
NONPARAMETRIC CAUSAL MEDIATION ANALYSIS FOR STOCHASTIC INTERVENTIONAL (IN)DIRECT EFFECTS

Nima S. Hejazi

Division of Biostatistics,
Department of Population Health Sciences,
Weill Cornell Medicine
nhejazi@berkeley.edu

Kara E. Rudolph

Department of Epidemiology,
Mailman School of Public Health,
Columbia University
kr2854@cumc.columbia.edu

Mark J. van der Laan

Division of Biostatistics,
School of Public Health, and
Department of Statistics,
University of California, Berkeley
laan@berkeley.edu

Iván Díaz

Division of Biostatistics,
Department of Population Health Sciences,
Weill Cornell Medicine
ild2005@med.cornell.edu

January 13, 2022

ABSTRACT

Causal mediation analysis has historically been limited in two important ways: (i) a focus has traditionally been placed on binary exposures and static interventions, and (ii) direct and indirect effect decompositions have been pursued that are only identifiable in the absence of intermediate confounders affected by exposure. We present a theoretical study of an (in)direct effect decomposition of the population intervention effect, defined by stochastic interventions jointly applied to the exposure and mediators. In contrast to existing proposals, our causal effects can be evaluated regardless of whether an exposure is categorical or continuous and remain well-defined even in the presence of intermediate confounders affected by exposure. Our (in)direct effects are identifiable without a restrictive assumption on cross-world counterfactual independencies, allowing for substantive conclusions drawn from them to be validated in randomized controlled trials. Beyond the novel effects introduced, we provide a careful study of nonparametric efficiency theory relevant for the construction of flexible, multiply robust estimators of our (in)direct effects, while avoiding undue restrictions induced by assuming parametric models of nuisance parameter functionals. To complement our nonparametric estimation strategy, we introduce inferential techniques for constructing confidence intervals and hypothesis tests, and discuss open source software, the `medshift` R package, implementing the proposed methodology. Application of our (in)direct effects and their nonparametric estimators is illustrated using data from a comparative effectiveness trial examining the direct and indirect effects of pharmacological therapeutics on relapse to opioid use disorder.

1 Introduction

In myriad applications, one is often interested in the effect of an exposure on an outcome only through a particular pathway between the two. Indeed, efforts in defining and identifying such *path-specific* effects have come to constitute a rich history in not only philosophy but also in the sciences of statistics, causal inference, epidemiology, economics, and psychology. In each of these disciplines, and in many others among the biomedical and social sciences, developing a mechanistic understanding of the complexities that admit representations as path-specific effects remains central; examples include elucidating the biological mechanism by which a vaccine reduces infection risk [e.g., Hejazi et al., 2020], assessing the effect on preterm birth of maternal exposure to environmental toxins, and ascertaining the effect of novel pharmacological therapies on substance abuse disorder relapse.

The latter serves as our motivating example as we consider how exposure to a buprenorphine dose schedule characterized by successive increases toward a maximum dose early in treatment (versus static dose) affects the risk of relapse to opioid use disorder, both directly and indirectly through mediating factors such as depression and pain. Developing a detailed mechanistic understanding of the process by which such therapeutics modulate intermediary states is necessarily a *causal* question — one central to designing and successively improving upon available therapies in a manner targeted towards the mitigation of the risk of substance abuse relapse. In comparative effectiveness trials of promising opioid use disorder therapeutics, detailed dissections of the complex neurological and psychiatric pathways involved in the development of addiction disorders is of clinical interest [Lee et al., 2018, Rudolph et al., 2020a]. The ability to define and evaluate causal effects along paths involving or avoiding mediating neuropsychiatric sequela would facilitate drug efficacy assessments; moreover, the ability to refine scientific conclusions based on statistical evidence through randomized controlled trials remains integral to furthering clinical progress.

To carefully study complex mediation relationships, a wealth of techniques rooted in statistical causal inference have been formulated. Path analysis [Wright, 1934], perhaps the earliest example of such methodology, directly inspired the development of subsequent techniques that leveraged parametric structural equation models [e.g., Goldberger, 1972, Baron and Kenny, 1986] for mediation analysis. More recently, the advent of modern frameworks and formalisms for causal inference, including nonparametric structural equation models, directed acyclic graphs, and their underlying do-calculus [Pearl, 1995, 2000], provided the necessary foundational tools to express causal mechanisms without reliance on more restrictive approaches tied to parametric modeling.

In tandem with the developments of Pearl [2000], similar approaches spearheaded by Robins [1986], Spirtes et al. [2000], and Dawid [2000] allowed nonparametric formulations of mediation analysis and uncovered significant limitations of the earlier efforts focused on structural equation models [Imai et al., 2010]. Recent applications of modern causal models have illustrated the failings of popular parametric modeling strategies [i.e., Baron and Kenny, 1986], in the presence of intermediate confounders of the mediator-outcome relationship [Cole and Hernán, 2002]. Consequently, the usually implausible assumptions that underlie such restrictive structural equation models make these approaches of limited applicability for examining complex phenomena in the biomedical and health sciences.

Modern approaches to causal inference have allowed for significant advances over the methodology of traditional mediation analysis, overcoming the significant restrictions imposed by the use of parametric structural equation modeling. For example, Robins and Greenland [1992] and Pearl [2006], using distinct frameworks, provided equivalent nonparametric decompositions of the average treatment effect (for binary exposures) into the *natural* direct and indirect effects, which quantify all effects of the exposure on the outcome through paths avoiding the mediator and all paths involving the mediator, respectively. Such advances were not without their limitations, however. A key assumption of the nonparametric decomposition of the average treatment effect is the requirement of *cross-world* counterfactual independencies (i.e., independence of counterfactuals indexed by distinct interventions). Unfortunately, such an assumption limits the scientific relevance of the natural (in)direct effects by making them unidentifiable in randomized trials, directly implying that corresponding scientific claims cannot be falsified through experimentation [Popper, 1934, Dawid, 2000]. Importantly, such cross-world independencies are also unsatisfied in the presence of intermediate confounders affected by exposure [Avin et al., 2005, Tchetgen Tchetgen and VanderWeele, 2014]. Given that such confounders are challenging to rule out in practice, the natural (in)direct effects are of limited applicability in real-world data analysis. This incompatibility motivated the recent development of a rich family of *interventional* (in)direct effects [Didelez et al., 2006, VanderWeele et al., 2014, Vansteelandt and Daniel, 2017, Rudolph et al., 2017, Nguyen et al., 2021], which utilize a flexible joint intervention strategy to retain identifiability in the presence of such confounding. Until quite recently, nonparametric effect decompositions and efficiency theory were unavailable for this class of effects, though recent efforts by Díaz et al. [2020] and Benkeser and Ran [2021] resolved this gap in the literature. Like their natural effect counterparts, the interventional effects are limited to settings with binary exposures. Our work outlines a general class of causal (in)direct effect estimands that do not require the cross-world independence condition and are robust to intermediate confounding (like the interventional effects), though our effect definitions are capable of readily accommodating exposure variables of all varieties, resolving a significant practical limitation of both classes of (in)direct effects.

A related thread of the literature has considered stochastic interventions, which generalize many intervention classes. For example, within this framework, static interventions result in post-intervention exposures that have degenerate distributions. Stock [1989] first considered the estimation of the total effects of stochastic interventions, while many others [e.g., Didelez et al., 2006, Haneuse and Rotnitzky, 2013, Young et al., 2014] provided careful studies that expanded the underlying theory of stochastic interventions and demonstrated their numerous applications. Within the population intervention models framework [Hubbard and van der Laan, 2008], Díaz and van der Laan [2012] formulated total causal effects attributable to continuous-valued exposures using a particular class of stochastic intervention. Conveniently, these causal effects of stochastic interventions carry an interpretation echoing that of standard regression adjustment. For example, Haneuse and Rotnitzky [2013] described modified treatment policies, which assign

post-intervention counterfactuals based on the natural value of the exposure; their methods were demonstrated in the context of reducing surgical time for non-small-cell lung cancer operations. Stochastic interventions have also successfully been applied to binary exposures: Kennedy [2019] proposed incremental propensity score interventions and demonstrated their use in longitudinal studies in order to circumvent identifiability and estimation issues arising from positivity violations. Building on this flexible framework, Díaz and Hejazi [2020] proposed a decomposition of the total effect of stochastic interventions [Díaz and van der Laan, 2012] into the *population intervention (in)direct effects*, which are endowed with interpretations analogous to that of the natural (in)direct effects. The (in)direct effects of Díaz and Hejazi [2020] do not require cross-world counterfactual independencies, apply to exposure variables of all types, and succeed in accommodating nonparametric estimation strategies. Consequently, their population intervention (in)direct effects may be estimated without restrictive assumptions and yield scientific results that can be tested through randomization of both the exposure and mediator. Unfortunately, the results of Díaz and Hejazi [2020] suffer a serious shortcoming — these effects cannot be identified in the presence of mediator–outcome confounders affected by exposure. In this vein, our work formulates alternative (in)direct effect estimands that retain the flexibility of the (in)direct effects of Díaz and Hejazi [2020]; however, our identification strategy emphasizes effects robust to this form of confounding, which is accomplished by leveraging joint stochastic interventions on the exposure and mediator.

In the present work, we outline a general framework encompassing many prior causal mediation analysis approaches, including the natural (in)direct effects, their interventional effect counterparts, and the stochastic (in)direct effects. Building upon the foundations laid by Díaz and Hejazi [2020], the introduced class of mediation effects originate from combining the novel lines of inquiry established in the distinct literatures on stochastic interventions and the interventional effects; accordingly, we denote these *stochastic interventional (in)direct effects*. Our proposed class of effects are the first to simultaneously avoid the requirement of cross-world counterfactual independencies; leverage stochastic interventions to be applicable to binary, categorical, and continuous-valued exposures; and remain identifiable despite intermediate confounding. Our contributions apply to a broader class of exposures than the interventional effects [e.g., Díaz et al., 2020, Benkeser and Ran, 2021] while generalizing stochastic (in)direct effects [i.e., Díaz and Hejazi, 2020] to accommodate the presence of intermediate confounders. While our robust and flexible causal mediation analysis framework subsumes prior classes of effect definitions, this is far from enough for the successful application of our proposed (in)direct effects. To this end, we develop novel efficiency theory and efficient nonparametric estimators of this broad class of causal mediation parameters, within the frameworks of one-step [Pfanzagl and Wefelmeyer, 1985, Bickel et al., 1993] and targeted minimum loss estimation [van der Laan and Rubin, 2006, van der Laan and Rose, 2011]. These flexible estimators have desirable asymptotic properties even when nuisance parameter functionals are estimated via machine learning; moreover, they are endowed with a form of multiple robustness producing consistent point estimates under several configurations of nuisance parameter misspecification. Lastly, we provide implementations of our methodological advances in our free and open source `medshift` [Hejazi and Díaz, 2020] package, for the R language and environment for statistical computing [R Core Team, 2022].

2 Mediation analysis for the population intervention effects

Let A denote a continuous or categorical exposure, Y denote a continuous or binary outcome, Z denote mediator(s), W denote a vector of observed pre-exposure covariates, and L denote an intermediate (mediator–outcome) confounder affected by exposure. The nonparametric structural equation model (NPSEM) formalizes the problem:

$$W = f_W(U_W); A = f_A(W, U_A); L = f_L(A, W, U_L); Z = f_Z(L, A, W, U_Z); Y = f_Y(Z, L, A, W, U_Y). \quad (1)$$

In the NPSEM (1), $U = (U_W, U_A, U_L, U_Z, U_Y)$ is a vector of exogenous factors, and the functions f are assumed deterministic but unknown. This mechanistic model is assumed to generate the observed data O ; it encodes several fundamental assumptions. First, an implicit temporal ordering $W \rightarrow A \rightarrow L \rightarrow Z \rightarrow Y$ is assumed. Second, each variable (i.e., $\{W, A, L, Z, Y\}$) is assumed to be generated from the corresponding deterministic function of the observed variables that precede it temporally, plus an exogenous variable denoted by U . Each exogenous variable is assumed to contain all unobserved causes of the corresponding observed variable. For a random variable X , let X_a denote the counterfactual outcome observed in a hypothetical world in which $P(A = a) = 1$. For example, we have $L_a = f_L(a, W, U_L)$, $Z_a = f_Z(L_a, a, W, U_Z)$, and $Y_a = f_Y(Z_a, L_a, a, W, U_Y)$. Likewise, we let $Y_{a,z} = f_Y(z, L_a, a, W, U_Y)$ denote the value of the outcome in a hypothetical world where $P(A = a, Z = z) = 1$. Figure 1 represents model (1) in terms of a directed acyclic graph (DAG).

Letting $O = (W, A, L, Z, Y)$ represent a random variable with distribution P , we denote by O_1, \dots, O_n a sample of n i.i.d. observations of O . We let $Pf = \int f(o)dP(o)$ for a given function $f(o)$. We use P_c to denote the joint distribution of (O, U) , and let E and E_c denote corresponding expectation operators. We use P_n to denote the empirical distribution of O_1, \dots, O_n , and assume $P \in \mathcal{M}$, where \mathcal{M} is the nonparametric statistical model defined as all continuous densities on O with respect to a dominating measure ν . Let p denote the corresponding probability density function. We use

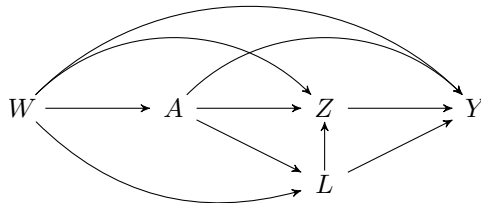


Figure 1: Directed Acyclic Graph of NPSEM (1).

$g(a | w)$ and $e(a | z, w)$ to denote the probability density function or the probability mass function of A conditional on $W = w$ and (Z, W) , respectively; $m(z, l, a, w)$ to denote the outcome regression function $E(Y | Z = z, L = l, A = a, W = w)$. Let $g(\cdot | w)$ and $e(\cdot | z, w)$ be dominated by a measure $\kappa(a)$ (e.g., the counting measure for binary A and the Lebesgue measure for continuous A). In constructing our estimators, we will use

$$\frac{p(z | w)}{p(z | a, w)} = \frac{g(a | w)}{e(a | z, w)}; \quad \frac{p(z | a, w)}{p(z | l, a, w)} = \frac{p(l | a, w)}{p(l | z, a, w)} \quad (2)$$

as such parameterizations allow for estimation and integration with respect to multivariate conditional densities on the mediator Z to be avoided. We use $\mathcal{W}, \mathcal{A}, \mathcal{L}, \mathcal{Z}$, and \mathcal{Y} to denote the support of the corresponding random variables.

Causal effects are defined in terms of hypothetical interventions on the NPSEM (1). In particular, consider an intervention in which the structural equation corresponding to A is removed, with the exposure drawn instead from a user-specified distribution $g_\delta(a | w)$, which may itself depend on the natural exposure distribution and a user-specified parameter δ . Going forward, we let A_δ denote a draw from $g_\delta(a | w)$. Alternatively, such modifications can occasionally be described in terms of an intervention in which the structural equation corresponding to A is removed and the exposure is set equal to a hypothetical regime $d(A, W)$. Regime d depends on the exposure level A that would have been assigned in the absence of the regime as well as on W . The latter intervention has been referred to as depending on the *natural value of treatment*, or as a *modified treatment policy* [MTP; Haneuse and Rotnitzky, 2013]. For such interventions, Haneuse and Rotnitzky [2013] introduced the assumption of *piecewise smooth invertibility*, which ensures that the change of variable formula can be used when computing integrals over A :

A1 (Piecewise smooth invertibility). For each $w \in \mathcal{W}$, assume that the interval $\mathcal{I}(w) = (l(w), u(w))$ may be partitioned into subintervals $\mathcal{I}_{\delta,j}(w) : j = 1, \dots, J(w)$ such that $d(a, w)$ is equal to some $d_j(a, w)$ in $\mathcal{I}_{\delta,j}(w)$ and $d_j(\cdot, w)$ has inverse function $h_j(\cdot, w)$ with derivative $h'_j(\cdot, w)$.

Assumption A1 can be used to show that the intervention drawing A_δ from the post-intervention distribution $g_\delta(a | w)$ can be interpreted on the individual level. Young et al. [2014] provide a discussion comparing and contrasting the interpretation and identification of these two interventions. Such stochastic interventions can be used to define the *population intervention effect (PIE)* of A on Y . To illustrate, suppose A to be continuous-valued and assume the distribution of A conditional on $W = w$ is supported in the interval $(l(w), u(w))$. Then, one may define

$$d(a, w) = \begin{cases} a - \delta & \text{if } a > l(w) + \delta \\ a & \text{if } a \leq l(w) + \delta, \end{cases} \quad (3)$$

where $0 < \delta < u(w)$ is an arbitrary prespecified value. We can alternatively define a tilted intervention distribution as

$$g_\delta(a | w) = \frac{\exp(\delta a) g(a | w)}{\int \exp(\delta a) g(a | w) d\kappa(a)}, \quad (4)$$

for $\delta \in \mathbb{R}$. Kennedy [2019] proposed a form of exponential tilting (4) under the parameterization $\delta' = \exp(\delta)$, appropriate for incremental interventions on the propensity score for binary A . Two key distinctions between the interventions defining modified treatment policies (3) and exponential tilting (4) can be helpful in differentiating between the two in practice. Firstly, through Assumption A1, intervention (3) defines counterfactual exposures equipped with an individual-level interpretation, whereas the intervention (4) only admits interpretation as a random draw from a post-intervention distribution. Secondly, intervention (3) has been historically studied in settings with continuous (or ordinal) exposures [e.g., Díaz and van der Laan, 2013, Hejazi et al., 2020], though it may apply more broadly. While intervention (4) applies readily to exposures of any type (e.g., continuous, categorical), Kennedy [2019] introduced it to construct total effects with weakened positivity requirements in longitudinal studies with time-varying binary exposures. Díaz and Hejazi [2020] provide a careful study of these interventions for mediation analysis, introducing novel (in)direct effects and efficiency theory. Their contributions assume the absence of intermediate confounding; our generalization remedies this inconvenient shortcoming.

2.1 Stochastic Mediation Effects

Díaz and Hejazi [2020] defined the (in)direct effect of A on Y in terms of a decomposition of the total effect of a stochastic intervention. In particular, the total effect $E(Y - Y_{A_\delta})$ may be decomposed as the sum of the population intervention direct and indirect effects (PIDE; PIIE):

$$\begin{aligned} \text{PIDE} &= E_c\{f_Y(Z, L, A, W, U_Y) - f_Y(Z, L_{A_\delta}, A_\delta, W, U_Y)\} \\ \text{PIIE} &= E_c\{f_Y(Z, L_{A_\delta}, A_\delta, W, U_Y) - f_Y(Z_{A_\delta}, L_{A_\delta}, A_\delta, W, U_Y)\}. \end{aligned} \quad (5)$$

Upon inspection, the definitions above reveal that the direct effect measures the effect through paths *not* involving the mediator (i.e., $A \rightarrow Y$ and $A \rightarrow L \rightarrow Y$), whereas the indirect effect measures the effect through paths involving the mediator (i.e., $A \rightarrow Z \rightarrow Y$ and $A \rightarrow L \rightarrow Z \rightarrow Y$).

Unfortunately, the population intervention (in)direct effects are not generally identified in the presence of an intermediate confounder affected by exposure such as in the DAG in Figure 1 [Díaz and Hejazi, 2020]. This is due to the dual role of L as a confounder of the relation between Z and Y , which requires adjustment, and a variable on the path from A to Y , which precludes adjustment. Note that Vansteelandt and VanderWeele [2012] attempted to circumvent these restrictions in the context of direct effects, while Fulcher et al. [2019] complemented their study by introducing an alternative effect decomposition that yielded indirect effects with weaker identification requirements. In particular, Fulcher et al. [2019] contrive an indirect effect definition that is identifiable under unmeasured baseline confounding of the exposure–outcome relationship and that is formulated by the application of a stochastic intervention to the mediator. By contrast, Díaz and Hejazi [2020] leverage joint stochastic interventions on the exposure and mediator in their effect decomposition, which can be made to achieve a similar identification property. (Coincidentally, Fulcher et al. [2019] refer to their effects as “population intervention (in)direct effects,” terminology also used by Díaz and Hejazi [2020], though the effects of the former are in some ways more restrictive than those proposed by the latter.) The interventional effects [VanderWeele et al., 2014] resolve the identification issue brought on by L , though their limitation to static interventions acting upon binary exposures is a significant hurdle to their use. Next, we present a solution to this complication using a joint stochastic intervention on the exposure A and mediator Z . We also show that the effects defined in this manuscript are a generalization of the effects of Díaz and Hejazi [2020] in the sense that the former reduce to the latter in the absence of intermediate confounding.

2.2 Stochastic Interventional Mediation Effects

To introduce (in)direct effects robust to the presence of intermediate confounders, we draw upon ideas first outlined by Didelez et al. [2006] and van der Laan and Petersen [2008], later formalized or subsumed by VanderWeele et al. [2014] and Vansteelandt and Daniel [2017]. Owing to their definition in terms of stochastic interventions on the mediator, these (in)direct effects have been collectively termed *interventional effects*. We leverage two types of stochastic interventions: one on the exposure A , which defines the intervention of interest, and one on the mediator Z , which is used to achieve identifiability of the effects. Following the convention of the literature, we term stochastic interventions on the mediator *interventional*, while reserving the label of *stochastic* to refer only to interventions on the exposure A . To proceed, let G_δ denote a random draw from the distribution of Z_{A_δ} conditional on (A_δ, W) , and let G denote a random draw from the distribution of Z conditional on (A, W) . To distinguish G_δ from the previously defined $g_\delta(a | w)$, note that the latter is the post-intervention distribution of the exposure, based on the user-specified scalar δ , and gives rise to the counterfactual A_δ , while the former is a counterfactual arising from a draw from the *interventional distribution* of the mediator, which breaks the dependence of Z on L . These draws are denoted G_δ (when A_δ is used) or G (when A is used). We consider the effect defined by $\psi_\delta = E_c\{Y_{A,G} - Y_{A_\delta,G_\delta}\}$. Note that the effect ψ_δ is distinct from the effect considered by Díaz and Hejazi [2020], which may be expressed $E_c\{Y_{A,Z} - Y_{A_\delta,Z_\delta}\}$. The effect ψ_δ arises from fixing the mediator to a random value chosen from its distribution among all those with a particular exposure level, rather than fixing it to what it would have been under a particular exposure level. The choice of intervention on Z (i.e., defining G_δ , which does not depend on L) allows for the contribution of the intermediate confounder L upon the mediator Z to be eliminated. The reasoning behind this intervention is as follows. Considering the DAG (see Figure 1) corresponding to model (1), L satisfies the “recanting witness” criterion of Avin et al. [2005]. This introduces unidentifiability of the natural (in)direct effects. Removal of the directed path from L to Z in the interventional distribution resolves this complication, analogously to solutions presented by those authors for the identification of simpler path-specific effects. Defining the effect in this way aids in achieving an identifiable decomposition into direct and indirect effects. In particular, we may decompose this effect in terms of stochastic interventional *direct effects* (DE) and *indirect effects* (IE):

$$\psi(\delta) = \overbrace{E\{Y_{A,G} - Y_{A_\delta,G}\}}^{\text{DE}} + \overbrace{E\{Y_{A_\delta,G} - Y_{A_\delta,G_\delta}\}}^{\text{IE}}. \quad (6)$$

Decomposition as the sum of direct and indirect effects affords an interpretation analogous to the corresponding standard decomposition of the average treatment effect into the natural direct and indirect effects [Pearl, 2006]. In particular, the direct effect arises from drawing a counterfactual value of A from a post-intervention distribution while keeping the distribution of Z fixed. The indirect effect arises from replacing the distribution of Z with a candidate post-intervention distribution while holding A fixed. Our proposed stochastic interventional effects have an interpretation similar to the interventional effects of VanderWeele et al. [2014]; moreover, while both effect definitions account for the presence of an intermediate confounder, our (in)direct effects utilize flexible, stochastic interventions on the exposure while those of VanderWeele et al. [2014] are limited to static interventions on binary exposures. By generalizing the effect definitions of Díaz and Hejazi [2020], our proposed (in)direct effects include, as special cases, the natural (in)direct effects (under a static intervention on binary A and no intermediate confounders L), the interventional (in)direct effects (under a static intervention on binary A and a stochastic intervention on Z , allowing intermediate confounders L), and the stochastic (in)direct effects (under a stochastic intervention on arbitrary-valued A and no intermediate confounders L).

2.3 Identification

To construct estimators of our proposed causal (in)direct effects, we turn to examining assumptions needed to estimate components of the post-intervention quantities corresponding to counterfactuals of interest. First, we consider Assumptions A2–A6, which allow identification of the stochastic interventional effects of Equation (6):

A2 (Common support). Assume $\text{supp}\{g_\delta(\cdot | w)\} \subseteq \text{supp}\{g(\cdot | w)\}$ for all $w \in \mathcal{W}$.

A3 (Mediator positivity). Assume $p(z | a, w) > 0$ and further assume $p(z | l, a, w) > 0$.

A4 (No unmeasured exposure-outcome confounder). Assume $Y_{a,z} \perp\!\!\!\perp A | W$.

A5 (No unmeasured mediator-outcome confounder). Assume $Y_{a,z} \perp\!\!\!\perp Z | (L, A, W)$.

A6 (No unmeasured exposure-mediator confounder). Assume $Z_a \perp\!\!\!\perp A | W$.

Assumption A4 states that, conditional on W , there is no unmeasured confounding of the relation between A and Y ; Assumption A6 states that conditional on W there is no unmeasured confounding of the relation between A and Z ; and Assumption A5 states that conditional on (W, A, L) there is no unmeasured confounding of the relation between Z and Y . These assumptions are standard in causal mediation analysis. In addition to these assumptions, standard mediation analyses [e.g., VanderWeele et al., 2014] require positivity assumptions on the exposure and mediation mechanisms (e.g., Assumption A3). The stochastic intervention framework we adopt does not require such assumptions, as positivity can be arranged by definition of g_δ . For example, the interventions in expressions (3) and (4) satisfy Assumption A2 by definition [see, e.g., Kennedy, 2019, Hejazi et al., 2020, Díaz and Hejazi, 2020]. Notably, our effects do not require any assumption on the independence of cross-world counterfactuals, required for identification of the natural (in)direct effects. The cross-world independence assumption is restrictive, as it cannot be tested under randomization, significantly limiting the scientific relevance of these effects [Díaz and Hejazi, 2020, Popper, 1934]. Our novel effect definitions are hardly the first to circumvent this assumption: earlier work [e.g., Didelez et al., 2006, van der Laan and Petersen, 2008, Vansteelandt and VanderWeele, 2012] provided strategies for loosening the reliance of the natural (in)direct effects on the cross-world assumption. As with the earlier effect definitions of Díaz and Hejazi [2020], our stochastic interventional effects similarly eschew this restrictive condition for their identification. Under these assumptions, the following identification results hold; proofs appear in the Supplementary Materials.

Theorem 1 (Identification). *Define*

$$\begin{aligned}\theta_{1,\delta} &= \int m(z, l, a, w) p(l | a, w) p(z | a, w) g_\delta(a | w) p(w) d\nu(a, z, l, w), \\ \theta_{2,\delta} &= \int m(z, l, a, w) p(l | a, w) p(z | w) g_\delta(a | w) p(w) d\nu(a, z, l, w).\end{aligned}$$

Under Assumptions A2–A6, the direct effect $\psi_{D,\delta} = E\{Y_{A,G} - Y_{A_\delta,G}\}$ and indirect effect $\psi_{I,\delta} = E\{Y_{A_\delta,G} - Y_{A_\delta,G_\delta}\}$ (6) are identified, respectively, by

$$\psi_{D,\delta} = \theta_{1,0} - \theta_{2,\delta} \quad \text{and} \quad \psi_{I,\delta} = \theta_{2,\delta} - \theta_{1,\delta}. \quad (7)$$

A consequence of this identification result is that the definitions reduce to the stochastic (in)direct effects of Díaz and Hejazi [2020] in the absence of intermediate confounders L . Importantly, this implies that our estimators can be safely used in the absence of intermediate confounders; furthermore, it implies that the corresponding estimates may be interpreted in terms of a decomposition of the population intervention effect $E_c\{Y - Y_{A_\delta}\}$, which, like the interventional effect $\psi_\delta = E_c\{Y_{A,G} - Y_{A_\delta,G_\delta}\}$, may be of scientific relevance.

Examination of Definition (7) reveals that evaluation of $\psi_{D,\delta}$ and $\psi_{I,\delta}$ requires access to $\theta_{2,\delta}$, which is based on $g_\delta(a \mid w)$. The quantities $\psi_{D,\delta}$ and $\psi_{I,\delta}$ further require access, respectively, to either $\theta_{1,0}$ or $\theta_{1,\delta}$, which draw upon $g_\delta(a \mid w)$. Notably, $\theta_{1,0}$ is *not* the population mean outcome, rather it is the “interventional” counterfactual mean arising from changing the distribution of the mediator from $p(Z \mid L, A, W)$ to $p(Z \mid A, W)$, resulting in the counterfactuals G and G_δ . In fact, this induced independence of L and Z , conditional on A and W , recovers the effects of Díaz and Hejazi [2020] when there is no post-exposure confounder L . We next turn our attention to developing efficiency theory for estimation of the statistical parameter $\theta_{j,\delta} : j = 1, 2$, which depends on the observed data distribution P .

3 Optimality theory for estimation of the direct effect

Thus far, we have discussed the decomposition of the effect of a stochastic intervention into direct and indirect effects and have provided identification results under standard assumptions. Next, we develop efficiency theory for estimating $\theta_{1,\delta}$ and $\theta_{2,\delta}$ in the nonparametric model \mathcal{M} . To do so, we introduce the *efficient influence function* (EIF), which characterizes the asymptotic behavior of all regular and asymptotically linear estimators [Bickel et al., 1993]. Three common frameworks exist for constructing locally efficient estimators based on the EIF: (i) estimating equations [e.g., van der Laan and Robins, 2003], (ii) one-step bias correction [e.g., Pfanzagl and Wefelmeyer, 1985, Bickel et al., 1993], and targeted minimum loss estimation [van der Laan and Rubin, 2006, van der Laan and Rose, 2011].

As a consequence of its representation in terms of orthogonal score equations, the EIF allows the construction of consistent estimators of the target parameter even when certain nuisance components are inconsistently estimated. Second-order bias terms may be derived from asymptotic analysis of estimators constructed based on the EIF — often, these estimators require slow convergence rates (e.g., $n^{-1/4}$) for the nuisance parameters involved, in order for the estimators to be regular and asymptotically linear (thereby achieving a Gaussian limit distribution), achieve \sqrt{n} -consistency, and exhibit asymptotic efficiency. Importantly, it is this rate-convergence property that enables the use of flexible, data adaptive regression techniques in estimating these quantities.

In Theorem 2, we present the EIF for a general stochastic intervention. Although the components of the EIF associated with (W, L, Z, Y) are the same, the component associated with the model for the distribution of A must be computed on a case-by-case basis, that is, for each intervention of interest. Lemmas 1 and 2 present such components for modified treatment policies satisfying Assumption A1 and for exponential tilting, respectively. In Theorem 2 below, we present a representation of the EIF that circumvents the challenging computation of multivariate integrals over Z . To introduce the EIF, we define the following auxiliary nuisance parameters:

$$\begin{aligned} u(z, a, w) &= \int m(z, l, a, w) dP(l \mid a, w); & \bar{u}(a, w) &= \int u(z, a, w) dP(z \mid a, w) \\ v(l, a, w) &= \int m(z, l, a, w) dP(z \mid a, w); & \bar{v}(a, w) &= \int v(l, a, w) dP(l \mid a, w) \\ s(l, a, w) &= \int m(z, l, a, w) dP(z \mid w); & \bar{s}(a, w) &= \int s(l, a, w) dP(l \mid a, w) \end{aligned} \quad (8)$$

Proofs for the following results are detailed in the Supplementary Materials.

Theorem 2 (Efficient influence functions). *Define*

$$H_{P,\delta}^1(a, z, l, w) := \frac{g_\delta(a \mid w)}{g(a \mid w)} \frac{p(z \mid a, w)}{p(z \mid a, l, w)}; \quad H_{P,\delta}^2(a, z, l, w) := \frac{g_\delta(a \mid w)}{g(a \mid w)} \frac{p(z \mid w)}{p(z \mid a, l, w)}.$$

The efficient influence functions for $\theta_{j,\delta} : j = 1, 2$ in the nonparametric model are equal to $D_{P,\delta}^j(o) - \theta_{j,\delta}$, where $D_{P,\delta}^j(o) = S_{P,\delta}^j(o) + S_{P,\delta}^{j,A}(o)$ and

$$S_{P,\delta}^1(o) = H_{P,\delta}^1(a, z, l, w) \{y - m(z, l, a, w)\} \quad (9)$$

$$+ \frac{g_\delta(a | w)}{g(a | w)} [v(l, a, w) - \bar{v}(a, w) + u(z, a, w) - \bar{u}(a, w)] \quad (10)$$

$$+ \int \bar{u}(a, w) g_\delta(a | w) d\kappa(a)$$

$$S_{P,\delta}^2(o) = H_{P,\delta}^2(a, z, l, w) \{y - m(z, l, a, w)\} \quad (11)$$

$$+ \frac{g_\delta(a | w)}{g(a | w)} \{s(l, a, w) - \bar{s}(a, w)\} \quad (12)$$

$$+ \int u(z, a, w) g_\delta(a | w) d\kappa(a),$$

and $S_{P,\delta}^{1,A}(o)$, $S_{P,\delta}^{2,A}(o)$ are the respective efficient score functions of the model for $g(a | w)$.

An immediate consequence of Theorem 2 is that, in a randomized trial, $S_{P,\delta}^{j,A}(o) = 0$ for $j = 1, 2$; however, even in such trials, covariate adjustment can improve the efficiency of the resultant estimator [van der Laan and Robins, 2003]. We now present the efficient scores $S_{P,\delta}^{j,A}(o)$ for modified treatment policies and exponentially tilted stochastic interventions. To do so, we define the parameter $q(a, w) = \int u(z, a, w) dP(z | w)$.

Lemma 1 (Modified treatment policies). *If the modified treatment policy $d(A, W)$ satisfies Assumption A1, then*

$$S_{P,\delta}^{1,A}(o) = \bar{u}(d(a, w), w) - \int \bar{u}(d(a, w), w) g(a | w) d\kappa(a) \quad (13)$$

$$S_{P,\delta}^{2,A}(o) = q(d(a, w), w) - \int q(d(a, w), w) g(a | w) d\kappa(a). \quad (14)$$

Lemma 2 (Exponential tilt). *If the stochastic intervention is the exponential tilt (4), then*

$$S_{P,\delta}^{1,A}(o) = \frac{g_\delta(a | w)}{g(a | w)} \left\{ \bar{u}(a, w) - \int \bar{u}(a, w) g_\delta(a | w) d\kappa(a) \right\} \quad (15)$$

$$S_{P,\delta}^{2,A}(o) = \frac{g_\delta(a | w)}{g(a | w)} \left\{ q(a, w) - \int q(a, w) g_\delta(a | w) d\kappa(a) \right\} \quad (16)$$

For binary exposures, the EIF for the incremental propensity score intervention may be simplified as follows.

Lemma 3 (EIF for incremental propensity score interventions). *Let $A \in \{0, 1\}$ and let the exponentially tilted intervention $g_{\delta,0}(1 | W)$ be based on (4) under the parameterization $\delta' = \exp(\delta)$. Then, the EIF of Lemma 2 may be simplified as follows. Specifically, we have*

$$S_{\eta,\delta}^{j,A}(o) = \frac{\delta q^j(w) \{a - g(1 | w)\}}{\{\delta g(1 | w) + 1 - g(1 | w)\}^2}, \quad \text{where}$$

$$q^1(w) = \bar{u}(1, w) - \bar{u}(0, w), \quad \text{and} \quad q^2(w) = E \{u(Z, 1, W) - u(Z, 0, W) | W = w\}. \quad (17)$$

In contrast to the efficient influence function for the interventional (in)direct effects [Díaz et al., 2020], the contribution of the exposure mechanism to the EIF for the stochastic interventional effects is nonzero. This is a direct consequence of the fact that the parameter of interest depends on $g(a | w)$; moreover, this implies that the efficiency bound in observational studies differs from the efficiency bound in randomized trials. Thus, it is not generally possible to obtain estimating equations robust to inconsistent estimation of $g(a | w)$. Such robustness will only be possible if the stochastic intervention is also a modified treatment policy satisfying Assumption A1.

The form of Theorem 2 makes it clear that estimation of multivariate or continuous conditional density functions on the mediators Z or intermediate confounders L , as well as integrals with respect to these density functions, is generally necessary for computation of the EIF. This poses a significant challenge from the perspective of estimation, due to both the curse of dimensionality and the practical computational complexity inherent in solving multivariate

numerical integrals. A simplification is possible when either Z or L is low-dimensional; this is achieved by re-parameterizing the densities as conditional expectations (or low-dimensional conditional densities) that take other nuisance parameters as pseudo-outcomes. In cases where L or Z is low-dimensional, our proposed re-parameterizations allow for the conditional density to be estimated via appropriate semiparametric density estimation procedures [e.g., Díaz and van der Laan, 2011, Hejazi et al., 2021].

Lemma 4 (EIF for low-dimensional L). *Let L be low-dimensional and Z multivariate. A representation of v , s , and \bar{u} in terms of conditional expectations may be chosen in order to simplify their estimation. Denote by $b(l | a, w)$ and $d(l | z, a, w)$ the density of L conditional on (A, W) and (Z, A, W) , respectively. Then, using (2), we have*

$$\begin{aligned} v(l, a, w) &= E \left[m(z, l, a, w) \frac{b(L | A, W)}{d(L | Z, A, W)} \middle| L = l, A = a, W = w \right], \\ s(l, a, w) &= E \left[m(z, l, a, w) \frac{b(L | A, W)}{d(L | Z, A, W)} \frac{g(A | W)}{e(A | Z, W)} \middle| L = l, A = a, W = w \right], \\ \bar{u}(a, w) &= E \left[u(Z, A, W) \middle| A = a, W = w \right]. \end{aligned} \quad (18)$$

Likewise,

$$H_{P, \delta}^1(a, z, l, w) = \frac{g_\delta(a | w)}{g(a | w)} \frac{b(l | a, w)}{d(l | z, a, w)}, \quad \text{and} \quad H_{P, \delta}^2(a, z, l, w) = \frac{g_\delta(a | w)}{e(a | z, w)} \frac{b(l | a, w)}{d(l | z, a, w)}, \quad \text{then}$$

$$q(a, w) = E \left\{ \frac{g(A | W)}{e(A | Z, W)} u(Z, A, W) \middle| A = a, W = w \right\}.$$

Analogous representations may be constructed for \bar{v} , \bar{s} , and \bar{u} based on the parameterizations (2) if L is multivariate and Z is of low dimension. We note, however, that at least one of Z or L must be of low dimensionality for its density to be easily estimable and integrals over its range computed with relative ease. Henceforth, denote by $\eta = (m, g, b, \bar{v}, \bar{s}, \bar{u}, d, e, s, q)$ and let $D_{P, \delta}^j(o) = D_{\eta, \delta}^j(o)$.

We note that the choice of parameterization in Lemma 4 has important consequences for the purpose of estimation, as it helps to bypass estimation of the (possibly high-dimensional) conditional density of the mediators, by requiring only that the intermediate confounders be of modest dimensionality for the purpose of estimation. In the particularly simple case that $L = l \in \{0, 1\}$ (as in our motivating application), the nuisance quantities $b(l | a, w)$ and $d(l | z, a, w)$ reduce to conditional expectations, allowing for regression methods, far more readily available throughout the statistics literature and software, to be used for their estimation. In addition to the expression for the EIF in Lemma 4, it is important to understand the behavior of the difference $PD_{\eta_1} - \theta$, which is expected to yield a second-order term in differences $\eta_1 - \eta$, so that consistent estimation of θ is possible under consistent estimation of certain configurations of the parameters in η . As we will see in Theorems 3 and 4, this second-order term is fundamental in the construction of asymptotically linear estimators. Lemmas S1 and S2, in the Supplementary Materials, delineate these second-order terms. The following lemma is a consequence.

Lemma 5 (Multiple robustness for modified treatment policies). *Let the modified treatment policy satisfy A1, and let η_1 be such that one of the following conditions hold:*

	$m_1 = m$	$g_1 = g$	$b_1 = b$	$\bar{u}_1 = \bar{u}$	$v_1 = v$	$d_1 = d$	$e_1 = e$	$s_1 = s$	$q_1 = q$
Cond. 1	×	×	×						
Cond. 2	×	×			×			×	
Cond. 3		×	×			×	×		
Cond. 4		×		×	×	×	×		
Cond. 5	×		×	×					×
Cond. 6	×			×	×			×	×

Table 1: Different configurations of consistency for nuisance parameters

Then $PD_{\eta_1, \delta}^1 = \theta_{1, \delta}$ and $PD_{\eta_1, \delta}^2 = \theta_{2, \delta}$, with $D_{\eta, \delta}^1$ and $D_{\eta, \delta}^2$ as defined in Theorem 2 and Lemma 1.

The above lemma implies that it is possible to construct consistent estimators for the (in)direct effects under consistent estimation of subsets of the nuisance parameters in η , in the configurations described in the lemma. Lemma 5 follows directly from Lemma S1, found in the Supplementary Materials. It may be surprising that estimation of $\theta_{j,\delta}$ can be robust to inconsistent estimation of g , even when the parameter definitions are explicitly dependent on g . We offer some intuition for this result by noting that Assumption A1 allows use of the change of variable formula to obtain $\theta_{2,\delta} = E \left\{ \int m(z, l, d(A, W), W) p(l | d(A, W), W) p(z | W) d\nu(z, l) \right\}$. Estimation of this parameter without relying on g may be carried out by consistently estimating $m(z, l, a, w)$, $p(l | a, w)$, and $p(z | w)$ and using the empirical distribution as an estimator of the outer expectation. This behavior has been previously observed for related modified treatment policy effects A1 [Díaz and van der Laan, 2012, Haneuse and Rotnitzky, 2013, Díaz and Hejazi, 2020].

Robustness for exponentially tilted interventions (1), not satisfying Assumption A1, appears in Lemma 6.

Lemma 6 (Multiple robustness for exponential tilting). *Let g_δ be defined as in (4) and η_1 be such that one or more of Cond. 1-4 in Table 1 holds, then $PD_{\eta_1,\delta}^1 = \theta_{1,\delta}$ and $PD_{\eta_1,\delta}^2 = \theta_{2,\delta}$, with $D_{\eta_1,\delta}^1$ and $D_{\eta_1,\delta}^2$ as defined in Theorem 2 and Lemma 2*

Lemma 6 is a direct consequence of Lemma S2 in the Supplementary Materials. The corresponding proof reveals that the EIF for the binary distribution is not robust to inconsistent estimation of g — that is, the intervention fails to satisfy Assumption A1 and integrals over the range of A cannot be computed using the change of variable formula. This behavior has been previously observed for other interventions that do not satisfy Assumption A1. Even though this lemma implies that consistent estimation of g is required, the bias terms remain second-order; thus, an estimator of g converging at rate $n^{1/4}$ or faster is sufficient.

4 Efficient estimation and statistical inference

We discuss two efficient estimators that rely on the efficient influence function $D_{\eta,\delta}$, in order to build an estimator that is both asymptotically efficient and robust to model misspecification. We discuss an asymptotic linearity result for the doubly robust estimator that allows computation of asymptotically accurate Wald-style confidence intervals and hypothesis tests. In the sequel, we assume that preliminary estimators of the components of η are available. These estimators may be obtained from flexible regression techniques such as neural networks, regression trees, boosting, splines, or ensembles thereof [Breiman, 1996, van der Laan et al., 2007]. The consistency of these estimators determines consistency of our estimators of $\theta_{j,\delta}$.

Both of our proposed efficient estimators make use of the EIF $D_{\eta,\delta}$ to revise an initial substitution estimator through a bias correction step. Estimation proceeds by first constructing initial estimators of the nuisance parameters in η ; then, each of the efficient estimators is constructed by application of distinct bias-correction steps. In this process, we advocate for the use of cross-fitting [Klaassen, 1987, Zheng and van der Laan, 2011] to avoid imposing entropy conditions on the initial estimators of the nuisance parameters in η . Let $\mathcal{V}_1, \dots, \mathcal{V}_J$ denote a random partition of the index set $\{1, \dots, n\}$ into J prediction sets of approximately the same size. That is, $\mathcal{V}_j \subset \{1, \dots, n\}$; $\bigcup_{j=1}^J \mathcal{V}_j = \{1, \dots, n\}$; and $\mathcal{V}_j \cap \mathcal{V}_{j'} = \emptyset$. For each j , the associated training sample is given by $\mathcal{T}_j = \{1, \dots, n\} \setminus \mathcal{V}_j$, and we let $j(i)$ denote the index of the validation set containing observation i . Denote by $\hat{\eta}_j$ the estimator of η obtained by training a prediction algorithm using only data in the sample \mathcal{T}_j .

4.1 Efficient One-Step Estimator

To construct a robust and efficient estimator using the efficient influence function $D_{\eta,\delta}$, the one-step bias correction [Pfanzagl and Wefelmeyer, 1985, Bickel et al., 1993] adds the empirical mean of the estimated EIF $D_{\hat{\eta},\delta}$ to an initial substitution estimator. The estimators are thus defined

$$\hat{\psi}_{D,\delta}^{os} = \frac{1}{n} \sum_{i=1}^n \{D_{\hat{\eta}_{j(i)},\delta}^1(O_i) - D_{\hat{\eta}_{j(i)},\delta}^2(O_i)\} \quad \text{and} \quad \hat{\psi}_{I,\delta}^{os} = \frac{1}{n} \sum_{i=1}^n \{D_{\hat{\eta}_{j(i)},\delta}^2(O_i) - D_{\hat{\eta}_{j(i)},\delta}^1(O_i)\}. \quad (19)$$

Asymptotic linearity and efficiency of estimators for modified treatment policies follows.

Theorem 3 (Weak convergence of one-step estimators). *Let $\|\cdot\|$ denote the $L_2(P)$ -norm defined as $\|f\|^2 = \int f^2 dP$. Define the following conditions.*

(C1) $P\{|D_{\eta,\delta}^j(O)| \leq C\} = P\{|D_{\hat{\eta},\delta}^j(O)| \leq C\} = 1$ for $j = 1, 2$ and for some $C < \infty$.

(C2) *The following second-order terms converge at the specified rate $\|\hat{m} - m\| \{\|\hat{g} - g\| + \|\hat{e} - e\| + \|\hat{d} - d\|\} = o_P(n^{-1/2})$, $\|\hat{g} - g\| \{\|\hat{u} - u\| + \|\hat{q} - q\|\} = o_P(n^{-1/2})$, $\|\hat{b} - b\| \{\|\hat{v} - v\| + \|\hat{s} - s\|\} = o_P(n^{-1/2})$.*

(C3) The effect is defined in terms of modified treatment policy $d(a, w)$, which is piecewise smooth invertible (A1).

(C4) The intervention g_δ is an exponential tilting intervention and $P \left\{ \int (\hat{g} - g) d\kappa \right\}^2 = o_P(n^{-1/2})$.

If Conditions (C1) and (C2) hold, and one of Conditions (C3) and (C4) holds, then: $\sqrt{n}\{\hat{\psi}_{D,\delta}^{os} - \psi_{D,\delta}\} \rightsquigarrow N(0, \sigma_{D,\delta}^2)$, and $\sqrt{n}\{\hat{\psi}_{I,\delta}^{os} - \psi_{I,\delta}\} \rightsquigarrow N(0, \sigma_{I,\delta}^2)$, where $\sigma_{D,\delta}^2 = \text{Var}\{D_{\eta,0}^1(O) - D_{\eta,\delta}^2(O)\}$ and $\sigma_{I,\delta}^2 = \text{Var}\{D_{\eta,\delta}^2(O) - D_{\eta,\delta}^1(O)\}$ are the respective efficiency bounds.

Theorem 3 establishes the weak convergence of $\hat{\psi}_{D,\delta}^{os}$ and $\hat{\psi}_{I,\delta}^{os}$ pointwise in δ . This convergence is useful to derive confidence intervals in situations where the MTP has a scientific interpretation for a given realization of δ . Under Theorem 3, an estimator $\hat{\sigma}_{D,\delta}^2$ of $\sigma_{D,\delta}^2$ may be obtained as the empirical variance of $D_{\hat{\eta}_{j(i)},0}^1(O_i) - D_{\hat{\eta}_{j(i)},\delta}^2(O_i)$, and a Wald-style confidence interval may be constructed as $\hat{\psi}_{D,\delta}^{os} \pm z_{1-\alpha/2} \hat{\sigma}_{D,\delta}^2 / \sqrt{n}$; the same applies to $\hat{\psi}_{I,\delta}^{os}$.

Although the one-step estimator has optimal asymptotic performance, its finite-sample behavior may be affected by the inverse probability weighting involved in the computation of the efficient influence functions $D_{\hat{\eta}}^j(O_i) : j = 1, 2$. In particular, it is not guaranteed that $\hat{\psi}_{D,\delta}^{os}$ and $\hat{\psi}_{I,\delta}^{os}$ will remain within the bounds of the parameter space. This issue may be attenuated by performing weight stabilization. The estimated EIF $D_{\hat{\eta}_{j(i)}}^1(O_i)$ can be weight-stabilized by dividing (9) and (11) by the empirical mean of $H_{\hat{\eta}_{j(i)},\delta}^1(A_i, Z_i, L_i, W_i)$ and $H_{\hat{\eta}_{j(i)},\delta}^2(A_i, Z_i, L_i, W_i)$, respectively; as well as dividing (10), (12), (15), and (16) by the empirical mean of $\hat{g}_{j(i),\delta}(A_i | W_i) / \hat{g}_{j(i)}(A_i | W_i)$.

4.2 Efficient Targeted Minimum Loss Estimator

Although corrections may be applied to the one-step estimator, a more principled way to obtain estimators that remain in the parameter space may be derived from the targeted minimum loss (TML) estimation framework. The TML estimator is constructed by tilting an initial data adaptive estimator $\hat{\eta}$ towards a solution $\tilde{\eta}$ of the estimating equations

$$P_n\{D_{\tilde{\eta},0}^1 - D_{\tilde{\eta},\delta}^2\} = \psi_{D,\delta}(\tilde{\eta}) \quad \text{and} \quad P_n\{D_{\tilde{\eta},\delta}^2 - D_{\tilde{\eta},\delta}^1\} = \psi_{I,\delta}(\tilde{\eta}), \quad (20)$$

where $\psi_{D,\delta}(\tilde{\eta})$ and $\psi_{I,\delta}(\tilde{\eta})$ are the substitution estimators in formula (21) obtained by plugging in the estimates $\tilde{\eta}$ in the parameter definition (7). Thus, a TML estimator is guaranteed to remain in the parameter space by virtue of its being a substitution estimator. The fact that the nuisance estimators solve the relevant estimating equation is used to obtain a weak convergence result analogous to Theorem 3. Thus, while the TML estimator is expected to attain the same optimal asymptotic behavior as the one-step estimator, its finite-sample behavior may be better. An algorithm to compute a TML estimator $\tilde{\eta}$ is presented in the Supplementary Materials. Roughly, the algorithm proceeds by projecting the EIF into score functions for the model of each nuisance parameter and fitting appropriate parametric submodels [van der Laan and Rose, 2011]. For example, the following model is fitted for m :

$$\text{logit } m_\beta(a, z, l, w) = \text{logit } \hat{m}(z, l, a, w) + \beta_I H_I(o) + \beta_D H_D(o), \quad \text{where}$$

$$H_D(o) = \frac{\hat{b}(l | a, w)}{\hat{d}(l | z, a, w)} \left\{ 1 - \frac{\hat{g}_\delta(a | w)}{\hat{e}(a | z, w)} \right\}, \quad \text{and} \quad H_I(o) = \frac{\hat{b}(l | a, w)}{\hat{d}(l | z, a, w)} \left\{ \frac{\hat{g}_\delta(a | w)}{\hat{e}(a | z, w)} - \frac{\hat{g}(a | w)}{\hat{g}(a | w)} \right\},$$

and $\text{logit}(p) = \log\{p(1-p)^{-1}\}$. Here, the initial estimator $\text{logit } \hat{m}(z, l, a, w)$ is considered a fixed offset variable (i.e., a variable with known parameter value equal to one). The score of these tilting models is equal to the corresponding component of the EIF. The parameter $\beta = (\beta_I, \beta_D)$ may be estimated via standard logistic regression of Y on $(H_D(O), H_I(O))$ with no intercept and an offset term equal to $\text{logit } \hat{m}(z, l, a, w)$. Let $\hat{\beta}$ denote the MLE, and let $\tilde{m} = m_{\hat{\beta}}$ denote the updated estimates. Fitting this regression model ensures that \tilde{m} solves the relevant score equations. Regression models like this are estimated iteratively for all parameters in a way that guarantees that the estimating equations (Equation (20)) are solved up to an error term that converges to zero in probability at rate faster than $n^{-1/2}$. Upon termination of the iterative process, the TML estimators are defined as

$$\begin{aligned} \hat{\psi}_{D,\delta}^{tmle} &= \frac{1}{n} \int \sum_{i=1}^n \{ \tilde{u}(a, W_i) \tilde{g}(a | W_i) - \tilde{u}(Z_i, a, W_i) \tilde{g}_\delta(a | W_i) \} d\kappa(a) \\ \hat{\psi}_{I,\delta}^{tmle} &= \frac{1}{n} \int \sum_{i=1}^n \{ \tilde{u}(Z_i, a, W_i) - \tilde{u}(a, W_i) \} \tilde{g}_\delta(a | W_i) d\kappa(a). \end{aligned} \quad (21)$$

The fact that the TML estimator solves estimating equations (Equation (20)) is fundamental to the following theorem.

Theorem 4 (Weak convergence of TML estimator). *Assuming that (C1) and (C2) hold, and one of (C3) and (C4) of Theorem 3 hold, then $\sqrt{n}\{\hat{\psi}_{D,\delta}^{tmle} - \psi_{D,\delta}\} \rightsquigarrow N(0, \sigma_{D,\delta}^2)$, $\sqrt{n}\{\hat{\psi}_{I,\delta}^{tmle} - \psi_{I,\delta}\} \rightsquigarrow N(0, \sigma_{I,\delta}^2)$, where $\sigma_{D,\delta}^2 = \text{Var}\{D_{\eta,0}^1(O) - D_{\eta,\delta}^2(O)\}$ and $\sigma_{I,\delta}^2 = \text{Var}\{D_{\eta,\delta}^2(O) - D_{\eta,\delta}^1(O)\}$.*

Using Theorem 4, asymptotically valid variance estimators, p-values, and confidence intervals for the (in)direct effects may be obtained in a manner analogous to those for the one-step estimator. The proof of the theorem proceeds using similar arguments as the proof of Theorem 3 for the one-step estimator, using empirical process theory and leveraging cross-fitting to avoid entropy conditions on the initial estimators of η . Since the estimators now depend on the full sample through the estimates of the parameters β of the logistic tilting models, the empirical process treatment differs slightly to that of Theorem 3; its proof is detailed in the Supplementary Materials.

In Section S1 of the Supplementary Materials, we present a simulation study comparing the two efficient estimators under different configurations of nuisance parameter misspecification. In brief, our findings illustrate that our estimators empirically satisfy the forms of robustness identified by our theoretical investigations; moreover, in keeping with prior investigations in other settings [e.g., van der Laan and Rose, 2011], the TML estimator generally outperforms the one-step estimator in terms of both bias and efficiency. Given the favorable evaluation of our proposed estimators in these experiments, we next demonstrate their application.

5 Application to the X:BOT trial

We now apply our stochastic interventional direct and indirect effects to decompose the causal effect of a strategy where buprenorphine dose is successively increased early in the treatment course (regardless of opioid use) on relapse among those with opioid use disorder (OUD). Data for our illustrative analysis come from the X:BOT trial, a 24-week, multi-site randomized controlled trial designed to examine the comparative effectiveness of extended-release naltrexone (XR-NTX) and sublingual buprenorphine-naloxone (BUP-NX) on relapse [Lee et al., 2018]. The X:BOT trial enrolled 570 participants, all of whom were 18 years or older, had OUD [as per the Diagnostic and Statistical Manual of Mental Disorders-5; American Psychiatric Association, 2013], and had used non-prescribed opioids in the 30 days preceding enrollment. Participants were randomized to receive either XR-NTX or BUP-NX using a stratified permuted block design; 287 of the 570 were randomized to receive BUP-NX. Prior analytic efforts have established a protective effect of BUP-NX administration (versus placebo) on OUD relapse [Mattick et al., 2014]. For each participant assigned to receive BUP-NX, the prescribed dose was based on both clinical indication [Lee et al., 2018] and clinician judgment. Some clinicians tended to hold dose constant over time (i.e., a static regimen), while others increased dose — either based on clinical assessment or on the hypothesis that higher doses would result in better outcomes [Comer et al., 2005]. We estimated stochastic interventional (in)direct effects to assess the mechanism by which universally ramping up BUP-NX dose early in treatment (i.e., three or more dose increases in the first four weeks of treatment) could mitigate the risk of OUD relapse.

Baseline covariates (W) available in the data included site; gender; age; race/ethnicity; homeless status; educational attainment; employment status; marital status; current intravenous drug use; alcohol use disorder; cocaine use disorder; age at start of heroin use; severity of current opioid use; indicator of prior OUD treatment; past withdrawal discomfort level; histories of amphetamine use, sedative use, and cannabis use; weekly cost of primary drug; whether or not living with an individual currently using drugs or with alcohol use disorder; histories of psychiatric illnesses; randomization timing; baseline pain level; baseline depression symptoms. The exposure (A) was taken to be successive increases in dose of BUP-NX versus static dose, measured during the first four weeks of treatment. Mediating factors (Z) included depression and pain, measured from week 6 until relapse or week 24 (end of follow-up). Abstinence from illicit opioid use early in the treatment schedule, measured between weeks 4 and 6, acted as an intermediate confounder affected by exposure (L). OUD relapse status at the X:BOT trial’s end of follow-up was the outcome of interest (Y). To examine the effect of exposure to successive increases in BUP-NX dose, we consider an incremental propensity score intervention, which, for binary A , replaces the propensity score $g(1 | w)$ with a shifted variant constructed from multiplying the odds of exposure by a user-specified degree parameter δ [Kennedy, 2019], which we vary along a grid $\log(\delta) \in \{-10.0, -9.5, \dots, 9.5, 10.0\}$ of the observed exposure odds. Across all such estimates in the odds δ of exposure, the stochastic interventional (in)direct effects that we estimated may be interpreted in terms of the overall effect of increasingly encouraging ramping up BUP-NX dose early in treatment on the counterfactual risk of OUD relapse; thus, the results of our analysis may be informative of the mechanisms by which increasing BUP-NX dose can alter the risk of OUD relapse. Figure 2 presents the direct and indirect effect estimates across the grid in δ . We applied both of our cross-fitted, efficient one-step and TML estimators to examine the stochastic interventional direct and indirect effects of increasing the odds of ramping up BUP-NX dose. Both strategies produced results generally in close agreement with the magnitude of the (in)direct effects. For each point estimate, standard error estimates and 95% Wald-style confidence intervals were constructed based on Theorem 4. To ensure flexibility of

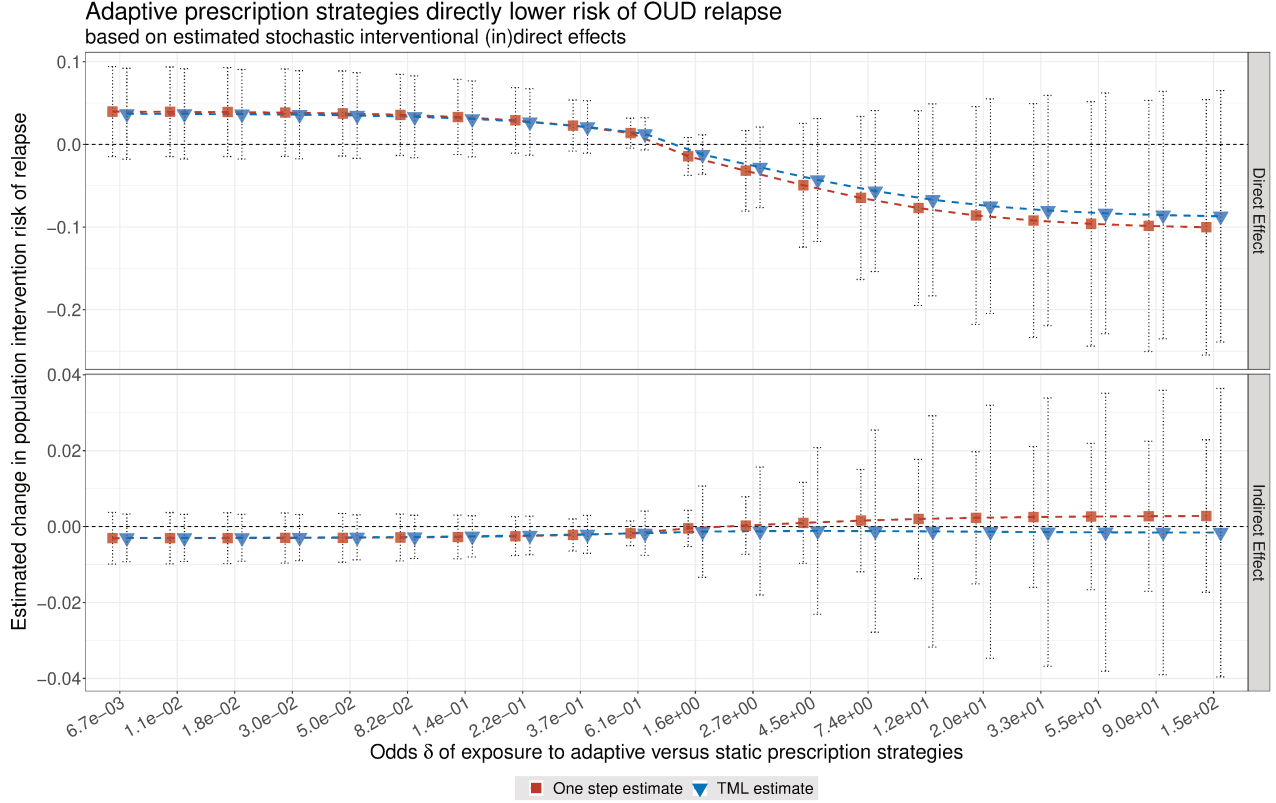


Figure 2: Stochastic interventional direct (upper panel) and indirect (lower panel) effect estimates of hypothetically increasing exposure odds to a BUP-NX dose schedule in which dose is increased early in OUD relapse treatment.

our estimators, all nuisance parameters $\eta = (e, m, d, g, b, u, v, s, q)$ were estimated using the Super Learner ensemble modeling algorithm [van der Laan et al., 2007, Coyle et al., 2021]. From examination of the point estimates and confidence intervals of the (in)direct effects in Figure 2, two conclusions may be drawn. Firstly, there appears to be little to no indirect effect of successively increasing BUP-NX dose on risk of OUD relapse, revealing that any effect of BUP-NX dose cannot be explained by actions on mediating factors such as pain or depression. Secondly, the direct effect of successively increasing BUP-NX dose varies considerably across changes in the odds of the introduction of such a dosing schedule. Importantly, lowered odds of dosage increases could lead to as much as a 5% increase in OUD relapse risk, with a plateau emerging at odds lower than $\approx 0.1\%$, suggesting that static dose can lead to comparatively heightened relapse risk. OUD relapse risk further appears to decrease by $\approx 10\%$ with increased odds of successive BUP-NX dose increases, with the risk plateauing at odds higher than 33%. This decrease in the counterfactual risk of OUD relapse suggests a protective effect of BUP-NX dose schedules when dose is successively increased early in the treatment course; however, the action mechanism is not readily apparent.

The conclusions that may be drawn from our re-analysis using the stochastic interventional (in)direct effects complement those previously reported in the investigations of Lee et al. [2018], who evaluated the total effect of BUP-NX (versus XR-NTX) treatment on OUD relapse, and Rudolph et al. [2020a], who used the interventional mediation analysis approach of Díaz et al. [2020] (for static interventions on A) to examine differences in relapse risk between homeless and non-homeless participants. Importantly, our substantive conclusion — that dosage increases directly lower the risk of relapse — agrees with those of Rudolph et al. [2020b], who found that dosage increases directly lowered risk of OUD relapse when such increases followed opioid use. Notably, our proposed (in)direct effects and estimation approach differ from prior efforts in three important ways: (i) our causal effect definitions remain unaltered in the presence of intermediate confounders affected by exposure and may be re-evaluated in randomized trials, (ii) the flexible estimators we introduce eschew restrictive modeling assumptions by incorporating modern machine learning in nuisance parameter estimation, and (iii) our strategy provides an analog to a dose-response analysis by allowing for the risk of OUD relapse to be traced out across changes in the odds of exposure to a schedule in which BUP-NX dose is increased repeatedly early in treatment.

6 Discussion

We have proposed a class of novel direct and indirect effect estimands for causal mediation analysis, as well as two efficient estimators of these effects in the nonparametric statistical model. Importantly, our proposed estimation framework allows for data adaptive estimation of nuisance parameters, while still preserving the benefits associated with similar classical techniques: our estimators are regular and asymptotically linear, provide unbiased point estimates, are multiply robust, allow the construction of asymptotically valid confidence intervals, and are capable of attaining the nonparametric efficiency bound. Notably, our (in)direct effects remain well-defined even in the presence of intermediate confounders affected by exposure. Further, any scientific conclusions drawn based upon our proposed (in)direct effects may be readily interrogated in trials that randomize both the exposure and mediators. Such flexible effect definitions and estimators appear necessary both to cope with the design complexity of modern epidemiological and biomedical studies and to take advantage of the ever-growing number of flexible, data adaptive regression techniques.

The challenge of leveraging data adaptive regression methodology to construct robust estimators that accommodate valid statistical inference is not a new one. It has been considered in great detail as early as the work of Pfanzagl and Wefelmeyer [1985] as well in numerous recent advances, most notably by van der Laan and Rose [2011, 2018] and Chernozhukov et al. [2018]; related work by these authors presents a wealth of extensions and applications. In the present work, we derive multiply robust, efficient estimators based on both the one-step and targeted minimum loss estimation frameworks. Following Klaassen [1987] and Zheng and van der Laan [2011], our estimators leverage cross-validation to avoid imposing possibly restrictive assumptions on nuisance function estimators. We demonstrated the properties of our estimators in simulation experiments that illustrated their ability to yield unbiased point estimates, attain the nonparametric efficiency bound, and build confidence intervals exhibiting coverage at the nominal rate across several nuisance parameter configurations — all within a context in which classical mediation effects are ill-defined. We demonstrated the application of our novel (in)direct effects in dissecting the mechanism by which increasing the odds of adopting a dosing schedule of universal successive increases in buprenorphine early in treatment affects OUD relapse [Lee et al., 2018, Rudolph et al., 2020a].

Several significant extensions and refinements are left for future consideration. Firstly, our proposed estimation strategy for the direct and indirect effects leverages re-parameterizations of factors of the likelihood in order to simplify the estimation of nuisance parameters. This approach works particularly well when either mediators or intermediate confounders are of modest dimension; however, improvements can be made to accommodate settings in which both mediators and intermediate confounders are of high dimensionality. When defining effects based upon stochastic interventions indexed by the user-specified parameter δ , an important consideration is choosing *a priori* a particular value of δ . One solution is to evaluate a set of causal effects indexed by a grid in δ . In such cases, aggregate effects (across δ) may be summarized via working marginal structural models [e.g., Hejazi et al., 2020] or the construction of uniform tests of the null hypothesis of no direct effect [e.g., Díaz and Hejazi, 2020]. Developments of these distinct summarization strategies would enrich the range of scientific problems to which these robust and flexible direct and indirect effects may be applied.

Supplementary Materials

The reader is referred to the on-line Supplementary Materials for technical appendices. R scripts used to conduct the simulation experiments and real-world data analysis have been made publicly available in a GitHub repository at https://github.com/nhejazi/pub_medshift_interv_biostats. While the estimation machinery is accessible from that GitHub repository, its integration into our open source `medshift` R package [Hejazi and Díaz, 2020] (<https://github.com/nhejazi/medshift>) is ongoing and will support wider long-term usage.

Acknowledgments

The authors thank John Rotrosen, Edward Nunes, and Marc Fishman for raising the research question in the Application and for helpful feedback. KER's time was supported by a grant from the National Institute on Drug Abuse (award no. R00-DA042127), MJvdL's time was supported by a grant from the National Institute of Allergy and Infectious Diseases (award no. R01-AI074345), and NSH's time was supported by a grant from the National Science Foundation (award no. DMS-2102840). The X:BOT trial was supported by the National Institute on Drug Abuse, Clinical Trials Network (award no.'s U10DA013046, UG1/U10DA013035, UG1/U10DA013034, U10DA013045, UG1/U10DA013720, UG1/U10DA013732, UG1/U10DA013714, UG1/U10DA015831, U10DA015833, HHSN271201200017C, and HHSN271201500065C).

References

- American Psychiatric Association. *Diagnostic and statistical manual of mental disorders (DSM-5®)*. American Psychiatric Pub, 2013.
- Chen Avin, Ilya Shpitser, and Judea Pearl. Identifiability of path-specific effects. In *IJCAI International Joint Conference on Artificial Intelligence*, pages 357–363, 2005.
- Reuben M Baron and David A Kenny. The moderator–mediator variable distinction in social psychological research: Conceptual, strategic, and statistical considerations. *Journal of Personality and Social Psychology*, 51(6):1173, 1986.
- David Benkeser and Jialu Ran. Nonparametric inference for interventional effects with multiple mediators. *Journal of Causal Inference*, 9(1):172–189, 2021. doi: 10.1515/jci-2020-0018.
- Peter J Bickel, Chris AJ Klaassen, YA’Acov Ritov, and Jon A Wellner. *Efficient and Adaptive Estimation for Semiparametric Models*. Johns Hopkins University Press, 1993.
- Leo Breiman. Stacked regressions. *Machine Learning*, 24(1):49–64, 1996.
- Victor Chernozhukov, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen, Whitney Newey, and James M Robins. Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal*, 21(1):C1–C68, 01 2018. doi: 10.1111/ectj.12097.
- Stephen R Cole and Miguel A Hernán. Fallibility in estimating direct effects. *International Journal of Epidemiology*, 31(1):163–165, 2002.
- Sandra D Comer, Ellen A Walker, and Eric D Collins. Buprenorphine/naloxone reduces the reinforcing and subjective effects of heroin in heroin-dependent volunteers. *Psychopharmacology*, 181(4):664–675, 2005.
- Jeremy R Coyle, Nima S Hejazi, Ivana Malenica, Rachael V Phillips, and Oleg Sofrygin. *sl3: Modern Pipelines for Machine Learning and Super Learning*, 2021. URL <https://doi.org/10.5281/zenodo.1342293>. R package version 1.4.4.
- A Philip Dawid. Causal inference without counterfactuals. *Journal of the American Statistical Association*, 95(450):407–424, 2000.
- Iván Díaz and Nima S Hejazi. Causal mediation analysis for stochastic interventions. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 82(3):661–683, 2020. doi: 10.1111/rssb.12362. URL <https://doi.org/10.1111/rssb.12362>.
- Iván Díaz and Mark J van der Laan. Super learner based conditional density estimation with application to marginal structural models. *International Journal of Biostatistics*, 7(1):1–20, 2011.
- Iván Díaz and Mark J van der Laan. Population intervention causal effects based on stochastic interventions. *Biometrics*, 68(2):541–549, 2012.
- Iván Díaz and Mark J van der Laan. Assessing the causal effect of policies: an example using stochastic interventions. *International Journal of Biostatistics*, 9(2):161–174, 2013.
- Iván Díaz, Nima S Hejazi, Kara E Rudolph, and Mark J van der Laan. Non-parametric efficient causal mediation with intermediate confounders. *Biometrika*, 108(3):627–641, 2020. doi: 10.1093/biomet/asaa085. URL <https://arxiv.org/abs/1912.09936>.
- Vanessa Didelez, Philip Dawid, and Sara Geneletti. Direct and indirect effects of sequential treatments. In *Proceedings of the 22nd Annual Conference on Uncertainty in Artificial Intelligence*, pages 138–146, 2006.
- Isabel R Fulcher, Ilya Shpitser, Stella Marealle, and Eric J Tchetgen Tchetgen. Robust inference on population indirect causal effects: the generalized front door criterion. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 82(1):199–214, 2019.

- Arthur S Goldberger. Structural equation methods in the social sciences. *Econometrica: Journal of the Econometric Society*, pages 979–1001, 1972.
- Sebastian Haneuse and Andrea Rotnitzky. Estimation of the effect of interventions that modify the received treatment. *Statistics in Medicine*, 32(30):5260–5277, 2013.
- Nima S Hejazi and Iván Díaz. *medshift: Causal mediation analysis for stochastic interventions*, 2020. URL <https://github.com/nhejazi/medshift>. R package version 0.1.4.
- Nima S Hejazi, Mark J van der Laan, Holly E Janes, Peter B Gilbert, and David C Benkeser. Efficient nonparametric inference on the effects of stochastic interventions under two-phase sampling, with applications to vaccine efficacy trials. *Biometrics*, 77(4):1241–1253, 2020. doi: 10.1111/biom.13375. URL <http://arxiv.org/abs/2003.13771>.
- Nima S Hejazi, David C Benkeser, and Mark J van der Laan. *haldensify: Highly adaptive lasso conditional density estimation*, 2021. URL <https://github.com/nhejazi/haldensify>. R package version 0.2.2.
- Alan E Hubbard and Mark J van der Laan. Population intervention models in causal inference. *Biometrika*, 95(1): 35–47, 2008.
- Kosuke Imai, Luke Keele, and Dustin Tingley. A general approach to causal mediation analysis. *Psychological Methods*, 15(4):309, 2010.
- Edward H Kennedy. Nonparametric causal effects based on incremental propensity score interventions. *Journal of the American Statistical Association*, 114(526):645–656, 2019.
- Chris AJ Klaassen. Consistent estimation of the influence function of locally asymptotically linear estimators. *Annals of Statistics*, pages 1548–1562, 1987.
- Joshua D Lee, Edward V Nunes Jr, Patricia Novo, Ken Bachrach, Genie L Bailey, Snehal Bhatt, Sarah Farkas, Marc Fishman, Phoebe Gauthier, Candace C Hodgkins, Jacquie King, Robert Lindblad, David Liu, Abigail G Matthews, Jeanine May, K Michelle Peavy, Stephen Ross, Dagmar Salazar, Paul Schkolnik, Dikla Shmueli-Blumberg, Don Stablein, Geetha Subramaniam, and John Rotrosen. Comparative effectiveness of extended-release naltrexone versus buprenorphine-naloxone for opioid relapse prevention (X:BOT): a multicentre, open-label, randomised controlled trial. *The Lancet*, 391(10118):309–318, 2018.
- Richard P Mattick, Courtney Breen, Jo Kimber, and Marina Davoli. Buprenorphine maintenance vs. placebo or methadone maintenance for opioid dependence. *Cochrane Database of Systematic Reviews*, (2), 2014.
- Trang Quynh Nguyen, Ian Schmid, and Elizabeth A Stuart. Clarifying causal mediation analysis for the applied researcher: Defining effects based on what we want to learn. *Psychological Methods*, 26(2), 2021.
- Judea Pearl. Causal diagrams for empirical research. *Biometrika*, 82(4):669–688, 1995.
- Judea Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, 2000.
- Judea Pearl. Direct and indirect effects. In *Proceedings of the 17th Annual Conference on Uncertainty in Artificial Intelligence*, 2006.
- J Pfanzagl and W Wefelmeyer. Contributions to a general asymptotic statistical theory. *Statistics & Risk Modeling*, 3 (3-4):379–388, 1985.
- Karl Popper. *The Logic of Scientific Discovery*. Routledge, 1934.
- R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2022. URL <https://www.R-project.org/>.
- James M Robins. A new approach to causal inference in mortality studies with sustained exposure periods — application to control of the healthy worker survivor effect. *Mathematical Modelling*, 7:1393–1512, 1986.

- James M Robins and Sander Greenland. Identifiability and exchangeability for direct and indirect effects. *Epidemiology*, pages 143–155, 1992.
- Kara E Rudolph, Oleg Sofrygin, Wenjing Zheng, and Mark J van der Laan. Robust and flexible estimation of stochastic mediation effects: a proposed method and example in a randomized trial setting. *Epidemiologic Methods*, 7(1), 2017.
- Kara E Rudolph, Iván Díaz, Nima S Hejazi, Mark J van der Laan, Sean X Luo, Matisyahu Shulman, Aimee Campbell, John Rotrosen, and Edward V Nunes. Explaining differential effects on opioid use disorder treatment using a novel causal approach incorporating mediating and intermediate variables. *Addiction*, 116(8):2094–2103, 2020a. doi: 10.1111/add.15377.
- Kara E Rudolph, Matisyahu Shulman, Marc Fishman, Iván Díaz, John Rotrosen, and Edward V Nunes. Association between dynamic dose adjustment of buprenorphine for treatment of opioid use disorder and risk of relapse. 2020b.
- Peter Spirtes, Clark N Glymour, Richard Scheines, David Heckerman, Christopher Meek, Gregory Cooper, and Thomas Richardson. *Causation, Prediction, and Search*. MIT Press, 2000.
- James H Stock. Nonparametric policy analysis. *Journal of the American Statistical Association*, 84(406):567–575, 1989.
- Eric J Tchetgen Tchetgen and Tyler J VanderWeele. On identification of natural direct effects when a confounder of the mediator is directly affected by exposure. *Epidemiology*, 25(2):282, 2014.
- Mark J van der Laan and Maya L Petersen. Direct effect models. *International Journal of Biostatistics*, 4(1), 2008.
- Mark J van der Laan and James M Robins. *Unified Methods for Censored Longitudinal Data and Causality*. Springer Science & Business Media, 2003.
- Mark J van der Laan and Sherri Rose. *Targeted Learning: Causal Inference for Observational and Experimental Data*. Springer Science & Business Media, 2011.
- Mark J van der Laan and Sherri Rose. *Targeted Learning in Data Science: Causal Inference for Complex Longitudinal Studies*. Springer Science & Business Media, 2018.
- Mark J van der Laan and Daniel Rubin. Targeted maximum likelihood learning. *International Journal of Biostatistics*, 2(1), 2006.
- Mark J van der Laan, Eric C Polley, and Alan E Hubbard. Super Learner. *Statistical Applications in Genetics and Molecular Biology*, 6(1), 2007.
- Tyler J VanderWeele, Stijn Vansteelandt, and James M Robins. Effect decomposition in the presence of an exposure-induced mediator-outcome confounder. *Epidemiology*, 25(2):300, 2014.
- Stijn Vansteelandt and Rhian M Daniel. Interventional effects for mediation analysis with multiple mediators. *Epidemiology*, 28(2):258, 2017.
- Stijn Vansteelandt and Tyler J VanderWeele. Natural direct and indirect effects on the exposed: effect decomposition under weaker assumptions. *Biometrics*, 68(4):1019–1027, 2012.
- Sewall Wright. The method of path coefficients. *The Annals of Mathematical Statistics*, 5(3):161–215, 1934.
- Jessica G Young, Miguel A Hernán, and James M Robins. Identification, estimation and approximation of risk under interventions that depend on the natural value of treatment using observational data. *Epidemiologic Methods*, 3(1): 1–19, 2014.
- Wenjing Zheng and Mark J van der Laan. Cross-validated targeted minimum-loss-based estimation. In *Targeted Learning: Causal Inference for Observational and Experimental Data*, pages 459–474. Springer, 2011.

Supplementary Materials for NONPARAMETRIC CAUSAL MEDIATION ANALYSIS FOR STOCHASTIC INTERVENTIONAL (IN)DIRECT EFFECTS

Nima S. Hejazi

Division of Biostatistics,
Department of Population Health Sciences,
Weill Cornell Medicine
nhejazi@berkeley.edu

Kara E. Rudolph

Department of Epidemiology,
Mailman School of Public Health,
Columbia University
kr2854@cumc.columbia.edu

Mark J. van der Laan

Division of Biostatistics,
School of Public Health, and
Department of Statistics,
University of California, Berkeley
laan@berkeley.edu

Iván Díaz

Division of Biostatistics,
Department of Population Health Sciences,
Weill Cornell Medicine
ild2005@med.cornell.edu

January 13, 2022

S1 Simulation study

We used simulation experiments to assess our two proposed efficient estimators of the (in)direct effects. On account of computational considerations, we focus on binary exposures and intermediate confounders in this example; however, as noted in the prior, our proposed methodology is general enough to be readily applicable in the presence of continuous-valued covariates, treatment, mediators, intermediate confounders, and outcome. We used the following data-generating mechanism for the joint distribution of O to generate synthetic data to evaluate our two estimators:

$$\begin{aligned}
 W_1 &\sim \text{Bernoulli}(p = 0.6); \\
 W_2 &\sim \text{Bernoulli}(p = 0.3); \\
 W_3 \mid (W_1, W_2) &\sim \text{Bernoulli}(p = 0.2 + 1/3 \cdot (W_1 + W_2)); \\
 A \mid W &\sim \text{Bernoulli}(p = \text{expit}(2 + (5/(W_1 + W_2 + W_3)))); \\
 L \mid (A, W) &\sim \text{Bernoulli}(p = \text{expit}(1/3(W_1 + W_2 + W_3) - A - \log(2) + 0.2)); \\
 Z \mid (L, A, W) &\sim \text{Bernoulli}(p = \text{expit}(\log(3) \cdot (W_1 + W_2) + A - L)); \\
 Y \mid (Z, L, A, W) &\sim \text{Bernoulli}\left(p = \text{expit}\left(1 - \frac{3 \cdot (3 - L - 3A + Z)}{2 + (W_1 + W_2 + W_3)}\right)\right),
 \end{aligned}$$

where $\text{expit}(x) := \{1 + \exp(x)\}^{-1}$. For each of the sample sizes $n \in \{200, 800, 1800, 3200, 5000, 7200, 9800, 12800, 16200\}$, 500 datasets were generated. For every dataset, six variations of each of the two efficient estimators was applied — five variants were based on misspecification of a single nuisance parameter among $\{e, m, d, g, b\}$ while the sixth variant was constructed based on consistent estimation of all five nuisance parameters. An intercept-only logistic regression model provided inconsistent estimation of each of the nuisance parameters $\{e, m, d, g, b\}$, while a Super Learner ensemble [van der Laan et al., 2007] was used to achieve consistent estimation. The Super Learner ensemble was constructed with a library of algorithms composed of intercept-only logistic regression; main-terms logistic regression; and several variants of the highly adaptive lasso [Benkeser and van der Laan, 2016, van der Laan, 2017, Coyle et al., 2020], a nonparametric regression approach capable of flexibly estimating arbitrary func-

tional forms at a fast convergence rate under only a global smoothness assumption [van der Laan and Bibaut, 2017, Bibaut and van der Laan, 2019]. Note that we do not consider cases of misspecified estimation of $\{v, s, q, \bar{u}\}$, as their consistent estimation depends on a subset of the nuisance parameters $\{e, m, d, g, b\}$. Generally, based on Lemmas 5 and 6, robustness of the direct and indirect effect estimators to misspecification of $\{e, m, d\}$ is to be expected, but the same is not true under misspecification of $\{g, b\}$.

Figure S1 summarizes the results of our investigations of the relative performance of the estimator variants enumerated above. Specifically, we assess the relative performance of our proposed estimators in terms of absolute bias, scaled (by $n^{1/2}$) bias, standard error and scaled (by n) mean squared error relative to the efficiency bound for the data-generating model, the empirical coverage of 95% confidence intervals, and relative efficiency. In terms of both raw (unscaled) bias and scaled bias, the estimator variants appear to conform to the predictions of Lemmas 5 and 6 — specifically, raw bias vanishes and scaled bias stabilizes to a small value (providing evidence for rate-consistency) under misspecification of any of $\{e, m, d\}$ as well as in the case of no nuisance parameter misspecification. In the same vein, when either of $\{g, b\}$ are estimated inconsistently, some of the estimator variants display diverging asymptotic (scaled) bias, in agreement with expectations based upon theory. The consistency of other estimator variants (e.g., the one-step estimator under misspecification of b) is likely an artifact of this data-generating mechanism, not to be taken as a general indication of robust performance. In terms of their relative mean squared error, the estimators of the (in)direct effects exhibit convergence to the efficiency bound under misspecification of $\{e, m, d\}$ and under no misspecification; this also appears to hold for a subset of the estimator variants under misspecification of $\{g, b\}$. We stress that aspects of this are likely to be a particularity of the given data-generating mechanism or on account of the irregularity of misspecified estimator variants, for the regularity and asymptotically linearity of the estimators is only to be expected under consistent estimation of all nuisance parameters. Finally, the empirical coverage of 95% confidence intervals is as expected: under a lack of nuisance parameter misspecification, both the one-step and TML estimators of the direct and indirect effect achieve 95% coverage in larger sample sizes. We note that misspecification of e leads to over-coverage for all estimator variants, implying an overly inflated variance estimate, while the confidence intervals fail to attain the nominal rate in most other instances. Notably, several of the estimator variants generate confidence intervals that are liable to converge to 0% coverage in larger samples under misspecification of $\{g, b\}$, very much in line with theoretical expectations.

Importantly, the TML estimator appears to generally outperform the one-step estimator throughout several scenarios. This comes in several forms, including lower bias, relative standard deviation, or relative mean squared error under misspecification of $\{e, m, d\}$ or under no misspecification; however, under inconsistent estimation of $\{g, b\}$, the irregularity of the estimators complicates this comparison. Interestingly, under misspecification of g , the TML estimators of the direct and indirect effects appear unbiased and efficient, a result unpredictable from theory given the irregularity of the estimators under this configuration. Altogether, results of our numerical experiments indicate that our proposed estimators exhibit properties that align with the theoretical results of Lemmas 5 and 6.

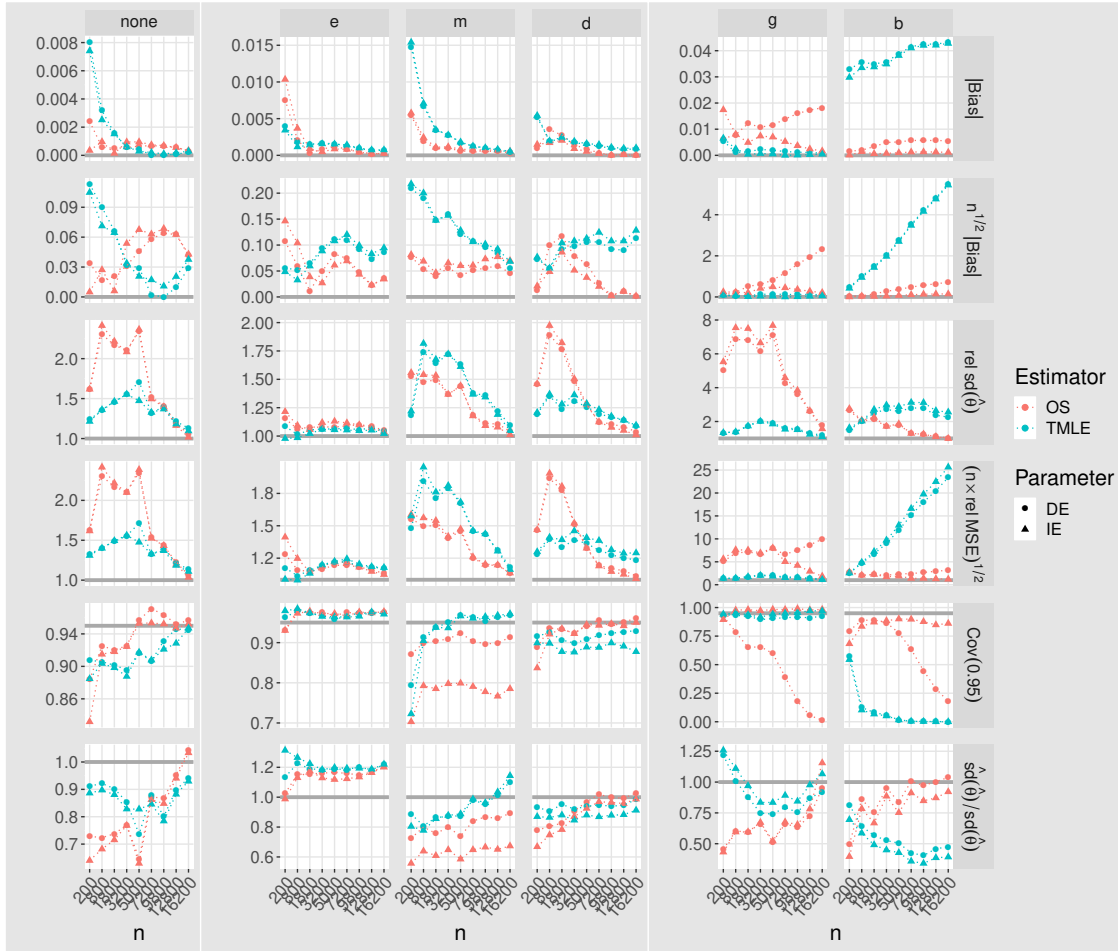


Figure S1: Comparison of efficient estimators across different nuisance parameter configurations.

In the interest of transparency and computational reproducibility, R scripts used to conduct these simulation experiments and summarize their results have been made publicly available in a curated GitHub repository located at https://github.com/nhejazi/pub_medshift_interv_biostats. The estimation machinery currently accessible via that repository is in the process of being integrated into the open source `medshift` R package, available at <https://github.com/nhejazi/medshift> [Hejazi and Díaz, 2020].

S2 Theorem 1

Proof First, we have

$$\begin{aligned} & \mathbb{E}\{Y_{A_\delta, G_\delta}\} \\ &= \int \mathbb{E}\{Y_{a,z} \mid A_\delta = a, G_\delta = z, W = w\} g_\delta(a \mid w) P(G_\delta = z \mid A_\delta = a, W = w) p(w) d\nu(a, z, w) \\ &= \int \mathbb{E}\{Y_{a,z} \mid W = w\} g_\delta(a \mid w) P(Z(a) = z \mid A_\delta = a, W = w) p(w) d\nu(a, z, w) \end{aligned} \quad (\text{S1})$$

$$= \int \mathbb{E}\{Y_{a,z} \mid A = a, W = w\} g_\delta(a \mid w) P(Z(a) = z \mid W = w) p(w) d\nu(a, z, w) \quad (\text{S2})$$

$$= \int \mathbb{E}\{Y_{a,z} \mid A = a, W = w\} g_\delta(a \mid w) P(Z(a) = z \mid A = a, W = w) p(w) d\nu(a, z, w) \quad (\text{S3})$$

$$\begin{aligned} &= \int \mathbb{E}\{Y_{a,z} \mid A = a, W = w, L = l\} b(l \mid a, w) g_\delta(a \mid w) p(z \mid a, w) p(w) d\nu(a, z, l, w) \\ &= \int m(a, z, l, w) b(l \mid a, w) g_\delta(a \mid w) p(z \mid a, w) p(w) d\nu(a, z, l, w), \end{aligned} \quad (\text{S4})$$

where (S1) follows by definition of (A_δ, G_δ) , (S2) follows by A4 and definition of A_δ , (S3) follows by A6, and (S4) follows by A5; the rearrangements made between (S3) and (S4) make use of A3. Similar arguments yield

$$\mathbb{E}\{Y_{A,G}\} = \int m(a, z, l, w) b(l \mid a, w) g(z \mid w) p(z \mid a, w) p(w) d\nu(a, z, l, w).$$

We also have

$$\begin{aligned} & \mathbb{E}\{Y_{A_\delta, G}\} \\ &= \int \mathbb{E}\{Y_{a,z} \mid A_\delta = a, G = z, W = w\} g_\delta(a \mid w) P(G = z \mid A_\delta = a, W = w) p(w) d\nu(a, z, w) \\ &= \int \mathbb{E}\{Y_{a,z} \mid W = w\} g_\delta(a \mid w) P(G = z \mid W = w) p(w) d\nu(a, z, w) \\ &= \int \mathbb{E}\{Y_{a,z} \mid A = a, W = w\} g_\delta(a \mid w) p(z \mid w) p(w) d\nu(a, z, w) \\ &= \int \mathbb{E}\{Y_{a,z} \mid A = a, W = w\} g_\delta(a \mid w) p(z \mid w) p(w) d\nu(a, z, w) \\ &= \int \mathbb{E}\{Y_{a,z} \mid A = a, W = w, L = l\} b(l \mid a, w) g_\delta(a \mid w) p(z \mid w) p(w) d\nu(a, z, l, w) \\ &= \int m(a, z, l, w) b(l \mid a, w) g_\delta(a \mid w) p(z \mid w) p(w) d\nu(a, z, l, w). \end{aligned}$$

Subtracting gives the expressions for the PIIE and PIDE in the theorem. \square

S3 Efficient influence functions (Theorem 2)

Proof In this proof we will use $\Theta_j(P) : j = 1, 2$ to denote a parameter as a functional that maps the distribution P in the model to a real number. We will assume that the measure ν is discrete so that integrals can be written as sums, and will omit the dependence on δ . It can be checked algebraically that the resulting influence function will also correspond to the influence function of a general measure ν . The true parameter value for θ_1 is thus given by

$$\theta_1 = \Theta_1(P) = \sum_{y,a,z,m,w} y p(y \mid a, z, l, w) p(l \mid a, w) p(z \mid a, w) g_\delta(a \mid w) p(w).$$

The non-parametric MLE of θ_1 in the model of g_δ known is given by

$$\Theta(P_n) = \sum_{y,a,z,m,w} y \frac{P_n f_{y,a,z,l,w}}{P_n f_{a,z,l,w}} \frac{P_n f_{l,a,w}}{P_n f_{a,w}} \frac{P_n f_{z,a,w}}{P_n f_{a,w}} g_\delta(a \mid w) P_n f_w, \quad (\text{S5})$$

where we remind the reader of the notation $Pf = \int f dP$. Here $f_{y,a,z,l,w} = \mathbb{1}(Y = y, A = a, Z = z, M = m, W = w)$, and $\mathbb{1}(\cdot)$ denotes the indicator function. The other functions f are defined analogously.

We will use the fact that the efficient influence function in a non-parametric model corresponds with the influence curve of the NPMLE. This is true because the influence curve of any regular estimator is also a gradient, and a non-parametric model has only one gradient. The Delta method [see, e.g., Appendix 18 of van der Laan and Rose, 2011] shows that if $\hat{\Theta}_1(P_n)$ is a substitution estimator such that $\theta_1 = \hat{\Theta}_1(P)$, and $\hat{\Theta}_1(P_n)$ can be written as $\hat{\Theta}_1^*(P_n f : f \in \mathcal{F})$ for some class of functions \mathcal{F} and some mapping Θ_1^* , the influence function of $\hat{\Theta}_1(P_n)$ is equal to

$$\text{IF}_P(O) = \sum_{f \in \mathcal{F}} \frac{d\hat{\Theta}_1^*(P)}{dP} \{f(O) - Pf\}.$$

Applying this result to (S5) with $\mathcal{F} = \{f_{y,a,z,l,w}, f_{a,z,l,w}, f_{z,a,w}, f_{a',w}, f_{l,a,w}, f_{a,w}, f_w : y, a, z, l, w\}$ and rearranging terms gives the result of the theorem. The algebraic derivations involved here are lengthy and not particularly illuminating, and are therefore omitted from the proof. Similar analyses may be performed for the model where only g_δ is unknown, as well as θ_2 . \square

S4 Targeted minimum loss estimation algorithm

To simplify notation, in the remaining of this section we will denote $\tilde{\eta}_{j(i)}(O_i)$ with $\tilde{\eta}(O_i)$. If L is binary, the efficient influence functions in Theorem 2 may be simplified using the following identity:

$$v(l, a, w) - \bar{v}(a, w) = \{v(1, a, w) - v(0, a, w)\} \{l - b(1 | a, w)\},$$

which also holds for v replaced by s and \bar{v} by \bar{s} .

Step 1. Initialize $\tilde{\eta} = \hat{\eta}$. Compute \tilde{v} , \tilde{s} , and \tilde{q}^j by plugging in \tilde{m} , \tilde{g} , \tilde{e} , \tilde{d} into equations (8), (18) and (17) if Z is multivariate, and fitting data-adaptive regression algorithms as appropriate.

Step 2. For each subject, compute the auxiliary covariates

$$\begin{aligned} H_{D,i} &= \frac{\tilde{b}(L_i | A_i, W_i)}{\tilde{d}(L_i | Z_i, A_i, W_i)} \left\{ 1 - \frac{\tilde{g}_\delta(A_i | W_i)}{\tilde{e}(A_i | Z_i, W_i)} \right\} \\ H_{I,i} &= \frac{\tilde{b}(L_i | A_i, W_i)}{\tilde{d}(L_i | Z_i, A_i, W_i)} \left\{ \frac{\tilde{g}_\delta(A_i | W_i)}{\tilde{e}(A_i | Z_i, W_i)} - \frac{\tilde{g}(A_i | W_i)}{\tilde{g}(A_i | W_i)} \right\} \\ K_{D,i} &= \tilde{v}(1, A_i, W_i) - \tilde{v}(0, A_i, W_i) - \frac{\tilde{g}_\delta(A_i | W_i)}{\tilde{g}(A_i | W_i)} \{\tilde{s}(1, A_i, W_i) - \tilde{s}(0, A_i, W_i)\} \\ K_{I,i} &= \frac{\tilde{g}_\delta(A_i | W_i)}{\tilde{g}(A_i | W_i)} \{\tilde{s}(1, A_i, W_i) - \tilde{s}(0, A_i, W_i) - \tilde{v}(1, A_i, W_i) + \tilde{v}(0, A_i, W_i)\} \\ M_{D,i} &= -\frac{\tilde{g}_\delta(1 | w)(1 - \tilde{g}_\delta(1 | w))}{\tilde{g}(1 | w)(1 - \tilde{g}(1 | w))} \tilde{q}^2(w) \\ M_{I,i} &= \frac{\tilde{g}_\delta(1 | w)(1 - \tilde{g}_\delta(1 | w))}{\tilde{g}(1 | w)(1 - \tilde{g}(1 | w))} \{\tilde{q}^2(w) - \tilde{q}^1(w)\} \end{aligned}$$

Step 3. Fit the logistic tilting models

$$\begin{aligned} \text{logit } m_\beta(A_i, Z_i, L_i, W_i) &= \text{logit } \tilde{m}(A_i, Z_i, L_i, W_i) + \beta_I H_{I,i} + \beta_D H_{D,i} \\ \text{logit } b_\alpha(1 | A_i, W_i) &= \text{logit } \tilde{b}(1 | A_i, W_i) + \alpha_I K_{I,i} + \alpha_D K_{D,i} \\ \text{logit } g_\gamma(1 | W_i) &= \text{logit } \tilde{g}(1 | W_i) + \gamma_I M_{I,i} + \gamma_D M_{D,i} \end{aligned}$$

where $\text{logit}(p) = \log\{p(1-p)^{-1}\}$. Here, $\text{logit } \tilde{m}(a, z, l, w)$ is an offset variable (i.e., a variable with known parameter value equal to one). The parameter $\beta = (\beta_I, \beta_D)$ may be estimated by running standard logistic regression of Y_i on $(H_{D,i}, H_{I,i})$ with no intercept and an offset term equal to $\text{logit } \tilde{m}(A_i, Z_i, L_i, W_i)$. Let $\hat{\beta}$ denote the estimate, and let $\tilde{m} = m_{\hat{\beta}}$ denote the updated estimates. Perform analogous computations for b and g .

Step 4. Compute \tilde{u} according to equation (8) by plugging in \tilde{m} and \tilde{b} . Compute the covariate

$$J_i = \frac{\tilde{g}_\delta(A_i | W_i)}{\tilde{g}(A_i | W_i)},$$

and fit the model

$$\text{logit } \bar{u}_\kappa(A_i, W_i) = \text{logit } \tilde{u}(A_i, W_i) + \kappa_D + \kappa_I J_i$$

by running a logistic regression of $\tilde{u}(Z_i, A_i, W_i)$ on J_i with an intercept and offset $\text{logit } \tilde{u}(A_i, W_i)$. Let $\hat{\kappa}$ denote the MLE, and update $\tilde{u} = \bar{u}_{\hat{\kappa}}$.

Step 5. The TMLE of the direct and indirect effects are defined as:

$$\begin{aligned}\hat{\psi}_{D,\delta}^{tmle} &= \frac{1}{n} \int \sum_{i=1}^n \{ \tilde{u}(a, W_i) \tilde{g}(a | W_i) - \tilde{u}(Z_i, a, W_i) \tilde{g}_\delta(a | W_i) \} d\kappa(a) \\ \hat{\psi}_{I,\delta}^{tmle} &= \frac{1}{n} \int \sum_{i=1}^n \{ \tilde{u}(Z_i, a, W_i) - \tilde{u}(a, W_i) \} \tilde{g}_\delta(a | W_i) d\kappa(a)\end{aligned}$$

S5 Proof of Theorem 3

Proof Let $P_{n,j}$ denote the empirical distribution of the prediction set \mathcal{V}_j , and let $G_{n,j}$ denote the associated empirical process $\sqrt{n/J}(P_{n,j} - P)$. For simplicity we denote a general parameter ψ with influence function D_η , the proof applies equally to the direct and indirect effect parameters. Note that

$$\hat{\psi}_\delta^{os} = \frac{1}{J} \sum_{j=1}^J P_{n,j} D_{\hat{\eta}_j, \delta}, \quad \psi_\delta = P D_\eta.$$

Thus,

$$\sqrt{n} \{ \hat{\psi}_\delta^{os} - \psi_\delta \} = G_n \{ D_{\eta, \delta} - \psi_\delta \} + R_{n,1}(\delta) + R_{n,2}(\delta),$$

where

$$R_{n,1}(\delta) = \frac{1}{\sqrt{J}} \sum_{j=1}^J G_{n,j} (D_{\hat{\eta}_j, \delta} - D_{\eta, \delta}), \quad R_{n,2}(\delta) = \frac{\sqrt{n}}{J} \sum_{j=1}^J P \{ D_{\hat{\eta}_j, \delta} - \psi_\delta \}.$$

It remains to show that $R_{n,1}(\delta)$ and $R_{n,2}(\delta)$ are $o_P(1)$. Lemmas 5 and 6 together with the Cauchy-Schwartz inequality and assumption (C2) of the theorem shows that $\|R_{n,2}\|_\Delta = o_P(1)$. For $\|R_{n,1}\|_\Delta$ we use empirical process theory to argue conditional on the training sample \mathcal{T}_j . In particular, Lemma 19.33 of van der Vaart [2000] applied to the class of functions $\mathcal{F} = \{D_{\hat{\eta}_j, \delta} - D_{\eta, \delta}\}$ (which consists of one element) yields

$$E \left\{ \left| G_{n,j} (D_{\hat{\eta}_j, \delta} - D_{\eta, \delta}) \right| \middle| \mathcal{T}_j \right\} \lesssim \frac{2C \log 2}{n^{1/2}} + \|D_{\hat{\eta}_j, \delta} - D_{\eta, \delta}\| (\log 2)^{1/2}$$

By assumption (C2), the left hand side is $o_P(1)$. Lemma 6.1 of Chernozhukov et al. [2018] may now be used to argue that conditional convergence implies unconditional convergence, concluding the proof. \square

S6 Theorem 4

Proof Let $P_{n,j}$ denote the empirical distribution of the prediction set \mathcal{V}_j , and let $G_{n,j}$ denote the associated empirical process $\sqrt{n/J}(P_{n,j} - P)$. For simplicity we denote a general parameter ψ with influence function D_η , the proof applies equally to the direct and indirect effect parameters. By definition, the sum of the scores of the sub-models $\{m_\beta, b_\alpha, g_\gamma, \bar{u}_\kappa : (\beta, \alpha, \gamma, \kappa)\}$ at the last iteration of the TMLE procedure is equal to $n^{-1} \sum_{i=1}^n D_{\hat{\eta}}(O_i) = o_P(n^{-1/2})$. Thus, we have

$$\hat{\psi}_\delta^{tmle} = \frac{1}{J} \sum_{j=1}^J P_{n,j} D_{\hat{\eta}_j} + o_P(n^{-1/2}).$$

Thus,

$$\sqrt{n} (\hat{\psi}_\delta^{tmle} - \theta) = G_n (D_\eta - \theta) + R_{n,1} + R_{n,2} + o_P(n^{-1/2}),$$

where

$$R_{n,1} = \frac{1}{\sqrt{J}} \sum_{j=1}^J G_{n,j}(D_{\hat{\eta}_j} - D_{\eta}), \quad R_{n,2} = \frac{\sqrt{n}}{J} \sum_{j=1}^J P(D_{\hat{\eta}_j} - \theta).$$

As in the proof of Theorem 3, Lemmas 5 and 6 together with the Cauchy-Schwartz inequality and the assumptions of the theorem shows that $R_{n,2} = o_P(1)$.

Since $D_{\hat{\eta}_j}$ depends on the full sample through the estimates of the parameters β of the logistic tilting models, the empirical process treatment of $R_{n,1}$ needs to be slightly from that in the proof of Theorem 3. To make this dependence explicit, we introduce the notation $D_{\hat{\eta}_j, \beta} = D_{\hat{\eta}_j}$ and $R_{n,1}(\beta)$. Let $\mathcal{F}_n^j = \{D_{\hat{\eta}_j, \beta} - D_{\eta} : \beta \in B\}$. Because the function $\hat{\eta}_j$ is fixed given the training data, we can apply Theorem 2.14.2 of van der Vaart and Wellner [1996] to obtain

$$E \left\{ \sup_{f \in \mathcal{F}_n^j} |G_{n,j} f| \mid \mathcal{T}_j \right\} \lesssim \|F_n^j\| \int_0^1 \sqrt{1 + N_{[]}(\epsilon \|F_n^j\|, \mathcal{F}_n^j, L_2(P))} d\epsilon,$$

where $N_{[]}(\epsilon \|F_n^j\|, \mathcal{F}_n^j, L_2(P))$ is the bracketing number and we take $F_n^j = \sup_{\beta \in B} |D_{\hat{\eta}_j, \beta} - D_{\eta}|$ as an envelope for the class \mathcal{F}_n^j . Theorem 2.7.2 of van der Vaart and Wellner [1996] shows

$$\log N_{[]}(\epsilon \|F_n^j\|, \mathcal{F}_n^j, L_2(P)) \lesssim \frac{1}{\epsilon \|F_n^j\|}.$$

This shows

$$\begin{aligned} \|F_n^j\| \int_0^1 \sqrt{1 + N_{[]}(\epsilon \|F_n^j\|, \mathcal{F}_n^j, L_2(P))} d\epsilon &\lesssim \int_0^1 \sqrt{\|F_n^j\|^2 + \frac{\|F_n^j\|}{\epsilon}} d\epsilon \\ &\leq \|F_n^j\| + \|F_n^j\|^{1/2} \int_0^1 \frac{1}{\epsilon^{1/2}} d\epsilon \\ &\leq \|F_n^j\| + 2\|F_n^j\|^{1/2}. \end{aligned}$$

Since $\|F_n^j\| = o_P(1)$, this shows $\sup_{f \in \mathcal{F}_n^j} G_{n,j} f = o_P(1)$ for each j , conditional on \mathcal{T}_j . Thus $\sup_{\beta \in B} R_{n,1}(\beta) = o_P(1)$. Lemmas 5 and 6 together with the Cauchy-Schwartz inequality and the assumptions of the theorem show that $R_{n,2} = o_P(1)$, concluding the proof of the theorem. \square

S7 Additional results

Lemma S1 (Second order terms for modified treatment policies). *Let $d\xi(o)$ denote $d\nu(a, l, z)dP(w)$, and let $r(z \mid a, w)$ denote $p(z \mid a, w)$, and let $h(z \mid w)$ denote $p(z \mid w)$. Let $d(a, w)$ denote a modified treatment policy satisfying A1. We have*

$$PD_{\eta_1, \delta}^1 - \psi_1(\delta) = \int \left(\frac{g}{g_1} \frac{d}{d_1} - 1 \right) (m - m_1) b_1 r g_{\delta, 1} d\xi \quad (S6)$$

$$- \int \left(\frac{g}{g_1} - 1 \right) (\bar{u}_1 - \bar{u}) g_{\delta, 1} d\xi \quad (S7)$$

$$+ \int \left(\frac{g}{g_1} - 1 \right) (m_1 - m) b_1 r g_{\delta, 1} d\xi \quad (S8)$$

$$- \int \frac{g}{g_1} (b_1 - b)(v_1 - v) g_{\delta, 1} d\xi \quad (S9)$$

$$- \int (\bar{u}_1 - \bar{u})(g_{\delta, 1} - g_{\delta}) d\xi \quad (S10)$$

and

$$PD_{\eta_1, \delta}^2 - \psi_2(\delta) = \int \left(\frac{e}{e_1} \frac{d}{d_1} - 1 \right) (m - m_1) b_1 h g_{\delta, 1} d\xi$$

$$+ \int \frac{g}{g_1} (b_1 - b)(s_1 - s) g_{\delta, 1} d\xi$$

$$- \int (q_1 - q)(g_{\delta, 1} - g_{\delta}) d\xi.$$

Proof Note that

$$\begin{aligned}
PS_{\eta_1, \delta}^1 - \psi_1(\delta) &= \int \left(\frac{\mathbf{g}}{\mathbf{g}_1} \frac{\mathbf{d}}{\mathbf{d}_1} - 1 \right) (\mathbf{m} - \mathbf{m}_1) \mathbf{b}_1 \mathbf{r} \mathbf{g}_{\delta, 1} d\xi + \int (\mathbf{m} - \mathbf{m}_1) \mathbf{b}_1 \mathbf{r} \mathbf{g}_{\delta, 1} d\xi \\
&\quad - \int \bar{\mathbf{u}}(\mathbf{g}_\delta - \mathbf{g}_{\delta, 1}) d\xi - \int \bar{\mathbf{u}} \mathbf{g}_{\delta, 1} d\xi \\
&\quad + \int \frac{\mathbf{g}}{\mathbf{g}_1} (\mathbf{b} - \mathbf{b}_1) \mathbf{v}_1 \mathbf{g}_{\delta, 1} d\xi + \int \frac{\mathbf{g}}{\mathbf{g}_1} (\mathbf{u}_1 \mathbf{r} - \bar{\mathbf{u}}_1) \mathbf{g}_{\delta, 1} d\xi + \int \bar{\mathbf{u}}_1 \mathbf{g}_{\delta, 1} d\xi \\
&= (S6) + \int (\bar{\mathbf{u}}_1 - \bar{\mathbf{u}}) \mathbf{g}_{\delta, 1} d\xi \\
&\quad + \int (\mathbf{m} - \mathbf{m}_1) \mathbf{b}_1 \mathbf{r} \mathbf{g}_{\delta, 1} d\xi + \int \frac{\mathbf{g}}{\mathbf{g}_1} (\mathbf{b} - \mathbf{b}_1) \mathbf{v}_1 \mathbf{g}_{\delta, 1} d\xi \\
&\quad + \int \frac{\mathbf{g}}{\mathbf{g}_1} \mathbf{u}_1 \mathbf{r} \mathbf{g}_{\delta, 1} d\xi - \int \frac{\mathbf{g}}{\mathbf{g}_1} \bar{\mathbf{u}}_1 \mathbf{g}_{\delta, 1} d\xi \\
&\quad - \int \bar{\mathbf{u}}(\mathbf{g}_\delta - \mathbf{g}_{\delta, 1}) d\xi \\
&= (S6) - \int \left(\frac{\mathbf{g}}{\mathbf{g}_1} - 1 \right) (\bar{\mathbf{u}}_1 - \bar{\mathbf{u}}) \mathbf{g}_{\delta, 1} d\xi \\
&\quad + \int \frac{\mathbf{g}}{\mathbf{g}_1} \mathbf{m}_1 \mathbf{b}_1 \mathbf{r} \mathbf{g}_{\delta, 1} d\xi - \int \frac{\mathbf{g}}{\mathbf{g}_1} \mathbf{m} \mathbf{b} \mathbf{r} \mathbf{g}_{\delta, 1} d\xi \\
&\quad + \int (\mathbf{m} - \mathbf{m}_1) \mathbf{b}_1 \mathbf{r} \mathbf{g}_{\delta, 1} d\xi + \int \frac{\mathbf{g}}{\mathbf{g}_1} (\mathbf{b} - \mathbf{b}_1) \mathbf{v}_1 \mathbf{g}_{\delta, 1} d\xi \\
&\quad - \int \bar{\mathbf{u}}(\mathbf{g}_\delta - \mathbf{g}_{\delta, 1}) d\xi \\
&= (S6) - (S7) \\
&\quad + \int (\mathbf{m} - \mathbf{m}_1) \mathbf{b}_1 \mathbf{r} \mathbf{g}_{\delta, 1} d\xi + \int \frac{\mathbf{g}}{\mathbf{g}_1} (\mathbf{b} - \mathbf{b}_1) \mathbf{v}_1 \mathbf{g}_{\delta, 1} d\xi \\
&\quad + \int \frac{\mathbf{g}}{\mathbf{g}_1} (\mathbf{m}_1 \mathbf{b}_1 + \mathbf{m} \mathbf{b}_1 - \mathbf{m} \mathbf{b}_1 - \mathbf{m} \mathbf{b}) \mathbf{r} \mathbf{g}_{\delta, 1} d\xi \\
&\quad - \int \bar{\mathbf{u}}(\mathbf{g}_\delta - \mathbf{g}_{\delta, 1}) d\xi \\
&= (S6) - (S7) \\
&\quad + \int (\mathbf{m} - \mathbf{m}_1) \mathbf{b}_1 \mathbf{r} \mathbf{g}_{\delta, 1} d\xi + \int \frac{\mathbf{g}}{\mathbf{g}_1} (\mathbf{b} - \mathbf{b}_1) \mathbf{v}_1 \mathbf{g}_{\delta, 1} d\xi \\
&\quad + \int \frac{\mathbf{g}}{\mathbf{g}_1} (\mathbf{m}_1 - \mathbf{m}) \mathbf{b}_1 \mathbf{r} \mathbf{g}_{\delta, 1} d\xi + \int \frac{\mathbf{g}}{\mathbf{g}_1} (\mathbf{b}_1 - \mathbf{b}) \mathbf{m} \mathbf{r} \mathbf{g}_{\delta, 1} d\xi \\
&\quad - \int \bar{\mathbf{u}}(\mathbf{g}_\delta - \mathbf{g}_{\delta, 1}) d\xi \\
&= (S6) - (S7) + (S8) - (S9) \\
&\quad - \int \bar{\mathbf{u}}(\mathbf{g}_\delta - \mathbf{g}_{\delta, 1}) d\xi.
\end{aligned} \tag{S11}$$

Using A1 we can change variables to obtain

$$PS_{\eta_1, \delta}^{A, 1} = \int \bar{\mathbf{u}}_1(\mathbf{g}_\delta - \mathbf{g}_{\delta, 1}) d\xi.$$

The proof for ψ_2 is analogous. This completes the proof of the theorem. \square

Lemma S2 (Second order terms for exponential tilting.). Define $c(w) = \{\int_a \exp(\delta a) g(a | w)\}^{-1}$, and let $c_1(w)$ be defined analogously. Let $b(a) = \exp(\delta a)$. Using the same notation as in Lemma 5, we have

$$\begin{aligned} PD_{\eta_1, \delta}^1 - \psi_1(\delta) &= \int \left(\frac{g}{g_1} \frac{d}{d_1} - 1 \right) (m - m_1) b_1 r g_{\delta, 1} d\xi \\ &\quad - \int \left(\frac{g}{g_1} - 1 \right) (\bar{u}_1 - \bar{u}) g_{\delta, 1} d\xi \\ &\quad + \int \left(\frac{g}{g_1} - 1 \right) (m_1 - m) b_1 r g_{\delta, 1} d\xi \\ &\quad - \int \frac{g}{g_1} (b_1 - b) (v_1 - v) g_{\delta, 1} d\xi \\ &\quad + \int (\bar{u}_1 - \bar{u}) (g_{\delta, 1} - g_\delta) d\xi \\ &\quad - \int \left\{ (c_1 - c)^2 \int b g_1 \bar{u}_1 d\kappa \int b g d\kappa \right\} d\xi \\ &\quad + \int \left\{ (c_1 - c) \int b \bar{u}_1 (g - g_1) d\kappa \right\} d\xi, \end{aligned}$$

and

$$\begin{aligned} PD_{\eta_1, \delta}^2 - \psi_2(\delta) &= \int \left(\frac{e}{e_1} \frac{d}{d_1} - 1 \right) (m - m_1) b_1 h g_{\delta, 1} d\xi \\ &\quad + \int \frac{g}{g_1} (b_1 - b) (s_1 - s) g_{\delta, 1} d\xi \\ &\quad - \int (q_1 - q) (g_{\delta, 1} - g_\delta) d\xi \\ &\quad - \int \left\{ (c_1 - c)^2 \int b g_1 \bar{q}_1 d\kappa \int b g d\kappa \right\} d\xi \\ &\quad + \int \left\{ (c_1 - c) \int b \bar{q}_1 (g - g_1) d\kappa \right\} d\xi. \end{aligned}$$

Proof In this proof, (S11) is also valid. We have

$$PS_{\eta_1, \delta}^{1, A} - \int \bar{u} (g_\delta - g_{\delta, 1}) d\xi = PS_{\eta_1, \delta}^{1, A} - \int \bar{u}_1 (g_\delta - g_{\delta, 1}) d\xi + \int (\bar{u}_1 - \bar{u}) (g_{\delta, 1} - g_\delta) d\xi$$

It thus remains to prove that

$$\begin{aligned} PS_{\eta_1, \delta}^{1, A} - \int \bar{u}_1 (g_\delta - g_{\delta, 1}) d\xi &= - \int \left\{ (c_1 - c)^2 \int b g_1 \bar{u}_1 d\kappa \int b g d\kappa \right\} d\xi \\ &\quad + \int \left\{ (c_1 - c) \int b \bar{u}_1 (g - g_1) d\kappa \right\} d\xi. \end{aligned}$$

We have

$$\begin{aligned}
& PS_{\eta_i}^{1,A} - \int \bar{u}_1(g_\delta - g_{\delta,1})d\xi \\
&= \int \left\{ \int \frac{g_{1,\delta}}{g_1} \bar{u}_1 g d\kappa - \int \frac{g_{1,\delta}}{g_1} g d\kappa \int \bar{u}_1 g_{1,\delta} d\kappa + \int (g_{1,\delta} - g_\delta) \bar{u}_1 d\kappa \right\} d\xi \\
&= \int \left\{ \frac{g_{1,\delta}}{g_1} g \bar{u}_1 d\kappa - \int g_\delta \bar{u}_1 d\kappa + \int g_{1,\delta} \bar{u}_1 d\kappa \left[1 - \int \frac{g_{1,\delta}}{g_1} g d\kappa \right] \right\} d\xi \\
&= \int \left\{ c_1 \int b \bar{u}_1 g d\kappa - c_1 \int b \bar{u}_1 g d\kappa + c_1 \int b g \bar{u}_1 d\kappa \int (c - c_1) b g d\kappa \right\} d\xi \\
&= \int (c_1 - c) \left\{ \int b \bar{u}_1 g d\kappa - c_1 \int b g_1 \bar{u}_1 d\kappa \int b g d\kappa \right\} d\xi \\
&= \int (c_1 - c) \left\{ \int b \bar{u}_1 g d\kappa - c \int b g_1 \bar{u}_1 d\kappa \int b g d\kappa - (c_1 - c) \int b g_1 \bar{u}_1 d\kappa \int b g d\kappa \right\} d\xi \\
&= \int \left\{ -(c_1 - c)^2 \int b g_1 \bar{u}_1 d\kappa \int b g d\kappa + (c_1 - c) \left[\int b \bar{u}_1 g d\kappa - \int b g_1 \bar{u}_1 d\kappa \right] \right\} d\xi \quad (S12) \\
&= \int \left\{ -(c_1 - c)^2 \int b g_1 \bar{u}_1 d\kappa \int b g d\kappa + (c_1 - c) \int b \bar{u}_1 (g - g_1) d\kappa \right\} d\xi,
\end{aligned}$$

where (S12) follows from $c \int b g d\kappa = 1$. The proof for ψ_2 is analogous. \square

References

- David Benkeser and Mark J van der Laan. The highly adaptive lasso estimator. In *2016 IEEE International Conference on Data Science and Advanced Analytics (DSAA)*, pages 689–696. IEEE, 2016.
- Aurélien F Bibaut and Mark J van der Laan. Fast rates for empirical risk minimization over càdlàg functions with bounded sectional variation norm. *arXiv preprint arXiv:1907.09244*, 2019.
- Victor Chernozhukov, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen, Whitney Newey, and James M Robins. Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal*, 21(1):C1–C68, 01 2018. doi: 10.1111/ectj.12097.
- Jeremy R Coyle, Nima S Hejazi, and Mark J van der Laan. *hal9001: The scalable highly adaptive lasso*, 2020. URL <https://doi.org/10.5281/zenodo.3558313>. R package version 0.2.6.
- Nima S Hejazi and Iván Díaz. *medshift: Causal mediation analysis for stochastic interventions*, 2020. URL <https://github.com/nhejazi/medshift>. R package version 0.1.4.
- Mark J van der Laan. A generally efficient targeted minimum loss based estimator based on the highly adaptive lasso. *International Journal of Biostatistics*, 13(2), 2017.
- Mark J van der Laan and Aurélien F Bibaut. Uniform consistency of the highly adaptive lasso estimator of infinite-dimensional parameters. *arXiv preprint arXiv:1709.06256*, 2017.
- Mark J van der Laan and Sherri Rose. *Targeted Learning: Causal Inference for Observational and Experimental Data*. Springer Science & Business Media, 2011.
- Mark J van der Laan, Eric C Polley, and Alan E Hubbard. Super Learner. *Statistical Applications in Genetics and Molecular Biology*, 6(1), 2007.
- Aad W van der Vaart. *Asymptotic statistics*, volume 3. Cambridge university press, 2000.
- Aad W van der Vaart and Jon A Wellner. *Weak convergence and empirical processes*. Springer, 1996.