

Quantifying an Adherence Path-Specific Effect
of Antiretroviral Therapy in the Nigeria
PEPFAR Program

Caleb Miles* Ilya Shpitser[†] Phyllis Kanki[‡]
Seema Meloni** Eric J. Tchetgen Tchetgen^{††}

*Harvard School of Public Health, chmiles@hsph.harvard.edu

[†]University of Southampton

[‡]Harvard School of Public Health, pkanki@hsph.harvard.edu

**Harvard School of Public Health, sthakore@hsph.harvard.edu

^{††}Harvard School of Public Health, etchetge@hsph.harvard.edu

This working paper is hosted by The Berkeley Electronic Press (bepress) and may not be commercially reproduced without the permission of the copyright holder.

<http://biostats.bepress.com/harvardbiostat/paper186>

Copyright ©2014 by the authors.

Quantifying an Adherence Path-Specific Effect of Antiretroviral Therapy in the Nigeria PEPFAR Program

Caleb Miles, Ilya Shpitser, Phyllis Kanki, Seema Meloni, and Eric Tchetgen
Tchetgen*

Abstract

Since the early 2000s, evidence has accumulated for a significant differential effect of first-line antiretroviral therapy (ART) regimens on human immunodeficiency virus (HIV) treatment outcomes, such as CD4 response and viral load suppression. This finding was replicated in our data from the Harvard President's Emergency Plan for AIDS Relief (PEPFAR) program in Nigeria. Investigators were interested in finding the source of these differences, i.e., understanding the mechanisms through which one regimen outperforms another, particularly via adherence. This amounts to a mediation question with adherence playing the role of mediator. Existing mediation analysis results, however, have relied on an assumption of no exposure-induced confounding of the intermediate variable, and generally require an assumption of no unmeasured confounding for nonparametric identification. Both assumptions are violated by the presence of drug toxicity. In this paper, we relax these assumptions and show that certain path-specific effects remain identified under weaker conditions. We focus on the path-specific effect solely mediated by adherence and not by toxicity and propose a suite of estimators for this effect, including a semiparametric-efficient, multiply-robust estimator. We illustrate with simulations and present results from a study applying the methodology to the Harvard PEPFAR data. Supplementary materials are available online.

Keywords: Human immunodeficiency virus, Mediation, Nonparametric identification, Unobserved confounding, Robustness

*Caleb Miles is Doctoral Student, Department of Biostatistics, Harvard School of Public Health, Boston, MA 02115. Ilya Shpitser is Lecturer in Statistics, Department of Mathematics, University of Southampton, Southampton, Hampshire, SO17 1BJ, United Kingdom. Phyllis Kanki is Professor and Seema Meloni is Research Associate, Department of Immunology and Infectious Diseases, Harvard School of Public Health, Boston, MA 02115. Eric Tchetgen Tchetgen is Associate Professor, Departments of Biostatistics and Epidemiology, Harvard School of Public Health, Boston, MA 02115. The authors gratefully acknowledge the hard work and dedication of the clinical, data, and laboratory staff at the PEPFAR supported Harvard/AIDS Prevention Initiative in Nigeria (APIN) hospitals that provided secondary data for this analysis. This work was funded, in part, by the US Department of Health and Human Services, Health Resources and Services Administration (U51HA02522) and by the National Institutes of Health (R01AI104459-01A1). The contents are solely the responsibility of the authors and do not represent the official views of the funding institutions.

1. INTRODUCTION

The President's Emergency Plan for AIDS Relief (PEPFAR) has been a highly successful program that has saved millions of lives worldwide since its inception in 2003. The Harvard School of Public Health was awarded one of the PEPFAR grants, receiving a total of \$362 million for work in Nigeria, Botswana, and Tanzania. The program has furnished these countries with invaluable medical infrastructure and provided AIDS care services in Nigeria for over 160,000 people and treatment to approximately 105,000 of those patients.

Our data set consists of previously antiretroviral therapy (ART)-naïve, human immunodeficiency virus (HIV)-1 infected, adult patients enrolled in the Harvard PEPFAR/AIDS Prevention Initiative in Nigeria (APIN) program between June 2004 and November 2010 who started ART in the program and were followed for at least 1 year after initiating ART. Upon entry into the Harvard/APIN PEPFAR HIV care program, all patients completed informed consent; all consent forms were approved by the institutional review boards at Harvard, APIN and all the corresponding Harvard/APIN PEPFAR HIV care and treatment sites. Patients not on one of 6 standard first-line regimens at baseline or seen at two of the hospitals without reliable viral load data were excluded from the data set. The analysis in this paper consists of only the complete cases, and results are given for all regimens but d4T+3TC+EFV (see Table 1 note for full drug names), due to the small sample of patients on this regimen as a consequence of it having been dropped mid-way through the program. (d4T+3TC+NVP was also dropped, but had a large enough sample to provide for stable inference.)

The significant funding support for AIDS treatment in resource-limited settings provided by PEPFAR and other international donor organizations relied on clinical trial data generated in resource rich settings. In order to maximize the benefit of providing ART to the largest number of patients, well-established drug regimens that were less costly were recommended and supported by the program. Studies dating back to the early 2000s have demonstrated evidence that these first-line regimens were not equally effective (Tang et al., 2012), and indeed, in the Harvard PEPFAR data, we have observed a significant differential effect of first-line ART regimens on virologic failure and, to a lesser extent, CD4 count. Since these were first-line regimens in use in most resource-limited settings, this difference could have widespread implications to the success

Table 1: Treatment regimen coding and their estimated average causal effects on risk of virologic failure (VF) and CD4 count

Code	ART regimen	Patients on regimen	RR of VF (s.e.)	log-RR of VF (s.e.)	Mean diff. in CD4 count
1	TDF + 3TC/FTC + EFV	1448 (14.6%)	0.65 (0.014)	-0.44 (0.12)	6.9 (7.4)
2	d4T + 3TC + NVP	854 (8.6%)	0.75 (0.017)	-0.29 (0.13)	-9.7 (6.8)
3	AZT + 3TC + EFV	1003 (10.1%)	0.78 (0.018)	-0.25 (0.14)	10.7 (8.7)
4	AZT + 3TC + NVP	4707 (47.4%)	0.82 (0.011)	-0.21 (0.078)	17.8 (4.9)
5	TDF + 3TC/FTC + NVP	1919 (19.3%)	-	-	-

NOTE: 3TC=lamivudine, AZT=zidovudine, d4T=stavudine, EFV=efavirenz, FTC=emtricitabine, NVP=nevirapine, TDF=tenofovir. Effects on risk of virologic failure are expressed on the risk ratio (RR) and log-risk ratio scale relative to treatment 5 and were estimated using inverse-probability weighted estimators. Effects on CD4 count are expressed on the mean difference scale relative to treatment 5 and were estimated using doubly-robust estimators. All effects adjusted for the confounders listed in Section 2.

of ART programs. These regimens and each of their corresponding total effects on virologic failure and CD4 count relative to a common reference treatment are reported in Table 1. The effects on virologic failure are reported as marginal and log-marginal risk ratios, and the effects on CD4 count and log CD4 count are reported on the mean-difference scale. Treatments were coded from strongest estimated effect on virologic failure to weakest, and the weakest treatment (TDF+3TC/FTC+NVP) was chosen as the reference for the purposes of the effects in Table 1. These effects are contrasts in the population between the risk of virologic failure had one intervened to assign everyone to a comparison-level treatment (1, 2, 3, or 4) and that if one had intervened to assign everyone to baseline treatment 5. The total effects of these regimens, however, do not quite tell the whole story. Investigators were interested in finding the source of these differences, i.e., understanding the mechanisms through which one regimen outperforms another. Mediation analysis serves to better explain these mechanisms that drive the differences in effects. This type of analysis has the potential to help target interventions to improve the performance of the less-robust regimens.

The total effect can be considered as a combination of effects, possibly in conflicting directions, through different pathways from the exposure to the outcome. Therefore, a weak total effect could be due to a combination of even weaker path-specific effects or several stronger

path-specific effects canceling one another out. One such path-specific effect could work strictly through biological pathways, in which case this population would benefit most from switching to a more favorable drug regimen. Alternatively, biological factors might play a comparatively smaller role relative to the effect of the treatment through nonbiological pathways, such as through adherence (Shpitser, 2013). We suspect a lack of adherence to treatment to be a driving mechanism of the observed differential effects, in which case it would be worth considering how to improve this mediating factor.

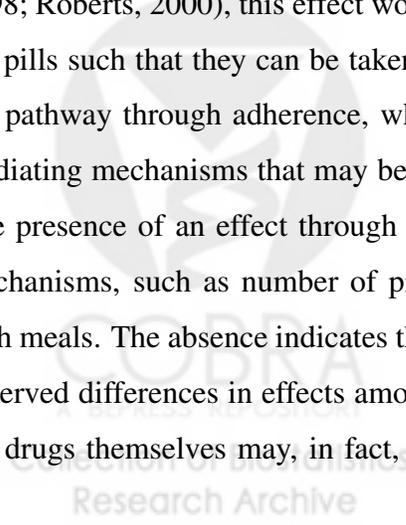
Adherence is widely accepted as a key factor for sustained viral suppression and is considered a prerequisite for maintenance on a prescribed drug regimen and optimal patient outcomes. However, the extent to which adherence to a given choice of first-line ART contributes to virologic failure (defined by the World Health Organization [WHO] as repeat viral load > 1000 copies/mL after 6 months of ART duration) is complex and still poorly understood and is a pressing mediation question in HIV research (Bangsberg et al., 2000). Understanding this issue is particularly important in resource-limited settings, where ART regimen options are few, and adherence to lifelong multi-drug daily dosing is challenging, but necessary. In such settings, quantifying to what degree differential rates of virologic failure are due to differences in adherence rates between therapies would inform the extent to which failure rates could be reduced by programs that improve adherence rates for certain ARTs, rather than changing the ART regimens themselves. Such adherence interventions have been very successful in the treatment of tuberculosis (China Tuberculosis Control Collaboration, 1996; Fujiwara et al., 1997; Suárez et al., 2001) and are considered similarly important in the treatment of HIV (Mills et al., 2006; Vranceanu et al., 2008; Pop-Eleches et al., 2011).

Remark. *Technically, the WHO also requires demonstration of adherence in their definition of virologic failure, which we avoid using in this paper since we cannot study the role of adherence as a mediator when it is part of the definition of the outcome.*

Among other potential mechanisms, the effect of treatment on virologic failure and CD4 count may be mediated by adherence, drug toxicity, or both. This study investigates the extent to which adherence, and not drug toxicity, mediates the effect, using the Harvard PEPFAR data set. That is, we focus on the role of adherence when it is differentially affected by the ways the drugs are obtained and taken, rather than by different levels of toxic side effects. The effect mediated by

nonadherence due to toxicity is unlikely to be appreciable, since toxicity in Nigeria is typically clinically recognizable and actionable. The magnitude of the roles of other drug-specific predictors of nonadherence, on the other hand, are less understood. These predictors also potentially point to lower-hanging fruit for development of adherence-promoting interventions. In mediation analysis terminology, we aim to estimate the effects of treatment assignment on virologic failure and CD4 count that are indirect with respect to adherence but direct with respect to toxicity. The definition, identification, and estimation of direct and indirect effects have received much attention in recent causal inference literature (Robins and Greenland, 1992; Robins, 1999, 2003; Pearl, 2001; Avin et al., 2005; Taylor et al., 2005; Petersen et al., 2006; Ten Have et al., 2007; Goetgeluk et al., 2008; van der Laan and Petersen, 2008; VanderWeele, 2009, 2011; VanderWeele and Vansteelandt, 2009, 2010; Imai et al., 2010a,b; Tchetgen Tchetgen, 2011; Tchetgen Tchetgen and Shpitser, 2014, 2012; Tchetgen Tchetgen, 2013).

The particular effect we are interested in can be classified as a *path-specific effect* (Pearl, 2001) – a class of estimands which can represent effects along any given causal pathway or collection of causal pathways. We consider the effect along the path from the provision of ART to virologic failure (or CD4 count) that goes through adherence, but not through toxicity. This effect is a measure of the change in risk of virologic failure (or mean CD4 count) were one to intervene on the mechanism by which the choice of treatment regimen directly, i.e., not through toxicity, affects adherence. For instance, if the difference in the effectiveness of ART through adherence were due to some regimens of ART having certain meal restrictions, posing a greater risk of patients missing dosages due to issues with food insecurity (Eldred et al., 1998; Gifford et al., 1998; Roberts, 2000), this effect would reflect the change in mean outcome if we were to modify the pills such that they can be taken without any meal restrictions. We emphasize our focus on the pathway through adherence, which does not involve toxicity, to learn about other possible mediating mechanisms that may be as important as toxicity, but are currently underappreciated. The presence of an effect through this pathway calls for closer investigation of these possible mechanisms, such as number of pills taken per dosage or the requirement that they be taken with meals. The absence indicates that differential effects through other pathways are driving the observed differences in effects among the treatment assignments. In particular, the efficacies of the drugs themselves may, in fact, differ, i.e., they may have a differential direct effect on the



outcome with respect to adherence, or they may have a differential effect on adherence due to their differing levels of toxicity.

Pearl (2001) defines path-specific effects, and Avin et al. (2005) provide general necessary and sufficient conditions for their identification for a single exposure and outcome, while Shpitser (2013) generalizes these definitions and conditions to settings with multiple exposures, multiple outcomes, and possible hidden variables. Our path-specific effect described above satisfies these identifying conditions, however an estimation strategy for its identifying functional does not yet exist. In this paper, we develop a suite of estimators (including a multiply-robust, semiparametric-efficient estimator) for the effect. The HIV case study detailed in this paper also functions as a guide for the application of this new method to analogous mediation settings where there is confounding that is affected by the exposure.

2. NOTATION & DEFINITIONS

To formalize our discussion, we begin by defining variables and counterfactuals. We will be considering pairwise comparisons of first-line ARTs prescribed to most HIV patients in Nigeria. Let E be an indicator of exposure to one of two such regimens of ART (coding given in Table 1). For notational simplicity, let e' denote the "reference level" treatment and e denote the "comparison level" treatment. Let C_1 be a bivariate vector of an indicator of any lab toxicities (alanine transaminase ≥ 120 UI/L, Creatinine ≥ 260 mmol/L, Hemoglobin ≤ 8 g/dL) observed six months after treatment initiation and an indicator that the patient's average percent adherence during the same six months, i.e., the total number of days that the patient had their drug supply divided by the number of days in the six month period, was no less than 95%. Let M be an indicator that the patient's average percent adherence during the subsequent six months was no less than 95%. Let Y be an indicator of whether the patient experienced virologic failure at the end of the year (based on viral load measurements at twelve and eighteen months for confirmation), or alternatively CD4 count at twelve months. Let C_0 be a vector of baseline confounders of the causal relationships between E , M , and Y not affected by exposure, viz. sex, age, marital status, WHO stage, hepatitis C virus, hepatitis B virus, CD4 count, and viral load. Throughout, we will assume that we observe i.i.d. sampling of $O = (C_0, E, C_1, M, Y)$.

We now consider counterfactuals under possible interventions on the variables (Rubin, 1974,

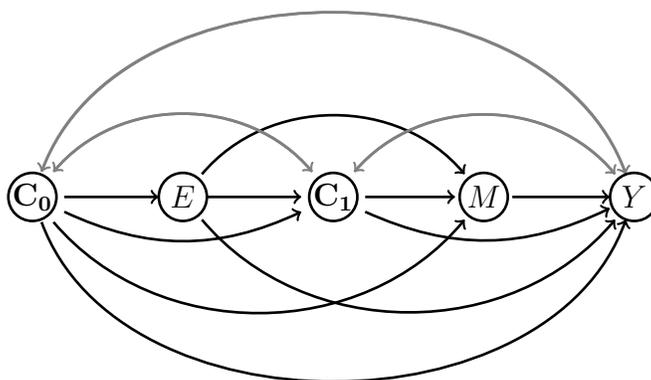


Figure 1: A causal graph with unobserved confounders that allows for identification of the \mathcal{P}_{EMY} -specific effect

1978). Let $Y(e^*)$ denote a patient's virologic suppression status or CD4 count if assigned, possibly contrary to fact, to the regimen of ART e^* . In the context of mediation, there will also be counterfactuals for intermediate variables. We define $C_1(e^*)$, $M(c_1, e^*)$ and $Y(m, e^*)$ similarly, and adopt the standard set of consistency assumptions (Robins, 1986) that if $E = e^*$, then $C_1(e^*) = C_1$ w.p.1, if $E = e^*$ and $C_1 = c_1$, then $M(c_1, e^*) = M$ w.p.1, if $E = e^*$ and $M = m$, then $Y(m, e^*) = Y$ w.p.1, and if $E = e^*$, then $Y(e^*) = Y$ w.p.1. Additionally, we adopt the standard set of positivity assumptions (Robins, 1986) that $f_{M|C_1, E, C_0}(m|C_1, E, C_0) > 0$ w.p.1 for each $m \in \text{supp}(M)$, $f_{C_1|E, C_0}(c_1|E, C_0) > 0$ w.p.1 for each $c_1 \in \text{supp}(C_1)$, $f_{E|C_0}(e^*|C_0) > 0$ w.p.1 for each $e^* \in \{e', e\}$.

To define the path-specific effect along the path $E \rightarrow M \rightarrow Y$, which we denote \mathcal{P}_{EMY} , we begin by discussing the graph in Figure 1. This is a complete graph of all observed variables in the sense that it includes all possible directed arrows that follow the natural temporal ordering. That is, any variable may directly affect any other variable succeeding it under this graph. The graph departs from the standard mediation graph (Baron and Kenny, 1986) in two important ways.

The first is with the presence of C_1 , which allows confounders of the effect of the mediator on the outcome to be affected by the exposure. In our HIV context, C_1 contains toxicity, which is clearly affected by the treatment assignment and may confound the effect of adherence on virologic failure. One way in which it might do this is on a biological level, toxicity might have an interactive effect with the drugs on the outcome, allowed for by the presence of the directed

arrow from C_1 to Y in conjunction with the directed arrow from E to Y . Thus, toxicity is a common cause of the outcome and adherence and, therefore, a confounder. Such a confounder is known as a *recanting witness*, due to its role in telling two conflicting “stories” about how E affects Y by being involved in two different pathway from E to Y – one involving M and the other not. Avin et al. (2005) showed the natural (or pure) direct and indirect effects (NDE and NIE, both highly popular in the mediation literature) (Robins and Greenland, 1992; Pearl, 2001) to be unidentified in the presence of a recanting witness.

The second way the graph in Figure 1 departs from the standard mediation graph is by the presence of the gray bidirected edges between C_0 , C_1 , and Y , each of which represents unobserved common causes between the two nodes to which it points. In the HIV application, these bidirected edges allow for the possibility of underlying biological factors which may be unobserved common causes of toxicity, the outcome, and biological baseline measurements such as viral load. The presence of these bidirected edges induces confounding of the effect of adherence on the outcome via toxicity, even if the arrow directed from C_1 to Y is absent. Since early adherence (during the first six months) may confound the effect of adherence at a later stage (during the subsequent six months) on the outcome, early adherence must be included in C_1 . Thus, \mathcal{P}_{EMY} involves only later adherence, and neither toxicity nor early-stage adherence.

As described above, we wish to quantify the mediating role of adherence along \mathcal{P}_{EMY} in Figure 1 which does not involve toxicity. Effects along such arbitrary (bundles of) causal pathways are known as path-specific effects (Pearl, 2000; Avin et al., 2005; Shpitser, 2013) and it is possible to define them inductively, which results in a quantity that is a function of a nested counterfactual (Shpitser, 2013). A general definition for the static-treatment and single-outcome case is given by Pearl (2000) and Avin et al. (2005). Defining

$$\begin{aligned}\beta_0 &\equiv \mathbb{E}[Y(M(e, C_1(e')), C_1(e'), e')] \\ \delta_0 &\equiv \mathbb{E}[Y(M(e', C_1(e')), C_1(e'), e')],\end{aligned}$$

the \mathcal{P}_{EMY} -specific effect, with respect to the comparison treatment value e and the baseline treatment value e' on the mean difference scale, is given by $\beta_0 - \delta_0$. δ_0 gives the mean outcome had everyone been assigned to the reference treatment regimen. β_0 gives the mean outcome had everyone been assigned to the reference treatment regimen, and adhered as they would have based

on the toxicity they experienced from this regimen, but otherwise as if they had been assigned to the comparison treatment. This is the \mathcal{P}_{EMY} -specific effect since it captures the impact of changing $M(e')$ to $M(e, \mathbf{C}_1(e'))$, which in turn would lead to an effect on Y only if M affects Y directly when all patients are assigned to e' .

3. IDENTIFICATION

Before introducing our identification result, we must first introduce a model that relaxes the assumption of independent errors of the Markovian model (Pearl, 2000) in a natural way. We will associate this model with the graph in Figure 1. This model consists of a set of equations, one for each variable in the graph. With each random variable on the graph is associated a distinct, arbitrary function, denoted g , and a distinct random disturbance, denoted ε , each with a subscript corresponding to its respective random variable. A component in a graph connected by bidirected edges (i.e., connected when ignoring directed edges) is known as a *district* (Richardson, 2009) or *c-component* (Tian and Pearl, 2002). The sets of random disturbances corresponding to each district are assumed to be mutually independent of one another. That is, $\{\varepsilon_{\mathbf{C}_0}, \varepsilon_{\mathbf{C}_1}, \varepsilon_Y\}$, ε_E , and ε_M are mutually independent; $\varepsilon_{\mathbf{C}_0}$, $\varepsilon_{\mathbf{C}_1}$, and ε_Y , however, are not. Each variable is generated by its corresponding function, which depends only on all variables that directly affect it (i.e., its parents on the graph), and its corresponding random disturbance, as follows: $\mathbf{C}_0 = \mathbf{g}_{\mathbf{C}_0}(\varepsilon_{\mathbf{C}_0})$, $E = g_E(\mathbf{C}_0, \varepsilon_E)$, $\mathbf{C}_1 = \mathbf{g}_{\mathbf{C}_1}(\mathbf{C}_0, E, \varepsilon_{\mathbf{C}_1})$, $M = g_M(\mathbf{C}_0, E, \mathbf{C}_1, \varepsilon_M)$, $Y = g_Y(\mathbf{C}_0, E, \mathbf{C}_1, M, \varepsilon_Y)$.

Just as the Markovian model, the model we introduce is especially useful for making counterfactual independence assumptions explicit. Take for instance the statement $\{Y(m, e'), \mathbf{C}_1(e')\} \perp\!\!\!\perp M(c_1, e) | \mathbf{C}_0$. To see whether this statement holds in the context of the graph in Figure 1, observe what occurs when we intervene on the mechanism in one case to force the exposure to be the comparison level, e , and set \mathbf{C}_1 to an arbitrary value \mathbf{c}_1 : $\mathbf{C}_0 = \mathbf{g}_{\mathbf{C}_0}(\varepsilon_{\mathbf{C}_0})$, $E = e$, $\mathbf{C}_1 = \mathbf{c}_1$, $M(\mathbf{c}_1, e) = g_M(\mathbf{C}_0, e, \mathbf{c}_1, \varepsilon_M)$, $Y(\mathbf{c}_1, e) = g_Y(\mathbf{C}_0, e, \mathbf{c}_1, M(\mathbf{c}_1, e), \varepsilon_Y)$; and in another case to force the exposure to be the reference level, e' , and set M to an arbitrary value m : $\mathbf{C}_0 = \mathbf{g}_{\mathbf{C}_0}(\varepsilon_{\mathbf{C}_0})$, $E = e'$, $\mathbf{C}_1(e') = \mathbf{g}_{\mathbf{C}_1}(\mathbf{C}_0, e', \varepsilon_{\mathbf{C}_1})$, $M = m$, $Y(m, e') = g_Y(\mathbf{C}_0, e', \mathbf{C}_1(e'), m, \varepsilon_Y)$. Note that the only sources of stochasticity in $M(\mathbf{c}_1, e)$ are \mathbf{C}_0 and ε_M , and the only sources of stochasticity in $\{Y(m, e'), \mathbf{C}_1(e')\}$ are \mathbf{C}_0 , $\varepsilon_{\mathbf{C}_1}$, and ε_Y . Hence the only source of dependence

between the two is \mathbf{C}_0 since $\varepsilon_M \perp\!\!\!\perp \{\varepsilon_{\mathbf{C}_1}, \varepsilon_Y\}$, and they are independent conditional on \mathbf{C}_0 . We are now prepared to present our identification result, whose proof is provided in the supplementary materials.

Theorem 1. *Suppose the data-generating mechanism from which the observed data \mathbf{O} are sampled follows the relaxation of the Markovian model that we introduce above, represented by the graph in Figure 1. Then β_0 is identified under this model by the following functional of $F_{\mathbf{O}}$:*

$$\beta_0 = \iiint_{m, \mathbf{c}_1, \mathbf{c}_0} \mathbb{E}(Y|m, \mathbf{c}_1, e', \mathbf{c}_0) dF(m|\mathbf{c}_1, e, \mathbf{c}_0) dF(\mathbf{c}_1|e', \mathbf{c}_0) dF(\mathbf{c}_0). \quad (1)$$

Remark. *The following conditions are sufficient for the same identification result, and are strictly weaker than those implied by our model: for all m, \mathbf{c}_1, e , and e' , $\{Y(m, e'), \mathbf{C}_1(e')\} \perp\!\!\!\perp E|\mathbf{C}_0$, $Y(m) \perp\!\!\!\perp M|\mathbf{C}_1, E, \mathbf{C}_0$, $M(\mathbf{c}_1, e) \perp\!\!\!\perp \{\mathbf{C}_1, E\}|\mathbf{C}_0$, $\{Y(m, e'), \mathbf{C}_1(e')\} \perp\!\!\!\perp M(\mathbf{c}_1, e)|\mathbf{C}_0$.*

Theorem 1, in conjunction with the standard g -formula result $\delta_0 = \int_{\mathbf{c}_0} \mathbb{E}(Y|e', \mathbf{c}_0) dF(\mathbf{c}_0)$ (Robins, 1986), which holds under the assumption encoded on the diagram that $Y(e') \perp\!\!\!\perp E|\mathbf{C}_0$, identifies the \mathcal{P}_{EMY} -specific effect, $\beta_0 - \delta_0$.

4. PATH-SPECIFIC INFERENCE

Thus far, we have only considered a nonparametric model \mathcal{M}_{nonpar} for the observed data, making our identifying functional of the \mathcal{P}_{EMY} -specific effect valid under any possible correct model for the data. Unfortunately, we will seldom have the luxury to continue using \mathcal{M}_{nonpar} through the estimation stage; because inference in \mathcal{M}_{nonpar} is rarely practical in situations with numerous or continuous confounders ($\mathbf{C}_0, \mathbf{C}_1$) (Robins et al., 1997), we will often be forced to posit parametric models. Which models we are to fit depend on how we choose to estimate (1). We now consider four estimators and the corresponding models needed to compute them. Note that, while these estimators are in fact asymptotically equivalent under a nonparametric model, they will have different asymptotic properties under parametric and semiparametric models (Tchetgen Tchetgen and Shpitser, 2012).

4.1 Maximum Likelihood Estimation

We first discuss the maximum likelihood estimator (MLE) for β_0 . By considering the identifying functional (1) as four nested expectations, it is clear that we can fit three appropriate regression models with parameters γ_1 , γ_2 , and γ_3 using maximum likelihood, and plug the predicted means under these models into the functional; the outermost mean can then be estimated empirically. If the conditional mean of Y is taken to be linear in M and \mathbf{C}_1 , and the conditional mean of M is taken as linear in \mathbf{C}_1 , then mean models can be fit for Y , M , and \mathbf{C}_1 . Thus, the MLE is

$$\hat{\beta}_{mle} \equiv \mathbb{P}_n \left\{ \hat{\mathbb{E}}(\hat{\mathbb{E}}(\hat{\mathbb{E}}(Y|M, \mathbf{C}_1, e', \mathbf{C}_0; \hat{\gamma}_1)|\mathbf{C}_1, e, \mathbf{C}_0; \hat{\gamma}_2)|e', \mathbf{C}_0; \hat{\gamma}_3) \right\},$$

where \mathbb{P}_n denotes the empirical mean.

Define $\gamma \equiv (\gamma_1, \gamma_2, \gamma_3)$, $g(\gamma) \equiv \hat{\mathbb{E}}(\hat{\mathbb{E}}(\hat{\mathbb{E}}(Y|M, \mathbf{C}_1, e', \mathbf{C}_0; \gamma_1)|\mathbf{C}_1, e, \mathbf{C}_0; \gamma_2)|e', \mathbf{C}_0; \gamma_3)$, $\mathbf{D}_\gamma \equiv \mathbb{E}[\nabla_\gamma g(\gamma)]$, and $\mathbf{U}(\gamma)$ and $\mathcal{I}(\gamma)$ to be the vector of score equations and block-diagonal matrix of expected informations, respectively, for γ . Let γ_0 be the true value of γ . Then $\hat{\beta}_{mle}$ is asymptotically normal with asymptotic variance equal to $\mathbb{E}[(g(\gamma_0) + \mathbf{D}_{\gamma_0}^T \mathcal{I}(\gamma_0) \mathbf{U}(\gamma_0) - \beta_0)^2]$, which can be estimated empirically, substituting $\hat{\gamma}$ and $\hat{\beta}_{mle}$ for γ_0 and β_0 . The MLE is asymptotically efficient when the three regression models are correctly specified, hence this is the minimum variance achievable by regular, asymptotically linear estimators under the choice of model \mathcal{M}_{par} of \mathbf{O} . $\hat{\beta}_{mle}$ will be consistent only under correct specification of the three models.

4.2 Multiply-Robust Estimation

The multiply-robust (MR) estimator, $\hat{\beta}_{mr}$, comes from an estimating equation involving the efficient influence function of β_0 in the model \mathcal{M}_{nonpar} placing no restriction on the observed data likelihood apart from the positivity assumptions given above. A derivation of this influence function is given in the supplementary materials. In order to express the estimator more succinctly, we introduce additional notation: $B(m, \mathbf{c}_1, e', \mathbf{c}_0) \equiv \mathbb{E}(Y|m, \mathbf{c}_1, e', \mathbf{c}_0)$, $B'(\mathbf{c}_1, e', e, \mathbf{c}_0) \equiv \mathbb{E}\{\mathbb{E}(Y|M, \mathbf{c}_1, e', \mathbf{c}_0)|\mathbf{c}_1, e, \mathbf{c}_0\}$, $B''(e', e, \mathbf{c}_0) \equiv \mathbb{E}[\mathbb{E}\{\mathbb{E}(Y|M, \mathbf{C}_1, e', \mathbf{C}_0)|\mathbf{C}_1, e, \mathbf{C}_0\}|e', \mathbf{C}_0]$, $M^{ratio} \equiv f(M|\mathbf{C}_1, e, \mathbf{C}_0)/f(M|\mathbf{C}_1, e', \mathbf{C}_0)$, and $C_1^{ratio} \equiv f(\mathbf{C}_1|e, \mathbf{C}_0)/f(\mathbf{C}_1|e', \mathbf{C}_0)$. The estimator is

then

$$\begin{aligned} \hat{\beta}_{mr} = & \mathbb{P}_n \left\{ \frac{1_{e'}(E)}{\hat{f}(e'|\mathbf{C}_0)} \hat{M}^{ratio} \left\{ Y - \hat{B}(M, \mathbf{C}_1, e', \mathbf{C}_0) \right\} \right. \\ & + \frac{1_e(E)}{\hat{f}(e|\mathbf{C}_0)} (\hat{C}_1^{ratio})^{-1} \left\{ \hat{B}(M, \mathbf{C}_1, e', \mathbf{C}_0) - \hat{B}'(\mathbf{C}_1, e', e, \mathbf{C}_0) \right\} \\ & \left. + \frac{1_{e'}(E)}{\hat{f}(e'|\mathbf{C}_0)} \left\{ \hat{B}'(\mathbf{C}_1, e', e, \mathbf{C}_0) - \hat{B}''(e', e, \mathbf{C}_0) \right\} + \hat{B}''(e', e, \mathbf{C}_0) \right\}, \end{aligned}$$

where $1_{e^*}(\cdot)$ is the indicator function, $\hat{B}'(\mathbf{C}_1, e', e, \mathbf{C}_0) = \hat{\mathbb{E}}\{\hat{B}(M, \mathbf{C}_1, e', \mathbf{C}_0)|\mathbf{C}_1, e, \mathbf{C}_0\}$, and $\hat{B}''(e', e, \mathbf{C}_0) = \hat{\mathbb{E}}[\hat{\mathbb{E}}\{\hat{B}(M, \mathbf{C}_1, e', \mathbf{C}_0)|\mathbf{C}_1, e, \mathbf{C}_0\}|e', \mathbf{C}_0]$.

Note that the estimator is only a function of estimates of $f_{M|\mathbf{C}_1, E, \mathbf{C}_0}$ and $f_{\mathbf{C}_1|E, \mathbf{C}_0}$ through the ratios M^{ratio} and C_1^{ratio} and mean functions $B'(\mathbf{C}_1, e', e, \mathbf{C}_0)$ and $B''(e', e, \mathbf{C}_0)$. When the mean of Y is linear in M , then $B'(\mathbf{C}_1, e', e, \mathbf{C}_0)$ only depends on the distribution of M through its conditional mean, $\mathbb{E}(M|\mathbf{C}_1, e, \mathbf{C}_0)$. Similarly, if in addition the means of Y and M are both linear in \mathbf{C}_1 , then $B''(e', e, \mathbf{C}_0)$ only depends on the distribution of \mathbf{C}_1 through its conditional mean, $\mathbb{E}(\mathbf{C}_1|e', \mathbf{C}_0)$. We denote $\boldsymbol{\theta}_M \equiv \{B'(\mathbf{C}_1, e', e, \mathbf{C}_0), M^{ratio}\}$ and $\boldsymbol{\theta}_{\mathbf{C}_1} \equiv \{B''(e', e, \mathbf{C}_0), C_1^{ratio}\}$.

B , $\boldsymbol{\theta}_M$, $\boldsymbol{\theta}_{\mathbf{C}_1}$, and $f_{E|\mathbf{C}_0}$ are estimated using low dimensional parametric working models, B^W , $\boldsymbol{\theta}_M^W = \{\mathbb{E}^W[B^W(M, \mathbf{C}_1, e', \mathbf{C}_0)|\mathbf{C}_1, e, \mathbf{C}_0], M^{ratio;W}\}$, $\boldsymbol{\theta}_{\mathbf{C}_1}^W = \{\mathbb{E}^W[B^W(\mathbf{C}_1, e, \mathbf{C}_0)|e', \mathbf{C}_0], C_1^{ratio;W}\}$, and $f_{E|\mathbf{C}_0}^W$, via standard maximum likelihood. Note that we are able to avoid estimating the densities for \mathbf{C}_1 and M by instead estimating their mean functions and density ratios directly. Mean functions can be estimated with standard regression techniques, and density ratios can be estimated using propensity score models since by Bayes' theorem,

$$\frac{f(\mathbf{C}_1|e, \mathbf{C}_0)}{f(\mathbf{C}_1|e', \mathbf{C}_0)} = \frac{f(e|\mathbf{C}_1, \mathbf{C}_0)}{f(e'|\mathbf{C}_1, \mathbf{C}_0)} \times \frac{f(e'|\mathbf{C}_0)}{f(e|\mathbf{C}_0)}$$

and

$$\frac{f(M|e, \mathbf{C}_1, \mathbf{C}_0)}{f(M|e', \mathbf{C}_1, \mathbf{C}_0)} = \frac{f(e|M, \mathbf{C}_1, \mathbf{C}_0)}{f(e'|M, \mathbf{C}_1, \mathbf{C}_0)} \times \frac{f(e'|\mathbf{C}_1, \mathbf{C}_0)}{f(e|\mathbf{C}_1, \mathbf{C}_0)}.$$

An attractive property of the multiply-robust estimator is its robustness to multiple types of potential model misspecification. Let \hat{B} , $\hat{\boldsymbol{\theta}}_M$, $\hat{\boldsymbol{\theta}}_{\mathbf{C}_1}$, and $\hat{f}_{E|\mathbf{C}_0}$ denote estimators of B^W , $\boldsymbol{\theta}_M^W$, $\boldsymbol{\theta}_{\mathbf{C}_1}^W$, and $f_{E|\mathbf{C}_0}^W$ consistent under correct specification. The mean functions in $\boldsymbol{\theta}_M$ and $\boldsymbol{\theta}_{\mathbf{C}_1}$ require correct specification of the functions of M and \mathbf{C}_1 based on the working models for Y and

$\{M, Y\}$, respectively, so that θ_M^W and $\theta_{C_1}^W$ can be correctly specified regardless of whether B^W is, and $\theta_{C_1}^W$ can be correctly specified regardless of whether θ_M^W is. The multiply-robust estimator is consistent and asymptotically normal (under standard regularity conditions) provided that one of the following holds: (a) $\{\theta_M, f_{E|C_0}\} \in \{\theta_M^W, f_{E|C_0}^W\}$, (b) $\{B, \theta_{C_1}, f_{E|C_0}\} \in \{B^W, \theta_{C_1}^W, f_{E|C_0}^W\}$, (c) $\{B, \theta_{C_1}, \theta_M\} \in \{B^W, \theta_{C_1}^W, \theta_M^W\}$. That is, $\hat{\beta}_{mr}$ offers three distinct opportunities to obtain valid inference about the path-specific effect. By contrast, $\hat{\beta}_{mle}$ will be consistent only if a slightly weaker form of (c) holds, where $M^{ratio;W}$ and $C_1^{ratio;W}$ need not be correctly specified.

For inference on $\hat{\beta}_{mr}$, we recommend the nonparametric bootstrap (Efron, 1979) or similar alternative resampling methods such as the wild bootstrap (Mammen, 1993) for nonparametric variance estimation. Due to its reliance on inverse-propensity-score weights, this estimator may suffer from instability in settings where the set of positivity assumptions is nearly violated (Kang and Schafer, 2007). A useful stabilization technique is to simply replace any propensity score $\hat{f}_{E|X}$ with $\hat{f}_{E|X}^\dagger$, where \mathbf{X} is some vector of covariates and $\text{logit} \hat{f}_{E|X}(e|\mathbf{X}) = \text{logit} \hat{f}_{E|X}(e|\mathbf{X}) - \log(1 - \mathbb{P}_n(1_e(E))) + \log(\mathbb{P}_n[1_e(E) \hat{f}_{E|X}(e'|\mathbf{X}) / \hat{f}_{E|X}(e|\mathbf{X})])$, which ensures the weights are bounded as discussed in Tchetgen Tchetgen and Shpitser (2012). An additional stabilization technique is given in the supplementary materials.

4.3 Other Estimators

We consider two additional estimators, both based on alternative representations of (1) as shown in the supplementary materials:

$$\hat{\beta}_a \equiv \mathbb{P}_n \left\{ \frac{1_{e'}(E)}{\hat{f}(e'|\mathbf{C}_0)} \hat{M}^{ratio} Y \right\}$$

$$\hat{\beta}_b \equiv \mathbb{P}_n \left\{ \frac{1_e(E)}{\hat{f}(e|\mathbf{C}_0)} (\hat{C}_1^{ratio})^{-1} \hat{\mathbb{E}}(Y|M, \mathbf{C}_1, e', \mathbf{C}_0) \right\},$$

which again involve plugging in estimated regression models and density curves $\hat{f}(e|M, \mathbf{C}_1, \mathbf{C}_0)$, $\hat{f}(e|\mathbf{C}_1, \mathbf{C}_0)$, and $\hat{f}(e|\mathbf{C}_0)$. Note that $\hat{\beta}_a$ and $\hat{\beta}_b$ depend only on a subset of the models in the multiple-robustness conditions (a) and (b), respectively. It follows that $\hat{\beta}_a$ will generally be consistent only if a slightly weaker form of (a) holds, where $B^{W'}$ need not be correctly specified. Similarly, $\hat{\beta}_b$ will be consistent only if a slightly weaker form of (b) holds, where $B^{W''}$ need not

be correctly specified. In settings with practical violations of positivity, stability of both estimators can be improved using the stabilization technique given in Section 4.2.

5. SIMULATION STUDY

We report results for a simulation study in which we generated 1000 data sets of size 1000 from the following models:

$$\begin{aligned}
 C_0 &\sim \mathcal{U}(0, 2) \\
 E|C_0 &\sim \text{Bernoulli} \left(1 - (1 + \exp(0.9 + 0.3C_0))^{-1} \right) \\
 \mathbf{C}_1 &= \begin{pmatrix} 0.8 \\ 0.6 \\ -0.3 \end{pmatrix} + \begin{pmatrix} 1 \\ 0.1 \\ 0.2 \end{pmatrix} C_0 + \begin{pmatrix} 0.5 \\ -0.4 \\ 0.5 \end{pmatrix} E + \begin{pmatrix} -0.1 \\ 0.8 \\ -0.2 \end{pmatrix} C_0 E + \mathcal{N}(0, I) \\
 M &= -0.5 - 0.2C_0 + 0.3E + [-0.2, 0.1, 0.5]\mathbf{C}_1 + [0.4, 0, 0]E\mathbf{C}_1 + N(0, 1) \\
 Y &= 0.2 + 0.2C_0 + 0.6E + [1, 0.7, 0.3]\mathbf{C}_1 - 0.9M - 0.8EM + N(0, 1).
 \end{aligned}$$

In order to investigate the impact of model misspecification, we computed each of the four estimators given above, $\hat{\beta}_{mr}$, $\hat{\beta}_{mle}$, $\hat{\beta}_a$, and $\hat{\beta}_b$, under the four parametric models, \mathcal{M}_a , \mathcal{M}_b , \mathcal{M}_c , and \mathcal{M}_{int} . Models \mathcal{M}_a , \mathcal{M}_b , and \mathcal{M}_c were specified such that statements (a)-(c) in Section 4.2 corresponding to their respective subscripts held, but the models for the remaining estimands were incorrectly specified. For instance, under \mathcal{M}_a , models θ_M^W and $f_{E|C_0}^W$ are correctly specified, while B^W and $\theta_{C_1}^W$ are not. The intersection model uses correctly-specified working models. All models were fit by maximum likelihood. The stabilization technique described in Section 4.2 was used to adjust propensity scores. We used the following working models, subscripted C for correctly specified and I for incorrectly specified:

$f_{E|C_0}^W$:

Correct: $\text{logit Pr}_C\{E = 1|C_0\} = [1, C_0]\alpha_C$

Incorrect: $\Phi^{-1}(\text{Pr}_I\{E = 1|C_0\}) = [1, C_0]\alpha_I$

B^W :

Correct: $\mathbb{E}_C[Y|M, \mathbf{C}_1, E, C_0] = [1, C_0, E, \mathbf{C}_1, M, EM]\eta_C$

Incorrect: $\mathbb{E}_I[Y|M, \mathbf{C}_1, E, C_0] = [1, C_0, E, \mathbf{C}_1, M]\boldsymbol{\eta}_I$

$\boldsymbol{\theta}_{C_1}^W$:

Correct: $C_1^{ratio;W} = \Pr_C(E = e|\mathbf{C}_1, C_0)/\Pr_C(E = e'|\mathbf{C}_1, C_0) \times \Pr_C(E = e'|C_0)\Pr_C(E = e|C_0)$, which depends on the correctly-specified $f_{E|C_0}^W$ model and the correctly-specified model $\text{logit } \Pr_C\{E = 1|\mathbf{C}_1, C_0\} = [1, C_0, C_0^2, \mathbf{C}_1, C_0\mathbf{C}_1]\boldsymbol{\lambda}_C$;

$B_C''(e', e, C_0) = \mathbb{E}_C\{\mathbb{E}_C\{\mathbb{E}_C(Y|M, \mathbf{C}_1, e', C_0)|\mathbf{C}_1, e, C_0\}|e', C_0\}$, which depends on the correctly-specified B^W model and the correctly-specified models $\mathbb{E}_C[C_{1j}|E, C_0] = [1, C_0, E, C_0E]\boldsymbol{\delta}_{j;C} \forall j \in \{1, 2, 3\}$ and $\mathbb{E}_C[M|\mathbf{C}_1, E, C_0] = [1, C_0, E, \mathbf{C}_1, EC_{11}]\boldsymbol{\zeta}_C$.

Incorrect: $C_1^{ratio;W,I} = \Pr_I(E = e|\mathbf{C}_1, C_0)/\Pr_I(E = e'|\mathbf{C}_1, C_0) \times \Pr_C(E = e'|C_0)/\Pr_C(E = e|C_0)$, which depends on the correctly-specified $f_{E|C_0}^W$ model and the incorrectly-specified model $\text{logit } \Pr_I\{E = 1|\mathbf{C}_1, C_0\} = [1, C_0, \mathbf{C}_1]\boldsymbol{\lambda}_I$;

$B_I''(e', e, C_0) = \mathbb{E}_I\{\mathbb{E}_C\{\mathbb{E}_I(Y|M, \mathbf{C}_1, e', C_0)|\mathbf{C}_1, e, C_0\}|e', C_0\}$, which depends on the incorrectly-specified B^W model, the correctly-specified working mean model for M used for $B_C''(e', e, C_0)$ above, and the incorrectly-specified model $\mathbb{E}_I[C_{1j}|E, C_0] = [1, C_0, E]\boldsymbol{\delta}_{j,I}$, since $\boldsymbol{\theta}_{C_1}^W$ is only misspecified in setting (a), under which B^W is also misspecified and $\boldsymbol{\theta}_M^W$ is correctly specified.

$\boldsymbol{\theta}_M^W$:

Correct: $M^{ratio;W,C} = \Pr_C(E = e|M, \mathbf{C}_1, C_0)/\Pr_C(E = e'|M, \mathbf{C}_1, C_0) \times \Pr_C(E = e'|\mathbf{C}_1, C_0)/\Pr_C(E = e|\mathbf{C}_1, C_0)$, which depends on the correctly-specified model $\text{logit } \Pr_C\{E = 1|M, \mathbf{C}_1, C_0\} = [1, C_0, C_0^2, \mathbf{C}_1, C_0\mathbf{C}_1, C_{11}\mathbf{C}_1, M, C_{11}M]\boldsymbol{\gamma}_C$ and the correctly-specified logistic model used for $C_1^{ratio;W}$ above;

$B_C'(C_1, e', e, C_0) = \mathbb{E}_C\{\mathbb{E}_C(Y|M, \mathbf{C}_1, e', C_0)|\mathbf{C}_1, e, C_0\}$ depends on the correctly-specified B^W model and the correctly-specified mean model for M used for $B_C''(e', e, C_0)$ above.

Incorrect: $M^{ratio;W,I} = \Pr_I(E = e|M, \mathbf{C}_1, C_0)/\Pr_I(E = e'|M, \mathbf{C}_1, C_0) \times \Pr_C(E = e'|\mathbf{C}_1, C_0)/\Pr_C(E = e|\mathbf{C}_1, C_0)$, which depends on the correctly-specified logistic model for $\Pr_C\{E = 1|\mathbf{C}_1, C_0\}$ and the incorrectly-specified model $\text{logit } \Pr_I\{E = 1|M, \mathbf{C}_1, C_0\} = [1, C_0, \mathbf{C}_1, M]\boldsymbol{\gamma}_I$; $B_I'(C_1, e', e, C_0) = \mathbb{E}_I\{\mathbb{E}_C(Y|M, \mathbf{C}_1, e', C_0)|\mathbf{C}_1, e, C_0\}$, which depends on the incorrectly-specified model $\mathbb{E}_I[M|\mathbf{C}_1, E, C_0] = [1, C_0, E, \mathbf{C}_1]\boldsymbol{\zeta}_I$ and the correctly-specified model $B^{W,C}$, since $\boldsymbol{\theta}_M^W$ is only misspecified in setting (c), under which B^W is correctly specified.

The results are summarized in the plot displayed in Figure 2 that shows the four point estimates under each model and their corresponding 95% confidence intervals. The point estimates

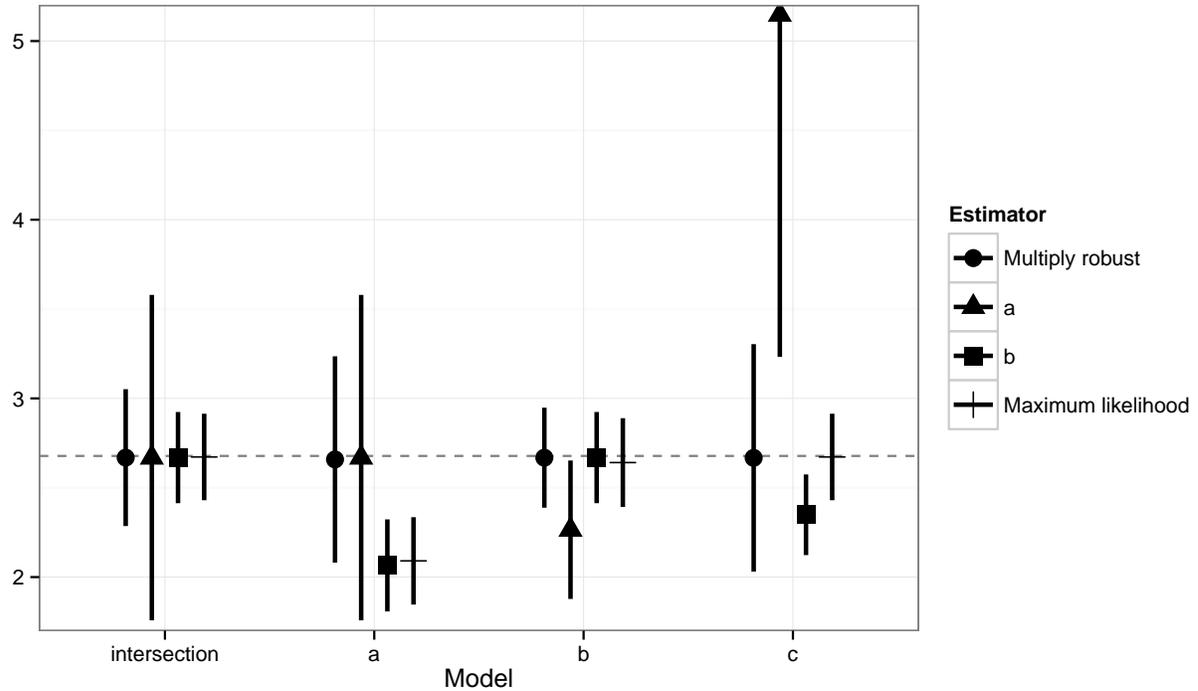


Figure 2: Simulation results for $n=1000$. Monte Carlo point estimates and confidence intervals of each of the four \mathcal{P}_{EMY} -specific estimators are given under \mathcal{M}_{int} , \mathcal{M}_a , \mathcal{M}_b , and \mathcal{M}_c . The horizontal dashed line is through the true parameter value, β_0 .

are the Monte Carlo means of the 1000 samples and the confidence intervals are the values within $t_{999,0.975}$ times the corresponding Monte Carlo standard errors of the point estimates. The confidence intervals correspond to t tests of $H_0 : \hat{\beta} = \beta_0 \equiv 2.678$, hence the confidence intervals not containing β_0 , represented by the horizontal dashed line, correspond to rejection of H_0 .

All estimators are consistent under \mathcal{M}_{int} . Besides $\hat{\beta}_{mr}$, $\hat{\beta}_a$ is the only consistent estimator under \mathcal{M}_a , $\hat{\beta}_b$ is the only consistent estimator under \mathcal{M}_b , and $\hat{\beta}_{mle}$ is the only consistent estimator under \mathcal{M}_c . $\hat{\beta}_{mr}$ is consistent under all models. Therefore, in moderate to large samples, we expect to reject H_0 at the nominal $\alpha = 0.05$ level for none of the estimators under \mathcal{M}_{int} , only for $\hat{\beta}_b$ and $\hat{\beta}_{mle}$ under \mathcal{M}_a , only for $\hat{\beta}_a$ and $\hat{\beta}_{mle}$ under \mathcal{M}_b , and only for $\hat{\beta}_a$ and $\hat{\beta}_b$ under \mathcal{M}_c .

The results illustrate quite well the multiple-robustness property of $\hat{\beta}_{mr}$. As predicted, while the other estimators failed to estimate β_0 without statistically-significant bias, the tests for $\hat{\beta}_{mr}$ failed to reject under every model. For the other estimators, the tests never rejected under \mathcal{M}_{int} and their corresponding models where the misspecified components did not factor into estimation, as expected. That is, the test for $\hat{\beta}_a$ did not reject under \mathcal{M}_a , the test for $\hat{\beta}_b$ did not reject under

\mathcal{M}_b , and the test for $\hat{\beta}_{mle}$ did not reject under \mathcal{M}_c . The tests did reject, however, under the other models, with the exception of $\hat{\beta}_c$ under \mathcal{M}_b . Thus, all estimators other than $\hat{\beta}_{mr}$ were significantly biased under at least one model.

We do see a tradeoff between efficiency and robustness; in all settings, $\hat{\beta}_{mle}$ and $\hat{\beta}_b$ perform best in terms of efficiency, with a slight advantage going to $\hat{\beta}_{mle}$. As such, a reasonable strategy may be to use these estimators in concert initially to diagnose model specification, and then select the most efficient estimator that appears to agree with the multiply-robust estimator, possibly $\hat{\beta}_{mr}$ itself.

Other sample sizes were also explored. At $n = 500$, asymptotic results began to come into focus, though a few of the tests expected to reject looked to be slightly underpowered due to the sample not being large enough. Still, the multiply-robust estimator outperformed the other estimators at this sample size in terms of robustness. At $n = 5000$, confidence intervals were tighter, as expected, but t test results were the same as those at $n = 1000$. The simulation study with $n = 1000$ is comparable to our data analysis in terms of sample size; every treatment comparison consisted of at least 1000 patients.

6. HARVARD PEPFAR NIGERIA ANALYSIS

We now present results of the Harvard PEPFAR data analysis. The data set consisted of 9968 complete observations, i.e., observations with no missing variables, which was 41.9% of the entire data set. We first consider the path-specific effect of treatment regimen assignment on virologic failure through adherence, expressed on the log-risk ratio scale. We used each of the four estimators for β_0 from Section 4, and in each case, δ_0 was estimated using only a subset of the models used to estimate β_0 . In particular, the doubly-robust estimator (Bang and Robins, 2005) was used to estimate δ_0 when contrasted with $\hat{\beta}_b$ and the multiply-robust estimator, $\hat{\beta}_{mr}$; the inverse-probability-weighted estimator (IPW) (Horvitz and Thompson, 1952) was used to estimate δ_0 when contrasted with $\hat{\beta}_a$; and the MLE was used to estimate δ_0 when contrasted with the MLE for the β_0 , $\hat{\beta}_{mle}$. Accordingly, let $\hat{\mathcal{P}}_{EMY;mle}$ denote the effect estimate using $\hat{\beta}_{mle}$, $\hat{\mathcal{P}}_{EMY;a}$ denote the effect estimate using $\hat{\beta}_a$, $\hat{\mathcal{P}}_{EMY;b}$ denote the effect estimate using $\hat{\beta}_b$, and $\hat{\mathcal{P}}_{EMY;mr}$ denote the effect estimate using $\hat{\beta}_{mr}$. We computed all four estimates and corresponding bootstrap confidence intervals for each pairwise comparison of treatments. The wild bootstrap (Mammen,

1993) with weights sampled from $\text{Exp}(1)$ was used to account for instability of resamples due to a small number of cases in some strata of E . Results are summarized in Figure 1 of the supplementary materials.

$\hat{\mathcal{P}}_{EMY;a}$ agreed with $\hat{\mathcal{P}}_{EMY;mr}$ across all comparisons, suggesting that it did not suffer much from model misspecification, or at least any worse than did $\hat{\mathcal{P}}_{EMY;mr}$. It also proved to be the more efficient estimator in this setting, with confidence intervals that were narrower than those of $\hat{\mathcal{P}}_{EMY;mr}$, and comparable to those of $\hat{\mathcal{P}}_{EMY;mle}$, which did not appear to be as robust. Thus, we chose to perform inference using $\hat{\mathcal{P}}_{EMY;a}$ for this portion of the analysis.

Recall that the treatment regimens were coded in descending order of magnitude of their total effects on risk of virologic failure, i.e., they were coded in ascending order of counterfactual risk of virologic failure had everyone been assigned to that treatment, since a lower counterfactual risk of failure corresponds to a higher magnitude of total effect. Because in practice we are more interested in learning how less-effective treatments can be improved, we only consider the higher-coded treatment in a pair as the baseline, e' . Using this ordering, the path-specific effect gives the improvement over the total effect of the less-effective treatment when intervening to make patients adhere as if they were on the more-effective treatment, but had the toxicity and direct effectiveness of the less-effective treatment.

We are primarily interested in the proportion of the total effect attributable to the mediated effect, i.e., the percent mediated by \mathcal{P}_{EMY} . If this proportion is close to or exceeds one, we can conclude that the drugs themselves likely have the same effectiveness on virologic failure, and that it is their differential effect on adherence not due to toxicity that is driving the difference in total effects. If, on the other hand, this proportion is small or negative, we can only say that the difference in total effects is not driven by a difference in effects through \mathcal{P}_{EMY} . It may be the case that the efficacies of the drugs themselves do, in fact, differ, or that the difference in total effects is driven by the differential effect on adherence due to toxicity, but we cannot confirm either. Table 2 shows $\hat{\mathcal{P}}_{EMY;a}$ divided by the total effect estimates, which are also on the log-risk ratio scale and are estimated by IPW. Superscripts indicate the comparisons with significant and marginally-significant path-specific effects. Due to the treatment coding, the denominators of the Table 2 values are always negative. Thus, a negative path-specific effect will be in the same direction as the total effect, and hence will explain a positive proportion of it.

Table 2: Proportion of total effect on virologic failure due to \mathcal{P}_{EMY} -specific effect

Comparison trt	Baseline treatment			
	2	3	4	5
1	0.41 [†]	0.21	-0.059	-0.068*
2	-	0.13	-0.49*	-0.20*
3	-	-	-0.57 [†]	-0.13*
4	-	-	-	-0.027 [†]

NOTE: *Significant path-specific effect ($\alpha = 0.05$). [†]Marginally-significant path-specific effect ($\alpha = 0.1$).

Note that all significant and marginally-significant proportions of total effects due to the effects through \mathcal{P}_{EMY} were negative apart from the one comparing treatment 1 (TDF+3TC/FTC+EFV) to baseline treatment 2 (d4T+3TC+NVP). This occurs when the \mathcal{P}_{EMY} -specific effect estimate is in the opposite direction of the total effect estimate, suggesting that directionally-opposite effects through other pathways overwhelm our estimated effect, and that the total effect would have been even greater if not for the \mathcal{P}_{EMY} -specific effect. For example, had the \mathcal{P}_{EMY} -specific effect been null in the case comparing treatment 3 (AZT+3TC+EFV) with treatment 5 (TDF+3TC/FTC+NVP), we estimate that the total effect would have been 13% larger. The effect of treatment 3 is stronger than 5 not because of its effect through \mathcal{P}_{EMY} , but in spite of it. All differences between treatment 5 and another treatment, and all differences between treatment 4 (AZT+3TC+NVP) and another treatment besides 1 were observed to exhibit this phenomenon as well.

Now consider the exception noted above: the comparison of treatment 1 (TDF+3TC/FTC+EFV) to baseline treatment 2 (d4T+3TC+NVP). We saw a marginally-negative effect, which would have the following interpretation: the effect of treatment 2 on the risk of virologic failure would be improved by patients adhering as if they were assigned to treatment 1, but still had the same toxicity that they did on treatment 2. Unfortunately, treatment 2 is known to have toxicities that were not measured in this data set that are likely to also be affected by underlying biological causes of virologic failure. This interpretation cannot even be considered to be valid for the effect through these unmeasured toxicities, since they induce unmeasured confounding that once again renders this effect unidentifiable. If there were no unmeasured toxicities, we would interpret this

effect as accounting for an estimated 41% of the differences in total effects between treatments 1 and 2.

In conclusion, of the significant \mathcal{P}_{EMY} -specific effects we observed, all apart from those involving unmeasured toxicities were countervailing to the total effect. This means that for these treatment pairs, the differences in their total effects on virologic failure would have been even greater if not for the effect along \mathcal{P}_{EMY} . Thus, the effect through \mathcal{P}_{EMY} does not explain the differential effects on virologic failure, and in some cases actually works against them. As mentioned above, the differential effects may instead be due to the drugs themselves differing in efficacy, or they may be driven by the differential effects on adherence due to toxicity, but such hypotheses require further investigation beyond the scope of our analysis.

We now consider the path-specific effect of treatment regimen assignment on log CD4 count, expressed on the mean difference scale. We again analyzed the four estimators given in Section 4. This time $\hat{\mathcal{P}}_{EMY;a}$ and $\hat{\mathcal{P}}_{EMY;b}$ were drastically less efficient than $\hat{\mathcal{P}}_{EMY;mle}$ and $\hat{\mathcal{P}}_{EMY;mr}$. One possible explanation for this is that the density of log CD4 count was less concentrated around zero, making $\hat{\mathcal{P}}_{EMY;a}$ and $\hat{\mathcal{P}}_{EMY;b}$ more sensitive to small weights. $\hat{\mathcal{P}}_{EMY;mle}$ disagreed with $\hat{\mathcal{P}}_{EMY;mr}$ on several occasions, so $\hat{\mathcal{P}}_{EMY;mr}$ was the best choice in terms of achieving both robustness and efficiency. It is worth noting that the linear outcome model for CD4 count did not seem to suffer too much from misspecification, while the logistic outcome model for virologic failure did. Results are summarized in Figure 2 of the supplementary materials.

Table 3 shows $\hat{\mathcal{P}}_{EMY;mr}$ divided by the total effect estimates, which are also on the log-risk ratio scale and are estimated using doubly-robust estimators. As before, we are interested in

Table 3: Proportion of total effect on CD4 count due to \mathcal{P}_{EMY} -specific effect

Comparison trt	Baseline treatment			
	3	1	5	2
4	0.48*	0.095	-0.045 [†]	0.036
3	-	-0.47	-0.13	0.030
1	-	-	-0.080	0.062
5	-	-	-	0.099

NOTE: *Significant path-specific effect ($\alpha = 0.05$). [†]Marginally-significant path-specific effect ($\alpha = 0.1$).

learning how the less-effective treatment can be improved, but now less-effective is in terms of CD4 count. Since the order of the effectiveness of the treatments for CD4 count is not the same as the order for virologic failure, the treatments which should be considered the comparison versus baseline level in a pair no longer correspond to the treatment coding. The order of the treatments in the margins of Table 3 is rearranged to reflect this different ordering of effectiveness. The denominator for each of the values in the table is positive, since a higher counterfactual CD4 count corresponds to a higher magnitude of total effect. Therefore, positive proportions correspond to positive path-specific effects, and negative proportions correspond to countervailing path-specific effects.

The path-specific effect was found to be significant for only one of the pairwise comparisons: treatment 4 (AZT+3TC+NVP) vs. treatment 3 (AZT+3TC+EFV). This effect is estimated to be in the positive direction, therefore we conclude that the effect of treatment 3 on CD4 count would be improved by patients adhering as if they were assigned to treatment 4 but without necessarily altering toxicity experienced under treatment 3 that they did on treatment 3. The effect through this pathway accounted for almost half of the total effect at an estimated 48%. Thus, if one were interested in improving the effect of AZT+3TC+EFV on CD4 count, it would be worthwhile to examine what mechanisms other than toxicity may be implicated in differential adherence rates between these two regimens. The \mathcal{P}_{EMY} -specific effect comparing treatments 4 and 5 (TDF+3TC/FTC+NVP) was found to be marginally-significantly less than zero. Thus, the difference in total effects of these two treatments is not attributable to their differential effect on adherence not due to toxicity, as the effect through this pathway was in fact in the opposite direction. Rather, this difference was due to differential effects through other pathways as previously described.

7. DISCUSSION

In the PEPFAR case study, we observed an interesting trend of countervailing effects along \mathcal{P}_{EMY} to the total effects on virologic failure for most treatment comparisons, meaning that the differences in the total effects of treatment assignment would have been even greater if not for the effects along \mathcal{P}_{EMY} . While this does not help explain why the treatment assignment effects are different (or at least different in the direction that we observe), it does suggest a method for im-

proving the regimens that we observed to have greater effects on virologic failure. For a treatment comparison with a significant \mathcal{P}_{EMY} -specific effect, if we could identify what is different about the more effective drug regimen that is causing people to not adhere to it as well, then we could potentially eliminate this mechanism in order to reduce the countervailing \mathcal{P}_{EMY} -specific effect and consequently improve its total effect on virologic failure.

A countervailing \mathcal{P}_{EMY} -specific effect on CD4 count was also observed between AZT+3TC+NVP and TDF+3TC/FTC+NVP, which has the same interpretation as the countervailing effects on virologic failure. On the other hand, almost half of the difference in the effects of AZT+3TC+EFV and AZT+3TC+NVP on CD4 count was found to be attributable to the effect through adherence, but not toxicity. This suggests that the effect of AZT+3TC+EFV on CD4 count could be improved up to that of AZT+3TC+NVP if one could identify and eliminate the mechanisms driving the difference in these treatments' effects on adherence. In the other treatment comparisons, none of the differences in total effects on CD4 count were found to be attributable to an effect through \mathcal{P}_{EMY} . Overall, we have achieved an enhanced understanding of the role of adherence in the effects of the five ART regimens considered on both virologic failure and CD4 count.

The most significant methodologic contribution of this paper is the extension of mediation analysis methods to settings in which the NDE and NIE may not be identified, viz. settings with unmeasured confounding and exposure-induced confounding of the mediator. We present conditions under which the \mathcal{P}_{EMY} -specific effect is nonparametrically identified as well as four estimators, including an efficient estimator that is multiply robust to model misspecification for settings where nonparametric estimation is not feasible.

Often effects of adherence are evaluated regarding the treatment assignment as an instrumental variable, relying on an assumption of no direct effect of assignment with respect to adherence. Furthermore, instrumental variable methods rely on an assumption of monotonicity in the effect of assignment on adherence. However, neither of these assumptions are reasonable in our setting where we are forced to compare treatments head-to-head rather than to a control exposure level.

This paper suffers from a few limitations. One is that our identifiability assumptions, though weaker than those of the Markovian model, are still untestable as stated. When possible, we can embed our mediation problem in a larger model represented by a larger graph where treatments

can be split into a component corresponding to the EMY pathway and a component corresponding to all other pathways. This can provide a testable reformulation of identifying assumptions, as was done in Robins and Richardson (2010) in simpler mediation contexts. Another limitation is that this method is not yet equipped to handle missing data. As such, only a complete-case analysis was conducted for the HIV data, allowing for the possibility of bias due to informative missingness. Additionally, for both virologic failure and CD4 count outcomes, it is possible that we are underestimating the effect of substantive interest if adherence over the first six months plays a large mediating role since we are forced to control for early adherence and can only estimate the effect through adherence over the second six months. Finally, not a limitation, but rather a caveat, is that the \mathcal{P}_{EMY} -specific effect is not a substitute for the NIE. The NIE is not fully captured by this effect and, in fact, even if the effects along both \mathcal{P}_{EMY} and $E \rightarrow C_1 \rightarrow M \rightarrow Y$ are in the same direction, the NIE does not necessarily have to be. Strong assumptions are needed to draw this conclusion. As such, while often practically meaningful, the \mathcal{P}_{EMY} -specific effect must be interpreted with care and not blindly substituted for the NIE.

Future directions for this work would, of course, include adjusting the method to account for missing data, which could improve the analysis conducted in this paper. Another important extension would be to the full longitudinal case, with repeated exposures, mediators, and confounders. Shpitser (2013) gives the identifying functional for the analog to the \mathcal{P}_{EMY} -specific effect in this setting, but no estimation strategy exists as of yet. Finally, it is not uncommon for a mediator to be measured with error, which tends to induce bias as shown by VanderWeele et al. (2012). It would be valuable to adapt the methods of Tchetgen Tchetgen and Lin (2012) for handling this problem to our setting. Alternatively, parametric approaches have been suggested (Valeri et al., 2014) that could also potentially be adapted.

8. SUPPLEMENTARY MATERIALS

All supplementary materials are contained in a single archive and can be obtained via a single download.

Proofs and theoretical results: We prove Theorem 1, derive the four estimation strategies, derive the efficient influence function of β_0 , and prove its robustness. (PDF file)

Additional stabilization technique for the multiply-robust estimator: We present an additional method to account for instability of the multiply-robust estimator due to near-positivity violations. (PDF file)

Plots comparing estimators in the PEPFAR Nigeria study: Plots summarizing the estimation results for the \mathcal{P}_{EMY} -specific effect on virologic failure and CD4 count using the four estimators. These plots were used to assess model misspecification and select the most appropriate estimator. (PDF file)



A. PROOFS AND THEORETICAL RESULTS

A.1 Identification Result Proof

Proof of Theorem 1.

$$\begin{aligned}
 \beta_0 &\equiv \mathbb{E}[Y(M(\mathbf{C}_1(e'), e), \mathbf{C}_1(e'), e')] \\
 &= \int_{\mathbf{c}_0, \mathbf{c}_1, m, y} y dF_{Y(M(\mathbf{C}_1(e'), e), \mathbf{C}_1(e'), e'), M(\mathbf{C}_1(e'), e), \mathbf{C}_1(e'), \mathbf{C}_0}(y, m, \mathbf{c}_1, \mathbf{c}_0) \\
 &= \iint_{\mathbf{c}_0, \mathbf{c}_1, m, y} y dF_{Y(m, e'), M(\mathbf{c}_1, e), \mathbf{C}_1(e') | \mathbf{C}_0}(y, m, \mathbf{c}_1 | \mathbf{c}_0) dF_{\mathbf{C}_0}(\mathbf{c}_0) \\
 &= \iiint_{\mathbf{c}_0, \mathbf{c}_1, m, y} y dF_{Y(m, e'), \mathbf{C}_1(e') | \mathbf{C}_0}(y, \mathbf{c}_1 | \mathbf{c}_0) dF_{M(\mathbf{c}_1, e) | \mathbf{C}_0}(m | \mathbf{c}_0) dF_{\mathbf{C}_0}(\mathbf{c}_0) \tag{2}
 \end{aligned}$$

$$= \iiint_{\mathbf{c}_0, \mathbf{c}_1, m, y} y dF_{Y(m, e'), \mathbf{C}_1(e') | E, \mathbf{C}_0}(y, \mathbf{c}_1 | e', \mathbf{c}_0) dF_{M(\mathbf{c}_1, e) | \mathbf{C}_0}(m | \mathbf{c}_0) dF_{\mathbf{C}_0}(\mathbf{c}_0) \tag{3}$$

$$= \iiint_{\mathbf{c}_0, \mathbf{c}_1, m, y} y dF_{Y(m), \mathbf{C}_1 | E, \mathbf{C}_0}(y, \mathbf{c}_1 | e', \mathbf{c}_0) dF_{M(\mathbf{c}_1, e) | \mathbf{C}_0}(m | \mathbf{c}_0) dF_{\mathbf{C}_0}(\mathbf{c}_0) \tag{4}$$

$$\begin{aligned}
 &= \iiint_{\mathbf{c}_0, \mathbf{c}_1, m, y} y dF_{Y(m) | \mathbf{C}_1, E, \mathbf{C}_0}(y | \mathbf{c}_1, e', \mathbf{c}_0) dF_{M(\mathbf{c}_1, e) | \mathbf{C}_0}(m | \mathbf{c}_0) dF_{\mathbf{C}_1 | E, \mathbf{C}_0}(\mathbf{c}_1 | e', \mathbf{c}_0) dF_{\mathbf{C}_0}(\mathbf{c}_0) \\
 &= \iiint_{\mathbf{c}_0, \mathbf{c}_1, m, y} y dF_{Y(m) | M, \mathbf{C}_1, E, \mathbf{C}_0}(y | m, \mathbf{c}_1, e', \mathbf{c}_0) dF_{M(\mathbf{c}_1, e) | \mathbf{C}_1, E, \mathbf{C}_0}(m | \mathbf{c}_1, e, \mathbf{c}_0) \\
 &\hspace{20em} \times dF_{\mathbf{C}_1 | E, \mathbf{C}_0}(\mathbf{c}_1 | e', \mathbf{c}_0) dF_{\mathbf{C}_0}(\mathbf{c}_0) \tag{5}
 \end{aligned}$$

$$\begin{aligned}
 &= \iiint_{\mathbf{c}_0, \mathbf{c}_1, m, y} y dF_{Y | M, \mathbf{C}_1, E, \mathbf{C}_0}(y | m, \mathbf{c}_1, e', \mathbf{c}_0) dF_{M | \mathbf{C}_1, E, \mathbf{C}_0}(m | \mathbf{c}_1, e, \mathbf{c}_0) \\
 &\hspace{20em} \times dF_{\mathbf{C}_1 | E, \mathbf{C}_0}(\mathbf{c}_1 | e', \mathbf{c}_0) dF_{\mathbf{C}_0}(\mathbf{c}_0), \tag{6}
 \end{aligned}$$

where (2) follows from $\{Y(m, e'), \mathbf{C}_1(e')\} \perp\!\!\!\perp M(\mathbf{c}_1, e) | \mathbf{C}_0$, (3) follows from $\{Y(m, e'), \mathbf{C}_1(e')\} \perp\!\!\!\perp E | \mathbf{C}_0$, (4) follows by consistency, (5) follows from $Y(m) \perp\!\!\!\perp M | \mathbf{C}_1, E, \mathbf{C}_0$ and $M(\mathbf{c}_1, e) \perp\!\!\!\perp \{\mathbf{C}_1, E\} | \mathbf{C}_0$, and (6) follows by consistency. \square

A.2 Derivation of Estimation Strategies

A.2.1 Maximum Likelihood Estimator

The maximum likelihood estimator arises from the alternative representation of (1):

$$\begin{aligned} & \iiint_{m, \mathbf{c}_1, \mathbf{c}_0} \mathbb{E}(Y|m, \mathbf{c}_1, e', \mathbf{c}_0) dF_{M|\mathbf{C}_1, E, \mathbf{C}_0}(m|\mathbf{c}_1, e, \mathbf{c}_0) dF_{\mathbf{C}_1|E, \mathbf{C}_0}(\mathbf{c}_1|e', \mathbf{c}_0) dF_{\mathbf{C}_0}(\mathbf{c}_0) \\ &= \mathbb{E}(\mathbb{E}(\mathbb{E}(\mathbb{E}(Y|M, \mathbf{C}_1, e', \mathbf{C}_0)|\mathbf{C}_1, e, \mathbf{C}_0)|e', \mathbf{C}_0)). \end{aligned}$$

We replace the inner three expectations with their arguments' means under the empirical laws $\hat{f}_{\mathbf{C}_1|e', \mathbf{C}_0}$, $\hat{f}_{M|\mathbf{C}_1, e, \mathbf{C}_0}$, and $\hat{f}_{Y|M, \mathbf{C}_1, e', \mathbf{C}_1}$ respectively, and compute the empirical mean. Thus, we have

$$\hat{\beta}_{mle} \equiv \mathbb{P}_n \left\{ \hat{\mathbb{E}}(\hat{\mathbb{E}}(\hat{\mathbb{E}}(Y|M, \mathbf{C}_1, e', \mathbf{C}_0)|\mathbf{C}_1, e, \mathbf{C}_0)|e', \mathbf{C}_0) \right\}.$$

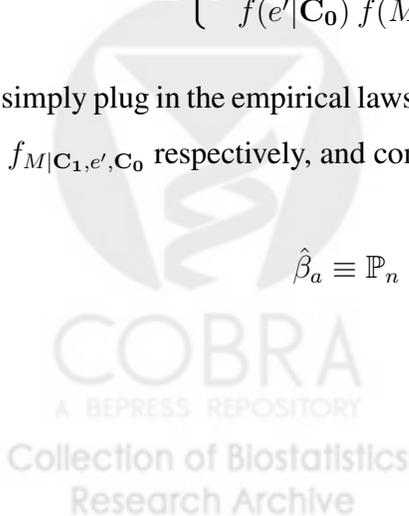
A.2.2 Estimator a

$\hat{\beta}_a$ arises from another alternative representation of (1):

$$\begin{aligned} & \iiint_{m, \mathbf{c}_1, \mathbf{c}_0} \mathbb{E}(Y|m, \mathbf{c}_1, e', \mathbf{c}_0) dF_{M|\mathbf{C}_1, E, \mathbf{C}_0}(m|\mathbf{c}_1, e, \mathbf{c}_0) dF_{\mathbf{C}_1|E, \mathbf{C}_0}(\mathbf{c}_1|e', \mathbf{c}_0) dF_{\mathbf{C}_0}(\mathbf{c}_0) \\ &= \sum_{e^* \in \{e', e\}} \int_{y, m, \mathbf{c}_1, \mathbf{c}_0} y \frac{1_{e'}(e^*)}{f(e'|\mathbf{C}_0)} \frac{f(m|\mathbf{c}_1, e, \mathbf{c}_0)}{f(m|\mathbf{c}_1, e^*, \mathbf{c}_0)} dF_{Y, M, \mathbf{C}_1, E, \mathbf{C}_0}(y, m, \mathbf{c}_1, e^*, \mathbf{c}_0) \\ &= \mathbb{E} \left\{ Y \frac{1_{e'}(E)}{f(e'|\mathbf{C}_0)} \frac{f(M|\mathbf{C}_1, e, \mathbf{C}_0)}{f(M|\mathbf{C}_1, e', \mathbf{C}_0)} \right\}. \end{aligned}$$

We simply plug in the empirical laws, $\hat{f}_{E=0|\mathbf{C}_0}$, $\hat{f}_{M|\mathbf{C}_1, e, \mathbf{C}_0}$, and $\hat{f}_{M|\mathbf{C}_1, e', \mathbf{C}_0}$ for $f_{E=0|\mathbf{C}_0}$, $f_{M|\mathbf{C}_1, e, \mathbf{C}_0}$, and $f_{M|\mathbf{C}_1, e', \mathbf{C}_0}$ respectively, and compute the empirical mean. Thus, we have

$$\hat{\beta}_a \equiv \mathbb{P}_n \left\{ Y \frac{1_{e'}(E)}{\hat{f}(e'|\mathbf{C}_0)} \frac{\hat{f}(M|\mathbf{C}_1, e, \mathbf{C}_0)}{\hat{f}(M|\mathbf{C}_1, e', \mathbf{C}_0)} \right\}.$$



A.2.3 Estimator b

$\hat{\beta}_b$ arises from a third representation of (1):

$$\begin{aligned} & \int \int \int_{m, \mathbf{c}_1, \mathbf{c}_0} \mathbb{E}(Y|m, \mathbf{c}_1, e', \mathbf{c}_0) dF_{M|\mathbf{C}_1, E, \mathbf{C}_0}(m|\mathbf{c}_1, e, \mathbf{c}_0) dF_{\mathbf{C}_1|E, \mathbf{C}_0}(\mathbf{c}_1|e', \mathbf{c}_0) dF_{\mathbf{C}_0}(\mathbf{c}_0) \\ &= \sum_{e^* \in \{e', e\}} \int_{m, \mathbf{c}_1, \mathbf{c}_0} \mathbb{E}(Y|M, \mathbf{C}_1, e', \mathbf{C}_0) \frac{1_e(e^*)}{f(e^*|\mathbf{C}_0)} \frac{f(\mathbf{c}_1|e', \mathbf{c}_0)}{f(\mathbf{c}_1|e^*, \mathbf{c}_0)} dF_{M, \mathbf{C}_1, E, \mathbf{C}_0}(m, \mathbf{c}_1, e^*, \mathbf{c}_0) \\ &= \mathbb{E} \left[\frac{1_e(E)}{f(e|\mathbf{C}_0)} \frac{f(\mathbf{C}_1|e', \mathbf{C}_0)}{f(\mathbf{C}_1|e, \mathbf{C}_0)} \mathbb{E}(Y|M, \mathbf{C}_1, e', \mathbf{C}_0) \right]. \end{aligned}$$

Again, we plug in the empirical laws $\hat{f}_{\mathbf{C}_1|e', \mathbf{C}_0}$, $\hat{f}_{\mathbf{C}_1|e, \mathbf{C}_0}$, and $\hat{f}_{E=1|\mathbf{C}_0}$ for $f_{\mathbf{C}_1|e', \mathbf{C}_0}$, $f_{\mathbf{C}_1|e, \mathbf{C}_0}$, and $f_{E=1|\mathbf{C}_0}$, respectively, replace $\mathbb{E}(Y|M, \mathbf{C}_1, e', \mathbf{C}_0)$ with $\hat{\mathbb{E}}(Y|M, \mathbf{C}_1, e', \mathbf{C}_0)$, the expectation of Y under the empirical law $\hat{f}_{Y|M, \mathbf{C}_1, e', \mathbf{C}_0}$, and compute the empirical mean. Thus, we have

$$\hat{\beta}_b \equiv \mathbb{P}_n \left\{ \frac{1_e(E)}{\hat{f}(e|\mathbf{C}_0)} \frac{\hat{f}(\mathbf{C}_1|e', \mathbf{C}_0)}{\hat{f}(\mathbf{C}_1|e, \mathbf{C}_0)} \hat{\mathbb{E}}(Y|M, \mathbf{C}_1, e', \mathbf{C}_0) \right\}.$$

We develop the multiply-robust estimator and prove its robustness properties in the following two sections.

A.3 Derivation of the Influence Function

Theorem 2. *The efficient influence function of β_0 in model \mathcal{M}_{nonpar} is given by*

$$\begin{aligned} V^{eff}(\beta_0) &= \frac{1_{e'}(E) f(M|e, \mathbf{C}_1, \mathbf{C}_0)}{f(M|e', \mathbf{C}_1, \mathbf{C}_0) f(e'|\mathbf{C}_0)} \{Y - B(M, \mathbf{C}_1, e', \mathbf{C}_0)\} \\ &\quad + \frac{1_e(E) f(\mathbf{C}_1|e', \mathbf{C}_0)}{f(\mathbf{C}_1|e, \mathbf{C}_0) f(e|\mathbf{C}_0)} \{B(M, \mathbf{C}_1, e', \mathbf{C}_0) - B'(\mathbf{C}_1, e', e, \mathbf{C}_0)\} \\ &\quad + \frac{1_{e'}(E)}{f(e'|\mathbf{C}_0)} \{B'(\mathbf{C}_1, e', e, \mathbf{C}_0) - B''(e', e, \mathbf{C}_0)\} + \{B''(e', e, \mathbf{C}_0) - \beta_0\}, \end{aligned}$$

implying that the asymptotic variance of a regular, asymptotically linear (RAL) estimator of β_0 in model \mathcal{M}_{nonpar} can be no smaller than $\mathbb{E}\{V^{eff}(\beta_0)^2\}^{-1}$, the semiparametric efficiency bound for the model.

$\hat{\beta}_{mr}$ is obtained simply by solving the estimating equation $V^{eff}(\beta_0)$ for β_0 . Since our model is nonparametric, the asymptotic variance is the same for any estimator in \mathcal{M}_{nonpar} so long as it is RAL. Furthermore, since all such estimators share the common influence function $V^{eff}(\beta_0)$, they also share a common asymptotic expansion, viz. $n^{1/2}(\hat{\beta}_0 - \beta_0) = n^{1/2}\mathbb{P}_n V^{eff}(\beta_0) + o_p(1)$, where \mathbb{P}_n denotes the empirical mean.

Proof. Let ν denote the appropriate dominating measure or product measure corresponding to each combination of random variables. Let $F_{\mathbf{O};t} = F_{Y|M,C,E,\mathbf{C}_0;t}F_{M|\mathbf{C}_1,E,\mathbf{C}_0;t}F_{\mathbf{C}_1|E,\mathbf{C}_0;t}F_{E|\mathbf{C}_0;t}F_{\mathbf{C}_0;t}$ denote a one-dimensional regular parametric submodel of \mathcal{M}_{nonpar} with $F_{\mathbf{O},0} = F_{\mathbf{O}}$, and let

$$\begin{aligned}\beta_t &= \beta_0(F_{\mathbf{O};t}) = \mathbb{E}_t(Y(M(e, \mathbf{C}_1(e')), \mathbf{C}_1(e'), e')) \\ &= \int_{m, \mathbf{c}_1, \mathbf{c}_0} \mathbb{E}_t(Y|m, \mathbf{c}_1, e', \mathbf{c}_0) f_t(M = m|\mathbf{c}_1, e, \mathbf{c}_0) f_t(\mathbf{C}_1 = \mathbf{c}_1|e', \mathbf{c}_0) f_t(\mathbf{C}_0 = \mathbf{c}_0) d\nu(m, \mathbf{c}_1, \mathbf{c}_0)\end{aligned}$$

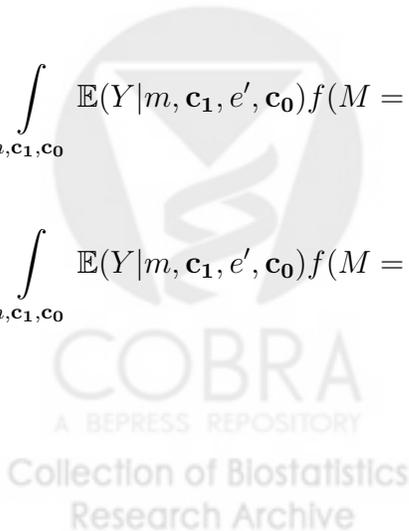
and $U_{\mathbf{O}} = \frac{\nabla_{t=0} f_t(\mathbf{O})}{f(\mathbf{O})}$ be the score for \mathbf{O} . Then

$$\begin{aligned}\frac{\partial \beta_t}{\partial t} \Big|_{t=0} &= \\ &\int_{m, \mathbf{c}_1, \mathbf{c}_0} \nabla_{t=0} \mathbb{E}_t(Y|m, \mathbf{c}_1, e', \mathbf{c}_0) f(M = m|\mathbf{c}_1, e, \mathbf{c}_0) f(\mathbf{C}_1 = \mathbf{c}_1|e', \mathbf{c}_0) f(\mathbf{C}_0 = \mathbf{c}_0) d\nu(m, \mathbf{c}_1, \mathbf{c}_0)\end{aligned}\tag{7}$$

$$+ \int_{m, \mathbf{c}_1, \mathbf{c}_0} \mathbb{E}(Y|m, \mathbf{c}_1, e', \mathbf{c}_0) \nabla_{t=0} f_t(M = m|\mathbf{c}_1, e, \mathbf{c}_0) f(\mathbf{C}_1 = \mathbf{c}_1|e', \mathbf{c}_0) f(\mathbf{C}_0 = \mathbf{c}_0) d\nu(m, \mathbf{c}_1, \mathbf{c}_0)\tag{8}$$

$$+ \int_{m, \mathbf{c}_1, \mathbf{c}_0} \mathbb{E}(Y|m, \mathbf{c}_1, e', \mathbf{c}_0) f(M = m|\mathbf{c}_1, e, \mathbf{c}_0) \nabla_{t=0} f_t(\mathbf{C}_1 = \mathbf{c}_1|e', \mathbf{c}_0) f(\mathbf{C}_0 = \mathbf{c}_0) d\nu(m, \mathbf{c}_1, \mathbf{c}_0)\tag{9}$$

$$+ \int_{m, \mathbf{c}_1, \mathbf{c}_0} \mathbb{E}(Y|m, \mathbf{c}_1, e', \mathbf{c}_0) f(M = m|\mathbf{c}_1, e, \mathbf{c}_0) f(\mathbf{C}_1 = \mathbf{c}_1|e', \mathbf{c}_0) \nabla_{t=0} f_t(\mathbf{C}_0 = \mathbf{c}_0) d\nu(m, \mathbf{c}_1, \mathbf{c}_0),\tag{10}$$



where

$$\begin{aligned}
(7) &= \int_{m, \mathbf{c}_1, \mathbf{c}_0} \nabla_{t=0} \mathbb{E}_t(Y|m, \mathbf{c}_1, e', \mathbf{c}_0) f(M = m|\mathbf{c}_1, e, \mathbf{c}_0) f(\mathbf{C}_1 = \mathbf{c}_1|e', \mathbf{c}_0) f(\mathbf{C}_0 = \mathbf{c}_0) d\nu(m, \mathbf{c}_1, \mathbf{c}_0) \\
&= \int_{m, \mathbf{c}_1, \mathbf{c}_0} \int_y \left\{ \frac{\nabla_{t=0} f_t(y, m, \mathbf{c}_1, e', \mathbf{c}_0)}{f(m, \mathbf{c}_1, e', \mathbf{c}_0)} - \frac{f(y, m, \mathbf{c}_1, e', \mathbf{c}_0) \nabla_{t=0} f_t(m, \mathbf{c}_1, e', \mathbf{c}_0)}{f(m, \mathbf{c}_1, e', \mathbf{c}_0)^2} \right\} d\nu(y) \\
&\quad \times f(M = m|\mathbf{c}_1, e, \mathbf{c}_0) f(\mathbf{C}_1 = \mathbf{c}_1|e', \mathbf{c}_0) f(\mathbf{C}_0 = \mathbf{c}_0) d\nu(m, \mathbf{c}_1, \mathbf{c}_0) \\
&= \int_{y, m, \mathbf{c}_1, \mathbf{c}_0} \left\{ y \frac{\nabla_{t=0} f_t(y, m, \mathbf{c}_1, e', \mathbf{c}_0)}{f(m, \mathbf{c}_1, e', \mathbf{c}_0)} - \frac{\nabla_{t=0} f_t(m, \mathbf{c}_1, e', \mathbf{c}_0)}{f(m, \mathbf{c}_1, e', \mathbf{c}_0)} \mathbb{E}(Y|m, \mathbf{c}_1, e', \mathbf{c}_0) f(y|m, \mathbf{c}_1, e', \mathbf{c}_0) \right\} \\
&\quad \times f(M = m|\mathbf{c}_1, e, \mathbf{c}_0) f(\mathbf{C}_1 = \mathbf{c}_1|e', \mathbf{c}_0) f(\mathbf{C}_0 = \mathbf{c}_0) d\nu(y, m, \mathbf{c}_1, \mathbf{c}_0) \\
&= \int_{y, m, \mathbf{c}_1, e^*, \mathbf{c}_0} \left\{ y \frac{\nabla_{t=0} f_t(y, m, \mathbf{c}_1, e^*, \mathbf{c}_0)}{f(m, \mathbf{c}_1, e', \mathbf{c}_0)} - \frac{\nabla_{t=0} f_t(m, \mathbf{c}_1, e^*, \mathbf{c}_0)}{f(m, \mathbf{c}_1, e', \mathbf{c}_0)} \mathbb{E}(Y|m, \mathbf{c}_1, e', \mathbf{c}_0) f(y|m, \mathbf{c}_1, e', \mathbf{c}_0) \right\} \\
&\quad \times 1_{e'}(e^*) f(M = m|\mathbf{c}_1, e, \mathbf{c}_0) f(\mathbf{C}_1 = \mathbf{c}_1|e', \mathbf{c}_0) f(\mathbf{C}_0 = \mathbf{c}_0) d\nu(y, m, \mathbf{c}_1, e^*, \mathbf{c}_0) \\
&= \mathbb{E} \left[\frac{1_{e'}(E) f(M|\mathbf{C}_1, e', \mathbf{C}_0) f(\mathbf{C}_1|e', \mathbf{C}_0) f(\mathbf{C}_0)}{f(Y, M, \mathbf{C}_1, E, \mathbf{C}_0)} \right. \\
&\quad \times \left. \left\{ Y \frac{\nabla_{t=0} f_t(Y, M, \mathbf{C}_1, E, \mathbf{C}_0)}{f(M, \mathbf{C}_1, e', \mathbf{C}_0)} - f(Y|M, \mathbf{C}_1, e', \mathbf{C}_0) \frac{\nabla_{t=0} f_t(M, \mathbf{C}_1, E, \mathbf{C}_0)}{f(M, \mathbf{C}_1, e', \mathbf{C}_0)} B(M, \mathbf{C}_1, e', \mathbf{C}_0) \right\} \right] \\
&= \mathbb{E} \left[\frac{1_{e'}(E) f(M|\mathbf{C}_1, e, \mathbf{C}_0) f(\mathbf{C}_1|e', \mathbf{C}_0) f(\mathbf{C}_0)}{f(Y, M, \mathbf{C}_1, E, \mathbf{C}_0) f(M, \mathbf{C}_1, e', \mathbf{C}_0)} \{Y \nabla_{t=0} f_t(Y, M, \mathbf{C}_1, E, \mathbf{C}_0) \right. \\
&\quad \left. - [\nabla_{t=0} f_t(Y, M, \mathbf{C}_1, E, \mathbf{C}_0) - f(M, \mathbf{C}_1, e', \mathbf{C}_0) \nabla_{t=0} f_t(Y|M, \mathbf{C}_1, E, \mathbf{C}_0)] B(M, \mathbf{C}_1, e', \mathbf{C}_0) \right] \\
&= \mathbb{E} \left[\frac{\nabla_{t=0} f_t(Y, M, \mathbf{C}_1, E, \mathbf{C}_0)}{f(Y, M, \mathbf{C}_1, E, \mathbf{C}_0)} \times \frac{1_{e'}(E) f(M|\mathbf{C}_1, e, \mathbf{C}_0)}{f(M|\mathbf{C}_1, e', \mathbf{C}_0) f(E = e'|\mathbf{C}_0)} \{Y - B(M, \mathbf{C}_1, e', \mathbf{C}_0)\} \right] \\
&\quad + \int_{m, \mathbf{c}_1, \mathbf{c}_0} f(\mathbf{c}_1|e', \mathbf{c}_0) f(m|\mathbf{c}_1, e, \mathbf{c}_0) f(\mathbf{c}_0) \nabla_{t=0} \left\{ \int_y f_t(y|m, \mathbf{c}_1, e', \mathbf{c}_0) d\nu(y) \right\} \\
&\quad \times B(m, \mathbf{c}_1, e', \mathbf{c}_0) d\nu(m, \mathbf{c}_1, \mathbf{c}_0) \\
&= \mathbb{E} \left[U_0 \frac{1_{e'}(E) f(M|\mathbf{C}_1, e, \mathbf{C}_0)}{f(M|\mathbf{C}_1, e', \mathbf{C}_0) f(E = e'|\mathbf{C}_0)} \{Y - B(M, \mathbf{C}_1, e', \mathbf{C}_0)\} \right], \\
(8) &= \int_{m, \mathbf{c}_1, \mathbf{c}_0} \mathbb{E}(Y|m, \mathbf{c}_1, e', \mathbf{c}_0) \left\{ \frac{\nabla_{t=0} f_t(m, \mathbf{c}_1, e, \mathbf{c}_0)}{f(\mathbf{c}_1, e, \mathbf{c}_0)} - \frac{\nabla_{t=0} f_t(\mathbf{c}_1, e, \mathbf{c}_0) f(m, \mathbf{c}_1, e, \mathbf{c}_0)}{f(\mathbf{c}_1, e, \mathbf{c}_0)^2} \right\} \\
&\quad \times f(\mathbf{c}_1|e', \mathbf{c}_0) f(\mathbf{c}_0) d\nu(m, \mathbf{c}_1, \mathbf{c}_0)
\end{aligned}$$

$$\begin{aligned}
&= \int_{m, \mathbf{c}_1, \mathbf{c}_0} \mathbb{E}(Y|m, \mathbf{c}_1, e', \mathbf{c}_0) \frac{\nabla_{t=0} f_t(m, \mathbf{c}_1, e, \mathbf{c}_0)}{f(\mathbf{c}_1, e, \mathbf{c}_0)} f(\mathbf{c}_1|e', \mathbf{c}_0) f(\mathbf{c}_0) d\nu(m, \mathbf{c}_1, \mathbf{c}_0) \\
&\quad - \int_{\mathbf{c}_1, \mathbf{c}_0} \frac{\nabla_{t=0} f_t(\mathbf{c}_1, e, \mathbf{c}_0)}{f(\mathbf{c}_1, e, \mathbf{c}_0)} \mathbb{E}(\mathbb{E}(Y|m, \mathbf{c}_1, e', \mathbf{c}_0)|\mathbf{c}_1, e, \mathbf{c}_0) f(\mathbf{c}_1|e', \mathbf{c}_0) f(\mathbf{c}_0) d\nu(\mathbf{c}_1, \mathbf{c}_0) \\
&= \int_{m, \mathbf{c}_1, \mathbf{c}_0} \frac{f(\mathbf{c}_1|e', \mathbf{c}_0) f(\mathbf{c}_0)}{f(\mathbf{c}_1, e, \mathbf{c}_0)} \left\{ \nabla_{t=0} f_t(m, \mathbf{c}_1, e, \mathbf{c}_0) \mathbb{E}(Y|m, \mathbf{c}_1, e', \mathbf{c}_0) \right. \\
&\quad \left. - \nabla_{t=0} f_t(\mathbf{c}_1, e, \mathbf{c}_0) f(m|\mathbf{c}_1, e, \mathbf{c}_0) B'(\mathbf{c}_1, e', e, \mathbf{c}_0) \right\} d\nu(m, \mathbf{c}_1, \mathbf{c}_0) \\
&= \int_{m, \mathbf{c}_1, \mathbf{c}_0} \frac{f(\mathbf{c}_1|e', \mathbf{c}_0)}{f(\mathbf{c}_1|e, \mathbf{c}_0) f(e|\mathbf{c}_0)} \left\{ \nabla_{t=0} f_t(m, \mathbf{c}_1, e, \mathbf{c}_0) \mathbb{E}(Y|m, \mathbf{c}_1, e', \mathbf{c}_0) \right. \\
&\quad \left. - [\nabla_{t=0} f_t(m, \mathbf{c}_1, e, \mathbf{c}_0) - f(\mathbf{c}_1, e, \mathbf{c}_0) \nabla_{t=0} f_t(m|\mathbf{c}_1, e, \mathbf{c}_0)] B'(\mathbf{c}_1, e', e, \mathbf{c}_0) \right\} d\nu(m, \mathbf{c}_1, \mathbf{c}_0) \\
&= \int_{m, \mathbf{c}_1, \mathbf{c}_0} \nabla_{t=0} f_t(m, \mathbf{c}_1, e, \mathbf{c}_0) \frac{f(\mathbf{c}_1|e', \mathbf{c}_0)}{f(\mathbf{c}_1|e, \mathbf{c}_0) f(e|\mathbf{c}_0)} \left\{ E(Y|m, \mathbf{c}_1, e', \mathbf{c}_0) - B'(\mathbf{c}_1, e', e, \mathbf{c}_0) \right\} d\nu(m, \mathbf{c}_1, \mathbf{c}_0) \\
&\quad + \int_{\mathbf{c}_1, \mathbf{c}_0} f(\mathbf{c}_1|e', \mathbf{c}_0) f(\mathbf{c}_0) \nabla_{t=0} \int_m f_t(m|\mathbf{c}_1, e, \mathbf{c}_0) d\nu(m) B'(\mathbf{c}_1, e', e, \mathbf{c}_0) d\nu(\mathbf{c}_1, \mathbf{c}_0) \\
&= \int_{m, \mathbf{c}_1, \mathbf{c}_0} \left\{ \int_y f(y|m, \mathbf{c}_1, e, \mathbf{c}_0) d\nu(y) \nabla_{t=0} f_t(m, \mathbf{c}_1, e, \mathbf{c}_0) \right. \\
&\quad \left. + \nabla_{t=0} \int_y f_t(y|m, \mathbf{c}_1, e, \mathbf{c}_0) d\nu(y) f(m, \mathbf{c}_1, e, \mathbf{c}_0) \right\} \frac{f(\mathbf{c}_1|e', \mathbf{c}_0)}{f(\mathbf{c}_1|e, \mathbf{c}_0) f(e|\mathbf{c}_0)} \\
&\quad \times \left\{ \mathbb{E}(Y|m, \mathbf{c}_1, e', \mathbf{c}_0) - B'(\mathbf{c}_1, e', e, \mathbf{c}_0) \right\} d\nu(m, \mathbf{c}_1, \mathbf{c}_0) \\
&= \int_{y, m, \mathbf{c}_1, \mathbf{c}_0} \nabla_{t=0} f_t(y, m, \mathbf{c}_1, e, \mathbf{c}_0) \frac{f(\mathbf{c}_1|e', \mathbf{c}_0)}{f(\mathbf{c}_1|e, \mathbf{c}_0) f(e|\mathbf{c}_0)} \\
&\quad \times \left\{ \mathbb{E}(Y|m, \mathbf{c}_1, e', \mathbf{c}_0) - B'(\mathbf{c}_1, e', e, \mathbf{c}_0) \right\} d\nu(y, m, \mathbf{c}_1, \mathbf{c}_0) \\
&= \int_{y, m, \mathbf{c}_1, e^*, \mathbf{c}_0} \nabla_{t=0} f_t(y, m, \mathbf{c}_1, e^*, \mathbf{c}_0) \frac{1_e(e^*) f(\mathbf{c}_1|e', \mathbf{c}_0)}{f(\mathbf{c}_1|e, \mathbf{c}_0) f(e|\mathbf{c}_0)} \\
&\quad \times \left\{ \mathbb{E}(Y|m, \mathbf{c}_1, e', \mathbf{c}_0) - B'(\mathbf{c}_1, e', e, \mathbf{c}_0) \right\} d\nu(y, m, \mathbf{c}_1, e^*, \mathbf{c}_0) \\
&= \mathbb{E} \left[U_0 \frac{1_e(E) f(\mathbf{C}_1|e', \mathbf{C}_0)}{f(\mathbf{C}_1|e, \mathbf{C}_0) f(e|\mathbf{C}_0)} \left\{ \mathbb{E}(Y|M, \mathbf{C}_1, e', \mathbf{C}_0) - B'(\mathbf{C}_1, e', e, \mathbf{C}_0) \right\} \right], \\
(9) &= \int_{m, \mathbf{c}_1, \mathbf{c}_0} \mathbb{E}(Y|m, \mathbf{c}_1, e', \mathbf{c}_0) f(m|\mathbf{c}_1, e, \mathbf{c}_0) \left\{ \frac{\nabla_{t=0} f_t(\mathbf{c}_1, e', \mathbf{c}_0)}{f(e', \mathbf{c}_0)} - \frac{\nabla_{t=0} f_t(e', \mathbf{c}_0) f(\mathbf{c}_1, e', \mathbf{c}_0)}{f(e', \mathbf{c}_0)^2} \right\}
\end{aligned}$$

$$\begin{aligned}
& \times f(\mathbf{c}_0) d\nu(m, \mathbf{c}_1, \mathbf{c}_0) \\
&= \int_{\mathbf{c}_1, \mathbf{c}_0} \mathbb{E}(\mathbb{E}(Y|M, \mathbf{C}_1, e', \mathbf{C}_0) | \mathbf{c}_1, e, \mathbf{c}_0) \left\{ \frac{\nabla_{t=0} f_t(\mathbf{c}_1, e', \mathbf{c}_0)}{f(e', \mathbf{c}_0)} - \frac{\nabla_{t=0} f_t(e', \mathbf{c}_0)}{f(e', \mathbf{c}_0)} f(\mathbf{c}_1 | e', \mathbf{c}_0) \right\} \\
& \quad \times f(\mathbf{c}_0) d\nu(\mathbf{c}_1, \mathbf{c}_0) \\
&= \int_{\mathbf{c}_1, \mathbf{c}_0} B'(\mathbf{c}_1, e', e, \mathbf{c}_0) \frac{\nabla_{t=0} f_t(\mathbf{c}_1, e', \mathbf{c}_0)}{f(e', \mathbf{c}_0)} f(\mathbf{c}_0) d\nu(\mathbf{c}_1, \mathbf{c}_0) \\
& \quad - \int_{\mathbf{c}_0} \mathbb{E}(\mathbb{E}(\mathbb{E}(Y|M, \mathbf{C}_1, e', \mathbf{C}_0) | \mathbf{C}_1, e, \mathbf{C}_0) | e', \mathbf{c}_0) \frac{\nabla_{t=0} f_t(e', \mathbf{c}_0)}{f(e', \mathbf{c}_0)} f(\mathbf{c}_0) d\nu(\mathbf{c}_0) \\
&= \int_{\mathbf{c}_1, \mathbf{c}_0} \frac{f(\mathbf{c}_0)}{f(e', \mathbf{c}_0)} \{ \nabla_{t=0} f_t(\mathbf{c}_1, e', \mathbf{c}_0) B'(\mathbf{c}_1, e', e, \mathbf{c}_0) \\
& \quad - \nabla_{t=0} f_t(e', \mathbf{c}_0) f(\mathbf{c}_1 | e', \mathbf{c}_0) B''(e', e, \mathbf{c}_0) \} d\nu(\mathbf{c}_1, \mathbf{c}_0) \\
&= \int_{\mathbf{c}_1, \mathbf{c}_0} \frac{1}{f(e' | \mathbf{c}_0)} \{ \nabla_{t=0} f_t(\mathbf{c}_1, e', \mathbf{c}_0) B'(\mathbf{c}_1, e', e, \mathbf{c}_0) \\
& \quad - [\nabla_{t=0} f_t(\mathbf{c}_1, e', \mathbf{c}_0) - \nabla_{t=0} f_t(\mathbf{c}_1 | e', \mathbf{c}_0) f(e', \mathbf{c}_0)] B''(e', e, \mathbf{c}_0) \} d\nu(\mathbf{c}_1, \mathbf{c}_0) \\
&= \int_{\mathbf{c}_1, \mathbf{c}_0} \frac{1}{f(e' | \mathbf{c}_0)} \nabla_{t=0} f_t(\mathbf{c}_1, e', \mathbf{c}_0) \{ B'(\mathbf{c}_1, e', e, \mathbf{c}_0) - B''(e', e, \mathbf{c}_0) \} d\nu(\mathbf{c}_1, \mathbf{c}_0) \\
& \quad + \int_{\mathbf{c}_0} f(\mathbf{c}_0) \nabla_{t=0} \int_{\mathbf{c}_1} f_t(\mathbf{c}_1 | e', \mathbf{c}_0) d\nu(\mathbf{c}_1) B''(e', e, \mathbf{c}_0) d\nu(\mathbf{c}_0) \\
&= \int_{\mathbf{c}_1, \mathbf{c}_0} \frac{1}{f(e' | \mathbf{c}_0)} \left\{ \int_{y, m} f(y, m | \mathbf{c}_1, e', \mathbf{c}_0) d\nu(y, m) \nabla_{t=0} f_t(\mathbf{c}_1, e', \mathbf{c}_0) \right. \\
& \quad \left. + \nabla_{t=0} \int_{y, m} f_t(y, m | \mathbf{c}_1, e', \mathbf{c}_0) d\nu(y, m) f(\mathbf{c}_1, e', \mathbf{c}_0) \right\} \{ B'(\mathbf{c}_1, e', e, \mathbf{c}_0) - B''(e', e, \mathbf{c}_0) \} d\nu(\mathbf{c}_1, \mathbf{c}_0) \\
&= \int_{y, m, \mathbf{c}_1, \mathbf{c}_0} \frac{\nabla_{t=0} f_t(y, m, \mathbf{c}_1, e', \mathbf{c}_0)}{f(e' | \mathbf{c}_0)} \{ B'(\mathbf{c}_1, e', e, \mathbf{c}_0) - B''(e', e, \mathbf{c}_0) \} d\nu(y, m, \mathbf{c}_1, \mathbf{c}_0) \\
&= \int_{y, m, \mathbf{c}_1, e^*, \mathbf{c}_0} \nabla_{t=0} f_t(y, m, \mathbf{c}_1, e^*, \mathbf{c}_0) \frac{1_{e'}(e^*)}{f(e' | \mathbf{c}_0)} \{ B'(\mathbf{c}_1, e', e, \mathbf{c}_0) - B''(e', e, \mathbf{c}_0) \} d\nu(y, m, \mathbf{c}_1, e^*, \mathbf{c}_0) \\
&= \mathbb{E} \left[U_{\mathbf{O}} \frac{1_e(E')}{f(e' | \mathbf{C}_0)} \{ B'(\mathbf{c}_1, e', e, \mathbf{C}_0) - B''(e', e, \mathbf{C}_0) \} \right],
\end{aligned}$$

and

$$\begin{aligned}
(10) &= \int_{\mathbf{c}_0} \mathbb{E}(\mathbb{E}(\mathbb{E}(Y|M, \mathbf{C}_1, e', \mathbf{C}_0)|\mathbf{C}_1, e, \mathbf{C}_0)|e', \mathbf{C}_0) \nabla_{t=0} f_t(\mathbf{c}_0) d\nu(\mathbf{c}_0) - \beta_0 \mathbb{E}U_{\mathbf{O}} \\
&= \int_{\mathbf{c}_0} \left\{ \int_{y, m, \mathbf{c}_1, e^*} f(y, m, \mathbf{c}_1, e^*|\mathbf{c}_0) d\nu(y, m, \mathbf{c}_1, e^*) \nabla_{t=0} f_t(\mathbf{c}_0) \right. \\
&\quad \left. + \nabla_{t=0} \int_{y, m, \mathbf{c}_1, e^*} f_t(y, m, \mathbf{c}_1, e^*|\mathbf{c}_0) d\nu(y, m, \mathbf{c}_1, e^*) f(\mathbf{c}_0) \right\} B''(e', e, \mathbf{c}_0) d\nu(\mathbf{c}_0) - \mathbb{E}[U_{\mathbf{O}}\beta_0] \\
&= \int_{y, m, \mathbf{c}_1, e^*, \mathbf{c}_0} \nabla_{t=0} f_t(y, m, \mathbf{c}_1, e^*, \mathbf{c}_0) B''(e', e, \mathbf{c}_0) d\nu(y, m, \mathbf{c}_1, e^*, \mathbf{c}_0) - \mathbb{E}[U_{\mathbf{O}}\beta_0] \\
&= \mathbb{E}[U_{\mathbf{O}} \{B''(e', e, \mathbf{C}_0) - \beta_0\}].
\end{aligned}$$

Thus, $\frac{\partial \beta_t}{\partial t} \Big|_{t=0} = \mathbb{E}[U_{\mathbf{O}} V^{eff}(\beta_0)]$ where

$$\begin{aligned}
V^{eff}(\beta_0) &= \frac{1_{e'}(E)f(M|e, \mathbf{C}_1, \mathbf{C}_0)}{f(M|e', \mathbf{C}_1, \mathbf{C}_0)f(e'|\mathbf{C}_0)} \{Y - B(M, \mathbf{C}_1, e', \mathbf{C}_0)\} \\
&\quad + \frac{1_e(E)f(\mathbf{C}_1|e', \mathbf{C}_0)}{f(\mathbf{C}_1|e, \mathbf{C}_0)f(e|\mathbf{C}_0)} \{B(M, \mathbf{C}_1, e', \mathbf{C}_0) - B'(\mathbf{C}_1, e', e, \mathbf{C}_0)\} \\
&\quad + \frac{1_{e'}(E)}{f(e'|\mathbf{C}_0)} \{B'(\mathbf{C}_1, e', e, \mathbf{C}_0) - B''(e', e, \mathbf{C}_0)\} + \{B''(e', e, \mathbf{C}_0) - \beta_0\},
\end{aligned}$$

so if a RAL estimator exists, then $V^{eff}(\beta_0)$ is the corresponding influence function. It is efficient because the model \mathcal{M}_{nonpar} is nonparametric. \square

A.4 Multiple-Robustness of the Efficient Influence Function

Let $\tilde{B}, \tilde{\theta}_M = \{\tilde{M}^{ratio}, \tilde{\mathbb{E}}[\tilde{B}(M, \mathbf{C}_1, e', \mathbf{C}_0)|\mathbf{C}_1, e, \mathbf{C}_0]\}$, $\tilde{\theta}_{\mathbf{C}_1} = \{\tilde{C}_1^{ratio}, \tilde{\mathbb{E}}[\tilde{B}'(\mathbf{C}_1, e, \mathbf{c}_0)|e', \mathbf{C}_0]\}$, and $\tilde{f}_{E|\mathbf{C}_0}$ denote limits of the estimators using the working models $B^W, \theta_M^W, \theta_{\mathbf{C}_1}^W$, and $f_{E|\mathbf{C}_0}^W$. We have established the following multiply-robust property of V^{eff} :

Theorem 3. *The estimating equation $V^{eff}(\beta_0, \tilde{B}, \tilde{\theta}_M, \tilde{\theta}_{\mathbf{C}_1}, \tilde{f}_{E|\mathbf{C}_0})$ is unbiased provided that one of the following holds:*

$$(a) \{\tilde{\theta}_M, \tilde{f}_{E|\mathbf{C}_0}\} = \{\theta_M, f_{E|\mathbf{C}_0}\},$$

- (b) $\{\tilde{B}, \tilde{\theta}_{C_1}, \tilde{f}_{E|C_0}\} = \{B, \theta_{C_1}, f_{E|C_0}\}$, or
(c) $\{\tilde{B}, \tilde{\theta}_{C_1}, \tilde{\theta}_M\} = \{B, \theta_{C_1}, \theta_M\}$.

Proof.

$$\begin{aligned} \mathbb{E}V^{eff}(\beta_0, \tilde{B}, \tilde{\theta}_M, \tilde{\theta}_{C_1}, \tilde{f}_{E|C_0}) = & \\ \mathbb{E} \left[\int_{m, \mathbf{c}_1} \frac{\tilde{M}^{ratio}}{\tilde{f}(e'|C_0)} \left\{ B(m, \mathbf{c}_1, e', C_0) - \tilde{B}(m, \mathbf{c}_1, e', C_0) \right\} \right. & \\ & \times f(m|\mathbf{c}_1, e', C_0) f(\mathbf{c}_1|e', C_0) f(e'|C_0) d\nu(m, \mathbf{c}_1) \\ & + \int_{\mathbf{c}_1} \frac{1}{\tilde{C}_1^{ratio} \tilde{f}(e|C_0)} \left\{ \mathbb{E} \left[\tilde{B}(M, \mathbf{c}_1, e', C_0) | \mathbf{c}_1, e, C_0 \right] - \tilde{\mathbb{E}} \left[\tilde{B}(M, \mathbf{c}_1, e', C_0) | \mathbf{c}_1, e, C_0 \right] \right\} \\ & \times f(\mathbf{c}_1|e, C_0) f(e|C_0) d\nu(\mathbf{c}_1) \\ & + \frac{f(e'|C_0)}{\tilde{f}(e'|C_0)} \left\{ \mathbb{E} \left[\tilde{\mathbb{E}} \left[\tilde{B}(M, C_1, e', C_0) | C_1, e, C_0 \right] | e', C_0 \right] \right. \\ & \quad \left. - \tilde{\mathbb{E}} \left[\tilde{\mathbb{E}} \left[\tilde{B}(M, C_1, e', C_0) | C_1, e, C_0 \right] | e', C_0 \right] \right\} \\ & \left. + \tilde{\mathbb{E}} \left[\tilde{\mathbb{E}} \left[\tilde{B}(M, C_1, e', C_0) | C_1, e, C_0 \right] | e', C_0 \right] - \mathbb{E} \left[\mathbb{E} \left[B(M, C_1, e', C_0) | C_1, e, C_0 \right] | e', C_0 \right] \right] \end{aligned}$$

Substituting under (a):

$$\begin{aligned} \mathbb{E}V^{eff}(\beta_0, \tilde{B}, \tilde{\theta}_M, \tilde{\theta}_{C_1}, \tilde{f}_{E|C_0}) = & \\ \mathbb{E} \left[\int_{m, \mathbf{c}_1} \left\{ B(m, \mathbf{c}_1, e', C_0) - \tilde{B}(m, \mathbf{c}_1, e', C_0) \right\} f(m|\mathbf{c}_1, e, C_0) f(\mathbf{c}_1|e', C_0) d\nu(m, \mathbf{c}_1) \right. & \\ & + \left\{ \mathbb{E} \left[\mathbb{E} \left[\tilde{B}(M, C_1, e', C_0) | C_1, e, C_0 \right] | e', C_0 \right] \right. \\ & \quad \left. - \tilde{\mathbb{E}} \left[\mathbb{E} \left[\tilde{B}(M, C_1, e', C_0) | C_1, e, C_0 \right] | e', C_0 \right] \right\} \\ & \left. + \tilde{\mathbb{E}} \left[\mathbb{E} \left[\tilde{B}(M, C_1, e', C_0) | C_1, e, C_0 \right] | e', C_0 \right] - \mathbb{E} \left[\mathbb{E} \left[B(M, C_1, e', C_0) | C_1, e, C_0 \right] | e', C_0 \right] \right] \end{aligned}$$

= 0

COBRA
 A BEPRESS REPOSITORY
 Collection of Biostatistics
 Research Archive

Substituting under (b):

$$\begin{aligned} \mathbb{E}V^{eff}(\beta_0, \tilde{B}, \tilde{\theta}_M, \tilde{\theta}_{C_1}, \tilde{f}_{E|C_0}) &= \\ & \int_{\mathbf{c}_1} \left\{ \mathbb{E}[B(M, \mathbf{c}_1, e', \mathbf{C}_0) | \mathbf{c}_1, e, \mathbf{C}_0] - \tilde{\mathbb{E}}[B(M, \mathbf{c}_1, e', \mathbf{C}_0) | \mathbf{c}_1, e, \mathbf{C}_0] \right\} f(\mathbf{c}_1 | e', \mathbf{C}_0) d\nu(\mathbf{c}_1) \\ & + \mathbb{E} \left[\tilde{\mathbb{E}}[B(M, \mathbf{C}_1, e', \mathbf{C}_0) | \mathbf{C}_1, e, \mathbf{C}_0] | e', \mathbf{C}_0 \right] - \mathbb{E}[\mathbb{E}[B(M, \mathbf{C}_1, e', \mathbf{C}_0) | \mathbf{C}_1, e, \mathbf{C}_0] | e', \mathbf{C}_0] \Big] \\ & = 0 \end{aligned}$$

Substituting under (c):

$$\mathbb{E}V^{eff}(\beta_0, \tilde{B}, \tilde{\theta}_M, \tilde{\theta}_{C_1}, \tilde{f}_{E|C_0}) = 0, \text{ trivially.} \quad \square$$

Thus, $\hat{\beta}_{mr}$ can be shown to be asymptotically normal under each of these scenarios using a Taylor expansion of $\mathbb{P}_n V^{eff}(\hat{\beta}_{mr}, \hat{B}, \hat{\theta}_M, \hat{\theta}_{C_1}, \hat{f}_{E|C_0})$ and applying the central limit theorem to $n^{-1/2} \sum_i V_i^{eff}(\beta_0, B^*, \theta_M^*, \theta_{C_1}^*, f_{E|C_0}^*)$.



B. ADDITIONAL STABILIZATION TECHNIQUE FOR THE MULTIPLY-ROBUST ESTIMATOR

This technique is an adaptation of the approach presented by Robins et al. (2007). The idea is to carefully select regression models and an estimation strategy such that the three terms in $\hat{\beta}_{mr}$ depending on weights are empirically evaluated as null, leaving the term $\hat{B}''(e', e, C_0)$, which does not depend on weights. This can be accomplished with the following steps. First, fit propensity score models to estimate $f_{E|C_0}$, M^{ratio} , and C_1^{ratio} . Substitute these estimates into the first term of $\hat{\beta}_{mr}$, and include the result in a set of estimating equation to solve for the Y -regression-model parameters. Next, plug in all parameters estimated thus far into the second term of $\hat{\beta}_{mr}$, and once again use the result in a set of estimating equations to solve for the M -regression-model parameters. Repeat this step with the third term of $\hat{\beta}_{mr}$ to solve for the C_1 -regression-model parameters. Finally, plugging all of these parameter estimates into $\hat{\beta}_{mr}$ leaves $\hat{B}''(e', e, C_0)$, as desired. If Y , M , and C_1 are all scalar, continuous random variables, this procedure is equivalent to repeatedly fitting regression models with intercepts using weighted least squares with appropriately-chosen weights.



C. PLOTS COMPARING ESTIMATORS IN THE PEPFAR NIGERIA STUDY

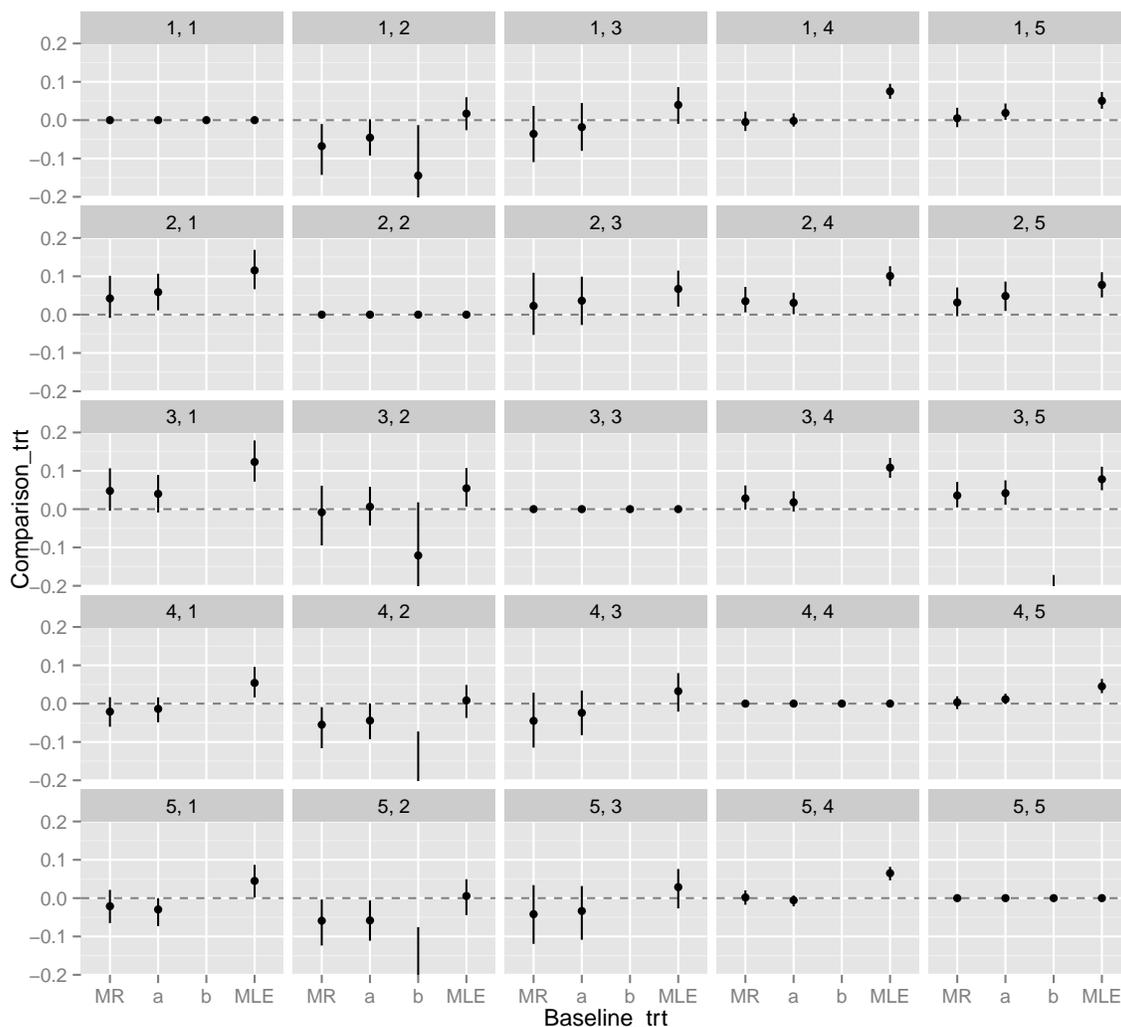
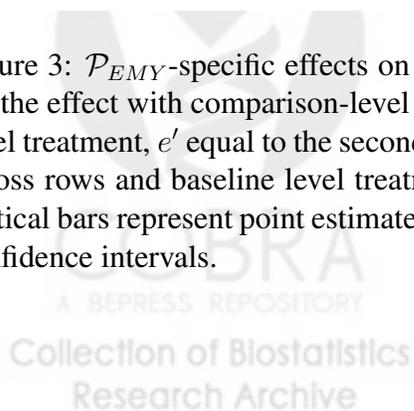


Figure 3: \mathcal{P}_{EMY} -specific effects on virologic failure. The plot in each cell represents estimates for the effect with comparison-level treatment, e , equal to the first index of the cell and baseline-level treatment, e' equal to the second index of the cell. That is, comparison level treatment varies across rows and baseline level treatment varies across columns. Within each plot, the dots and vertical bars represent point estimates using the four estimators and their corresponding bootstrap confidence intervals.



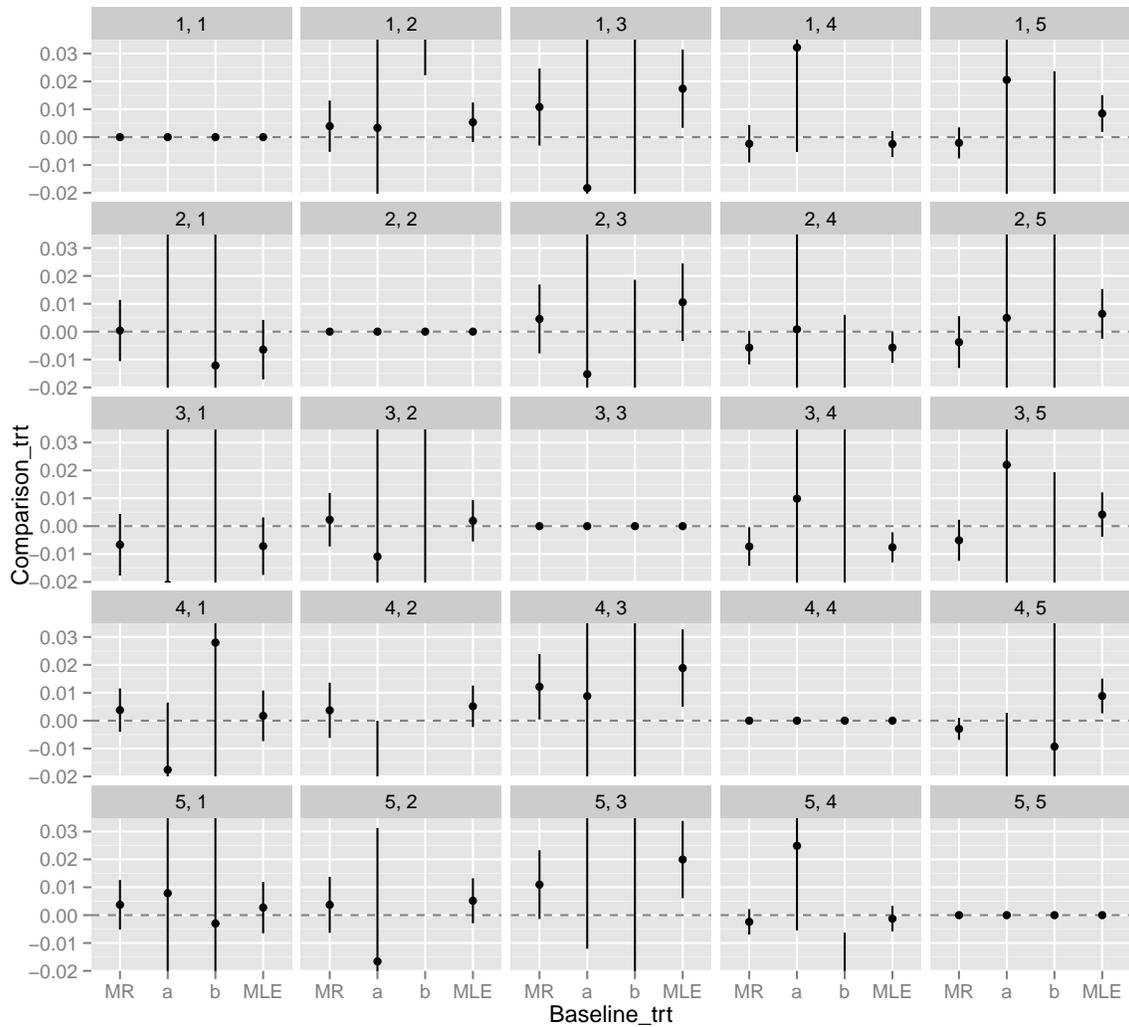
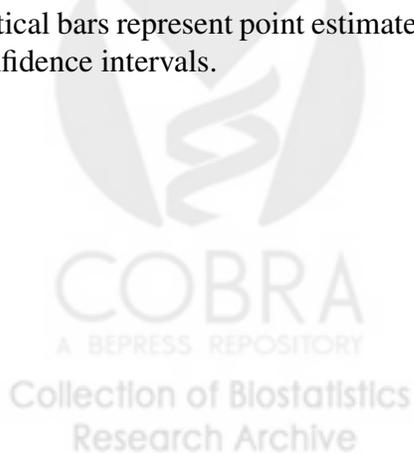


Figure 4: \mathcal{P}_{EMY} -specific effects on CD4 count. The plot in each cell represents estimates for the effect with comparison-level treatment, e , equal to the first index of the cell and baseline-level treatment, e' equal to the second index of the cell. That is, comparison level treatment varies across rows and baseline level treatment varies across columns. Within each plot, the dots and vertical bars represent point estimates using the four estimators and their corresponding bootstrap confidence intervals.



REFERENCES

- Avin, C., Shpitser, I., and Pearl, J. (2005). Identifiability of path-specific effects. In *IJCAI-05, Proceedings of the Nineteenth International Joint Conference on Artificial Intelligence*, pages 357–363.
- Bang, H. and Robins, J. M. (2005). Doubly robust estimation in missing data and causal inference models. *Biometrics*, 61(4):962–973.
- Bangsberg, D. R., Hecht, F. M., Charlebois, E. D., Zolopa, A. R., Holodniy, M., Sheiner, L., Bamberger, J. D., Chesney, M. A., and Moss, A. (2000). Adherence to protease inhibitors, HIV-1 viral load, and development of drug resistance in an indigent population. *AIDS*, 14(4):357–366.
- Baron, R. M. and Kenny, D. A. (1986). The moderator–mediator variable distinction in social psychological research: Conceptual, strategic, and statistical considerations. *Journal of Personality and Social Psychology*, 51(6):1173.
- China Tuberculosis Control Collaboration (1996). Results of directly observed short-course chemotherapy in 112 842 Chinese patients with smear-positive tuberculosis. *The Lancet*, 347(8998):358–362.
- Efron, B. (1979). Bootstrap methods: Another look at the jackknife. *The Annals of Statistics*, pages 1–26.
- Eldred, L. J., Wu, A. W., Chaisson, R. E., and Moore, R. D. (1998). Adherence to antiretroviral and pneumocystis prophylaxis in HIV disease. *Journal of Acquired Immune Deficiency Syndromes*, 18(2):117–125.
- Fujiwara, P. I., Larkin, C., and Frieden, T. R. (1997). Directly observed therapy in New York City: History, implementation, results, and challenges. *Clinics in Chest Medicine*, 18(1):135–148.
- Gifford, A., Shively, M., Bormann, J., Timberlake, D., and Bozzette, S. (1998). Self-reported adherence to combination antiretroviral medication regimens in a community-based sample of HIV-infected adults. In *12th World AIDS Conference*.

- Goetgeluk, S., Vansteelandt, S., and Goetghebeur, E. (2008). Estimation of controlled direct effects. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 70(5):1049–1066.
- Horvitz, D. G. and Thompson, D. J. (1952). A generalization of sampling without replacement from a finite universe. *Journal of the American Statistical Association*, 47(260):663–685.
- Imai, K., Keele, L., and Tingley, D. (2010a). A general approach to causal mediation analysis. *Psychological Methods*, 15(4):309.
- Imai, K., Keele, L., and Yamamoto, T. (2010b). Identification, inference and sensitivity analysis for causal mediation effects. *Statistical Science*, pages 51–71.
- Kang, J. D. and Schafer, J. L. (2007). Demystifying double robustness: A comparison of alternative strategies for estimating a population mean from incomplete data. *Statistical Science*, pages 523–539.
- Mammen, E. (1993). Bootstrap and wild bootstrap for high dimensional linear models. *The Annals of Statistics*, pages 255–285.
- Mills, E. J., Nachega, J. B., Buchan, I., Orbinski, J., Attaran, A., Singh, S., Rachlis, B., Wu, P., Cooper, C., Thabane, L., et al. (2006). Adherence to antiretroviral therapy in Sub-Saharan Africa and North America: A meta-analysis. *Journal of the American Medical Association*, 296(6):679–690.
- Pearl, J. (2000). *Causality: Models, Reasoning and Inference*. Cambridge University Press, New York. 2nd edition, 2009.
- Pearl, J. (2001). Direct and indirect effects. In *Proceedings of the Seventeenth Conference on Uncertainty in Artificial Intelligence*, pages 411–420. Morgan Kaufmann Publishers Inc.
- Petersen, M. L., Sinisi, S. E., and van der Laan, M. J. (2006). Estimation of direct causal effects. *Epidemiology*, 17(3):276–284.
- Pop-Eleches, C., Thirumurthy, H., Habyarimana, J. P., Zivin, J. G., Goldstein, M. P., De Walque, D., Mackeen, L., Haberer, J., Kimaiyo, S., Sidle, J., et al. (2011). Mobile phone technology

- gies improve adherence to antiretroviral treatment in a resource-limited setting: A randomized controlled trial of text message reminders. *AIDS (London, England)*, 25(6):825.
- Richardson, T. S. (2009). A factorization criterion for acyclic directed mixed graphs. In *Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence*, pages 462–470. AUAI Press.
- Roberts, K. J. (2000). Barriers to and facilitators of HIV-positive patients' adherence to antiretroviral treatment regimens. *AIDS Patient Care and STDs*, 14(3):155–168.
- 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(9):1393–1512.
- Robins, J., Sued, M., Lei-Gomez, Q., and Rotnitzky, A. (2007). Comment: Performance of double-robust estimators when "inverse probability" weights are highly variable. *Statistical Science*, pages 544–559.
- Robins, J. M. (1999). Testing and estimation of direct effects by reparameterizing directed acyclic graphs with structural nested models. *Computation, Causation, and Discovery*, pages 349–405.
- Robins, J. M. (2003). Semantics of causal DAG models and the identification of direct and indirect effects. *Highly Structured Stochastic Systems*, pages 70–81.
- Robins, J. M. and Greenland, S. (1992). Identifiability and exchangeability for direct and indirect effects. *Epidemiology*, pages 143–155.
- Robins, J. M. and Richardson, T. S. (2010). Alternative graphical causal models and the identification of direct effects. *Causality and Psychopathology: Finding the Determinants of Disorders and Their Cures*, pages 103–158.
- Robins, J. M., Ritov, Y., et al. (1997). Toward a curse of dimensionality appropriate (CODA) asymptotic theory for semi-parametric models. *Statistics in Medicine*, 16(3):285–319.

- Rubin, D. B. (1974). Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology*, 66(5):688.
- Rubin, D. B. (1978). Bayesian inference for causal effects: The role of randomization. *The Annals of Statistics*, pages 34–58.
- Shpitser, I. (2013). Counterfactual graphical models for longitudinal mediation analysis with unobserved confounding. *Cognitive Science*, 37(6):1011–1035.
- Suárez, P. G., Watt, C. J., Alarcón, E., Portocarrero, J., Zavala, D., Canales, R., Luelmo, F., Espinal, M. A., and Dye, C. (2001). The dynamics of tuberculosis in response to 10 years of intensive control effort in Peru. *Journal of Infectious Diseases*, 184(4):473–478.
- Tang, M. W., Kanki, P. J., and Shafer, R. W. (2012). A review of the virological efficacy of the 4 World Health Organization–recommended tenofovir-containing regimens for initial HIV therapy. *Clinical Infectious Diseases*, 54(6):862–875.
- Taylor, J. M., Wang, Y., and Thiébaud, R. (2005). Counterfactual links to the proportion of treatment effect explained by a surrogate marker. *Biometrics*, 61(4):1102–1111.
- Tchetgen Tchetgen, E. J. (2011). On causal mediation analysis with a survival outcome. *The International Journal of Biostatistics*, 7(1):1–38.
- Tchetgen Tchetgen, E. J. (2013). Inverse odds ratio-weighted estimation for causal mediation analysis. *Statistics in Medicine*, 32(26):4567–4580.
- Tchetgen Tchetgen, E. J. and Lin, S. H. (2012). Robust estimation of pure/natural direct effects with mediator measurement error. *Harvard University Biostatistics Paper Series*, Working Paper 152.
- Tchetgen Tchetgen, E. J. and Shpitser, I. (2012). Semiparametric theory for causal mediation analysis: Efficiency bounds, multiple robustness and sensitivity analysis. *The Annals of Statistics*, 40(3):1816–1845.
- Tchetgen Tchetgen, E. J. and Shpitser, I. (2014). Semiparametric estimation of models for natural direct and indirect effects. *Biometrika (In press)*.

- Ten Have, T. R., Joffe, M. M., Lynch, K. G., Brown, G. K., Maisto, S. A., and Beck, A. T. (2007). Causal mediation analyses with rank preserving models. *Biometrics*, 63(3):926–934.
- Tian, J. and Pearl, J. (2002). A general identification condition for causal effects. In *AAAI/IAAI*, pages 567–573.
- Valeri, L., Lin, X., and VanderWeele, T. J. (2014). Mediation analysis when a continuous mediator is measured with error and the outcome follows a generalized linear model. *Statistics in medicine*.
- van der Laan, M. J. and Petersen, M. L. (2008). Direct effect models. *The International Journal of Biostatistics*, 4(1):1–27.
- VanderWeele, T. and Vansteelandt, S. (2009). Conceptual issues concerning mediation, interventions and composition. *Statistics and its Interface*, 2:457–468.
- VanderWeele, T. J. (2009). Marginal structural models for the estimation of direct and indirect effects. *Epidemiology*, 20(1):18–26.
- VanderWeele, T. J. (2011). Causal mediation analysis with survival data. *Epidemiology (Cambridge, Mass.)*, 22(4):582.
- VanderWeele, T. J., Valeri, L., and Ogburn, E. L. (2012). The role of measurement error and misclassification in mediation analysis. *Epidemiology (Cambridge, Mass.)*, 23(4):561.
- VanderWeele, T. J. and Vansteelandt, S. (2010). Odds ratios for mediation analysis for a dichotomous outcome. *American Journal of Epidemiology*, 172(12):1339–1348.
- Vranceanu, A. M., Safren, S. A., Lu, M., Coady, W. M., Skolnik, P. R., Rogers, W. H., and Wilson, I. B. (2008). The relationship of post-traumatic stress disorder and depression to antiretroviral medication adherence in persons with HIV. *AIDS Patient Care and STDs*, 22(4):313–321.