

Some Doubly and Multiply Robust Estimators of Controlled Direct Effects*

Xiang Zhou

November 18, 2020

Abstract

This letter introduces several doubly, triply, and quadruply robust estimators of the controlled direct effect. Among them, the triply and quadruply robust estimators are locally semiparametric efficient, and well suited to the use of data-adaptive methods for estimating their nuisance functions.

1 Introduction

Over the past decade, causal mediation analysis has grown popular in social and biomedical sciences. A common approach to assessing causal mediation involves decomposing the total effect of a treatment on an outcome into the so-called natural direct and indirect effects (NDE and NIE; Robins and Greenland 1992; Pearl 2001). The NDE and NIE, however, are not nonparametrically identified in the presence of posttreatment confounders, i.e., when confounders of the mediator-outcome relationship may be affected by the treatment itself (Avin et al. 2005; VanderWeele and Vansteelandt 2009). In such cases, researchers have often focused on estimating the controlled direct effect (CDE), a quantity that measures the effect of treatment when a mediator is fixed at a given value for all units (Pearl 2001). Thus a nonzero CDE implies that the effect of treatment on the outcome does not operate exclusively through the mediator of interest. Unlike the NDE and NIE, the CDE is identified provided that all confounders for the treatment-outcome relationship

*Direct all correspondence to Xiang Zhou, Department of Sociology, Harvard University, 33 Kirkland Street, Cambridge MA 02138; email: xiang_zhou@fas.harvard.edu. The author thanks Aleksei Opacic for helpful comments.

and for the mediator-outcome relationship are observed, even if some of the mediator-outcome confounders are affected by the treatment itself.

Estimators of the CDE typically rely on correct specification of (at least) two nuisance functions about the conditional means/densities of the treatment, mediator, outcome, or posttreatment confounders. For example, the weighting estimator proposed by VanderWeele (2009) involves fitting two propensity score models, one for the treatment and one for the mediator, and the sequential g-estimator proposed by Vansteelandt (2009) involves fitting two outcome models, one for the observed outcome given all its antecedent variables and one for a “demediated” outcome given pretreatment confounders and the treatment (see also Joffe and Greene 2009). To alleviate bias due to model misspecification, Goetgeluk et al. (2009) proposed a doubly robust estimator of the CDE that depends on correct specification of (a) a model for treatment assignment, and (b) either an outcome model or a propensity score model for the mediator.

The causal structure underlying identification and estimation of the CDE is akin to that of estimating treatment effects with longitudinal data in the presence of time-varying confounders (Robins 1999). For the latter problem, Bang and Robins (2005) have proposed a doubly robust estimator for the mean of a potential outcome that depends on correct specification of either (a) propensity score models for treatment status at all time points or (b) models for an iteratively imputed outcome at all time points. In a recent paper, Rotnitzky et al. (2017) point out that the Bang-Robins estimator is actually “multiply robust” because it is consistent whenever the first k propensity score models and the last $K - k$ “outcome models” are correctly specified, where $0 \leq k \leq K$, and K is the number of time points (see also Molina et al. 2017). Moreover, these authors show that the Bang-Robins procedure can be further improved with a 2^K -robust estimator that requires correct specification of either the propensity score model or the outcome model at each time point.¹

Capitalizing on the above work, this letter introduces a set of doubly robust, triply robust,

¹In a separate strand of literature, the term “multiple robustness” has been used to characterize a class of estimators for the mean of incomplete data in a cross-sectional setting that are consistent if one of several models for the propensity score or one of several models for the outcome is correctly specified (e.g., Han and Wang 2013). Following Molina et al. (2017) and Rotnitzky et al. (2017), we use “multiple robustness” to characterize estimators that require modeling *multiple parts of the observed data likelihood* and are consistent if one of several sets of the corresponding models are correctly specified.

and quadruply robust estimators of the CDE, which, to the best of the author’s knowledge, are new to the causal mediation literature. While some of these estimators ($\hat{\psi}_{am}^{\text{tr1}}$ and $\hat{\psi}_{am}^{\text{qr}}$; see Section 5) are closely related to those proposed in Bang and Robins (2005) and Rotnitzky et al. (2017) for estimating time-varying treatment effects, the rest of them have not been discussed elsewhere. These estimators all involve estimating more than two nuisance functions; yet, under suitable regularity conditions, they are consistent and asymptotically normal (CAN) when only two of these nuisance functions are correctly specified and their estimates are \sqrt{n} -consistent. The triply and quadruply robust estimators are locally efficient, i.e., when all of the nuisance functions are correctly specified, they attain the semiparametric efficiency bound in the nonparametric model over observed data. Moreover, their estimating equations are Neyman orthogonal, encouraging the use of machine learning methods and cross-fitting to estimate the nuisance functions (Zheng and van der Laan 2011; Chernozhukov et al. 2018), in which case estimates of the CDE are semiparametric efficient when estimates of the nuisance functions, for example, all converge at faster-than- $n^{-1/4}$ rates.

2 Notation, Assumptions, and Identification

Let A denote treatment, M the mediator, Y the observed outcome, and $Y(a, m)$ the potential outcome under treatment status a and mediator value m . We focus on the simple setting where the treatment A and the mediator M are both discrete with finite support. In addition, we denote by X a vector of pretreatment variables that may confound the causal effect of (A, M) on Y , and denote by Z a vector of posttreatment variables that may confound the causal effect of M on Y . Note that the posttreatment confounders Z may themselves be affected by the treatment.

The controlled direct effect (CDE) is defined as the average effect of switching treatment status from a' to a while fixing the mediator at a given level m :

$$\text{CDE}(a, a', m) = \mathbb{E}[Y(a, m) - Y(a', m)]$$

By definition, the CDE is identified when the expected potential outcome $\mathbb{E}[Y(a, m)]$ is identified for any a and m . Thus, we focus on the latter estimand throughout the paper and denote it

as ψ_{am} . Since it is the expected potential outcome when both the treatment and the mediator are “controlled” at given values, we may refer to it as the controlled response function (CRF). The CRF can also be used to construct other estimands such as the controlled mediator effect $\text{CME}(a, m, m') = \psi_{am} - \psi_{am'}$ (Zheng and Zhou 2015) and the treatment-mediator interaction effect $(\psi_{am} - \psi_{am'}) - (\psi_{a'm} - \psi_{a'm'})$.

The CRF is identified under the assumptions of consistency, sequential ignorability, and positivity:

1. consistency: for any unit, if $A = a$ and $M = m$, then $Y = Y(a, m)$;
2. sequential ignorability: $Y(a, m) \perp\!\!\!\perp A|X, \forall a, m$, and $Y(a, m) \perp\!\!\!\perp M|X, A, Z, \forall a, m$.
3. positivity: $p_{A|X}(a|x) > \epsilon > 0$ and $p_{M|X,A,Z}(m|x, a, z) > \epsilon > 0 \forall a, m, x \in \text{supp}(X)$, and $z \in \text{supp}(Z|X = x, A = a)$,

where $p(\cdot)$ denotes a probability mass/density function. Under assumptions 1-3, the CRF (and hence the CDE) can be identified via Robins’s (1986) g-computation formula:

$$\psi_{am} = \iiint y dP(y|x, a, z, m) dP(z|x, a) dP(x), \quad (1)$$

where $P(u|v)$ denotes the cumulative distribution function of U given V .

3 G-Computation, Imputation, and Weighting

Using the law of iterated expectations, equation (1) can be written in several different forms, each of which points to a different way of estimating the CRF:

$$\psi_{am} = \iint \mathbb{E}[Y|x, a, z, m] dP(z|x, a) dP(x) \quad (\text{g-computation}) \quad (2)$$

$$= \mathbb{E}_X \mathbb{E}_{Z|X, A=a} \mathbb{E}[Y|X, A, Z, M = m] \quad (\text{pure imputation}) \quad (3)$$

$$= \mathbb{E} \left[\frac{\mathbb{I}(A = a) \mathbb{E}[Y|X, A, Z, M = m]}{\Pr[A = a|X]} \right] \quad (\text{imputation-then-weighting}) \quad (4)$$

$$= \mathbb{E} \left[\frac{\mathbb{I}(A = a) \mathbb{I}(M = m) Y}{\Pr[A = a|X] \Pr[M = m|X, A, Z]} \right] \quad (\text{pure weighting}) \quad (5)$$

$$= \mathbb{E}_X \mathbb{E} \left[\frac{\mathbb{I}(M = m) Y}{\Pr[M = m|X, A, Z]} \middle| X, A = a \right] \quad (\text{weighting-then-imputation}) \quad (6)$$

Equation (2) suggests a procedure akin to Robins’s g-computation algorithm: (1) fit a parametric model for the conditional distribution of Z given X and A ; (2) fit a parametric or semiparametric model for the conditional mean of Y given X , A , Z , and M ; and (3) evaluate the inner integral via Monte Carlo simulation and the outer integral via the empirical distribution of X . In the particular case where the models for $\mathbb{E}[Z|x, a]$ and $\mathbb{E}[Y|x, a, z, m]$ are both linear, equation (2) can be evaluated using a simple “regression-with-residuals” procedure (Zhou and Wodtke 2019). Equation (3) suggests a “pure imputation” procedure: (1) fit a model for the conditional mean of Y given X , A , Z , and M and obtain predicted values for all units at $M = m$, $\hat{\mathbb{E}}[Y|X, A, Z, M = m]$; (2) fit a model for the conditional mean of $\hat{\mathbb{E}}[Y|X, A, Z, M = m]$ given X and A and obtain its predicted values for all units at $A = a$; (3) average these predicted values over all units. This procedure is similar to the sequential g-estimator proposed in Vansteelandt (2009) and Joffe and Greene (2009). Equation (4) suggests an imputation-then-weighting procedure: (1) fit a model for the conditional mean of Y given X , A , Z , and M and obtain predicted values for all units at $M = m$, $\hat{\mathbb{E}}[Y|X, A, Z, M = m]$; (2) fit a propensity score model for treatment status and obtain fitted values $\widehat{\Pr}[A = a|X]$; (3) compute a weighted average of the predicted outcomes $\hat{\mathbb{E}}[Y|X, A, Z, M = m]$ with inverse-probability weights $\mathbb{I}(A = a)/\widehat{\Pr}[A = a|X]$.

Equation (5) suggests a “pure weighting” estimator (VanderWeele 2009): (1) fit a propensity score model for treatment status and obtain fitted values $\widehat{\Pr}[A = a|X]$; (2) fit a propensity score model for the mediator and obtain fitted values $\widehat{\Pr}[M = m|X, A, Z]$; and (3) compute a weighted average of observed outcomes with inverse-probability weights $\mathbb{I}(A = a)\mathbb{I}(M = m)/(\widehat{\Pr}[A = a|X] \cdot \widehat{\Pr}[M = m|X, A, Z])$. Finally, equation (6) suggests a weighting-then-imputation procedure: (1) fit a propensity score model for the mediator and obtain fitted values $\widehat{\Pr}[M = m|X, A, Z]$; (2) fit a model for the conditional mean of the inverse-probability-weighted outcome $\mathbb{I}(M = m)Y/\widehat{\Pr}[M = m|X, A, Z]$ given X and A and obtain predicted values for all units at $A = a$; (3) average these predicted values over all units.

All of the above estimators involve estimating two nuisance functions about the conditional means/distributions of the treatment, mediator, outcome, or posttreatment confounders. Specifically, the g-computation procedure requires correctly specified models for $\mathbb{E}[Y|x, a, z, m]$ and $P(z|x, a)$; the pure imputation estimator requires correctly specified models for $\mathbb{E}[Y|x, a, z, m]$ and $\mathbb{E}_{Z|x,a}\mathbb{E}[Y|X, A, Z, M = m]$; the imputation-then-weighting estimator requires correctly specified

models for $\mathbb{E}[Y|x, a, z, m]$ and $\Pr[A = a|x]$; the pure weighting estimator requires correctly specified models for $\Pr[A = a|x]$ and $\Pr[M = m|x, a, z]$; and the weighting-then-imputation estimator requires correctly specified models for $\Pr[M = m|x, a, z]$ and $\mathbb{E}\left[\frac{\mathbb{I}(M=m)Y}{\Pr[M=m|X,A,Z]}|x, a\right]$. When either of the two requisite models is misspecified, the resulting estimator will be inconsistent. Thus, in empirical applications where the confounders X and Z have many components, these estimators can be highly prone to model misspecification bias. In the following section, we introduce four “doubly robust” estimators, each of which requires correct specification of one particular nuisance function and *either* of two other nuisance functions.

4 Doubly Robust Estimators

Before proceeding, we introduce the following functions (treating a and m as fixed):

$$\begin{aligned}\mu_y(x, z) &:= \mathbb{E}[Y|x, a, z, m] \\ \nu_y(x) &:= \mathbb{E}_{Z|x,a}\mu_y(X, Z) \\ \pi_a(x) &:= \Pr[A = a|x] \\ \pi_m(x, z) &:= \Pr[M = m|x, a, z],\end{aligned}$$

Under assumptions 1-3, $\mu_y(x, z) = \mathbb{E}[Y(a, m)|x, a, z]$ and $\nu_y(x) = \mathbb{E}[Y(a, m)|x]$. Thus $\mu_y(x, z)$ reflects how the potential outcome $Y(a, m)$ depends on pretreatment confounders X and posttreatment confounders Z among units with treatment status a , and $\nu_y(x)$ reflects how the potential outcome $Y(a, m)$ depends on pretreatment confounders X . Let $\mu_y^w(x, z)$, $\nu_y^w(x)$, $\pi_a^w(x)$, and $\pi_m^w(x, z)$ denote a set of working models for these nuisance functions, and let $\hat{\mu}_y^w(x, z)$, $\hat{\nu}_y^w(x)$, $\hat{\pi}_a^w(x)$, and $\hat{\pi}_m^w(x, z)$ denote their estimates. In particular, consider three different two-step estimators of $\nu_y(x)$:

$$\hat{\nu}_y^w(x; \hat{\mu}_y^w) = \hat{\mathbb{E}}[\hat{\mu}_y^w(X, Z)|x, a] \tag{7}$$

$$\hat{\nu}_y^w(x; \hat{\pi}_m^w) = \hat{\mathbb{E}}\left[\frac{\mathbb{I}(M = m)Y}{\hat{\pi}_m^w(X, Z)}|x, a\right] \tag{8}$$

$$\hat{\nu}_y^w(x; \hat{\mu}_y^w, \hat{\pi}_m^w) = \hat{\mathbb{E}}\left[\hat{\mu}_y^w(X, Z) + \frac{\mathbb{I}(M = m)}{\hat{\pi}_m^w(X, Z)}(Y - \hat{\mu}_y^w(X, Z))|x, a\right], \tag{9}$$

where $\hat{\mathbb{E}}[U|x, a]$ denotes estimates of the conditional mean of a random variable U given $X = x$ and $A = a$. In the above equations, the notation $\hat{\nu}_y^w(x; \hat{\mu}_y^w)$ indicates that this quantity depends on previous estimates of $\mu_y^w(x, z)$, and the same applies to $\hat{\nu}_y^w(x; \hat{\pi}_m^w)$ and $\hat{\nu}_y^w(x; \hat{\mu}_y^w, \hat{\pi}_m^w)$. The last expression $\hat{\nu}_y^w(x; \hat{\mu}_y^w, \hat{\pi}_m^w)$ can be seen as a doubly robust estimator of $\nu_y(x)$: when $\nu_y(x)$ is correctly specified and either $\mu_y^w(x, z)$ or $\pi_m^w(x, z)$ is correctly specified, $\hat{\nu}_y^w(x; \hat{\mu}_y^w, \hat{\pi}_m^w)$ will be consistent for $\nu_y(x)$.

Now consider the following estimators of ψ_{am} :

$$\begin{aligned}\hat{\psi}_{am}^{\text{dr1}} &= \mathbb{P}_n \left[\hat{\nu}_y^w(X; \hat{\mu}_y^w) + \frac{\mathbb{I}(A = a)}{\hat{\pi}_a^w(X)} (\hat{\mu}_y^w(X, Z) - \hat{\nu}_y^w(X; \hat{\mu}_y^w)) \right] \\ \hat{\psi}_{am}^{\text{dr2}} &= \mathbb{P}_n \left[\hat{\nu}_y^w(X; \hat{\pi}_m^w) + \frac{\mathbb{I}(A = a)}{\hat{\pi}_a^w(X)} \left(\frac{\mathbb{I}(M = m)Y}{\hat{\pi}_m^w(X, Z)} - \hat{\nu}_y^w(X; \hat{\pi}_m^w) \right) \right] \\ \hat{\psi}_{am}^{\text{dr3}} &= \mathbb{P}_n \left[\frac{\mathbb{I}(A = a)}{\hat{\pi}_a^w(X)} \left(\hat{\mu}_y^w(X, Z) + \frac{\mathbb{I}(M = m)}{\hat{\pi}_m^w(X, Z)} (Y - \hat{\mu}_y^w(X, Z)) \right) \right] \\ \hat{\psi}_{am}^{\text{dr4}} &= \mathbb{P}_n \left[\hat{\nu}_y^w(X; \hat{\mu}_y^w, \hat{\pi}_m^w) \right],\end{aligned}$$

where $\mathbb{P}_n[\cdot] = n^{-1} \sum_i [\cdot]_i$. $\hat{\psi}_{am}^{\text{dr1}}$ can be seen as a combination of the pure imputation estimator and the imputation-then-weighting estimator; $\hat{\psi}_{am}^{\text{dr2}}$ a combination of the pure weighting estimator and the weighting-then-imputation estimator; $\hat{\psi}_{am}^{\text{dr3}}$ a combination of the pure weighting estimator and the imputation-then-weighting estimator; and $\hat{\psi}_{am}^{\text{dr4}}$ a combination of the pure imputation estimator and the weighting-then-imputation estimator. Their double robustness is given in Proposition 1.

Proposition 1. *Under assumptions 1-3 and suitable regularity conditions,*

1. $\hat{\psi}_{am}^{\text{dr1}}$ is CAN if $\mu_y^w(x, z)$ is correctly specified and either $\nu_y^w(x)$ or $\pi_a^w(x)$ is correctly specified.
2. $\hat{\psi}_{am}^{\text{dr2}}$ is CAN if $\pi_m^w(x, z)$ is correctly specified and either $\nu_y^w(x)$ or $\pi_a^w(x)$ is correctly specified.
3. $\hat{\psi}_{am}^{\text{dr3}}$ is CAN if $\pi_a^w(x)$ is correctly specified and either $\pi_m^w(x, z)$ or $\mu_y^w(x, z)$ is correctly specified.
4. $\hat{\psi}_{am}^{\text{dr4}}$ is CAN if $\nu_y^w(x)$ is correctly specified and either $\pi_m^w(x, z)$ or $\mu_y^w(x, z)$ is correctly specified.

The double robustness of these estimators is due to a similar logic to that of standard doubly robust estimators for the mean of incomplete data (Scharfstein et al. 1999; Robins et al. 2007). For example, for $\hat{\psi}_{am}^{\text{dr1}}$, when $\mu_y^w(x, z)$ and $\nu_y^w(x)$ are correctly specified, the second term inside $\mathbb{P}_n[\cdot]$ will have a zero mean (asymptotically), leaving only $\mathbb{P}_n[\hat{\nu}_y^w(X; \hat{\mu}_y^w)]$, i.e., the pure imputation estimator;

and when $\mu_y^w(x, z)$ and $\pi_a^w(x)$ are correctly specified, the terms involving $\hat{\nu}_y^w(X; \hat{\mu}_y^w)$ will have a zero mean, leaving only $\mathbb{P}_n[(\mathbb{I}(A = a)/\hat{\pi}_a^w(X))\hat{\mu}_y^w(X, Z)]$, i.e., the imputation-then-weighting estimator.

Among these doubly robust estimators, $\hat{\psi}_{am}^{\text{dr}_3}$ can be particularly useful in randomized trials where the treatment assignment mechanism is known. In this case, $\hat{\psi}_{am}^{\text{dr}_3}$ is consistent as long as either $\pi_m^w(X, Z)$ or $\mu_y^w(X, Z)$ is correctly specified. In observational studies, however, none of these nuisance functions is known a priori, and the relative utility of these estimators will depend on the subject matter knowledge the investigator might have about the data generating process. For example, if the investigator has a better understanding of the mediator model than of the outcome models, $\hat{\psi}_{am}^{\text{dr}_2}$ may be preferred. Yet, in many applications, little information is available about any part of the data generating process. In those cases, the multiply robust estimators presented below will be more useful as they do not hinge on correct specification of any particular nuisance function. Moreover, as we will see, they are more amenable to the use of flexible machine learning methods for estimating the nuisance functions.

5 Multiply Robust and Semiparametric Efficient Estimators

Henceforth, let $O = (X, A, Z, M, Y)$ denote the observed data, and \mathcal{P}_{np} a nonparametric model over O wherein all laws satisfy the positivity assumption described in Section 2. Define the following of submodels of \mathcal{P}_{np} :

- $\mathcal{P}_1 = \{P \in \mathcal{P}_{\text{np}}: \mu_y^w(x, z) \text{ and } \pi_a^w(x) \text{ are correctly specified}\}$
- $\mathcal{P}_2 = \{P \in \mathcal{P}_{\text{np}}: \mu_y^w(x, z) \text{ and } \nu_y^w(x) \text{ are correctly specified}\}$
- $\mathcal{P}_3 = \{P \in \mathcal{P}_{\text{np}}: \pi_m^w(x, z) \text{ and } \pi_a^w(x) \text{ are correctly specified}\}$
- $\mathcal{P}_4 = \{P \in \mathcal{P}_{\text{np}}: \pi_m^w(x, z) \text{ and } \nu_y^w(x) \text{ are correctly specified}\}$

Consider the following estimators of ψ_{am} :

$$\begin{aligned}\hat{\psi}_{am}^{\text{tr}_1} &= \mathbb{P}_n \left[\hat{\nu}_y^w(X; \hat{\mu}_y^w) + \frac{\mathbb{I}(A = a)}{\hat{\pi}_a^w(X)} (\hat{\mu}_y^w(X, Z) - \hat{\nu}_y^w(X; \hat{\mu}_y^w)) + \frac{\mathbb{I}(A = a)\mathbb{I}(M = m)}{\hat{\pi}_a^w(X)\hat{\pi}_m^w(X, Z)} (Y - \hat{\mu}_y^w(X, Z)) \right] \\ \hat{\psi}_{am}^{\text{tr}_2} &= \mathbb{P}_n \left[\hat{\nu}_y^w(X; \hat{\pi}_m^w) + \frac{\mathbb{I}(A = a)}{\hat{\pi}_a^w(X)} (\hat{\mu}_y^w(X, Z) - \hat{\nu}_y^w(X; \hat{\pi}_m^w)) + \frac{\mathbb{I}(A = a)\mathbb{I}(M = m)}{\hat{\pi}_a^w(X)\hat{\pi}_m^w(X, Z)} (Y - \hat{\mu}_y^w(X, Z)) \right] \\ \hat{\psi}_{am}^{\text{qr}} &= \mathbb{P}_n \left[\hat{\nu}_y^w(X; \hat{\mu}_y^w, \hat{\pi}_m^w) + \frac{\mathbb{I}(A = a)}{\hat{\pi}_a^w(X)} (\hat{\mu}_y^w(X, Z) - \hat{\nu}_y^w(X; \hat{\mu}_y^w, \hat{\pi}_m^w)) + \frac{\mathbb{I}(A = a)\mathbb{I}(M = m)}{\hat{\pi}_a^w(X)\hat{\pi}_m^w(X, Z)} (Y - \hat{\mu}_y^w(X, Z)) \right]\end{aligned}$$

The triple robustness of $\hat{\psi}_{am}^{\text{tr1}}$ and $\hat{\psi}_{am}^{\text{tr2}}$ and the quadruple robustness of $\hat{\psi}_{am}^{\text{qr}}$ are given below.

Proposition 2. *Under assumptions 1-3 and suitable regularity conditions, $\hat{\psi}_{am}^{\text{tr1}}$ is CAN in $\mathcal{P}_1 \cup \mathcal{P}_2 \cup \mathcal{P}_3$, $\hat{\psi}_{am}^{\text{tr2}}$ is CAN in $\mathcal{P}_1 \cup \mathcal{P}_3 \cup \mathcal{P}_4$, and $\hat{\psi}_{am}^{\text{qr}}$ is CAN in $\mathcal{P}_1 \cup \mathcal{P}_2 \cup \mathcal{P}_3 \cup \mathcal{P}_4$. In addition, $\hat{\psi}_{am}^{\text{tr1}}$, $\hat{\psi}_{am}^{\text{tr2}}$, and $\hat{\psi}_{am}^{\text{qr}}$ are all locally efficient in the sense that they attain the semiparametric efficiency bound of \mathcal{P}_{np} at $\mathcal{P}_1 \cap \mathcal{P}_3$, i.e., when all of the four nuisance functions are correctly specified.*

The multiple robustness of these estimators is due to a similar logic to that of the doubly robust estimators given previously. For example, for $\hat{\psi}_{am}^{\text{qr}}$, when $\mu_y^{\text{w}}(x, z)$ and $\pi_a^{\text{w}}(x)$ are correctly specified (\mathcal{P}_1), the terms involving $\hat{\nu}_y^{\text{w}}(X; \hat{\mu}_y^{\text{w}}, \hat{\pi}_m^{\text{w}})$ and the third term inside $\mathbb{P}_n[\cdot]$ will both have a zero mean (asymptotically), leaving only $\mathbb{P}_n[(\mathbb{I}(A = a)/\hat{\pi}_a^{\text{w}}(X))\hat{\mu}_y^{\text{w}}(X, Z)]$, the imputation-then-weighting estimator; when $\mu_y^{\text{w}}(x, z)$ and $\nu_y^{\text{w}}(x)$ are correctly specified (\mathcal{P}_2), both the second and third terms inside $\mathbb{P}_n[\cdot]$ will have a zero mean, leaving only $\mathbb{P}_n[\hat{\nu}_y^{\text{w}}(X; \hat{\mu}_y^{\text{w}}, \hat{\pi}_m^{\text{w}})]$, i.e., the doubly robust estimator $\hat{\psi}_{am}^{\text{dr4}}$; when $\pi_m^{\text{w}}(x, z)$ and $\pi_a^{\text{w}}(x)$ are correctly specified (\mathcal{P}_3), the terms involving $\hat{\mu}_y^{\text{w}}(X, Z)$ and $\hat{\nu}_y^{\text{w}}(X; \hat{\mu}_y^{\text{w}}, \hat{\pi}_m^{\text{w}})$ will both have a zero mean, leaving only $\mathbb{P}_n[(\mathbb{I}(A = a)\mathbb{I}(M = m))Y/(\hat{\pi}_a^{\text{w}}(X)\hat{\pi}_m^{\text{w}}(X, Z))]$, i.e., the pure weighting estimator; and when $\pi_m^{\text{w}}(x, z)$ and $\nu_y^{\text{w}}(x)$ are correctly specified (\mathcal{P}_4), the terms involving $\hat{\mu}_y^{\text{w}}(X, Z)$ and $\hat{\pi}_a^{\text{w}}(X)$ will both have a zero mean, leaving only $\mathbb{P}_n[\hat{\nu}_y^{\text{w}}(X; \hat{\mu}_y^{\text{w}}, \hat{\pi}_m^{\text{w}})]$, i.e., the doubly robust estimator $\hat{\psi}_{am}^{\text{dr4}}$.

The asymptotic efficiency of these estimators is due to the fact that they all solve the estimating equation formed by the efficient influence function of ψ_{am} , which is

$$\varphi_{am}(O) = \nu_y(X) + \frac{\mathbb{I}(A = a)}{\pi_a(X)}(\mu_y(X, Z) - \nu_y(X)) + \frac{\mathbb{I}(A = a)\mathbb{I}(M = m)}{\pi_a(X)\pi_m(X, Z)}(Y - \mu_y(X, Z)) - \psi_{am}, \quad (10)$$

and the fact that $\mathbb{E}[\varphi_{am}(O; \eta)]$ has a zero derivative with respect to the nuisance functions $\eta = (\mu_y^{\text{w}}(x, z), \nu_y^{\text{w}}(x), \pi_a^{\text{w}}(x), \pi_m^{\text{w}}(x, z))$ at the truth (for a derivation of this influence function in the context of time-varying treatments, see van der Laan and Gruber 2012). The latter property implies that first step estimation of the nuisance functions will have no (first-order) effect on the influence function of $\hat{\psi}_{am}^{\text{tr1}}$, $\hat{\psi}_{am}^{\text{tr2}}$, and $\hat{\psi}_{am}^{\text{qr}}$. In practice, the nuisance functions can be estimated via data-adaptive methods combined with cross-fitting (Zheng and van der Laan 2011; Chernozhukov et al. 2018), in which case estimates of ψ_{am} (and hence CDE) are semiparametric efficient when estimates of the nuisance functions, for example, all converge at faster-than- $n^{-1/4}$ rates.²

²More precisely, $\hat{\psi}_{am}^{\text{tr1}}$, $\hat{\psi}_{am}^{\text{tr2}}$, and $\hat{\psi}_{am}^{\text{qr}}$ are semiparametric efficient if $R_n(\hat{\pi}_a)R_n(\hat{\nu}_y) +$

Among the above estimators, $\hat{\psi}_{am}^{\text{tr1}}$ is akin to the estimator proposed by Bang and Robins (2005) for the mean of a potential outcome with time-varying treatments and time-varying confounders. Specifically, they suggest that $\mathbb{I}(A = a)\mathbb{I}(M = m)Y/(\hat{\pi}_a^{\text{w}}(X)\hat{\pi}_m^{\text{w}}(X, Z))$ be included as a covariate in a generalized linear model (with canonical link) for $\hat{\mu}_y^{\text{w}}(x, z)$, and $\mathbb{I}(A = a)/\hat{\pi}_a^{\text{w}}(X)$ be included as a covariate in a generalized linear model (with canonical link) for $\hat{\nu}_y^{\text{w}}(x; \hat{\mu}_y^{\text{w}})$, in which case the score equations ensure that both the second and third terms inside $\mathbb{P}_n[\cdot]$ have a *zero sample mean*, thus leaving only $\mathbb{P}_n[\hat{\nu}_y^{\text{w}}(X; \hat{\mu}_y^{\text{w}})]$, i.e., the pure imputation estimator. Because this procedure estimates ψ_{am} as a sample average of $\hat{\nu}_y^{\text{w}}(X; \hat{\mu}_y^{\text{w}})$, which typically resides in the parameter space of ψ_{am} , it tends to be more stable in finite samples than the unadjusted estimator $\hat{\psi}_{am}^{\text{tr1}}$ (Robins et al. 2007). $\hat{\psi}_{am}^{\text{tr2}}$ and $\hat{\psi}_{am}^{\text{qr}}$ differ from $\hat{\psi}_{am}^{\text{tr1}}$ only in the response variable they use to model $\nu_y(X)$: $\hat{\psi}_{am}^{\text{tr2}}$ uses $\mathbb{I}(M = m)/\hat{\pi}_m^{\text{w}}(X, Z)$ whereas $\hat{\psi}_{am}^{\text{qr}}$ uses $\hat{\mu}_y^{\text{w}}(X, Z) + \mathbb{I}(M = m)(Y - \hat{\mu}_y^{\text{w}}(X, Z))/\hat{\pi}_m^{\text{w}}(X, Z)$, which adds another layer of robustness. In fact, $\hat{\psi}_{am}^{\text{qr}}$ constitutes a special case of the 2^K -robust estimator proposed by Rotnitzky et al. (2017) in the context of time-varying treatments. In practice, the Bang-Robins procedure can also be applied to $\hat{\psi}_{am}^{\text{tr2}}$ and $\hat{\psi}_{am}^{\text{qr}}$ to improve their finite-sample performance.

When flexible machine learning methods (instead of generalized linear models) are used to estimate the nuisance functions, the Bang-Robins procedure can no longer ensure a zero sample mean of the second and third terms inside $\mathbb{P}_n[\cdot]$. In this case, the method of targeted maximum likelihood estimation (TMLE; van Der Laan and Rubin 2006) can be used to adjust the first step estimates of $\mu_y(x, z)$ and $\nu_y(x)$ such that the second and third terms inside $\mathbb{P}_n[\cdot]$ have a zero sample mean. This approach may yield better finite-sample performance than the unadjusted estimators and more robustness than the Bang-Robins procedure based on generalized linear models.

For inference of $\hat{\psi}_{am}^{\text{tr1}}$, $\hat{\psi}_{am}^{\text{tr2}}$, $\hat{\psi}_{am}^{\text{qr}}$, and the corresponding estimates of the CDE, the nonparametric bootstrap can be used when the nuisance functions are estimated using parametric models. When data-adaptive methods are used to estimate the nuisance functions, it will be reasonable to use the empirical analog of the efficient influence function to construct standard errors and Wald-type confidence intervals. For example, the variance of $\widehat{\text{CDE}}(a, a', m)$ can be estimated by $\mathbb{P}_n[(\hat{\varphi}_{am}(O) - \hat{\varphi}_{a'm}(O))^2]/n$. When the CDE is defined on the risk ratio or odds ratio scale, corresponding variance estimates can be obtained using the delta method.

$R_n(\hat{\pi}_m)R_n(\hat{\mu}_y) = o(n^{-1/2})$, where $R_n(\cdot)$ maps a nuisance function to its $L_2(P)$ convergence rate with respect to the true distribution P . See Supporting Material C or Rotnitzky et al. (2017).

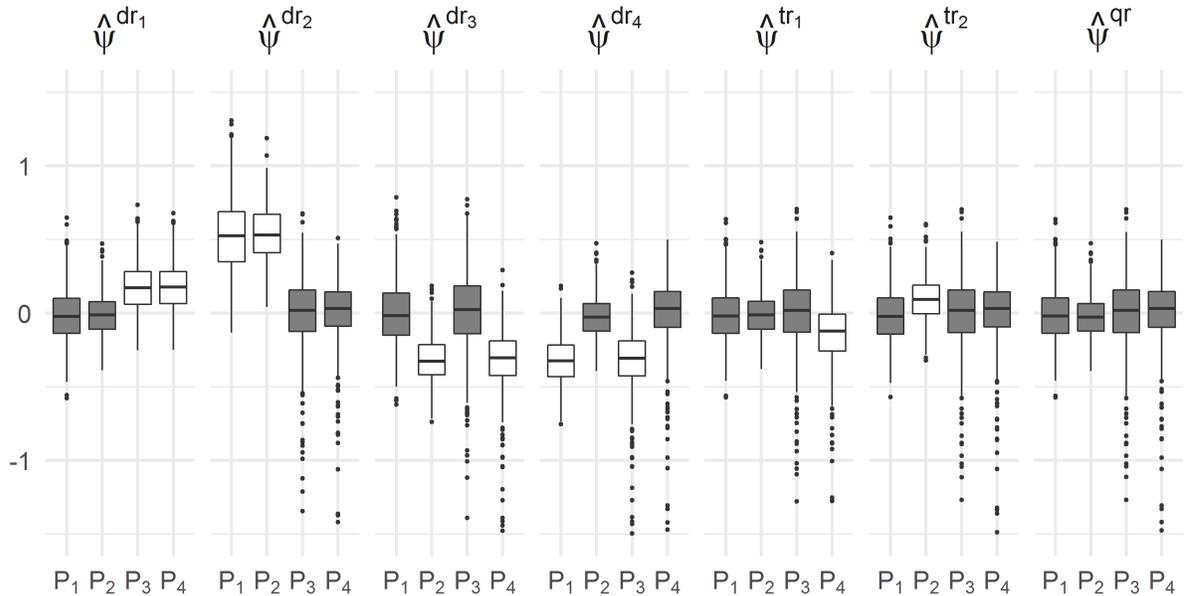


Figure 1: Sampling distributions of the doubly and multiply robust estimators for $n = 2,000$.

6 Monte Carlo Evidence

We now present a simulation study to demonstrate the multiple robustness of the proposed estimators. The data generating process is similar to that used in Miles et al. (2020) and is described in greater detail in Supporting Material D. We generate 1,000 Monte Carlo samples of size 2,000, and, without loss of generality, focus on the estimand $\psi_{01} = \mathbb{E}[Y(0, 1)]$. We examine the sampling distributions of all of the doubly and multiply robust estimators described above under conditions associated with submodels $\mathcal{P}_1, \mathcal{P}_2, \mathcal{P}_3, \mathcal{P}_4$. The results are shown in Figure 1, where each panel corresponds to an estimator, and the y axis is recentered at the true value of ψ_{01} . The shaded box plots highlight the cases under which a given estimator should be consistent. We can see that all of the doubly and multiply robust estimators behave as expected. They center around the true value if and only if the requisite nuisance functions are all correctly specified.

References

Avin, C., Shpitser, I., Pearl, J., 2005. Identifiability of path-specific effects, in: Proceedings of the 19th International Joint Conference on Artificial Intelligence, Morgan Kaufmann Publishers Inc..

- pp. 357–363.
- Bang, H., Robins, J.M., 2005. Doubly robust estimation in missing data and causal inference models. *Biometrics* 61, 962–973.
- Chernozhukov, V., Chetverikov, D., Demirer, M., Duflo, E., Hansen, C., Newey, W., Robins, J., 2018. Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal* 21, C1–C68.
- van Der Laan, M.J., Rubin, D., 2006. Targeted maximum likelihood learning. *The International Journal of Biostatistics* 2.
- Goetgeluk, S., Vansteelandt, S., Goetghebeur, E., 2009. Estimation of controlled direct effects. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 70, 1049–1066.
- Han, P., Wang, L., 2013. Estimation with missing data: Beyond double robustness. *Biometrika* 100, 417–430.
- Joffe, M.M., Greene, T., 2009. Related causal frameworks for surrogate outcomes. *Biometrics* 65, 530–538.
- van der Laan, M.J., Gruber, S., 2012. Targeted minimum loss based estimation of causal effects of multiple time point interventions. *The International Journal of Biostatistics* 8.
- Miles, C.H., Shpitser, I., Kanki, P., Meloni, S., Tchetgen Tchetgen, E.J., 2020. On semiparametric estimation of a path-specific effect in the presence of mediator-outcome confounding. *Biometrika* 107, 159–172.
- Molina, J., Rotnitzky, A., Sued, M., Robins, J., 2017. Multiple robustness in factorized likelihood models. *Biometrika* 104, 561–581.
- Pearl, J., 2001. Direct and indirect effects, in: *Proceedings of the Seventeenth Conference on Uncertainty in Artificial Intelligence*, Morgan Kaufmann Publishers Inc.. pp. 411–420.
- Robins, J., 1986. A new approach to causal inference in mortality studies with a sustained exposure period-application to control of the healthy worker survivor effect. *Mathematical Modelling* 7, 1393–1512.

- Robins, J., Sued, M., Lei-Gomez, Q., Rotnitzky, A., 2007. Comment: Performance of double-robust estimators when “inverse probability” weights are highly variable. *Statistical Science* 22, 544–559.
- Robins, J.M., 1999. Marginal structural models versus structural nested models as tools for causal inference. *Statistical Models in Epidemiology: The Environment and Clinical Trials* .
- Robins, J.M., Greenland, S., 1992. Identifiability and exchangeability for direct and indirect effects. *Epidemiology* 3, 143–155.
- Rotnitzky, A., Robins, J., Babino, L., 2017. On the multiply robust estimation of the mean of the g-functional. arXiv preprint arXiv:1705.08582 .
- Scharfstein, D.O., Rotnitzky, A., Robins, J.M., 1999. Adjusting for nonignorable drop-out using semiparametric nonresponse models. *Journal of the American Statistical Association* 94, 1096–1120.
- VanderWeele, T.J., 2009. Marginal structural models for the estimation of direct and indirect effects. *Epidemiology* 20, 18–26.
- VanderWeele, T.J., Vansteelandt, S., 2009. Conceptual issues concerning mediation, interventions and composition. *Statistics and its Interface* 2, 457–468.
- Vansteelandt, S., 2009. Estimating direct effects in cohort and case–control studies. *Epidemiology* 20, 851–860.
- Zheng, C., Zhou, X.H., 2015. Causal mediation analysis in the multilevel intervention and multicomponent mediator case. *Journal of the Royal Statistical Society: Series B: Statistical Methodology* 77, 581–615.
- Zheng, W., van der Laan, M.J., 2011. Cross-validated targeted minimum-loss-based estimation, in: *Targeted Learning*. Springer, New York, NY, pp. 459–474.
- Zhou, X., Wodtke, G.T., 2019. A regression-with-residuals method for estimating controlled direct effects. *Political Analysis* 27, 360–369.

A Proof of Equations (2-5).

Starting from equation (6), we have

$$\begin{aligned}
& \mathbb{E}_X \mathbb{E} \left[\frac{\mathbb{I}(M = m)Y}{\Pr[M = m|X, A, Z]} \middle| X, A = a \right] \\
&= \mathbb{E}_X \left[\mathbb{E} \left[\frac{\mathbb{I}(A = a)\mathbb{I}(M = m)Y}{\Pr[A = a|X] \Pr[M = m|X, A, Z]} \middle| X, A = a \right] \cdot \Pr[A = a|X] + 0 \cdot \Pr[A \neq a|X] \right] \\
&= \mathbb{E} \left[\frac{\mathbb{I}(A = a)\mathbb{I}(M = m)Y}{\Pr[A = a|X] \Pr[M = m|X, A, Z]} \right] \tag{11}
\end{aligned}$$

$$\begin{aligned}
&= \mathbb{E} \left[\frac{\mathbb{I}(A = a)}{\Pr[A = a|X]} \mathbb{E} \left[\frac{\mathbb{I}(M = m)Y}{\Pr[M = m|X, A, Z]} \middle| X, A, Z \right] \right] \\
&= \mathbb{E} \left[\frac{\mathbb{I}(A = a)}{\Pr[A = a|X]} \mathbb{E} \left[\frac{\mathbb{I}(M = m)Y}{\Pr[M = m|X, A, Z]} \middle| X, A, Z, M = m \right] \Pr[M = m|X, A = a, Z] \right] \\
&= \mathbb{E} \left[\frac{\mathbb{I}(A = a)\mathbb{E}[Y|X, A, Z, M = m]}{\Pr[A = a|X]} \right] \tag{12}
\end{aligned}$$

$$\begin{aligned}
&= \mathbb{E}_{X,A} \mathbb{E} \left[\frac{\mathbb{I}(A = a)\mathbb{E}[Y|X, A, Z, M = m]}{\Pr[A = a|X]} \middle| X, A \right] \\
&= \mathbb{E}_{X,A} \left[\mathbb{E} \left[\frac{\mathbb{I}(A = a)\mathbb{E}[Y|X, A, Z, M = m]}{\Pr[A = a|X]} \middle| X, A = a \right] \Pr[A = a|X] \right] \\
&= \mathbb{E}_X \mathbb{E} [\mathbb{E}[Y|X, A, Z, M = m] | X, A = a] \tag{13}
\end{aligned}$$

$$= \iiint \mathbb{E}[Y|x, a, z, m] dP(z|x, a) dP(x). \tag{14}$$

Equations (11), (12), (13), and (14) correspond to equations (5), (4), (3), and (2), respectively.

B Proof of Proposition 1

Below we show that $\hat{\psi}_{am}^{\text{dr1}}$ is CAN when (a) $\mu_y^w(x, z)$ is correctly specified and (b) either $\nu_y^w(x)$ or $\pi_a^w(x)$ is correctly specified. The double robustness of $\hat{\psi}_{am}^{\text{dr2}}$, $\hat{\psi}_{am}^{\text{dr3}}$, and $\hat{\psi}_{am}^{\text{dr4}}$ can be verified analogously.

A first-order Taylor expansion of $\hat{\psi}_{am}^{\text{dr1}}$ implies that

$$\hat{\psi}_{am}^{\text{dr1}} = \mathbb{P}_n [\nu_y^*(X) + \frac{\mathbb{I}(A = a)}{\pi_a^*(X)} (\mu_y^*(X, Z) - \nu_y^*(X))] + o_p(1),$$

where $\nu_y^*(X)$, $\pi_a^*(X)$ and $\mu_y^*(X, Z)$ denote the probability limits of $\hat{\nu}_y^w(X)$, $\hat{\pi}_a^w(X)$, and $\hat{\mu}_y^w(X, Z)$. Hence it suffices to show $\mathbb{E} [\nu_y^*(X) + \frac{\mathbb{I}(A = a)}{\pi_a^*(X)} (\mu_y^*(X, Z) - \nu_y^*(X))] = \psi_{am}$ if $\mu_y^*(X, Z) = \mu_y(X, Z)$ and

either $\pi_a^*(X) = \pi_a(X)$ or $\nu_y^*(X) = \nu_y(X)$. Consistency follows from the law of large numbers, and asymptotic normality follows from standard regularity conditions for M-estimators.

When $\mu_y^*(X, Z) = \mu_y(X, Z)$ and $\pi_a^*(X) = \pi_a(X)$,

$$\begin{aligned}
\text{plim } \hat{\psi}_{am}^{\text{dr}_1} &= \mathbb{E} \left[\nu_y^*(X) + \frac{\mathbb{I}(A = a)}{\pi_a(X)} (\mu_y(X, Z) - \nu_y^*(X)) \right] \\
&= \mathbb{E}_X \mathbb{E} \left[\nu_y^*(X) + \frac{\mathbb{I}(A = a)}{\pi_a(X)} (\mu_y(X, Z) - \nu_y^*(X)) \mid X \right] \\
&= \mathbb{E}_X \left[\nu_y^*(X) + \mathbb{E}[\mu_y(X, Z) - \nu_y^*(X) \mid X, A = a] \right] \\
&= \mathbb{E}_X \left[\nu_y^*(X) + \mathbb{E}[\mu_y(X, Z) \mid X, A = a] - \nu_y^*(X) \right] \\
&= \mathbb{E}_X \mathbb{E}[\mu_y(X, Z) \mid X, A = a] \\
&= \psi_{am}.
\end{aligned}$$

When $\mu_y^*(X, Z) = \mu_y(X, Z)$ and $\nu_y^*(X) = \nu_y(X)$,

$$\begin{aligned}
\text{plim } \hat{\psi}_{am}^{\text{dr}_1} &= \mathbb{E} \left[\nu_y(X) + \frac{\mathbb{I}(A = a)}{\pi_a^*(X)} (\mu_y(X, Z) - \nu_y(X)) \right] \\
&= \mathbb{E}_X \mathbb{E} \left[\nu_y(X) + \frac{\mathbb{I}(A = a)}{\pi_a^*(X)} (\mu_y(X, Z) - \nu_y(X)) \mid X \right] \\
&= \mathbb{E}_X \left[\nu_y(X) + \mathbb{E} \left[\frac{\mathbb{I}(A = a)}{\pi_a^*(X)} (\mu_y(X, Z) - \nu_y(X)) \mid X, A = a \right] \pi_a(X) \right] \\
&= \mathbb{E}_X \left[\nu_y(X) + \frac{\pi_a(X)}{\pi_a^*(X)} \mathbb{E}[\mu_y(X, Z) - \nu_y(X) \mid X, A = a] \right] \\
&= \mathbb{E}_X [\nu_y(X) + 0] \\
&= \psi_{am}.
\end{aligned}$$

C Proof of Proposition 2

Below we show that $\hat{\psi}_{am}^{\text{tr}_1}$ is CAN in $\mathcal{P}_1 \cup \mathcal{P}_2 \cup \mathcal{P}_3$ and locally efficient in $\mathcal{P}_1 \cap \mathcal{P}_3$. The multiple robustness and local efficiency of $\hat{\psi}_{am}^{\text{tr}_2}$ and $\hat{\psi}_{am}^{\text{qr}}$ can be verified analogously.

A first-order Taylor expansion of $\hat{\psi}_{am}^{\text{tr}_1}$ implies that

$$\hat{\psi}_{am}^{\text{tr}_1} = \mathbb{P}_n \left[\nu_y^*(X) + \frac{\mathbb{I}(A = a)}{\pi_a^*(X)} (\mu_y^*(X, Z) - \nu_y^*(X)) + \frac{\mathbb{I}(A = a)\mathbb{I}(M = m)}{\pi_a^*(X)\pi_m^*(X, Z)} (Y - \mu_y^*(X, Z)) \right] + o_p(1),$$

where $\mu_y^*(X, Z)$, $\nu_y^*(X)$, $\pi_a^*(X)$, and $\pi_m^*(X, Z)$ denote the probability limits of $\hat{\mu}_y^w(X, Z)$, $\hat{\nu}_y^w(X)$, $\hat{\pi}_a^w(X)$, and $\hat{\pi}_m^w(X, Z)$. Hence it suffices to show that the expectation of the quantity inside $\mathbb{P}_n[\cdot]$ equals ψ_{am} in $\mathcal{P}_1 \cup \mathcal{P}_2 \cup \mathcal{P}_3$. Consistency follows from the law of large numbers, and asymptotic normality follows from standard regularity conditions for M-estimators. First, consider submodel \mathcal{P}_2 , under which we have $\mu_y^*(X, Z) = \mu_y(X, Z)$ and $\nu_y^*(X) = \nu_y(X)$. From the proof of proposition 1, we know that

$$\mathbb{E}\left[\nu_y(X) + \frac{\mathbb{I}(A = a)}{\pi_a^*(X)}(\mu_y(X, Z) - \nu_y(X))\right] = \psi_{am}.$$

Thus

$$\begin{aligned} \text{plim } \hat{\psi}_{am}^{\text{tr}_1} &= \psi_{am} + \mathbb{E}\left[\frac{\mathbb{I}(A = a)\mathbb{I}(M = m)}{\pi_a^*(X)\pi_m^*(X, Z)}(Y - \mu_y(X, Z))\right] \\ &= \psi_{am} + \mathbb{E}_X \mathbb{E}\left[\frac{\mathbb{I}(A = a)\mathbb{I}(M = m)}{\pi_a^*(X)\pi_m^*(X, Z)}(Y - \mu_y(X, Z)) \mid X\right] \\ &= \psi_{am} + \mathbb{E}_X \left[\mathbb{E}\left[\frac{\mathbb{I}(A = a)\mathbb{I}(M = m)}{\pi_a^*(X)\pi_m^*(X, Z)}(Y - \mu_y(X, Z)) \mid X, A = a\right] \pi_a(X) \right] \\ &= \psi_{am} + \mathbb{E}_X \left[\frac{\pi_a(X)}{\pi_a^*(X)} \mathbb{E}\left[\frac{\mathbb{I}(M = m)}{\pi_m^*(X, Z)}(Y - \mu_y(X, Z)) \mid X, A = a\right] \right] \\ &= \psi_{am} + \mathbb{E}_X \left[\frac{\pi_a(X)}{\pi_a^*(X)} \mathbb{E}_{Z \mid X, A = a} \mathbb{E}\left[\frac{\mathbb{I}(M = m)\pi_m(X, Z)}{\pi_m^*(X, Z)}(Y - \mu_y(X, Z)) \mid X, A = a, Z, M = m\right] \right] \\ &= \psi_{am} + \mathbb{E}_X \left[\frac{\pi_a(X)}{\pi_a^*(X)} \mathbb{E}_{Z \mid X, A = a} \frac{\pi_m(X, Z)}{\pi_m^*(X, Z)} \underbrace{\mathbb{E}[(Y - \mu_y(X, Z)) \mid X, A = a, Z, M = m]}_{=0} \right] \\ &= \psi_{am}. \end{aligned}$$

Then, under \mathcal{P}_1 , we have $\mu_y^*(X, Z) = \mu_y(X, Z)$ and $\pi_a^*(X) = \pi_a(X)$. From the proof of proposition 1, we know that

$$\mathbb{E}\left[\nu_y^*(X) + \frac{\mathbb{I}(A = a)}{\pi_a(X)}(\mu_y(X, Z) - \nu_y^*(X))\right] = \psi_{am}.$$

Thus $\text{plim } \hat{\psi}_{am}^{\text{tr}_1} = \psi_{am} + \mathbb{E}\left[\frac{\mathbb{I}(A = a)\mathbb{I}(M = m)}{\pi_a^*(X)\pi_m^*(X, Z)}(Y - \mu_y(X, Z))\right] = 0$ (directly from the above proof for submodel \mathcal{P}_2). Finally, under \mathcal{P}_3 , we have $\pi_a^*(X) = \pi_a(X)$ and $\pi_m^*(X, Z) = \pi_m(X, Z)$.

$$\begin{aligned} \text{plim } \hat{\psi}_{am}^{\text{tr}_1} &= \mathbb{E}\left[\nu_y^*(X) + \frac{\mathbb{I}(A = a)}{\pi_a(X)}(\mu_y^*(X, Z) - \nu_y^*(X)) + \frac{\mathbb{I}(A = a)\mathbb{I}(M = m)}{\pi_a(X)\pi_m(X, Z)}(Y - \mu_y^*(X, Z))\right] \\ &= \mathbb{E}_X[\nu_y^*(X)] + \mathbb{E}_X \left[\mathbb{E}\left[\frac{\mathbb{I}(A = a)}{\pi_a(X)}(\mu_y^*(X, Z) - \nu_y^*(X)) \mid X, A = a\right] \pi_a(X) \right] \\ &\quad + \mathbb{E}_X \mathbb{E}_{Z \mid X, A = a} \left[\mathbb{E}\left[\frac{\mathbb{I}(A = a)\mathbb{I}(M = m)}{\pi_a(X)\pi_m(X, Z)}(Y - \mu_y^*(X, Z)) \mid X, A = a, Z\right] \pi_a(X) \right] \end{aligned}$$

$$\begin{aligned}
&= \mathbb{E}_X[\nu_y^*(X)] + \mathbb{E}_X \mathbb{E}[(\mu_y^*(X, Z) - \nu_y^*(X)) | X, A = a] \\
&+ \mathbb{E}_X \mathbb{E}_{Z|X, A=a} \left[\mathbb{E} \left[\frac{\mathbb{I}(A = a) \mathbb{I}(M = m)}{\pi_a(X) \pi_m(X, Z)} (Y - \mu_y^*(X, Z)) | X, A = a, Z, M = m \right] \pi_m(X, Z) \pi_a(X) \right] \\
&= \mathbb{E}_X \left[\mu_y^*(X, Z) + \mathbb{E}_{Z|X, A=a} \mathbb{E}[(Y - \mu_y^*(X, Z)) | X, A = a, Z, M = m] \right] \\
&= \mathbb{E}_X \left[\mu_y^*(X, Z) + \mu_y(X, Z) - \mu_y^*(X, Z) \right] \\
&= \psi_{am}.
\end{aligned}$$

To show that $\hat{\psi}_{am}^{\text{tr1}}$ is locally efficient, we first verify that equation (10) is the efficient influence function of ψ_{am} in \mathcal{P}_{np} , i.e.,

$$\left. \frac{\partial \psi_{am}(t)}{\partial t} \right|_{t=0} = \mathbb{E}[\varphi_{a,m}^{\text{eff}}(O) S_0(O)], \quad (15)$$

where $S_0(O)$ is the score function for any one-dimensional submodel $P_t(O)$ evaluated at $t = 0$. We first note that $S_t(O)$ can be written as $S_t(O) = S_t(X) + S_t(A|X) + S_t(Z|X, A) + S_t(M|X, A, Z) + S_t(Y|X, A, Z, M)$, where $S_t(u|v) = \partial \log p_t(u|v) / \partial t$ and $p_t(u|v)$ is the conditional probability density/mass function of U given V . Using equation (1) and the product rule, the left hand side of equation (15) can be written as

$$\begin{aligned}
\left. \frac{\partial \psi_{am}(t)}{\partial t} \right|_{t=0} &= \left. \frac{\partial \iiint y dP_t(y|x, a, z, m) dP_t(z|x, a) dP_t(x)}{\partial t} \right|_{t=0} \\
&= \underbrace{\iiint y S_0(x) dP_0(y|x, a, z, m) dP_0(z|x, a) dP_0(x)}_{=:\phi_1} \\
&+ \underbrace{\iiint y S_0(z|x, a) dP_0(y|x, a, z, m) dP_0(z|x, a) dP_0(x)}_{=:\phi_2} \\
&+ \underbrace{\iiint y S_0(y|x, a, z, m) dP_0(y|x, a, z, m) dP_0(z|x, a) dP_0(x)}_{=:\phi_3} \\
&= \phi_1 + \phi_2 + \phi_3
\end{aligned}$$

where the second equality follows from the fact that $\partial dP_t(u|v) / \partial t = S_t(u|v) dP_t(u|v)$.

Before evaluating the right hand side of equation (15), we introduce the following shorthands:

$$\varphi_1(X) = \nu_y(X),$$

$$\begin{aligned}\varphi_2(X, A, Z) &= \frac{\mathbb{I}(A = a)}{\pi_a(X)} (\mu_y(X, Z) - \nu_y(X)), \\ \varphi_3(X, A, Z, M, Y) &= \frac{\mathbb{I}(A = a)\mathbb{I}(M = m)}{\pi_a(X)\pi_m(X, Z)} (Y - \mu_y(X, Z)).\end{aligned}$$

Thus $\varphi_{a,m}^{\text{eff}}(O) = \varphi_1(X) + \varphi_2(X, A, Z) + \varphi_3(X, A, Z, M, Y) - \psi_{am}$. We first observe

$$\begin{aligned}& \mathbb{E}[\varphi_1(X)S_0(O)] \\ &= \mathbb{E}[\varphi_1(X)(S_0(X) + S_0(A|X) + S_0(Z|X, A) + S_0(M|X, A, Z) + S_0(Y|X, A, Z, M))] \\ &= \mathbb{E}[\varphi_1(X)S_0(X)] + \mathbb{E}[\varphi_1(X)S_0(A|X)] + \dots + \mathbb{E}[\varphi_1(X)S_0(Y|X, A, Z, M)] \\ &= \mathbb{E}[\varphi_1(X)S_0(X)] + \mathbb{E}[\varphi_1(X)\underbrace{\mathbb{E}[S_0(A|X)|X]}_{=0}] + \dots + \mathbb{E}[\varphi_1(X)\underbrace{\mathbb{E}[S_0(Y|X, A, Z, M)|X, A, Z, M]}_{=0}] \\ &= \mathbb{E}[\varphi_1(X)S_0(X)] \\ &= \mathbb{E}[\mathbb{E}_{Z|X, A=a}\mathbb{E}[Y|X, A = a, Z, M = m]S_0(X)] \\ &= \phi_1\end{aligned}\tag{16}$$

where we used the fact that $\mathbb{E}[S(U|V)|V] = 0$ for any score function $S(U, V)$. Second,

$$\begin{aligned}& \mathbb{E}[\varphi_2(X, A, Z)S_0(O)] \\ &= \mathbb{E}[\varphi_2(X, A, Z)(S_0(X) + S_0(A|X) + S_0(Z|X, A) + S_0(M|X, A, Z) + S_0(Y|X, A, Z, M))] \\ &= \mathbb{E}[\varphi_2(X, A, Z)S_0(X)] + \mathbb{E}[\varphi_2(X, A, Z)S_0(A|X)] + \mathbb{E}[\varphi_2(X, A, Z)S_0(Z|X, A)] \\ &\quad + \mathbb{E}[\varphi_2(X, A, Z)\underbrace{\mathbb{E}[S_0(M|X, A, Z)|X, A, Z]}_{=0}] + \mathbb{E}[\varphi_2(X, A, Z)\underbrace{\mathbb{E}[S_0(Y|X, A, Z, M)|X, A, Z, M]}_{=0}] \\ &= \mathbb{E}[\varphi_2(X, A, Z)S_0(X)] + \mathbb{E}[\varphi_2(X, A, Z)S_0(A|X)] + \mathbb{E}[\varphi_2(X, A, Z)S_0(Z|X, A)] \\ &= \mathbb{E}[S_0(X)\underbrace{\mathbb{E}[\varphi_2(X, A, Z)|X, A]}_{=0}] + \mathbb{E}[S_0(A|X)\underbrace{\mathbb{E}[\varphi_2(X, A, Z)|X, A]}_{=0}] + \mathbb{E}[\varphi_2(X, A, Z)S_0(Z|X, A)] \\ &= \mathbb{E}[\varphi_2(X, A, Z)S_0(Z|X, A)] \\ &= \mathbb{E}\left[\frac{\mathbb{I}(A = a)}{\Pr[A = a|X]} (\mathbb{E}[Y|X, A = a, Z, M = m] - \varphi_1(X))S_0(Z|X, A)\right] \\ &= \mathbb{E}_X\left\{\mathbb{E}\left[\frac{\mathbb{I}(A = a)}{\Pr[A = a|X]} (\mathbb{E}[Y|X, A = a, Z, M = m] - \varphi_1(X))S_0(Z|X, A)\middle|X, A = a\right] \cdot \Pr[A = a|X]\right\} \\ &= \mathbb{E}_X\mathbb{E}\left[(\mathbb{E}[Y|X, A = a, Z, M = m] - \varphi_1(X))S_0(Z|X, A)\middle|X, A = a\right]\end{aligned}$$

$$\begin{aligned}
&= \mathbb{E}_X \mathbb{E}_{Z|X, A=a} [\mathbb{E}[Y|X, A=a, Z, M=m] S_0(Z|X, A)] \\
&= \phi_2
\end{aligned} \tag{17}$$

where the fifth line follows from the fact that

$$\begin{aligned}
\mathbb{E}[\varphi_2(X, A, Z)|X, A] &= \mathbb{E}\left[\frac{\mathbb{I}(A=a)}{\Pr[A=a|X]} (\mathbb{E}[Y|X, A=a, Z, M=m] - \varphi_1(X)) |X, A\right] \\
&= \frac{\mathbb{I}(A=a)}{\Pr[A=a|X]} \underbrace{\left[\mathbb{E}_{Z|X, A} \mathbb{E}[Y|X, A=a, Z, M=m] - \varphi_1(X)\right]}_{= \varphi_1(X)} \\
&= 0.
\end{aligned}$$

Third,

$$\begin{aligned}
&\mathbb{E}[\varphi_3(X, A, Z, M, Y) S_0(O)] \\
&= \mathbb{E}[\varphi_3(X, A, Z, M, Y) (S_0(X) + S_0(A|X) + S_0(Z|X, A) + S_0(M|X, A, Z) + S_0(Y|X, A, Z, M))] \\
&= \mathbb{E}\left[\left(S_0(X) + S_0(A|X) + S_0(Z|X, A) + S_0(M|X, A, Z)\right) \underbrace{\mathbb{E}[\varphi_3(X, A, Z, M, Y)|X, A, Z, M]}_{=0}\right] \\
&\quad + \mathbb{E}[\varphi_3(X, A, Z, M, Y) S_0(Y|X, A, Z, M)]
\end{aligned} \tag{18}$$

$$\begin{aligned}
&= \mathbb{E}[\varphi_3(X, A, Z, M, Y) S_0(Y|X, A, Z, M)] \\
&= \mathbb{E}\left[\frac{\mathbb{I}(A=a)\mathbb{I}(M=m)Y S_0(Y|X, A, Z, M)}{\Pr[A=a|X]\Pr[M=m|X, A=a, Z]}\right] \\
&\quad - \mathbb{E}\left[\frac{\mathbb{I}(A=a)\mathbb{I}(M=m)\mathbb{E}[Y|X, A=a, Z, M=m]}{\Pr[A=a|X]\Pr[M=m|X, A=a, Z]}\right] \underbrace{\mathbb{E}[S_0(Y|X, A, Z, M)|X, A, Z, M]}_{=0}
\end{aligned} \tag{19}$$

$$\begin{aligned}
&= \mathbb{E}\left\{\mathbb{E}\left[\frac{\mathbb{I}(A=a)\mathbb{I}(M=m)Y S_0(Y|X, A, Z, M)}{\Pr[A=a|X]\Pr[M=m|X, A=a, Z]}\right] \middle| X, A, Z, M=m\right\} \Pr[M=m|X, A, Z] \\
&= \mathbb{E}_{X, A} \mathbb{E}\left[\frac{\mathbb{I}(A=a)\Pr[M=m|X, A, Z]\mathbb{E}[Y S_0(Y|X, A, Z, M)|X, A, Z, M=m]}{\Pr[A=a|X]\Pr[M=m|X, A=a, Z]}\right] \middle| X, A \\
&= \mathbb{E}_X \left\{ \mathbb{E}\left[\frac{\mathbb{I}(A=a)\Pr[M=m|X, A, Z]\mathbb{E}[Y S_0(Y|X, A, Z, M)|X, A, Z, M=m]}{\Pr[A=a|X]\Pr[M=m|X, A=a, Z]}\right] \middle| X, A=a\right\} \Pr[A=a|X] \\
&= \mathbb{E}_X \left\{ \mathbb{E}\left[\mathbb{E}[Y S_0(Y|X, A, Z, M)|X, A, Z, M=m] \middle| X, A=a\right]\right\} \\
&= \phi_3
\end{aligned} \tag{20}$$

where the first equality follows from the fact that

$$\begin{aligned}
& \mathbb{E}[\varphi_3(X, A, Z, M, Y)|X, A, Z, M] \\
&= \mathbb{E}\left[\frac{\mathbb{I}(A = a)\mathbb{I}(M = m)(Y - \mathbb{E}[Y|X, A = a, Z, M = m])}{\Pr[A = a|X] \Pr[M = m|X, A = a, Z]}|X, A, Z, M\right] \\
&= \frac{\mathbb{I}(A = a)\mathbb{I}(M = m)(\mathbb{E}[Y|X, A, Z, M] - \mathbb{E}[Y|X, A = a, Z, M = m])}{\Pr[A = a|X] \Pr[M = m|X, A = a, Z]} \\
&= 0.
\end{aligned}$$

Summing up equations (16-20), we have

$$\begin{aligned}
& \mathbb{E}[\varphi_{a,m}^{\text{eff}}(O)S_0(O)] \\
&= \mathbb{E}[(\varphi_1(X) + \varphi_2(X, A, Z) + \varphi_3(X, A, Z, M, Y) - \psi_{am})S_0(O)] \\
&= \mathbb{E}[\varphi_1(X)S_0(O)] + \mathbb{E}[\varphi_2(X, A, Z)S_0(O)] + \mathbb{E}[\varphi_3(X, A, Z, M, Y)S_0(O)] - \psi_{am} \cdot \mathbb{E}[S_0(O)] \\
&= \phi_1 + \phi_2 + \phi_3 - 0 \\
&= \left. \frac{\partial \psi_{am}(t)}{\partial t} \right|_{t=0}.
\end{aligned}$$

One way to show the local efficiency of $\hat{\psi}_{am}^{\text{tr}_1}$ is to verify that $\mathbb{E}[\varphi_{am}(O; \eta)]$ has a zero derivative with respect to the nuisance functions $\eta = (\mu_y^w(X, Z), \nu_y^w(X), \pi_a^w(X), \pi_m^w(X, Z))$ at the truth. Suppose these nuisance functions are parameterized by different components of a vector-valued parameter β , where β_0 denotes the truth. We then have

$$\begin{aligned}
\left. \frac{\partial \varphi_{am}(O; \eta)}{\partial \beta} \right|_{\beta=\beta_0} &= \frac{\mathbb{I}(A = a)}{\pi_a(X)} \left[1 - \frac{\mathbb{I}(M = m)}{\pi_m(X, Z)}\right] \dot{\mu}_y^w(X, Z) \\
&\quad + \left[1 - \frac{\mathbb{I}(A = a)}{\pi_a(X)}\right] \dot{\nu}_y^w(X) \\
&\quad - \frac{\mathbb{I}(A = a)}{\pi_a^2(X)} \left[\mu_y(X, Z) - \nu_y(X) + \frac{\mathbb{I}(M = m)}{\pi_m(X, Z)} (Y - \mu_y(X, Z))\right] \dot{\pi}_a(X) \\
&\quad - \frac{\mathbb{I}(A = a)\mathbb{I}(M = m)}{\pi_a(X)\pi_m^2(X, Z)} [Y - \mu_y(X, Z)] \dot{\pi}_m(X, Z)
\end{aligned}$$

where $\dot{\mu}_y^w(X, Z)$, $\dot{\nu}_y^w(X)$, $\dot{\pi}_a(X)$, and $\dot{\pi}_m(X, Z)$ denote the derivatives of the corresponding functions with respect to β at β_0 . It is easy to verify that these components all have a zero mean. Thus $\mathbb{E}[\partial \varphi_{am}(O; \beta_0)/\partial \beta] = 0$, implying that the influence function of $\hat{\psi}_{am}^{\text{tr}_1}$ at $\mathcal{P}_1 \cap \mathcal{P}_3$ is $\varphi_{a,m}^{\text{eff}}(O)$.

Alternatively, we can also analyze the asymptotic expansion of $\hat{\psi}_{am}^{\text{tr}_1}$ to establish weaker conditions for its semiparametric efficiency. Denote $m(O; \eta) = \varphi_{am}(O; \eta) + \psi_{am}$, we have

$$\begin{aligned}\hat{\psi}_{am}^{\text{tr}_1} &= \mathbb{P}_n[m(O; \hat{\eta})] \\ &= \mathbb{P}_n[m(O; \eta)] + P[m(O; \hat{\eta}) - m(O; \eta)] + (\mathbb{P}_n - P)[m(O; \hat{\eta}) - m(O; \eta)],\end{aligned}\quad (21)$$

where $Pg = \int gdP$ denotes the expectation of function $g(O)$ taken at the truth. In equation (21), the first term can be analyzed with the standard central limit theorem and has an asymptotic variance of $\mathbb{E}[(\varphi_{am}(O; \eta))^2]$. The last term is an empirical process term that will be $o_p(1/\sqrt{n})$ if either the nuisance functions fall in a Donsker class or if cross-fitting is used to induce independence between $\hat{\eta}$ and O (Chernozhukov et al. 2018). By rearranging terms, using the law of iterated expectations, and applying the Cauchy-Schwartz inequality, we can rewrite the second term as

$$\begin{aligned}P[m(O; \hat{\eta}) - m(O; \eta)] &= P\left[\frac{(\hat{\pi}_a(X) - \pi_a(X))(\hat{\nu}_y(X) - \nu_y(X))}{\hat{\pi}_a(X)}\right] \\ &\quad + P\left[\frac{\mathbb{I}(A = a)(\hat{\pi}_m(X, Z) - \pi_m(X, Z))(\hat{\mu}_y(X, Z) - \mu_y(X, Z))}{\hat{\pi}_a(X)\hat{\pi}_m(X, Z)}\right]. \\ &\leq C_1\|\hat{\pi}_a(X) - \pi_a(X)\| \cdot \|\hat{\nu}_y(X) - \nu_y(X)\| \\ &\quad + C_2\|\hat{\pi}_m(X, Z) - \pi_m(X, Z)\| \cdot \|\hat{\mu}_y(X, Z) - \mu_y(X, Z)\|,\end{aligned}$$

where C_1 and C_2 are on the order of $O_p(1)$, and $\|g\| = (\int g^T g dP)^{1/2}$. The last line is due to the positivity assumption that $\pi_a(X)$ and $\pi_m(X, Z)$ are bounded away from zero. Thus the second term in equation (21) is asymptotically negligible if $R_n(\hat{\pi}_a)R_n(\hat{\nu}_y) + R_n(\hat{\pi}_m)R_n(\hat{\mu}_y) = o(n^{-1/2})$, where $R_n(\cdot)$ maps a nuisance function to its $L_2(P)$ convergence rate. This result implies that if all nuisance functions are consistently estimated and converge at faster than $n^{1/4}$ rates, then $\hat{\psi}_{am}^{\text{tr}_1}$ is semiparametric efficient.

D More Details of the Simulation Study

The variables X, A, Z, M, Y in the simulation study are generated via the following model:

$$(U_{XA}, U_{XZ}, U_{XM}, U_{XY}) \sim N(0, I_4)$$

$$\begin{aligned}
X &\sim N((U_{XA}, U_{XZ}, U_{XM}, U_{XY})\beta_X, 1) \\
A &\sim \text{Bernoulli}(\text{logit}^{-1}[(1, U_{XA}, |X|)\beta_A]) \\
Z &\sim N((1, U_{XZ}, X, X^2, A)\beta_Z, 1) \\
M &\sim \text{Bernoulli}(\text{logit}^{-1}[(1, U_{XM}, X, X^2, A, Z, XA, XZ)\beta_M]) \\
Y &\sim N((1, U_{XY}, X, X^2, A, Z, XZ, M, AM)\beta_Y, 1).
\end{aligned}$$

The coefficients $\beta_X, \beta_A, \beta_Z, \beta_M, \beta_Y$ are generated using a set of uniform distributions with certain constraints designed to create nontrivial degrees of model misspecification.

It can be shown that under the above model, the nuisance functions $\mu_y(x, z)$, $\nu_y(x)$, $\pi_a(x)$, $\pi_m(x, z)$ can be consistently estimated via the following GLMs:

$$\begin{aligned}
\mathbb{E}[Y|X, A, Z, M] &= (1, X, X^2, A, Z, XZ, M, AM)\theta; & \mu_y(X, Z) &= \mathbb{E}[Y|X, A = a, Z, M = m] \\
\mathbb{E}[U|X, A] &= (1, X, X^2, X^3, A, XA)\eta; & \nu_y(X) &= \mathbb{E}[U|X, A = a] \\
\pi_a(X) &= \text{logit}^{-1}[(1, |X|)\alpha] \\
\pi_m(X, Z) &= \text{logit}^{-1}[(1, X, X^2, A, Z, XA, XZ)\gamma].
\end{aligned}$$

Here U is the outcome variable used to fit the model for $\nu_y(X)$, as shown in equations (7-9). To demonstrate the multiple robustness of the proposed estimators, we also fit a misspecified model for each of the nuisance functions:

$$\begin{aligned}
\mathbb{E}[Y|X, A, Z, M] &= (1, X, A, Z, M)\tilde{\theta}; & \mu_y(X, Z) &= \mathbb{E}[Y|X, A = a, Z, M = m] \\
\mathbb{E}[U|X, A] &= (1, X, A)\tilde{\eta}; & \nu_y(X) &= \mathbb{E}[U|X, A = a] \\
\pi_a(X) &= \text{logit}^{-1}[(1, X)\tilde{\alpha}] \\
\pi_m(X, Z) &= \text{logit}^{-1}[(1, X)\tilde{\gamma}].
\end{aligned}$$

Each of the four cases described in Figure 1 reflects a combination of estimated nuisance functions from these correctly and incorrectly specified models. For example, for submodel \mathcal{P}_1 , we use correctly specified models for $\mu_y(x, z)$ and $\pi_a(x)$ and incorrectly specified models for $\nu_y(x)$ and $\pi_m(x, z)$ for all estimators.