

# Propensity Score-based Methods Versus MTE-based Methods in Causal Inference: Identification, Estimation, and Application

Sociological Methods &amp; Research

2016, Vol. 45(1) 3-40

© The Author(s) 2014

Reprints and permission:

[sagepub.com/journalsPermissions.nav](http://sagepub.com/journalsPermissions.nav)

DOI: 10.1177/0049124114555199

[smr.sagepub.com](http://smr.sagepub.com)

Xiang Zhou<sup>1</sup> and Yu Xie<sup>1</sup>

## Abstract

Since the seminal introduction of the propensity score (PS) by Rosenbaum and Rubin, PS-based methods have been widely used for drawing causal inferences in the behavioral and social sciences. However, the PS approach depends on the ignorability assumption: there are no unobserved confounders once observed covariates are taken into account. For situations where this assumption may be violated, Heckman and his associates have recently developed a novel approach based on marginal treatment effects (MTEs). In this article, we (1) explicate the consequences for PS-based methods when aspects of the ignorability assumption are violated, (2) compare PS-based methods and MTE-based methods by making a close examination of their identification assumptions and estimation performances, (3) apply these two approaches in estimating the economic return to college using data from the National Longitudinal Survey of Youth (NLSY) of 1979 and discuss their discrepancies in results. When there is a sorting gain but no systematic baseline difference between treated and untreated units given observed covariates,

---

<sup>1</sup>Institute for Social Research, University of Michigan, Ann Arbor, MI, USA

## Corresponding Author:

Xiang Zhou or Yu Xie, Institute for Social Research, University of Michigan, 426 Thompson Street, Ann Arbor, MI 48104, USA.

Email: [xiangzh@umich.edu](mailto:xiangzh@umich.edu); [yuxie@umich.edu](mailto:yuxie@umich.edu)

PS-based methods can identify the treatment effect of the treated (TT). The MTE approach performs best when there is a valid and strong instrumental variable (IV). In addition, this article introduces the “smoothing-difference PS-based method,” which enables us to uncover heterogeneity across people of different PSs in both counterfactual outcomes and treatment effects.

### **Keywords**

causal effects, exclusion restriction, heterogeneity, ignorability, instrumental variable, marginal treatment effect, propensity score, selection bias

### **Introduction**

Since the seminal introduction of the propensity score (PS) by Rosenbaum and Rubin (1983), PS-based methods, including matching, stratification, and weighting, have become a mainstay strategy for drawing causal inferences in the behavioral and social sciences. By reducing a large array of confounding variables to a univariate measure that preserves all the relevant information of potential confounders, the PS provides a more effective tool than covariate adjustment does for eliminating confounder bias (Rosenbaum and Rubin 1984). Furthermore, social science researchers have recently utilized PS methods to study heterogeneous treatment effects across individuals with different propensities of being treated. For example, Brand and Xie (2010) recently found that those students who are least likely to obtain a college education benefit most from college.

Like all other attempts to resolve confounding problems in causal inference, the PS approach is by no means a panacea. The primary limitation of this approach lies in the impossibility of capturing unobserved individual and contextual confounders. In fact, the whole justification of PS-based methods hinges on the common “ignorability assumption”: through control of a given set of relevant observed covariates, treatment status is assumed to be independent of potential outcomes. This assumption is unverifiable, indeed unlikely to be true, in practice. For instance, economic theory predicts that attainment of college education may be selective because it may attract young persons who are motivated by economic gain from college education (Carneiro, Heckman, and Vytlacil 2011; Willis and Rosen 1979). This example illustrates the effect of “sorting on gain” that may not be captured by observed covariates such as family background and cognitive abilities.

Despite the aforementioned limitation, PS-based methods are still widely used by empirical researchers in a variety of disciplines. Not only is the PS approach simple and straightforward, but methods of addressing unobserved selection would require either additional data unavailable to the researcher or strong assumptions implausible in a research setting. However, the performance of PS-based methods is questionable when the ignorability assumption breaks down. Although sensitivity analysis is usually employed to assess the plausibility of findings (DiPrete and Gangl 2004; Harding 2003), systematic investigation is also needed to directly examine the consequences for PS-based methods when ignorability is violated. A related discussion can be found in Heckman and Navarro-Lozano (2004), who compared matching, instrumental variables (IVs), and control functions in the estimation of economic choice models. Blundell, Dearden, and Sianesi (2005) also compared least squares, matching, control functions, and IV from a methodological point of view within a common framework. More recently, Shadish, Clark, and Steiner (2008) explored the performances of ordinary least squares (OLS) adjustment and PS adjustment with different sets of predictors in an experimental setting. Inspired by these studies, this article aims to provide another examination of PS-based methods in a variety of plausible situations.

Our article goes beyond PS-based methods by evaluating a structural approach, developed by James Heckman and his associates, for situations in which the ignorability assumption is violated (Heckman, Urzua, and Vytlačil 2006a, 2006b; Heckman and Vytlačil 1999, 2001, 2005). Different from traditional IV-based methods, this approach is based on the building block of marginal treatment effect (MTE), which enables us to derive various parameters of interest within a single framework. However, MTE-based methods have not been widely used in empirical research (for a few exceptions, see Carneiro et al. 2011; Moffitt 2008; Tsai and Xie 2011), partly due to their complexity and demands on data. In fact, the aforementioned literature suggests that the utility of MTE hinges heavily on the validity of the exclusion restriction as well as the strength of IVs.<sup>1</sup> The properties of this approach are not yet well known, when either the exclusion restriction is violated or the IV is too “weak.”

In this article, we evaluate and compare the widely adopted PS-based methods and the less popular MTE-based approach, as follows. In the second section, we revisit population heterogeneity and two types of selection bias in causal inference, presenting PS-based methods and their implications in settings where the ignorability assumption is partially or completely violated. In particular, we propose a PS-based method by modeling counterfactual

outcomes as nonparametric functions of the PS, which we call the “smoothing-difference method.” In the third section, we introduce the MTE-based approach as a remedy for situations where the ignorability assumption may be violated. In the fourth section, we compare PS- and MTE-based methods by examining their identification assumptions and estimation performances. We use numerical simulation to explore (1) the relative efficiency of these two approaches when both ignorability and the exclusion restriction hold true and (2) the potential biases from using methods based on the two approaches when neither ignorability nor the exclusion restriction is guaranteed. In the fifth section, we illustrate both methods in analyzing the economic return to college education using a sample of white males from the National Longitudinal Survey of Youth of 1979 (NLSY) and discuss their discrepancies in results. In the sixth section, we conclude the article.

## **Population Heterogeneity, Ignorability, and PS-based Methods**

Population sciences, including economics, demography, epidemiology, psychology, and sociology, treat individual-level variation as a part of reality subject to scientific inquiry, rather than a mere nuisance or measurement error (Angrist and Krueger 1999; Ansari and Kamel 2000; Bauer and Curran 2003; Greenland and Poole 1988; Heckman 2001, 2005; Heckman and Robb 1985; Heckman and Vytlacil 2005; Lubke and Muthén 2005; Manski 2007; Moffitt 1996; Rothman and Greenland 1998; Winship and Morgan 1999; Xie 2007). The recognition of inherent individual-level heterogeneity has important consequences for research designs in the social sciences. Because individuals differ from one another and differ in their responses to a common treatment, results can vary widely depending on population composition.

The large methodological literature on causal inference using statistical methods recognizes the importance of and consequently allows for population heterogeneity (Heckman and Vytlacil 2005; Holland 1986; Manski 1995; Rubin 1974; Winship and Morgan 1999). Suppose that a population,  $U$ , is being studied. Let  $Y$  denote an outcome variable of interest (its realized value being  $y$ ) that is defined for each member in  $U$ . Let us define treatment as an externally induced intervention that can, at least in principle, be given to or withheld from a unit under study. For simplicity, we consider only dichotomous treatments and use  $D$  to denote the treatment status (its realized value being  $d$ ), with  $D = 1$  if a member is treated and  $D = 0$  if a member is not treated. Let subscript  $i$  represent the  $i$ th member in  $U$ . We further denote

$y_i^1$  as the  $i$ th member's potential outcome if treated (i.e., when  $d_i = 1$ ), and  $y_i^0$  as the  $i$ th member's potential outcome if untreated (i.e., when  $d_i = 0$ ). The framework for counterfactual reasoning in causal inference (Heckman 2005; Holland 1986; Manski 1995; Morgan and Winship 2007; Rubin 1974; Sobel 2000; Winship and Sobel 2004) states that we should conceptualize a treatment effect as the difference in potential outcomes associated with different treatment states for the same member in  $U$ :

$$\delta_i = y_i^1 - y_i^0, \quad (1)$$

where  $\delta_i$  represents the hypothetical treatment effect for the  $i$ th member.<sup>2</sup> The fundamental problem of causal inference (Holland 1986) is that, for a given unit  $i$ , we observe either  $y_i^1$  (if  $d_i = 1$ ) or  $y_i^0$  (if  $d_i = 0$ ), but not both. Given this fundamental problem, Holland describes two possible solutions, the "scientific solution" and the "statistical solution." The scientific solution capitalizes on homogeneity in assuming that all members in  $U$  are the same, in either the treated state or the control state:  $y_i^1 = y_j^1$  and  $y_i^0 = y_j^0$ , where  $j \neq i$  in  $U$ . This strong homogeneity assumption would enable a researcher to identify individual-level treatment effects by as few as two cases in  $U$ . However, as we discussed previously, pervasive heterogeneity across units is the norm rather than the exception in a population science. Thus, in general, the scientific solution has no practical value in the social and behavioral sciences.

### Quantities of Interest

For a population science, the statistical solution is a necessity. The statistical approach is to compute quantities of interest that reveal treatment effects only at the group level. For example, we may evaluate the average difference between a set of members in  $U$  that were randomly selected for treatment and another set of members that were randomly selected for control. This comparison yields a quantity that is called the average treatment effect (ATE):

$$\text{ATE} = E(Y^1 - Y^0).$$

While *ATE* is defined for the whole population, the researcher may wish to focus on and define a treatment effect for a well-defined subpopulation. In contexts of program evaluation, for example, researchers may be primarily interested in the treatment effect of the treated (TT; Heckman and Robb

1985), which refers to the average difference by treatment status among those individuals who are actually treated:

$$TT = E(Y^1 - Y^0|D = 1).$$

Although various statistical quantities of interest can easily be defined theoretically with the statistical “solution,” estimating these quantities in social research can be very difficult, due to two types of selection bias, a topic to be discussed subsequently.

### *Two Types of Selection Bias*

In the preceding subsection, we established the need to conduct group-level comparisons for causal inference, because causal inference is impossible at the individual level. However, due to population heterogeneity, there is no guarantee that the group that actually receives the treatment is comparable, in observed and particularly in unobserved contextual and individual characteristics, to the group that does not receive the treatment.<sup>3</sup> Individuals may self-select into treatment based on their anticipated monetary and nonmonetary benefits and costs of treatment. To see this, let us partition the total population  $U$  into the subpopulation of the treated  $U_1$  (for which  $D = 1$ ) and the subpopulation of the untreated  $U_0$  (for which  $D = 0$ ). We can thus decompose the expectations for the two counterfactual outcomes as follows:

$$E(Y^1) = E(Y^1|D = 1)P(D = 1) + E(Y^1|D = 0)P(D = 0),$$

and

$$E(Y^0) = E(Y^0|D = 1)P(D = 1) + E(Y^0|D = 0)P(D = 0).$$

The issue of selection stems from the scenario

$$E(Y^1|D = 1) \neq E(Y^1|D = 0) \neq E(Y^1), \quad (2)$$

and

$$E(Y^0|D = 1) \neq E(Y^0|D = 0) \neq E(Y^0). \quad (3)$$

Note that what we observe from data are  $\hat{E}(Y^1|D = 1)$ ,  $\hat{E}(Y^0|D = 0)$ ,  $\hat{P}(D = 1)$ , and  $\hat{P}(D = 0)$ . Due to inequalities (2) and (3), the simple-comparison estimator  $\hat{E}(Y^1|D = 1) - \hat{E}(Y^0|D = 0)$ , as a naive estimator for

ATE, may be contaminated by selection bias. Denoting this estimator by  $\hat{\beta}_{\text{naive}}$ , we can decompose its expectation as follows (Winship and Morgan 1999:667):

$$\begin{aligned}
 E\left(\hat{\beta}_{\text{naive}}\right) &= E\left(Y^1|D=1\right) - E\left(Y^0|D=0\right) \\
 &= E\left(Y^1 - Y^0|D=1\right) + E\left(Y^0|D=1\right) - E\left(Y^0|D=0\right) \quad (4) \\
 &= \text{TT} + E\left(Y^0|D=1\right) - E\left(Y^0|D=0\right) \\
 &= \text{ATE} + (\text{TT} - \text{ATE}) + E\left(Y^0|D=1\right) - E\left(Y^0|D=0\right).
 \end{aligned}$$

From equation (4), we see two sources of selection bias:

1. The difference in average outcome between the treatment and control groups if neither group receives treatment:  $E(Y^0|D=1) - E(Y^0|D=0)$ . We call this the “pretreatment heterogeneity bias” or “type I selection bias.”
2. The difference in average treatment effect between the treated group and the entire population. We call this the “treatment-effect heterogeneity bias” or “type II selection bias.” There is treatment-effect heterogeneity bias if and only if  $\text{TT} \neq \text{ATE}$ .

We now illustrate the two different sources of selection bias with two concrete examples. First, preschool children from poor families are selected into head start programs and thus would compare unfavorably to other children who do not attend head start programs without an adequate control for family socioeconomic resources (Xie 2000). Second, economic theory predicts that attainment of college education may be selective because it may attract young persons who are more motivated than their peers to gain from college education (Willis and Rosen 1979). While the first example reflects the importance of pretreatment heterogeneity bias that may be represented by “covariates” or “fixed effects,” the second example underscores the possibility of treatment-effect heterogeneity bias—sorting on the treatment effects—which might not be captured by “covariates” or “fixed effects.”

### *Ignorability and PS*

In observational studies, to overcome the two types of selection bias resulting from nonrandomness in treatment assignment, a natural idea is to control for observed pretreatment covariates. While it is not possible for a researcher to claim that he or she has controlled for all of the variables that may affect the outcome, it is more plausible to assume that the researcher has controlled for

almost all of the relevant pretreatment covariates that may affect *both* the treatment assignment *and* the outcome, a subset of all the variables affecting the outcome. In fact, only the covariates that meet the condition of affecting both the treatment assignment and the outcome may potentially confound the observed relationship between treatment and outcome (Rubin 1997). Thus, if we assume that all these *relevant* pretreatment variables are observed, the treatment status will be independent of potential outcomes through control of these covariates. This conditional independence assumption is called “ignorability,” “unconfoundedness,” or “selection on observables.” If we let  $\mathbf{X}$  be the vector of these observed covariates, the ignorability assumption states:

$$(Y^1, Y^0) \perp\!\!\!\perp D | \mathbf{X}. \quad (5)$$

Because we can never be sure after inclusion of which covariates relation (5) would hold true, the ignorability condition is always held as an assumption, indeed an unverifiable assumption. Substantive knowledge about the subject matter needs to be brought in before a researcher can entertain the ignorability assumption. Measurement of theoretically meaningful confounders makes ignorability tentatively plausible but not necessarily true. However, the researcher can always consider the ignorability assumption and then assess its plausibility in a concrete setting through sensitivity or auxiliary analyses (Cornfield et al. 1959; DiPrete and Gangl 2004; Harding 2003; Rosenbaum 2002; Xie and Wu 2005).

If the ignorability assumption of equation (5) holds true, we can change inequalities (2) and (3) into two equations by conditioning on  $\mathbf{X}^4$ :

$$E(Y^1 | D = 1, \mathbf{X}) = E(Y^1 | D = 0, \mathbf{X}) = E(Y^1 | \mathbf{X}); \quad (6)$$

$$E(Y^0 | D = 1, \mathbf{X}) = E(Y^0 | D = 0, \mathbf{X}) = E(Y^0 | \mathbf{X}). \quad (7)$$

For now, we are concerned only with identification and postpone inference issues to a later discussion. Let us define quantities of interest for causal inference, conditioning on  $\mathbf{X}$ , as follows:

$$\text{ATE}(\mathbf{X}) = E(Y^1 - Y^0 | \mathbf{X}),$$

$$\text{TT}(\mathbf{X}) = E(Y^1 - Y^0 | D = 1, \mathbf{X}),$$

Similarly, we define the naive estimator conditioning on  $\mathbf{X}$  as:



$$\hat{\beta}_{\text{naive}}(\mathbf{X}) = \hat{E}(Y^1|D = 1, \mathbf{X}) - \hat{E}(Y^0|D = 0, \mathbf{X}).$$

Then equations (6) and (7) imply the following identity,

$$E\left(\hat{\beta}_{\text{naive}}(\mathbf{X})\right) = \text{ATE}(\mathbf{X}) = \text{TT}(\mathbf{X}). \quad (8)$$

As a result, the ignorability assumption enables the naive estimator to identify both ATE and TT through control of  $\mathbf{X}$ . Conditioning on  $\mathbf{X}$ , however, can be difficult in applied research due to the “curse of dimensionality.” Rosenbaum and Rubin (1983, 1984) show that, when the ignorability assumption holds true, it is sufficient to condition on the PS as a function of  $\mathbf{X}$ . That is to say, relation (5) implies:

$$(Y^1, Y^0) \perp\!\!\!\perp D | P(D = 1 | \mathbf{X}),$$

where  $P(D = 1 | \mathbf{X})$  is the PS, the conditional probability of treatment given all the relevant information in covariates  $\mathbf{X}$ . In other words, only through the PS  $P(D = 1 | \mathbf{X})$  may covariates  $\mathbf{X}$  confound the observed relationship between treatment  $D$  and outcome  $Y$ . In empirical settings, however, the PS first needs to be estimated. Because a fully nonparametric estimation of the PS would also suffer from the curse of dimensionality, the estimation is conventionally accomplished by a logit or probit regression. From here on, we denote by  $p$  the PS and by  $\hat{\beta}_{\text{naive}}(p)$  the naive estimator of treatment effect conditional on  $p$ . Note that  $\hat{\beta}_{\text{naive}}(p)$  here is a function of  $p$  but not necessarily a linear function of  $p$ . There are a variety of methods for constructing  $\hat{\beta}_{\text{naive}}(p)$ , such as matching and stratification (see Morgan and Winship 2007:chap. 4). Subsequently, we introduce a new PS-based method by modeling counterfactual outcomes as nonparametric functions of the PS, which we call the “smoothing-difference method.”<sup>5</sup>

### *The Smoothing-difference PS-based Method*

If treatment effects are heterogeneous, there are many possible ways to characterize the heterogeneity (Heckman and Robb 1985; Pearl 2009; Winship and Morgan 1999; Xie, Brand, and Jann 2012). The basic idea of the smoothing-difference method is fitting two nonparametric functions for  $E(Y^1|p)$  and  $E(Y^0|p)$  and taking their differences as estimates of treatment effects. Specifically, it consists of the following four steps:

1. From observed  $Y^1$  and  $Y^0$  and estimated PS  $\hat{p}$ , fit two univariate functions,  $f_1(p)$  and  $f_0(p)$ , to approximate  $E(Y^1|p)$  and  $E(Y^0|p)$ .

2. Use  $f_1(p)$  and  $f_0(p)$  to predict counterfactual outcomes  $Y_i^1$  and  $Y_i^0$  for each individual  $i$  in the sample.
3. Obtain the estimated treatment effect  $\delta_i$  for each individual  $i$  by taking the difference between the predicted counterfactual outcomes.
4. Average estimated  $\delta_i$  over the entire sample as the estimate of ATE or over a specific subsample to estimate a corresponding group-level causal effect. For example, we may estimate TT by averaging  $\delta_i$  over those subjects who are actually treated.

Note that in the first step,  $f_1(p)$  and  $f_0(p)$  could be fitted via different estimation methods. One simple possibility is to fit two linear models of  $Y^1$  and  $Y^0$  on  $p$  through OLS. However, this strategy imposes too strong a parametric assumption concerning the relationship between the outcome variables and the PS. As we will see in the Results section, empirical data could suggest a nonmonotone treatment effect of college education as a function of PS. In this case, imposing a linear structure would mask interesting patterns of treatment effect heterogeneity, making it impossible to identify population subgroups that benefit differently from the treatment. Therefore, we propose to fit two smoothing splines (Hastie, Tibshirani, and Friedman 2008) of  $Y^1$  and  $Y^0$  on  $p$  with the smoothing parameter determined by generalized cross validation or some other criterion,<sup>6</sup> an approach that we will adopt in our simulation studies in the fourth section as well as the analysis of NLSY data in the fifth section.

Other PS-based methods include matching, stratification, and weighting, all of which have been widely used in empirical research. Compared to these traditional methods, our smoothing-difference method has three distinct advantages.<sup>7</sup> First and foremost, the research interest may lie in the trend of  $\Delta f(p)$ , that is,  $f_1(p) - f_0(p)$ , which characterizes how the treatment effect varies across subjects with different propensities of being treated (Brand and Xie 2010; Xie et al. 2012; Xie and Wu 2005). Such a pattern may identify population subgroups that benefit the most, or the least, from treatment, thus offering meaningful policy implications. For example, many countries are now rapidly expanding the size of college enrollment. The expansion of college education would be more effective if it were targeted at those individuals who are likely to benefit most from attending college. As the propensity of attending college could be directly estimated from individual characteristics for every potential college goer, the expected returns to college could likewise be estimated. Hence, a PS-based cost-benefit analysis would enable policy makers to fine-tune their strategies in expanding higher education. Furthermore, the researcher may be interested in the trends of the

outcome as a function of the PS among either treated or untreated subjects, that is,  $f_0(p)$  or  $f_1(p)$ , which could not be extracted from results using matching or weighting methods. As we will illustrate in the Results section, these trends can enhance our understanding of the social processes that may be masked by too exclusive a focus on the estimation of causal effects. In addition, after we pool information through nonparametric regressions across adjacent cases within either treated or untreated groups, we are able to derive the estimated treatment effect  $\delta_i$  for each individual  $i$ , before we calculate any group-level treatment effect  $\delta$  through averaging over an appropriate subsample.

### Decomposition of Ignorability

Equations (4) and (8) reveal that, under the ignorability assumption, controlling for observed covariates eliminates both types of selection bias. This suggests that the ignorability assumption should contain two components corresponding to the two types of bias. Indeed, ignorability expressed as relation (5) can be rewritten as

$$(Y^0, Y^1 - Y^0) \perp\!\!\!\perp D | \mathbf{X}.$$

This expression reveals the two conditions underlying the ignorability assumption:

**Condition 1:**

$$Y^0 \perp\!\!\!\perp D | \mathbf{X},$$

that is, given observed covariates, treatment status is independent of the baseline outcome.

**Condition 2:**

$$Y^1 - Y^0 \perp\!\!\!\perp D | \mathbf{X},$$

that is, given observed covariates, treatment status is independent of the treatment effect. We may call condition 1 the *ignorability of type I selection bias* and condition 2 the *ignorability of type II selection bias*.<sup>8</sup> The ignorability assumption, in essence, means the ignorability of both type I selection bias and type II selection bias.

As mentioned in the Ignorability and PS subsection, Rosenbaum and Rubin (1983) derive the sufficiency of PS under the assumption of complete ignorability for eliminating confounding bias:

If  $(Y^1, Y^0) \perp\!\!\!\perp D | \mathbf{X}$ , then  $(Y^1, Y^0) \perp\!\!\!\perp D | P(D = 1 | \mathbf{X})$ .

Indeed, the proof provided by Rosenbaum and Rubin (1983) implies the following two propositions:

**Proposition 1:**

$$\text{If } Y^0 \perp\!\!\!\perp D | \mathbf{X}, \text{ then } Y^0 \perp\!\!\!\perp D | P(D = 1 | \mathbf{X}); \quad (9)$$

**Proposition 2:**

$$\text{If } Y^1 - Y^0 \perp\!\!\!\perp D | \mathbf{X}, \text{ then } Y^1 - Y^0 \perp\!\!\!\perp D | P(D = 1 | \mathbf{X}). \quad (10)$$

As before, we denote by  $p$  the PS  $P(D = 1 | \mathbf{X})$  and by  $\hat{\beta}_{\text{naive}}(p)$  the naive estimator of treatment effect conditional on  $p$ . In light of equation (4), the expectation of  $\hat{\beta}_{\text{naive}}(p)$  could be decomposed as:

$$\begin{aligned} E\left(\hat{\beta}_{\text{naive}}(p)\right) &= \text{ATE}(p) + \text{TT}(p) - \text{ATE}(p) + E(Y^0 | D = 1, p) \\ &\quad - E(Y^0 | D = 0, p). \end{aligned} \quad (11)$$

Equation (11), combined with Propositions 1 and 2, shows that the naive estimator  $\hat{\beta}_{\text{naive}}(p)$  identifies different quantities under different conditions. First, if condition 1 holds true, type I selection bias thus becomes ignorable, that is,

$$E(Y^0 | D = 1, p) = E(Y^0 | D = 0, p).$$

In this scenario, we have

$$E\left(\hat{\beta}_{\text{naive}}(p)\right) = \text{TT}(p).$$

Second, if condition 2 holds true, type II selection bias becomes ignorable, that is,

$$E(Y^1 - Y^0 | p, D = 1) = E(Y^1 - Y^0 | p),$$

or

$$\text{TT}(p) = \text{ATE}(p).$$

In this scenario, we have

$$E\left(\hat{\beta}_{\text{naive}}(p)\right) = \text{ATE}(p) + E(Y^0|D = 1, p) - E(Y^0|D = 0, p).$$

As a result, we conclude that:

1. If the ignorability of type I selection bias (i.e., condition 1) holds true, the naive estimator conditional on the PS,  $\hat{\beta}_{\text{naive}}(p)$ , confronts only type II selection bias (treatment-effect heterogeneity bias), that is,  $\text{TT}(p) \neq \text{ATE}(p)$ . However, if our quantity of interest is TT rather than ATE,  $\hat{\beta}_{\text{naive}}(p)$  is an unbiased estimator.
2. If the ignorability of type II selection bias (i.e., condition 2) holds true, the naive estimator conditional on the PS,  $\hat{\beta}_{\text{naive}}(p)$ , is subject only to type I selection bias (pretreatment heterogeneity bias), that is,  $E(Y^0|D = 1, p) \neq E(Y^0|D = 0, p)$ .<sup>9</sup>

## Marginal-treatment-effect-based Approach

So far, we have seen that PS-based methods are subject to biases when the ignorability assumption is violated. Unfortunately, the ignorability assumption can never be verified. What recourse is available to a researcher who finds the ignorability assumption implausible in a research setting? In this section, we introduce the MTE approach developed by Heckman and his associates (Heckman and Vytlačil 1999, 2001, 2005; Heckman et al. 2006a, 2006b). Essentially, these researchers show that consistent estimates of a wide range of treatment parameters (including ATE and TT) can be obtained through different weighted averages of MTE (Björklund and Moffitt 1987). MTE could be estimated either parametrically or semiparametrically. In the following, we briefly review this class of methods, under the heading “MTE-based approach.”

### The Definition of MTE

It is most convenient to explicate the MTE-based approach with three equations: outcome equations under two counterfactual regimes ( $D = 0, D = 1$ ) and a treatment selection equation. In writing out each of the equations, we assume separability of the outcome variable into a structural component due to a linear function of pretreatment covariates and residual components due to unobserved variables<sup>10</sup>:

$$Y^0 = \beta_0' \mathbf{X} + \epsilon,$$

$$Y^1 = \beta_1' \mathbf{X} + \epsilon + \eta.$$

Here, corresponding to our decomposition of the two sources of ignorability, the error term  $\epsilon$  captures the unobserved factors that affect only the baseline outcome, while the error term  $\eta$  represents the unobserved factors that affect only units that are treated. Thus, the treatment effect contains both the structural component ( $\beta'_0\mathbf{X}$  vs.  $\beta'_1\mathbf{X}$ ) and the residual component  $\eta$ . This setup changes equation (1) so that heterogeneous treatment effect can be written in a structural form:

$$\delta(\mathbf{X}) = Y^1(\mathbf{X}) - Y^0(\mathbf{X}) = (\beta_1 - \beta_0)' \mathbf{X} + \eta.$$

Note that the treatment effect  $\delta$  depends on covariates  $\mathbf{X}$ . If we denote by  $Y$  the observed outcome and by  $D$  the treatment status, the previous model could be written in the notation of switching regression models:

$$\begin{aligned} Y &= (1 - D)Y^0 + DY^1 \\ &= Y^0 + D(Y^1 - Y^0) \\ &= \beta'_0\mathbf{X} + (\beta_1 - \beta_0)' \mathbf{X}D + \epsilon + \eta D. \end{aligned}$$

We further specify a model for selection into treatment. Let  $D^*$  be the latent tendency to be treated:

$$\begin{aligned} D^* &= \gamma' \mathbf{Z} - V, \\ D &= 1(D^* > 0). \end{aligned} \tag{12}$$

Here,  $\mathbf{Z}$  is a vector of variables that predict the treatment probability,  $\gamma$  is a vector of coefficients, and  $V$  is a latent random variable representing disturbance. As in standard regression models, we assume that error terms ( $\epsilon, \eta, V$ ) have zero means and are jointly independent of  $\mathbf{X}$  and  $\mathbf{Z}$ . In practice,  $\mathbf{Z}$  consists of all observed predictors of treatment probability, including all the components in  $\mathbf{X}$  as well as some additional variables that predict only the treatment status  $D$ . These additional variables are called instrumental variables (IVs). The assumption that IVs affect only the treatment status  $D$  but not the outcome variable  $Y$  directly is called the exclusion restriction.

We can easily rewrite the treatment selection model, equation (12), in the following form:

$$\begin{aligned} \tilde{D}^* &= p(\mathbf{Z}) - U_D \\ D &= 1(\tilde{D}^* > 0), \end{aligned}$$

where  $p(\mathbf{Z}) = P(D = 1|\mathbf{Z}) = F_V(\gamma'\mathbf{Z})$  denotes the PS of being treated given  $\mathbf{Z}$  and  $U_D = F_V(V)$  follow a standard uniform distribution on  $[0,1]$ .  $U_D$

represents a catch-all unobserved selection component, interpretable as the level of unobserved resistance to receiving treatment, normalized between 0 and 1. We see that  $\mathbf{Z}$  enters the treatment selection model only through the PS  $p(\mathbf{Z})$ .

Based on the earlier specification of the outcome models and treatment selection model, the MTE is defined as follows:

$$\begin{aligned} \text{MTE}(\mathbf{x}, u_D) &= E(\delta | \mathbf{X} = \mathbf{x}, U_D = u_D) \\ &= E\left((\beta_1 - \beta_0)' \mathbf{x} + \eta | \mathbf{X} = \mathbf{x}, U_D = u_D\right) \quad (13) \\ &= (\beta_1 - \beta_0)' \mathbf{x} + E(\eta | V = F_V^{-1}(u_D)). \end{aligned}$$

Thus, MTE is essentially the expected treatment effect conditional on observed covariates  $\mathbf{X} = \mathbf{x}$  as well as the unobserved selection component  $U_D = u_D$ .

As mentioned previously, Heckman et al. (2006a, 2006b) have shown that group-level treatment effects such as ATE and TT can be expressed as weighted averages of  $\text{MTE}(\mathbf{x}, u_D)$ .<sup>11</sup> However, the estimation of  $\text{MTE}(\mathbf{x}, u_D)$  is not straightforward since neither the counterfactual outcome nor the latent variable  $u_D$  is observed. Now we briefly sketch the two approaches to estimating MTE: (1) the parametric method and (2) the semiparametric method.

### *Parametric and Semiparametric Estimation of MTE*

First, we write out the expectation of the observed outcome  $Y$  given covariates  $\mathbf{X} = \mathbf{x}$  and the PS  $p(\mathbf{Z}) = p$ :

$$\begin{aligned} E(Y | \mathbf{X} = \mathbf{x}, p(\mathbf{Z}) = p) &= E\left(\beta_0' \mathbf{X} + (\beta_1 - \beta_0)' \mathbf{X}D + \epsilon + \eta D | \mathbf{X} = \mathbf{x}, p(\mathbf{Z}) = p\right) \\ &= \beta_0' \mathbf{x} + (\beta_1 - \beta_0)' \mathbf{x}p + E(\eta | D = 1, p(\mathbf{Z}) = p)p \\ &= \beta_0' \mathbf{x} + (\beta_1 - \beta_0)' \mathbf{x}p + E(\eta | V < F_V^{-1}(p))p \\ &= \beta_0' \mathbf{x} + (\beta_1 - \beta_0)' \mathbf{x}p + \int_0^p E(\eta | V = F_V^{-1}(u_D)) du_D. \end{aligned} \quad (14)$$

Incorporating equation (13), the previous expression can be simplified:

$$E(Y | \mathbf{X} = \mathbf{x}, p(\mathbf{Z}) = p) = \beta_0' \mathbf{x} + \int_0^p \text{MTE}(\mathbf{x}, u_D) du_D.$$

Differentiating the previous equation with respect to  $p$ , we obtain MTE:

$$\text{MTE}(\mathbf{x}, p) = \frac{\partial E(Y|\mathbf{X} = \mathbf{x}, p(\mathbf{Z}) = p)}{\partial p}. \tag{15}$$

This expression relates  $\text{MTE}(\mathbf{x}, p)$  to  $E(Y|\mathbf{X} = \mathbf{x}, p(\mathbf{Z}) = p)$  and thus provides a possible route for estimating MTE.

In general, the third term in equation (14),  $\int_0^p E(\eta|V = F_V^{-1}(u_D))du_D$ , is an unknown function of  $p$ . However, if we assume that error terms  $(\epsilon, \eta, V)$  follow a joint Gaussian distribution  $N(0, \Sigma)$ ,  $E(Y|\mathbf{X} = \mathbf{x}, p(\mathbf{Z}) = p)$  would become a linear combination of  $\mathbf{x}$ ,  $\mathbf{x}p$ , and  $\phi(\Phi^{-1}(p))$ . Accordingly, the expression of MTE reduces to:

$$\text{MTE}(\mathbf{x}, u_D) = (\beta_1 - \beta_0)' \mathbf{x} + \sigma_{\eta V} \Phi^{-1}(u_D),$$

where  $\sigma_{\eta V}$  represents the covariance between  $\eta$  and  $V$ . With this parametric specification, we can estimate its unknown parameters  $(\beta_1, \beta_0, \sigma_{\eta V})$  via maximum likelihood (ML). This is the parametric MTE-based method, which is also called the “control function approach.” In fact, given the parametric assumption, identification of  $\text{MTE}(\mathbf{x}, u_D)$  is theoretically possible even without the presence of IVs.

In empirical settings, the assumption of joint normality is rarely justifiable. This motivated Heckman et al. (2006b) to develop a semiparametric method to identify  $\text{MTE}(\mathbf{x}, p)$ , using equation (15), after first estimating equation (14) under more flexible assumptions. Their method involves four steps:<sup>12</sup>

1. Fit local linear regressions of  $Y$ ,  $\mathbf{X}$ , and  $\mathbf{X}p$  on  $p$  and extract their residuals  $R_Y$ ,  $\mathbf{R}_X$ , and  $\mathbf{R}_{Xp}$ .
2. Regress  $R_Y$  on  $\mathbf{R}_X$  and  $\mathbf{R}_{Xp}$  using least squares to estimate the parametric component of equation (14), that is,  $\beta_0$  and  $\beta_1 - \beta_0$ , and denote its residuals by  $R_Y^*$ .
3. Regress  $R_Y^*$  on  $p$  using standard nonparametric techniques (such as local polynomial regression) to model the third term in equation (14) as well as its derivative, that is,  $E(\eta|V = F_V^{-1}(p))$ .
4. Construct  $\text{MTE}(\mathbf{x}, u_D)$ , expressed in equation (13), using  $\hat{\beta}_1 - \hat{\beta}_0$  from step 2 and the estimate of  $E(\eta|V = F_V^{-1}(p))$  from step 3.

Since this method capitalizes on the net relationship of  $E(Y|\mathbf{X} = \mathbf{x}, p(\mathbf{Z}) = p)$  with  $p(\mathbf{Z})$  after all covariates in  $\mathbf{X}$  are controlled for,



**Table 1.** Assumptions for Identification.

Method	Exclusion restriction (A valid IV)	Distributional form of error terms	Ignorability assumption
PS-based methods	No	No	Yes
Parametric MTE	No	Yes	No
Semiparametric MTE (LIV)	Yes	No	No

Note: For the parametric MTE-based method, multivariate normality is conventionally assumed for error terms. For PS-based methods, when the parameter of interest is TT, only the ignorability of type I selection bias is necessary. MTEs = marginal treatment effects; LIV = local instrumental variable; PS = propensity score; IV = instrumental variable; TT = treatment effect of the treated.

the presence of at least a valid IV in  $Z$  is indispensable for identification. For this reason, the semiparametric approach is also called the “local instrumental variable” (LIV) method. In spite of its flexibility, the LIV method has not been widely adopted in empirical research, due partly to its high data demand pertaining to IVs. Although a detailed discussion on the asymptotic variance of the ML estimator for the parametric MTE method can be found in Heckman (1979) and Puhani (2000), statistical properties of the semiparametric LIV method are not yet familiar to a wider research community.<sup>13</sup> This motivates us to evaluate the performance of MTE-based methods through numerical simulation.

## Evaluation of Different Methods

### Assumptions for Identification

The preceding discussion indicates that different methods require different assumptions to identify group-level treatment effects. Table 1 summarizes the assumptions that are required for the PS-based methods, the parametric MTE-based method, and the semiparametric MTE-based method to identify ATE and TT. Generally speaking, both the PS- and the MTE-based methods rely on strong and unverifiable assumptions, the former on ignorability, and the latter on exclusion restriction or the distribution of error terms. In particular, several facts deserve our attention. First, as we note in the second section, the applicability of PS-based methods is not limited to settings in which complete ignorability is satisfied. As long as the ignorability of type I selection bias is satisfied, that is, there is no systematic baseline difference between treated and untreated units given observed covariates, the PS-based models yield good estimates of TT, even in the presence of a heterogeneous treatment effect bias.

Second, in contrast to the semiparametric LIV approach, the parametric MTE-based method requires the assumption of joint normality for error

terms. Although in principle the parametric MTE method does not need an IV for identification, estimation based on the parametric distribution can be highly imprecise. In practice, availability of IVs that satisfy the exclusion restriction would greatly improve the precision of the ML estimation.<sup>14</sup> We shall study this facet of the MTE-based methods through numerical simulation in the next subsection.

In comparison with the PS-based methods, the MTE-based methods require a more explicit micro-level model. From the discussion in the Parametric and Semiparametric Estimation of MTE subsection, we see that the expression of MTE in equation (13) requires the separability of observables and unobservables in the outcome equations. Also, to estimate the parametric component of the model, we need to specify a functional form (such as linearity) to characterize the dependence of  $Y^0$  and  $Y^1$  on  $X$ . For the PS-based methods, however, the assumption of ignorability allows us to model counterfactual outcomes (or treatment effects) along the single dimension of PS. As this can be conducted in a purely nonparametric manner (as is commonly done in practice), the PS-based methods do not require an explicit model specification (except for the PS model).

Ultimately, a choice between the PS- and the MTE-based methods is driven by the plausibility of the ignorability assumption versus the exclusion restriction assumption. On one hand, if we suspect violation of ignorability but have valid IVs, the MTE-based approach is preferable to the PS-based methods. However, as the exclusion restriction is also a strong and unverifiable assumption, determining its degree of plausibility requires substantial knowledge in an actual research setting. On the other hand, if we have no satisfactory instruments that would satisfy the exclusion restriction but have sufficient information on relevant individual and contextual characteristics so that the ignorability assumption becomes plausible, the PS-based approach is a reasonable choice.

From the previous discussion, the general guideline is quite clear when one assumption is more plausible than the other. What, however, happens if the ignorability and the exclusion restriction assumptions are equally plausible—or equally implausible? In the rest of this section, we examine the performances of different methods for the following two scenarios: (1) when both the ignorability and the exclusion restriction assumptions hold true and (2) when both the ignorability and the exclusion restriction assumptions break down.

### *When Both Ignorability and the Exclusion Restriction Hold True*

As Table 1 shows, when both the ignorability and the exclusion restriction assumptions hold true, both the PS- and the MTE-based methods can

correctly identify group-level causal effects (as long as the data-generating model is correctly specified). In this case, they both provide estimates of ATE and TT that are asymptotically unbiased. Nonetheless, their statistical efficiency may differ. Although estimation uncertainty may converge to zero as sample size goes to infinity, we cannot avoid the limitation of sample size in empirical research. At present, we know little about the asymptotic variance of the estimators produced by the semiparametric MTE method. It is therefore of practical relevance for us to explore the statistical efficiency of the methods being evaluated in this article.

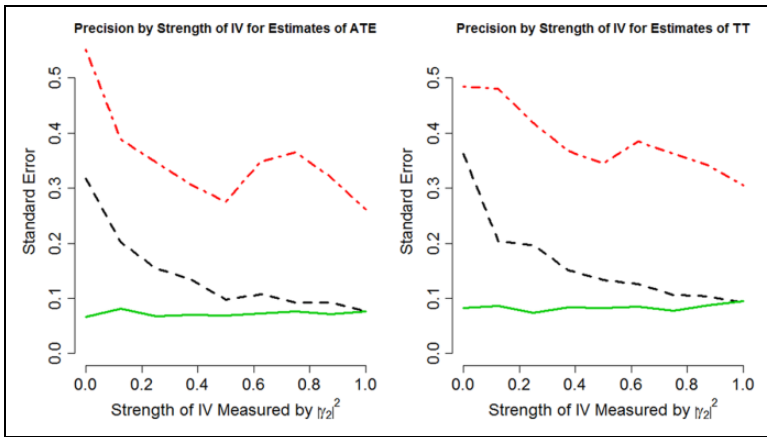
To achieve this end, we utilize simulated data. First, we generate data through the two potential outcome models and the treatment selection model described in the previous section:

$$\begin{aligned}
 Y^0 &= \beta_{00} + \beta_{01}X + \epsilon, \\
 Y^1 &= \beta_{10} + \beta_{11}X + \epsilon + \eta, \\
 D^* &= \gamma_0 + \gamma_1X + \gamma_2Z - V \\
 D &= 1(D^* > 0),
 \end{aligned}$$

with the following parameterization:

$$\begin{aligned}
 \beta_0 &= [\beta_{00}, \beta_{01}] = [0, 1], \beta_1 = [\beta_{10}, \beta_{11}] = [3, 2], \\
 \gamma &= [\gamma_0, \gamma_1, \gamma_2] = [0, \gamma_1, \gamma_2], \gamma_1^2 + \gamma_2^2 = 1 \\
 X, Z &\sim N(0, 1), X \perp\!\!\!\perp Z \\
 \epsilon, \eta, V &\sim N(0, 1), \epsilon, \eta, V \text{ mutually independent,} \\
 \epsilon, \eta, V &\perp\!\!\!\perp X, Z.
 \end{aligned}$$

Note that in this parameterization, the mutual independence of error terms  $\epsilon$ ,  $\eta$ , and  $V$  implies the validity of the ignorability assumption. The exclusion restriction is made true by the joint independence of error terms and  $Z$ , which serves here as an IV. However, the values of  $\gamma_1$  and  $\gamma_2$  are not fixed. We manipulate the value of  $\gamma_2$  to vary the relative importance of  $Z$  in determining the treatment status. Evidently, when  $\gamma_2$  is small, the strength of the IV,  $Z$ , is weak. Meanwhile, since  $X$  and  $Z$  are distributed as independent standard normal, we fix  $\|\gamma\|$  at 1 to keep the importance of the observables ( $X$  and  $Z$ ) relative to the unobservable  $V$  roughly constant regardless of the value of  $\gamma_1$  or  $\gamma_2$ .



**Figure 1.** Standard error by strength of instrumental variable (IV) for different estimators of average treatment effect (ATE; left) and treatment effect of the treated (TT; right). The standard error for each estimator is calculated from 100 random samples of size 2,500. Solid line: estimators using the smoothing-difference propensity score (PS)-based method; dashed line: estimators using the parametric marginal treatment effects (MTEs) method; dot dash line: estimators using the semiparametric MTE method.

In our simulation, we alter the value of  $\gamma_2^2$  from 0 to 1 with a step size of 0.125, thus generating nine scenarios with a gradual change in the strength of the IV. For each of these scenarios, we conduct a Monte Carlo experiment as follows: first, we generate a hypothetical population of size 100,000. Next, we draw 100 samples of size 2,500 from each of these populations. Then, for each sample, we estimate the causal parameters of ATE and TT with the three methods that we discussed in previous sections: (1) the smoothing-difference PS-based method (using smoothing splines), (2) the parametric MTE-based method, and (3) the semiparametric LIV method. For the first method, we construct the PS using only  $X$ .<sup>15</sup> Finally, for each estimator, we report its standard error as an indicator of statistical efficiency.

In Figure 1, we plot the trends of the standard error for estimates of ATE and TT as we vary the explanatory power of the IV in the treatment selection model. First of all, we see that the semiparametric MTE method (dot-dash line) generally yields estimates with much larger standard errors than those from the other two methods. Indeed, when  $\gamma_2$  is very small (weak IV), the standard error of the semiparametric LIV method (around 0.5) is more than

five times as large as that of the PS-based method (less than 0.1). More importantly, Figure 1 shows that the relative strength of IV matters greatly for the efficiency of the two MTE-based methods. In fact, both the parametric and the semiparametric MTE-based estimates undergo a substantial decline in standard error when the IV becomes a stronger predictor of treatment selection. In summary, two findings emerge from the results. On one hand, when the IV is relatively weak, the PS-based method outperforms both MTE-based methods. On the other hand, when treatment selection is dominated by the IV ( $\gamma_2 = 1$ ), the parametric MTE method and the PS-based method converge in their estimation uncertainty (around 0.1), whereas the semiparametric MTE approach still suffers from an inefficiency penalty with a significantly larger standard error (around 0.3) for either ATE or TT.

From this simulation, we observe that when both ignorability and the exclusion restriction hold true, the PS-based method is generally preferable to the MTE-based methods, especially when the IV is relatively weak. Although the parametric MTE-based approach does not require the exclusion restriction for identification, its estimation efficiency depends heavily on the availability of a strong IV. The semiparametric LIV estimation depends on an IV for identification and a strong IV for efficiency. As Figure 1 shows, when the IV strengthens as a predictor of treatment selection, estimates from the MTE-based methods become less uncertain.

### *When Both Ignorability and the Exclusion Restriction Break Down*

When the ignorability and the exclusion restriction assumptions are both violated, neither the PS-based method nor the MTE-based methods produces theoretically unbiased estimators of ATE and TT. Unfortunately, this is a likely situation in actual settings of empirical research. This motivates us to explore potential patterns of under-/overestimation due to the violation of both ignorability (for PS-based methods) and exclusion restriction (for MTE-based methods).

In the Decomposition of Ignorability subsection, we showed some implications of PS-based estimation when ignorability is violated. Rewriting equation (11), we can express the biases of PS-specific estimators as:

$$\text{BiasATE}(p) = E(Y^0|D = 1, p) - E(Y^0|D = 0, p) + \text{TT}(p) - \text{ATE}(p);$$

$$\text{BiasTT}(p) = E(Y^0|D = 1, p) - E(Y^0|D = 0, p).$$

Therefore, the bias of ATE is an aggregate of type I and type II selection biases due to unobservables, whereas the bias of TT is due only to

**Table 2.** Biases of ATE and TT due to Unobserved Selection for PS-based Estimators.

Unobserved selection			
Type I	Type II	BiasATE	BiasTT
+	+	+	+
+	0	+	+
+	-	Uncertain	+
0	+	+	0
0	0	0	0
0	-	-	0
-	+	Uncertain	-
-	0	-	-
-	-	-	-

Note: ATE = average treatment effect; TT = treatment effect of the treated; PS = propensity score.

unobserved type I selection. Neither of them depends on the validity of the exclusion restriction since IVs play no role in PS-based methods. According to the previous two expressions, we summarize the directions of BiasATE and BiasTT and under different scenarios of unobserved selection in Table 2. Hence, if we postulate with some confidence the underlying pattern of unobserved selection, we may surmise the direction of over-/underestimation of ATE and TT in PS-based estimation. For example, when there is a sorting on gain but no selection on level, PS-based methods are likely to overestimate ATE but not TT. If unobserved type I and type II selection biases are in the opposite direction, the sign of BiasATE, but not of BiasTT, would be indeterminate.

In contrast, when the exclusion restriction breaks down, it is difficult to know the direction of the bias for MTE-based estimators, because their analytical expressions are not readily available. To complicate things, the breakdown of the exclusion restriction may take different forms. The IV may be correlated, either positively or negatively, with either of the two error terms ( $\epsilon, \eta$ ). In any of these scenarios, the potential biases for ATE and TT may also be affected by the specific pattern of unobserved selection, that is, the dependence of  $V$  on  $\epsilon$  and  $\eta$ . In sum, there is no simple guideline that helps us decide whether ATE or TT will be overestimated or underestimated by MTE-based methods when the exclusion restriction is violated.

However, for any particular case, we can still explore its consequences through numerical simulation. As an illustration, we explore subsequently

a concrete case that merits detailed investigation in which (1) there is a negative (unobserved) type I selection and a positive (unobserved) type II selection and (2) the IV used for estimating MTE is correlated with the treatment effect  $Y^1 - Y^0$ . We use the same simulation setup as that in the previous subsection with the following parameterization:

$$\beta_0 = [\beta_{00}, \beta_{01}] = [0, 1], \beta_1 = [\beta_{10}, \beta_{11}] = [3, 2],$$

$$\gamma = [\gamma_0, \gamma_1, \gamma_2] = [0, 1, 0.2],$$

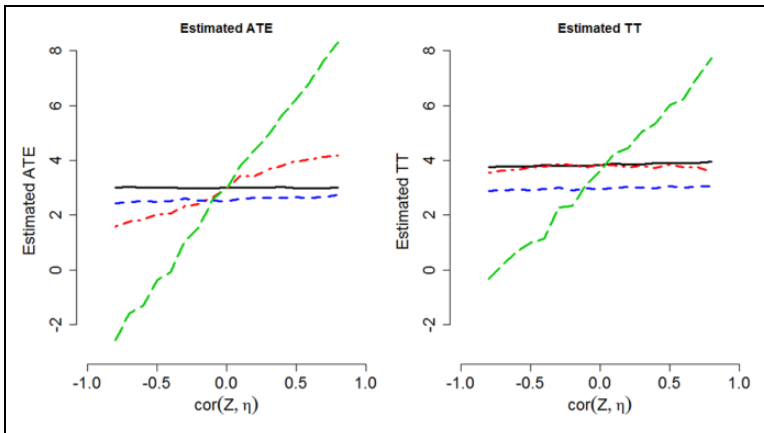
$$X, Z \sim N(0, 1), X \perp\!\!\!\perp Z,$$

$$\epsilon, \eta, V \sim N(0, 1), \epsilon \perp\!\!\!\perp \eta, \text{cor}(\epsilon, V) = 0.5, \text{cor}(\eta, V) = -0.5,$$

$$\epsilon, \eta, V \perp\!\!\!\perp X, \epsilon, V \perp\!\!\!\perp Z.$$

Note that in this specification, we assume a positive correlation between  $\epsilon$  and  $V$  (type I selection) but a negative correlation between  $\eta$  and  $V$  (type II selection). Since  $V$  represents the latent resistance to receiving treatment, this setup is one of “negative sorting on level” and “positive sorting on gain.” This pattern reflects the literature in estimating returns to schooling. Under the conventional common effects model, the argument for “ability bias” (Griliches 1977) predicted a positive type I selection due to unobserved ability, that is, more capable individuals tend to acquire more education as well as to earn more money. However, more recent research considering heterogeneous effects has found support for the “comparative advantage” argument, which implies negative sorting on level as well as positive sorting on gain (Cunha, Heckman, and Navarro 2005; Willis and Rosen 1979). In other words, it is predicted that individuals who actually went to college would be worse off than those who did not if they had not attended college, although the former group has benefited more from college education than the latter group would have had they attended college. Patterns of self-selection have been widely observed in other contexts (Winship and Mare 1992). For example, Smock, Manning, and Gupta (1999: 809) found that “divorced women would not fare as well economically as married women had they remained married instead of divorcing,” indicating some self-selection into divorce among women who have less to lose from divorce.

Further, we assume independence between  $Z$  and  $\epsilon$  but not between  $Z$  and  $\eta$ . The dependence between  $Z$  and  $\eta$  implies that  $Z$  is not truly exogenous but “moderates” the treatment effect. In the literature estimating earnings returns to college education, researchers have used the distance from a



**Figure 2.** Estimated average treatment effect (ATE; left) and treatment effect of the treated (TT; right) when exclusion restriction breaks down. This figure shows the estimates of ATE and TT using different methods as the correlation between  $Z$  and  $\eta$  changes from  $-0.8$  to  $0.8$ . Solid line: actual values of ATE and TT; dashed line: estimates using the smoothing-difference propensity score (PS)-based method; dot dash line: estimates using the parametric marginal treatment effects (MTEs) method; long dash line: estimates using the semiparametric MTE method.

youth's home to college as an IV (Cameron and Taber 2004; Card 1995; Currie and Moretti 2003; Kane and Rouse 1995). However, if returns to college vary by the distance measure, for example, students living closer to colleges would benefit more from college than those who live further away, the IV would be correlated with  $\eta$ . Here, we vary the correlation between  $Z$  and  $\eta$  from  $-0.8$  to  $0.8$  with a step size of  $0.1$ , generating 17 scenarios.<sup>16</sup> For each of these scenarios, we simulate a hypothetical sample of 20,000 and estimate the causal parameters of ATE and TT using the same three methods as specified in the previous subsection.<sup>17</sup> Finally, we display the results in Figure 2.

The left panel of Figure 2 shows the estimates of ATE, along with its actual values (solid line). First of all, we can see that the MTE-based estimates of ATE are upwardly biased when  $\text{cor}(Z, \eta) > 0$  and downwardly biased when  $\text{cor}(Z, \eta) < 0$ . In fact, the larger the correlation between  $Z$  and  $\eta$ , the higher the estimates from the MTE-based methods, especially the semiparametric LIV estimates (long dash line). For example, when  $\text{cor}(Z, \eta)$  is larger than  $0.5$ , the semiparametric LIV estimates are greater than  $6.0$ , twice as large as its actual value ( $3.0$ ), whereas the parametric MTE-based estimates (dot-dash line) are upwardly biased by a smaller



magnitude, at about 4.0. In comparison, the PS-based estimates (dashed line) show a moderate downward bias in this setup. As expected, the magnitude of bias for the PS-based estimates does not depend on  $\text{cor}(Z, \eta)$ , because the PS-based estimates do not rely on the exclusion restriction for estimation.

The right panel compares estimates of TT. Similar to the case of ATE, the semiparametric LIV approach yields estimates that are significantly upwardly biased when  $\text{cor}(Z, \eta) > 0$ , and downwardly biased when  $\text{cor}(Z, \eta) < 0$ . Nonetheless, the parametric MTE-based estimates are almost equal to the true value of TT across the entire range of  $\text{cor}(Z, \eta)$ . Finally, the PS-based estimates of TT show a significant underestimation. In fact, we may infer this last result from an earlier discussion, as Table 2 indicates that TT is underestimated as long as there is a negative sorting on level.

Overall, the previous simulation reveals that, when there is a negative sorting on level (type I selection) and positive sorting on gain (type II selection) due to unobservables, the MTE-based methods, especially the semiparametric LIV method, may severely overestimate or underestimate ATE and TT due to the use of an improper IV. As expected, the same causal parameters may be underestimated by the PS-based method. As we will see in the next section, these results can reasonably explain apparent discrepancies in an assessment of returns to college.

## Empirical Example

To illustrate the three methods we discussed earlier, we applied them to the data used in the Carneiro et al. (2011) study of returns to college education using MTE. In the subsections that follow: we (1) describe the data, (2) demonstrate the use of the smoothing-difference PS-based method, (3) replicate the Carneiro et al. (2011) results using MTE, and (4) compare MTE- and PS-based estimates of ATE and TT.

### Data Description

Following Carneiro et al. (2011), we reanalyze a sample of white males ( $N = 1,747$ ) who were 16–22 years old in 1979, drawn from the NLSY 1979. Treatment is college attendance measured by having attained any postsecondary education by 1991. By this definition, the treated group consists of 865 subjects and the control group consists of 882 subjects. The wage variable is measured as an average of deflated (to 1983 constant dollars) nonmissing hourly wages reported between 1989 and 1993. Pretreatment covariates ( $\mathbf{X}$ ) are urban residence at 14, the Armed Forces Qualification Test

**Table 3.** Propensity Score Probit Model Predicting College Attendance.

Predictors	Coefficient
Urban residence at 14	0.127 (0.084)
Corrected AFQT	0.667*** (0.045)
Corrected AFQT square	0.196*** (0.039)
Mother's years of schooling	-0.110 (0.089)
Mother's years of schooling square	0.010** (0.004)
Number of siblings	-0.090† (0.053)
Number of siblings square	0.002 (0.006)
Permanent local log earnings at 17	-43.9* (17.2)
Permanent local log earnings at 17 square	2.15* (0.84)
Permanent state unemployment rate at 17	0.240 (0.369)
Permanent state unemployment rate at 17 square	-0.018 (0.029)
Model $\chi^2$	684.4 ( $df = 18$ )

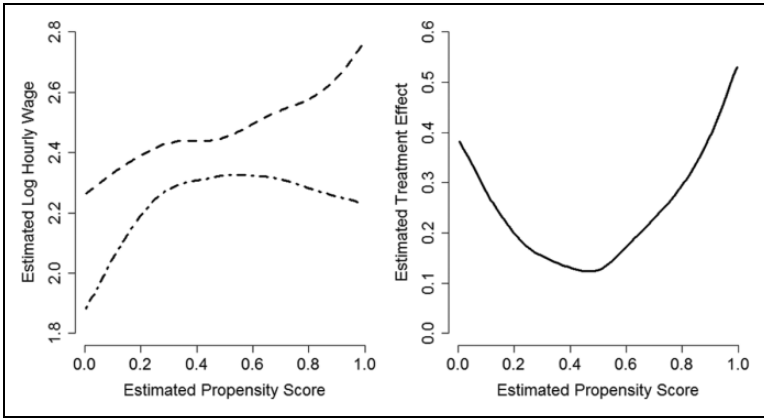
Note: Numbers in parentheses are standard errors. † $p < .1$ . \* $p < .05$ . \*\* $p < .01$ . \*\*\* $p < .001$ .

(AFQT) score adjusted by years of schooling, mother's years of schooling, number of siblings, permanent local log earnings at 17 (county log earnings averaged between 1973 and 2000), permanent local unemployment rate at age 17 (state unemployment rate averaged between 1973 and 2000), and cohort dummies. IVs ( $Z/X$ ) include (a) the presence of a four-year college in the county of residence at age 14, (b) local wage in the county of residence at age 17, (c) local unemployment rate in the state of residence at age 17, and (d) average tuition in public four-year colleges in the county of residence at age 17. More detailed description of the data set is provided in Carneiro et al. (2011).

### *The Smoothing-difference PS-based Results*

Subsequently, we show results from the smoothing-difference PS-based method. First of all, we estimate the PS of attending college for each subject in the sample given  $X$  using a probit regression model. Table 3 presents the fitted PS model. We can see that the likelihood of attending college is predicted positively by corrected AFQT score and negatively by number of siblings and permanent local log earnings at age 17.

In the next step, we fit two separate nonparametric models regressing the log hourly wage on the estimated PS, one for the treated group that went to college and one for the untreated group that did not go to college. Here, we use smoothing splines with five equivalent degrees of freedom.<sup>18</sup> Figure 3



**Figure 3.** The smoothing-difference propensity score (PS)-based method for estimating returns to college. The left panel shows the expected annual wages, respectively, for those who attended college (dashed line) and for those who did not attend college (dot dash line). The right panel demonstrates the expected return to college for people with different propensity scores.

displays the resulting curves, evaluated over the entire interval (with a small portion being extrapolated). In the left panel, the dashed line and the dot-dash line show the expected log hourly wage, respectively, for those who went to college and for those who did not. Two patterns emerge from this figure. First, for persons who attended college, the expected wage increases steadily with the PS. That is, labor market outcomes differ systematically among college goers, as those with a higher propensity to attend college earn more than those with a lower propensity. Second, for persons who did not attend college, expected wage shows a rapid increase at the lower end of PS but flattens out thereafter. Hence, individuals who are very unlikely to go to college on the basis of their observed covariates included in the PS are truly disadvantaged. If they do not go to college, then they earn much lower wages than their peers with a higher propensity to attend college ( $e^{1.9} = 6.7$  at  $p \approx 0$ , compared to  $e^{2.2} = 9.0$  at  $p \approx 0.2$ ). However, they also stand to gain a lot from attending college ( $e^{2.2} = 9.0$  at  $p \approx 0$ ), although their wages would be still substantially lower than those of other college goers with a higher propensity of attending college (e.g.,  $e^{2.8} = 16.4$  at  $p \approx 1.0$ ).

We now turn to the right panel, which depicts estimated heterogeneous treatment effects by PS.<sup>19</sup> This curve is obtained directly by differencing the two functions in the left panel. The nonmonotonic pattern suggests that two groups of individuals exist who benefit most from college: those most unlikely to go to

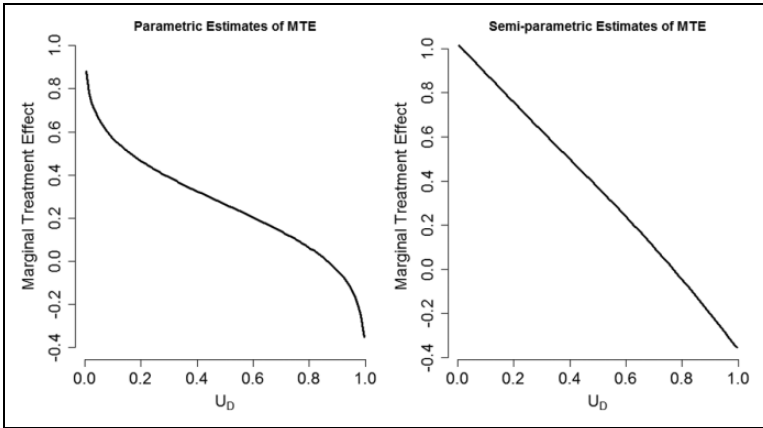
college and those most likely to go to college.<sup>20</sup> Thus, college education seems to be more valuable for persons at either the low end or the high end of the PS than for those in the middle. Therefore, from these data, we observe a mix of positive selection and negative selection into college using the PS-based approach.

Next, we use the earlier curve to predict treatment effect  $\delta_i$  for each individual  $i$  in the sample. We then average these  $\delta_i$ 's over the entire sample to obtain ATE, and over those who actually attended college to estimate TT. We will discuss these summary results in the next subsection, comparing them to those produced by the MTE-based methods.

### *MTE-based Results*

We now give up the ignorability assumption and thus the PS approach. Instead, we use the MTE-based methods, with covariates  $X$  and IVs  $Z/X$  specified in the Data Description subsection. We first estimate two sets of MTEs, one from the parametric model, and the other from the semiparametric LIV method. Figure 4 plots these two sets of  $MTE(x, u_D)$ , both evaluated at mean values of  $X$ . Both the parametric and the semiparametric estimates of MTE show a declining trend with respect to  $u_D$ , that is, the unobserved resistance to attending college. These results show that individuals with higher returns to college are more likely to go to college (in having lower  $u_D$ ). Furthermore, the magnitude of the heterogeneity in MTE is substantial: returns can vary from as high as 80 percent – 100 percent (for low  $u_D$  persons who would double their wages from attending college) to as low as –40 percent (for high  $u_D$  persons who would lose from attending college).

Using weights provided by Heckman et al. (2006a), we construct standard treatment parameters from the two sets of estimated MTE. Columns 1–4 of Table 4 show the final estimates of ATE and TT from different methods, with bootstrapped standard errors. We observe that MTE-based estimates of ATE and TT are less precise than those from the PS-based method. The lack of precision for MTE-based estimates is expected since the IVs we use are relatively weak compared to  $X$  in determining treatment selection (see When Both Ignorability and the Exclusion Restriction Hold True subsection). More importantly, MTE- and PS-based results differ in magnitude. For ATE, the differences are not statistically significant, although the semiparametric LIV method seems to give a larger point estimate than do the other two methods. For TT, the difference between MTE- and PS-based results is more substantial. Both the parametric and the semiparametric MTE-based methods yield significantly higher estimates of TT than the PS-based estimate. Specifically,



**Figure 4.** Estimated marginal treatment effects (averaged over  $\mathbf{X}$ ) from marginal treatment effects (MTEs)-based methods.

**Table 4.** Estimates for Returns to College from NLSY Data.

Causal parameters	Smoothing-difference PS-based method	MTE-based methods		MTE-based methods (without <i>tuition</i> )	
		Parametric	Semiparametric	Parametric	Semiparametric
ATE	0.242 (0.067)	0.264 (0.159)	0.356 (0.174)	0.231 (0.147)	0.277 (0.202)
TT	0.278 (0.093)	0.567 (0.156)	0.736 (0.226)	0.540 (0.144)	0.604 (0.243)

Note: Numbers in parentheses are bootstrapped standard errors with 250 repetitions. MTEs = marginal treatment effects; NLSY = National Longitudinal Survey of Youth.

we obtain the TT estimate of college returns at 73.6 percent by the semiparametric MTE method but only 27.8 percent by the PS-based method.

### Discussion on the Discrepancy in Estimates of TT

From Table 4 a natural question arises, why is there such a large discrepancy between PS-based estimate and MTE-based estimate of TT? In light of our discussion in When Both Ignorability and the Exclusion Restriction Break Down subsection, we can offer some speculations. On one hand, there could be an underestimation by the PS-based method due to the breakdown of the

ignorability assumption. One way to examine the potential pattern of unobserved selection is from the MTE-based analysis. In fact, the parametric MTE approach provides the following estimates:

$$\hat{\sigma}_{\epsilon V} = 0.08, \hat{\sigma}_{\eta V} = -0.24,$$

where  $\sigma_{\epsilon V}$  and  $\sigma_{\eta V}$  denote the covariances, respectively, between  $\epsilon$  and  $V$  and between  $\eta$  and  $V$ . Since  $V$  represents a latent resistance to receiving treatment, these estimates suggest a negative sorting on level and a positive sorting on gain. This finding accords well with Willis and Rosen's (1979) model of comparative advantage, which argues that college goers would do worse if they did not go to college but benefit more from college education than persons who do not go to college. If we accept these estimates as evidence for a negative type I selection and a positive type II selection (due to unobservables), our earlier discussion around Table 2 would suggest indeed a downward bias for the PS-based estimate of TT. Hence, the discrepancy in TT estimates between PS- and MTE-based methods could be attributed to a negative sorting on precollege earnings.

On the other hand, there might be an overestimation by MTE-based methods due to the violation of the exclusion restriction. The numerical simulation results in the When Both Ignorability and the Exclusion Restriction Break Down subsection suggest a potentially upward bias when there is a positive correlation between IV and the treatment effect, that is,  $\text{cor}(Z, \eta) > 0$  (for  $\gamma > 0$ ). Unfortunately, such a correlation is empirically unverifiable, since  $\eta$  is an unobserved attribute that cannot be individually recovered from the data. Nonetheless, our results provide good grounds for questioning the theoretical validity of IV in concrete settings. In our example, one of the IVs is average tuition in public four-year colleges in the county of residence. This variable, however, is likely to be correlated with college quality and thus could influence the returns to college. To test this conjecture, we excluded average tuition as an IV and reanalyzed the data.<sup>21</sup> Columns 5 and 6 of Table 4 show the results for ATE and TT after this modification. Compared with column 3, the new results by the parametric MTE-based method are largely unchanged. However, for the semiparametric LIV approach, the large estimates reported earlier, especially of TT, are markedly reduced.

In sum, the large discrepancy in TT between PS- and MTE-based estimates may be caused by one of the three underlying causal mechanisms: (1) the negative sorting on precollege earnings, (2) the use of an improper IV, or (3) a mixture of the previous two. Because of this uncertainty, neither

the PS-based method nor the (semiparametric) MTE-based method yields the true ATE of college education for college goers. Most likely, the truth lies somewhere in between.

## Concluding Remarks

In this study, we have examined certain statistical properties of PS- and MTE-based methods through an exposition of identification issues, two simulation analyses, and an empirical application. We showed that the applicability of PS-based methods is not limited to settings in which complete ignorability is satisfied. In fact, it is useful to decompose ignorability into two components: (1) ignorability of type I selection bias or baseline difference between treated and untreated units and (2) ignorability of type II selection bias or difference in treatment effects between treated and untreated units. We have shown that as long as the ignorability of type I selection bias is satisfied, PS-based methods can still identify TT, even in the presence of a heterogeneous treatment effect bias. Furthermore, when type I selection bias cannot be ignored, the bias for TT is in the same direction as the type I selection bias. For example, in the evaluation of returns to college, a negative type I selection bias is part of the model of “comparative advantage.” An underestimation of TT by PS-based methods would occur under this situation.

By comparison, MTE-based methods are robust to different types of violation of the ignorability assumption. However, they require strong IVs to achieve statistical efficiency. This is true for both the parametric model and the semiparametric method. Furthermore, when the exclusion restriction is violated, MTE-based methods, especially the semiparametric LIV approach, can be subject to severe overestimation or underestimation of treatment effects. In practice, the plausibility of the exclusion restriction assumption cannot be verified but can be evaluated based on substantive knowledge about the research setting. If substantive concerns suggest the violation of the exclusion restriction, we could, as we did in Discussion on the Discrepancy in Estimates of TT subsection, exclude the susceptible IV in the treatment selection model and reanalyze the data. In addition, we may directly assess the consequence of a violation of the assumption through sensitivity analyses (e.g., see Angrist 1990; Angrist, Imbens, and Rubin 1996).

This article has also proposed a PS-based method based on first smoothing two counterfactual outcomes, which we call the smoothing-difference method. Compared to traditional matching and stratification methods, the smoothing-difference method has two distinct advantages. On one hand, it enables the researcher to examine the nonparametric trends of counterfactual

outcomes by treatment status across the spectrum of PS. In our empirical example, we have shown the variations in wages by both the PS of attending college and the status of college attendance. On the other hand, this method produces a nonparametric pattern of treatment effect heterogeneity across individuals with different PSs. Such an observed pattern of heterogeneity is of interest to social science researchers, although its interpretation is still ambiguous, depending on the validity of the ignorability assumption (Brand and Xie 2010; Xie et al. 2012). For example, if the ignorability assumption holds true, observed results reveal the pattern of heterogeneous treatment effects. If one accepts only the ignorability of type I selection bias, heterogeneous treatment effects along the PS should be interpreted only for those who are actually treated. If one does not embrace any form of ignorability, the observed pattern may reveal an underlying selection process sorting out treated units from untreated units (Xie and Wu 2005).

### **Acknowledgments**

We thank three anonymous referees for their helpful comments on the previous version of this article.

### **Authors' Note**

The authors benefited from communications with Jennie Brand, James Heckman, and Ben Jann. The ideas expressed herein are those of the authors.

### **Declaration of Conflicting Interests**

The author(s) declared no potential conflicts of interest with respect to the research, authorship, and/or publication of this article.

### **Funding**

The author(s) disclosed receipt of the following financial support for the research, authorship, and/or publication of this article: Financial support for this research was provided by the National Institutes of Health, Grant 1 R21 NR010856-01 and by the Population Studies Center at the University of Michigan, which receives core support from the National Institute of Child Health and Human Development, Grant R24HD041028.

### **Notes**

1. The exact meaning of the strength of an IV will be defined in the fourth section.
2. An implicit condition for defining causal effects within the counterfactual framework is the stable-unit-treatment-value assumption (SUTVA), which requires



- that the value of  $\delta_i$  does not depend on what mechanism is used to assign the treatment to subject  $i$ , or what treatments the other subjects receive (Rubin 1986).
3. There is a guarantee of comparability of the treated group and the control group in an experiment. In this article, we restrict our attention to observational studies.
  4. These two equations constitute a necessary but not sufficient condition for the ignorability assumption of equation (5). In the literature, they are usually called the “weak ignorability assumption” or “conditional mean independence.” See Woodridge (2001).
  5. The smoothing-difference method, as an improvement upon the approach adopted by Brand and Xie (2010), is a by-product of this research. It has also been incorporated into Xie et al. (2012).
  6. Alternative nonparametric regression techniques, such as kernel methods and local polynomial regression, could also be applied here.
  7. Xie et al. (2012) also provide a comparison between this method and other PS-based methods.
  8. The first condition is also called “unconfoundedness for controls” (Imbens 2004).
  9. In this case, type I selection bias due to unobservables could be reduced by other methods such as the conventional IV approach, the fixed-effect model, and difference-in-difference methods (after conditioning on the PS).
  10. The linearity assumption is convenient but not necessary. We can generally assume  $Y^0 = \mu_0(\mathbf{X}) + \epsilon$  and  $Y^1 = \mu_1(\mathbf{X}) + \epsilon + \eta$  for any given functions  $\mu_0(\mathbf{X})$  and  $\mu_1(\mathbf{X})$ .
  11. Weights for different parameters of interest are given in Heckman et al. (2006a).
  12. For specific issues on the implementation of the semiparametric MTE method, see Heckman et al. (2006b).
  13. An illustrative simulation study for the MTE-based approach is given in Heckman et al. (2006b). However, the authors considered only the situation where variables in  $\mathbf{Z}$  and variables in  $\mathbf{X}$  are mutually exclusive and independent, that is, the treatment selection equation is purely determined by IVs, which is rather unrealistic.
  14. For the parametric model, a detailed discussion on the asymptotic variance of the maximum likelihood (ML) estimator could be found in Heckman (1979) and Puhani (2000).
  15. For a discussion of whether to include IV in estimating the PS, see Pearl (2009).
  16. For the covariance matrix of  $(\epsilon, \eta, V, Z)$  to be positive definite, the correlation between  $\eta$  and  $Z$  cannot exceed 0.8.
  17. Here, we include  $Z$  in estimating the PS, because it is correlated with unobservables in the outcome equations.
  18. Alternative choices of the smoothing parameter do not substantially alter our results.

19. For a discussion of why such heterogeneity is of special interest, see Xie et al. (2012).
20. The second part of this finding, that is, a larger return at the very high level of the PS, is inconsistent with Brand and Xie's (2010) main conclusion. Future research is needed to explain this inconsistency.
21. We also reanalyzed the data after excluding the presence of college, another controversial IV, from the set of instruments. The corresponding estimates of ATE and TT, however, do not change much.

## References

- Angrist, Joshua D. 1990. "Lifetime Earnings and the Vietnam Era Draft Lottery: Evidence from Social Security Administrative Records." *American Economic Review* 80:313-35.
- Angrist, Joshua D., Guido W. Imbens, and Donald B. Rubin. 1996. "Identification of Causal Effects Using Instrumental Variables." *Journal of the American Statistical Association* 91:444-55.
- Angrist, Joshua D. and Alan B. Krueger. 1999. "Empirical Strategies in Labor Economics." Pp. 1277-366 in *Handbook of Labor Economics*, vol. 3A, edited by O. Ashenfelter and D. Card. Amsterdam, the Netherlands: Elsevier.
- Ansari, Asim and Jedidi Kamel. 2000. "Bayesian Factor Analysis for Multilevel Binary Observations." *Psychometrika* 65:475-96.
- Bauer, Daniel J. and Patrick J. Curran. 2003. "Distributional Assumptions of Growth Mixture Models: Implications for Overextraction of Latent Trajectory Classes." *Psychological Methods* 8:338-63.
- Björklund, Anders and Robert Moffitt. 1987. "The Estimation of Wage Gains and Welfare Gains in Self-selection Models." *The Review of Economics and Statistics* 69:42-49.
- Blundell, Richard, Lorraine Dearden, and Barbara Sianesi. 2005. "Evaluating the Effect of Education on Earnings: Models, Methods, and Results from the National Child Development Survey." *Journal of the Royal Statistical Society: Series A* 168:473-512.
- Brand, Jennie E. and Yu Xie. 2010. "Who Benefits Most from College? Evidence for Negative Selection in Heterogeneous Economic Returns to Higher Education." *American Sociological Review* 75:273-302.
- Cameron, Stephen V. and Christopher Taber. 2004. "Estimation of Educational Borrowing Constraints Using Returns to Schooling." *Journal of Political Economy* 112:132-82.
- Card, David. 1995. "Using Geographic Variation in College Proximity to Estimate the Return to Schooling." Pp. 201-22 in *Aspects of Labour Market Behavior: Essays in Honor of John Vanderkamp*, edited by Louis N. Christofides,

- E. Kenneth Grant, and Robert Swidinsky. Toronto, Canada: University of Toronto Press.
- Carneiro, Pedro, James Heckman, and Edward Vytlacil. 2011. "Estimating Marginal Returns to Education." *American Economic Review* 101:2754-81.
- Cornfield, Jerome, William Haenszel, E. Cuyler Hammond, Abraham M. Lilienfeld, Michael B. Shimkin, and Ernst L. Wynder. 1959. "Smoking and Lung Cancer: Recent Evidence and a Discussion of Some Questions." *Journal of the National Cancer Institute* 22:173-203.
- Cunha, Flavio, James J. Heckman, and Salvador Navarro. 2005. "Separating Uncertainty from Heterogeneity in Life Cycle Earnings, the 2004 Hicks Lecture." *Oxford Economic Papers* 57:191-261.
- Currie, Janet and Enrico Moretti. 2003. "Mother's Education and the Intergenerational Transmission of Human Capital: Evidence from College Openings." *Quarterly Journal of Economics* 118:1495-532.
- DiPrete, Thomas and Markus Gangl. 2004. "Assessing Bias in the Estimation of Causal Effects: Rosenbaum Bounds on Matching Estimators and Instrumental Variables Estimation with Imperfect Instruments." *Sociological Methodology* 34:271-310.
- Greenland, Sander and Charles Poole. 1988. "Invariants and Noninvariants in the Concept of Interdependent Effects." *Scandinavian Journal of Work, Environment & Health* 14:125-29.
- Griliches, Zvi. 1977. "Estimating the Returns to Schooling: Some Econometric Problems." *Econometrica* 45:1-22.
- Harding, David J. 2003. "Counterfactual Models of Neighborhood Effects: The Effect of Neighborhood Poverty on High School Dropout and Teenage Pregnancy." *American Journal of Sociology* 109:676-719.
- Hastie, Trevor, Robert Tibshirani, and Jerome Friedman. 2008. *The Elements of Statistical Learning: Data Mining, Inference and Prediction*. New York: Springer.
- Heckman, James J. 1979. "Sample Selection Bias as a Specification Error." *Econometrica* 47:153-61.
- Heckman, James J. 2001. "Micro Data, Heterogeneity, and the Evaluation of Public Policy: Nobel Lecture." *Journal of Political Economy* 109:673-748.
- Heckman, James J. 2005. "The Scientific Model of Causality." *Sociological Methodology* 35:1-98.
- Heckman, James J. and Salvador Navarro-Lozano. 2004. "Using Matching, Instrumental Variables, and Control Functions to Estimate Economic Choice Models." *The Review of Economics and Statistics* 86:30-57.
- Heckman, James J. and Richard Robb. 1985. "Alternative Methods for Evaluating the Impact of Interventions." Pp.156-245 in *Longitudinal Analysis of Labor Market Data*, edited by James Heckman and Burton Singer. Cambridge, UK: Cambridge University Press.

- Heckman, James J., Sergio Urzua, and Edward Vytlacil. 2006a. "Understanding Instrumental Variables in Models with Essential Heterogeneity." *The Review of Economics and Statistics* 88:389-432.
- Heckman, James J., Sergio Urzua, and Edward Vytlacil. 2006b. "Estimation of Treatment Effects under Essential Heterogeneity." Retrieved October 12, 2010 ([http://jenni.uchicago.edu/underiv/documentation\\_2006\\_03\\_20.pdf](http://jenni.uchicago.edu/underiv/documentation_2006_03_20.pdf)).
- Heckman, James J. and Edward J. Vytlacil. 1999. "Local Instrumental Variables and Latent Variable Models for Identifying and Bounding Treatment Effects." *Proceedings of the National Academy of Sciences of the United States of America* 96:4730-34.
- Heckman, James J. and Edward J. Vytlacil. 2001. "Local Instrumental Variables." Pp. 1-46 in *Nonlinear Statistical Modeling: Proceedings of the Thirteenth International Symposium in Economic Theory and Econometrics: Essays in Honor of Takeshi Amemiya*, edited by Cheng Hsiao, Kimio Morimune, and James L. Powel. New York: Cambridge University Press.
- Heckman, James J. and Edward J. Vytlacil. 2005. "Structural Equations, Treatment Effects, and Econometric Policy Evaluation." *Econometrica* 73:669-738.
- Holland, Paul W. 1986. "Statistics and Causal Inference." *Journal of American Statistical Association* 81:945-60.
- Imbens, Guido W. 2004. "Nonparametric Estimation of Average Treatment Effects under Exogeneity: A Review." *Review of Economics and Statistics* 86:4-30.
- Kane, Thomas J. and Cecilia E. Rouse. 1995. "Labor-market Returns to Two- and Four-year College." *American Economic Review* 85:600-14.
- Lubke, Gitta H. and Bengt Muthén. 2005. "Investigating Population Heterogeneity with Factor Mixture Models." *Psychological Methods* 10:21-39.
- Manski, Charles. 1995. *Identification Problems in the Social Sciences*. Boston, MA: Harvard University Press.
- Manski, Charles. 2007. *Identification for Prediction and Decision*. Cambridge, MA: Harvard University Press.
- Moffitt, Robert. 1996. "Identification of Causal Effects Using Instrumental Variables: Comment." *Journal of the American Statistical Association* 91:462-65.
- Moffitt, Robert. 2008. "Estimating Marginal Treatment Effects in Heterogeneous Populations." *Annals of Economics and Statistics* 91/92:239-61.
- Morgan, Stephen and Christopher Winship. 2007. *Counterfactuals and Causal Inference: Methods and Principles for Social Research*. Cambridge, UK: Cambridge University Press.
- Pearl, Judea. 2009. *Causality: Models, Reasoning, and Inference*. 2nd ed. New York: Cambridge University Press.
- Puhani, Patrick. 2000. "The Heckman Correction for Sample Selection and Its Critique." *Journal of Economic Surveys* 14:53-68.

- Rosenbaum, Paul R. 2002. *Observational Studies*. New York: Springer.
- Rosenbaum, Paul R. and Donald B. Rubin. 1983. "The Central Role of the Propensity Score in Observational Studies for Causal Effects." *Biometrika* 70:41-55.
- Rosenbaum, Paul R. and Donald B. Rubin. 1984. "Reducing Bias in Observational Studies Using Subclassification on the Propensity Score." *Journal of the American Statistical Association* 79:516-24.
- Rothman, Kenneth J. and Sander Greenland, eds. 1998. *Modern Epidemiology*. 2nd ed. Philadelphia, PA: Lippincott-Raven Publishers.
- Rubin, Donald B. 1974. "Estimating Causal Effects of Treatments in Randomized and Nonrandomized Studies." *Journal of Educational Psychology* 66:688-701.
- Rubin, Donald B. 1986. "What Ifs Have Causal Answers?" *Journal of American Statistical Association* 81:961-62.
- Rubin, Donald B. 1997. "Estimating Causal Effects from Large Data Sets Using Propensity Scores." *Annals of Internal Medicine* 5:757-63.
- Shadish, William R., M. H. Clark, and Peter M. Steiner. 2008. "Can Nonrandomized Experiments Yield Accurate Answers? A Randomized Experiment Comparing Random and Nonrandom Assignments." *Journal of American Statistical Association* 103:1334-44.
- Smock, Pamela J., Wendy D. Manning, and Sanjiv Gupta. 1999. "The Effect of Marriage and Divorce on Women's Economic Well-being." *American Sociological Review* 64:794-812.
- Sobel, Michael E. 2000. "Causal Inference in the Social Science." *Journal of the American Statistical Association* 95:647-51.
- Tsai, Shu-Ling and Yu Xie. 2011. "Heterogeneity in Returns to College Education: Selection Bias in Contemporary Taiwan." *Social Science Research* 40:796-810.
- Willis, Robert J. and Sherwin Rosen. 1979. "Education and Self-selection." *Journal of Political Economy* 87:S7-36.
- Winship, Christopher and Robert. D. Mare. 1992. "Models for Sample Selection Bias." *Annual Review of Sociology* 18:327-50.
- Winship, Christopher and Stephen L. Morgan. 1999. "The Estimation of Causal Effects from Observational Data." *Annual Review of Sociology* 25:659-707.
- Winship, Christopher and Michael Sobel. 2004. "Causal Inference in Sociological Studies." Pp. 481-503 in *Handbook of Data Analysis*, edited by Melissa Hardy and Alan Bryman. London, UK: Sage.
- Woodridge, Jeffery M. 2001. *Econometric Analysis of Cross Section and Panel Data*. 1st ed. Cambridge: The MIT Press.
- Xie, Yu. 2000. "Assessment of the Long-term Benefits of Head Start." Pp.139-67 in *Into Adulthood: A Study of the Effects of Head Start*, edited by Sherri Oden, Lawrence J. Schweinhart, and David P. Weikart. Ypsilanti, MI: High/Scope Press.

- Xie, Yu. 2007. "Otis Dudley Duncan's Legacy: The Demographic Approach to Quantitative Reasoning in Social Science." *Research in Social Stratification and Mobility* 25:141-56.
- Xie, Yu, Jennie Brand, and Ben Jann. 2012. "Estimating Heterogeneous Treatment Effects with Observational Data." *Sociological Methodology* 42:314-47.
- Xie, Yu and Xiaogang Wu. 2005. "Market Premium, Social Process, and Statisticism." *American Sociological Review* 70:865-70.

### Author Biographies

**Xiang Zhou** is a PhD candidate in Sociology and Statistics at the University of Michigan. His research interests include economic inequality, social stratification, quantitative methods, and Chinese studies. His recent publications appear in *Social Forces*, *Sociological Methodology*, and *Proceedings of the National Academy of Sciences*.

**Yu Xie** is Otis Dudley Duncan distinguished university professor of Sociology, Statistics, and Public Policy at the University of Michigan. His main areas of interest are social stratification, demography, statistical methods, Chinese studies, and sociology of science. His recently published works include: *Marriage and Cohabitation* (University of Chicago Press 2007) with Arland Thornton and William Axinn, *Statistical Methods for Categorical Data Analysis* with Daniel Powers (Emerald 2008, second edition), and *Is American Science in Decline?* (Harvard University Press, 2012) with Alexandra Killewald.