

# Identification and Estimation of Causal Mechanisms and Net Effects of a Treatment\*

Carlos A. Flores<sup>†</sup>

Alfonso Flores-Lagunes<sup>‡</sup>

September 26, 2007

## Abstract

An important goal in the analysis of the causal effect of a treatment on an outcome is to understand the mechanisms through which the treatment causally works. In the economics literature, however, there seems to be no available framework to estimate the relative importance of different causal mechanisms of a treatment. We fill this void by precisely defining a causal mechanism effect of a treatment and the causal effect of a treatment net of that mechanism using the potential outcomes framework, and by considering their identification and estimation. The definition of our parameters results in an intuitive decomposition of the total effect of a treatment that is useful for policy purposes. We offer conditions under which these causal effects can be estimated for the case of a randomly assigned treatment and when selection into the treatment is random conditional on a set of covariates. We close with two empirical applications that illustrate the concepts and methods introduced in this paper.

Key words and phrases: causal inference, post-treatment variables, principal stratification.

JEL classification: C13, C21, C14

---

\*We thank Jon Guryan, Kei Hirano, Hilary Hoynes, Sabine Kroger, Oscar Mitnik, participants at the 2006 annual meeting of the Society of Labor Economists, the University of Miami Labor Lunch, 2006 Economic Science Association meetings, 2006 Midwest Econometrics Group meeting, 2006 Latin American meetings of the Econometric Society, and seminar participants at the Industrial Relations Section (Princeton), Laval, Purdue, and Virginia Tech for useful comments and discussions. All errors are our own.

<sup>†</sup>Department of Economics, University of Miami. Email: caflores@miami.edu

<sup>‡</sup>Food and Resource Economics Department, University of Florida, and Department of Economics, University of Arizona. Email: alfonsofl@ufl.edu

# 1 Introduction

The main purpose in the estimation of causal effects of a treatment or intervention is to estimate its total impact on a particular outcome. Commonly estimated parameters are the average treatment effect, the average treatment effect on the treated, and the marginal treatment effect.<sup>1</sup> In addition, it is of interest to estimate causal mechanisms through which the treatment or intervention works, and/or causal effects of the treatment on the outcome *net* of these mechanisms. Knowledge of these causal parameters allows a better understanding of the treatment and, as a result, can be used for policy purposes in the design, development, and evaluation of interventions. This paper analyzes identification and estimation of the average causal mechanism through which a treatment or intervention affects an outcome and the average causal effect of the treatment net of this mechanism. Using the potential outcomes framework, we precisely define our estimands of interest, consider different assumptions that can be employed in their identification and estimation, and analyze other related parameters mentioned elsewhere in the literature.

To briefly motivate the importance of understanding the mechanism through which a treatment works, consider the following example that is later employed as empirical illustration of the methods developed in the paper. When analyzing the causal effect of smoking during pregnancy on birth weight, it is of particular interest to determine what part of this causal effect works through a shorter gestation. If it were determined that the causal effect of smoking during pregnancy on birth weight works mainly through a shorter gestation time (as opposed to working through a low intrauterine growth), then medical procedures that help delay birth may be deemed helpful.

Many studies in economics are concerned with estimating mechanism effects and effects net of one or more mechanisms. For example, the literature on the effect of school quality on labor market outcomes recognizes that part of this effect may work through increasing years of schooling. To address this, Dearden et al. (2002) present results of the effect of school quality on wages with and without controlling for schooling to measure the total impact of school quality on wages and the effect that works through higher educational attainment. Similarly, Black and Smith (2004) use propensity score matching methods with and without including years of education in the propensity score specification. Another example is Simonsen and Skipper (2006), who estimate the effect of motherhood on wages in Denmark with emphasis on the total effect of having children on wages, taking into account various mechanisms through which motherhood may affect wages.<sup>2</sup> They use a propensity score matching approach and discuss

---

<sup>1</sup>See, for instance, Heckman, Lalonde and Smith (1999) for a detailed discussion of these parameters.

<sup>2</sup>For instance, they consider as a possible channel the sector of employment because, as they point out, Denmark's public sector is known to have higher benefits regarding maternity leave and more flexible working conditions than the private sector. Other channels they consider are working experience and occupation.

assumptions needed to estimate the total effect of motherhood on wages and the effect of motherhood on wages net of the mechanisms. As a final example, Ehrenberg et al. (2006) look at the channels or mechanisms through which the Andrew W. Mellon Foundation’s Graduate Education Initiative (GEI) affected the attrition and graduation probabilities of PhD students in various academic departments during the 1990s.<sup>3</sup> In general, every time causal effects of a treatment are estimated, it is natural to ask about the relative importance of different potential causal mechanisms through which the treatment works.<sup>4</sup>

A common problem in the existing literature attempting to estimate causal mechanisms of a treatment is that the parameters are not clearly defined or are defined within the context of the estimation procedure used (e.g., OLS, matching) and, most importantly, the assumptions needed for a causal interpretation of the estimates are not always made explicit.<sup>5</sup> To avoid this problem, we use the potential outcomes framework (Neyman, 1923; Rubin, 1974) to clearly define our parameters of interest and decompose the average (total) treatment effect into the average causal mechanism and the average causal effect net of that mechanism. In addition, to give a causal interpretation to our parameters we use the concept of principal stratification introduced in Frangakis and Rubin (2002) for estimating treatment effects controlling for a post-treatment variable (in our case the mechanism variable). The basic idea behind Frangakis and Rubin (2002) is to compare treated and control individuals based on the potential values of the post-treatment variable. As stressed in Rubin (2005), when drawing causal inferences it is very important to keep the distinction between observed values of a variable (e.g., observed gestation) and the potential values it represents (e.g., gestation if smoked during pregnancy or not). In the previously mentioned papers this important distinction is missing.

The following ideal situation provides intuition for the definition of our parameters and the challenges faced in their estimation. Suppose we are interested on the effect of a randomly assigned treatment  $T$  on an outcome  $Y$ , and want to learn what part of that effect is through a mechanism  $S$ . Ideally, we would perform a new experiment in which the new (counterfactual) treatment is the same as the original one but blocks the effect of  $T$  on  $S$ ; or, in other words, sets the value of the mechanism variable  $S$  at the level it would have been if this individual were a

---

<sup>3</sup>In 1991 the Andrew W. Mellon Foundation launched the GEI to improve the structure and organization of the PhD programs in the humanities and social sciences. Some of the channels considered by Ehrenberg et al. (2006) are more financial support to graduate students, more course and seminar requirements, higher quality advising, among others.

<sup>4</sup>Another situation in which it is relevant estimating causal mechanisms is in the evaluation of government programs that are a combination of different services (e.g. they are “bundled treatments”), especially when policymakers aim at reforming them. For example, see the discussions in Meyer (1995) and Currie and Neidell (2007) in the context of unemployment insurance reforms and the Head Start program, both in the U.S. A related study is Card and Hyslop (2005) on the different incentives provided by Canada’s Self Sufficiency Project.

<sup>5</sup>Simonsen and Skipper (2006) define the parameter to estimate before discussing its estimation. However, as we will discuss later, they do not acknowledge explicitly the fact that the mechanism variable represents two different potential variables, and the relationship of their parameter to the total average treatment effect is not discussed.

control under the original treatment. We define the net average treatment effect as the difference in mean outcomes of this new experiment and the control treatment. If this counterfactual experiment is available, estimation of the net average treatment effect is straightforward by comparing the average outcomes of the individuals that took this new treatment and those in the control group. Therefore, intuitively, estimation of treatment effects net of a mechanism requires learning about a different treatment from the one we have at hand. This motivates the difficulty in estimating this kind of treatment effects since, unfortunately, the commonly available data may provide limited information about this counterfactual experiment. Given the usual trade-offs between data availability and assumptions, it is not surprising that estimation of causal net treatment effects requires stronger assumptions than estimation of total average effects.

We present two different approaches for estimation of our parameters. The first is based on a functional form assumption relating the potential outcomes of interest, plus a selection-on-observables assumption analogous to that used for estimation of (total) average treatment effects. The second is based on estimation of the causal net average treatment effect for a particular subpopulation: those individuals for which the treatment does not affect the mechanism variable. We present each of these approaches for the case in which the treatment is randomly assigned and for the case in which selection into the treatment is based on a set of observable covariates. For comparison, we also discuss a set of assumptions under which the usual approach of controlling for the observed value of the mechanism variable can be interpreted as a causal net average treatment effect.

The paper is organized as follows. Section 2 reviews related literature. Section 3 presents the general framework and defines our parameters of interest. Section 4 analyzes the identification and estimation of our parameters. Section 5 presents the results from two empirical applications that illustrate the methods discussed in this paper. The first application analyzes the importance of the “lock-in” effect of a major training program on participant’s earnings, while the second analyzes the importance of gestation as a mechanism for the effect of smoking during pregnancy on the incidence of low birth weight. Concluding remarks are provided in the last section of the paper.

## 2 Related Literature

Our goal is to analyze two related effects: a causal mechanism through which a treatment affects an outcome, and the causal effect of the treatment net of this mechanism. To achieve this goal we employ the potential outcomes framework and, more specifically, build on literature related to the estimation of causal effects adjusting for covariates that are affected by the treatment. This literature relates to our goal since estimating the causal mechanism of a

treatment implies accounting for variables that are observed after the treatment and that are affected by it.<sup>6</sup>

Rosenbaum (1984) analyzes the consequences of adjusting for covariates that are affected by the treatment using the potential outcomes framework. He concludes that estimators adjusting for these variables are generally biased, and specifies sufficient conditions under which controlling for such variables yields the average treatment effect (*ATE*).<sup>7</sup> Trivially, these conditions imply that the *ATE* can be identified when the post-treatment variables are not affected by the treatment, in which case they can be regarded as pre-treatment variables. Rosenbaum (1984) also defines the “net treatment difference” (*NTD*), a parameter that is estimated by simply adjusting for the observed value of the post-treatment variable and is argued to “provide insight into the treatment mechanism”, even though it lacks causal interpretation.

More recently, Frangakis and Rubin (2002) introduced the concept of principal stratification to define causal effects when controlling for post-treatment variables in a variety of settings. Principal stratification, which will be further discussed in the following section, builds on the concept of potential values and the central idea of comparing “comparable” individuals to obtain causal effects. They define causal effects by comparing individuals with the same potential values of the post-treatment variable under each of the treatment arms. In this paper, we define our estimands to have causal interpretation based on principal stratification.

Some of the work closer to ours is in Mealli and Rubin (2003) and Rubin (2004). Both papers motivate the use of principal stratification to clarify and analyze the discussion of “direct” versus “indirect” causal effects, which answer questions similar to the ones we consider here. A direct effect corresponds to a causal effect of a treatment net of a post-treatment variable, while an indirect effect corresponds to the causal effect of a treatment that is mediated by another variable (i.e., a mechanism). The main goal in both papers is to illustrate that the use of principal stratification clarifies the concepts of causality when controlling for post-treatment variables, and that other methods that ignore potential values of variables influenced by the treatment can potentially lead to misleading causal conclusions.<sup>8</sup>

Even though the concepts of direct and indirect effects in the previous papers are similar to the causal mechanism and causal net effects we define and analyze here, there are important

---

<sup>6</sup>Note that to assess the importance of a potential mechanism through which a treatment works one needs to have a measurement of it. In practice, the mechanisms considered may depend on the availability of a variable measuring them. Since such variables are measured after the treatment, we indistinctly refer to them as “post-treatment variables” or “mechanisms”.

<sup>7</sup>More recently, Imbens (2004) also warns about similar pitfalls when controlling for post-treatment variables affected by a treatment, while Lechner (2005) specifies more explicit conditions to assess the endogeneity bias introduced when controlling for variables influenced by the treatment. Both deal with situations in which interest lies on identification of the *ATE*.

<sup>8</sup>Mealli and Rubin (2003) discuss the application of principal stratification to analyze the assumptions needed to estimate direct versus indirect effects in the context of the temporal causal relationships between health and socioeconomic status analyzed by Adams et al. (2003). Similar discussions and illustrations are provided in Rubin (2004).

differences between those papers and ours. First, the relationship of the concepts of direct and indirect effects to the (total)  $ATE$  is not discussed in those papers, while the parameters to be presented here intuitively decompose the  $ATE$  into two effects (mechanism and net effect) that are relevant for policy purposes. Second, as we explain later, the concept of direct effect is a special case of our causal net average treatment effect for a specific subpopulation. Finally, we formally discuss identification and estimation under different assumptions and present empirical applications, which none of the other papers do.

Another strand of literature related to our work is that of Robins and Greenland (1992) and Petersen, Sinisi and van der Laan (2006) in the field of epidemiology (see also references therein). Robins and Greenland (1992) make a similar distinction of direct and indirect effects and present conditions under which they can be estimated; whereas Petersen, Sinisi and van der Laan (2006) discuss the related concepts of “controlled” and “natural” direct effects.<sup>9</sup> While these concepts are related to ours in some respects, there are important differences as well. Most notably, these papers do not employ the concept of principal stratification we employ here and do not distinguish the potential values of the post-treatment variable in their assumptions for identification. In our view, this obscures the assessment of the plausibility of the assumptions and, as discussed in Rubin (2005), may lead to invalid causal conclusions.<sup>10</sup>

### 3 The Estimands of Interest

#### 3.1 Definition of Estimands

We employ the potential outcomes framework (Neyman, 1923; Rubin, 1974) now widely used in the program evaluation literature. Assume we have a random sample of size  $N$  from a large population. For each unit  $i$  in the sample, let  $T_i \in \{0, 1\}$  indicate whether the unit received the treatment of interest ( $T_i = 1$ ) or the control treatment ( $T_i = 0$ ). We are interested on the effect of the treatment  $T$  on an outcome  $Y$ . As usual in this framework, let  $Y_i(1)$  denote the potential outcome for individual  $i$  under treatment and  $Y_i(0)$  denote the potential outcome under the control treatment. The (population) average treatment effect is hence given by  $ATE = E[Y(1) - Y(0)]$ .<sup>11</sup> We are interested on analyzing the part of the  $ATE$  that works through a mechanism variable  $S$ , and the causal effect of  $T$  on  $Y$  net of the effect through  $S$ .

---

<sup>9</sup>These two concepts are discussed later in section 3.3.

<sup>10</sup>Robins and Greenland (1992) is actually an application of a more general literature on the estimation of dynamic causal effects (e.g. Robins (1986) in epidemiology and more recently Lechner and Miquel (2005) in economics). In this literature, the identification of causal effects from sequences of interventions is analyzed. Accounting for the possibility of a dynamic selection process implies making assumptions about the dependence of both the sequence of treatments and the final outcome of interest on intermediate outcomes. We abstract from modeling dynamics explicitly, so we concentrate on a static model of causal effects as in Robins and Greenland (1992).

<sup>11</sup>Another treatment effect usually analyzed in the literature is the average treatment effect on the treated, which is given by  $ATT = E[Y(1) - Y(0)|T = 1]$ . For ease of exposition we focus on decomposing the  $ATE$ , but the discussion and results can easily be extended to the  $ATT$  and other parameters, as is the case in section 5.1.

Since  $S$  is affected by the treatment, we must consider its potential values, denoted by  $S_i(1)$  and  $S_i(0)$ . Hence,  $S_i(1)$  represents the value of the post-treatment variable individual  $i$  would get if exposed to the treatment, and  $S_i(0)$  represents the value she would get if exposed to the control treatment.<sup>12</sup> For each unit  $i$ , we observe the vector  $(T_i, Y_i^{obs}, S_i^{obs})$ , where  $Y_i^{obs} = T_i Y_i(1) + (1 - T_i) Y_i(0)$  and  $S_i^{obs} = T_i S_i(1) + (1 - T_i) S_i(0)$ . That is, we only observe  $Y_i(1)$  and  $S_i(1)$  for units that receive treatment and  $Y_i(0)$  and  $S_i(0)$  for units in the control group. This is the so-called “fundamental problem of causal inference” (Holland, 1986). It is important to stress the fact that  $S^{obs}$  represents two different potential variables:  $S(1)$  for treated units and  $S(0)$  for controls.<sup>13</sup>

In our case, it is convenient to let the potential outcomes be a function of the mechanism variable  $S$ . For each individual  $i$ , define the “composite” potential outcomes  $Y_i(\tau, \zeta)$ , where the first argument refers to one of the treatment arms ( $\tau \in \{0, 1\}$ ) and the second argument represents one of the potential values of the post-treatment variable  $S$  ( $\zeta \in \{S_i(0), S_i(1)\}$ ). Using this notation, we can consider the following composite potential outcomes for any given individual:

1.  $Y_i(1, S_i(1))$ : this is the potential outcome the individual would obtain if she received treatment and post-treatment variable level  $S_i(1)$ . It includes the total effect of receiving treatment on  $Y$  (i.e., through  $S$  or not). This is exactly the potential outcome  $Y_i(1)$  under the treatment.
2.  $Y_i(0, S_i(0))$ : this is the potential outcome when no treatment is received and the post-treatment variable value is  $S_i(0)$ . It is the outcome an individual would obtain if the treatment is not given to her and if the value of her post-treatment variable is not altered either. This is exactly the potential outcome  $Y_i(0)$  under the control treatment.
3.  $Y_i(1, S_i(0))$ : this is the potential outcome the individual would receive if she were exposed to the treatment but kept the level of  $S$  she would obtain had not been treated. In other words, it is the outcome the individual would get if we were to give her the treatment but held the value of her post-treatment variable at  $S_i(0)$ . As a result, this potential outcome includes the effect of  $T$  on  $Y$  that is *not* through  $S$ . This is the key potential outcome we use to define net and mechanism effects below.<sup>14</sup>

Based on these composite potential outcomes, the following three individual-level comparisons are of interest for our purposes:

---

<sup>12</sup>Note that  $S$  is not restricted to be binary.

<sup>13</sup>We also adopt the stable unit treatment value assumption (SUTVA) following Rubin (1980). This assumption is common throughout the literature, and it implies that the treatment effects at the individual level are not affected either by the mechanism used to assign the treatment or by the treatment received by other units. In practice, this assumption rules out general equilibrium effects of the treatment that may impact individuals.

<sup>14</sup>For completeness, note that  $Y_i(0, S_i(1))$  is the potential outcome the individual would obtain when the treatment is not given to her but she receives a value of the post-treatment variable equal to  $S_i(1)$ .

- (a)  $Y_i(1, S_i(1)) - Y_i(0, S_i(0))$ : this represents the usual individual total treatment effect or *ITTE*. For example, the total effect of smoking during pregnancy on birth weight.
- (b)  $Y_i(1, S_i(1)) - Y_i(1, S_i(0))$ : this difference gives the effect of a change in  $S$ , which is *due* to  $T$ , on the outcome  $Y$ . Here we hold constant all other ways in which  $T$  may affect  $Y$ , since  $Y_i(1, S_i(0))$  already considers the effect of  $T$  on  $Y$  through other channels. For example, this difference shows the effect of a change in gestation time due to smoking on birth weight, holding all other effects of smoking during pregnancy fixed. We call this the *individual causal mechanism effect*.
- (c)  $Y_i(1, S_i(0)) - Y_i(0, S_i(0))$ : this difference gives the effect of  $T$  on  $Y$  when the value of the post-treatment variable is held constant at  $S_i(0)$ . Hence, it is the part of the effect of  $T$  on  $Y$  that is *not* due to a change in  $S$  caused by the treatment. For example, the effect of smoking during pregnancy on birth weight that is *not* due to a change in gestation time caused by smoking. We call this the *individual causal net effect*.

Given these comparisons, we can decompose the individual total treatment effect in (a) as:

$$ITTE = [Y_i(1, S_i(1)) - Y_i(1, S_i(0))] + [Y_i(1, S_i(0)) - Y_i(0, S_i(0))]. \quad (1)$$

Hence, at the individual level, the total effect is decomposed into the part of the effect due to a change in  $S$  because of a change in  $T$  (first term in (1) or mechanism effect); and the part of the effect holding  $S$  fixed at  $S(0)$  (second term in (1) or net effect).

The population (total) average treatment effect (*ATE*) can be decomposed in a similar way as:

$$ATE = E[Y(1, S(1)) - Y(1, S(0))] + E[Y(1, S(0)) - Y(0, S(0))]. \quad (2)$$

As in (1), the first term reflects the part of the average treatment effect that is due only to a change in  $S$  because of a change in  $T$ , and the second term shows the part of the average effect holding  $S$  fixed at  $S(0)$ .

It is clear from the decomposition in (2) that we need to make treatment comparisons adjusting for the post-treatment variable  $S$  that is affected by the treatment. In order to causally interpret our parameters of interest, we employ the concept of principal stratification developed in Frangakis and Rubin (2002).

In the potential outcomes framework, a causal effect must be a comparison of potential outcomes for the same group of individuals under treatment and control.<sup>15</sup> Frangakis and Rubin (2002) (hereafter FR) introduce the concept of principal stratification for defining causal effects

---

<sup>15</sup>For example, assuming that selection into treatment is random conditional on a set of covariates  $X$ , we can write  $ATE = E\{E[Y(1) - Y(0)|X]\}$ . In this case, a causal effect is defined by a comparison of the potential outcomes  $Y(1)$  and  $Y(0)$  for those units with the same vector of covariates  $X$ —units with the same value of  $X$  are comparable.



in the presence of post-treatment variables. The idea is to define the “same group of individuals” based on the potential values of the post-treatment variable. In FR terminology, the basic principal stratification with respect to post-treatment variable  $S$  is a partition of individuals into groups such that within each group all individuals have the same vector  $\{S(0) = s_0, S(1) = s_1\}$ , where  $s_0$  and  $s_1$  are generic values of  $S(0)$  and  $S(1)$ , respectively. A principal effect with respect to a principal strata is defined as a comparison of potential outcomes within that strata. Since principal strata are not affected by treatment assignment, individuals in that group are indeed comparable and thus principal effects are causal effects.<sup>16</sup>

Based on FR, we condition on the principal strata  $\{S(0) = s_0, S(1) = s_1\}$  in order to give a causal interpretation to our parameters. Write the ATE controlling for  $S(0)$  and  $S(1)$  as

$$ATE = E\{E[Y(1, S(1)) - Y(0, S(0)) | S(0) = s_0, S(1) = s_1]\} = E[\tau(s_0, s_1)], \quad (3)$$

where the outer expectation is taken over  $S(0)$  and  $S(1)$  and we let  $\tau(s_0, s_1) = E[Y(1, S(1)) - Y(0, S(0)) | S(0) = s_0, S(1) = s_1]$ . Then, using the same decomposition as in (2) we have:

$$ATE = E\{E[Y(1, S(1)) - Y(1, S(0)) | S(0) = s_0, S(1) = s_1]\} \\ + E\{E[Y(1, S(0)) - Y(0, S(0)) | S(0) = s_0, S(1) = s_1]\}. \quad (4)$$

We define the (causal) net average treatment effect or  $NATE$  as:

$$NATE = E\{E[Y(1, S(0)) - Y(0, S(0)) | S(0) = s_0, S(1) = s_1]\} \quad (5)$$

and the (causal) mechanism average treatment effect or  $MATE$  as:<sup>17</sup>

$$MATE = E\{E[Y(1, S(1)) - Y(1, S(0)) | S(0) = s_0, S(1) = s_1]\}. \quad (6)$$

In the remaining of this section we discuss  $NATE$  and  $MATE$ , and in the next section we concentrate on different assumptions and conditions upon which they can be identified and estimated.

### 3.2 Discussion of Estimands

An intuitive way to think about our estimands is to consider  $Y(1, S(0))$  as the potential outcome of an alternative counterfactual experiment in which the treatment is the same as the

---

<sup>16</sup>FR’s idea of principal stratification is closely related to the local average treatment effect interpretation of instrumental variables (e.g., Imbens and Angrist, 1994). For example, in the terminology of Imbens and Angrist (1994), the group of “compliers” is the set of individuals that always comply with their treatment assignment regardless of whether their assignment is to treatment ( $T = 1$ ) or control group ( $T = 0$ ). Therefore, for this group  $\{S(0) = 0, S(1) = 1\}$ , where  $S$  is an indicator of actual treatment reception.

<sup>17</sup>Although we could have used the terms “direct effect” and “indirect effect” to define our parameters, we prefer our names for two reasons. First, they differ in important ways from direct effects as defined in Mealli and Rubin (2003) and Rubin (2004), as discussed later in section 3.3. Second, our names make clear that these effects are considered with respect to a particular mechanism  $S$ . Strictly speaking, a “pure” direct effect would have to net out all possible mechanisms through which the treatment may affect the outcome.

original one but blocks the effect of  $T$  on  $S$  by holding  $S$  fixed at  $S_i(0)$  for each individual  $i$ . The causal net treatment effect or  $NATE$  for individual  $i$  is the difference between the outcome of this alternative treatment,  $Y_i(1, S_i(0))$ , and  $Y_i(0, S_i(0))$  from the original control treatment. Similarly, her causal mechanism treatment effect or  $MATE$  is given by the difference in the potential outcomes of the original treatment and the alternative one.

An important property of  $NATE$  in (5) is that it includes not only the part of the  $ATE$  that is totally unrelated to the mechanism variable  $S$ , but also the part of the  $ATE$  that results from a change *in the way*  $S$  affects  $Y$ . That is, even though the level of  $S$  is held fixed at  $S(0)$ , the treatment may still affect the way  $S$  affects the outcome, and this is counted as part of  $NATE$ . To illustrate this point, consider one of our empirical applications in which we analyze the lock-in effect as a causal mechanism of a job training program, i.e. the labor market experience lost due to participation in the program. If participants lose substantial labor market experience due to the program, and this negatively affects their future earnings, a policy maker may want to change the program to be as the original one but without affecting labor market experience (i.e., holding experience fixed at  $S(0)$ ). In this case the policy maker would like to know the average effect of this alternative training program on future earnings. This effect would include not only the part of the effect of the program on earnings that is totally unrelated to experience, but it would also include the effect of the program on how experience affects wages, i.e., the program’s effect on the returns to experience.  $NATE$  takes this into account, correctly measuring what the effect of this alternative treatment (training program) would be.

We argue that including the effect of  $T$  on how  $S$  affects  $Y$  (i.e., returns to  $S$ ) in  $NATE$  is more relevant from a policy perspective, compared to a different parameter that holds constant the way  $S$  affects  $Y$ . The reason is that a policy maker typically has some degree of control over  $S$ , while very rarely over how  $S$  affects  $Y$ . In the previous example, the administrators of a training program have some degree of control over the level of labor market experience that might be lost due to the time spent in training (e.g. by offering training while on the job or shortening the time to completion of the program), but it seems unlikely that they could influence the (potentially) different returns to experience that the market awards to trained versus non-trained individuals.<sup>18</sup> Our argument is consistent with the notion of a “treatment” being an intervention that can be potentially applied to each individual (e.g., Holland, 1986).

As a final remark about  $NATE$  and  $MATE$ , we note that their definitions conform to intuition in the following two extreme cases. First, consider the situation in which all the effect of  $T$  on  $Y$  works exclusively through  $S$  for the entire population. In this case  $Y(1, S(0)) = Y(0, S(0))$  and, as expected,  $NATE = 0$  and  $MATE = ATE$  from equations (5) and (6),

---

<sup>18</sup>One potentially interesting case where the policymaker might have some degree of influence on how  $S$  affects  $Y$  is when general equilibrium effects due to the treatment are present.

respectively. Second, consider the situation in which none of the effect of  $T$  on  $Y$  is through  $S$ , in which case  $NATE$  should equal  $ATE$  and  $MATE$  should be zero. This can arise due to two reasons: either  $S$  does not affect  $Y$  (even though  $S$  may be affected by  $T$ ) and thus  $\{S(1), S(0)\}$  is independent of  $\{Y(1), Y(0)\}$ ; or  $T$  simply does not affect  $S$  and thus  $S(1) = S(0)$ . Regardless of the reason, the consequence is that  $Y(1, S(1)) = Y(1, S(0))$  and thus (5) and (6) imply  $NATE = ATE$  and  $MATE = 0$ , respectively.

### 3.3 Relation of the Estimands to Other Parameters in the Literature

As discussed in section 2, Rosenbaum (1984) defines the  $NTD$ . This parameter is characterized by conditioning on the observed post-treatment variable and, without further assumptions, has no causal interpretation if the post-treatment variable is affected by the treatment. It can be written as  $NTD = E\{E[Y(1) - Y(0) | S^{obs}]\}$ . The reason for  $NTD$ 's lack of causal interpretation is that it compares individuals with the same values of  $S^{obs}$ . Since  $S^{obs}$  represents two different potential variables,  $S(1)$  and  $S(0)$ , units with the same value of  $S^{obs}$  are generally not comparable. This point is further discussed and illustrated in Mealli and Rubin (2003), Rubin (2004), and Rubin (2005).<sup>19</sup> In contrast, by conditioning on principal strata,  $NATE$  explicitly accounts for the possibility that the post-treatment variable is affected by the treatment. Furthermore, our parameters effectively decompose the  $ATE$  into causal mechanism and net effects (see (4)).

Although both our estimands and the concepts of direct and indirect effects in Mealli and Rubin (2003) and Rubin (2004) rely on the idea of principal stratification and thus can be interpreted as causal effects, they differ in other aspects. Mealli and Rubin (2003) define a direct effect as a comparison of  $Y(1)$  and  $Y(0)$  within the stratum for which  $S(0) = S(1) = s$ , which implies  $Y(1, S(1)) = Y(1, S(0))$ . Using our notation in (3), we can write their direct effect as  $DE(s) = \tau(s, s)$ , which corresponds to  $NATE$  in (5) defined for *this particular subpopulation or strata*. More generally, we can define the direct average effect as  $DAE = E[\tau(s, s)]$ , which is the average of the direct effects over the possible values  $s$  of  $S$ . Note that, unless  $NATE$  is constant in the population,  $DAE$  will differ from  $NATE$ . Moreover,  $DAE$  does not decompose the  $ATE$  in the way  $NATE$  does because  $DAE$  ignores all the individuals for which  $S_i(1) \neq S_i(0)$ . Finally, note that the definition of the direct effect effectively rules out a mechanism effect, since it is only defined for subpopulations for which there is no mechanism effect. For these reasons,  $NATE$  and  $MATE$  are more general and, in our view, more relevant for policy purposes.

Finally, there are other parameters related to  $NATE$  that have been used in the epidemiol-

---

<sup>19</sup>Another way to see the problem of conditioning by  $S^{obs}$  is to note that when estimating the  $NTD$  based on the observed data we are implicitly assuming that the treatment is “randomly assigned” conditional on  $S^{obs}$  so that we can write  $E[Y(1) | S^{obs}] = E[Y^{obs} | T = 1, S^{obs}]$ . However, in general, we can infer something about the treatment assignment  $T$  based on  $S^{obs}$  and hence the assumption fails. See Rubin (2005) for further discussion.

ogy literature: the controlled direct effect (*CDE*) and the natural direct effect (*NDE*)<sup>20</sup>. The *CDE* at a specific value  $\bar{s}$  of  $S$  can be written as  $CDE = E[Y(1, S(1) = \bar{s}) - Y(0, S(0) = \bar{s})]$ . The *CDE* gives the average difference between the counterfactual outcome under the two treatment arms controlling for the value of the mechanism variable at  $\bar{s}$ . While this parameter may be informative in some applications, in our view has some undesirable features for the estimation of net effects. First, it does not decompose the *ATE* into a net and a mechanism effect in the way *NATE* and *MATE* do.<sup>21</sup> Second, since neither of the two potential outcomes used in the definition of *CDE* necessarily correspond to the observed outcome ( $Y^{obs}$ ) for any particular individual, its estimation surely requires stronger assumptions than the ones used for estimation of *NATE*, where at least one of the potential outcomes ( $Y(0, S(0))$ ) is observed for some individuals.<sup>22</sup> Lastly, using the *CDE* to estimate net effects has the undesirable property that, even if in fact the treatment does not affect the mechanism variable  $S$ , the *ATE* may be different from the *CDE* if there is heterogeneity in the effect of  $T$  on  $Y$  along the values of  $S$ . Conversely, as previously discussed for our parameters,  $NATE = ATE$  and  $MATE = 0$  in this case.

The *NDE* used in epidemiology can be written as  $E[Y(1, S(0)) - Y(0, S(0))]$ . Hence, this parameter is similar to *NATE* in (5) with the subtle but important difference that *NATE* conditions on principal strata in order to retain causal interpretation. This distinction becomes crucial when stating assumptions for estimation. For example, Robins and Greenland (1992) and Petersen, Sinisi and van der Laan (2006) do not keep the distinction between  $S(1)$  and  $S(0)$  when stating their assumptions for estimation of *CDE* and *NDE*. Since  $S$  represents two potential variables, this makes difficult assessing the plausibility of the assumptions needed for estimation.<sup>23</sup>

## 4 Identification and Estimation of the Parameters of Interest

In this section we discuss identification and estimation of the parameters *NATE* and *MATE* defined in section 3.1. We focus our attention on *NATE* since, by definition, we can obtain  $MATE = ATE - NATE$ . We start by discussing the type of assumptions needed to interpret the standard approach of directly controlling for  $S^{obs}$  as an estimate of *NATE*. Unfortunately, these assumptions are too strong to be useful in practice. Next, we present

<sup>20</sup>See, for instance, Robins and Greenland (1992) and Petersen, Sinisi and van der Laan (2006).

<sup>21</sup>For example, we could write the *ATE* as:  $ATE = E[Y(1, S(1)) - Y(1, S(1) = \bar{s})] + CDE + E[Y(0, S(0) = \bar{s}) - Y(0, S(0))]$ . The first term gives the average effect of giving the treatment to the individuals and moving the value of the post-treatment variable from  $\bar{s}$  to  $S(1)$ . The second term represents the average effect of giving the control treatment to the individuals and moving the value of the post-treatment variable from  $S(0)$  to  $\bar{s}$ . These two effects are hard to interpret as mechanism effects of  $T$  on  $Y$  through  $S$ .

<sup>22</sup>See following section for details.

<sup>23</sup>For further discussion on the importance of conditioning on principal strata see Rubin (2004, 2005) and Mealli and Rubin (2003).

two different estimation strategies for each of two treatment-assignment mechanisms. We first consider the situation in which the treatment is randomly assigned. This case is important in its own right given the existence of social experiments in economics, such as the one used in our first empirical application. We then discuss the case in which the treatment is assumed to be random given a set of observed covariates.

Regardless of the mechanism used to assign the treatment, identification and estimation of *NATE* faces two challenges. First, we have to take into account that for each unit under study only one of the potential values of the post-treatment variable is observed:  $S^{obs}$  represents  $S(1)$  for treated units and  $S(0)$  for controls units. This implies that the principal strata  $\{S(0) = s_0, S(1) = s_1\}$ , necessary for a causal interpretation of *NATE*, is never observed. Note that  $S$  can be regarded as an outcome, and thus the distribution of the principal strata equals the joint distribution of the potential outcomes  $\{S(1), S(0)\}$ , which is not easily identifiable (e.g., Heckman, Smith and Clements, 1997). The second challenge is that a key potential outcome needed for estimation of *NATE*,  $Y_i(1, S_i(0))$ , is generally not observed—this is in contrast to the case of estimation of the *ATE*, where only one of the relevant potential outcomes is missing for every unit. In an ideal situation in which we could perform the counterfactual experiment and observe  $Y(1, S_i(0))$  for some units, none of these two challenges would arise and estimation of *NATE* would be straightforward.<sup>24</sup>

Despite the missing data challenges that result from the unavailability of the counterfactual experiment, we can still impose assumptions under which *NATE* can be identified from the available data. We present below two strategies for its estimation. Given the challenges faced, it is not surprising that the assumptions needed for estimation of *NATE* are stronger than those typically needed for estimation of *ATE*. This reflects the difficulty of answering the question at hand given the available data, and illustrates the usual trade off between data quality and assumptions needed for estimation. Importantly, one must have a clear understanding of the assumptions needed to identify and estimate our parameters in practice in order to be able to evaluate their plausibility in any particular setting.

#### 4.1 Assumptions under which controlling directly for $S^{obs}$ yields *NATE*

It is important to state conditions under which the standard approach of controlling for the observed value of the post-treatment variable ( $S^{obs}$ ), and possibly a set of covariates  $X$ , yields *NATE*. Importantly, the kind of assumptions needed are very strong—certainly stronger than those we present in the following sections to identify our parameters.

Consider the following parameter that is representative of the standard approach, where the

---

<sup>24</sup>Under this counterfactual treatment we have that  $S(1) = S(0)$  for all units (by construction of the counterfactual treatment), and the potential outcome  $Y(1, S(0))$  would be observed for those who received this treatment.

second line uses the fact that  $S^{obs}$  represents  $S(0)$  or  $S(1)$  depending on the treatment received:

$$\begin{aligned}\gamma &= E\{E[Y^{obs}|T = 1, S^{obs} = s, X = x] - E[Y^{obs}|T = 0, S^{obs} = s, X = x]\} \\ &= E\{E[Y^{obs}|T = 1, S(1) = s, X = x] - E[Y^{obs}|T = 0, S(0) = s, X = x]\}. \quad (7)\end{aligned}$$

A set of sufficient conditions under which  $\gamma = ATE$  are (Rosenbaum, 1984): (i)  $S(1) = S(0)$  for all subjects in the population (“unaffected post-treatment variable”), and (ii) the treatment assignment is ignorable in the sense that  $\{Y(1), Y(0)\} \perp T|X$  and  $0 < \Pr(T = 1|X) < 1$  for all  $X$ .<sup>25</sup> Intuitively, the issue when estimating the  $ATE$  based on (7) is that the outer expectation should be taken with respect to the distribution  $\Pr(S(1)|X)$  for the first term and with respect to  $\Pr(S(0)|X)$  for the second. As a result, if  $S$  is affected by  $T$ , bias will arise from averaging both terms over  $\Pr(S^{obs}|X)$  instead. In other words, looking at units with the same values of  $S^{obs}$  in fact compares treated units with  $S(1) = s$  to control units with  $S(0) = s$ , which are in general not comparable.<sup>26</sup> Condition (i) implies that  $S^{obs} = S(1) = S(0)$ , ensuring that the averaging is over the correct distribution; nonetheless, this condition is too strong since it rules out an effect of  $T$  on  $S$ .

Unfortunately, the same conditions are needed to have  $\gamma = NATE$ . Even if we were to assume that people in different strata are comparable conditional on  $X$ ,<sup>27</sup> we still need to assume that  $S_i(1) = S_i(0) = s$  for all units in order to have  $Y_i(1, S_i(0)) = Y_i(1, S_i(1))$  and thus deal with the problem that  $Y_i(1, S_i(0))$  is unobserved. Only then could we have  $\gamma = NATE$ . However, the condition that  $S_i(1) = S_i(0) = s$  for all units again rules out the role of  $S$  as a mechanism by assumption, rendering it too strong to be used in practice. In the following sub-sections we present weaker assumptions that allow us to estimate  $NATE$  and  $MATE$ .

## 4.2 Identification and Estimation based on $Y(1, S(1))$

We first consider the case in which individuals are randomly assigned to the treatment. We keep the following assumption until our discussion of non-random treatment assignment in section 4.4.

**Assumption 1**  $Y(1, S(1)), Y(0, S(0)), Y(1, S(0)), S(1), S(0) \perp T$

Under this assumption, the treatment received by each individual is independent of her potential outcomes and potential values of the post-treatment variable. Note that Assumption 1 implies  $Y(1, S(1)), Y(0, S(0)), Y(1, S(0)) \perp T|\{S(1), S(0)\}$ , so that potential outcomes are independent of the treatment given the principal strata.<sup>28</sup>

<sup>25</sup>As in Dawid (1979), we write  $X \perp Y$  to denote independence of  $X$  and  $Y$ .

<sup>26</sup>Yet another way to see the problem of estimating  $ATE$  controlling for  $S^{obs}$  is to regard  $S^{obs}$  as an endogenous control variable since it is affected by the treatment. See Lechner (2005).

<sup>27</sup>In which case the groups with  $\{T = 1, S(1) = s, X = x\}$  and  $\{T = 0, S(0) = s, X = x\}$  would be comparable.

<sup>28</sup>See, for instance, Lemma 4 in Dawid (1979).

Let us start by considering the challenge that the principal strata  $\{S(0) = s_0, S(1) = s_1\}$  is not observed. Note that identification of the principal strata is challenging since it entails determining the effect of the treatment  $T$  on the intermediate outcome  $S$  for every individual using only the marginal distributions of  $S(1)$  and  $S(0)$  for treated and controls, respectively. We follow an approach analogous to the commonly-used selection on observables framework in program evaluation (e.g., Imbens, 2004) and assume that the principal strata is independent of the potential outcomes given a rich set of covariates  $X$ .

**Assumption 2**  $Y(1, S(1)), Y(0, S(0)), Y(1, S(0)) \perp \{S(1), S(0)\} | X$

Assumption 2 implies that individuals in different strata are comparable once we condition on a set of covariates  $X$ , ruling out the existence of variables not included in  $X$  that simultaneously affect the principal strata an individual belongs to and her potential outcomes (i.e., confounders). Assumption 2 also implies that by conditioning on  $X$  we rule out confounders of the relationship between (i) each of the potential values of  $S$  ( $S(1)$  and  $S(0)$ ) and the potential outcomes, and (ii) any function of  $S(1)$  and  $S(0)$  and the potential outcomes. In particular, note that individuals with different treatment effects of  $T$  on  $S$  (i.e.,  $S(1) - S(0)$ ) are comparable conditional on  $X$ . Finally, Assumptions 1 and 2 imply that  $Y(1, S(1)), Y(0, S(0)), Y(1, S(0)) \perp \{T, S(1), S(0)\} | X$ , so that control and treated units in different strata, but with the same values of covariates, are comparable.<sup>29</sup>

The other challenge in the estimation of *NATE* is making inferences about a potential outcome that is not usually observed,  $Y(1, S(0))$ . One approach is to use the information in  $Y(1, S(1))$ , and possibly  $Y(0, S(0))$ , to learn about  $Y(1, S(0))$ . This can be done in many different ways, with the specific assumption to be made depending on what is judged plausible in the particular application at hand. We present here one assumption to illustrate the approach. Suppose that the conditional expectations of the potential outcomes  $Y(1, S(0))$  and  $Y(1, S(1))$  have the same functional form in terms of  $\{X, S(0)\}$  and  $\{X, S(1)\}$ , respectively, but the former sets  $S(1) = S(0)$ . As a simple example, let  $E[Y(1, S(1))]$  be of the form  $E[Y(1, S(1)) | S(1), X] = a_1 + b_1 S(1) + c_1 X$ . Then, this assumption implies that  $E[Y(1, S(0)) | S(0), X] = a_1 + b_1 S(0) + c_1 X$ . We can state this assumption more generally as follows:

**Assumption 3** Let  $E[Y(1, S(1)) | S(1) = s_1, X = x] = f_1(S(1), X)$ . Then,

$$E[Y(1, S(0)) | S(0) = s_0, X = x] = f_1(S(0), X).$$

---

<sup>29</sup>Note the importance of stating the assumptions used in terms of principal strata as opposed to using simply  $S$ , as commonly done in the literature (e.g., Petersen, Sinisi and van der Laan, 2006; Simonsen and Skipper, 2006). Principal strata is not affected by  $T$ , and it acknowledges the fact that  $S$  represents two potential variables:  $S(1)$  and  $S(0)$ . If we were to use  $S$  instead of the principal strata in Assumption 2, its interpretation (which is needed to gauge its plausibility in practice) would be obscured by the fact that  $S$  is affected by the treatment.

We offer a few comments on this assumption. First, Assumption 3 directly acknowledges that we are trying to learn about a counterfactual treatment based on the information available on the original treatment. Second, we clarify that regarding  $Y(1, S(0))$  as the outcome of the counterfactual treatment does not imply Assumption 3. The definition of  $Y_i(1, S_i(0))$  implies that  $Y_i(1, S_i(0)) = Y_i(1, S_i(1))$  for those units with  $S_i(1) = S_i(0)$ ; however, for those with  $S_i(1) \neq S_i(0)$  it is not necessarily the case that  $Y_i(1, S_i(0))$  has the same functional form as  $Y_i(1, S_i(1))$  but setting  $S_i(1) = S_i(0)$ . Finally, note that Assumption 3 implies that the covariates  $X$  and the mechanism variable  $S$  affect the outcome in the same way in both the original and counterfactual treatments. Hence, as discussed in section 3.2,  $NATE$  will include the effect of  $T$  on the way  $S$  affects  $Y$ .

Under Assumptions 1-3 we can identify  $NATE$  by writing it as a function of observed variables as:

$$\begin{aligned}
NATE &= E \{E[Y(1, S(0)) - Y(0, S(0)) | S(0) = s_0, S(1) = s_1, X = x]\} \\
&= E \{E[Y(1, S(0)) | S(0) = s_0, S(1) = s_1, X = x]\} \\
&\quad - E \{E[Y(0, S(0)) | S(0) = s_0, S(1) = s_1, X = x]\} \\
&= E \{E[Y(1, S(0)) | S(1) = s_1, X = x]\} - E \{E[Y(0, S(0)) | S(0) = s_0, X = x]\} \\
&= E \{f_1(S(0), X)\} - E \left\{ E \left[ Y^{obs} | T = 0, S^{obs} = s_0, X = x \right] \right\} \tag{8}
\end{aligned}$$

where we have used Assumption 2 in the third equality, Assumptions 1 and 3 in the last equality, and we have that  $E[Y(1, S(1)) | S(1) = s_1, X = x] = f_1(S(1), X) = E[Y^{obs} | T = 1, S^{obs} = s_1, X = x]$ .

In practice, this identification strategy can be implemented as follows: (i) estimate a model for  $E[Y^{obs} | T = 1, S^{obs} = s_1, X = x] = f_1(S(1), X)$ ; (ii) compute  $E[Y(1, S(0)) | S(0) = s_0, X = x] = f_1(S(0), X)$  based on the model in (i); (iii) estimate  $NATE$  based on (8) and  $MATE = ATE - NATE$ . For steps (i) and (ii) a simple way to proceed is to run a linear regression of  $Y^{obs}$  on  $S(1)$  and  $X$  for treated units and evaluate this estimated model on  $S(1) = \widehat{E}[S_i(0)]$ . One may allow this function to be more flexible by employing a polynomial series expansion of  $S(1)$  and interactions with the covariates, for instance.

### 4.3 Estimation of $NATE$ Based on a Specific Subpopulation

In this section we present the second approach to estimate  $NATE$  by focusing on a particular subpopulation or principal strata: those for which  $T$  does not affect  $S$ , so that  $S_i(1) = S_i(0)$ . For them we have that  $Y_i(1, S_i(0)) = Y_i(1, S_i(1))$  and hence  $Y_i(1, S_i(0))$  is observed for those receiving treatment. Therefore, any non-zero causal effect  $Y_i(1, S_i(1)) - Y_i(0, S_i(0))$  in this subpopulation is due to factors different from the change in  $S$  caused by  $T$ . For this particular



subpopulation with  $S_i(1) = S_i(0)$ , we define its local *NATE* (hereafter *LNATE*) as:

$$LNATE = E \{E[Y(1, S(0)) - Y(0, S(0)) | S(0) = s, S(1) = s]\} = E \{\Delta(s)\} \quad (9)$$

where  $\Delta(s) = E[Y(1, S(0)) - Y(0, S(0)) | S(0) = s, S(1) = s]$  is the local *NATE* in the stratum  $S(1) = S(0) = s$ .<sup>30</sup> The key to regard *LNATE* as a causal effect is to note that it is defined within a principal strata, and thus has a causal interpretation (Frangakis and Rubin, 2002). *LNATE* is similar to the “direct effect” discussed in Mealli and Rubin (2003) and Rubin (2004). More precisely, *LNATE* equals the direct average effect (*DAE*) discussed in section 3.3, which is simply the average of the direct effects for all the stratum for which  $S(1) = S(0) = s$ . It is important to remark that the direct effect is a local *NATE* since it is defined for a specific subpopulation, and without further assumptions it does not decompose the population *ATE* into a net and mechanism effect. Note also that the *LNATE* equals the local average treatment effect for this subpopulation (*LATE<sup>sub</sup>*), since  $Y_i(1, S_i(0)) = Y_i(1, S_i(1))$  and there is no mechanism effect by definition. There is precedent in the literature on the estimation and importance of local average treatment effects. In particular, Imbens and Angrist (1994) interpret instrumental variables (*IV*) estimators as estimators of a local average treatment effect (*LATE*).

Suppose for the moment that we know which individuals belong to each of the strata  $S(1) = S(0) = s$  for all values of  $s$ . Under Assumption 1, conditioning on the principal strata  $\{S(1), S(0)\}$  controls for all observed and unobserved individual characteristics reflected in the post-treatment variable  $S$ . Then, *LNATE* is identified in this subpopulation as

$$\begin{aligned} LNATE &= E \{E[Y(1, S(0)) - Y(0, S(0)) | S(0) = s, S(1) = s]\} \\ &= E \{E[Y(1, S(1)) | S(0) = s, S(1) = s] - E[Y(0, S(0)) | S(0) = s, S(1) = s]\} \\ &= E \left\{ E[Y^{obs} | T = 1, S(0) = s, S(1) = s] - E[Y^{obs} | T = 0, S(0) = s, S(1) = s] \right\} \\ &= E \left\{ E[Y^{obs} | T = 1, S^{obs} = S(0) = S(1) = s] \right. \\ &\quad \left. - E[Y^{obs} | T = 0, S^{obs} = S(0) = S(1) = s] \right\} \end{aligned} \quad (10)$$

Note that the last equation in (10) resembles the standard approach that controls for the observed value of the post-treatment variable (see (7)). The difference is that (10) is calculated for the subpopulation for which  $S(1) = S(0)$ , and therefore, it can be interpreted as a causal effect.<sup>31</sup> In addition, since in this subpopulation  $LNATE = LATE^{sub}$ , under Assumption 1 we can alternatively write *LNATE* as the simple difference in mean outcomes between treated and controls.

---

<sup>30</sup>Note that the outer expectation in *LNATE* is taken over all strata with  $S(1) = S(0) = s$ . For instance, in the context of a binary post-treatment variable, *LNATE* would be the average of the net effects effects for the stratum with  $S(0) = S(1) = 0$  and  $S(0) = S(1) = 1$ .

<sup>31</sup>See also the discussion in section 4.1

Unfortunately, in practice we usually do not have knowledge about a subpopulation for which  $S_i(1) = S_i(0)$ . In some specific situations, though, the nature of the treatment and post-treatment variable may convey some knowledge about a subpopulation for which  $T$  does not affect  $S$ . In such cases one can restrict attention to that subpopulation for estimation of  $LNATE$ . An illustration of this is when a law or regulation restricts the effect of the treatment on the post-treatment variable and results on  $S(1)$  being equal to  $S(0)$  for a known group, allowing the estimation of  $LNATE$  for this group.<sup>32</sup>

When this additional knowledge is not available, which may be the usual case, we propose to use the covariates to find a subpopulation for which there is no effect of  $T$  on  $S$  (or for which such effect is close to zero), and then estimate the corresponding local  $NATE$  for this subpopulation. The approach we follow to find such subpopulation is to rely on predicted values of the potential values of the post-treatment variable based on the covariates  $X$ . Let  $\widehat{S}(1)$  and  $\widehat{S}(0)$  be the estimators of the potential values of  $S$  based on  $X$ .<sup>33</sup> Then, we focus on estimation of the local  $NATE$  for the subpopulation with  $\widehat{S}_i(1) = \widehat{S}_i(0)$ :

$$LNATE_{\widehat{S}} = E\{E[Y(1, S(0)) - Y(0, S(0)) | S(0) = s_0, S(1) = s_1, \widehat{S}(1) = \widehat{S}(0)]\} \quad (11)$$

We condition on the principal strata in the definition of  $LNATE_{\widehat{S}}$  in order to interpret it as a causal effect. Importantly, although the parameter in (11) depends on the unobserved principal strata, we expect that  $S(1) \approx S(0)$  within this subpopulation since  $\widehat{S}(1) = \widehat{S}(0)$ ; hence, we estimate it as if  $S(1) = S(0) = s$ . Similar to the case of  $LNATE$ , two possible estimators of  $LNATE_{\widehat{S}}$  are (i) the estimator of the local average treatment effect for the subpopulation with  $\widehat{S}_i(1) = \widehat{S}_i(0)$  ( $\widehat{LNATE}_{\widehat{S}}^{sub}$ ),<sup>34</sup> and (ii) the estimator based on the last line in (10) for the same subpopulation as in (i) ( $LNATE_{\widehat{S}}$ ).

Clearly, since the subpopulation for which  $\widehat{S}(1) = \widehat{S}(0) = s$  will likely have units for which  $S_i(1) \neq S_i(0)$ , the estimators of  $LNATE_{\widehat{S}}$  previously mentioned will be generally biased. For  $\widehat{LNATE}_{\widehat{S}}^{sub}$  it is not hard to derive an expression for such bias. The key insight is that the bias associated with estimating the  $NATE$  for a population by using an unbiased estimator of  $ATE$  equals  $ATE - NATE = MATE$ . The same holds true for any subpopulation when estimating a local  $NATE$  using an unbiased estimator of the local  $ATE$ . Therefore, the bias associated with  $\widehat{LNATE}_{\widehat{S}}^{sub}$  in estimating  $LNATE_{\widehat{S}}$  equals the “local  $MATE$ ” for the subpopulation with

<sup>32</sup>Note that in this case information on  $S$  is not even necessary to estimate  $LNATE$ .

<sup>33</sup>One can construct estimators  $\widehat{S}(1)$  and  $\widehat{S}(0)$  in different ways. For example, we could use a single matching approach and let  $\widehat{S}_i(0) = S_k^{obs}$  and  $\widehat{S}_i(1) = S_k^{obs}$  if unit  $i$  is a control, and  $\widehat{S}_i(1) = S_k^{obs}$  and  $\widehat{S}_i(0) = S_k^{obs}$  if unit  $i$  is treated, where  $S_k^{obs}$  is the observed value of  $S$  for the closest unit to  $i$  in terms of a given distance measure  $\|X_i - X_j\|$ , with  $T_i \neq T_j$ . Alternatively, we could use a regression function approach to predict  $S(1)$  and  $S(0)$ . Let  $\mu_t(x) = E[S(t) | X = x]$  for  $t = \{0, 1\}$  be the regression functions of the post-treatment potential values on  $X$ . Then, given the estimators  $\widehat{\mu}_t(x)$  of these regression functions, we would define  $\widehat{S}(1)$  and  $\widehat{S}(0)$  for each unit  $i$  as  $\widehat{\mu}_1(x)$  and  $\widehat{\mu}_0(x)$ , respectively.

<sup>34</sup>Under Assumption 1 this estimator equals the difference in mean outcomes between treated and control units for those units with  $\widehat{S}_i(1) = \widehat{S}_i(0)$ .

$\widehat{S}_i(1) = \widehat{S}_i(0)$ . More precisely, letting  $Z_i = 1$  if indeed  $S_i(1) = S_i(0)$  and zero otherwise, the bias associated with estimating  $LNATE_{\widehat{S}}$  using the unbiased estimator  $\widehat{LATE}_{\widehat{S}}^{sub}$  can be written as:<sup>35</sup>

$$\begin{aligned}
Bias(\widehat{LATE}_{\widehat{S}}^{sub}) &= LATE_{\widehat{S}}^{sub} - LNATE_{\widehat{S}} \\
&= E \left[ Y(1, S(1)) - Y(1, S(0)) \mid \widehat{S}(1) = \widehat{S}(0) \right] \\
&= E \left\{ E \left[ Y(1, S(1)) - Y(1, S(0)) \mid Z = z, \widehat{S}(1) = \widehat{S}(0) \right] \mid \widehat{S}(1) = \widehat{S}(0) \right\} \\
&= \Pr(Z = 0 \mid \widehat{S}(1) = \widehat{S}(0)) E[Y(1, S(1)) - Y(1, S(0)) \mid Z = 0, \widehat{S}(1) = \widehat{S}(0)]
\end{aligned} \tag{12}$$

Not surprisingly, the first term states that the closer  $\Pr(S(1) \neq S(0) \mid \widehat{S}(1) = \widehat{S}(0))$  is to zero, the smaller the bias associated with the estimation of  $LNATE_{\widehat{S}}$ . Consequently, the better we predict  $S(1)$  and  $S(0)$ , the smaller the bias will be. In the limit, if we are able to perfectly predict  $S(1)$  and  $S(0)$ , the bias will be zero. The second term equals the local average mechanism effect for those units with  $S(0) \neq S(1)$  and  $\widehat{S}(1) = \widehat{S}(0)$ . The notion that for this subpopulation  $S(0) \approx S(1)$  supports the idea that the second term is likely to be close to zero and thus the bias is likely to be small. Finally, note that the sign of the bias is given by the second term in (12). This may be useful to determine the direction of the bias if information is available about the sign of the mechanism effect for this subpopulation.

In practice, we can assess the plausibility that  $S_i(1) = S_i(0)$  for all units in the subpopulation with  $\widehat{S}_i(1) = \widehat{S}_i(0)$  by employing, for example, a Fisher randomization test for the sharp null hypothesis that the treatment effect of  $T$  on  $S$  is zero for all units (Fisher, 1935). Even though failure to reject this null hypothesis does not mean that the treatment effect is zero for all units, rejecting it is a clear sign that the subpopulation characterized by  $\widehat{S}_i(1) = \widehat{S}_i(0)$  is not appropriate to estimate  $LNATE_{\widehat{S}}$ .<sup>36</sup> In any event, it is important to keep in mind that even if there are small differences between  $S_i(1)$  and  $S_i(0)$  for some units in the subpopulation characterized by  $\widehat{S}_i(1) = \widehat{S}_i(0)$ , this will generally cause bias in the estimation of  $LNATE_{\widehat{S}}$ .

So far we have discussed estimation of  $LNATE_{\widehat{S}}$ . The following assumption allows the interpretation of this estimator as an estimator of  $NATE$  as well.

**Assumption 4** *NATE is constant over the population.*

Under this assumption,  $NATE = LNATE_{\widehat{S}}$ . In addition, the part of the  $ATE$  that is due

<sup>35</sup>To simplify notation we omit the conditioning on the principal strata in the expression below.

<sup>36</sup>Another possibility to informally gauge the plausibility that the treatment does not affect  $S$  within the subpopulation is to estimate  $LNATE_{\widehat{S}}$  both controlling for  $S^{obs}$  as suggested by (10) and not controlling for it (i.e., using  $\widehat{LATE}_{\widehat{S}}^{sub}$ ). Since for any subpopulation with  $S(1) = S(0)$  we have  $LNATE = LATE^{sub}$ , a statistically significant difference between the two can be regarded as evidence that the corresponding subpopulation is not appropriate to estimate a local  $NATE$ .

to the mechanism  $S$  is given by  $MATE = ATE - LNATE_{\hat{S}}$ .<sup>37</sup>

A few observations about Assumption 4 are in order. First, note that the need for this assumption is analogous to the need of the assumption of a constant average treatment effect when estimating  $ATE$  using instrumental variables. In that case we can only identify  $LATE$  for the group of individuals who change treatment status in response to a change in the instrumental variable. However, under the assumption of a constant  $ATE$  we have that  $LATE = ATE$ . Second, we point out that Assumption 4 is weaker than assuming a constant  $ATE$ , which is a relatively common assumption in the literature (see, e.g., Heckman, LaLonde and Smith, 1999). Assumption 4 allows for heterogeneous effects of the treatment on the outcome variable, but such heterogeneity is restricted to work through the mechanism or post-treatment variable  $S$  (i.e. through  $MATE$ ). The plausibility of this assumption can be gauged in light of this observation. Third, we note that the standard approach of controlling directly for  $S^{obs}$  implicitly assumes a stronger condition than Assumption 4, since it imposes a zero mechanism effect ( $MATE$ ) for the population. Finally, when Assumption 4 is judged to be untenable in a particular application,  $LNATE_{\hat{S}}$  can still be an informative parameter for policymakers, just as  $LATE$  commonly is.

For the case of a randomly assigned treatment, one way to implement this estimation strategy is as follows: (i) specify a model (based on  $X$ ) to estimate  $S(1)$  and  $S(0)$ ; (ii) identify the subpopulation for which  $\hat{S}(1) = \hat{S}(0)$ ;<sup>38,39</sup> (iii) For that subpopulation, estimate  $LATE_{\hat{S}}^{sub}$  (without controlling for  $S^{obs}$ ) and/or estimate  $LNATE_{\hat{S}}$  based on (10); (iv) under Assumption 4, estimate  $MATE = ATE - LNATE_{\hat{S}}$ .

In sum, given the reliance on predicted instead of actual potential values of  $S$ , the estimate of  $LNATE_{\hat{S}}$  will be biased in general; although the magnitude and direction of the bias can be judged in light of (12). In addition, with this approach we need the assumption of a constant  $NATE$  in order to estimate  $NATE$  and  $MATE$ . Conversely, it is important to note that under this strategy Assumptions 2 and 3 are not needed; that is, we do not need to assume a specific relation between  $Y(1, S(1))$  and  $Y(1, S(0))$ , nor that selection into principal strata is based on observables. Therefore, the approach described in this section allows for unobservables affecting selection into principal strata and potential outcomes. There are thus trade-offs in employing either one of the two approaches presented, and the specific one to use should depend on the plausibility of the different assumptions in the application at hand.

---

<sup>37</sup>Note that under Assumption 4 we also have  $LNATE = LNATE_{\hat{S}} = NATE$ , where  $LNATE$  is defined as in (9).

<sup>38</sup>If the post-treatment variable under consideration is continuous or if the procedure used to estimate  $S(1)$  and  $S(0)$  yields a continuous variable (and hence the probability of finding someone with  $\hat{S}(1) = \hat{S}(0)$  is zero), one could consider a window within which values of  $\hat{S}(1)$  and  $\hat{S}(0)$  are considered to be equal. As usual, such window will tend to zero as the sample size goes to infinity. Alternatively, one could use a kernel function to give higher weight to observations for which  $\hat{S}(1)$  is closer to  $\hat{S}(0)$ .

<sup>39</sup>At this point, one can use a test like the Fisher randomization test to gauge whether the subpopulation meets the minimum requirements for the estimation of  $LNATE_{\hat{S}}$ .

## 4.4 Identification and Estimation under Non-random Assignment

In the previous sections we analyzed the problem of estimating *NATE* and *MATE* when the treatment  $T$  is randomly assigned. In the absence of an experiment, a common approach in the literature is to assume that selection into treatment is based on a set of observed covariates ( $X$ ) and on unobserved components not correlated with the potential outcomes. This assumption is known in the literature as unconfoundedness, conditional independence, or selection on observables.<sup>40</sup> We extend the framework discussed in sections 4.2 and 4.3 to the case when  $T$  is not randomly assigned using the following unconfoundedness assumption:

**Assumption 5**  $Y(1, S(1)), Y(0, S(0)), Y(1, S(0)), S(1), S(0) \perp T | X$ .

Assumption 5 implies that the treatment received by each individual is independent of her potential outcomes and potential values of the post-treatment variable given the set of covariates  $X$ . Hence, the covariates will now have the additional role of controlling for selection into the treatment.<sup>41</sup>

We also add the following overlap assumption:

**Assumption 6**  $0 < \Pr(T = 1 | X = x) < 1$ , for all  $x$ .

Assumption 6 ensures that in infinite samples we are able to compare treated and control units for all values of  $X$ . When Assumptions 5 and 6 hold, the treatment assignment is said to be strongly ignorable (Rosenbaum and Rubin, 1983).

We start by discussing the identification strategy in section 4.2. In this case, as before, we need Assumptions 2 and 3. Assumptions 2 and 5 imply that  $Y(1, 1), Y(0, 0), Y(1, 0) \perp \{T, S(1), S(0)\} | X$ . Thus, the covariates correct for selection not only into the treatment, but also into the principal strata.<sup>42</sup> Finally, Assumption 3 allows using  $Y(1, S(1))$  to learn about  $Y(1, S(0))$ . Then, under Assumptions 2, 3, 5, and 6 we identify *NATE* as in (8). This estimator can be implemented using the same approach outlined in section 4.2.

Now consider estimating *NATE* with a strongly ignorable treatment assignment by focusing on the subpopulation for which  $T$  does not affect  $S$ . The main difference from the random assignment case is that now focus is on the subpopulation for which  $S(1) = S(0) = s$  and that also has the same values of  $X$ . Therefore, *LNATE* can be defined as in (9) including  $X$  in the conditioning set. If there is knowledge about a subpopulation for which  $T$  does not affect  $S$ , *LNATE* can be estimated as in the previous section (with the additional conditioning on  $X$ ).

<sup>40</sup>See for example Heckman, Lalonde and Smith (1999) and Imbens (2004).

<sup>41</sup>Similar to Assumption 1, this assumption implies  $Y(1, S(1)), Y(0, S(0)), Y(1, S(0)) \perp T | (X, S(1), S(0))$ .

<sup>42</sup>Assumptions 2 and 5 also imply (again using Lemma 4 in Dawid, 1979) that  $Y(1, 1), Y(0, 0), Y(1, 0) \perp \{S(1), S(0)\} | \{T, X\}$ , so that individuals in different strata but with the same values of the treatment and covariates are comparable.

If we do not have knowledge of such subpopulation (as is usually the case), we follow the same approach as in section 4.3 and focus on the subpopulation for which  $\widehat{S}_i(1) = \widehat{S}_i(0)$ , facing the same consequences discussed there. Specifically, estimation is focused on:

$$LNATE_{\widehat{S}} = E\{E[Y(1, S(0)) - Y(0, S(0)) | S(0) = s_0, S(1) = s_1, \widehat{S}(1) = \widehat{S}(0), X = x]\} \quad (13)$$

where the term inside the outer expectation is the local  $NATE$  for the strata with  $S(0) = s_0, S(1) = s_1, \widehat{S}(1) = \widehat{S}(0)$  and  $X = x$ . As before, (13) is a causal effect because it is an average over effects defined within a principal stratum. Note that in this case the covariates both remove selection into the treatment and predict  $S(1)$  and  $S(0)$ .

When the treatment is not randomly assigned we need to add an overlap assumption in order to ensure that, for sufficiently large samples, there will be both treated and control individuals at each value of  $X, S(0)$  and  $S(1)$  where we have  $\widehat{S}(1) = \widehat{S}(0)$ . Specifically:

**Assumption 7**  $0 < \Pr(T = 1 | S(0) = s_0, S(1) = s_1, X = x, \widehat{S}(1) = \widehat{S}(0)) < 1$ , for all  $s_0, s_1$  and  $x$ .

Assumption 7 is similar to the common overlap condition in Assumption 6, except that it includes  $S(1)$  and  $S(0)$  as additional covariates and only has to hold for the subpopulation of interest.<sup>43</sup>

Finally, under Assumption 4,  $NATE = LNATE_{\widehat{S}}$  and  $MATE = ATE - LNATE_{\widehat{S}}$ . The implementation of this approach follows closely the one outlined when  $T$  is randomly assigned. The only distinction is the inclusion of the covariates  $X$  in the conditioning set when estimating  $LNATE_{\widehat{S}}$  based on (10) or when estimating  $LATE^{sub}$ , both for the subpopulation with  $\widehat{S}_i(1) = \widehat{S}_i(0)$ .

## 5 Empirical Applications

In this section, we present two empirical applications that illustrate the implementation of our strategies to estimate the causal net average treatment effect ( $NATE$ ), from which the causal mechanism average treatment effect ( $MATE$ ) can also be estimated. The first application illustrates the case of a randomly assigned treatment using data from the social experiment undertaken in the National Job Corps Study (NJCS), while the second implements our estimators to observational data from the Natality Data Sets of Pennsylvania (1989-1991).

---

<sup>43</sup>In practice, since we expect  $S_i(0) \approx S_i(1)$  within the subpopulation with  $\widehat{S}_i(1) = \widehat{S}_i(0)$ , one way to check the overlap condition is to look at the overlap in the distribution of the propensity scores for treated and control units within this subpopulation, where the propensity score includes  $S^{obs}$  as an additional covariate.

## 5.1 Random Assignment

If receipt of treatment is randomly assigned, we can focus on the estimation of *NATE* and *MATE* without the additional complication of controlling for self-selection into the treatment. Our data comes from the National Job Corps Study (NJCS), a randomized experiment to evaluate the effectiveness and social value of the Job Corps (JC) training program. JC provides low-skilled young people (ages 16-24) with marketable skills to enhance their labor market outcomes by offering academic, vocational, and social skills training at JC centers throughout the United States, where most students reside during their enrollment period.<sup>44</sup>

An important finding of the NJCS was that, 16 quarters after randomization, individuals in the treatment group earned a statistically significant 12% more per week (\$25.2) than individuals in the control group (Burghardt et al., 2001). However, upon looking at different race and ethnic groups, it was found that Hispanics in the treatment group earned 10% less (a statistically insignificant -\$15.1) than those in the control group during the same period of time. In contrast, black and white treatment-group members experienced a statistically significant earnings increase of 14% (\$22.8) and 24% (\$46.2) over their control group members, respectively (Schochet, Burghardt and Glazerman, 2001).<sup>45</sup>

The bold differential impact on Hispanics was labeled the most prominent “failure” of JC and it could not be explained by individual and institutional variables (Burghardt et al., 2001). In a recent paper, Flores-Lagunes, Gonzalez and Neumann (2006) (hereafter FGN) document that Hispanics in the control group earned a significant amount of labor market experience during the study compared to treated Hispanics and also to control-group blacks and whites. They show, using several pieces of evidence, that such accumulated experience resulted in an earnings advantage that treated Hispanics were not able to overcome by the end of the study. Thus, this post-treatment variable (experience) potentially accounts for the lack of earnings gain of Hispanics in JC; nevertheless, the methods employed in FGN (2006) come up short of having clean causal interpretations. This setup offers a situation in which the estimation of our parameters are policy-relevant: if lost labor market experience (i.e. the lock-in effect) is a relevant causal mechanism through which JC fails to increase the earnings of Hispanics, policies that reduce the lock-in effect of JC on Hispanics can be judged beneficial. At the same time, by focusing on subgroups that seem to differ in terms of their lock-in effect, this application provides an interesting setting in which our parameters should result in distinct inferences for these groups.

Table 1 presents estimates of different parameters of interest using data from the NJCS.

---

<sup>44</sup>For more information on Job Corps and the NJCS see Burghardt et al. (2001).

<sup>45</sup>These estimated effects reported by the NJCS were computed using differences-in-means estimates adjusted for non-compliance, identifying a *LATE* on those who comply with their treatment assignment (Imbens and Angrist, 1994). The proportion of those in the treatment group who enrolled in Job Corps was 73%, and the proportion of those in the control group that managed to enroll in Job Corps was 1.4%.

Specifically, we employ the same sample used by FGN (2006) that includes individuals with information on pre-treatment covariates plus the post-treatment variable “average hours worked per week during the study”, and that report being Hispanic, white or black.<sup>46,47</sup> In contrast to FGN (2006), however, we pool the samples of blacks and whites for simplicity, since for both of these groups it is found that post-treatment experience is not a relevant mechanism (we present further evidence of this below), making unnecessary to present a separate analysis for them.<sup>48</sup>

Rows 1 through 4 in Table 1 report estimates of the average intention-to-treat (*ITT*) parameter.<sup>49</sup> Row 1 presents unadjusted differences in means between treatment- and control-group individuals. These estimates are qualitatively similar to the originally reported NJCS estimates: for the full sample and white/black samples the estimates are positive and statistically significant (\$15.6 and \$23.8, respectively), while for Hispanics the effects show a loss of \$19.7 that is marginally statistically significant.

Next, we present *ITT* estimates that control for pre-treatment variables through weighting by the estimated propensity score (pscore) in order to improve precision. In particular, we employ the estimator due to Robins and Rotnitzky (1995), as described in Imbens (2004), that combines weighting and regression and has the desirable property of being asymptotically efficient. Its implementation amounts to applying weighted least squares (WLS) to a regression of the outcome on the treatment indicator and optional additional covariates, with weights given by  $\lambda_i = \sqrt{\frac{T_i}{p(X_i)} + \frac{1-T_i}{1-p(X_i)}}$  where  $p(X_i)$  is the estimated pscore.<sup>50</sup> Rows 2, 3, and 4 differ in the way the WLS regression is specified: using no additional covariates, the pscore, and up to a cubic pscore as additional regressors, respectively.

The estimates in rows 2-4 are fairly comparable to the unadjusted estimates in row 1, except for Hispanics. This might be due to the smaller sample size of this group and the fact that the group shows some pre-treatment imbalances in the covariates.<sup>51</sup> For this reason the remaining estimates adjust for covariates. These estimates are in line with the conclusions in the NJCS although, as documented in FGN (2006), the negative effects on Hispanics are less dramatic

---

<sup>46</sup>The original NJCS sample contains 11,313 individuals. Out of these, 219 individuals do not complete the baseline interview, 1,295 more have missing values in any of the variables we control for, and 694 are individuals whose race or ethnicity is not white, black or Hispanic. The resulting sample contains 9,105 individuals. FGN (2006) report that this sample is consistent with the overall profile of the JC population.

<sup>47</sup>The pre-treatment variables include: indicators for a high school diploma or GED, speaks English as a native language, married or cohabitating, household head, one or more children, gender, vocational degree, ever been convicted, employed, unemployed, not in the labor force, resides in a PMSA, MSA, pre-treatment weekly earnings, age, and indicators for race and ethnicity.

<sup>48</sup>We have estimated all parameters for the black and white samples separately as well, corroborating that this is indeed the case.

<sup>49</sup>Given the presence of non-compliance in the sample, we estimate the intention-to-treat (*ITT*) parameter. This parameter is commonly estimated in the program evaluation literature and allows relying on the random assignment as much as possible. Consequently, in this application, our parameters decompose the *ITT* and not the *ATE*.

<sup>50</sup>The propensity score (pscore) is estimated using all pre-treatment variables, their squares, and interactions in a logit model.

<sup>51</sup>The misalignment of pre-treatment variables for Hispanics is documented and discussed in FGN (2006).



once covariates are controlled for. These “total effect” estimates will be the benchmark to compare our estimated effects net of the lock-in mechanism effect.

The next set of estimates in Table 1 are of Rosenbaum’s (1984) *NTD* parameter. All of them are obtained controlling for the observed value of post-treatment labor market experience employing the WLS approach described above with a pscore that includes a flexible specification of experience ( $S^{obs}$ ) in its estimation.<sup>52</sup> This way of controlling for a post-treatment variable by including it in the estimation of the propensity score is followed by Black and Smith (2004), although they use a matching approach to control for the estimated pscore, as opposed to weighting.

Recall that the *NTD* estimates typically lack causal interpretation as estimates of the total effect, and correspond to *NATE* under very stringent conditions (see section 4.1). We report them for comparison to our estimates below. The *NTD* estimates for the full sample are less than 20% larger compared to the *ITT* estimates, while for whites/blacks they are less than 10% larger. For Hispanics, the two sets of estimates are starkly different: more than 150% larger. In sum, despite the fact that all these effects are statistically insignificant for Hispanics and the lack of causal interpretation of *NTD*, the point estimates are suggestive of a relevant lock-in effect for Hispanics (contrary to whites/blacks) that would seem to explain an important portion of the lack of effects of JC on them.

Two sets of estimates of *NATE* appear in rows 8-9, obtained using the estimation strategy outlined in section 4.2. To implement it, we model (under Assumption 3) the first term in (8) as a linear function of  $S(1)$  and all available pre-treatment covariates. The second term in (8) is similarly predicted, but using  $S(0)$  instead. The two *NATE* estimates differ on the specification of the experience variable and the covariates included: row 8 includes experience up to a cubic term and all covariates, while row 9 adds interactions between the experience variable and the covariates to this specification. The *NATE* estimates for whites/blacks and for the full sample are closer to the *ITT* estimates than the *NTD* estimates, which is consistent with a non-existent lock-in effect for them. Among the two *NATE* estimates, the richer specification (row 9) is closer to the *ITT* estimates. For Hispanics, however, the *NATE* estimates are very different from the *ITT* estimates (as was the case with *NTD*), strengthening the notion of a relevant lock-in effect for them. Unfortunately, as before, the estimates are imprecisely estimated.

In the last panel of Table 1, we present estimates of  $LNATE_{\hat{S}}$  that differ in the way they are implemented. In all of them, the potential values of post-treatment labor market experience ( $S(0)$  and  $S(1)$ ) are estimated based on covariates  $X$  employing the matching approach described in footnote 33, using a single match on the estimated pscore that does not include experience.<sup>53</sup> Given that  $S$  is defined as the average number of hours worked per week during

---

<sup>52</sup>In particular, we use the same specification of the pscore as in *ITT*, but include experience up to a cubic term, and interactions of this variable with the pre-treatment covariates.

<sup>53</sup>We estimated  $S(1)$  and  $S(0)$  based on  $X$  separately for the full sample, whites/blacks, and Hispanics.

the study, it is difficult to find individuals for which  $\widehat{S}(1) = \widehat{S}(0)$ . We approach this feature by defining a window around  $\widehat{S}(1) - \widehat{S}(0) = 0$  using a Silverman-type bandwidth to characterize the subpopulation with  $\{\widehat{S}(1) = \widehat{S}(0)\}$ .<sup>54</sup> The proportional size of the resulting  $LNATE_{\widehat{S}}$  subpopulation is similar across the three samples.

As mentioned in section 4.3, we can assess the plausibility that the subpopulation found for the estimation of  $LNATE_{\widehat{S}}$  satisfies the requirement that experience ( $S$ ) is not affected by the treatment. To do this, we implement several versions of the Fisher randomization test. This test provides evidence on the sharp null hypothesis  $H_0 : S_i(1) = S_i(0)$  for every  $i$ . In its simplest form, the implementation consists of simulating the distribution under  $H_0$  for the observed test statistic  $\frac{\sum T_i S_i(1)}{\sum T_i} - \frac{\sum (1-T_i) S_i(0)}{\sum (1-T_i)}$  (or other quantity measuring the effect of  $T$  on  $S$ ) by randomizing the treatment indicator to the units in the sample and computing the test statistic in each repetition. Then, an approximate p-value is constructed by comparing the observed test statistic with the simulated distribution. A rejection of the test indicates that the post-treatment variable  $S$  is affected by the treatment.

Panel A of Table 3 presents results of Fisher randomization tests for the three groups under analysis. For each group, tests are applied to the population of the group (for comparison) and to the subpopulation characterized by  $\{\widehat{S}(1) = \widehat{S}(0)\}$ . We present five versions of the test, which turn out to yield the same conclusion. The first three tests are based on comparing the coefficient on an OLS regression of  $S$  on the treatment indicator ( $T$ ) in row 1, adding the pscore in row 2, and further adding the square and cube of the pscore in row 3.<sup>55</sup> The last two rows are based on applying the simple form of the test described above to the residuals from OLS regressions of  $S$  on the pscore and then adding the pscore square and cube, respectively.

The results of the tests are in line with what would be expected. Whites/blacks, the group that shows the least amount of lock-in effect, have p-values for their population that range from 0.31 to 0.5, not rejecting the null hypothesis of no effect of the treatment on post-treatment experience. For the full sample, the null hypothesis is rejected at the 10% level in all cases, and at the 5% level in one. However, the subpopulation cannot reject the test with a p-value of 0.77 or higher. Finally, as expected, Hispanics show the strongest rejections of the null hypothesis for their population, consistent with the documented relevance of the lock-in effect for them. Importantly, the characterized subsample substantially decreases the strong relationship between  $S$  and  $T$ , with p-values that range from 0.57 to 0.99. In sum, for the three groups under analysis, the statistical evidence cannot reject the notion that the corresponding subpopulations characterized by  $\{\widehat{S}(1) = \widehat{S}(0)\}$  have a zero effect of  $T$  on  $S$  for all units.<sup>56</sup>

<sup>54</sup>The Silverman-type bandwidth employed is equal to  $0.79 * IQR * N^{-1/5}$ , where  $IQR$  is the interquartile range and  $N$  is the sample size. This bandwidth has the advantage of being more robust to outliers than the usual one based on the standard deviation (see, e.g., Pagan and Ullah, 1999).

<sup>55</sup>Given that in this case  $S$  is the outcome of interest, in all cases the pscore is the specification that does not include experience on it. Importantly, we re-estimate the pscore within the corresponding subpopulation.

<sup>56</sup>Note that the Fisher randomization test could also be used at the outset of the estimation of mechanism and

We now discuss the  $LNATE_{\hat{S}}$  estimates based on the subpopulations described above. We start by estimating  $LNATE_{\hat{S}}$  in rows 10-12 as the local average treatment effect ( $LATE_{\hat{S}}^{sub}$ ) for this subpopulation. For each group, we estimate a propensity score without including experience in this subpopulation and use WLS with similar specifications as those used in the estimation of  $ITT$  and  $NTD$ . For the full sample, the  $LNATE_{\hat{S}}$  estimates average \$29.9, substantially higher than  $ITT$  (\$19.2),  $NTD$  (\$23), and also  $NATE$  (\$21.7). In contrast, for whites/blacks the  $LNATE_{\hat{S}}$  estimates are around \$24, which is about the same magnitude as  $ITT$ . This result reinforces the observation that for whites/blacks experience is not a mechanism through which JC affects wages. For Hispanics, however, the  $LNATE_{\hat{S}}$  estimates (about \$10.9 on average) are larger than the  $NTD$  and  $NATE$ , and substantially larger than the  $ITT$ . Unfortunately, it is also estimated very imprecisely and these differences are not statistically significant. Note that if we were to assume that the average mechanism effect is negative for all three subpopulations, then the estimates in rows 10-12 would be downward biased (see (12)).

As a robustness check, and following our discussion in section 4.3, rows 13-15 present estimates of  $LNATE_{\hat{S}}$  controlling for  $S^{obs}$  based on equation (10). These estimates differ from those in rows 10-12 in that they include experience in the specification of the pscore.<sup>57</sup> Overall, the estimates for the full population are close to the ones presented in rows 10-12, which give us confidence in our  $LNATE_{\hat{S}}$  results for the full sample. For whites and blacks the results are not as robust as for the full sample; however, they remain below the  $NTD$  estimates in rows 5-7. Note that for this group there is a considerable decrease in the precision of our estimates by introducing experience as an additional control, since now the  $LNATE_{\hat{S}}$  estimates in rows 13-15 are not statistically different from zero. Something similar occurs for Hispanics, for which the  $LNATE_{\hat{S}}$  estimates fall when including experience in the pscore specification, although none of their estimates are statistically significant.

We gather the following conclusions from this empirical illustration. First, our estimates of  $NATE$  and  $LNATE_{\hat{S}}$  suggest that the lock-in effect results in a negative causal mechanism for the effect of JC training on Hispanic’s earnings, although the estimates remain statistically insignificant. Second, the full set of estimates corroborate the high degree of heterogeneity that exist among whites/blacks and Hispanics, which results in very different inferences in terms of their estimated total, net and mechanism effects from JC training. Lastly and unfortunately, this application appears less than ideal since many of the differences in the estimates are not statistically significant. This may be due to lack of differences among the true parameters in this application, or the need for larger sample sizes to increase precision. Fortunately, in the next empirical application, we eliminate the sample size as a source of concern.

---

net causal effects to gauge the importance of a potential mechanism of a treatment in the population.

<sup>57</sup>Given that we are within the subpopulation with  $\hat{S}_i(1) = \hat{S}_i(0)$ , for which  $S(1) \approx S(0)$ , we can include  $S^{obs}$  in the estimation of the propensity score.

## 5.2 Non-random Assignment

When the treatment is not randomly assigned we face the additional issue of controlling for self-selection. To deal with this, we employ the selection on observables assumption and regard the treatment as randomly assigned conditional on a rich set of observed covariates. For this application, the data comes from Pennsylvania’s Natality Data Sets from 1989 to 1991, which includes all births (although we focus on single births) and has been previously used and documented by Chay, Flores and Torelli (2005). The availability of a wide range of observable characteristics, including characteristics of both parents and previous birth history, makes the assumption of selection on observables more plausible.

Our focus is on evaluating the extent to which smoking during pregnancy (treatment) affects the incidence on low birth weight (outcome) through a shorter gestation time (a mechanism). The outcome “low birth weight” (LBW) has the standard definition in the medical literature of birth weight below 2,500 grams, and is widely associated with a myriad of health, behavioral and socioeconomic problems in later stages of individual development (e.g. UNICEF and WHO, 2004). For instance, LBW has been negatively associated to educational attainment, self-reported health status and employment (Currie and Hyson, 1999). The consensus in the literature (e.g., Stein et al., 1983; Center for Disease Control and Prevention, 2001) is that smoking during pregnancy causally reduces birth weight and thus increases the probability of incidence on LBW, but the importance of specific mechanisms is not completely understood. In general, there might be two ways in which smoking during pregnancy affects birth weight: a shorter gestation time and intrauterine growth retardation (IGR). The importance of determining the causal relative importance of a channel is that particular policies aimed at minimizing the negative effects of smoking during pregnancy may be considered. For instance, if gestation time is an important causal mechanism, drugs that lengthen gestation time may be deemed useful.

Table 2 presents the results for this application. Given the importance of satisfying the support condition in observational studies using the selection-on-observables assumption (e.g. Heckman, Ichimura and Todd, 1997; Dehejia and Wahba, 1999), we concentrate on a sample in the overlap region of the estimated p-score between the 1 percentile of the p-score values for the treated and 99 percentile of the p-score values for controls.<sup>58</sup> For reference, the average LBW incidence in the sample employed is 58.3 per 1,000 births, and 20.8% of women smoked during pregnancy. The incidence on LBW is 48.4 per 1,000 births for non-smokers (control) and 95.7 per 1,000 births for smokers (treatment), yielding the unadjusted difference of 47.3 (per 1,000 births) shown in the first row of Table 2. Thus, mothers who smoked during pregnancy were about twice as likely to deliver a LBW baby than those mothers who did not.

Rows 2 through 4 present estimates of the total effect (*ATE*) of smoking on LBW incidence,

---

<sup>58</sup>The sample consists of 496,212 individuals, of which 425,219 are contained within the overlap region.

controlling for self-selection using an estimated pscore.<sup>59</sup> Rows 2-4 employ the same WLS approach and specifications as in the previous empirical application. The three estimates are close to each other, reflecting an effect of smoking on the incidence on LBW of 33 per 1,000 births. The fact that this figure is smaller than the unadjusted difference in row 1 implies a selection bias of about 14 in the unadjusted figure. Still, the  $ATE$  estimate suggests a sizable effect of smoking during pregnancy on LBW incidence, as the probability increases 68%.

The second panel of Table 2 (rows 5-7) presents estimates of the  $NTD$  parameter that control directly for gestation time in weeks ( $S^{obs}$ ) using specifications similar to those used in the previous application. For these estimates the pscore includes observed gestation (in a flexible way) in its estimation. The  $NTD$  is precisely estimated at about 28. This suggests that 15.2% of the effect of smoking during pregnancy on the incidence on LBW (5 of 33 per 1,000 births) can potentially be attributed to gestation time, although these estimates correspond to  $NATE$  under very stringent conditions.

Rows 8 and 9 present estimates of  $NATE$ , implemented in the same way as in the previous application. Both estimates are essentially identical to each other at 26.5 per 1,000 births, and slightly smaller than  $NTD$ . Based on the  $NATE$  estimates and under the assumptions discussed in section 4, about 20% of the effect of smoking during pregnancy on the incidence on LBW can be causally attributed to gestation time. We note that the difference between the  $NATE$  and the  $ATE$  estimates is statistically significant, whereas the difference between the  $NATE$  and the  $NTD$  estimates is not.

Finally, the last panel in Table 2 presents estimates of  $LNATE_{\hat{S}}$ . The subpopulation of interest is obtained, as in the previous application, by estimating the potential values of gestation time using a single match on the estimated pscore and selecting those units with  $\{\hat{S}(1) = \hat{S}(0)\}$ , resulting in a sample of about 15% of the one used for estimation of the  $ATE$ .<sup>60</sup> We test whether the individual treatment effect of  $T$  on  $S$  is zero for all mothers in this subpopulation using different versions of the Fisher randomization test. Panel B of Table 3 shows that for the population the null hypothesis of no effect of  $T$  on  $S$  is soundly rejected, while in the subpopulation it is not, with p-values ranging from 0.35 to 0.92. Rows 10-12 estimate  $LNATE_{\hat{S}}$  by estimating

---

<sup>59</sup>The propensity score is estimated with a logit model. The covariates used are mother's age, education, race, ethnicity, marital status, foreign-born status; father's age, education, race and ethnicity; dummies for trimester of first prenatal care visits, adequacy of care, number of prenatal visits, number of drinks per week, alcohol use, live birth order, number of previous births were newborn died, parity indicator, interval since last birth, indicators for previous birth over 4000 grams and previous birth preterm or small for gestational age; maternal medical risk factors that are not believed to be affected by smoking during pregnancy: anemia, cardiac disease, lung disease, diabetes, genital herpes, hydramnios/oligohydramnios, hemoglobinopathy, chronic hypertension, eclampsia, incompetent cervix, renal disease, Rh sensitization, uterine bleeding; indicators for: month of birth, county of residence at birth, state of occurrence and residence different, each variable that is missing for some mothers. The particular specification used includes nonlinear functions and interactions, and is similar to the one used in Chay, Flores and Torelli (2005) and Almond, Chay and Lee (2005).

<sup>60</sup>Note that, contrary to the previous application, the post-treatment variable gestation time, measured in weeks, is sufficiently discrete to allow identifying a population for which  $\{\hat{S}(1) = \hat{S}(0)\}$  exactly.

the  $ATE$  within the relevant subpopulation ( $LATE_{\hat{S}}^{sub}$ ) employing a pscore without gestation in its specification; whereas rows 13-15 estimate  $LNATE_{\hat{S}}$  by including gestation in the pscore specification. For this application, all six  $LNATE_{\hat{S}}$  estimates are essentially the same at about 22.8 per 1,000 births, an amount smaller than both the  $NTD$  and  $NATE$  estimates. Under Assumption 4 (constant  $NATE$ ), these estimates imply that 31% of the effect of smoking during pregnancy on the incidence on LBW can be causally attributed to gestation time.<sup>61</sup> Note that in this case whether we include gestation or not in the estimation of the pscore does not affect the results. This supports the notion that within this subpopulation  $LNATE_{\hat{S}} \approx LATE_{\hat{S}}^{sub}$  and thus the bias in (12) is close to zero. Remarkably, the difference between  $ATE$  and each of  $NTD$ ,  $NATE$  and  $LNATE_{\hat{S}}$  is highly statistically significant, speaking to the relevance of the mechanism. In addition, the difference between  $LNATE_{\hat{S}}$  and  $NTD$  is statistically significant at the 7% level, while the difference between  $LNATE_{\hat{S}}$  and  $NATE$  is not.

In sum, the results of this empirical application are consistent with a causal role of gestation time as a channel through which smoking during pregnancy increases the incidence on LBW. While the total effect is 33 per 1,000 births (or about 70% higher than non-smokers), our results indicate that between 20 to 30 percent of this effect works causally through a shorter gestation time. Importantly, we also find that the  $NTD$  understates the importance of gestation time as a causal mechanism by between 5 to 15 percentage points. Clearly, an advantage of this empirical application is that the sample sizes allow us to estimate our parameters of interest very precisely.

Summarizing the implementation of our methods in the two empirical applications, some patterns emerge. First, our methods are feasible to implement in estimating the causal average net effect and the causal average mechanism effect. Second, the estimation of the parameters introduced in this paper yields new insights about the treatment at hand; although we remark that carefully evaluating the plausibility of the assumptions made by each estimation strategy in particular applications cannot be overstated. Finally, in both applications, the (non-causal) estimates for  $NTD$ , which is the parameter commonly used in the literature to estimate net effects, are different from our causal estimates (especially in the second empirical application). This underscores the potentially misleading conclusions that can be reached by controlling for the observed values of the post-treatment variable.

---

<sup>61</sup>If we further assume that the average mechanism effect is positive for the subpopulation where  $LNATE_{\hat{S}}$  is estimated, then the estimates in rows 10-12 are upward biased (see (12)) and thus the estimated mechanism effect of 31% is downward biased.

## 6 Conclusion

This paper analyzes the identification and estimation of an average causal mechanism through which a treatment or intervention affects an outcome, and the average causal effect of the treatment net of this mechanism. These causal effects are of interest since they allow a better understanding of the treatment and, as a result, can be used for policy purposes in the design, development, and evaluation of interventions. Not surprisingly, it is common in the literature to informally analyze potential mechanisms of a treatment as a natural step after estimating the “total” effect of the treatment. Unfortunately, these analyses are typically based on a standard approach that controls for observed values of a variable representing a mechanism, resulting in estimates that generally lack causal interpretation.

We avoid this pitfall by using the concept of principal stratification (Frangakis and Rubin, 2002) to define causal parameters of interest. These parameters intuitively decompose the total effect of a treatment into the part that is causally due to a particular mechanism (mechanism average treatment effect or *MATE*) and the part that is net of such mechanism (net average treatment effect or *NATE*). In addition, we show that for interpreting the standard approach that controls for observed values of the potential mechanism as *NATE* we need to rely on assumptions that are typically too strong to be useful in practice.

We develop two strategies for estimation of our parameters, both for the case of a randomly assigned treatment and the case of non-random assignment. The first strategy is based on an assumption in the spirit of the familiar selection on observables approach (Rosenbaum and Rubin, 1983; Imbens, 2004), along with an assumption relating the partially observable potential outcome  $Y(1, S(1))$  to the unobserved  $Y(1, S(0))$  that is necessary for the estimation of *NATE*. The second approach estimates *NATE* by estimating a local *NATE* for the subpopulation for which  $T$  does not affect  $S$ , where the covariates are used to find such population. We recognize that our assumptions, although weaker than the implicit assumptions needed to interpret the standard approach as an estimate of *NATE*, remain strong. This is because estimation of causal net and mechanism effects is a difficult task given the data typically available to a researcher. Therefore, it is important in applied work to explicitly discuss the plausibility of the assumptions we present. Finally, we present two different empirical applications that illustrate the implementation of our methods.

Several natural extensions are left for future work. For instance, it is of interest to develop a set of alternative assumptions that lead to an identification and estimation strategy that allows for selection into the treatment based on unobservables. A couple of possibilities come to mind, such as the construction of bounds for our parameters in the spirit of Manski (1990). Similarly, an analysis of the way in which additional information can be used to estimate our parameters is of interest, such as the availability of instrumental variables.

## References

- [1] Adams, P., Hurd, M., McFadden, D., Merrill, A. and Ribeiro, T. (2003) "Healthy, Wealthy and Wise? Tests for Direct Causal Paths Between Health and Socioeconomic Status" *Journal of Econometrics*, 112, 3-56.
- [2] Almond, D., Chay, K. Y. and Lee, D. S. (2005) "The Cost of Low Birth Weight". *Quarterly Journal of Economics*, 120 (3), 1031-1083.
- [3] Black, D. and Smith, J. (2004), "How Robust is the Evidence on the Effects of College Quality? Evidence from Matching." *Journal of Econometrics*, 121, 99-124.
- [4] Burghardt, J., Schochet, P., McConnell, S., Johnson, T., Gritz, R., et. al. (2001) "Does Job Corps Work? Summary of the National Job Corps Study" 8140-530. Mathematica Policy Research, Inc., Princeton, NJ.
- [5] Card, D. and Hyslop, D. (2005) "Estimating the Effects of a Time-Limited Earnings Subsidy for Welfare-Leavers" *Econometrica*, 73, 1723-70.
- [6] Center for Disease Control and Prevention (2001) *Women and Smoking: A Report of the Surgeon General*.
- [7] Chay, K.; Flores, C. A. and Torelli, P. (2005) "The Association between Maternal Smoking during Pregnancy and Fetal and Infant health: New Evidence from United States Birth Records", mimeo, University of California, Berkeley.
- [8] Currie, J. and Hyson, R. (1999) "Is the Impact of Health Shocks Cushioned by Socioeconomic Status? The Case of Low Birthweight " *American Economic Review*, 89, 245-250.
- [9] Currie, J. and Neidell, M. (2007) "Getting Inside the 'Black Box' of Head Start Quality: What Matters and What Doesn't" *Economics of Education Review*, 26, 83-99.
- [10] Dawid, A. (1979) "Conditional Independence in Statistical Theory (with Discussion)" *Journal of the Royal Statistical Society, Series B*, 41, 1-31.
- [11] Dearden, L. Ferri, J. and Meguir, C. (2002), "The Effect of School Quality on Educational Attainment and Wages." *Review of Economics and Statistics*, 84, 1-20.
- [12] Dehejia, R. and Wahba, S. (1999), "Causal Effects in Nonexperimental Studies: Reevaluating the Evaluation of Training Programs", *Journal of the American Statistical Association*, 94, 1053-1062.
- [13] Ehrenberg, R., Jakubson, G., Groen, J., So, E., and Price, J. (2006), "Inside the Black Box of Doctoral Education: What Program Characteristics Influence Doctoral Students' Attrition and Graduation Probabilities?" NBER Working Paper 12065.
- [14] Fisher, R.A. (1935). *The Design of Experiments*. Edingburgh: Oliver and Boyd.
- [15] Flores-Lagunes, A., Gonzalez, A., and Neumann, T. (2006) "Learning But Not Earning? The Impact of Job Corps Training on Hispanic Youths", mimeo, University of Arizona.
- [16] Frangakis, C.E. and Rubin D. (2002) "Principal Stratification in Causal Inference" *Biometrics*, 58, 21-29.
- [17] Heckman, J.; Ichimura, H. and Todd, P. (1997), "Matching as an Econometric Evaluation Estimator: Evidence from Evaluating a Job Training Programme", *Review of Economic Studies*, 64(4), 605-654.
- [18] Heckman, J.; Smith, J. and Clements, N. (1997), "Making the Most Out of Programme Evaluations and Social Experiments: Accounting for Heterogeneity in Programme Impacts", *Review of Economic Studies*, 64(3), 487-535.



- [19] Heckman, J., LaLonde, R. and Smith, J. (1999) "The Economics and Econometrics of Active Labor Market Programs" in O. Ashenfelter and D. Card (eds.) *Handbook of Labor Economics*. Elsevier Science North Holland, 1865-2097.
- [20] Holland, P. (1986) "Statistics and Causal Inference" *Journal of the American Statistical Association*, 81, 945-70.
- [21] Imbens, G. (2004) "Nonparametric Estimation of Average Treatment Effects under Exogeneity: A Review" *Review of Economics and Statistics*, 84, 4-29.
- [22] Imbens, G. and Angrist, J. (1994) "Identification and Estimation of Local Average Treatment Effects" *Econometrica*, 62, 467-75.
- [23] Lechner, M. (2005) "A Note on Endogenous Control Variables in Evaluation Studies" Discussion paper 2005-16, University of St. Gallen.
- [24] Lechner, M. and R. Miquel (2005) "Identification of the Effects of Dynamic Treatments by Sequential Conditional Independence Assumptions" Discussion paper, University of St. Gallen.
- [25] Manski, C. (1990) "Nonparametric Bounds on Treatment Effects" *American Economic Review Papers and Proceedings*, 80, 319-23.
- [26] Meyer, B. (1995) "Lessons from the U.S. Unemployment Insurance Experiments" *Journal of Economic Literature*, XXXIII, 91-131.
- [27] Mealli, F. and Rubin, D. (2003) "Assumptions Allowing the Estimation of Direct Causal Effects" *Journal of Econometrics*, 112, 79-87.
- [28] Neyman, J. (1923) "On the Application of Probability Theory to Agricultural Experiments: Essays on Principles" Translated in *Statistical Science*, 5, 465-80.
- [29] Pagan, A. and Ullah A. (1999) *Nonparametric Econometrics*. Cambridge university Press.
- [30] Petersen, M., Sinisi, S., and van der Laan, M. (2006) "Estimation of Direct Causal Effects" *Epidemiology*, 17, 276-284.
- [31] Robins, J. (1986) "A New Approach to Causal Inference in Mortality Studies with Sustained Exposure Periods—Application to Control of the Healthy Worker Survivor Effect" *Mathematical Modeling*, 7, 1393-1512.
- [32] Robins, J. and Greenland, S. (1992) "Identifiability and Exchangeability for Direct and Indirect Effects" *Epidemiology*, 3, 143-155.
- [33] Robins, J. and Rotnitzky, A. (1995) "Semiparametric Efficiency in Multivariate Regression Models with Missing Data" *Journal of the American Statistical Association*, 90, 122-129.
- [34] Rosenbaum, P. (1984) "The Consequences of Adjustment for a Concomitant Variable That Has Been Affected by the Treatment" *Journal of the Royal Statistical Society, Series A*, 147, 656-66.
- [35] Rosenbaum, P. and Rubin, D. (1983). "The Central Role of the Propensity Score in Observational Studies for Causal Effects", *Biometrika*, 70, 41-55.
- [36] Rubin, D. (1974) "Estimating Causal Effects of Treatments in Randomized and Nonrandomized Studies" *Journal of Educational Psychology*, 66, 688-701.
- [37] Rubin, D. (1980) "Discussion of 'Randomization Analysis of Experimental Data in the Fisher Randomization Test' by Basu" *Journal of the American Statistical Association*, 75, 591-93.

- [38] Rubin, D. (2004) "Direct and Indirect Causal Effects via Potential Outcomes" *Scandinavian Journal of Statistics*, 31, 161-70.
- [39] Rubin, D. (2005) "Causal Inference Using Potential Outcomes: Design, Modeling, Decisions" *Journal of the American Statistical Association*, 100, 322-331.
- [40] Stein et al. (1983). "Smoking, Alcohol and Reproduction" *American Journal of Public Health*, 73 (10), 1154-1156.
- [41] Schochet, P., Burghardt, J. and Glazerman, S. (2001) "National Job Corps Study: The Impacts of Job Corps on Participants' Employment and Related Outcomes." 8140-530. Mathematica Policy Research, Inc., Princeton, NJ.
- [42] Simonsen, M. and Skipper, L. (2006), "The Costs of Motherhood: An Analysis Using Matching Estimators", *Journal of Applied Econometrics*, 21, 919-934.
- [43] UNICEF and WHO (2004), *Low Birthweight: Country, Regional and Global Estimates*, New York.

**Table 1. Random Assignment Application: Estimation of the effect of the Job Corps training program on weekly earnings during quarter 16 after randomization. Mechanism analyzed: post-treatment labor market experience<sup>1</sup>**

	Full Sample			White and Black			Hispanic			
	(N=9,105)	t-stat	N	(N=7,412)	t-stat	N	(N=1,693)	t-stat	N	
<i>Estimation of Intention to Treat Effects, ITT</i>										
1	Unadjusted Difference	15.6	(3.40)	23.8	(4.70)		-19.7	(-1.79)		
2	WLS using pscore in weights. No pscore in regression	19.2	(4.21)	24.8	(4.90)		-7.6	(-0.69)		
3	WLS using pscore in weights. Linear pscore as regressor	19.3	(4.24)	24.8	(4.90)		-7.5	(-0.69)		
4	WLS using pscore in weights. Up to cubic pscore as regressor	19.2	(4.23)	25.0	(4.96)		-7.3	(-0.67)		
<i>Estimation of "Net Treatment Difference" controlling for observed post-treatment experience, NTD</i>										
5	WLS using pscore in weights. No pscore in regression	23.1	(5.14)	27.0	(5.42)		5.2	(0.47)		
6	WLS using pscore in weights. Linear pscore as regressor	23.1	(5.21)	27.0	(5.43)		5.8	(0.55)		
7	WLS using pscore in weights. Up to cubic pscore as regressor	22.7	(5.13)	27.0	(5.46)		4.5	(0.44)		
<i>Estimation of NATE using <math>E[Y(I,S(I))   S(I), X]</math> to predict <math>E[Y(I,S(0))   S(0), X]</math>. <math>E[Y(0,0)   S(0), X]</math> is similarly predicted.</i>										
8	OLS outcome on experience and its polynomials up to degree 3 plus linear covariates	22.7	(6.33)	26.8	(6.68)		3.6	(0.71)		
9	OLS outcome on experience and its polynomials up to degree 3, linear covariates, and interactions	20.7	(5.29)	23.6	(5.55)		6.7	(0.69)		
<i>Estimation of the local NATE for the subpopulation with (predicted) <math>S(0)=S(1)</math>, where predicted <math>S(0)</math> and <math>S(1)</math> are based on matching on the pscore.<sup>2</sup></i>										
<i>Using pscore that does not include experience in its estimation (Estimates of Local ATE)</i>										
10	WLS using pscore in weights. No pscore in regression	28.9	(2.45)	1273	23.9	(2.04)	1072	10.5	(0.36)	344
11	WLS using pscore in weights. Linear pscore as regressor	28.9	(2.49)	1273	24.0	(2.05)	1072	9.6	(0.34)	344
12	WLS using pscore in weights. Up to cubic pscore as regressor	31.9	(3.08)	1273	24.1	(2.07)	1072	12.7	(0.44)	344
<i>Using pscore that includes experience in its estimation</i>										
13	WLS using pscore in weights. No pscore in regression	30.1	(2.48)	1273	18.2	(1.52)	1072	1.6	(0.05)	344
14	WLS using pscore in weights. Linear pscore as regressor	30.1	(2.52)	1273	18.2	(1.52)	1072	2.9	(0.09)	344
15	WLS using pscore in weights. Up to cubic pscore as regressor	34.4	(3.28)	1273	18.0	(1.49)	1072	3.6	(0.12)	344

<sup>1</sup> All estimates use a sample that contains those who completed both a 48-month and baseline interviews, and with non-missing information on the covariates employed by the estimators. The sample sizes are indicated at the top of each column, unless otherwise indicated for particular estimators. Standard errors do not take into account the estimation of the propensity score.

<sup>2</sup> Given that S is defined as the average number of hours worked during the study, the predicted values S(1) and S(0) are continuous. The subpopulation with predicted S(1)=S(0) is obtained employing a window around (predicted) S(1)-S(0)=0 using a Silverman-type bandwidth based on the inter-quantile range (IQR):  $h=0.79*IQR*N^{(-1/5)}$ .

**Table 2. Non-Random Assignment Application: Estimation of the effect of smoking during pregnancy on the incidence of low birth weight (less than 2,500 grams) per 1,000 births. Mechanism analyzed: weeks of gestation (single births in Pennsylvania from 1989 to 1991)**

	<i>Estimate</i>	<i>t-statistic</i>
<i>Estimation of Average Treatment Effects, ATE. Focus on a population with overlap region of pscore between the 1 percentile of pscore for treated and 99 percentile of pscore for controls (N=425,219).</i>		
1 Unadjusted Difference	47.3	(44.82)
2 WLS using pscore in weights. No pscore in regression	33.1	(26.85)
3 WLS using pscore in weights. Linear pscore as regressor	32.8	(26.57)
4 WLS using pscore in weights. Up to cubic pscore as regressor	32.8	(26.56)
<i>Estimation of "Net Treatment Difference" controlling for observed gestation, NTD. Focus on a population with and overlap region of corresponding pscore between the 1 percentile of pscore for treated and 99 percentile of pscore for controls (N=424,677).</i>		
5 WLS using pscore in weights. No pscore in regression	28.2	(23.56)
6 WLS using pscore in weights. Linear pscore as regressor	27.9	(23.29)
7 WLS using pscore in weights. Up to cubic pscore as regressor	27.9	(23.28)
<i>Estimation of NATE using <math>E[Y(1,S(1))   S(1), X]</math> to predict <math>E[Y(1,S(0))   S(0), X]</math>. <math>E[Y(0,0)   S(0), X]</math> is similarly predicted. Focus on subpopulation with overlap region of pscore between the 1 percentile of pscore for treated and 99 percentile of pscore for controls. (N=425,219)</i>		
8 OLS outcome on gestation and its polynomials up to degree 3 plus covariates	26.6	(23.46)
9 OLS outcome on gestation and its polynomials up to degree 3, covariates, and interactions	26.5	(23.36)
<i>Estimation of the local NATE for the subpopulation with (predicted) <math>S(0)=S(1)</math> and overlap region of pscore between the 1 percentile of pscore for treated and 99 percentile of pscore for controls. Predicted values of <math>S(0)</math> and <math>S(1)</math> are based on matching on the pscore.</i>		
<i>Using pscore that does not include experience in its estimation (Estimates of Local ATE). N=63,666</i>		
10 WLS using pscore in weights. No pscore in regression	22.9	(9.51)
11 WLS using pscore in weights. Linear pscore as regressor	22.9	(9.47)
12 WLS using pscore in weights. Up to cubic pscore as regressor	22.9	(9.52)
<i>Using pscore that includes experience in its estimation. N=63,748</i>		
13 WLS using pscore in weights. No pscore in regression	22.8	(9.50)
14 WLS using pscore in weights. Linear pscore as regressor	22.7	(9.45)
15 WLS using pscore in weights. Up to cubic pscore as regressor	22.8	(9.49)

Note: The standard errors do not take into account the estimation of the propensity score.

**Table 3. Simulated p-values for Fisher's Randomization Test for the presence of individual effects of T on the post-treatment variable S. Based on 10,000 repetitions**

	PANEL A						PANEL B	
	<i>Random Assignment Application: testing the effect of Job Corps training on post-treatment experience</i>						<i>Non-Random Assignment Application: testing the effect of smoking on gestation</i>	
	<i>Full Sample</i>		<i>Hispanics</i>		<i>Whites and Blacks</i>		<i>Population</i>	<i>Subpopulation</i>
<i>Test based on:</i>	<i>Population</i>	<i>Subpopulation</i>	<i>Population</i>	<i>Subpopulation</i>	<i>Population</i>	<i>Subpopulation</i>	<i>Population</i>	<i>Subpopulation</i>
1 OLS coefficient, no covariates	0.01	0.96	0.00	0.57	0.31	0.25	--	--
2 OLS coefficient, including pscore	0.06	0.78	0.00	0.99	0.38	0.90	0.00	0.35
3 OLS coefficient, including up to cubic pscore	0.06	0.77	0.01	0.94	0.39	0.91	0.00	0.45
4 OLS residuals, including pscore	0.06	0.79	0.01	0.99	0.50	0.84	0.00	0.77
5 OLS residuals, including up to pscore	0.06	0.79	0.01	0.95	0.39	0.91	0.00	0.92

Note: "Subpopulation" refers to that subpopulation for which (predicted)  $S(0)=S(1)$ , which is used in the estimation of the local NATE.