

# IDENTIFICATION AND INFERENCE FOR MARGINAL AVERAGE TREATMENT EFFECT ON THE TREATED WITH AN INSTRUMENTAL VARIABLE

Lan Liu, Wang Miao, Baoluo Sun, James Robins and Eric Tchetgen Tchetgen

*University of Minnesota, Peking University, National University of Singapore,  
Harvard University and University of Pennsylvania*

*Abstract:* In observational studies, treatments are typically not randomized and, therefore, estimated treatment effects may be subject to a confounding bias. The instrumental variable (IV) design plays the role of a quasi-experimental handle because the IV is associated with the treatment and only affects the outcome through the treatment. In this paper, we present a novel framework for identification and inferences, using an IV for the marginal average treatment effect amongst the treated (ETT) in the presence of unmeasured confounding. For inferences, we propose three semiparametric approaches: (i) an inverse probability weighting (IPW); (ii) an outcome regression (OR); and (iii) a doubly robust (DR) estimation, which is consistent if either (i) or (ii) is consistent, but not necessarily both. A closed-form locally semiparametric efficient estimator is obtained in the simple case of a binary IV, and outcome, and the efficiency bound is derived for the more general case.

*Key words and phrases:* Counterfactuals, double robustness, effect of treatment on the treated, instrumental variable, unmeasured confounding.

## 1. Introduction

Epidemiology studies and social sciences often aim to evaluate the effect of a treatment. For practical reasons, the average treatment effect among treated individuals (ETT) is sometimes of greater interest than the treatment effect in the population. In epidemiology studies concerning the toxic effects of a new drug or the treatment effect only on those who take the treatment, the ETT is the parameter of interest, and is known as “the effect of exposure on the exposed,” or “standardized morbidity” (Miettinen (1974); Greenland and Robins (1986)). In econometrics, ETT is often used to evaluate the effects of a policy on those to whom it applies. For example, Angrist (1995) evaluated the average effect of military service on the civilian earnings for veterans. Heckman, Ichimura and Todd (1997, 1998) evaluated the average effect of job training on the program participants.

In observational or randomized studies with noncompliance, a primary challenge is the presence of unmeasured confounding, that is, the outcomes between treatment groups may differ, not only because of the treatment effect, but also because of unmeasured factors that may affect the treatment selection.

Instrumental variables (IV) are useful in addressing unmeasured confounding. An IV is associated with the treatment and affects the outcome only through the treatment. The key idea of the IV method is to extract exogenous variation in the treatment that is unconfounded with the outcome, and to take advantage of this bias-free component to make a causal inference about the treatment effect (Robins (1989); Angrist, Imbens and Rubin (1996); Heckman (1997)).

The development of the IV approach can be traced back to Wright (1928) and Goldberger (1972) under linear structural equations in econometrics. Imbens and Angrist (1994), Angrist, Imbens and Rubin (1996), and Heckman (1997) formalized the IV approach within the framework of potential outcomes or counterfactuals. Under additive and multiplicative structural nested models (SNMs), Robins (1989) and Robins (1994) evaluated the corresponding average treatment effect among treated individuals (ETT), conditional on the IV and observed covariates. Identification is achieved by assuming a certain degree of homogeneity with regard to the IV in an SNM of the conditional ETT (Hernán and Robins (2006)). Mainly, the assumption states that the magnitude of the conditional ETT does not vary with the IV. This is also referred to as the no-current treatment value interaction assumption. Under a similar identifying assumption, Vansteelandt and Goetghebeur (2003), Robins and Rotnitzky (2004), Tan (2010), Clarke, Palmer and Windmeijer (2015), and Matsouaka and Tchetgen Tchetgen (2014) investigated estimations of this conditional causal effect using additive, multiplicative and logistic SNMs.<sup>1</sup>

The literature mentioned above has several limitations. First, the literature focuses on the ETT conditional on the IV and observed covariates. The identification of such conditional ETT was achieved by specifying a functional form of the treatment causal effect. However, this is not an appealing solution, because it places constraints directly on the main parameter of interest, and a misspecification of this functional form would lead to biased result. Second, the available inference methods require that the treatment propensity score be

---

<sup>1</sup> In another line of research, Imbens and Angrist (1994) and Angrist, Imbens and Rubin (1996) defined the treatment effect on individuals who would comply to their assigned treatment. Under a monotonicity assumption about the effect of the IV on exposure, the complier average treatment effect can be identified. Further research along these lines include fully parametric estimation strategies (Tan (2006); Barnard et al. (2003); Frangakis et al. (2004)) and semiparametric methods (Abadie (2003); Abadie, Angrist and Imbens (2002); Tan (2006); Ogburn, Rotnitzky and Robins (2015)).

correctly specified, even for an outcome regression-based estimator (Tan (2010)).

In this study, we remedy these limitations in a novel framework for identification and estimation using an IV of the marginal ETT in the presence of unmeasured confounding. By targeting directly the marginal ETT, we allow the conditional causal effect to remain unrestricted. Our methods are particularly valuable when the primary goal is to obtain an accurate estimate of the treatment effect. Additionally, we propose a new identification strategy that is applicable to any type of outcome, and provides necessary and sufficient global identification conditions. Moreover, for inference purposes, we propose three semiparametric estimators, allowing for flexible covariate adjustment: (i) an inverse probability weighting (IPW); (ii) an outcome regression (OR), and (iii) a doubly robust (DR) estimation, which is consistent if either (i) or (ii) is consistent, but not necessarily both.

The remainder of the paper proceeds as follows. In Section 2, we introduce our notation and state our main assumptions. We study the nonparametric identification of the ETT in Section 3. We introduce the IPW, OR, and DR estimators in Section 4. In Section 5, we assess the performance of the various estimators in a simulation study. In Section 6, we further illustrate the methods using a study on the impact of participation in a 401(k) retirement program on savings. We conclude with a brief discussion in Section 7.

## 2. Preliminary Results

Suppose that one observes independently and identically distributed (i.i.d) data  $O = (A, Y, Z, C)$ , where  $A$  is a binary treatment,  $Y$  is the outcome of interest, which may be dichotomous, polytomous, discrete, or continuous, and the candidate IV  $Z$  and covariates  $C$  are both pre-exposure variables. Let  $a, y, z$ , and  $c$  denote the possible values of  $A, Y, Z$ , and  $C$ , respectively. Let  $Y_{az}$  denote the potential outcome if  $A$  and  $Z$  are set to  $a$  and  $z$ , respectively and let  $Y_a$  denote the potential outcome if only  $A$  is set to  $a$ . We formalize the IV assumptions using potential outcomes as follows:

(IV.1) Stochastic exclusion restriction:

$$Y_{az} = Y_a \text{ almost surely for all } a \text{ and } z;$$

(IV.2) Unconfounded IV-outcome relation:

$$f_{Y_0|Z,C}(y|z, c) = f_{Y_0|C}(y|c), \text{ for all } z \text{ and } c;$$

(IV.3) IV relevance:

$$\Pr(A = 1|Z = z, C = c) \neq \Pr(A = 1|Z = 0, C = c), \text{ for all } z \neq 0 \text{ and } c.$$

Assumption (IV.1) states that  $Z$  does not have a direct effect on the outcome  $Y$ . Thus, we use  $Y_a$  to denote the potential outcome under treatment  $a$  for  $a = 0, 1$ . Assumption (IV.2) is ensured under physical randomization, but will hold more generally if  $C$  includes all common causes of  $Z$  and  $Y$ . Assumptions (IV.1)–(IV.2) together imply that, conditional on  $C$ , the IV is independent of the potential outcome for the unexposed; that is,  $Y_0 \perp\!\!\!\perp Z|C$ . Assumption (IV.3) states that  $A$  and  $Z$  have a non-null association, conditional on  $C$ , even if the association is not causal. If assumptions (IV.1)–(IV.3) are satisfied,  $Z$  is said to be a valid IV.

We make the consistency assumption  $Y = AY_1 + (1 - A)Y_0$ . The marginal treatment effect on the treated is  $\text{ETT} = E(Y_1 - Y_0|A = 1)$ . Because  $E(Y_1|A = 1) = E(Y|A = 1)$  can be consistently estimated from the average observed outcome of treated individuals, we focus on making inferences about  $\psi$ , where

$$\psi = E(Y_0|A = 1).$$

Suppose there exist unmeasured variables, denoted by  $U$ , such that controlling for  $(U, Z, C)$  suffices to account for confounding; that is,  $Y_0 \perp\!\!\!\perp A|(U, Z, C)$ ; however,

$$Y_0 \not\perp\!\!\!\perp A|(Z, C), \tag{2.1}$$

where  $\perp\!\!\!\perp$  denotes statistical independence. As pointed out by Robins, Rotnitzky and Scharfstein (2000), potential outcomes can be viewed as the ultimate unmeasured confounders. This is because, by the consistency assumption, the observed outcome  $Y$  is a deterministic function of the treatment and the potential outcomes. Thus, given  $(Y_0, Y_1)$ ,  $U$  does not contain any further information about  $Y$ . To make explicit use of (2.1), we define the extended propensity score

$$\pi(Y_0, Z, C) = \Pr(A = 1|Y_0, Z, C),$$

as a function of  $Y_0$ .

### 3. Nonparametric Identification

While assumptions (IV.1)–(IV.3) suffice to obtain a valid test of the sharp null hypothesis of no treatment effect (Robins (1994)), and can also be used

to test for the presence of a confounding bias (Pearl (1995)), the ETT is not uniquely determined by the observed data without any additional restrictions. For simplicity, we first consider the situation where covariates are omitted, and the outcome and the IV are both binary. From the observed data, one can identify the quantities  $\Pr(Y_0, Z|A = 0)$ ,  $\Pr(Z|A = 1)$ , and  $\Pr(A = 0)$ . These quantities are functions of the unknown parameters  $\Pr(Z = 1)$ ,  $\Pr(Y_0 = 1)$ , and  $\Pr(A = 0|Y_0, Z)$ . Without imposing any additional assumption, there are six unknown parameters (one for  $\Pr(Z = 1)$ , one for  $\Pr(Y_0 = 1)$ , and four for  $\Pr(A = 0|Y_0, Z)$ ). However, only five degrees of freedom are available from the observed data (one for  $\Pr(A = 0)$ , one for  $\Pr(Z|A = 1)$ , and three for  $\Pr(Y, Z|A = 0)$ ). As a result, the joint distribution  $f(A, Y_0, Z)$  is not uniquely identified. In particular,  $\psi$  is not identified.

For identification purposes, additional assumptions, such as Robins' no-current treatment value interaction assumption (Hernán and Robins (2006)), must be imposed to reduce the set of candidate models for the joint distribution  $f(A, Y_0, Z, C)$ . Below, we give a general necessary and sufficient condition for identification. Let  $\mathcal{P}_{A|Y_0, Z, C}$  and  $\mathcal{P}_{Y_0|C}$  denote the collections of candidates for  $\Pr(A = 0|Y_0, Z, C)$  and  $f(Y_0|C)$ , respectively, which are known to satisfy (IV.1) and (IV.2).

**Condition 1.** Any two distinct elements  $\Pr_1(A = 0|Y_0, Z, C)$ ,  $\Pr_2(A = 0|Y_0, Z, C) \in \mathcal{P}_{A|Y_0, Z, C}$  and  $f_1(Y_0|C)$ ,  $f_2(Y_0|C) \in \mathcal{P}_{Y_0|C}$  satisfy the inequality:

$$\frac{\Pr_1(A = 0|Y_0, Z, C)}{\Pr_2(A = 0|Y_0, Z, C)} \neq \frac{f_2(Y_0|C)}{f_1(Y_0|C)}.$$

The following proposition states that condition 1 is a necessary and sufficient condition for the identifiability of the joint distribution of  $(A, Y_0, Z, C)$ , where  $Y_0$  and  $Z$  may be dichotomous, polytomous, discrete, or continuous.

**Proposition 1.** *The joint distribution of  $(A, Y_0, Z, C)$  is identified in the model defined by  $\mathcal{P}_{A|Y_0, Z, C}$  and  $\mathcal{P}_{Y_0|C}$  if and only if condition 1 holds.*

It is convenient to check condition 1 for parametric models, but it may be more difficult for semiparametric and nonparametric models, because  $\mathcal{P}_{A|Y_0, Z, C}$  and  $\mathcal{P}_{Y_0|C}$  can be complicated. The following corollary gives a more convenient condition.

**Corollary 1.** *Suppose that for any two candidates  $\Pr_1(A = 0|Y_0, Z, C)$ ,  $\Pr_2(A = 0|Y_0, Z, C) \in \mathcal{P}_{A|Y_0, Z, C}$ , the ratio  $\Pr_1(A = 0|Y_0, Z, C)/\Pr_2(A = 0|Y_0, Z, C)$  is*

either a constant or varies with  $Z$ . Then, the joint distribution of  $(A, Y_0, Z, C)$  is identified.

Although the condition provided in Corollary 1 is a sufficient condition for identification, it allows identification of a large class of models. The proofs of Proposition 1 and Corollary 1 are given in the Supplementary Material. We further illustrate Proposition 1 and Corollary 1 with several examples. For simplicity, we again omit the covariates; however, we show at the end of this section that similar results with covariates can be derived. For simplicity, we first consider the case of a binary outcome with a binary IV.

**Example 1.** Consider a model  $\mathcal{P}_{A|Y_0,Z} = \{\Pr(A = 0|Y_0, Z) : \text{logit } \Pr(A = 0|Y_0, Z; \theta_1, \theta_2, \eta_1, \eta_2) = \theta_1 + \theta_2 Z + \eta_1 Y_0 + \eta_2 Y_0 Z, \theta_1, \theta_2, \eta_1, \eta_2 \in (-\infty, \infty)\}$ . The model is saturated because  $\mathcal{P}_{A|Y_0,Z}$  contains all possible treatment mechanisms. It can be shown that neither the joint distribution nor  $\psi$  is identified, even under assumptions (IV.1)–(IV.3).

Example 1 shows that the joint density  $f(A, Y_0, Z)$  is not identified when the treatment selection mechanism is left unrestricted under (IV.1)–(IV.3). However, we show that the joint density  $f(A, Y_0, Z)$  is identified, assuming a separable treatment mechanism on the additive scale.

**Example 2.** Consider a model  $\mathcal{P}_{A|Y_0,Z} = \{\Pr(A = 0|Y_0, Z) : \text{logit } \Pr(A = 0|Y_0, Z; \theta_1, \theta_2, \eta_1) = \theta_1 + \theta_2 Z + \eta_1 Y_0; \theta_1, \theta_2, \eta_1 \in (-\infty, \infty)\}$ . The model is separable because  $\mathcal{P}_{A|Y_0,Z}$  excludes an interaction between  $Y_0$  and  $Z$ . It can be shown that both the joint distribution and  $\psi$  are identified under assumptions (IV.1)–(IV.3).

Example 2 agrees with the intuition that identification follows from having fewer parameters than the saturated model. Under the assumed model, we have five unknown parameters and five available degrees of freedom from the empirical distribution. We show in the next example that the joint distribution and  $\psi$  can be identified in a general separable model when the outcome and instrument are both continuous.

**Example 3.** Consider the logistic separable treatment mechanism:  $\mathcal{P}_{A|Y_0,Z} = \{\Pr(A = 0|Y_0, Z) : \text{logit } \Pr(A = 0|Y_0, Z) = q(Z) + h(Y_0)\}$ , where  $q$  and  $h$  are unknown differentiable functions, with  $h(0) = 0$ . It can be shown that  $\mathcal{P}_{A|Y_0,Z}$  satisfies condition 1, and thus the joint distribution is identified under (IV.1)–(IV.3).

These results can be generalized to include covariates  $C$ . For instance, by allowing both  $q$  and  $h$  to depend on  $C$  in example 3:

$$\mathcal{P}_{A|Y_0,Z,C} = \{\Pr(A = 0|Y_0, Z, C) : \text{logit } \Pr(A = 0|Y_0, Z, C) = q(Z, C) + h(Y_0, C)\},$$

where  $h(0, C) = 0$ , the joint distribution is identified whenever the interaction term of  $Y_0$  and  $Z$  is absent.

In the Supplementary Material, we present proofs for the above examples. We also provide additional examples, such as the case of a continuous outcome with a binary IV, probit link, and separable treatment mechanism.

#### 4. Estimation

Although nonparametric identification conditions are provided in Section 3, such conditions will seldom suffice for reliable statistical inferences. Typically, in observational studies, the set of covariates  $C$  is too large for a nonparametric inference, owing to the curse of dimensionality (Robins and Ritov (1997)). Therefore, we posit parametric models for various nuisance parameters, and provide three possible approaches for a semiparametric inference that depend on different subsets of models. We describe an IPW, an OR, and a DR estimator of the marginal ETT under assumptions (IV.1)–(IV.2) and condition 1. Throughout, we posit a parametric model  $f_{Z|C}(z|c) = \Pr(Z = z|C = c; \rho)$  for the conditional density of  $Z$ , given  $C$ . Let  $\hat{\rho}$  denote the maximum likelihood estimator (MLE) of  $\rho$ . Let  $\mathbb{P}_n$  denote the empirical measure; that is,  $\mathbb{P}_n f(O) = n^{-1} \sum_{i=1}^n f(O_i)$ . Let  $\hat{E}$  denote the expectation taken under the empirical distribution of  $C$ , and let  $\widehat{\Pr}(A = 1) = \sum_{i=1}^n A_i/n$  denote the empirical probability of receiving treatment.

##### 4.1. IPW estimator

For the estimation, we first propose an IPW IV approach that extends the standard IPW estimation of an ETT to an IV setting. We make the positivity assumption that for all values of  $Y_0$ ,  $Z$ , and  $C$ , the probability of not being exposed to treatment is bounded away from zero. When  $A = 0$ , by the consistency assumption,  $Y = Y_0$ ; thus, we could use  $Y_0$  and  $Y$  interchangeably for the  $A = 0$  group. The IPW approach relies on the crucial assumption that the extended propensity score model  $\pi(Y_0, Z, C; \gamma)$  is correctly specified, with unknown finite-dimensional parameter  $\gamma$ , and the following representation of ETT:

$$E(Y_0|A = 1) = E \left\{ \frac{\pi(Y, Z, C)Y(1 - A)}{\Pr(A = 1)\{1 - \pi(Y, Z, C)\}} \right\}. \quad (4.1)$$

A derivation of the above equation is given in the Supplementary Material. We solve the following equations to obtain an estimator  $\hat{\gamma}$  of  $\gamma$ :

$$\mathbb{P}_n \left\{ \frac{1 - A}{1 - \pi(Y, Z, C; \hat{\gamma})} - 1 \right\} = 0, \quad (4.2)$$

$$\mathbb{P}_n \left[ \frac{1 - A}{1 - \pi(Y, Z, C; \hat{\gamma})} \{h_1(Z, C) - E(h_1(Z, C)|C; \hat{\rho})\} \right] = 0, \quad (4.3)$$

$$\mathbb{P}_n \left[ \frac{1 - A}{1 - \pi(Y, Z, C; \hat{\gamma})} \{h_2(C) - \hat{E}(h_2(C))\} \right] = 0, \quad (4.4)$$

$$\mathbb{P}_n \left[ \frac{1 - A}{1 - \pi(Y, Z, C; \hat{\gamma})} t(Y, C) \{l(Z, C) - E(l(Z, C)|C; \hat{\rho})\} \right] = 0, \quad (4.5)$$

where  $(h_1^T, h_2^T, l^T)^T$  satisfies the regularity condition (A.1) described in the Supplementary Material. Equations (4.3) and (4.4) identify the association between  $(Z, C)$  and  $A$  in  $\pi(0, Z, C)$ . If there is no selection bias, equations (4.2)–(4.4) are adequate to estimate the propensity score. By utilizing the IV property (IV.1)–(IV.2), equation (4.5) identifies the degree of selection bias encoded in the dependence of  $\pi$  on  $Y_0$ . Equations (4.3) and (4.5) both require the conditional density of IV  $Z$ ,  $f_{Z|C}(z|c; \rho)$ , to be correctly modeled. By equation (4.1), an extended propensity score estimator leads to an estimator of  $\psi$ . We have the following result:

**Proposition 2.** *Under (IV.1)–(IV.2) and condition 1, suppose the extended propensity score model  $\pi(Y_0, Z, C; \gamma)$  and  $f_{Z|C}(z|c; \rho)$  are correctly specified. Then, the IPW estimator*

$$\hat{\psi}^{ipw} = \mathbb{P}_n \widehat{\frac{\pi(Y, Z, C; \hat{\gamma})Y(1 - A)}{\Pr(A = 1)\{1 - \pi(Y, Z, C; \hat{\gamma})\}}}$$

*is consistent for  $\psi$ .*

Note that the extended propensity score model can use any well-defined link function (such as, logit, probit), and if condition 1 holds, Proposition 2 still holds. The functions  $h_1$ ,  $h_2$ ,  $t$ , and  $l$  can be chosen based on the model for the extended propensity score. For example, assuming logit  $\pi(Y_0, Z, C; \gamma) = \theta_0 + \theta_1 Z + \theta_2 C + \eta Y_0$ , where  $\tilde{\eta} = (\theta_1, \theta_2, \eta)^T$  is a  $k$ -dimensional parameter vector. The  $k$ -dimensional function  $(h_1, h_2, t)^T$  can be chosen as  $(h_1, h_2, t)^T =$

$\partial \log \pi(Y_0, Z, C; \gamma) / \partial \tilde{\eta} = (Z, C, Y_0)^T$ , and  $l$  can be chosen as any scalar function of  $(Z, C)$ , such as  $l(Z, C) = Z$ . Thus, we have exactly  $k + 1$  estimating equations. The choice of  $h_1, h_2, t$ , and  $l$  will generally impact efficiency but should not affect consistency, as long as the identification conditions hold and the required models are correctly specified. The choices of  $h_1, h_2, t$ , and  $l$  that lead to the most efficient IPW estimator can be derived using the results in Newey and McFadden (1994). Owing to space constraints, we illustrate in details the choice of similar functions for efficient DR estimator in the section 4.3; a similar derivation could be made here.

The asymptotic variance of the IPW estimator can be derived using standard M-estimation theory (van der Vaart (1998)). Specifically, let  $G_{\psi}^{ipw}(O) = \hat{\psi}^{ipw} - \psi$ , and let  $G_{\gamma}^{ipw}(O)$  and  $G_{\rho}^{ipw}(O)$  denote the score functions for  $\gamma$  and  $\rho$ , respectively. Then,  $\theta = (\psi, \gamma, \rho)$  is the solution to  $\int G^{ipw}(o; \theta) dF(o; \theta) = 0$ , where  $G^{ipw}(O) = \{G_{\psi}^{ipw}(O), G_{\gamma}^{ipw}(O), G_{\rho}^{ipw}(O)\}$ . Thus,  $\hat{\theta}^{ipw} = (\hat{\psi}^{ipw}, \hat{\gamma}, \hat{\rho})$  is the solution to  $\mathbb{P}_n G^{ipw}(O; \theta) = 0$ . By M-estimation theory, under regularity conditions,  $\sqrt{n}(\hat{\theta}^{ipw} - \theta)$  converges in distribution to  $N(0, \Sigma^{ipw})$  when  $n$  goes to infinity, where  $\Sigma^{ipw} = U^{-1} V U^{-T}$ ,  $U = -E\{\partial G^{ipw}(O_i; \theta) / \partial \theta\}$  and  $V = E\{G^{ipw}(O_i; \theta)^{\otimes 2}\}$ . A consistent estimator of the asymptotic variance of  $\hat{\psi}^{ipw}$  can be constructed by replacing the expectations with their empirical counterparts, and by replacing the parameters with their estimates. Such a variance estimator is also referred to as a sandwich estimator.

## 4.2. OR and DR estimators

Because  $Y_0$  is never observed for the treated group, we use the following equation to decompose  $E[Y_0 | A = 1, Z, C]$  into two parts: one can be estimated directly using a restricted MLE, and the other can be computed by solving an estimating equation. Specifically, we have

$$E\{g(Y_0, C) | A = 1, Z, C\} = \frac{E[\exp\{\alpha(Y, Z, C)\}g(Y, C) | A = 0, Z, C]}{E[\exp\{\alpha(Y, Z, C)\} | A = 0, Z, C]}, \quad (4.6)$$

where  $g$  is any function of  $Y_0$  and  $C$ , and  $\alpha(Y_0, Z, C)$  is the generalized odds ratio function relating  $A$  and  $Y_0$ , conditional on  $Z$  and  $C$ , as

$$\alpha(Y_0, Z, C) = \log \frac{f(Y_0 | A = 1, Z, C) f(Y_0 = 0 | A = 0, Z, C)}{f(Y_0 | A = 0, Z, C) f(Y_0 = 0 | A = 1, Z, C)}.$$

Because the association between  $Y_0$  and  $A$  is attributed to unmeasured confounding,  $\alpha(Y_0, Z, C)$  can be interpreted as the selection bias function. Thus,

we express the conditional mean function  $E\{g(Y_0, C)|A = 1, Z, C\}$  in terms of  $f(Y|A = 0, Z, C)$  and  $\alpha(Y_0, Z, C)$ . We prove equation (4.6) in the Supplementary Material.

Let  $f(Y|A = 0, Z, C; \xi)$  denote a model for the density of the outcome among the unexposed, conditional on  $Z$  and  $C$ , and let  $\hat{\xi}$  denote the restricted MLE of  $\xi$  obtained using only the data for the unexposed. Let  $\eta$  denote the parameter indexing a parametric model for the selection bias function  $\alpha$  as  $\alpha(Y_0, Z, C; \eta)$ . We obtain an estimator for  $\eta$  by solving:

$$\mathbb{P}_n \left[ \left\{ w(Z, C) - E(w(Z, C)|C; \hat{\rho}) \right. \right. \\ \left. \left. \left\{ AE[g(Y_0, C)|A = 1, Z, C; \eta, \hat{\xi}] + (1 - A)g(Y, C) \right\} \right\} \right] = 0, \quad (4.7)$$

for any choice of functions  $w$  and  $g$ , such that the regularity condition (A.2) stated in the Supplementary Material holds. Intuitively, the left-hand side of equation (4.7) is an empirical estimator of the expected conditional covariance between  $w(Z, C)$  and  $g(Y_0, C)$ , given  $C$ , which should be zero, by (IV.1)–(IV.2). Equation (4.7) requires that the conditional density of IV  $Z$ ,  $f_{Z|C}(z|c; \rho)$ , be correctly modeled. Based on equation (4.6), we can construct an estimator for  $\psi$  based on  $\hat{\eta}$ ,  $\hat{\xi}$ , and  $\hat{\rho}$ .

**Proposition 3.** *Under (IV.1)–(IV.2) and condition 1, suppose  $\alpha(Y_0, Z, C; \eta)$ ,  $f_{Z|C}(z|c; \rho)$ , and  $f(Y|A = 0, Z, C; \xi)$  are correctly specified. Then, the OR estimator*

$$\hat{\psi}^{reg} = \mathbb{P}_n \frac{A}{\Pr(A = 1)} \frac{E[\exp\{\alpha(Y, Z, C; \hat{\eta})\}Y|A = 0, Z, C; \hat{\xi}]}{E[\exp\{\alpha(Y, Z, C; \hat{\eta})\}|A = 0, Z, C; \hat{\xi}]}$$

*is consistent for  $\psi$ .*

Functions  $g$  and  $\omega$  in equation (4.7) can be chosen based on the model we posit for  $\alpha(Y_0, Z, C)$ . For example, assuming

$$\alpha(Y_0, Z, C; \eta) = \eta Y_0, \quad (4.8)$$

$g$  can be chosen as  $g(Y_0, C) = \partial\alpha(Y_0, Z, C; \eta)/\partial\eta = Y_0$ , and  $\omega$  can be chosen as any scalar function of  $(Z, C)$ , such as,  $\omega(Z, C) = Z$ . The choices of  $g$  and  $\omega$  may impact efficiency, but do not affect consistency, as long as the identification conditions hold and the required models are correctly specified. The choices of  $g$  and  $\omega$  that lead to the most efficient OR estimator can be derived using Newey and McFadden (1994).

Tan (2010) proposed an OR estimator for the conditional ETT, which requires correctly specified models for both the treatment propensity score and the outcome regression function. In contrast, we circumvent the dependence of the regression estimator on the propensity score.

The proposed estimator for the nuisance parameter  $\eta$  is closely related to the regression estimator proposed by Vansteelandt and Goetghebeur (2003) when  $Y$  is binary. Vansteelandt and Goetghebeur (2003) developed a two-stage logistic estimator that combines a logistic SMM at the first stage and a logistic regression association model at the second stage. Specifically, they focused on estimating  $\zeta(Z, C) = \text{logit } \Pr(Y_1 = 1|A = 1, Z, C) - \text{logit } \Pr(Y_0 = 1|A = 1, Z, C)$ , which encodes the conditional ETT, given  $Z$  and  $C$ . Let  $\nu$  denote the parameter indexing a model for  $\zeta(Z, C)$  as  $\zeta(Z, C; \nu)$ . They proposed estimating  $\nu$  in the estimating equation

$$\mathbb{P}_n \left[ \left\{ w(Z, C) - E(w(Z, C)|C; \hat{\rho}) \right. \right. \\ \left. \left. \left\{ A \text{expit}\{\vartheta(Z, C; \hat{\varrho}) - \zeta(Z, C; \nu)\} + (1 - A)Y \right\} \right\} \right] = 0, \quad (4.9)$$

where  $\text{expit}(x) = \exp(x)/\{1 + \exp(x)\}$  and  $\vartheta(Z, C; \varrho) = \text{logit } \Pr(Y = 1|A = 1, Z, C; \varrho)$ .

Recall that we obtain an estimator of  $\eta$  indexing  $\alpha(Y_0, Z, C; \eta)$  in equation (4.7), which can be re-expressed as

$$\mathbb{P}_n \left[ \left\{ w(Z, C) - E(w(Z, C)|C; \hat{\rho}) \right. \right. \\ \left. \left. \left\{ A \text{expit}\{\delta(Z, C; \hat{\xi}) + \alpha(1, Z, C; \eta)\} + (1 - A)Y \right\} \right\} \right] = 0, \quad (4.10)$$

where  $\delta(Z, C; \hat{\xi}) = \text{logit } \Pr(Y_0 = 1|A = 0, Z, C)$ . Equations (4.9) and (4.10) differ mainly in the way  $\Pr(Y_0 = 1|A = 1, Z, C)$  is estimated. More specifically, (4.9) obtains  $\Pr(Y_0 = 1|A = 1, Z, C)$ , using  $\Pr(Y_1 = 1|A = 1, Z, C)$  as a baseline risk for the model, whereas (4.10) uses  $\Pr(Y_0 = 1|A = 0, Z, C)$  as a baseline risk. This difference is important, because Vansteelandt and Goetghebeur (2003) failed to obtain a DR estimator of  $\zeta(Z, C)$ , whereas, as we show next, our choice of parameterization yields a DR estimator of the marginal ETT.

Heretofore, we have constructed estimators using two different approaches. Both approaches assume correct models for  $\alpha(Y_0, Z, C; \eta)$  and  $f_{Z|C}(z|c; \rho)$ . The IPW approach further relies on a consistent estimator of the baseline extended

propensity score  $\beta(Z, C) = \text{logit } \Pr(A = 1|Y_0 = 0, Z, C)$ , which under the logit link, and together with  $\alpha(Y_0, Z, C; \eta)$ , provides a consistent estimator of the extended propensity score  $\pi(Y_0, Z, C; \gamma) = \text{expit } \{\alpha(Y_0, Z, C; \eta) + \beta(Z, C; \theta)\}$ . The OR approach further relies on a consistent estimator of  $f(Y|A = 0, Z, C)$ , which, together with  $\alpha(Y_0, Z, C; \eta)$ , provides a consistent estimator of  $\Pr(Y_0 = 1|A = 1, Z, C)$  by (4.6). Define  $\mathcal{M}_a$  as the collection of laws with parametric models  $f_{Z|C}(z|c; \rho)$ ,  $\alpha(Y_0, Z, C; \eta)$ , and  $\beta(Z, C; \theta)$ , while  $f(Y|A = 0, Z, C)$  is unrestricted. Likewise, define  $\mathcal{M}_y$  as the collection of laws with parametric models  $f_{Z|C}(z|c; \rho)$ ,  $\alpha(Y_0, Z, C; \eta)$ , and  $f(Y|A = 0, Z, C; \xi)$ , while  $\beta(Z, C)$  is unrestricted. The main appeal of a doubly robust estimator is that it remains consistent if either  $\beta(Z, C; \theta)$  or  $f(Y|A = 0, Z, C; \xi)$  is correctly specified. To derive a DR estimator for  $\psi$  in the union space  $\mathcal{M}_a \cup \mathcal{M}_y$ , we first propose a DR estimator for the parameter  $\eta$  of the selection bias model  $\alpha(Y_0, Z, C; \eta)$ . For notational convenience, let

$$\begin{aligned} & Q_g(Y, A, Z, C; \gamma, \xi) \\ &= \frac{(1 - A)\pi(Y, Z, C; \gamma)}{1 - \pi(Y, Z, C; \gamma)} \left[ g(Y, C) - \frac{E[\exp\{\alpha(Y, Z, C; \eta)\}g(Y, C)|A = 0, Z, C; \xi]}{E[\exp\{\alpha(Y, Z, C; \eta)\}|A = 0, Z, C; \xi]} \right] \\ &+ A \frac{E[\exp\{\alpha(Y, Z, C; \eta)\}g(Y, C)|A = 0, Z, C; \xi]}{E[\exp\{\alpha(Y, Z, C; \eta)\}|A = 0, Z, C; \xi]}. \end{aligned} \quad (4.11)$$

Consider the estimating equation for the selection bias parameter  $\tilde{\eta}$

$$\mathbb{P}_n \left[ [\omega(Z, C) - E\{\omega(Z, C)|C; \hat{\rho}\}] \tilde{Q}_g(Y, A, Z, C; \tilde{\gamma}, \hat{\xi}) \right] = 0, \quad (4.12)$$

where

$$\begin{aligned} & \tilde{Q}_g(Y, A, Z, C; \tilde{\gamma}, \hat{\xi}) \\ &= Q_g(Y, A, Z, C; \tilde{\gamma}, \hat{\xi}) + (1 - A)g(Y, C) \\ &= \frac{1 - A}{1 - \pi(Y, Z, C; \gamma)} g(Y, C) \\ &+ \frac{A - \pi(Y, Z, C; \gamma)}{1 - \pi(Y, Z, C; \gamma)} \frac{E[\exp\{\alpha(Y, Z, C; \eta)\}g(Y, C)|A = 0, Z, C; \xi]}{E[\exp\{\alpha(Y, Z, C; \eta)\}|A = 0, Z, C; \xi]}. \end{aligned}$$

Equation (4.12) is key to obtaining a DR estimation of the selection bias function, and thus of the ETT. Intuitively, the left-hand side of equation (4.12) is also an empirical estimator of the expected conditional covariance between  $w(Z, C)$  and  $g(Y_0, C)$ , given  $C$ , which should be zero, by (IV.1)–(IV.2). In ad-

dition to the model  $f_{Z|C}(z|c; \rho)$  for IV, equation (4.7) only involves an outcome regression model, whereas equation (4.12) involves both outcome regression and propensity score models. Hence, the parameter  $\hat{\eta}$  obtained from (4.7) depends on the correct specification of the outcome regression. In contrast, as shown in the following proposition, the parameter estimate for  $\eta$  obtained from (4.12) is doubly robust. We solve equation (4.12) jointly with equations (4.2)–(4.4), with  $\hat{\gamma}$  replaced by  $\tilde{\gamma} = (\hat{\eta}^{DR}, \tilde{\theta})$ . The choices of  $h_1, h_2, g$ , and  $w$  can be decided as in Sections 4.1 and 4.2.

**Proposition 4.** *Under (IV.1)–(IV.2) and condition 1,  $\hat{\eta}^{DR}$  and  $\hat{\psi}^{DR}$  are consistent in the union model  $\mathcal{M}_a \cup \mathcal{M}_y$ , where  $\hat{\psi}^{DR} = \hat{\mathbb{P}}_n Q_{\tilde{g}}(Y, A, Z, C; \tilde{\gamma}, \hat{\xi}) / \widehat{\Pr}(A = 1)$  and  $\tilde{g}(Y, C) = Y$ .*

Proposition 4 implies that  $\hat{\eta}^{DR}$  and  $\hat{\psi}^{DR}$  are both DR estimators, because their consistency requires that either the extended propensity score or the outcome regression model be correctly specified, but not necessarily both.

The asymptotic variance of the OR and DR estimators and the corresponding sandwich variance estimators can be derived similarly to those of the IPW estimator. We omit the details here owing to space constraints.

### 4.3. Local efficiency

The large sample variance of the doubly robust estimators  $\hat{\eta}^{DR}$  and  $\hat{\psi}^{DR}$  at the intersection submodel  $\mathcal{M}_a \cap \mathcal{M}_y$ , where all models are correctly specified, is determined by the choice of  $g(Y, C)$  and  $\omega(Z, C)$  in equation (4.12). In the Supplementary Material, we derive the semiparametric efficient score of  $(\eta, \psi)$  in a model  $\mathcal{M}_{np}$  that assumes only that  $Z$  is a valid IV and the selection bias function  $\alpha(Y_0, Z, C; \eta)$  is correctly specified. As discussed in the Supplementary Material, the efficient score is generally not available in closed form, except in special cases, such as when  $Z$  and  $Y$  are both polytomous. Here, we illustrate the result by constructing a locally efficient estimator of  $(\eta, \psi)$  when  $Z$  and  $Y$  are both binary. Here, similarly to the definition of  $\tilde{Q}_g(Y, A, Z, C; \gamma, \xi)$ , define

$$\begin{aligned} \tilde{Q}_v(Y, A, Z, C; \gamma, \xi) &= \frac{(1 - A)v(Y, Z, C)}{1 - \pi(Y, Z, C; \gamma)} \\ &+ \frac{A - \pi(Y, Z, C; \gamma)}{1 - \pi(Y, Z, C; \gamma)} \frac{E[\exp\{\alpha(Y, Z, C; \eta)\}v(Y, Z, C)|A = 0, Z, C; \xi]}{E[\exp\{\alpha(Y, Z, C; \eta)\}|A = 0, Z, C; \xi]}, \end{aligned}$$

where  $v$  is any function of  $(Y_0, Z, C)$ .

A one-step locally efficient estimator of  $\eta$  in  $\mathcal{M}_{np}$  is given by

$$\hat{\eta}^{eff} = \hat{\eta}^{DR} - \left\{ E\left(\nabla_{\eta} \widehat{S}_{\eta}^{eff} | \hat{\gamma}, \hat{\xi}\right) \right\}^{-1} E(\widehat{S}_{\eta}^{eff} | \hat{\gamma}, \hat{\xi}),$$

where  $\bar{v}(Y, Z, C) = \{Y - E(Y|C)\}\{Z - E(Z|C)\}$ ,  $\Delta(\eta) = \tilde{Q}_{\bar{v}}(Y, A, Z, C; \gamma, \xi)$  and

$$\widehat{S}_{\eta}^{eff} = E\{\Delta(\eta)\Delta(\eta)^T | C; \hat{\gamma}, \hat{\xi}\}^{-1} E\left\{ \frac{\partial \Delta(\eta)}{\partial \eta^T} | C; \hat{\gamma}, \hat{\xi} \right\} \Delta(\hat{\eta}^{eff})$$

is the efficient score of  $\eta$  evaluated at the estimated intersection submodel  $\mathcal{M}_a \cap \mathcal{M}_y$ . Further, let  $\hat{\psi}^{DR}(\hat{\eta}^{eff})$  denote a DR estimator for  $\psi$  evaluated at the estimated intersection submodel  $\mathcal{M}_a \cap \mathcal{M}_y$ , with  $\hat{\eta}^{eff}$  replacing  $\hat{\eta}^{DR}$ . Then, the efficient estimator of  $\psi$  is given by

$$\hat{\psi}^{eff} = \hat{\psi}^{DR}(\hat{\eta}^{eff}) - E\{\Delta^2(\hat{\eta}^{eff}) | C; \hat{\gamma}, \hat{\xi}\}^{-1} E\{\hat{\psi}^{DR}(\hat{\eta}^{eff})\Delta(\hat{\eta}^{eff}) | C; \hat{\gamma}, \hat{\xi}\} \Delta(\hat{\eta}^{eff}).$$

## 5. Simulations

Simulations for both binary and continuous outcomes were conducted to evaluate the finite-sample performance of the causal effect estimators derived in Sections 4.1 and 4.2. Let  $\mathcal{M}_a^c$  denote the complement space of  $\mathcal{M}_a$ , and define  $\mathcal{M}_y^c$  in a similar manner. Simulations were conducted under three scenarios: (i)  $\mathcal{M}_a \cap \mathcal{M}_y$ , that is, both the outcome regression and the extended propensity score are correctly specified; (ii)  $\mathcal{M}_a \cap \mathcal{M}_y^c$  that is only the extended propensity score is correctly specified; and (iii)  $\mathcal{M}_a^c \cap \mathcal{M}_y$ , that is, only the outcome regression model is correctly specified.

Simulations were first carried out for a binary outcome. For scenario (i), the simulation study was conducted in the following steps:

Step 1: A hypothetical study population of size  $n = 1,000$  (or  $n = 5,000$ ) was generated, and each individual had baseline covariates  $C_1$  and  $C_2$  generated independently from Bernoulli distributions with probability 0.4 and 0.6, respectively. Then, the IV  $Z$  was generated from the model  $\text{logit } \Pr(Z = 1|C) = 0.2 + 0.4C_1 - 0.5C_2$  and potential outcomes  $Y_0, Y_1$  were generated from the models  $\text{logit } \Pr(Y_0 = 1|Z, C) = 0.6 + 0.8C_1 - 2C_2$  and  $\text{logit } \Pr(Y_1 = 1|Z, C) = 0.7 - 0.3C_1$ , respectively. The treatment variable  $A$  was generated from  $\text{logit } \Pr(A = 1|Y_0, Z, C) = 0.4 + 2Z + 0.8C_1 - 0.6Y_0 - 1.6C_1Z$ , and the observed outcome was  $Y = Y_0(1 - A) + Y_1A$ .

Step 2: The following extended propensity score model was estimated, and the parameters  $\gamma = (\theta_1, \theta_2, \theta_3, \theta_4, \eta)$  in the model

$$\text{logit } \Pr(A = 1|Y_0, Z, C; \gamma) = \theta_1 + \theta_2 Z + \theta_3 C_1 + \theta_4 C_1 Z + \eta Y_0 \quad (5.1)$$

were estimated using estimating equations (4.2)–(4.5), with  $h_1(Z, C) = (Z, C_1 Z)^T$ ,  $h_2(C) = C_1$ ,  $t(Y, C) = Y$ , and  $l(Z, C) = Z$ . Then,  $\hat{\psi}^{ipw}$  was evaluated.

Step 3: The selection bias function was correctly specified as in (4.8),  $\xi$  in the regression outcome model

$$\text{logit } E(Y|A = 0, Z, C; \xi) = \xi_1 + \xi_2 C_1 + \xi_3 C_2 + \xi_4 Z + \xi_5 C_1 Z \quad (5.2)$$

was estimated using a restricted MLE, and  $\alpha$  was estimated by solving equation (4.7), with  $\omega(Z, C) = Z$  and  $g(Y, C) = Y$ . Then,  $\hat{\psi}^{reg}$  was evaluated.

Step 4: The selection bias function was correctly specified as in (4.8),  $\xi$  in equation (5.2) was estimated using a restricted MLE, the parameters  $\gamma$  in (5.1) were estimated using (4.2)–(4.4) and (4.12), where  $h, t, l, \omega$ , and  $g$  are chosen as in Step 2 and Step 3. Then,  $\hat{\psi}^{DR}$  was evaluated.

Step 5: Steps 1–4 were repeated 1,000 times.

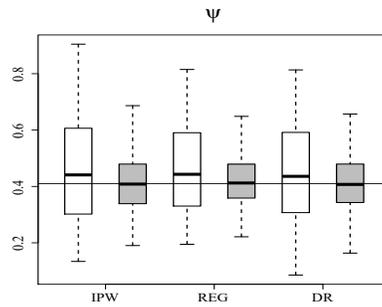
The data-generating mechanism described in Step 1 satisfies assumptions (IV.1)–(IV.2) for both  $a = 0, 1$ . As shown in example 1,  $\psi$  is identified from the observed data because the treatment mechanism is a separable logit model. In addition, in the Supplementary Material, we verify that model (5.2) for  $E(Y|A = 0, C, Z)$  contains the true data-generating mechanism. Simulations for scenario (ii) were similar to scenario (i), except that (5.1) was replaced with

$$\text{logit } \Pr(A = 1|Y_0, Z, C; \gamma) = \theta_1 + \theta_2 Z + \theta_3 C_1 + \eta Y_0, \quad (5.3)$$

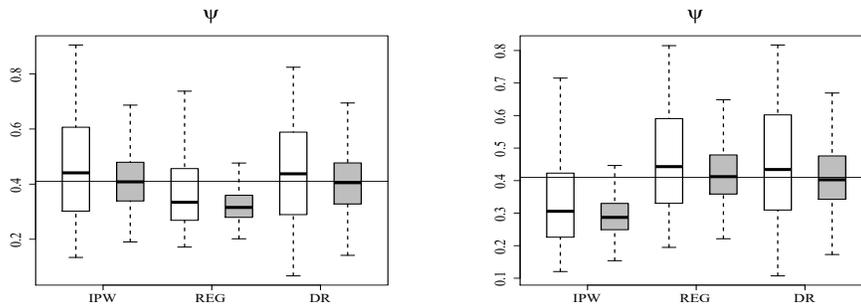
which is misspecified if  $\theta_4 \neq 0$  in equation (5.1). For scenario (iii), the potential outcome model (5.2) was replaced with

$$\text{logit } E(Y|A = 0, Z, C; \xi) = \xi_1 + \xi_2 C_1 + \xi_4 Z, \quad (5.4)$$

which is misspecified if  $\xi_3 \neq 0$  and  $\xi_5 \neq 0$  in equation (5.2). We use the R package BB (Varadhan and Gilbert (2009)) to solve the nonlinear estimating equations.



(a) Both outcome regression and extended propensity score are correctly specified.



(b) Only the extended propensity score is correctly specified.

(c) Only the outcome model is correctly specified.

Figure 1. Performance of the IPW, OR and DR estimators of  $\psi$  with binary outcomes.

Note: In each boxplot, the true value  $\psi_0$  is marked by the horizontal lines, white boxes are for  $n = 1,000$  and grey boxes are for  $n = 5,000$ .

The simulation results for 1,000 Monte Carlo samples are reported in Figure 1 (see also Table A.4 in the Supplementary Material), and the empirical coverage rates of 95% Wald-type confidence intervals are presented in Table 1. The confidence intervals are constructed using the sandwich variance estimators given in Section 4. Under a correct model specification, all estimators have negligible bias, which diminishes with an increasing sample size. The finite-sample biases are slightly larger, relative to the variability of the estimators in the case of binary outcomes than for continuous outcomes, resulting in lower coverages (around 91%) for the confidence intervals. The detailed empirical bias, Monte Carlo standard error (MCSE), and average estimated standard error (ASE) for the IPW, regression, and DR estimators are given in the Supplementary Material. In agreement with

Table 1. Empirical coverage rates based on 95% Wald-type confidence intervals for both binary and continuous outcomes.

	Binary $Y$		Cont. $Y$	
sample size ( $n$ )	1,000	5,000	1,000	5,000
(i) both $\pi$ and $\mu$ are correct				
$\hat{\psi}^{ipw}$	0.86	0.90	0.96	0.95
$\hat{\psi}^{reg}$	0.84	0.92	0.97	0.95
$\hat{\psi}^{DR}$	0.85	0.91	0.97	0.96
(ii) only $\pi$ is correct				
$\hat{\psi}^{ipw}$	0.86	0.90	0.96	0.95
$\hat{\psi}^{reg}$	0.79	0.60	0.39	0.00
$\hat{\psi}^{DR}$	0.86	0.91	0.97	0.95
(iii) only $\mu$ is correct				
$\hat{\psi}^{ipw}$	0.78	0.53	0.39	0.00
$\hat{\psi}^{reg}$	0.84	0.92	0.97	0.95
$\hat{\psi}^{DR}$	0.85	0.92	0.96	0.96

The coverage was evaluated under three scenarios: (i) both the outcome regression and the extended propensity score are correctly specified; (ii) only the extended propensity score is correct, and (iii) only the outcome regression model is correct.

our theoretical results, the IPW and regression estimators are biased, with poor empirical coverages when the extended propensity score or the outcome model, respectively, is misspecified. The DR estimator performs well in terms of bias and coverage when either model is misspecified but the other is correct. When all models are correctly specified, the relative efficiency of the locally semiparametric efficient estimator compared with that of the DR estimator of  $\eta$  and  $\psi$  are 0.840 and 0.810, respectively, based on Monte Carlo standard errors at a sample size  $n = 5,000$ . This shows that a substantial efficiency gain may be possible at the intersection submodel when using the locally efficient score.

Simulations for a continuous outcome were conducted similarly as for the binary outcome as follows:

Step 1\*: Covariates  $C_1$  and  $C_2$  were generated as in Step 1,  $Z$  was generated from the model  $\text{logit } \Pr(Z = 1|C) = 0.7 + 0.8C_1 - C_2$ ,  $Y_0, Y_1$  were generated from models  $Y_0|Z, C \sim N(0.5 + C_1 + 3C_2, 1)$  and  $Y_1|Z, C \sim N(1.1 - 1.3C_1, 1)$ , respectively,  $A$  was generated from  $\text{logit } \Pr(A = 1|Y_0, Z, C) = -0.2 - 3Z - 3C_1 + 0.3Y_0 + 4C_1Z$ , and  $Y = Y_0(1 - A) + Y_1A$ .

Step 2\*: See Step 2.

Step 3\*: As in Step 3, except that the following regression outcome models were fitted to the data:

$$E\{Y \exp(\eta Y) | A = 0, Z, C; \xi\} = \xi_1 + \xi_2 C_1 + \xi_3 C_2 + \xi_4 Z + \xi_5 C_1 Z + \xi_6 C_2 Z + \xi_7 C_1 C_2 + \xi_8 C_1 C_2 Z, \quad (5.5)$$

$$E\{\exp(\eta Y) | A = 0, Z, C; \xi\} = \xi_9 + \xi_{10} C_1 + \xi_{11} C_2 + \xi_{12} Z + \xi_{13} C_1 Z + \xi_{14} C_2 Z + \xi_{15} C_1 C_2 + \xi_{16} C_1 C_2 Z. \quad (5.6)$$

Step 4\*: As in Step 4, except that (5.2) was replaced by (5.5) and (5.6).

Step 5\*: See Step 5.

A simulation for a continuous outcome under scenario (ii) was carried out similarly to that for scenario (i), except that (5.1) was replaced by (5.3). For scenario (iii), the potential outcome models (5.5) and (5.6) were replaced with the linear models

$$E\{Y \exp(\eta Y) | A = 0, Z, C; \xi\} = \xi_1 + \xi_2 C_1 + \xi_4 Z, \quad (5.7)$$

$$E\{\exp(\eta Y) | A = 0, Z, C; \xi\} = \xi_9 + \xi_{10} C_1 + \xi_{12} Z. \quad (5.8)$$

We use the R package `nleqslv` (Hasselmann (2014)) to solve the nonlinear estimating equations.

We verify in Example A.1 of the Supplementary Material that  $\psi$  is identified from the observed data. The simulation results for 1,000 Monte Carlo samples are reported in Figure 2 (see also Table A.3 in the Supplementary Material), and the empirical coverage rates for the 95% Wald-type confidence intervals are presented in Table 1. The results are similar to the those for the binary outcome. Under a correct model specification, all estimators have negligible bias, which diminishes with an increasing sample size. The IPW and OR estimators are biased, with poor empirical coverages when the corresponding model is misspecified. The DR estimator performs well in terms of bias and coverage when either the extended propensity score or the outcome regression model is correctly specified.

## 6. Application

Since the 1980s, tax-deferred programs such as individual retirement accounts (IRAs) and the 401(k) plan have played an important role as a channel for personal savings in the United States. Aiming to encourage investment for fu-

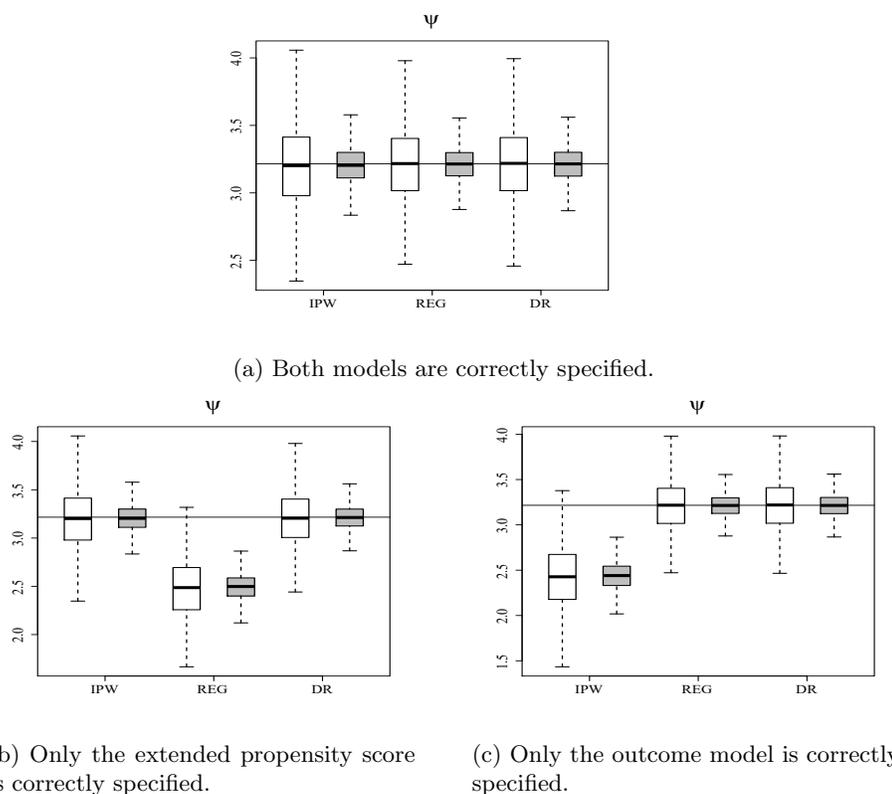


Figure 2. Performance of the IPW, OR, and DR estimators of  $\psi$  with continuous outcomes.

Note: In each box plot, the true value  $\psi_0$  is marked by the horizontal lines, white boxes are for  $n = 1,000$ , and grey boxes are for  $n = 5,000$ .

ture retirement, the 401(k) plan offers tax deductions on deposits into retirement accounts and a tax-free accrual of interest. The 401(k) plan shares similarities with IRAs in that both are deferred compensation plans for wage earners; however, the 401(k) plan is only provided by employers. The study includes 9,275 people. Once offered the 401(k) plan, individuals decide whether to participate in the program. However, participants usually have a stronger preference for savings, which suggests the presence of a selection bias. This was addressed as individual heterogeneity by Abadie (2003), and it has been pointed out that a simple comparison of personal savings between participants and non-participants may yield results that are biased upward. It has also been postulated that for a given income, 401(k) eligibility is unrelated to individuals' preferences for savings,

and thus can be used as an instrument for participation in the program (Poterba and Venti (1994); Poterba, Venti and Wise (1995)). The complier causal effect for the 401(k) plan was studied by Abadie (2003). Here, we reanalyze these data to illustrate the proposed estimators of the marginal ETT.

We illustrate the methods in the context of a dichotomous outcome, defined as the indicator that a person falls in the first quartile of net savings of the observed sample (equal to  $-\$500$ ). The treatment variable is a binary indicator of participation in a 401(k) plan, and the IV is a binary indicator of 401(k) eligibility. The covariates are standardized log family income ( $\log_{10}(\text{income}) - 4.5$ ), standardized age ( $\text{age} - 41$ ) and its square, marital status, and family size. Age ranges from 25 to 64 years, marital status is a binary indicator variable, and family size ranges from 1 to 13 people. These covariates are thought to be associated with unobserved preferences for savings. For a family that participated in the 401(k) program, let  $\psi = E(Y_0|A = 1)$  denote the probability that they would have had net financial assets above the first quartile, had they been forced not to participate in the program. The  $\text{ETT} = E(Y_1 - Y_0|A = 1)$  is the effect of the 401(k) plan on the difference scale for the probability of a family's net financial assets being above the first quartile among participants. Equivalently, the ETT can also be interpreted as the effect of an intervention in reducing a person's risk for poor savings performance, as measured by falling below the first quartile of the empirical distribution of savings for the sample. Before implementing our IV estimators, we first obtained a standard IPW estimator of the ETT under an assumption of no unmeasured confounding; that is,  $\hat{\psi}_0^{ipw}$  is defined as  $\hat{\psi}^{ipw}$ , with  $\alpha = 0$ . Thus, the propensity score was modeled as:

$$\text{logit Pr}(A = 1|Z, C) = 1 + Z + \log(\text{income}) + \text{married} + \text{age} + \text{fsize} + \text{age}^2,$$

and estimated using the standard maximum likelihood. The IPW estimate of  $\psi$  was  $\hat{\psi}_0^{ipw} = 0.688$ , with standard error (SE) 0.014, where SE was evaluated using the sandwich estimator, accounting for all sources of variability. In comparison, the estimator based on the empirical estimate of  $E(Y|A = 1)$  was 0.883 (SE = 0.006). Thus, the estimate of ETT was  $\widehat{\text{ETT}} = 0.194$  (SE = 0.016), which suggests that the 401(k) plan may have a significant effect on increasing the family net financial assets among participants.

However, this result may be spurious, owing to the suspicion that even after controlling for observed covariates, there may still exist unmeasured factors that confound the relationship between the 401(k) plan and the family net financial

assets. Assuming assumptions (IV.1)–(IV.2) and condition 1, we applied the methods proposed in Section 4 to estimate the ETT in the presence of unmeasured confounders. The following parametric models were considered:

$$\text{logit } \Pr(Z = 1|C) = 1 + \log(\text{income}) + \text{married} + \text{age} + \text{fsize} + \text{age}^2,$$

$$\text{logit } \Pr(Y = 1|A = 0, Z, C) = 1 + Z + \log(\text{income}) + \text{married} + \text{age} + \text{fsize} + \text{age}^2.$$

We specified the selection bias function as in (4.8). Thus, the selection bias function were assumed to depend on  $Y_0$  linearly. Possible deviations from this simple model was explored by allowing for potential interactions of  $Y_0$  with observed covariates in the extended propensity score. Thus, we posited the following parametric model for the extended propensity score, which satisfies the identifying condition 1 as a submodel of the separable model:

$$\text{logit } \Pr(A = 1|Y_0, Z, C) = 1 + Z + Y_0 + \log(\text{income}) + \text{married} + \text{age} + \text{fsize} + \text{age}^2.$$

Table 2 reports the point estimates and estimated standard errors for the IV, extended propensity score, and outcome regression models. Although the DR estimator also involves an outcome regression model among the unexposed, it is the same model required for the regression estimator. Therefore, these estimates are repeated only once. The instrument is strongly associated with family income (log OR = 2.823, SE = 0.106), age (log OR = 0.007, SE = 0.002), and age square (log OR = -0.002, SE =  $2e^{-4}$ ). The selection bias parameter was estimated to be 0.320 (SE = 0.115) by the IPW, 0.385 (SE = 0.135) by the OR and 0.280 (SE = 0.101) by the DR estimation. This provides strong evidence that unmeasured confounding may be present, and that a stronger saving preference means a person is more likely to participate in the 401(k) plan. All three estimators of the marginal ETT agree: they are significant, but with a smaller Z-score value than when the selection bias is ignored (e.g., the IPW estimator suggests  $\widehat{\text{ETT}} = 0.134$ , SE = 0.013). The efficient estimator for the selection bias parameter is 0.273, and for the ETT is 0.137, both in agreement with the other three estimators. Thus, we may conclude that even after adjustment for unobserved preferences for savings, the 401(k) plan can still increase net financial assets among participants.

These findings roughly agree with the results obtained by Abadie in the sense that the IV estimate corrects the observational estimate toward the null. However, it may be difficult to directly compare our findings to those of Abadie, who reported the compliers average treatment effect under a monotonicity assumption

Table 2. Point estimates and estimated SE (in parentheses) of IPW, OR and DR estimators for ETT of the 401(k) plan, as well as the parameters for the IV, extended propensity score, and outcome regression outcome models required by those estimators.

	IV model	IPW propensity	regression	DR propensity
Intercept	-0.180 (0.058)	-8.685 (1.832)	1.307 (0.073)	-8.629 (1.796)
linc	2.695 (0.107)	1.626 (0.210)	0.618 (0.128)	1.633 (0.209)
age	0.007 (0.002)	-0.009 (0.005)	0.035 (0.003)	-0.009 (0.005)
fsize	-0.037 (0.019)	-0.004 (0.033)	-0.127 (0.022)	-0.005 (0.033)
marr	-0.145 (0.063)	-0.032 (0.108)	-0.133 (0.075)	-0.031 (0.108)
age <sup>2</sup>	-0.002 (2e-04)	0.001 (4e-04)	6e-04 (3e-04)	0.001 (4e-04)
Z		9.150 (1.820)	-0.210 (0.074)	9.126 (1.781)
$\alpha$		0.320 (0.115)	0.385 (0.135)	0.280 (0.101)
$\psi = E(Y_0 A = 1)$		0.749 (0.012)	0.746 (0.012)	0.750 (0.012)
ETT		0.134 (0.013)	0.137 (0.014)	0.132 (0.014)

of the IV-exposure relationship, and assumed no unmeasured confounding of this first-stage relation. Our approaches rely on neither assumption, but instead rely on condition 1, encoded in the functional form of the extended propensity score model for identification. In order to assess the robustness of the selection bias model, additional functional forms were explored. We considered adding to  $\alpha$  an interaction between  $Y_0$  and each of the covariate: log income, marriage status, and family size. However, there was no evidence in favor of any such interaction.

## 7. Discussion

In this study, we establish that access to an IV allows us to identify an association between the exposure to a treatment and the potential outcome when unexposed, which directly encodes the magnitude of the selection bias in the treatment due to confounding. We propose IPW, OR, and DR estimators for the treatment effect amongst treated individuals. Vansteelandt and Goetghebuer (2003) and Robins (1994) proposed identification and inference approaches under a no-current treatment value interaction assumption. Thus, their estimators remain consistent under the null hypothesis of no ETT. In contrast, the identification and inference approaches proposed here may be particularly valuable when an ITT analysis indicates a non-null treatment effect, in which case, Robins' identification assumption may be violated.

When condition 1 does not hold, the ETT is not identified. However, sharp bounds could be derived for the ETT under monotonicity and dominance assumptions (Huber, Laffers and Mellace (2017)). A sensitivity analysis should be

carried out to investigate the parameter estimates within the sharp bounds.

The proposed methods assume the treatment is binary. They can be generalized without much effort to a categorical treatment. However, when the treatment is continuous (e.g.,  $A$  is a treatment dose), then a parametric model for the treatment effect and a model for the density of  $A$  may be unavoidable for an estimation. We leave this as a topic for future research.

## Supplementary Material

The online Supplementary Material contains proofs of the propositions, proofs of the examples in the main text, and additional examples related to identifying the models. It also presents more derivations mentioned in the main text, and the presents derivations of the semiparametric efficiency theory.

## Acknowledgements

The content is solely the responsibility of the authors. Lan Liu was supported by NSF DMS 1916013. Professor Eric Tchetgen Tchetgen was supported by R01 AI032475, R21 AI113251, R01 ES020337, and R01 AI104459. Wang Miao was supported by the China Scholarship Council.

## References

- Abadie, A. (2003). Semiparametric instrumental variable estimation of treatment response models. *Journal of Econometrics* **113**, 231–263.
- Abadie, A., Angrist, J. and Imbens, G. (2002). Instrumental variables estimates of the effect of subsidized training on the quantiles of trainee earnings. *Econometrica* **70**, 91–117.
- Angrist, J. (1995). Using social security data on military applicants to estimate the effect of voluntary military service on earnings.
- Angrist, J. D., Imbens, G. W. and Rubin, D. B. (1996). Identification of causal effects using instrumental variables. *Journal of the American Statistical Association* **91**, 444–455.
- Barnard, J., Frangakis, C. E., Hill, J. L. and Rubin, D. B. (2003). Principal stratification approach to broken randomized experiments: A case study of school choice vouchers in New York City. *Journal of the American Statistical Association* **98**, 299–323.
- Clarke, P. S., Palmer, T. M. and Windmeijer, F. (2015). Estimating structural mean models with multiple instrumental variables using the generalised method of moments. *Statistical Science* **30**, 96–117.
- Frangakis, C. E., Brookmeyer, R. S., Varadhan, R., Safaeian, M., Vlahov, D. and Strathdee, S. A. (2004). Methodology for evaluating a partially controlled longitudinal treatment using principal stratification, with application to a needle exchange program. *Journal of the American Statistical Association* **99**, 239–249.
- Goldberger, A. S. (1972). Structural equation methods in the social sciences. *Econometrica*:

- Journal of the Econometric Society* **40**, 979–1001.
- Greenland, S. and Robins, J. M. (1986). Identifiability, exchangeability, and epidemiological confounding. *International Journal of Epidemiology* **15**, 413–419.
- Hasselmann, B. (2014). *nleqslv: Solve Systems of Non Linear Equations*. R package version 2.1.1.
- Heckman, J. (1997). Instrumental variables: A study of implicit behavioral assumptions used in making program evaluations. *Journal of Human Resources* **32**, 441–462.
- Heckman, J. J., Ichimura, H. and Todd, P. E. (1997). Matching as an econometric evaluation estimator: Evidence from evaluating a job training programme. *The review of economic studies* **64**, 605–654.
- Heckman, J. J., Ichimura, H. and Todd, P. (1998). Matching as an econometric evaluation estimator. *The Review of Economic Studies* **65**, 261–294.
- Hernán, M. A. and Robins, J. M. (2006). Instruments for causal inference: an epidemiologist's dream?. *Epidemiology* **17**, 360–372.
- Huber, M., Laffers, L. and Mellace, G. (2017). Sharp IV bounds on average treatment effects on the treated and other populations under endogeneity and noncompliance. *Journal of Applied Econometrics* **32**, 56–79.
- Imbens, G. W. and Angrist, J. D. (1994). Identification and estimation of local average treatment effects. *Econometrica: Journal of the Econometric Society* **62**, 467–475.
- Matsouaka, R. A. and Tchetgen Tchetgen, E. J. (2014). Likelihood based estimation of logistic structural nested mean models with an instrumental variable.
- Miettinen, O. S. (1974). Proportion of disease caused or prevented by a given exposure, trait or intervention. *American Journal of Epidemiology* **99**, 325–332.
- Newey, W. K. and McFadden, D. (1994). Large sample estimation and hypothesis testing. *Handbook of Econometrics* **4**, 2111–2245.
- Ogburn, E. L., Rotnitzky, A. and Robins, J. M. (2015). Doubly robust estimation of the local average treatment effect curve. *Journal of the Royal Statistical Society, Series B (Statistical Methodology)* **77**, 373–396.
- Pearl, J. (1995). On the testability of causal models with latent and instrumental variables. In *Proceedings of the Eleventh Conference on Uncertainty in Artificial Intelligence*, 435–443. Morgan Kaufmann Publishers Inc.
- Poterba, J. M. and Venti, S. F. (1994). 401 (k) plans and tax-deferred saving. In *Studies in the Economics of Aging*, 105–142. University of Chicago Press.
- Poterba, J. M., Venti, S. F. and Wise, D. A. (1995). Do 401 (k) contributions crowd out other personal saving? *Journal of Public Economics* **58**, 1–32.
- Robins, J. M. (1989). The analysis of randomized and non-randomized AIDS treatment trials using a new approach to causal inference in longitudinal studies. *Health Service Research Methodology: A Focus on AIDS*, 113–159.
- Robins, J. M. (1994). Correcting for non-compliance in randomized trials using structural nested mean models. *Communications in Statistics-Theory and Methods* **23**, 2379–2412.
- Robins, J. M. and Ritov, Y. (1997). Toward a curse of dimensionality appropriate (CODA) asymptotic theory for semi-parametric models. *Statistics in Medicine* **16**, 285–319.
- Rotnitzky, A. and Robins, J. M. (1997). Analysis of semi-parametric regression models with non-ignorable non-response. *Statistics in Medicine* **16**, 81–102.
- Robins, J. M. and Rotnitzky, A. (2004). Estimation of treatment effects in randomised trials with

- non-compliance and a dichotomous outcome using structural mean models. *Biometrika* **91**, 763–783.
- Robins, J. M., Rotnitzky, A. and Scharfstein, D. (2000). Sensitivity analysis for selection bias and unmeasured confounding in missing data and causal inference models. In *Statistical Models in Epidemiology, the Environment, and Clinical Trials*, 1–94. Springer.
- Tan, Z. (2006). Regression and weighting methods for causal inference using instrumental variables. *Journal of the American Statistical Association* **101**, 1607–1618.
- Tan, Z. (2010). Marginal and nested structural models using instrumental variables. *Journal of the American Statistical Association* **105**, 157–169.
- van der Vaart, A. (1998). *Asymptotic Statistics*. Cambridge University Press, Cambridge.
- Vansteelandt, S. and Goetghebeur, E. (2003). Causal inference with generalized structural mean models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **65**, 817–835.
- Varadhan, R. and Gilbert, P. (2009). BB: An R package for solving a large system of nonlinear equations and for optimizing a high-dimensional nonlinear objective function. *Journal of Statistical Software* **32**, 1–26.
- Wright, S. (1928). Appendix to the tariff on animal and vegetable oils. *New York: MacMillan. (1934), "The Method of Path Coefficients," Annals of Mathematical Statistics* **5**, 161–215.

School of Statistics, University of Minnesota, Minneapolis, MN 55455, USA.

E-mail: liux3771@umn.edu

Guanghua School of Management, Peking University, Beijing 100871, China.

E-mail: mwfy@gsm.pku.edu.cn

National University of Singapore, 119077, Singapore.

E-mail: stasb@nus.edu.sg

Harvard T.H. Chan School of Public Health, Harvard University, MA, 02115, USA.

E-mail: robins@hsph.harvard.edu

The Wharton School, University of Pennsylvania, PA, 19104, USA.

E-mail: ett@wharton.upenn.edu

(Received April 2017; accepted September 2018)