

Decomposition and Interpretation of Treatment Effects in Settings with Delayed Outcomes

Federico A. Bugni

Department of Economics

Northwestern University

federico.bugni@northwestern.edu

Ivan A. Canay

Department of Economics

Northwestern University

iacanay@northwestern.edu

Steve McBride

Head of Economy Science

Roblox

February 23, 2023

Abstract

This paper studies settings where there is interest in identifying and estimating an average causal effect of a binary treatment on an outcome of interest, under complete randomization or selection on observables assumptions. The outcome does not get immediately realized after treatment assignment, a feature that is ubiquitous in empirical settings, creating a time window in between the treatment and the realization of the outcome. The existence of such a time window, in turn, opens up the possibility of other observed endogenous actions to take place and affect the interpretation of popular parameters, including the average treatment effect. In this context, we study several regression-based estimands that are routinely used in empirical work, and present five results that shed light on how to interpret them in terms of *ceteris paribus* effects, indirect causal effects, and selection terms. Our three main takeaways are the following. First, the three most popular estimands do not satisfy what we call *strong sign preservation*, in the sense these estimands may be negative even when the treatment positively affects the outcome for any possible combination of other actions. Second, the by-far most popular estimand that “controls” for the other actions in the regression does not improve upon a simple comparisons of means in the sense that negative weights multiplying relevant *ceteris paribus* effects become more prevalent. Finally, while non-parametric identification of the effects we study is straightforward under our assumptions and follows from saturated regressions, we also show that linear regressions that correctly control for the other actions by stratifying lead to estimands that always satisfy *strong sign preservation*.

KEYWORDS:

JEL classification codes: C12, C14

1 Introduction

We study settings where the analyst is interested in identifying and estimating an average causal effect of a binary treatment on an outcome of interest, where the treatment status could be determined in the context of a randomized controlled experiment or in the context of an observational study under conditional independence assumptions. We focus on settings where the outcome of interest does not get immediately realized after treatment assignment, a feature that is ubiquitous in empirical settings, including applications in development economics (Beaman et al., 2013), clinical trials (Moderna, 2021), and a variety of applications in the industry (Akhtari et al., 2021). The delay in the realization of the outcomes creates a time window in between the treatment assignment and the realization of the outcome that, in turn, opens up the possibility for other observed endogenous actions to take place before the outcome is finally realized; see Figure 1 for a graphical representation. In this context, we study the interpretation of several popular estimands that arise from running regressions of the outcome on the treatment and different ways of “controlling” for the other actions. Some of these estimands are not only popular in the economics literature, see, e.g., Fagereng et al. (2021); Heckman et al. (2013); Chernozhukov et al. (2021), but are also widely used across other social sciences, like psychology and political science, as shown by the large number of citations associated with the regression approach popularized by Baron and Kenny (1986). For each of these estimands, our results present a decomposition that facilitates their interpretation in terms of ceteris paribus effects of the treatment on the outcomes, indirect effects caused by the other actions, and selection terms; and provide a framework that allows us to clarify under what type of conditions the practice of “controlling” for the presence of other actions leads to estimands that admit the desired interpretation.

The main findings of this paper can be grouped into three sets of results. First, the standard practice of studying estimands that arise from a regression of an outcome on the treatment, with or without “controlling” for the other actions in such regressions, does not satisfy what we call *strong sign preservation*. Strong sign preservation, formally defined in Definition 3.3, is satisfied when an estimand that intends to capture a ceteris paribus causal effect of a treatment on an outcome is positive when the effect of the treatment on the outcome is positive conditional on *all* possible values of the other actions. Failure to satisfy strong sign preservation introduces a Simpson’s Paradox-like sign reversal where the estimands may be negative even when the treatment positively affects the outcome for any possible combination of other actions. Second, the most popular estimand that linearly controls for the other actions in the regression, and that we label the long regression, does not generally provide benefits relative to the short regression that includes no controls whatsoever. More concretely, while both the short and long regressions do not satisfy strong sign preservation, the estimand associated with the long regression admits a decomposition in terms of weighted averages of well

defined causal effects but where the weights could potentially be negative. This feature introduces yet another source that may separate the sign of the estimand from the sign of *ceteris paribus* causal effects. Notably, this feature does not depend on whether the regression includes interaction terms between the treatment and the other actions. Finally, while non-parametric identification of the effects we study is straightforward under the stronger version of our identifying assumptions and immediately follows from a saturated regressions, we also show that linear regressions that correctly control for the other actions by proper stratification (an approach we label as the strata fixed effects regression due to its connection with the standard practice of including strata fixed effects in randomized controlled trials with covariate adaptive randomization; see [Bugni et al. \(2018, 2019\)](#)) always lead to estimands that automatically satisfy strong sign preservation.

The decompositions we derive for each of the five estimands we study can be interpreted as decomposing a “total” effect into a “direct” and an “indirect” effect (and possibly “selection” effect depending on the assumptions), and so our results are linked to the vast literature on mediation analysis, see, e.g., [Baron and Kenny \(1986\)](#), [Pearl \(2001\)](#), [Robins \(2003\)](#), [Imai et al. \(2010\)](#), [Glynn \(2012\)](#), and [Remark 2.1](#) for a discussion. However, as opposed to the literature on mediation that studies the type of identifying assumptions that would identify the causal effects of the so-called *mediators*, which in our context would simply be the other actions taken before the outcome is realized, here our goal is not to identify these indirect effects but rather to gain a better understanding of how to properly interpret certain popular estimands of the effect of the treatment on the outcome. Despite the different goals, our results directly speak to the dominant empirical practice in the mediation literature and point to limitations in the scope of such a practice. In particular, we show that the validity of the long regression (and the long regression with interaction terms) not only depends on identifying assumptions like sequential ignorability as shown by [Imai et al. \(2010\)](#), but also on the mediators being scalar valued and a correctly specified linear model for potential outcomes.

Beyond the literature on mediation analysis, we are also not the first to acknowledge the importance of the distinctions between “partial” and “total” causal effects in the economics literature, where early discussions include those in [Manski \(1997\)](#) and [Heckman \(2000\)](#); see [Remark 3.1](#). For example, [Manski \(1997, page 1321 and 1323\)](#) provides two interpretations of potential outcomes (one that keeps other actions fixed and another one that lets the other actions change in response to treatment) and clarifies that the interpretation of treatment effects depends on how we think about potential outcomes. Our goal in this paper is not dwell on discussions about relative merits of partial or total effects but rather seek to understand when and how commonly used estimands in empirical work admit either one of these interpretations. Finally, an important characteristic of the setting we consider is that the other actions are observed

by the analyst and this separates the types of concerns we focus on here from those that are related with unobserved factors that may affect the outcome and may be affected by the treatment, see, e.g., [Rosenzweig and Wolpin \(2000\)](#).

The remainder of the paper is organized as follows. Section 2 introduces the basic notation and provides motivating examples. Section 3 defines the main concepts we use throughout the paper, including partial causal effects, direct causal effects, and strong sign preservation. Section 4 introduces the five estimands we study, the short regression, the long regression, the long regression with interactions, the strata fixed effects regression, and the saturated regression, and then presents the main results on how each of these estimands admit different decompositions into direct, indirect, and selection effects. Finally, Section 5 concludes.

2 Setup and Notation

Consider a setting where Y denotes the (observed) outcome of interest, \tilde{A} denotes a vector of actions that the individuals or units under study may take, and X denotes observed covariates that include features beyond actions. We partition the vector of actions \tilde{A} into the main action of interest, denoted by D , and “other” actions, denoted by A , i.e.,

$$\tilde{A} = (D, A) \in \mathcal{D} \times \mathcal{A}. \quad (1)$$

All actions are assumed to be discrete, with the main action, D , further assumed to be binary, i.e., $\mathcal{D} \equiv \{0, 1\}$. The other actions, A , are a K dimensional vector taking values in

$$\mathcal{A} \equiv \{a = (a_1, \dots, a_K) : a_j \in \mathcal{A}_j \text{ and } j = 1, \dots, K\}, \quad (2)$$

where $\mathcal{A}_j \equiv \{0, 1, \dots, \bar{a}_j\}$ for $\bar{a}_j \geq 1$.

The setting we study in this paper is one with the following characteristics. First, the analyst controls the action of interest D via a randomized controlled experiment (or, alternatively, by an exogeneity assumption like selection on observables). We therefore alternatively call this action the “treatment”. Second, the outcome Y is not instantaneous and takes some time to get realized within the timeline of the experiment. In the period in-between the treatment assignment and the realization of the outcome, the other actions contained in A get chosen by the units participating in the experiment. Figure 1 illustrates the setting and Examples 2.1-2.3 provide concrete empirical situations where the setting we study currently applies.

Example 2.1 (Agriculture). [Beaman et al. \(2013\)](#) conduct a field experiment that provides free fertilizer to women rice farmers in Mali to measure how farmers choose to use the fertilizer, what changes they make to their agricultural practices, and the overall

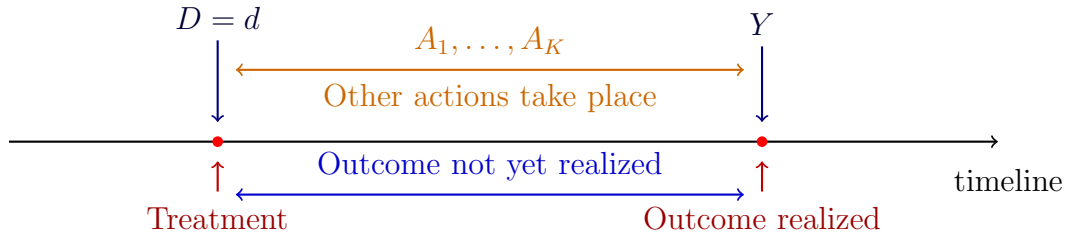


Figure 1: Timeline of actions. The first action, D , is randomly assigned in the context of the experiment (what we call the treatment). The outcome is not instantaneous and may take a short or long period of time to get realized. In the meantime, units choose the value of the other actions D, \dots, A_K .

in part on profitability. The authors distributed the fertilizer in May 2020 and conducted two follow up surveys, one in August 2020 and one in December 2020, right after the harvest. Using the notation in Figure 1, in this example D would be an indicator of whether the farmer received free fertilizer, Y would be a measure of output like crop yield or just profits, and A would include all relevant complementary agricultural inputs, such as labor, herbicides, and water usage. ■

Example 2.2 (Covid). The clinical trial run by Moderna (2021) studied the efficacy of the Moderna COVID-19 vaccine against SARS-CoV-2 infections. Participants in the study were randomized to Immediate Vaccination Group 1 (receiving the Moderna COVID-19 Vaccine on Day 1 and Day 29) or Standard of Care Group 2, with vaccination given at months 4 and 5. During the months following vaccinations, participants received visits that checked for infections and could include blood collection, nasal swab, SARS-CoV-2 screening, COVID-19 symptom check, and questionnaires. Using the notation in Figure 1, in this example D would be an indicator of whether the participant received a vaccine, Y would be an indicator of whether the participant got infected within the 4 months of the study, and A would include other actions taken by the participants that could affect infection rates, like whether the participants wear masks in public, whether the participants avoid large gatherings, etc. ■

Example 2.3 (AirBnB). Online platforms often allow customers to take a variety of actions and it becomes important to understand how much value to the company specific actions may bring (e.g, buying an item, leaving a review, making a reservation, ordering a delivery, renting a movie, etc); see Jain and Singh (2002) for a review of the literature on life-time valuation in marketing. For example, Akhtari et al. (2021) discuss how AirBnB measures the short and long-term value of actions and events that take place on their platform. These actions could be a guest making a booking or a host adding amenities to their listing, among many others. While it is often possible to rely on randomized experiments to measure the causal effects of some of these actions, see, e.g., Huang et al. (2018), there are others that are difficult to evaluate using experiments due to ethical, legal, or user experience concerns. In this cases it is common practice

to rely on selection on observables assumptions, as those discussed in the next section, and focus on a given action of interest at a time. For example, the main case discussed in [Akhtari et al. \(2021\)](#) is the valuation of the long-term impact of a guest making a booking at AirBnB. Using the notation in [Figure 1](#), in this example D would be an indicator for whether the customer made a booking on the platform, Y would be the revenue crated by that customer over 365 days, and A would capture other relevant actions, most notably, cancellations, leaving a review on the platform, etc. ■

Remark 2.1. What we call the other actions in [Figure 1](#) can be alternatively labeled as “mediator” variables since one could define mediators as any post-treatment variables that occurs before the outcome is realized, see, e.g., [Baron and Kenny \(1986\)](#), [Pearl \(2001\)](#), [Robins \(2003\)](#), [Imai et al. \(2010\)](#), [Glynn \(2012\)](#), and [Hernan and Robins \(2023, Ch. 23\)](#) for a recent book treatment. However, our work deviates from this literature in two important ways. First, while the literature on causal mediation analysis focuses on the identification of causal effects induced by mediators, our focus in this paper is to understand whether common estimands that are used to capture causal effects of main action D on the outcome Y admit clear interpretations through the lens of total and direct effects. Second, our decompositions in terms of direct and indirect effects are defined in terms of potential values for all of the actions, including those that may be labeled as mediators, and this implies “indirect” effects in our context do not coincide with the definition of indirect effects in the mediation literature but rather with the so-called “controlled” effects discussed by [Pearl \(2001\)](#) and [Robins \(2003\)](#); see [Remark 3.2](#) for additional discussion on this distinction. It is worth noting, however, that several of our results have implications for the causal mediation literature and we discuss these implications as we present our main results. ■

We denote potential outcomes by $Y(d, a)$ and their expectation by

$$\mu(d, a) \equiv E[Y(d, a)] . \tag{3}$$

Depending on the setting, we may expand a into (a_1, \dots, a_K) and write $Y(d, a_1, \dots, a_K)$ instead of $Y(d, a)$, although we prioritize the more concise notation whenever possible. We also introduce the concept of a pooled potential outcome to isolate the counterfactual outcome associated with the main action of interest (the treatment),

$$Y(d) = \sum_{a \in \mathcal{A}} Y(d, a) I\{A = a\} . \tag{4}$$

Finally, the observed outcome Y is related to potential outcomes by the relationship

$$Y = \sum_{(d,a) \in \mathcal{D} \times \mathcal{A}} Y(d, a) I\{(D, A) = (d, a)\} . \tag{5}$$

With this notation, we can state our basic maintained assumption as follows, where we denote by X the covariates or pre-determined variables.

Assumption 2.1. For all $d \in \mathcal{D}$ it follows that $D \perp Y(d) \mid X$.

Assumption 2.1 can be obtained by the design of the experiment (as in Examples 2.1 or 2.2) or by relying on a rich set of covariates that would make the exogeneity requirement credible (as in Example 2.3). In general, and as we will discuss in the next sections, this assumption will not be enough to identify causal effects of interest but we take it as the natural starting point of our analysis. For this reason, we also consider a stronger version of this assumption that requires conditional exogeneity of potential outcomes with respect to A as well. Formally, we state this assumption as follows.

Assumption 2.2. For all $(d, a) \in \mathcal{D} \times \mathcal{A}$ it follows that $(D, A) \perp Y(d, a) \mid X$.

Assumption 2.2 essentially re-interprets the problem as a problem of multiple conditionally exogenous treatments where, out of all possible treatments (D, A) , the analyst is interested in the effect of D only. As such, it may not be credible in settings where a randomized controlled experiment randomly assigned D across units but not A , as in Examples 2.1 or 2.2, but it is often invoked in settings where the main identification argument relies on selection on observables, as in Example 2.3. The assumption is strong in the sense that it is sufficient to non-parametrically identify $\mu(d, a)$ from the data, but the results in this paper show that the assumption is not strong enough to deliver a clean interpretation to popular estimands that are often used in practice. In addition, this assumption is implied by the so-called sequential ignorability assumption, a commonly used assumption in the literature on mediation analysis; see Section 4.2.

Remark 2.2. All the parameters we consider in this paper are characterized by conditional means. It would then be sufficient for all of our formal results to work with versions of Assumptions 2.1 and 2.2 that only require conditional mean independence, as opposed to full independence. We choose to maintain full independence for clarity.

3 Causal Treatment Effects

We start by discussing the type of counterfactual treatment effects that could be of interest to the researcher in the canonical setting where D is binary. Viewing $Y(d, a)$ as a function of two types of actions immediately suggests that there could be partial effects, total effects, direct effects, and indirect effects, all which may or may not be of interest in the context of a concrete application. Understanding the variety of causal effects that one could describe in turn will help us provide representations and interpretations of commonly used target parameters, like the average treatment effect (ATE), in terms of

these types of causal effects. We start with what is perhaps one of the most natural types of *ceteris paribus* effects in Definition 3.1.

Definition 3.1 (Average Partial Causal Effect). An **average partial causal effect** of D on the outcome Y is any difference of the form $\mu(d, a) - \mu(d', a)$, where the value $a \in \mathcal{A}$ is kept constant.

Definition 3.1 defines an average partial causal effect of the main action as a mean comparison that keeps the value of the other actions unchanged in both states of the comparison. In the context of Example 2.1 it would capture the average causal effect of using fertilizer on the crop yield, while keeping other inputs, like labor, herbicides, and water usage, constant in the counterfactual comparison.

The *ceteris paribus* effect in Definition 3.1 could also be defined conditional on certain events or subpopulations. In order to account for this, we also consider the concept in Definition 3.1 conditional on some conditioning set Ω , i.e.,

$$E[Y(d, a) - Y(d', a) \mid \Omega] , \tag{6}$$

where Ω is a function of (D, A, X) . For example, $\Omega = I\{D = 1\}$ would lead to an average partial causal effect on the treated and $\Omega = I\{X = x\}$ would lead to an average partial causal effect for units with covariates x .

The definition of a partial causal effect for the main action delivers a potentially different causal effect for each possible value of the other actions or, alternatively, provides a collection of partial causal effects indexed by $a \in \mathcal{A}$. While the goal could just be to identify such a collection of effects, in many settings it may be natural to aggregate this collection of partial effects in a way that summarizes the effect of the main action on the outcome of interest. The following definition defines a direct causal effect as any weighted average of partial causal effects.

Definition 3.2 (Average Direct Causal Effect). The **average direct causal effect** of D on the outcome Y is any convex combination of average partial causal effects of D on Y . That is,

$$\sum_{a \in \mathcal{A}} \omega(a) (\mu(d', a) - \mu(d, a)) , \tag{7}$$

where $\omega(a) \in [0, 1]$ for all $a \in \mathcal{A}$ and $\sum_{a \in \mathcal{A}} \omega(a) = 1$.

In the context of Example 2.1 with A only capturing low and high water usage for simplicity, the parameter (7) combines the average causal effect of using fertilizer on the crop yield for units with high water usage, say $A = 1$, and units with low water usage, say $A = 0$. Average direct causal effects could be defined conditional on a set Ω by weighting terms like those in (6). The definition does not determine how these two

groups are weighted, but it does require that no group gets a negative weight. In this sense, any average direct causal effect satisfies what we call *strong sign preservation*, as defined below.

Definition 3.3 (Strong Sign Preservation). A parameter Δ that measures a causal effect of the main action D on the outcome Y satisfies **strong sign preservation** if

$$\mu(d', a) - \mu(d, a) > 0 \text{ for all } a \in \mathcal{A} \text{ implies } \Delta > 0 .$$

In the context of Example 2.1 with A only capturing low and high water usage, strong sign preservation of a parameter Δ implies that whenever fertilizers improve the expected crop yield both for units with high water usage and units with low water usage, Δ should be positive as well. As the name suggests, strong sign preservation does not allow for the possibility of what it is typically referred to as *sign reversal*, understood as a situation where $\Delta < 0$ when $\mu(d', a) - \mu(d, a) > 0$ for all $a \in \mathcal{A}$.

Remark 3.1. While strong sign preservation may be perceived as a key requirement for parameters that intend to identify partial causal effects, it may not be a reasonable requirement in settings where the counter-factual question of interest involves total effects, as introduced and discussed in the next section. The distinctions between “partial” and “total” causal effects have appeared in the literature in a variety of contexts, even beyond the mediation analysis literature discussed in Remark 2.1, where Pearl (2001) and Robins (2003) provide comprehensive treatments on these distinctions. For example, Heckman (2000) defines a causal effect as a partial derivative and states that while the assumption that an isolated action can be varied independently of others is strong but “...essential to the definition of a causal parameter”. In fact, Heckman (2000, page 47) writes “Defining causality within a model is relatively straightforward when the causes can be independently varied. Defining causality when the causes are interrelated is less straightforward and is a major achievement of econometrics”. Manski (1997, page 1321 and 1323), in turn, provides two interpretations of potential outcomes (one that keeps other actions fixed and another one that lets the other actions change in response to the main action) and clarifies that the interpretation of treatment effects depends on how we think about potential outcomes. Here, we do not dwell on discussions about relative merits of partial or total effects but rather seek to understand whether commonly used estimands in empirical work admit either one of these interpretations under different assumptions.

Remark 3.2. Our definitions of average causal partial effects and average direct causal effects are not analogous to the notions of total causal effects, causal mediation effects, and natural direct effects that are commonly used in the mediation analysis literature. For example, the average natural direct effect corresponds to $E[Y(1, A(1)) - Y(0, A(0))]$ in our notation, where $A(d)$ denotes potential outcomes for the actions A as a function of

the treatment d . Contrary to these type of effects that are defined in terms of potential actions, $A(d)$, the effects we focus on in this paper are defined in terms of specific values of the actions A , say $A = a$ for any $a \in \mathcal{A}$, and are therefore analogous to the notions of a controlled direct (or total) effect that have been discussed in Pearl (2001); Robins (2003), among others.

4 Decomposing Common Estimands

In this section we analyze five natural and highly popular estimands that are intended to capture treatment effects of D on Y . For each of these estimands we derive a decomposition in terms of parameters that can be labeled according to Definitions 3.1 and 3.2 and discuss under what assumptions they can be interpreted as intended. In order to keep our exposition as simple as possible, from here on we ignore the role of the covariates, X , in the type of regressions we consider. This could be interpreted as a situation where the covariates are discrete, and the regressions are viewed as within cell regressions with cells given by $X = x$, or more generally where the covariates have been already accounted for by other means, like clustering or via a partially linear model, among many possibilities.

The first such estimand is the usual difference in means, which we write here as the slope coefficient Δ_{short} in a regression (projection) of Y on D and a constant term,

$$\text{short regression: } Y = \beta + \Delta_{\text{short}}D + U, \quad (8)$$

where $E[UD] = 0$ by properties of projections and $E[U|D] = 0$ follows from D being binary. We call this the short regression.

The second estimand is the slope coefficient D in a linear regression of Y on D , a constant term, and the K actions A_1, \dots, A_K ,

$$\text{long regression: } Y = \Delta_{\text{long}}D + \theta_0 + \sum_{j=1}^K \theta_j A_j + V, \quad (9)$$

where $E[VD] = E[VA_j] = 0$ by properties of projections. We call this the long regression.

The third estimand is the slope coefficient D in a linear regression of Y on D , a constant term, the K actions A_1, \dots, A_K , and their interactions with D ,

$$\text{interaction regression: } Y = \Delta_{\text{inter}}D + \theta_0 + \sum_{j=1}^K \theta_j A_j + \sum_{j=1}^K \lambda_j A_j D + e, \quad (10)$$

where $E[eD] = E[eA_j] = E[eA_j D] = 0$ by properties of projections. We call this the

interaction regression. Note that this is not a fully saturated regression in general, since the random variables A_j take values in $\mathcal{A}_j \equiv \{0, 1, \dots, \bar{a}_j\}$ and \bar{a}_j is allowed to be greater than 1.

The fourth estimand is the slope coefficient D in a regression of Y on D and a set of indicator functions for all the values that A takes,

$$\textbf{sfe regression: } Y = \Delta_{\text{sfe}}D + \sum_{a \in \mathcal{A}} \theta(a)I\{A = a\} + \nu, \quad (11)$$

where $E[\nu D] = E[\nu I\{A = a\}] = 0$ by properties of projections. Note that this is a regression of Y on D with “strata fixed effects”, where the event $\{A = a\}$ defines a stratum for each value of a . As a result, we call this the strata fixed effect (sfe) regression. The regression in (11) can also be interpreted as a regression of Y on D and A where the linear component that captures the effect of A is fully saturated; i.e.,

$$\sum_{a \in \mathcal{A}} \theta(a)I\{A = a\} = \sum_{a_1 \in \mathcal{A}_1} \sum_{a_2 \in \mathcal{A}_2} \cdots \sum_{a_K \in \mathcal{A}_K} \theta(a)I\{A_1 = a_1\}I\{A_2 = a_2\} \cdots I\{A_K = a_K\}.$$

The last set of estimands are the slope coefficients $\Delta_{\text{sat}}(a)$, for $a \in \mathcal{A}$, in a saturated regression of Y on a set of indicator functions for all the values that A takes and their interactions with D ,

$$\textbf{sat regression: } Y = \sum_{a \in \mathcal{A}} \gamma(a)I\{A = a\} + \sum_{a \in \mathcal{A}} \Delta_{\text{sat}}(a)I\{A = a\}D + \epsilon, \quad (12)$$

where $E[\epsilon D I\{A = a\}] = 0$ by properties of projections and $E[\epsilon | D, I\{A = a\}] = 0$ follows from D and $I\{A = a\}$ being binary for all $a \in \mathcal{A}$. We call this the saturated regression.

Remark 4.1. The use of short, long, and interaction regressions in the social science literature is ubiquitous. When [Glynn \(2012\)](#) discusses the popularity of these regressions, he writes that long regressions are so pervasive within the social science and empirical mediation literature that “examples are too numerous to cite.” Indeed, [Baron and Kenny \(1986\)](#), the paper that largely established the use of these and related regressions, has over 115,000 citation as of 2022.

4.1 Short regression

The short regression is algebraically very simple, so we build up towards the main result introducing the main concepts and notation along the way. The other regressions, on the contrary, have derivations that are more opaque, and so in those cases we first present the formal results and then discuss their interpretation.

The slope coefficient Δ_{short} in (8) equals $\Delta_{\text{short}} = E[Y|D = 1] - E[Y|D = 0]$ by

elementary arguments. If we define

$$\pi_d(a) \equiv P\{A = a|D = d\} , \quad (13)$$

and note that

$$E[Y|D = d] = \sum_{a \in \mathcal{A}} E[Y(d, a)|D = d, A = a]\pi_d(a) ,$$

we can decompose Δ_{short} into the following three terms,

$$\Delta_{\text{short}} = \Delta_{\text{dce}}^{\text{s}} + \Delta_{\text{ind}}^{\text{s}} + \Delta_{\text{sel}}^{\text{s}} \quad (14)$$

where

$$\Delta_{\text{dce}}^{\text{s}} \equiv \sum_{a \in \mathcal{A}} \pi_1(a) E[Y(1, a) - Y(0, a)|D = 1, A = a] \quad (15)$$

$$\Delta_{\text{ind}}^{\text{s}} \equiv \sum_{a \in \mathcal{A}} (\pi_1(a) - \pi_0(a)) (E[Y(0, a)|D = 0, A = a] - E[Y(0, 0)|D = 0, A = 0]) \quad (16)$$

$$\Delta_{\text{sel}}^{\text{s}} \equiv \sum_{a \in \mathcal{A}} \pi_1(a) (E[Y(0, a)|D = 1, A = a] - E[Y(0, a)|D = 0, A = a]) . \quad (17)$$

Each of the three terms in the above decomposition for Δ_{short} have a clear interpretation and show that there are two channels of endogeneity that are introduced by the fact that the other actions, A , take place in-between the treatment assignment and the realization of the outcome of interest. The term in (15), $\Delta_{\text{dce}}^{\text{s}}$, captures an average direct causal effect on the treated, as in Definition 3.2. Note that this term conditions on $\Omega = I\{D = 1, A = a\}$ and so it is a conditional effect like those defined in (6). The term in (16), $\Delta_{\text{ind}}^{\text{s}}$, admits a clean interpretation for each value $a \in \mathcal{A}$ under additional assumptions we introduce below. Without additional assumptions, this term is a type of “indirect” effect that contains the product of the difference in conditional probabilities, $\pi_1(a) - \pi_0(a)$, and a term that confounds the average partial causal effect of A moving from 0 to a on Y , with selection that arises from the distinct conditioning sets $\{D = 0, A = a\}$ and $\{D = 0, A = 0\}$. Finally, the term in (17), $\Delta_{\text{sel}}^{\text{s}}$, is a selection term that captures the fact that the action $A = a$ may not be independent of $Y(0, a)$ and D .

Two points are worth highlighting. First, the above decomposition does not invoke either Assumption 2.1 or Assumption 2.2. Importantly, while Assumption 2.1 guarantees that

$$\Delta_{\text{short}} = E[E[Y|D = 1, X] - E[Y|D = 0, X]] = E[Y(1) - Y(0)] , \quad (18)$$

where $Y(d)$ are the pooled potential outcomes in (4), it is not enough to characterize Δ_{short} as an average direct causal effect or as a parameter that satisfies strong sign preservation. In particular, the two endogeneity channels, $\Delta_{\text{ind}}^{\text{s}}$ and $\Delta_{\text{sel}}^{\text{s}}$, in the decom-

position of Δ_{short} could be positive or negative and, more importantly, lead to Δ_{short} to have opposite sign to $\Delta_{\text{dce}}^{\text{s}}$. Second, the two endogeneity channels, $\Delta_{\text{ind}}^{\text{s}}$ and $\Delta_{\text{sel}}^{\text{s}}$, are conceptually different. While the channels entering the term $\Delta_{\text{ind}}^{\text{s}}$ are difficult to shut down, the selection term $\Delta_{\text{sel}}^{\text{s}}$ can be set equal to zero under Assumption 2.2.

Under Assumption 2.2 the three terms entering the decomposition for Δ_{short} simplify in the following way,

$$\Delta_{\text{dce}}^{\text{s}} = \sum_{a \in \mathcal{A}} \pi_1(a) (\mu(1, a) - \mu(0, a)) \quad (19)$$

$$\Delta_{\text{ind}}^{\text{s}} = \sum_{a \in \mathcal{A}} (\pi_1(a) - \pi_0(a)) (\mu(0, a) - \mu(0, 0)) \quad (20)$$

$$\Delta_{\text{sel}}^{\text{s}} = 0 . \quad (21)$$

That is, the first two terms are now a function of the unconditional expectations $\mu(d, a)$ defined in (3), and the selection term $\Delta_{\text{sel}}^{\text{s}}$ is no longer present. Importantly, the term $\Delta_{\text{ind}}^{\text{s}}$ is still part of the decomposition since Assumption 2.2 does not restrict how A may affect outcomes, so that $\mu(0, a) - \mu(0, 0) \neq 0$, nor does it affect how the main action may affect the other ones, so that $\pi_1(a) - \pi_0(a) \neq 0$. Aside from removing the term capturing selection bias, Assumption 2.2 also delivers a clean interpretation to the indirect effects captured by $\Delta_{\text{ind}}^{\text{s}}$. Each summand in $\Delta_{\text{ind}}^{\text{s}}$ contains the average partial causal effect of A moving from 0 to a on Y , $\mu(0, a) - \mu(0, 0)$, multiplied by the difference $\pi_1(a) - \pi_0(a)$, which admits a causal interpretation of an average direct causal effect of D on A under the additional assumption $A(d) \perp D$ - where here we use $A(d)$ to denote potential outcomes for the other actions. As an illustrative example, suppose that

$$\begin{aligned} \mu(1, a) - \mu(0, a) &> 0 \quad \forall a \in \mathcal{A} \\ \mu(0, a) - \mu(0, 0) &< 0 \quad \forall a \in \mathcal{A} \\ \pi_1(a) - \pi_0(a) &> 0 \quad \forall a \in \mathcal{A} . \end{aligned}$$

That is, D increases the mean outcome for any value of a and also increases the probability that the other actions take the value a , while the other actions decrease the mean outcome relative to $a = 0$ under no treatment ($D = 0$). In this case, $\Delta_{\text{ind}}^{\text{s}}$ is immediately negative and the sign of Δ_{short} gets determined by the relative magnitudes of $\Delta_{\text{ind}}^{\text{s}}$ and $\Delta_{\text{dce}}^{\text{s}}$. In the other extreme where $\mu(0, a) - \mu(0, 0) > 0$ for all $a \in \mathcal{A}$, the decomposition shows that Δ_{short} amplifies the average direct causal effect of D on Y by taking a piece $\pi_1(a) - \pi_0(a)$ of the effect of A on Y , $\mu(0, a) - \mu(0, 0)$. It follows that Δ_{short} measures a total causal effect of D on Y .

We can interpret the terms entering the decomposition of Δ_{short} in (14) in the context of Examples 2.1-2.3. For example, consider the case of Example 2.1 where Y is crop yield, D is an indicator for the use of fertilizer, and A is for simplicity an indicator

for high water usage. In this setting, $\Delta_{\text{dce}}^{\text{s}}$ captures the average direct causal effect of using fertilizer on the crop yield, where the effect weights units with high and low water usage according to the respective probabilities of these actions happening for the treated, $\pi_1(a)$. The term $\Delta_{\text{ind}}^{\text{s}}$, in turn, captures a piece of the causal effect of water usage on crop yield that depends on the magnitude of differential water usage between the treated and the untreated. If water usage causally improves crop yield in the absence of fertilizer, and getting an exogenous fertilizer incentivizes units to increase their water usage, this term would be positive.

The following theorem summarizes our discussion above.

Theorem 4.1. Consider the short regression in (8) and assume $P\{D = 1\} \in (0, 1)$. Then, Δ_{short} can be decomposed as in (14)-(17). If Assumption 2.2 holds, then $\Delta_{\text{sel}}^{\text{s}} = 0$ and $\Delta_{\text{dce}}^{\text{s}}$ and $\Delta_{\text{ind}}^{\text{s}}$ simplify to the expressions in (19) and (20).

Remark 4.2. It is important to note that, even under the stronger exogeneity condition in Assumption 2.2, the parameter Δ_{short} does not satisfy strong sign preservation as defined in Definition 3.3. Indeed, it is certainly possible that $\mu(1, a) - \mu(0, a) > 0$ for all $a \in \mathcal{A}$ and yet $\Delta_{\text{short}} < 0$ due to $\Delta_{\text{ind}}^{\text{s}} < -\Delta_{\text{dce}}^{\text{s}} < 0$. ■

Remark 4.3. Under Assumption 2.2 and $A(d) \perp D$, Δ_{short} is a linear combination of average partial causal effects and it captures a “total” effect rather than a “partial” effect, as discussed in Remark 3.1. To understand this, notice that Δ_{ind} in (20) is the product of $\pi_1(a) - \pi_0(a)$ and $\mu(0, a) - \mu(0, 0)$ for each $a \in \mathcal{A}$. Both of these terms are partial effects, where $\pi_1(a) - \pi_0(a) = E[A(1) - A(0)]$ is the average partial effect of D on A and $\mu(0, a) - \mu(0, 0)$ is the average partial effect of moving A from 0 to a on the outcome Y for units with $D = 0$. With this interpretation, Δ_{short} captures a total effect of D on Y that adds up the direct effect of D on Y , captured by $\Delta_{\text{dce}}^{\text{s}}$, and the indirect effect that D has on Y via its effect on A and how A affects Y . This distinction between partial and total effects mimics the usual one associated with total and partial derivatives in mathematical analysis. Whether total or partial effects are relevant in the context of a given application has been already discussed elsewhere; see, for example, Manski (1997); Heckman (2000); Imai et al. (2010); Glynn (2012). Our main goal here is to clarify the interpretation of estimands like Δ_{short} in terms of these notions. ■

In what follows we prioritize results that hold under Assumption 2.2, with discussions on how the main implications would be affected if Assumption 2.2 is replaced by its weaker analog, Assumption 2.1. In general, moving from Assumption 2.2 to Assumption 2.1 in all of the cases we study below leads to the same implication: interpreting the estimands under consideration becomes difficult as a selection term, like $\Delta_{\text{sel}}^{\text{s}}$ above, becomes present. Indeed, Robins and Greenland (1992) argued early on in the empirical mediation literature that direct and indirect effects cannot be separated in randomized controlled trials without additional assumptions; a problem that at least within the

literature in development economics appears to be well understood, see, for example, Mel et al. (2009); Duflo et al. (2011); Beaman et al. (2013).

Settings where D is randomized in the context of an RCT, as in Examples 2.1 and 2.2, are the ones where Assumption 2.2 may be particularly difficult to defend. We note, however, that selection terms could disappear under alternative assumptions to Assumption 2.2, and we discuss some of these alternatives below. On the other hand, settings where the identification argument relies on selection on observables, as it is typically the case in applications in the industry, as in Example 2.3, are such that Assumption 2.2 may become more natural and may even be implied by the symmetric nature of D and A - as these are all actions that customers can take in the platform.

4.2 Long Regression

A seemingly natural, and certainly popular, way to mitigate the presence of indirect effects and obtain an estimand that satisfies strong sign preservation is to control for the other actions linearly as in (9); an approach we call the long regression. Our main result below shows that the slope coefficient Δ_{long} in (9) admits a decomposition similar to that derived by Δ_{short} , and thus includes a combination of direct effects and indirect effects. However, except in some special cases, the coefficients multiplying each average partial causal effect, as in Definition 3.1, could be negative and so Δ_{long} may be negative even in the absence of indirect effects. We formalize this below and provide a proof in Appendix A.

Theorem 4.2. Let Assumption 2.2 hold and assume that the covariance matrix of (D, A) is positive definite. Then, the coefficient Δ_{long} in (9) admits the decomposition

$$\Delta_{\text{long}} = \Delta_{\text{dce}}^1 + \Delta_{\text{ind}}^1, \quad (22)$$

where

$$\Delta_{\text{dce}}^1 \equiv \sum_{a \in \mathcal{A}} \omega_{\text{dce}}^1(a) (\mu(1, a) - \mu(0, a)) \quad (23)$$

$$\Delta_{\text{ind}}^1 \equiv \sum_{a \in \mathcal{A}} \omega_{\text{ind}}^1(a) (\mu(0, a) - \mu(0, 0)), \quad (24)$$

and $\{\omega_{\text{dce}}^1(a) : a \in \mathcal{A}\}$ and $\{\omega_{\text{ind}}^1(a) : a \in \mathcal{A}\}$ are as defined in Theorem A.1 and satisfy $\sum_{a \in \mathcal{A}} \omega_{\text{dce}}^1(a) = 1$ and $\sum_{a \in \mathcal{A}} \omega_{\text{ind}}^1(a) = 0$. Furthermore, the following statements are equivalent:

- (a) A are mutually exclusive binary variables, i.e., $\mathcal{A}_j = \{0, 1\}$ for $j = 1, \dots, K$ and $A_j A_l = 0$ for all $j, l = 1, \dots, K$ with $j \neq l$.

(b) For any distribution of (A, D) , $\omega_{\text{dce}}^1(a) \geq 0$ for all $a \in \mathcal{A}$.

(c) For any distribution of (A, D) , $\omega_{\text{ind}}^1(a) = 0$ for all $a \in \mathcal{A}$.

Theorem 4.2 shows that Δ_{long} can be decomposed into direct and indirect effects, but it leaves open the possibility that the coefficients entering each of these terms could, in general, be negative. An immediate implication is that, except in the special case where the actions in A are all mutually exclusive binary variables, which includes the case where A is a scalar binary variable as a special case, the term Δ_{dce}^1 could be negative even if $\mu(1, a) - \mu(0, a) > 0$ for all $a \in \mathcal{A}$. This is because $\omega_{\text{dce}}(a)$ may be negative for some $a \in \mathcal{A}$. As a result, Δ_{long} in general does not satisfy strong sign preservation for the following two reasons. First, it may be possible that $\Delta_{\text{ind}}^1 < -\Delta_{\text{dce}}^1$, so that the indirect effect dominates the direct effects. This phenomenon is the same as the one we discussed for the short regression. Second, even in the absence of indirect effects, where $\Delta_{\text{ind}}^1 = 0$, the term Δ_{dce}^1 could be negative by itself. This second possibility represents a stark distinction between the long regression estimand, Δ_{long} , and the short regression estimand, Δ_{short} .

Remark 4.4. Replacing Assumption 2.2 with Assumption 2.1 leads to a decomposition of Δ_{long} that introduces three changes relative to the one in Theorem 4.2. First, the term Δ_{dce}^1 becomes a linear combination of expectations that condition on $\{D = 1, A = a\}$. Second, the interpretation of Δ_{ind}^1 becomes convoluted for the same reasons discussed for $\Delta_{\text{ind}}^{\text{s}}$ before. Finally, the decomposition additionally includes a selection term that is conceptually identical to $\Delta_{\text{ind}}^{\text{s}}$ in the short regression. The details of these expressions are presented in Theorem A.1 in the Appendix. ■

The possibility of $\omega_{\text{dce}}(a)$ being negative for some $a \in \mathcal{A}$ does not rely on pathological data generating processes and may arise in rather simple settings under reasonable distributions for (A, D) . This raises an important red flag for the use of linear in A regressions, as they may lead to results that are quite difficult to interpret and, in general, do not offer an improvement relative to the short regression in (8). Below we illustrate this situation with two canonical simple cases: one where A is a scalar random variable taking multiple values, and another one where $A = (A_1, A_2)$ with A_1 and A_2 being binary. These same examples appear in the proof of Theorem 4.2.

Consider first the case where $A = A_1$ is a scalar random variable taking values in $\mathcal{A}_1 = \{0, 1, 2, \dots, \bar{a}_1\}$. The regression in (9) simplifies to

$$Y = \Delta_{\text{long}}D + \theta_0 + \theta_1 A_1 + V . \quad (25)$$

Theorem A.1 in the appendix provides general closed-form expressions for $\{\omega_{\text{dce}}^1(a) : a \in$

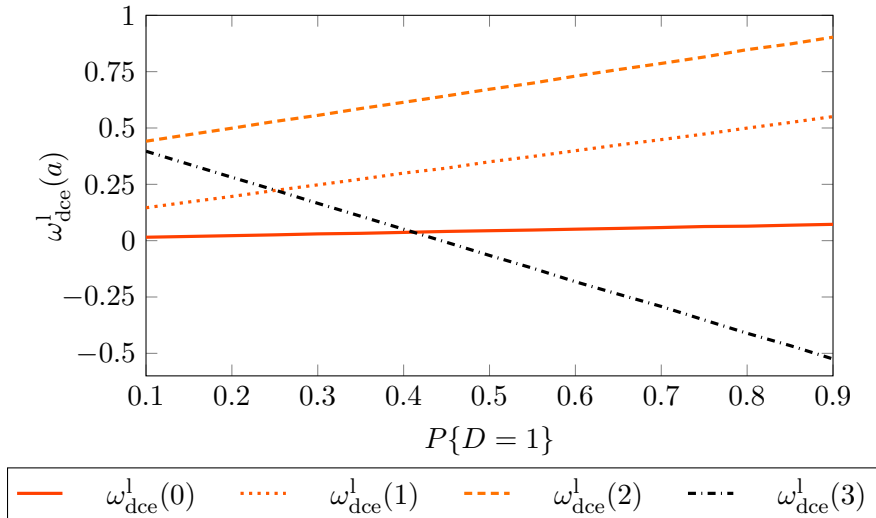


Figure 2: Weights $\omega_{\text{dce}}^1(a)$ as a function of p when $\{A|D=1\} \sim \text{Bi}(3, 0.8)$, and $\{A|D=0\} \sim \text{Bi}(3, 0.2)$

\mathcal{A} and $\{\omega_{\text{ind}}^1(a) : a \in \mathcal{A}\}$ that, when applied to this specific example, lead to

$$\omega_{\text{dce}}^1(a) \propto \left(\pi_1(a) - \frac{\text{Cov}(D, A_1)(a - E[A_1])}{\text{Var}(A_1)(1-p)} \right), \quad (26)$$

where $p = P\{D=1\}$. From this expression it follows that any distribution of (A, D) for which

$$\frac{\text{Cov}(D, A_1)(a - E(A_1))}{\text{Var}(A_1)(1-p)} > \pi_1(a),$$

would lead to negative weights. For example, consider the case where $p = 0.8$, $\bar{a}_1 = 3$, $\{A|D=1\} \sim \text{Bi}(3, 0.8)$, and $\{A|D=0\} \sim \text{Bi}(3, 0.2)$, where $\text{Bi}(n, \pi)$ denotes a Binomial distribution with n trials and success probability π . In this case, $\omega_{\text{dce}}^1(3) = -0.41 < 0$. Figure 2 plots the weights $\omega_{\text{dce}}^1(a)$ as a function of p and shows that $\omega_{\text{dce}}^1(3)$ is negative for any $p > 0.4$ in this example.

Next, consider the case where $A = (A_1, A_2)$ with A_1 and A_2 both being binary variables, so that $\mathcal{A} = \{(0, 0), (1, 0), (0, 1), (1, 1)\}$. The regression in (9) simplifies to

$$Y = \Delta_{\text{long}} D + \theta_0 + \theta_1 A_1 + \theta_2 A_2 + V.$$

The closed-form expressions for $\{\omega_{\text{dce}}^1(a) : a \in \mathcal{A}\}$ and $\{\omega_{\text{ind}}^1(a) : a \in \mathcal{A}\}$ derived in Theorem A.1 also simplify to this case and lead to simple conditions for which

$$\omega_{\text{dce}}^1(1, 1) = -\omega_{\text{dce}}^1(1, 0). \quad (27)$$

That is, whenever one of the average partial causal effects gets a positive weight, the other one necessarily gets a negative one. As an illustrative example, consider the case

where $\text{Cov}[A_1, A_2] = 0$,

$$P\{D = 1\} = P\{A_2 = 1\} = \frac{1}{2}, \quad P\{A_1 = 1 \mid D = 1\} = 2P\{A_1 = 1\}, \quad (28)$$

and

$$P\{A_1 = A_2 = 1 \mid D = 1\} = P\{A_1 = 1, A_2 = 0 \mid D = 1\} = \frac{1}{4}. \quad (29)$$

Using the expressions in Theorem A.1, we obtain $\omega_{\text{dce}}^1(1, 0) = -\omega_{\text{dce}}^1(1, 1) = -0.30$, which one more time illustrates that negative weights arise naturally in settings with non-pathological DGPs. In the proof of Theorem 4.2 we present even simpler counterexamples that also illustrate how the weights $\{\omega_{\text{ind}}^1(a) : a \in \mathcal{A}\}$ are generally non-zero and potentially negative, as well as how $\omega_{\text{dce}}^1(a)$ may be negative without necessarily satisfying (27).

Remark 4.5. It is important to note that, even under the stronger exogeneity condition in Assumption 2.2, the parameter Δ_{long} does not generally satisfy strong sign preservation as defined in Definition 3.3. Indeed, it is certainly possible that $\mu(1, a) - \mu(0, a) > 0$ for all $a \in \mathcal{A}$ and yet $\Delta_{\text{long}} < 0$ due to either $\Delta_{\text{ind}}^1 < -\Delta_{\text{dce}}^1 < 0$ or simply $\Delta_{\text{dce}}^1 < 0$ because of negative weights $\{\omega_{\text{dce}}^1(a) : a \in \mathcal{A}\}$. This second condition implies that *not even* Δ_{dce}^1 satisfies strong sign preservation and thus, in general, Δ_{long} does not offer much of a benefit relative to Δ_{short} , as Δ_{short} can at least be interpreted as a kind of “total” effect, as discussed in Remark 4.3. ■

Remark 4.6. As discussed in Remark 4.1, the long regression in (9) is used extensively in the social sciences and the mediation literature. In economics, Heckman et al. (2013, Eq. (6)) consider a long regression in the context of a more restrictive model for potential outcomes that are linear and separable in (d, a) . More recently, Fagereng et al. (2021, Eq. (7)) use the same mediation model from Heckman et al. (2013), in combination with the long regression in (9), to disentangle the average causal effect on outcomes into direct and indirect effects. The causal interpretation assigned to the estimands in these last set of papers is correct under the modeling assumptions for potential outcomes, despite both applications involving actions A that are multidimensional and non-mutually exclusive. Our results, however, imply that the main conclusions from such an analysis delicately rely on a linear model for $\mu(d, a)$ and do not generally extend to more general models for $\mu(d, a)$. ■

The results in Theorem 4.2 are novel to the best of our knowledge, though there are related results that differ in focus and scope. For example, Imai et al. (2010) study the interpretation of the long regression popularized by Baron and Kenny (1986) under an assumption they refer to as *sequential ignorability* and a linear model for potential outcomes. We state this assumption below.

Assumption 4.1 (Sequential Ignorability). Let $A(d)$ denote the potential outcome for

A and assume that

$$(Y(d', a), A(d)) \perp D \mid X = x \quad (30)$$

$$Y(d', a) \perp A(d) \mid D = d, X = x, \quad (31)$$

for $d, d' = 0, 1$ and all x , where in addition $0 < P\{D = d \mid X = x\}$ and $0 < P\{A(d) = a \mid D = d, X = x\}$ for all d, x , and a .

The results in [Imai et al. \(2010\)](#) about the long regression in (9) invoke

- (a) sequential ignorability,
- (b) a scalar random variable A (though not necessarily binary), and
- (c) a linear model for $\mu(d, a)$ in (d, a) .

Under (a)-(c) above, [Imai et al. \(2010, Theorem 2\)](#) shows that Δ_{long} identifies $\bar{\zeta} = \bar{\zeta}(1) = \bar{\zeta}(0)$ where

$$\bar{\zeta}(d) \equiv E[Y(1, A(d))] - E[Y(0, A(d))] = \sum_{a \in \mathcal{A}} (\mu(1, a) - \mu(0, a)) \pi_d(a), \quad (32)$$

and the equality follows from Assumption 4.1 implying $Y(d', a) \perp A(d) \mid X = x$; see Lemma B.2 in Appendix B. The linear model for $\mu(d, a)$ implies that the difference $\mu(1, a) - \mu(0, a)$ does not depend on the value of a , and so it is just a constant that we can denote by $\bar{\zeta}$ without loss of generality. The fact that $\bar{\zeta} = \bar{\zeta}(1) = \bar{\zeta}(0)$ then follows from $\sum_{a \in \mathcal{A}} \pi_d(a) = 1$ for $d \in \{0, 1\}$.

The additional assumptions (a)-(c) mentioned above have implications on the conclusions of Theorem 4.2, which does not invoke any of these assumptions. In particular, the linearity of $\mu(d, a)$ implies that Δ_{dce}^1 in (23) equals $\bar{\zeta} \sum_{a \in \mathcal{A}} \omega_{\text{dce}}^1(a) = \bar{\zeta}$ by the weights adding up to one according to Theorem 4.2. The same linearity assumption also implies that

$$\Delta_{\text{ind}}^1 \equiv \sum_{a \in \mathcal{A}} \omega_{\text{ind}}^1(a) (\mu(0, a) - \mu(0, 0)) \propto \sum_{a \in \mathcal{A}} \omega_{\text{ind}}^1(a) a = 0,$$

where the last equality follows from Theorem A.1 in the appendix. We conclude that Theorem 4.2 coincides with the results in [Imai et al. \(2010\)](#) in delivering Δ_{long} being equal to $\bar{\zeta}$ under the additional assumption that $\mu(d, a)$ is linear in (d, a) . This means that, while sequential ignorability is a stronger assumption than Assumption 2.2 (we prove this claim in Lemma B.1 in Appendix B), the main driving force of this result is the linear model for the potential outcomes or, equivalently, the linear model for $\mu(d, a)$. As we discussed in Remark 4.6, this linearity assumption has been used in economic applications, e.g., [Heckman et al. \(2013\)](#); [Fagereng et al. \(2021\)](#), and while it imposes

enough restrictions to provide a clean interpretation to the coefficient Δ_{long} , our results show that such a clean interpretation generally breaks down when $\mu(d, a)$ is not linear in (d, a) .

Remark 4.7. We note that while Theorem 4.2 is a result on how to properly interpret Δ_{long} in the context of a long regression, Imai et al. (2010, Theorem 2) is a result on the identification of *natural* indirect effects via the same type of regression and the additional conditions in (a)-(c) above. Related to our discussion in Remark 3.2, the *natural* indirect effect does not coincide with our notion of indirect effect in Theorem 4.2. To see the difference, note that the *natural* indirect effect, defined as $\bar{\delta}(d) = E[Y(d, A(1)) - Y(d, A(0))]$, can be written as

$$\bar{\delta}(d) = \sum_{a \in \mathcal{A}} (\mu(d, a) - \mu(d, 0))(\pi_1(a) - \pi_0(a)) , \quad (33)$$

under sequential ignorability, and is distinct from $\Delta_{\text{ind}}^!$ in (24) because $\omega_{\text{ind}}^!(a) \neq \pi_1(a) - \pi_0(a)$. Conceptually, the literature on mediation analysis defines an indirect effect as a target parameter and then determines conditions under which such indirect effects could be identified from the data. In contrast, we characterize the decomposition of estimands in terms of average direct causal effects, as defined in Definition 3.2, and then group the reminding terms as indirect or selection terms, depending on the case. We note, however, that our indirect effects coincide with those characterized by $\bar{\delta}(d)$ in the case of the short regression from Section 4.1. That is, $\Delta_{\text{ind}}^{\text{s}}$ in (20) equals $\delta(0)$.

4.3 Long regression with interactions

A common variant of the long regression we just study is the regression that additionally includes the interactions between the K actions, A_1, \dots, A_K , and the treatment D ; i.e., the slope coefficient Δ_{inter} in (10). We call this the long regression with interactions. Our main result below shows that the slope coefficient Δ_{inter} in (10) admits a decomposition with the same shortcomings of the one we derived for Δ_{long} , including the possibility of Δ_{inter} being negative even in the absence of indirect effects. We formalize this below and provide a proof in Appendix A.

Theorem 4.3. Let Assumption 2.2 hold and assume that the covariance matrix of (D, A, AD) is positive definite. Then, the coefficient Δ_{inter} in (10) admits the decomposition

$$\Delta_{\text{inter}} = \Delta_{\text{dce}}^{\text{i}} + \Delta_{\text{ind}}^{\text{i}} , \quad (34)$$

where

$$\begin{aligned}\Delta_{\text{dce}}^{\text{i}} &\equiv \sum_{a \in \mathcal{A}} \omega_{\text{dce}}^{\text{i}}(a) (\mu(1, a) - \mu(0, a)) , \\ \Delta_{\text{ind}}^{\text{i}} &\equiv \sum_{a \in \mathcal{A}} \omega_{\text{ind}}^{\text{i}}(a) (\mu(0, a) - \mu(0, 0)) ,\end{aligned}$$

and $\{\omega_{\text{dce}}^{\text{i}}(a) : a \in \mathcal{A}\}$ and $\{\omega_{\text{ind}}^{\text{i}}(a) : a \in \mathcal{A}\}$ are as defined in Theorem A.2 and satisfy $\sum_{a \in \mathcal{A}} \omega_{\text{dce}}^{\text{i}}(a) = 1$ and $\sum_{a \in \mathcal{A}} \omega_{\text{ind}}^{\text{i}}(a) = 0$. Furthermore, the following statements are equivalent:

- (a) A are mutually exclusive binary variables, i.e., $\mathcal{A}_j = \{0, 1\}$ for $j = 1, \dots, K$ and $A_j A_l = 0$ for all $j, l = 1, \dots, K$ with $j \neq l$.
- (b) For any distribution of (A, D) , $\omega_{\text{dce}}^{\text{i}}(a) \geq 0$ for all $a \in \mathcal{A}$.
- (c) For any distribution of (A, D) , $\omega_{\text{ind}}^{\text{i}}(a) = 0$ for all $a \in \mathcal{A}$.

Theorem 4.3 is analogous to Theorem 4.2 and has very similar implications. Except in the special case where the actions in A are all mutually exclusive binary variables, which includes the case where A is a scalar binary variable as a special case, the term $\Delta_{\text{dce}}^{\text{i}}$ could be negative even if $\mu(1, a) - \mu(0, a) > 0$ for all $a \in \mathcal{A}$. This is because $\omega_{\text{dce}}^{\text{i}}(a)$ may be negative for some $a \in \mathcal{A}$. As a result, Δ_{inter} in general does not satisfy strong sign preservation for the same two reasons Δ_{long} did not satisfy it either. That is, (a) it is possible that $\Delta_{\text{ind}}^{\text{i}} < -\Delta_{\text{dce}}^{\text{i}}$, and (b) even if $\Delta_{\text{ind}}^{\text{i}} = 0$ the term $\Delta_{\text{dce}}^{\text{i}}$ could be negative by itself due to negative weights. Again, this second possibility separates Δ_{long} and Δ_{inter} from the short regression estimand, Δ_{short} .

Remark 4.8. Replacing Assumption 2.2 with Assumption 2.1 leads to a decomposition of Δ_{inter} that introduces three changes relative to the one in Theorem 4.3. First, the term $\Delta_{\text{dce}}^{\text{i}}$ becomes a linear combination of expectations that condition on $\{D = 1, A = a\}$. Second, the interpretation of $\Delta_{\text{ind}}^{\text{i}}$ becomes convoluted for the same reasons discussed for $\Delta_{\text{ind}}^{\text{s}}$. Finally, the decomposition additionally includes a selection term that is conceptually identical to $\Delta_{\text{ind}}^{\text{s}}$ in the short regression. The details of these expressions are presented in Theorem A.2 in Appendix A. ■

The possibility of $\omega_{\text{dce}}^{\text{i}}(a)$ being negative for some $a \in \mathcal{A}$ does not rely on pathological data generating processes and may arise in rather simple settings under reasonable distributions for (A, D) . We present simple examples in the proof of Theorem 4.3. Since these examples are analogous to those we described in Section 4.2, we do not describe them in detail here. We note, however, that in the example illustrated in Figure 2, the weights $\omega_{\text{dce}}^{\text{i}}(3) = -1.02$ for all values of p .

Remark 4.9. Similar to the long regression in (9) that we discussed in Remarks 4.1 and 4.6, the interaction regression is used extensively in the mediation literature. In the context of mediation analysis, this variant has been popularized and advocated by Judd and Kenny (1981); Kraemer et al. (2002, 2008). However, the main goal in that particular setting has been to test for the existence of mediation effects using the estimated coefficients in (10), see Kraemer et al. (2008) for details on the proposed test and Imai et al. (2010) for a result that shows that, under Assumption 4.1, such a test does not provide evidence in favor or against the parameter $\bar{\delta}(d)$ in (33) being zero. Our results, on the other hand, imply that even in settings where mediation effects are nuisance and the main goal is to interpret the coefficients directly related to the treatment D , the main conclusions depend on the distribution of (A, D) . ■

While Theorem 4.3 focuses on the properties of the estimand Δ_{inter} , in settings with interactions terms it is most often the case that the analyst would rather focus on the estimand $\Delta_{\text{inter}} + \sum_{j=1}^K \lambda_j E[A_j]$ (or, simply, $\Delta_{\text{inter}} + \sum_{j=1}^K \lambda_j a_j$ for given values a_j , $j = 1, \dots, K$). In Lemma B.4 in Appendix B we show that

$$\Delta_{\text{inter}} + \sum_{j=1}^K \lambda_j E[A_j] = \sum_{a \in \mathcal{A}} \omega_{\text{dce}}^{\text{i}\star}(a) (E[Y(1, a) - Y(0, a)]) + \omega_{\text{ind}}^{\text{i}\star}(a) (E[Y(0, a) - Y(0, 0)]) ,$$

where the “weights” $\{\omega_{\text{dce}}^{\text{i}\star}(a), \omega_{\text{ind}}^{\text{i}\star}(a) : a \in \mathcal{A}\}$ may be negative in general, and thus leading to an estimand with similar properties to those of Δ_{inter} . The one special case where the estimand $\Delta_{\text{inter}} + \sum_{j=1}^K \lambda_j E[A_j]$ works well is when $\mu(d, a)$ is assumed to take the functional form

$$\mu(d, a) = \kappa_0 + \kappa_1 d + \sum_{j=1}^K \kappa_{2,a} a_j + d \sum_{j=1}^K \kappa_{3,j} a_j , \quad (35)$$

which is equivalent to assume that the conditional mean of the observed outcome, Y , is correctly specified in the interaction regression in (10). Lemma B.5 in Appendix B shows that $\Delta_{\text{inter}} + \sum_{j=1}^K \lambda_j a_j = \mu(1, a) - \mu(0, a)$ in this case, delivering an average partial causal effect given $a = (a_1, \dots, a_K)$, as defined in Definition 3.1. It follows from these results that a clean interpretation of Δ_{inter} or $\Delta_{\text{inter}} + \sum_{j=1}^K \lambda_j a_j$ in terms of the definitions introduced in Section 3 essentially depends on a correctly specified linear model for potential outcomes and does not generally apply to non-parametric models, similarly to our results about the long regression in Section 4.2.

4.4 Strata fixed effects (SFE) regression

A lesson from Theorems 4.2 and 4.3 is that adding the other actions linearly in the regression is attractive only when the actions are mutually exclusive. This suggests

that if we could mechanically make the actions mutually exclusive, we could obtain an estimand that is free from indirect effects and that satisfies strong sign preservation. Perhaps not surprisingly, this is possible by considering the slope coefficient Δ_{sfe} in (11), where the regression controls for all possible values of A , i.e., $\{I\{A = a\} : a \in \mathcal{A}\}$. We refer to this regression as a strata-fixed effects regression, given its connection with the standard practice of including strata fixed effects in randomized controlled trials with covariate adaptive randomization; see Bugni et al. (2018, 2019). Our main result below shows that Δ_{sfe} in (11) admits a decomposition similar to the ones previously derived for Δ_{short} and Δ_{long} that is free from indirect effects. This is formalized below.

Theorem 4.4. Let Assumption 2.2 hold, $\pi_1(a)$ be as in (13), and assume that $P\{A = a\} > 0$ for all $a \in \mathcal{A}$ where

$$p_a(d) \equiv P\{D = d \mid A = a\} . \quad (36)$$

Then

$$\Delta_{\text{sfe}} = \sum_{a \in \mathcal{A}} \omega_{\text{sfe}}(a) (\mu(1, a) - \mu(0, a)) , \quad (37)$$

where the weights $\{\omega_{\text{sfe}}(a) : a \in \mathcal{A}\}$ are given by

$$\omega_{\text{sfe}}(a) \equiv \frac{\pi_1(a)\pi_0(a)}{\sum_{a' \in \mathcal{A}} \pi_1(a')\pi_0(a')} , \quad (38)$$

and satisfy $\sum_{a \in \mathcal{A}} \omega_{\text{sfe}}(a) = 1$ and $\omega_{\text{sfe}}(a) \geq 0$.

Theorem 4.4 shows that Δ_{sfe} identifies an average direct causal effect as in Definition 3.2. Importantly, it does not contain indirect effects and, as a result, Δ_{sfe} satisfies strong sign preservation as in Definition 3.3. The weights $\omega_{\text{sfe}}(a)$ admits a simple representation and depend only on the conditional probabilities that the action a happens for the treated and control group, $\pi_1(a)$ and $\pi_0(a)$. These weights are generally different than the weights associated with the direct effect in the short regression, $\Delta_{\text{dce}}^{\text{s}}$, which are simply $\pi_1(a)$, unless D and A are independent. We emphasize that the result in Theorem 4.4 does not require the other actions, A , to be singled-valued or mutually exclusive.

Remark 4.10. Replacing Assumption 2.2 with Assumption 2.1 in Theorem 4.4 leads to

$$\Delta_{\text{sfe}} = \Delta_{\text{dce}}^{\text{f}} + \Delta_{\text{sel}}^{\text{f}} , \quad (39)$$

where, for $\omega_{\text{sfe}}(a)$ as in (38),

$$\Delta_{\text{dce}}^{\text{f}} \equiv \sum_{a \in \mathcal{A}} \omega_{\text{sfe}}(a) E[Y(1, a) - Y(0, a) \mid D = 1, A = a] \quad (40)$$

$$\Delta_{\text{sel}}^{\text{f}} \equiv \sum_{a \in \mathcal{A}} \omega_{\text{sfe}}(a) \left(E[Y(0, a) \mid D = 1, A = a] - E[Y(0, a) \mid D = 0, A = a] \right) . \quad (41)$$

The two terms, $\Delta_{\text{dce}}^{\text{f}}$ and $\Delta_{\text{sel}}^{\text{f}}$, are directly comparable to the terms $\Delta_{\text{dce}}^{\text{s}}$ and $\Delta_{\text{sel}}^{\text{s}}$ in (14). In particular, $\Delta_{\text{sel}}^{\text{f}}$ is a selection term similar to $\Delta_{\text{sel}}^{\text{s}}$ in (17), with the only difference being the weights $\omega_{\text{sfe}}(a)$ replacing $\pi_1(a)$. ■

4.5 Saturated (SAT) Regression

We now turn our attention to the last set of estimands we study in this paper; the slope coefficient Δ_{sat} in (12). As we have stated in the introduction, under Assumption 2.2 it follows that $\mu(d, a)$ is immediately identified from $E[Y|D = d, A = a]$ for any $d \in \{0, 1\}$ and $a \in \mathcal{A}$ and so identification of any contrast of means $\mu(d, a)$ is straightforward. From this, it immediately follows that the same result could be achieved by running a saturated regression, as in (12), that we re-write here for readability,

$$Y = \sum_{a \in \mathcal{A}} \gamma(a) I\{A = a\} + \sum_{a \in \mathcal{A}} \Delta_{\text{sat}}(a) I\{A = a\} D + \epsilon .$$

Standard results on saturated regressions imply that $\Delta_{\text{sat}}(a) = \mu(1, a) - \mu(0, a)$ for all $a \in \mathcal{A}$, and so $\Delta_{\text{sat}}(a)$ captures an average causal partial effect of D on Y for each value of the other actions, $a \in \mathcal{A}$, aligned with Definition 3.1. For completeness, we state and prove this result formally in Theorem A.4 in the appendix. The same theorem also shows that replacing Assumption 2.2 with Assumption 2.1 leads to a decomposition of $\Delta_{\text{sat}}(a)$ that includes a selection term, as it was also the case for the other estimands we considered.

5 Concluding Remarks

In this paper we study settings where the analyst is interested in identifying and estimating an average causal effect of a binary treatment on an outcome of interest, in instances where the outcomes are “delayed” in the sense that they do not get immediately realized after treatment assignment. This delay in the realization of the outcomes creates a time window in between the treatment assignment and the realization of the outcome that, in turn, opens up the possibility for other observed endogenous actions to take place before the outcome is realized. In this context, we present formal results on how we can decompose popular estimands that arise from running regressions of the outcome on the treatment and different ways of “controlling” for the other actions and show that our decompositions have immediate implications on how these estimands can be interpreted in applications. All in all, our results provide a framework that allows analysts to understand under what type of conditions the practice of “controlling” for the presence of other actions leads to estimands that admit causal and ceteris paribus interpretations. Perhaps our most salient result is the one that shows that the most popular estimand

that linearly controls for the other actions in the regression, with or without interactions with the treatment, does not generally provide benefits relative to a simple regression of outcome on treatment. Under our assumptions, however, identification of partial causal effects immediately follows from saturated regressions.

A Proofs

Proof of Theorem 4.1. This proof follows from derivations in Section 4.1 and basic algebraic manipulations. ■

Theorem A.1. Consider the long regression in (8) and let Σ_{long} denote the variance-covariance matrix of (A, D) . Assume Σ_{long} is positive definite and let $M = \text{Cov}(D, A) \text{Var}(A)^{-1}$. Then,

$$\begin{aligned} \Delta_{\text{long}} &= \sum_{a \in \mathcal{A}} \omega_{\text{dce}}^1(a) E[Y(1, a) - Y(0, a) | D = 1, A = a] \\ &\quad + \sum_{a \in \mathcal{A}} \omega_{\text{ind}}^1(a) (E[Y(0, a) | D = 0, A = a] - E[Y(0, 0) | D = 0, A = 0]) \\ &\quad + \sum_{a \in \mathcal{A}} \omega_{\text{dce}}^1(a) (E[Y(0, a) | D = 1, A = a] - E[Y(0, a) | D = 0, A = a]) , \end{aligned} \quad (\text{A-1})$$

where

$$\begin{aligned} \omega_{\text{dce}}^1(a) &\equiv \frac{\pi_1(a) [\text{Var}(D) - P\{D = 1\} \sum_{j=1}^K M_j (a_j - E[A_j])]}{\text{Var}(D) - \text{Cov}(D, A) \text{Var}(A)^{-1} \text{Cov}(A, D)} \\ \omega_{\text{ind}}^1(a) &\equiv \frac{\text{Var}(D) [\pi_1(a) - \pi_0(a)] - P\{A = a\} \sum_{j=1}^K M_j (a_j - E[A_j])}{\text{Var}(D) - \text{Cov}(D, A) \text{Var}(A)^{-1} \text{Cov}(A, D)} , \end{aligned} \quad (\text{A-2})$$

and $\pi_d(a)$ is defined in (13). Furthermore, $\sum_{a \in \mathcal{A}} \omega_{\text{dce}}^1(a) = 1$, $\sum_{a \in \mathcal{A}} a \omega_{\text{ind}}^1(a) = \mathbf{0}$, and $\sum_{a \in \mathcal{A}} \omega_{\text{ind}}^1(a) = 0$.

Proof. Let $\theta \equiv (\theta_j : j = 1, \dots, K)$. By properties of projections,

$$E[(1, D, A')'(Y - (\Delta_{\text{long}} D + \theta_0 + \theta' A))] = \mathbf{0} . \quad (\text{A-3})$$

Profiling θ_0 leads to

$$\text{Cov}(D, Y) = \text{Var}(D) \Delta_{\text{long}} + \text{Cov}(A, D)' \theta \quad (\text{A-4})$$

$$\text{Cov}(A, Y) = \text{Cov}(A, D) \Delta_{\text{long}} + \text{Var}(A) \theta . \quad (\text{A-5})$$

Since Σ_{long} is positive definite, $\text{Var}(A)$ is positive definite. Then, (A-5) implies that $\theta = \text{Var}(A)^{-1} (\text{Cov}(A, Y) - \text{Cov}(A, D) \Delta_{\text{long}})$. If we plug this into (A-4), we get

$$(\text{Var}(D) - \text{Cov}(D, A) \text{Var}(A)^{-1} \text{Cov}(A, D)) \Delta_{\text{long}} = \text{Cov}(D, Y) - M \text{Cov}(A, Y) . \quad (\text{A-6})$$

Since Σ_{long} is positive definite, $\text{Cov}(D, A) \text{Var}(A)^{-1} \text{Cov}(A, D) > 0$, and so (A-6) implies that

$$\Delta_{\text{long}} = \frac{\text{Var}(D)\Delta_{\text{short}} - \sum_{j=1}^K M_j \text{Cov}(A_j, Y)}{\text{Var}(D) - \text{Cov}(D, A) \text{Var}(A)^{-1} \text{Cov}(A, D)}, \quad (\text{A-7})$$

where we used that $\text{Var}(D)\Delta_{\text{short}} = \text{Cov}(D, Y)$. For any $j = 1, \dots, K$, some algebra shows that

$$\begin{aligned} \text{Cov}(A_j, Y) &= \sum_{a \in \mathcal{A}} E[Y(1, a) - Y(0, a) | D = 1, A = a](a_j - E[A_j])\pi_1(a)E[D] \\ &+ \sum_{a \in \mathcal{A}} (E[Y(0, a) | D = 0, A = a] - E[Y(0, 0) | D = 0, A = 0])(a_j - E[A_j])P\{A = a\} \\ &+ \sum_{a \in \mathcal{A}} (E[Y(0, a) | D = 1, A = a] - E[Y(0, a) | D = 0, A = a])(a_j - E[A_j])\pi_1(a)E[D]. \end{aligned} \quad (\text{A-8})$$

Then, (A-1) follows from combining (14), (A-7), and (A-8).

To show $\sum_{a \in \mathcal{A}} \omega_{\text{dce}}^1(a) = 1$, consider the following derivation.

$$\begin{aligned} \sum_{a \in \mathcal{A}} \omega_{\text{dce}}^1(a) &\stackrel{(1)}{=} \frac{\text{Var}(D) - P\{D = 1\} \sum_{j=1}^K M_j (E[A_j | D = 1] - E[A_j])}{\text{Var}(D) - \text{Cov}(D, A) \text{Var}(A)^{-1} \text{Cov}(A, D)} \\ &\stackrel{(2)}{=} \frac{\text{Var}(D) - M \text{Cov}(A, D)}{\text{Var}(D) - \text{Cov}(D, A) \text{Var}(A)^{-1} \text{Cov}(A, D)} \stackrel{(3)}{=} 1, \end{aligned}$$

where (1) holds by $\sum_{a \in \mathcal{A}} \pi_1(a) = 1$ and $\sum_{a \in \mathcal{A}} \pi_1(a)a_j = E[A_j | D = 1]$, and (2) holds by $P\{D = 1\}(E[A_j | D = 1] - E[A_j]) = \text{Cov}(A_j, D)$, and (3) holds by definition of M .

We show $\sum_{a \in \mathcal{A}} \omega_{\text{ind}}^1(a) = 0$ by the following derivation applied to its numerator:

$$\text{Var}(D) \sum_{a \in \mathcal{A}} [\pi_1(a) - \pi_0(a)] - \sum_{a \in \mathcal{A}} P\{A = a\} \sum_{j=1}^K M_j (a_j - E[A_j]) = 0,$$

where the equality holds by $\sum_{a \in \mathcal{A}} \pi_1(a) = \sum_{a \in \mathcal{A}} \pi_0(a) = \sum_{a \in \mathcal{A}} P\{A = a\} = 1$ and $\sum_{a \in \mathcal{A}} P\{A = a\} a_j = E[A_j]$.

Finally, we show $\sum_{a \in \mathcal{A}} a_u \omega_{\text{ind}}^1(a) = 0$ for any $u = 1, \dots, K$. Once again, we focus on following derivation applied to its numerator:

$$\begin{aligned} &\sum_{a \in \mathcal{A}} a_u \text{Var}(D)[\pi_1(a) - \pi_0(a)] - \sum_{a \in \mathcal{A}} a_u P\{A = a\} \sum_{j=1}^K M_j (a_j - E[A_j]) \\ &\stackrel{(1)}{=} \text{Var}(D)[E[A_u | D = 1] - E[A_u | D = 0]] - M \text{Cov}(A_u, A) \\ &\stackrel{(2)}{=} \text{Cov}(D, A_u) - \text{Cov}(D, A) \text{Var}(A)^{-1} \text{Cov}(A, A_u) \stackrel{(3)}{=} 0, \end{aligned}$$

where (1) holds by $\sum_{a \in \mathcal{A}} \pi_d(a)a_j = E[A_j | D = d]$ for $d = 0, 1$, and $\sum_{a \in \mathcal{A}} a_u P\{A = a\} (a_j - E[A_j]) = \text{Cov}(A_u, A_j)$, (2) holds by $\text{Var}(D)[E[A_u | D = 1] - E[A_u | D = 0]] = \text{Cov}(D, A_u)$ and the definition of M , and (3) holds by the fact that $\text{Var}(A)^{-1} \text{Cov}(A, A_u)$ equals a column vector with zeros except for a one in the u th position. ■

Proof of Theorem 4.2. The first part follows from Theorem A.1, which also yields $\sum_{a \in \mathcal{A}} \omega_{\text{dce}}^1(a) = 1$ and $\sum_{a \in \mathcal{A}} \omega_{\text{ind}}^1(a) = 0$. To complete the proof, we now show the equivalence between (a), (b),

and (c).

First, we show that (a) implies (b) and (c). To this end, assume (a) holds. Then, the long regression in (9) is equivalent to an SFE regression in (11). To see why, note that (a) implies that A is a K dimensional vector that is either equal to 0 or equal to a canonical vector (i.e., a vector with a 1 in only one of its coordinates and zeroes otherwise). If we then let $\theta(a) = \theta_0$ for $a = 0$ and $\theta(a) = \theta_j$ for a being the canonical vector with j th coordinate equal to one, we get

$$\theta_0 + \theta' A = \sum_{a \in \mathcal{A}} \theta(a) I\{A = a\} .$$

Therefore, $\Delta_{\text{long}} = \Delta_{\text{sfe}}$ and Theorem A.3 imply (b) (with $\omega_{\text{dce}}^{\text{l}}(a) = \omega_{\text{dce}}^{\text{f}}(a)$) and (c).

Second, we show that (b) or (c) implies (a) or, equivalently, the negation of (a) implies the negation of (b) and the negation of (c).

Start by considering the case when $K = 1$. Then, (a) fails when $\mathcal{A} \neq \{0, 1\}$. For example, consider the case where $\mathcal{A} = \{0, 1, 2\}$ with $\{A|D = 1\} \sim \text{Bi}(2, 0.3)$, $\{A|D = 0\} \sim \text{Bi}(2, 0.9)$, and $P\{D = 1\} = 0.5$, where $\text{Bi}(n, p)$ denotes a Binomial distribution with n trials and probability p . With this distribution of (A, D) , the weights in (A-2) become $\omega_{\text{dce}}^{\text{l}} \approx [-0.1, 0.76, 0.34]$ and $\omega_{\text{ind}}^{\text{l}} \approx [-0.14, 0.28, -0.14]$, and so (b) and (c) fail.

Next, consider the case when $K = 2$. In this case (a) can fail when (i) $\mathcal{A}_j \neq \{0, 1\}$ for some $j = 1, 2$ or (ii) $\mathcal{A}_j = \{0, 1\}$ for all $j = 1, 2$ but $A_1 A_2 \neq 0$. For (i), consider $\{A_1|D = 1\} \sim \text{Bi}(2, 0.3)$, $\{A_1|D = 0\} \sim \text{Bi}(2, 0.9)$, $P\{D = 1\} = 0.5$, $A_2 \perp \{D, A_1\}$, and $\text{Var}(A_2) > 0$. The fact that $A_2 \perp \{D, A_1\}$ and $\text{Var}(A_2) > 0$ implies that A_2 drops out of the expressions in (A-2), and the example becomes identical to the one considered when $K = 1$, where (b) and (c) fail. For (ii), let $\text{Ber}(p)$ denote a Bernoulli distribution with parameter p and consider $\{A_j|D = 0\} \sim \text{Ber}(0.1)$ and $\{A_j|D = 1\} \sim \text{Ber}(0.7)$ for $j = 1, 2$, with $P\{D = 1\} = 0.5$, so that $\mathcal{A}_j = \{0, 1\}$ for $j = 1, 2$ and $P\{A_1 A_2 = 0\} \approx 0.45$. With this distribution of (A, D) the weights in (A-2) become $\omega_{\text{dce}}^{\text{l}} \approx [0.34, 0.38, 0.48, -0.10]$ and $\omega_{\text{ind}}^{\text{l}} \approx [-0.14, 0.14, 0.14, -0.14]$, and so (b) and (c) fail.

Finally, consider the case $K > 2$. Then, (a) can fail when (i) $\mathcal{A}_j \neq \{0, 1\}$ for some $j = 1, \dots, K$ or (ii) $\mathcal{A}_j = \{0, 1\}$ for all $j = 1, \dots, K$ but $A_j A_l \neq 0$ for some $j, l = 1, \dots, K$ with $j \neq l$. In either case, we can repeat the examples used for $K = 2$ by adding coordinates $j = 3, \dots, K$ with $\{A_j : j > 2\} \perp \{D, \{A_j : j \leq 2\}\}$, and $\text{Var}(A_j) > 0$ for $j > 2$. By construction, $\{A_j : j > 2\}$ drops out of the expressions in (A-2), and the examples considered with $K = 2$ imply the failure of (b) and (c). ■

Theorem A.2. Consider the interaction regression in (10) and let Σ_{inter} denote the variance-covariance matrix of (A, D, AD) . Assume Σ_{inter} is positive definite and let $M = \text{Cov}(D, W) \text{Var}(W)^{-1}$ with $W \equiv (A', A'D)'$. Then,

$$\begin{aligned} \Delta_{\text{inter}} &= \sum_{a \in \mathcal{A}} \omega_{\text{dce}}^{\text{i}}(a) E[Y(1, a) - Y(0, a)|D = 1, A = a] \\ &\quad + \sum_{a \in \mathcal{A}} \omega_{\text{ind}}^{\text{i}}(a) (E[Y(0, a)|D = 0, A = a] - E[Y(0, 0)|D = 0, A = 0]) \\ &\quad + \sum_{a \in \mathcal{A}} \omega_{\text{dce}}^{\text{i}}(a) (E[Y(0, a)|D = 1, A = a] - E[Y(0, a)|D = 0, A = a]) , \end{aligned} \quad (\text{A-9})$$

where

$$\omega_{\text{dce}}^i(a) \equiv \frac{\pi_1(a) \left[\sigma_D^2 - p \sum_{j=1}^K M_j (a_j - E[A_j]) - p \sum_{j=1}^K M_{j+K} (a_j - pE[A_j|D=1]) \right]}{\sigma_D^2 - M \text{Cov}(W, D)} \quad (\text{A-10})$$

$$\omega_{\text{ind}}^i(a) \equiv \frac{\sigma_D^2 (\pi_1(a) - \pi_0(a)) - \sum_{j=1}^K M_j p_a (a_j - E[A_j]) - p \sum_{j=1}^K M_{j+K} (\pi_1(a) a_j - p_a E[A_j|D=1])}{\sigma_D^2 - M \text{Cov}(W, D)},$$

$p = P\{D = 1\}$, $p_a = P\{A = a\}$, $\sigma_D^2 = \text{Var}(D)$, and $\pi_d(a)$ is defined in (13). Furthermore, $\sum_{a \in \mathcal{A}} \omega_{\text{dce}}^i(a) = 1$, $\sum_{a \in \mathcal{A}} a \omega_{\text{ind}}^i(a) = \mathbf{0}$, and $\sum_{a \in \mathcal{A}} \omega_{\text{ind}}^i(a) = 0$.

Proof. Let $\theta = (\theta_j : j = 1, \dots, K)$, $\lambda = (\lambda_j : j = 1, \dots, K)$, and $\alpha = (\theta', \lambda')'$. By properties of projections,

$$E[(1, D, A', DA')'(Y - (\Delta_{\text{inter}} D + \theta_0 + \alpha' W))] = \mathbf{0}. \quad (\text{A-11})$$

Profiling θ_0 leads to,

$$\text{Cov}(D, Y) = \text{Var}(D) \Delta_{\text{inter}} + \text{Cov}(W, D)' \alpha \quad (\text{A-12})$$

$$\text{Cov}(W, Y) = \text{Cov}(W, D) \Delta_{\text{inter}} + \text{Var}(W) \alpha. \quad (\text{A-13})$$

Since Σ_{inter} is positive definite, $\text{Var}(W)$ is positive definite. Then, (A-13) implies that $\alpha = \text{Var}(W)^{-1} (\text{Cov}(W, Y) - \text{Cov}(W, D) \Delta_{\text{inter}})$. If we plug this into (A-12), we get

$$(\text{Var}(D) - M \text{Cov}(W, D)) \Delta_{\text{inter}} = \text{Cov}(D, Y) - M \text{Cov}(W, Y). \quad (\text{A-14})$$

Since Σ_{inter} is positive definite, $\text{Var}(D) - \text{Cov}(W, D)' \text{Var}(W)^{-1} \text{Cov}(W, D) > 0$ and so (A-14) implies that

$$\Delta_{\text{inter}} = \frac{\text{Cov}(D, Y) - \sum_{j=1}^K M_j \text{Cov}(A_j, Y) - \sum_{j=1}^K M_{j+K} \text{Cov}(DA_j, Y)}{\text{Var}(D) - M \text{Cov}(W, D)}, \quad (\text{A-15})$$

where we used that $\text{Var}(D) \Delta_{\text{short}} = \text{Cov}(D, Y)$. For any $j = 1, \dots, K$, some algebra shows that

$$\begin{aligned} \text{Cov}(A_j, Y) &= \sum_{a \in \mathcal{A}} E[Y(1, a) - Y(0, a) | D = 1, A = a] (a_j - E[A_j]) \pi_1(a) p \\ &+ \sum_{a \in \mathcal{A}} (E[Y(0, a) | D = 0, A = a] - E[Y(0, 0) | D = 0, A = 0]) (a_j - E[A_j]) p_a \\ &+ \sum_{a \in \mathcal{A}} (E[Y(0, a) | D = 1, A = a] - E[Y(0, a) | D = 0, A = a]) (a_j - E[A_j]) \pi_1(a) p, \end{aligned} \quad (\text{A-16})$$

and

$$\begin{aligned} \text{Cov}(DA_j, Y) &= p \sum_{a \in \mathcal{A}} E[Y(1, a) - Y(0, a) | D = 1, A = a] \pi_1(a) (a_j - pE[A_j|D=1]) \\ &+ p \sum_{a \in \mathcal{A}} (E[Y(0, a) | D = 1, A = a] - E[Y(0, a) | D = 0, A = a]) \pi_1(a) (a_j - pE[A_j|D=1]) \\ &+ p \sum_{a \in \mathcal{A}} (E[Y(0, a) | D = 0, A = a] - E[Y(0, 0) | D = 0, A = 0]) (\pi_1(a) a_j - p_a E[A_j|D=1]). \end{aligned} \quad (\text{A-17})$$

By plugging in (14), (A-16), and (A-17) into (A-15), (A-9) follows.

To show $\sum_{a \in \mathcal{A}} \omega_{\text{dce}}^i(a) = 1$, consider the following derivation.

$$\begin{aligned} \sum_{a \in \mathcal{A}} \omega_{\text{dce}}^i(a) &\stackrel{(1)}{=} \frac{\sigma_D^2 - p \sum_{j=1}^K M_j (E[A_j|D=1] - E[A_j]) - \sum_{j=1}^K M_{j+K} \sigma_D^2 E[A_j|D=1]}{\sigma_D^2 - M \text{Cov}(W, D)} \\ &\stackrel{(2)}{=} \frac{\sigma_D^2 - \sum_{j=1}^K M_j \text{Cov}(D, A_j) - \sum_{j=1}^K M_{j+K} \text{Cov}(D, DA_j)}{\sigma_D^2 - M \text{Cov}(W, D)} \stackrel{(3)}{=} 1, \end{aligned}$$

where (1) holds by $\sum_{a \in \mathcal{A}} \pi_1(a) = 1$, $\sum_{a \in \mathcal{A}} \pi_1(a) a_j = E[A_j|D=1]$, and $\sigma_D^2 = p(1-p)$, (2) holds by $p(E[A_j|D=1] - E[A_j]) = \text{Cov}(D, A_j)$ and $\sigma_D^2 E[A_j|D=1] = \text{Cov}(DA_j, D)$, and (3) holds by definition of M .

We show $\sum_{a \in \mathcal{A}} \omega_{\text{ind}}^i(a) = 0$ by the following derivation applied to its numerator:

$$\sigma_D^2 \sum_{a \in \mathcal{A}} (\pi_1(a) - \pi_0(a)) - \sum_{j=1}^K M_j \sum_{a \in \mathcal{A}} p_a (a_j - E[A_j]) - p \sum_{j=1}^K M_{j+K} \sum_{a \in \mathcal{A}} (\pi_1(a) a_j - p_a E[A_j|D=1]) = 0,$$

where the equality holds by $\sum_{a \in \mathcal{A}} \pi_1(a) = \sum_{a \in \mathcal{A}} \pi_0(a) = \sum_{a \in \mathcal{A}} p_a = 1$, $\sum_{a \in \mathcal{A}} p_a a_j = E[A_j]$, and $\sum_{a \in \mathcal{A}} \pi_1(a) a_j = E[A_j|D=1]$.

Finally, we show $\sum_{a \in \mathcal{A}} a_u \omega_{\text{ind}}^i(a) = 0$ for any $u = 1, \dots, K$. Once again, we focus on following derivation applied to its numerator:

$$\begin{aligned} &\sum_{a \in \mathcal{A}} a_u \sigma_D^2 (\pi_1(a) - \pi_0(a)) - \sum_{j=1}^K M_j \sum_{a \in \mathcal{A}} a_u p_a (a_j - E[A_j]) - p \sum_{j=1}^K M_{j+K} \sum_{a \in \mathcal{A}} a_u (\pi_1(a) a_j - p_a E[A_j|D=1]) \\ &\stackrel{(1)}{=} \sigma_D^2 [E[A_u|D=1] - E[A_u|D=0]] - \sum_{j=1}^K M_j \text{Cov}(A_u, A_j) - \sum_{j=1}^K M_{j+K} p (E[A_u A_j|D=1] - E[A_u] E[A_j|D=1]) \\ &\stackrel{(2)}{=} \text{Cov}(D, A_u) - \sum_{j=1}^K M_j \text{Cov}(A_j, A_u) - \sum_{j=1}^K M_{j+K} \text{Cov}(DA_j, A_u) \\ &\stackrel{(3)}{=} \text{Cov}(D, A_u) - \text{Cov}(D, W) \text{Var}(W)^{-1} \text{Cov}(W, A_u) \stackrel{(4)}{=} 0, \end{aligned}$$

where (1) holds by $\sum_{a \in \mathcal{A}} \pi_d(a) a_j = E[A_j|D=d]$ for $d = 0, 1$, $\sum_{a \in \mathcal{A}} a_u p \{A = a\} (a_j - E[A_j]) = \text{Cov}(A_u, A_j)$, and $\sum_{a \in \mathcal{A}} a_u p_a = E[A_u]$, (2) holds by $\text{Var}(D)[E[A_u|D=1] - E[A_u|D=0]] = \text{Cov}(D, A_u)$ and $p(E[A_u A_j|D=1] - E[A_u] E[A_j|D=1]) = \text{Cov}(DA_j, A_u)$, (3) holds by the definition of M , and (4) holds by the fact that $\text{Var}(W)^{-1} \text{Cov}(W, A_u)$ equals a column vector with zeros except for a one in the u th position. ■

Proof of Theorem 4.3. The first part follows from Theorem A.2, which also yields that $\sum_{a \in \mathcal{A}} \omega_{\text{dce}}^i(a) = 1$ and $\sum_{a \in \mathcal{A}} \omega_{\text{ind}}^i(a) = 0$. To complete the proof, we now show the equivalence between (a), (b), and (c).

First, we show that (a) implies (b) and (c). To this end, assume (a) holds. Then, the long with interactions regression in (10) is equivalent to a SAT regression in (11). To see why, note that (a) implies that $\mathcal{A} = \{\mathbf{0}_{K \times 1}, \{e_j : j = 1, \dots, K\}\}$, where $e_j \in \mathbb{R}^{K \times 1}$ has a one in the j 'th coordinate and zero otherwise. By defining $A_0 = 1 - \sum_{j=1}^K A_j$, $\gamma(a) = \theta_0$ and $\Delta_{\text{sat}}(a) = \Delta_{\text{inter}}$

for $a = 0$, and $\gamma(a) = \theta_0 + \theta_j$ and $\Delta_{\text{sat}}(a) = \Delta_{\text{inter}} + \lambda_j$ for $a = e_j$ with $j = 1, \dots, K$, we get

$$\Delta_{\text{inter}}D + \theta_0 + \theta' A + \lambda' AD = \sum_{a \in \mathcal{A}} \gamma(a) I\{A = a\} + \sum_{a \in \mathcal{A}} \Delta_{\text{sat}}(a) I\{A = a\} D .$$

Therefore, $\Delta_{\text{inter}} = \Delta_{\text{sat}}(0)$ and Theorem A.4 imply (b) (with $\omega_{\text{dce}}(0) = 1$ and $\omega_{\text{dce}}(e_j) = 0$ for $j = 1, \dots, K$) and (c).

To conclude, we now show that (b) or (c) implies (a) or, equivalently, the negation of (a) implies the negation of (b) and the negation of (c).

First, consider the case when $K = 1$. Then, (a) fails when $\mathcal{A} \neq \{0, 1\}$. For example, if $\{A|D = 0\} \sim \text{Bi}(2, 0.3)$, $\{A|D = 1\} \sim \text{Bi}(2, 0.9)$, and $P\{D = 1\} = 0.5$, and so $\mathcal{A} = \{0, 1, 2\}$. By evaluating this information on (A-10), we get $\omega_{\text{dce}} \approx [0.19, 1.62, -0.81]$ and $\omega_{\text{ind}} \approx [-0.72, 1.44, -0.72]$, i.e., (b) and (c) fail.

Second, consider the case when $K = 2$. Then, (a) can fail when (i) $\mathcal{A}_j \neq \{0, 1\}$ for some $j = 1, 2$ or (ii) $\mathcal{A}_j = \{0, 1\}$ for all $j = 1, 2$ but $A_1 A_2 \neq 0$. For (i), consider $\{A_1|D = 0\} \sim \text{Bi}(2, 0.3)$, $\{A_1|D = 1\} \sim \text{Bi}(2, 0.9)$, $P\{D = 1\} = 0.5$, $A_2 \perp \{D, A_1\}$, and $\text{Var}(A_2) > 0$. The fact that $A_2 \perp \{D, A_1\}$ and $\text{Var}(A_2) > 0$ implies that A_2 drops out of the expressions in (A-10), and the example becomes identical to the one considered when $K = 1$ and, thus, (b) and (c) fail. For (ii), consider $\{A_j|D = 0\} \sim \text{Be}(0.3)$ and $\{A_j|D = 1\} \sim \text{Be}(0.9)$ for $j = 1, 2$, and $P(D = 1) = 0.5$, and so $\mathcal{A}_j = \{0, 1\}$ for $j = 1, 2$ and $P(A_1 A_2 = 0) \approx 0.25$. By evaluating this information on (A-10), we get $\omega_{\text{dce}} \approx [0.19, 0.81, 0.81, -0.81]$ and $\omega_{\text{ind}} \approx [-0.72, 0.72, 0.72, -0.72]$, i.e., (b) and (c) fail.

Finally, consider $K > 2$. Then, (a) can fail when (i) $\mathcal{A}_j \neq \{0, 1\}$ for some $j = 1, \dots, K$ or (ii) $\mathcal{A}_j = \{0, 1\}$ for all $j = 1, \dots, K$ but $A_j A_l \neq 0$ for some $j, l = 1, \dots, K$ with $j \neq l$. In either case, we can repeat the examples used for $K = 2$ by adding coordinates $j = 3, \dots, K$ with $\{A_j : j > 2\} \perp \{D, \{A_j : j \leq 2\}\}$, and $\text{Var}(A_j) > 0$ for $j > 2$. By construction, $\{A_j : j > 2\}$ drops out of the expressions in (A-10), and the examples considered with $K = 2$ imply the failure of (b) and (c). ■

Theorem A.3. Consider the SFE regression in (11), and assume that $P\{A = a\} > 0$ and $P\{D = 1|A = a\} \in (0, 1)$ for all $a \in \mathcal{A}$. Then,

$$\Delta_{\text{sfe}} = \Delta_{\text{dce}}^{\text{f}} + \Delta_{\text{sel}}^{\text{f}} , \quad (\text{A-18})$$

where

$$\omega_{\text{sfe}}(a) \equiv \frac{P\{D = 0|A = a\}P\{D = 1|A = a\}P\{A = a\}}{\sum_{\tilde{a} \in \mathcal{A}} P\{D = 1|A = \tilde{a}\}P\{D = 0|A = \tilde{a}\}P\{A = \tilde{a}\}} \quad \text{for all } a \in \mathcal{A} \quad (\text{A-19})$$

$$\Delta_{\text{dce}}^{\text{f}} \equiv \sum_{a \in \mathcal{A}} \omega_{\text{sfe}}(a) E[Y(1, a) - Y(0, a)|D = 1, A = a] \quad (\text{A-20})$$

$$\Delta_{\text{sel}}^{\text{f}} \equiv \sum_{a \in \mathcal{A}} \omega_{\text{sfe}}(a) (E[Y(0, a)|D = 1, A = a] - E[Y(0, a)|D = 0, A = a]) . \quad (\text{A-21})$$

Furthermore, note that $\sum_{a \in \mathcal{A}} \omega_{\text{sfe}}(a) = 1$ and $\omega_{\text{sfe}}(a) \geq 0$.

Proof. By properties of projections,

$$E[YD] = \Delta_{\text{sfe}}E[D] + \sum_{a \in \mathcal{A}} \theta(a)E[I\{A = a\}D] \quad (\text{A-22})$$

$$E[YI\{A = a\}] = \Delta_{\text{sfe}}E[DI\{A = a\}] + \theta(a)P\{A = a\} \text{ for all } a \in \mathcal{A} . \quad (\text{A-23})$$

By $P\{A = a\} > 0$ for all $a \in \mathcal{A}$, (A-23) implies that

$$\theta(a) = E[Y|A = a] - \Delta_{\text{sfe}}E[D|A = a] \quad \text{for all } a \in \mathcal{A} . \quad (\text{A-24})$$

Then, (A-22), (A-24), and some algebra imply that

$$\begin{aligned} & E[Y|D = 1] - \sum_{a \in \mathcal{A}} E[Y|A = a]P\{A = a|D = 1\} \\ &= \Delta_{\text{sfe}} \sum_{a \in \mathcal{A}} \frac{P\{D = 1|A = a\}P\{D = 0|A = a\}P\{A = a\}}{P\{D = 1\}} , \end{aligned} \quad (\text{A-25})$$

Under $P\{A = a\} > 0$ and $P\{D = 1|A = a\} \in (0, 1)$ for all $a \in \mathcal{A}$, (A-25) implies that

$$\begin{aligned} \Delta_{\text{sfe}} &= \frac{P\{D = 1\}E[Y|D = 1] - \sum_{a \in \mathcal{A}} E[Y|A = a]P\{A = a, D = 1\}}{\sum_{a \in \mathcal{A}} P\{D = 1|A = a\}P\{D = 0|A = a\}P\{A = a\}} \\ &= \sum_{a \in \mathcal{A}} \omega_{\text{sfe}}(a)(E[Y|A = a, D = 1] - E[Y|A = a, D = 0]) . \end{aligned} \quad (\text{A-26})$$

By doing algebra on (A-26), (A-18) follows. Finally, verifying $\sum_{a \in \mathcal{A}} \omega_{\text{sfe}}(a) = 1$ and $\omega_{\text{sfe}}(a) \geq 0$ is straightforward given the definition in (A-19). ■

Proof of Theorem 4.4. This result follows immediately from Theorem A.3. ■

Theorem A.4. Consider the SAT regression in (12), and assume that $P\{A = a\} > 0$ and $P\{D = 1|A = a\} \in (0, 1)$ for all $a \in \mathcal{A}$. Then, for all $a \in \mathcal{A}$,

$$\Delta_{\text{sat}}(a) = \Delta_{\text{dce}}^{\text{t}}(a) + \Delta_{\text{sel}}^{\text{t}}(a) , \quad (\text{A-27})$$

where

$$\Delta_{\text{dce}}^{\text{t}}(a) \equiv E[Y(1, a) - Y(0, a)|D = 1, A = a] \quad (\text{A-28})$$

$$\Delta_{\text{sel}}^{\text{t}}(a) \equiv E[Y(0, a)|D = 1, A = a] - E[Y(0, a)|D = 0, A = a] . \quad (\text{A-29})$$

Furthermore, under Assumption 2.2, $\Delta_{\text{sel}}^{\text{t}}(a) = 0$ and

$$\Delta_{\text{sat}}(a) = \Delta_{\text{dce}}^{\text{t}}(a) = \mu(1, a) - \mu(0, a) . \quad (\text{A-30})$$

Proof. Fix $a \in \mathcal{A}$ arbitrarily throughout this proof. By projection,

$$\begin{aligned} E[YI\{A = a\}] &= \gamma(a)P\{A = a\} + \Delta_{\text{sat}}(a)E[DI\{A = a\}] \\ E[YDI\{A = a\}] &= (\gamma(a) + \Delta_{\text{sat}}(a))E[DI\{A = a\}] . \end{aligned} \quad (\text{A-31})$$

By $P\{A = a\} > 0$, (A-31) implies that

$$\gamma(a) = E[Y|A = a] - \Delta_{\text{sat}}(a)P\{D = 1|A = a\} \quad (\text{A-32})$$

$$E[YD|A = a] = (\gamma(a) + \Delta_{\text{sat}}(a))P\{D = 1|A = a\}. \quad (\text{A-33})$$

By plugging in (A-32) on (A-33), we get

$$E[YD|A = a] - E[Y|A = a]P\{D = 1|A = a\} = \Delta_{\text{sat}}(a)P\{D = 1|A = a\}P\{D = 0|A = a\}. \quad (\text{A-34})$$

By (A-34) and $P\{D = 1|A = a\} \in (0, 1)$, we get that

$$\Delta_{\text{sat}}(a) = E[Y(1, a)|D = 1, A = a] - E[Y(0, a)|D = 0, A = a]. \quad (\text{A-35})$$

The desired result follows from adding and subtracting $E[Y(0, a)|D = 1, A = 1]$ to (A-35). ■

B Auxiliary Lemmas

Lemma B.1. The following statements are true.

- (a) Assumption 4.1 implies Assumption 2.2.
- (b) Assumption 2.2 does not imply Assumption 4.1.
- (c) Assumption 4.1 implies that $Y(\tilde{d}, a) \perp A(d) | X$ for $(\tilde{d}, d, a) \in \mathcal{D} \times \mathcal{D} \times \mathcal{A}$.

Proof. Part (a). For any $(\tilde{d}, \tilde{a}, y, a, d)$, we have

$$\begin{aligned} P\{Y(\tilde{d}, \tilde{a}) \leq y, A(d) = a, D = d|X\} &\stackrel{(1)}{=} P\{Y(\tilde{d}, \tilde{a}) \leq y|X\} P\{A(d) = a|X\} P\{D = d|X\} \\ &\stackrel{(2)}{=} P\{Y(\tilde{d}, \tilde{a}) \leq y|X\} P\{A(d) = a, D = d|X\} \\ &\stackrel{(3)}{=} P\{Y(\tilde{d}, \tilde{a}) \leq y|X\} P\{A = a, D = d|X\}, \end{aligned} \quad (\text{B-36})$$

where (1) holds by (30) and (31), (2) holds by (30), and (3) holds by $A(D) = A$. Since $(\tilde{d}, \tilde{a}, y, a, d)$ is arbitrary, (B-36) implies Assumption 2.2.

Part (b). Consider the following example. Assume $X \perp (D, (A(d) : d \in \mathcal{D})', (Y(\tilde{d}, a) : (\tilde{d}, a) \in \mathcal{D} \times \mathcal{A})')'$, $Y(d, a) = 0$ for all (d, a) , $(A(1), A(0)) = (D, D)$, and $D \sim \text{Be}(0.5)$. Since $Y(d, a) = 0$, it is independent of $(D, A(D)) = (D, D)$. Thus, Assumption 2.2 holds. By $Y(d, a) = 0$ for all (d, a) and also $A(d) = D$, we have $Y(\tilde{d}, a) \perp A(d)|D$, so (31) holds. However, $(Y(\tilde{d}, a), A(d)) = (0, D) \not\perp D$, and so (30) and Assumption 4.1 fail.

Part (c). For any $(\tilde{d}, \tilde{a}, y, a, d)$, we have

$$\begin{aligned} P\{Y(\tilde{d}, \tilde{a}) \leq y, A(d) = a|X\} &\stackrel{(1)}{=} P\{Y(\tilde{d}, \tilde{a}) \leq y, A(d) = a|X, D\} \\ &\stackrel{(2)}{=} P\{Y(\tilde{d}, \tilde{a}) \leq y|X, D\} P\{A(d) = a|X, D\} \\ &\stackrel{(3)}{=} P\{Y(\tilde{d}, \tilde{a}) \leq y|X\} P\{A(d) = a|X\}, \end{aligned} \quad (\text{B-37})$$

where (1) and (3) hold by (30), and (2) holds by (31). Since $(\tilde{d}, \tilde{a}, y, a, d)$ is arbitrary, (B-37) implies Assumption 2.2. ■

Lemma B.2. Assume the conditions in Theorem 4.2, and that

$$\mu(d, a) = \kappa_0 + \kappa_1 d + \kappa_2' a \quad \text{for all } (d, a) \in \{0, 1\} \times \mathcal{A} \quad (\text{B-38})$$

for some constants $\kappa_0, \kappa_1, \kappa_2$. First, the coefficients in (9) satisfy $\Delta_{\text{long}} = \kappa_1$, $\theta_0 = \kappa_0$, and $\theta_1 = \kappa_2$. Second, the terms in the decomposition in (22) are $\Delta_{\text{dce}}^1 = \kappa_1$ and $\Delta_{\text{ind}}^1 = 0$.

Proof. Assumption 2.2 implies that $E(Y|D = d, A = a) = \mu(d, a)$ which, combined with (B-38), implies that the conditional expectation of Y is linear in $(1, a, d)$. From here, the first result follows from the fact that the linear regression consistently estimates the parameters of a linear conditional expectation. The second part follows immediately from combining (B-38) with $\sum_{a \in \mathcal{A}} a \omega_{\text{ind}}^1(a) = \mathbf{0}$ and $\sum_{a \in \mathcal{A}} \omega_{\text{dce}}^1(a) = 1$ (both shown in Theorem A.1). ■

Lemma B.3. The examples used in the proofs of Theorem 4.2 and 4.3 can be completed to satisfy Assumption 4.1.

Proof. For brevity, we focus on the example in the proof of Theorem 4.2 when $K = 1$. A similar argument can be made for all other examples.

Recall that the example in the proof of Theorem 4.2 when $K = 1$ is as follows: $\{A|D = 0\} \sim \text{Bi}(2, 0.3)$, $\{A|D = 1\} \sim \text{Bi}(2, 0.9)$, and $P\{D = 1\} = 0.5$, and so $\mathcal{A} = \{0, 1, 2\}$. The example is silent about X or $\{Y(d, a) : (d, a) \in \mathcal{D} \times \mathcal{A}\}$, and so it is unclear whether Assumption 4.1 holds or not. We now provide one way to complete the specification of the example in a manner compatible with Assumption 4.1.

Assume that $X \perp (\{Y(d, a) : (d, a) \in \mathcal{D} \times \mathcal{A}\}, D, \{A(\tilde{d}) : \tilde{d} \in \mathcal{D}\}, \{Y(d, a) : (d, a) \in \mathcal{D} \times \mathcal{A}\})$ non-stochastic and equal to $\{\mu(d, a) : (d, a) \in \mathcal{D} \times \mathcal{A}\}$, $A(0) \sim \text{Bi}(2, 0.3)$, $A(1) \sim \text{Bi}(2, 0.9)$, $D \sim \text{Be}(0.5)$, and $\{A(1), A(0), D\}$ are independent random variables. These conditions imply that $A(0) \stackrel{d}{=} \{A(0)|D = 0\} = \{A|D = 0\} \sim \text{Bi}(2, 0.3)$, $A(1) \stackrel{d}{=} \{A(1)|D = 1\} = \{A|D = 1\} \sim \text{Bi}(2, 0.9)$, and $P\{D = 1\} = 0.5$, as required by the example. Next, we show that the completed example satisfies Assumption 4.1. First, we have that (30) holds from the fact that X is independent of the rest of the problem, $\{Y(d, a) : (d, a) \in \mathcal{D} \times \mathcal{A}\}$ is non-stochastic, and $A(d) \perp D$. Second, we have that (31) holds from the fact that X is independent of the rest of the problem and $\{Y(d, a) : (d, a) \in \mathcal{D} \times \mathcal{A}\}$ is non-stochastic. ■

Lemma B.4. Consider the setup in Theorem 4.3 and that A is scalar. Then, the coefficients in (10) satisfy the following decomposition:

$$\Delta_{\text{inter}} + E[A]\lambda = \sum_{a \in \mathcal{A}} \omega_{\text{dce}}^{i*}(a)(E[Y(1, a) - Y(0, a)]) + \omega_{\text{ind}}^{i*}(a)(E[Y(0, a) - Y(0, 0)]) , \quad (\text{B-39})$$

where

$$\begin{aligned}
\Delta &= \text{Var}(AD) \text{Var}(A) - (\text{Cov}(DA, A))^2 \\
\Psi &= 1 + \frac{E[A]}{\Delta} (\text{Cov}(A, DA) \text{Cov}(A, D) - \text{Var}(A) \text{Cov}(DA, D)) \\
\omega_{\text{dce}}^{\text{i}*}(a) &= \Psi \omega_{\text{dce}}^{\text{i}}(a) + \frac{E[A]}{\Delta} (\text{Var}(A) p \pi_1(a) (a - p E[A|D=1]) - \text{Cov}(A, DA) (a - E[A]) \pi_1(a) p) \\
\omega_{\text{ind}}^{\text{i}*}(a) &= \Psi \omega_{\text{ind}}^{\text{i}}(a) + \frac{E[A]}{\Delta} (\text{Var}(A) p (\pi_1(a) a - p_a E[A|D=1]) - \text{Cov}(A, DA) (a - E[A]) p_a) .
\end{aligned} \tag{B-40}$$

Moreover, $\sum_{a \in \mathcal{A}} \omega_{\text{dce}}^{\text{i}*}(a) = 1$ and $\sum_{a \in \mathcal{A}} \omega_{\text{ind}}^{\text{i}*}(a) = 0$. Furthermore, it is possible to have $\omega_{\text{dce}}^{\text{i}*}(a) < 0$ and $\omega_{\text{ind}}^{\text{i}*}(a) \neq 0$ for some $a \in \mathcal{A}$.

Proof. By properties of projection,

$$(\text{Var}(W))^{-1} (\text{Cov}(W, Y) - \text{Cov}(W, D) \Delta_{\text{inter}}) = \alpha = (\theta', \lambda')' .$$

We can use the fact that A is scalar to obtain an explicit formula for $(\text{Var}(W))^{-1}$. With this expression in hand, we get

$$\Delta_{\text{inter}} + E[A] \lambda = \Psi \Delta_{\text{inter}} + \frac{E[A]}{\Delta} (\text{Var}(A) \text{Cov}(DA, Y) - \text{Cov}(A, DA) \text{Cov}(A, Y)) , \tag{B-41}$$

By plugging in the expressions for (A-9), (A-16), (A-17) on the right-hand side of (B-41), imposing Assumption 2.2, we obtain (B-39) and (B-40).

By the definition of $\{(\omega_{\text{dce}}^{\text{i}*}(a), \omega_{\text{ind}}^{\text{i}*}(a)) : a \in \mathcal{A}\}$ in (B-40) and repeating arguments used in the proof of Theorem A.2, it is immediate to show that $\sum_{a \in \mathcal{A}} \omega_{\text{dce}}^{\text{i}*}(a) = 1$ and $\sum_{a \in \mathcal{A}} \omega_{\text{ind}}^{\text{i}*}(a) = 0$.

To conclude, it suffices to find an example in which $\omega_{\text{dce}}^{\text{i}*}(a) < 0$ and $\omega_{\text{ind}}^{\text{i}*}(a) \neq 0$ for some $a \in \mathcal{A}$. To this end, consider an example with $\{A|D=0\} \sim \text{Bi}(2, 0.9)$, $\{A|D=1\} \sim \text{Bi}(2, 0.1)$, and $P\{D=1\} = 0.3$, and so $\mathcal{A} = \{0, 1, 2\}$. By evaluating this information on (B-40), we get $\omega_{\text{dce}} \approx [-0.2, 1.08, 0.12]$ and $\omega_{\text{ind}} \approx [-0.26, 0.52, -0.26]$, i.e., (b) and (c) fail. ■

Lemma B.5. Assume the conditions in Theorem 4.3, and that

$$\mu(d, a) = \kappa_0 + \kappa_1 d + \kappa_2 a + \kappa_3' a d \quad \text{for all } (d, a) \in \{0, 1\} \times \mathcal{A} \tag{B-42}$$

for some constants $\kappa_0, \kappa_1, \kappa_2, \kappa_3$. Then, the coefficient in (10) satisfies $\Delta_{\text{inter}} = \kappa_1$, $\theta_0 = \kappa_0$, $\theta = \kappa_2$, and $\lambda = \kappa_3$. Furthermore, the decomposition in (34) are $\Delta_{\text{dce}}^{\text{i}} = \kappa_1$ and $\Delta_{\text{ind}}^{\text{i}} = 0$.

Proof. Assumption 2.2 implies that $E(Y|D=d, A=a) = \mu(d, a)$ which, combined with (B-42), implies that the conditional expectation of Y is linear in $(1, a, d, ad)$. From here, the first result follows from the fact that the linear regression consistently estimates the parameters of a linear conditional expectation. The second part follows from combining $\sum_{a \in \mathcal{A}} a \omega_{\text{ind}}^{\text{i}}(a) = \mathbf{0}$ (shown in Theorem A.2) and (B-42). ■

References

- AKHTARI, M., CHEN, J., LEMIONET, A., NGUYEN, D., OBEID, H. and ZHU, Y. (2021). How airbnb measures future value to standardize tradeoffs. *Medium.com*. URL <https://medium.com/airbnb-engineering/how-airbnb-measures-future-value-to-standardize->
- BARON, R. M. and KENNY, D. A. (1986). The moderator–mediator variable distinction in social psychological research: Conceptual, strategic, and statistical considerations. *Journal of personality and social psychology*, **51** 1173.
- BEAMAN, L., KARLAN, D., THUYSBAERT, B. and UDRY, C. (2013). Profitability of Fertilizer: Experimental Evidence from Female Rice Farmers in Mali. *American Economic Review*, **103** 381–386.
- BUGNI, F. A., CANAY, I. A. and SHAIKH, A. M. (2018). Inference under covariate adaptive randomization. *Journal of the American Statistical Association (Theory & Methods)*, **113** 1741–1768.
- BUGNI, F. A., CANAY, I. A. and SHAIKH, A. M. (2019). Inference under covariate-adaptive randomization with multiple treatments. *Quantitative Economics*, **10** 1741–1768.
- CHERNOZHUKOV, V., KASAHARA, H. and SCHRIMPF, P. (2021). Causal impact of masks, policies, behavior on early covid-19 pandemic in the us. *Journal of econometrics*, **220** 23–62.
- DUFLO, E., KREMER, M. and ROBINSON, J. (2011). Nudging Farmers to Use Fertilizer: Theory and Experimental Evidence from Kenya. *American Economic Review*, **101** 2350–2390.
- FAGERENG, A., MOGSTAD, M. and RØNNING, M. (2021). Why do wealthy parents have wealthy children? *Journal of Political Economy*, **129** 703–756.
- GLYNN, A. N. (2012). The product and difference fallacies for indirect effects. *American Journal of Political Science*, **56** 257–269.
- HECKMAN, J., PINTO, R. and SAVELYEV, P. (2013). Understanding the mechanisms through which an influential early childhood program boosted adult outcomes. *American Economic Review*, **103** 2052–86.
- HECKMAN, J. J. (2000). Causal Parameters and Policy Analysis in Economics: A Twentieth Century Retrospective*. *Quarterly Journal of Economics*, **115** 45–97.
- HERNAN, M. and ROBINS, J. (2023). *Causal Inference: What If*. SCRC Press. <https://doi.org/10.1201/9781315374932>.

- HUANG, J., REILEY, D. and RIABOV, N. (2018). Measuring consumer sensitivity to audio advertising: A field experiment on pandora internet radio. *Available at SSRN 3166676*.
- IMAI, K., KEELE, L. and YAMAMOTO, T. (2010). Identification, inference and sensitivity analysis for causal mediation effects. *Statistical science*, **25** 51–71.
- JAIN, D. and SINGH, S. S. (2002). Customer lifetime value research in marketing: A review and future directions. *Journal of interactive marketing*, **16** 34–46.
- JUDD, C. M. and KENNY, D. A. (1981). Process analysis: Estimating mediation in treatment evaluations. *Evaluation review*, **5** 602–619.
- KRAEMER, H., KIERNAN, M., ESSEX, M. and KUPFER, D. J. (2008). How and why criteria defining moderators and mediators differ between the baron & kenny and macarthur approaches. *Health Psychology*, **27** S101.
- KRAEMER, H. C., WILSON, G. T., FAIRBURN, C. G. and AGRAS, W. S. (2002). Mediators and moderators of treatment effects in randomized clinical trials. *Archives of general psychiatry*, **59** 877–883.
- MANSKI, C. F. (1997). Monotone Treatment Response. *Econometrica*, **65** 1311.
- MEL, S. D., MCKENZIE, D. and WOODRUFF, C. (2009). Returns to Capital in Microenterprises: Evidence from a Field Experiment. *The Quarterly Journal of Economics*, **124** 423–423.
- MODERNA (2021). A Study of SARS CoV-2 Infection and Potential Transmission in Individuals Immunized With Moderna COVID-19 Vaccine (CoVPN 3006). *ClinicalTrials.gov Identifier: NCT04811664*. URL <https://clinicaltrials.gov/ct2/show/study/NCT04811664>.
- PEARL, J. (2001). Direct and indirect effects. In *Proceedings of the Seventeenth Conference on Uncertainty in Artificial Intelligence*. UAI’01, Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 411–420.
- ROBINS, J. M. (2003). Semantics of causal dag models and the identification of direct and indirect effects. In *Highly Structured Stochastic Systems* (N. L. H. P. J. Green and S. Richardson, eds.). Oxford University Press, 70–81.
- ROBINS, J. M. and GREENLAND, S. (1992). Identifiability and exchangeability for direct and indirect effects. *Epidemiology* 143–155.
- ROSENZWEIG, M. R. and WOLPIN, K. I. (2000). Natural “Natural Experiments” in Economics. *Journal of Economic Literature*, **38** 827–874.