

## NOTES AND COMMENTS

### MATCHING ON THE ESTIMATED PROPENSITY SCORE

BY ALBERTO ABADIE AND GUIDO W. IMBENS<sup>1</sup>

Propensity score matching estimators (Rosenbaum and Rubin (1983)) are widely used in evaluation research to estimate average treatment effects. In this article, we derive the large sample distribution of propensity score matching estimators. Our derivations take into account that the propensity score is itself estimated in a first step, prior to matching. We prove that first step estimation of the propensity score affects the large sample distribution of propensity score matching estimators, and derive adjustments to the large sample variances of propensity score matching estimators of the average treatment effect (ATE) and the average treatment effect on the treated (ATET). The adjustment for the ATE estimator is negative (or zero in some special cases), implying that matching on the estimated propensity score is more efficient than matching on the true propensity score in large samples. However, for the ATET estimator, the sign of the adjustment term depends on the data generating process, and ignoring the estimation error in the propensity score may lead to confidence intervals that are either too large or too small.

KEYWORDS: Matching estimators, propensity score matching, average treatment effects, causal inference, program evaluation.

#### 1. INTRODUCTION

PROPENSITY SCORE MATCHING ESTIMATORS (Rosenbaum and Rubin (1983)) are widely used to estimate treatment effects.<sup>2</sup> Rosenbaum and Rubin (1983) defined the propensity score as the conditional probability of assignment to a treatment given a vector of covariates. Suppose that adjusting for a set of covariates is sufficient to eliminate confounding. The key insight of Rosenbaum and Rubin (1983) is that adjusting only for the propensity score is also sufficient to eliminate confounding. Relative to matching directly on the covariates, propensity score matching has the advantage of reducing the dimensionality of matching to a single dimension. This greatly facilitates the matching process

<sup>1</sup>We are grateful to the editor and three referees for helpful comments, to Ben Hansen, Judith Lok, James Robins, Paul Rosenbaum, Donald Rubin, and participants in many seminars for comments and discussions, and to Jann Spiess for expert research assistance. Financial support by the NSF through Grants SES 0820361 and SES 0961707 is gratefully acknowledged.

<sup>2</sup>Following the terminology in Abadie and Imbens (2006), the term “matching estimator” is reserved in this article to estimators that match each unit (or each unit of some sample subset, e.g., the treated) to a small number of units with similar characteristics in the opposite treatment arm. Thus, our discussion does not refer to regression imputation methods, like the kernel matching method of Heckman, Ichimura, and Todd (1998), which use a large number of matches per unit and nonparametric smoothing techniques to consistently estimate unit-level regression values under counterfactual treatment assignments. See Hahn (1998), Heckman, Ichimura, and Todd (1998), Imbens (2004), and Imbens and Wooldridge (2009) for a discussion of such estimators.

because units with dissimilar covariate values may nevertheless have similar values for their propensity scores.

In observational studies, propensity scores are not known, so they have to be estimated prior to matching. In spite of the great popularity that propensity score matching methods have enjoyed since they were proposed by Rosenbaum and Rubin in 1983, their large sample distribution has not yet been derived for the case when the propensity score is estimated in a first step.<sup>3</sup> A possible reason for this void in the literature is that matching estimators are non-smooth functionals of the distribution of the matching variables, which makes it difficult to establish an asymptotic approximation to the distribution of matching estimators when a matching variable is estimated in a first step. This has motivated the use of bootstrap standard errors for propensity score matching estimators. However, recently it has been shown that the bootstrap is not, in general, valid for matching estimators (Abadie and Imbens (2008)).<sup>4</sup>

In this article, we derive large sample approximations to the distribution of propensity score matching estimators. Our derivations take into account that the propensity score is itself estimated in a first step. We show that propensity matching estimators have approximately Normal distributions in large samples. We demonstrate that first step estimation of the propensity score affects the large sample distribution of propensity score matching estimators, and derive adjustments to the large sample variance of propensity score matching estimators that correct for first step estimation of the propensity score. We do this for estimators of the average treatment effect (ATE) and the average treatment effect on the treated (ATET). The adjustment for the ATE estimator is negative (or zero in some special cases), implying that matching on the estimated propensity score is more efficient than matching on the true propensity score in large samples. As a result, treating the estimated propensity score as it was the true propensity score for estimating the variance of the ATE estimator leads to conservative confidence intervals. However, for the ATET estimator, the sign of the adjustment depends on the data generating process, and ignoring the estimation error in the propensity score may lead to confidence intervals that are either too large or too small.

## 2. MATCHING ESTIMATORS

The setup in this article is a standard one in the program evaluation literature, where the focus of the analysis is often the effect of a binary treatment,

<sup>3</sup>Influential papers using matching on the estimated propensity score include Heckman, Ichimura, and Todd (1997), Dehejia and Wahba (1999), and Smith and Todd (2005).

<sup>4</sup>In contexts other than matching, Heckman, Ichimura, and Todd (1998), Hirano, Imbens, and Ridder (2003), Abadie (2005), Wooldridge (2007), and Angrist and Kuersteiner (2011) derived large sample properties of statistics based on a first step estimator of the propensity score. In all these cases, the second step statistics are smooth functionals of the propensity scores and, therefore, standard stochastic expansions for two-step estimators apply (see, e.g., Newey and McFadden (1994)).

represented in this paper by the indicator variable  $W$ , on some outcome variable,  $Y$ . More specifically,  $W = 1$  indicates exposure to the treatment, while  $W = 0$  indicates lack of exposure to the treatment. Following [Rubin \(1974\)](#), we define treatment effects in terms of potential outcomes. We define  $Y(1)$  as the potential outcome under exposure to treatment, and  $Y(0)$  as the potential outcome under no exposure to treatment. Our goal is to estimate the average treatment effect,

$$\tau = E[Y(1) - Y(0)],$$

where the expectation is taken over the population of interest. Alternatively, the goal may be estimation of the average effect for the treated,

$$\tau_t = E[Y(1) - Y(0)|W = 1].$$

Estimation of these average treatment effects is complicated by the fact that for each unit in the population, we observe at most one of the potential outcomes:

$$Y = \begin{cases} Y(0) & \text{if } W = 0, \\ Y(1) & \text{if } W = 1. \end{cases}$$

Let  $X$  be a vector of covariates of dimension  $k$ . The propensity score is  $p(X) = \Pr(W = 1|X)$ , and  $p^* = \Pr(W = 1)$  is the probability of being treated. The following assumption is often referred to as “strong ignorability” ([Rosenbaum and Rubin \(1983\)](#)). It means that adjusting for  $X$  is sufficient to eliminate all confounding.

ASSUMPTION 1: (i)  $Y(1), Y(0) \perp\!\!\!\perp W|X$  *almost surely*; (ii)  $\underline{p} \leq p(X) \leq \bar{p}$  *almost surely, for some*  $\underline{p} > 0$  *and*  $\bar{p} < 1$ .

Assumption 1(i) uses the conditional independence notation in [Dawid \(1979\)](#). This assumption is often referred to as “unconfoundedness.” It will hold, for example, if all confounders are included in  $X$ , so that after controlling for  $X$ , treatment exposure is independent of the potential outcomes. [Hahn \(1998\)](#) derived asymptotic variance bounds and studied asymptotically efficient estimation under Assumption 1(i). Assumption 1(ii) implies that, for almost all values of  $X$ , the population includes treated and untreated units. Moreover, Assumption 1(ii) bounds the values of the propensity score away from zero and 1. [Khan and Tamer \(2010\)](#) have shown that this condition is necessary for root- $N$  consistent estimation of the average treatment effect.

Let  $\mu(w, x) = E[Y|W = w, X = x]$  and  $\sigma^2(w, x) = \text{var}(Y|W = w, X = x)$  be the conditional mean and variance of  $Y$  given  $W = w$  and  $X = x$ . Similarly, let  $\bar{\mu}(w, p) = E[Y|W = w, p(X) = p]$  and  $\bar{\sigma}^2(w, p) = \text{var}(Y|W = w, p(X) = p)$  be the conditional mean and variance of  $Y$  given  $W = w$  and

$p(X) = p$ . Under Assumption 1,

$$\tau = E[\mu(1, X) - \mu(0, X)]$$

and

$$\tau_i = E[\mu(1, X) - \mu(0, X) | W = 1]$$

(see [Rubin \(1974\)](#)). Therefore, adjusting for differences in the distribution of  $X$  between treated and nontreated removes all confounding and, therefore, allows identification of ATE and ATET. [Rosenbaum and Rubin \(1983\)](#) proved that  $W$  and  $X$  are independent conditional on the propensity score,  $p(X)$ , which implies that under Assumption 1:

$$\tau = E[\bar{\mu}(1, p(X)) - \bar{\mu}(0, p(X))]$$

and

$$\tau_i = E[\bar{\mu}(1, p(X)) - \bar{\mu}(0, p(X)) | W = 1].$$

In other words, under Assumption 1, adjusting for the propensity score only is enough to remove all confounding. This result motivates the use of propensity score matching estimators. A propensity score matching estimator for the average treatment effect can be defined as

$$\hat{\tau}_N^* = \frac{1}{N} \sum_{i=1}^N (2W_i - 1) \left( Y_i - \frac{1}{M} \sum_{j \in \mathcal{J}_M(i)} Y_j \right),$$

where  $M$  is a fixed number of matches per unit and  $\mathcal{J}_M(i)$  is the set of matches for unit  $i$ .<sup>5</sup> (The superscript  $*$  on  $\hat{\tau}_N^*$  indicates that matching is done on the true propensity score.) For concreteness, in this article we will consider matching with replacement, so each unit in the sample can be used as a match multiple times. In the absence of matching ties, the set of matches  $\mathcal{J}_M(i)$  can formally

<sup>5</sup>In typical applications,  $M$  is small, often  $M = 1$ . Choosing a small  $M$  reduces finite sample biases caused by matches of poor quality, that is, matches between individuals with substantial differences in their propensity score values. Larger values of  $M$  produce lower large sample variances (see [Abadie and Imbens \(2006, Section 3.4\)](#)), and one could consider increasing  $M$  in a particular application if such increase has a small effect on the size of the matching discrepancies, which are observed in the data ([Abadie and Imbens \(2011\)](#)). Similarly to [Yatchew's \(1997\)](#) work on semiparametric differencing estimators, we derive a large sample approximation to the distribution of matching estimators for fixed values of  $M$ . Large sample approximations based on fixed values of smoothing parameters have been shown to increase accuracy in other contexts (see, in particular, [Kiefer and Vogelsang \(2005\)](#)).

be defined as

$$\mathcal{J}_M(i) = \left\{ j = 1, \dots, N : W_j = 1 - W_i, \right. \\ \left. \left( \sum_{k: W_k = 1 - W_i} 1_{|p(X_i) - p(X_k)| \leq |p(X_i) - p(X_j)|} \right) \leq M \right\},$$

where  $1_{[\cdot]}$  is a binary indicator that takes value 1 if the event inside brackets is true, and value zero if not. For the average effect on the treated,  $\tau_t$ , the corresponding estimator is

$$\hat{\tau}_{t,N}^* = \frac{1}{N_1} \sum_{i=1}^N W_i \left( Y_i - \frac{1}{M} \sum_{j \in \mathcal{J}_M(i)} Y_j \right),$$

where  $N_1 = \sum_{i=1}^N W_i$  is the number of treated units in the sample.

The large sample distributions of  $\hat{\tau}_N^*$  and  $\hat{\tau}_{t,N}^*$  can be easily derived from results on matching estimators in [Abadie and Imbens \(2006\)](#), applied to the case where the only matching variable is the known propensity score. However, one of the assumptions in [Abadie and Imbens \(2006\)](#) requires that the density of the matching variables is bounded away from zero. Although this assumption may be appropriate in settings where the matching is carried out directly on the covariates,  $X$ , it is much less appealing for propensity score matching estimators. For example, if the propensity score has the form  $p(X) = F(X'\theta)$ , then even if the density of  $X$  is bounded away from zero on its support, the density of  $F(X'\theta)$  will generally not be bounded away from zero on its support. We therefore generalize the results in [Abadie and Imbens \(2006\)](#) to allow the density of propensity score to take values that are arbitrarily close to zero.

**ASSUMPTION 2:** (i) *The propensity score  $p(X)$  is continuously distributed, has interval support  $[\underline{p}, \overline{p}]$ , and has a density that is a continuous function on  $[\underline{p}, \overline{p}]$ ;* (ii) *for  $w = 0, 1$ ,  $\bar{\mu}(w, p)$  and  $\bar{\sigma}^2(w, p)$  are Lipschitz-continuous and continuous in  $p$ , respectively;* (iii) *for  $w = 0, 1$ , there exists  $\delta > 0$  such that  $E[|Y|^{2+\delta} | W = w, p(X) = p]$  is uniformly bounded.*

**ASSUMPTION 3:**  $\{(Y_i, W_i, X_i)\}_{i=1}^N$  *are independent draws from the distribution of  $(Y, W, X)$ .*

The next proposition presents the large sample distributions of  $\hat{\tau}_N^*$  and  $\hat{\tau}_{t,N}^*$  under Assumptions 1–3.

**PROPOSITION 1:** *Suppose Assumptions 1–3 hold. Then, (i)*

$$\sqrt{N}(\hat{\tau}_N^* - \tau) \xrightarrow{d} N(0, \sigma^2),$$

where

$$\begin{aligned}\sigma^2 = & E[(\bar{\mu}(1, p(X)) - \bar{\mu}(0, p(X)) - \tau)^2] \\ & + E\left[\bar{\sigma}^2(1, p(X))\left(\frac{1}{p(X)} + \frac{1}{2M}\left(\frac{1}{p(X)} - p(X)\right)\right)\right] \\ & + E\left[\bar{\sigma}^2(0, p(X))\right. \\ & \left. \times \left(\frac{1}{1-p(X)} + \frac{1}{2M}\left(\frac{1}{1-p(X)} - (1-p(X))\right)\right)\right],\end{aligned}$$

and (ii)

$$\sqrt{N}(\hat{\tau}_{t,N}^* - \tau_t) \xrightarrow{d} N(0, \sigma_t^2),$$

where

$$\begin{aligned}\sigma_t^2 = & \frac{1}{E[p(X)]^2} E[p(X)(\bar{\mu}(1, p(X)) - \bar{\mu}(0, p(X)) - \tau_t)^2] \\ & + \frac{1}{E[p(X)]^2} E[\bar{\sigma}^2(1, p(X))p(X)] \\ & + \frac{1}{E[p(X)]^2} E\left[\bar{\sigma}^2(0, p(X))\right. \\ & \left. \times \left(\frac{p^2(X)}{1-p(X)} + \frac{1}{M}p(X) + \frac{1}{2M}\frac{p^2(X)}{1-p(X)}\right)\right].\end{aligned}$$

The proof of this proposition is available in the Supplemental Material (Abadie and Imbens (2016)).

Motivated by the fact that, in observational studies, propensity scores are not known, we are interested in the case where matching is not on the true propensity score  $p(X)$ , but on an estimate of the propensity score. Following Rosenbaum and Rubin (1983) and most of the empirical literature, we consider a generalized linear specification for the propensity score,  $p(x) = F(x'\theta)$ . In empirical research, the link function  $F$  is usually specified as Logit or Probit. It is straightforward to extend our results to more general parametric models for the propensity score.<sup>6</sup> For unit  $i$ , and for arbitrary values for  $\theta$ , let  $\mathcal{J}_M(i, \theta)$

<sup>6</sup>Misspecification of the propensity score typically leads to inconsistency of the treatment effect estimator, unless the misspecified propensity score constitutes a balancing score, that is, a function,  $b(X)$ , of the covariates such that  $X \perp\!\!\!\perp W|b(X)$  (see Rosenbaum and Rubin (1983)). Moti-

denote the set of  $M$  matches where we match on  $F(X'\theta)$ :

$$\mathcal{J}_M(i, \theta) = \left\{ j = 1, \dots, N : W_j = 1 - W_i, \left( \sum_{k: W_k = 1 - W_i} 1_{|F(X'_i\theta) - F(X'_k\theta)| \leq |F(X'_i\theta) - F(X'_j\theta)|} \right) \leq M \right\}.$$

The matching estimator for the average treatment effect where we match on  $F(X'\theta)$  is then

$$\hat{\tau}_N(\theta) = \frac{1}{N} \sum_{i=1}^N (2W_i - 1) \left( Y_i - \frac{1}{M} \sum_{j \in \mathcal{J}_M(i, \theta)} Y_j \right).$$

Let  $\theta^*$  denote the true value of the propensity score model parameter vector, so that  $p(X) = F(X'\theta^*)$ . Then, the estimator based on matching on the true propensity score can be written as  $\hat{\tau}_N^* = \hat{\tau}_N(\theta^*)$ . We are interested in the case where  $\hat{\tau}_N(\theta)$  is evaluated at an estimator  $\hat{\theta}_N$  of  $\theta^*$ , based on a sample  $\{Y_i, W_i, X_i\}_{i=1}^N$ . We focus on the case where  $\hat{\theta}_N$  is the maximum likelihood estimator of  $\theta$ :<sup>7</sup>

$$\hat{\theta}_N = \arg \max_{\theta} L(\theta | W_1, X_1, \dots, W_N, X_N),$$

where the log-likelihood function is

$$\begin{aligned} L(\theta | W_1, X_1, \dots, W_N, X_N) \\ = \sum_{i=1}^N W_i \ln F(X'_i\theta) + (1 - W_i) \ln(1 - F(X'_i\theta)). \end{aligned}$$

The propensity score matching estimator of  $\tau$  that matches on the estimated propensity score can now be written as

$$\hat{\tau}_N = \hat{\tau}_N(\hat{\theta}_N) = \frac{1}{N} \sum_{i=1}^N (2W_i - 1) \left( Y_i - \frac{1}{M} \sum_{j \in \mathcal{J}_M(i, \hat{\theta}_N)} Y_j \right).$$

---

vated by this consideration, empirical researchers routinely use measures of balance in the distribution of the covariates between treated and nontreated, conditional on the estimated propensity score, to perform specification searches on the propensity score (see, e.g., Dehejia and Wahba (1999)). An alternative safeguard against misspecification of the propensity score is the use “doubly robust” matching estimators, like the bias-corrected matching estimator of Abadie and Imbens (2011).

<sup>7</sup>It is straightforward to extend our results to other asymptotically linear estimators of  $\theta^*$ .

Similarly, the propensity score matching estimator of  $\tau_t$  that matches on the estimated propensity score can be written as

$$\hat{\tau}_{t,N} = \hat{\tau}_{t,N}(\hat{\theta}_N) = \frac{1}{N_1} \sum_{i=1}^N W_i \left( Y_i - \frac{1}{M} \sum_{j \in \mathcal{J}_M(i, \hat{\theta}_N)} Y_j \right).$$

Whenever confusion is possible, we will be explicit in the dependence of the matching estimators on  $\theta$ . If the argument is omitted,  $\hat{\tau}_N$  and  $\hat{\tau}_{t,N}$  are used as shorthand for  $\hat{\tau}_N(\hat{\theta}_N)$  and  $\hat{\tau}_{t,N}(\hat{\theta}_N)$ , respectively.

The two main questions addressed in this article are (i) do the estimators based on matching on the estimated propensity score have Normal large sample distributions, and (ii) if so, how does their large sample variance compare to that of the estimators that match on the true propensity score, given in Proposition 1? In the next section, we answer these two questions and derive the large sample distribution of  $\hat{\tau}_N(\hat{\theta}_N)$  and  $\hat{\tau}_{t,N}(\hat{\theta}_N)$ . Conventional linearization methods for two-step statistics are difficult to apply in the context of matching estimators because matching estimators are complicated functionals of the distribution of the data. We therefore follow a different route, building on work by [Andreou and Werker \(2012\)](#) on residual based statistics, and the martingale representations for matching estimators derived in [Abadie and Imbens \(2012\)](#).

### 3. LARGE SAMPLE DISTRIBUTION

In the first part of this section, we derive the large sample approximation to the sampling distribution of  $\hat{\tau}_N(\hat{\theta}_N)$ , and in the second part, we present the results for  $\hat{\tau}_{t,N}(\hat{\theta}_N)$ .

Let  $P^\theta$  be the distribution of  $Z = \{Y, W, X\}$  induced by the propensity score,  $F(X'\theta)$ , the marginal distribution of  $X$ , and the conditional distribution of  $Y$  given  $X$  and  $W$ . We index this distribution  $P^\theta$  by  $\theta$ , and will consider properties of estimators for different values of  $\theta$ , under the same marginal distribution for  $X$ , and the same conditional distribution for  $Y$  given  $W$  and  $X$ . Given Assumption 1, the average treatment effect is equal to

$$\tau = E[Y(1) - Y(0)] = E[E[Y|W = 1, X] - E[Y|W = 0, X]].$$

From this equation, it can be seen that ATE does not depend on the propensity score; it only depends on the conditional distribution of  $Y$  given  $W$  and  $X$  and the marginal distribution of  $X$ . The average treatment effect for the treated is

$$\begin{aligned} \tau_t &= E[Y(1) - Y(0)|W = 1] \\ &= E[E[Y|W = 1, X] - E[Y|W = 0, X]|W = 1] \\ &= \frac{E[F(X'\theta^*)(E[Y|W = 1, X] - E[Y|W = 0, X])]}{E[F(X'\theta^*)]}. \end{aligned}$$



In contrast to the average treatment effect,  $\tau$ , the average treatment effect for the treated,  $\tau_t$ , depends on the propensity score, and we make this dependence explicit by indexing  $\tau_t$  by  $\theta$  wherever appropriate. In particular,  $\tau_t = \tau_t(\theta^*)$  is the average effect of the treatment on the treated.

To derive the large sample distribution of  $\hat{\tau}_N$  and  $\hat{\tau}_{t,N}$ , we invoke some additional regularity conditions. First, we extend Assumption 2 to hold for all  $\theta$  in a neighborhood of  $\theta^*$ .

ASSUMPTION 4: (i)  $\theta^* \in \text{int}(\Theta)$  with  $\Theta$  compact,  $X$  has a bounded support, and  $E[XX']$  is nonsingular; (ii)  $F: \mathbb{R} \mapsto (0, 1)$  is continuously differentiable with strictly positive and bounded derivative  $f$ ; (iii) there exists a component of  $X$  that is continuously distributed, has nonzero coefficient in  $\theta^*$ , and admits a continuous density function conditional on the rest of  $X$ ; and (iv) there exists  $\varepsilon > 0$  such that, for all  $\theta$  with  $\|\theta - \theta^*\| \leq \varepsilon$ ,  $E[Y|W = w, F(X'\theta) = p]$  is Lipschitz-continuous in  $p$ ,  $\text{var}(Y|W = w, F(X'\theta) = p)$  is continuous in  $p$ , and there is  $\delta > 0$  such that  $E[|Y|^{2+\delta}|W = w, F(X'\theta) = p]$  is uniformly bounded.

The case of only discrete regressors is left out from Assumption 4, but, as noted in Abadie and Imbens (2006), it is a simple case to treat separately. With only discrete regressors, each observation (or each treated observation in the case of ATET) can be matched to every observation with identical estimated propensity score value in the opposite treatment arm. In this case, the propensity score matching estimator is identical to the subclassification estimator in Cochran (1968) and Angrist (1998), and valid analytical and bootstrap standard errors can be easily derived.

We will study the behavior of certain statistics under sequences  $\theta_N$  that are local to  $\theta^*$ . Consider  $Z_{N,i} = \{Y_{N,i}, W_{N,i}, X_{N,i}\}$  with distribution given by the local “shift”  $P^{\theta_N}$  with  $\theta_N = \theta^* + h/\sqrt{N}$ , where  $h$  is a conformable vector of constants. Let  $\Delta_N(\theta|\theta')$  be the difference between the value of the log-likelihood function evaluated at  $\theta$  and the value of the log-likelihood function evaluated at  $\theta'$ :

$$\Delta_N(\theta|\theta') = L(\theta|Z_{N,1}, \dots, Z_{N,N}) - L(\theta'|Z_{N,1}, \dots, Z_{N,N}).$$

Let  $\Delta_N(\theta)$  be the normalized score function, or central sequence,

$$\begin{aligned} \Delta_N(\theta) &= \frac{1}{\sqrt{N}} \frac{\partial}{\partial \theta} L(\theta|Z_{N,1}, \dots, Z_{N,N}) \\ &= \frac{1}{\sqrt{N}} \sum_{i=1}^N X_{N,i} \frac{W_{N,i} - F(X'_{N,i}\theta)}{F(X'_{N,i}\theta)(1 - F(X'_{N,i}\theta))} f(X'_{N,i}\theta). \end{aligned}$$

Finally, let

$$I_{\theta} = E \left[ \frac{f(X'\theta)^2}{F(X'\theta)(1 - F(X'\theta))} XX' \right]$$

be the Fisher Information Matrix for  $\theta$ . The expectation in this equation is taken over the marginal distribution of  $X$ , which does not depend on  $\theta$ , so the indexing by  $\theta$  solely reflects the value of  $\theta$  where  $f(X'\theta)$  and  $F(X'\theta)$  are evaluated. The following intermediate lemma derives some important regularity properties of the propensity score model that will be needed for our derivations.

LEMMA 1: *Suppose that Assumptions 3, 4(i), and 4(ii) hold. Then, under  $P^{\theta_N}$ ,*

$$(1) \quad \Delta_N(\theta^*|\theta_N) = -h'\Delta_N(\theta_N) - \frac{1}{2}h'I_{\theta^*}h + o_p(1),$$

$$(2) \quad \Delta_N(\theta_N) \xrightarrow{d} N(0, I_{\theta^*}),$$

where  $I_{\theta^*}$  is not singular, and

$$(3) \quad \sqrt{N}(\widehat{\theta}_N - \theta_N) = I_{\theta^*}^{-1}\Delta_N(\theta_N) + o_p(1).$$

The proof of this proposition is available in the Supplemental Material (Abadie and Imbens (2016)). For regular parametric models, equation (1) can be established using Proposition 2.1.2 in Bickel, Klaassen, Ritov, and Wellner (1998). Also for regular parametric models, equation (2) is derived in the proof of Proposition 2.1.2 in Bickel et al. (1998). Equation (3) can be established using the same set of results in combination with classical conditions for asymptotic linearity of maximum likelihood estimators (see, e.g., van der Vaart (1998, Theorem 5.39), Lehmann and Romano (2005, Theorem 12.4.1)).

Let  $E_{\theta}$  be the expectation operator with respect to the distributions  $P^{\theta}$ . The following assumption is a regularity condition that will be used later in this section.

ASSUMPTION 5:  $E_{\theta_N}[r(Y, W, X)|W, F(X'\theta_N)]$  converges to  $E[r(Y, W, X)|W, F(X'\theta^*)]$  almost surely, for any  $\mathbb{R}^{k+2}$ -to- $\mathbb{R}$  bounded and measurable function,  $r(y, w, x)$ , continuous in  $x$ , and any sequence,  $\theta_N \rightarrow \theta^*$ .

Primitive conditions for this assumption are provided in the Supplemental Material.

### 3.1. Large Sample Distribution for $\widehat{\tau}_N(\widehat{\theta}_N)$

Our derivation of the limit distribution of  $\sqrt{N}(\widehat{\tau}_N - \tau)$  is based on the techniques developed in [Andreou and Werker \(2012\)](#) to analyze the limit distribution of residual-based statistics. We proceed in three steps. First, we derive the joint limit distribution of  $(\sqrt{N}(\widehat{\tau}_N(\theta_N) - \tau), \sqrt{N}(\widehat{\theta} - \theta_N), \Lambda_N(\theta^*|\theta_N))$  under  $P^{\theta_N}$ .

PROPOSITION 2: Suppose that Assumptions 1–5 hold. Then, under  $P^{\theta_N}$ ,

$$\begin{pmatrix} \sqrt{N}(\widehat{\tau}_N(\theta_N) - \tau) \\ \sqrt{N}(\widehat{\theta}_N - \theta_N) \\ \Lambda_N(\theta^*|\theta_N) \end{pmatrix} \xrightarrow{d} N \left( \begin{pmatrix} 0 \\ 0 \\ -h'I_{\theta^*}h/2 \end{pmatrix}, \begin{pmatrix} \sigma^2 & c'I_{\theta^*}^{-1} & -c'h \\ I_{\theta^*}^{-1}c & I_{\theta^*}^{-1} & -h \\ -h'c & -h' & h'I_{\theta^*}h \end{pmatrix} \right),$$

where

$$(4) \quad c = E \left[ \left( \frac{\text{cov}(X, \mu(1, X)|F(X'\theta^*))}{F(X'\theta^*)} + \frac{\text{cov}(X, \mu(0, X)|F(X'\theta^*))}{1 - F(X'\theta^*)} \right) f(X'\theta^*) \right].$$

All proofs for the results in this section are provided in the [Appendix](#).

Asymptotic Normality of the first component,  $\sqrt{N}(\widehat{\tau}_N(\theta_N) - \tau)$ , under  $P^{\theta^*}$  follows from Proposition 1. Asymptotic joint Normality of the last two components,  $\sqrt{N}(\widehat{\theta}_N - \theta_N)$  and  $\Lambda_N(\theta^*|\theta_N)$ , follows from Lemma 1. Proposition 2 derives the joint large sample distribution of the three components under  $P^{\theta_N}$ . The proof extends the martingale techniques of [Abadie and Imbens \(2012\)](#) to derive the result of the proposition.

In the second step of our argument, we use Le Cam's third lemma (e.g., [van der Vaart \(1998, p. 90\)](#)). Given the result of Proposition 2, Le Cam's third lemma implies that, under  $P^{\theta^*}$ ,

$$\begin{pmatrix} \sqrt{N}(\widehat{\tau}_N(\theta_N) - \tau) \\ \sqrt{N}(\widehat{\theta}_N - \theta_N) \end{pmatrix} \xrightarrow{d} N \left( \begin{pmatrix} -c'h \\ -h \end{pmatrix}, \begin{pmatrix} \sigma^2 & c'I_{\theta^*}^{-1} \\ I_{\theta^*}^{-1}c & I_{\theta^*}^{-1} \end{pmatrix} \right).$$

Substituting  $\theta_N = \theta^* + h/\sqrt{N}$ , this implies that (still under  $P^{\theta^*}$ ), for any  $h \in \mathbb{R}^k$ ,

$$(5) \quad \begin{pmatrix} \sqrt{N}(\widehat{\tau}_N(\theta^* + h/\sqrt{N}) - \tau) \\ \sqrt{N}(\widehat{\theta}_N - \theta^*) \end{pmatrix} \xrightarrow{d} N \left( \begin{pmatrix} -c'h \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma^2 & c'I_{\theta^*}^{-1} \\ I_{\theta^*}^{-1}c & I_{\theta^*}^{-1} \end{pmatrix} \right).$$

If equation (5) was an exact result rather than an approximation based on convergence in distribution, it would directly lead to the result of interest. In that case, it would follow that

$$\begin{aligned} & \sqrt{N}(\widehat{\tau}_N(\theta^* + h/\sqrt{N}) - \tau) | \sqrt{N}(\widehat{\theta}_N - \theta^*) \\ &= h \sim N(0, \sigma^2 - c' I_{\theta^*}^{-1} c) \end{aligned}$$

(see, e.g., [Goldberger \(1991, p. 197\)](#)). Because  $\sqrt{N}(\widehat{\theta}_N - \theta^*) = h$  implies  $\theta^* + h/\sqrt{N} = \widehat{\theta}_N$ , and thus implies that  $\widehat{\tau}_N(\theta^* + h/\sqrt{N}) = \widehat{\tau}_N(\widehat{\theta}_N) = \widehat{\tau}_N$ , the last displayed equation can also be written as

$$\sqrt{N}(\widehat{\tau}_N - \tau) | \sqrt{N}(\widehat{\theta}_N - \theta) = h \sim N(0, \sigma^2 - c' I_{\theta^*}^{-1} c).$$

Because this conditional distribution does not depend on  $h$ , this in turn implies that, under  $P^{\theta^*}$ , unconditionally,

$$\sqrt{N}(\widehat{\tau}_N - \tau) \sim N(0, \sigma^2 - c' I_{\theta^*}^{-1} c),$$

which is the result we are looking for: the distribution of the matching estimator based on matching on the estimated propensity score.

A challenge formalizing this argument is that convergence of  $\sqrt{N}(\widehat{\tau}_N - \tau) | \sqrt{N}(\widehat{\theta}_N - \theta^*) = h$  involves convergence in a conditioning event. To overcome this challenge, in the third step of the argument, we employ a Le Cam discretization device, as proposed in [Andreou and Werker \(2012\)](#). Consider a grid of cubes in  $\mathbb{R}^k$  with sides of length  $d/\sqrt{N}$ , for arbitrary positive  $d$ . Then  $\bar{\theta}_N$  is the discretized estimator, defined as the midpoint of the cube  $\widehat{\theta}_N$  belongs to. If  $\widehat{\theta}_{N,j}$  is the  $j$ th component of the  $k$ -vector  $\widehat{\theta}_N$ , then the  $j$ th component of the  $k$ -vector  $\bar{\theta}_N$  is  $\bar{\theta}_{N,j} = (d/\sqrt{N})[\sqrt{N}\widehat{\theta}_{N,j}/d]$ , where  $[\cdot]$  is the nearest integer function. Now we can state the main result of the paper.

**THEOREM 1:** *Suppose Assumptions 1–5 hold. Then, under  $P^{\theta^*}$ ,*

$$\begin{aligned} & \lim_{d \downarrow 0} \lim_{N \rightarrow \infty} \Pr(\sqrt{N}(\sigma^2 - c' I_{\theta^*}^{-1} c)^{-1/2} (\widehat{\tau}_N(\bar{\theta}_N) - \tau) \leq z) \\ &= \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}x^2\right) dx. \end{aligned}$$

An implication of Theorem 1 is that we can approximate the distribution of  $\sqrt{N}(\widehat{\tau}_N(\bar{\theta}_N) - \tau)$  by a Normal distribution with mean zero and variance  $\sigma^2 - c' I_{\theta^*}^{-1} c$ . The result in Theorem 1 indicates that the adjustment to the stan-

dard error of the propensity score matching estimator for first step estimation of the propensity score is always negative, or zero in some special cases as discussed below. This implies that using matching on the estimated propensity score, rather than on the true propensity score, to estimate ATE increases precision in large samples. As we will see later, this gain in precision from using the estimated propensity score does not necessarily hold for the estimation of parameters different than ATE.

Equation (4) and the fact that  $X$  and  $W$  are independent conditional on the propensity score imply that if the covariance of  $X$  and  $\mu(W, X)$  given  $F(X'\theta^*)$  and  $W$  is equal to zero, then  $c = 0$  and first step estimation of the propensity score does not affect the large sample variance of  $\sqrt{N}(\hat{\theta}_N - \theta^*)$ . This would be the case if the propensity score provides no “dimension reduction,” that is, if the propensity score is a bijective function of  $X$ . In that case, each value of the propensity score corresponds to only one value of  $X$ , so  $\text{cov}(X, \mu(W, X)|W, F(X'\theta^*)) = 0$  and, therefore,  $c = 0$ .

For concreteness, our derivations focus on the case of matching with replacement. However, as shown in [Abadie and Imbens \(2012\)](#), martingale representations analogous to the one employed in the proof of Proposition 2 exist for alternative matching estimators (e.g., estimators that construct the matches without using replacement). An inspection of the proof of Proposition 2 reveals that the adjustment term,  $-c'I_{\theta^*}c$ , does not depend on the type of matching employed to obtain the estimators. Therefore, the result in Theorem 1 translates easily to other matching settings. A different type of matching scheme may change the form of  $\sigma^2$  in the result of Theorem 1, but not the adjustment term,  $-c'I_{\theta^*}c$ .

### 3.2. Large Sample Distribution for $\hat{\tau}_{t,N}$

In this section, we consider the asymptotic distribution for  $\sqrt{N}(\hat{\tau}_{t,N} - \tau_t(\theta^*))$ . The derivations are similar to those for the case of ATE, so we relegate details to the Supplemental Material. The following theorem provides the result.

**THEOREM 2:** *Suppose Assumptions 1–5 hold, and let*

$$\begin{aligned} c_t = & \frac{1}{E[F(X'\theta^*)]} E[Xf(X'\theta^*)(\bar{\mu}(1, F(X'\theta^*)) - \bar{\mu}(0, F(X'\theta^*)) - \tau_t)] \\ & + \frac{1}{E[F(X'\theta^*)]} E\left[\left(\text{cov}(X, \mu(1, X)|F(X'\theta^*))\right.\right. \\ & \left.\left.+ \frac{F(X'\theta^*)}{1 - F(X'\theta^*)} \text{cov}(X, \mu(0, X)|F(X'\theta^*))\right)f(X'\theta^*)\right]. \end{aligned}$$

Then, under  $P^{\theta^*}$ ,

$$\begin{aligned} & \lim_{d \downarrow 0} \lim_{N \rightarrow \infty} \Pr \left( \sqrt{N} \left( \sigma_t^2 - c_t' I_{\theta^*}^{-1} c_t + \frac{\partial \tau_t(\theta^*)'}{\partial \theta} I_{\theta^*}^{-1} \frac{\partial \tau_t(\theta^*)}{\partial \theta} \right)^{-1/2} \right. \\ & \quad \left. \times (\hat{\tau}_{t,N}(\bar{\theta}_N) - \tau_t) \leq z \right) \\ &= \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}x^2\right) dx. \end{aligned}$$

Notice that, in contrast to the ATE case, the adjustment for first step estimation of the propensity score for the ATET estimator may result in a decrease *or* an increase in the standard error. The adjustment to the standard error of the ATET estimator will be positive, for example, if the propensity score does not provide “dimension reduction,” so  $c_t = 0$ , and  $\partial \tau_t(\theta^*)/\partial \theta \neq 0$  (which will typically be the case if the average effect of the treatment varies with the covariates,  $X$ ). In contrast, if  $c_t \neq 0$  and  $\partial \tau_t(\theta^*)/\partial \theta = 0$ , the adjustment is negative. Like for the case of ATE, it can be shown that the adjustment term does not depend on the particular type of matching estimator of ATET.

#### 4. ESTIMATION OF THE ASYMPTOTIC VARIANCE

In this section, we discuss estimation of the large sample variances of ATE and ATET adjusting for first step estimation of the propensity score. As shown in the previous section, the asymptotic variance for  $\sqrt{N}(\hat{\tau}_N - \tau)$  is

$$(6) \quad \sigma_{\text{adj}}^2 = \sigma^2 - c' I_{\theta^*}^{-1} c,$$

and the asymptotic variance for  $\sqrt{N}(\hat{\tau}_{t,N} - \tau_t)$  is

$$(7) \quad \sigma_{t,\text{adj}}^2 = \sigma_t^2 - c_t' I_{\theta^*}^{-1} c_t + \frac{\partial \tau_t(\theta^*)'}{\partial \theta} I_{\theta^*}^{-1} \frac{\partial \tau_t(\theta^*)}{\partial \theta}.$$

To estimate  $\sigma_{\text{adj}}^2$  and  $\sigma_{t,\text{adj}}^2$ , we define estimators of each of the components of the right-hand sides of equations (6) and (7). First, estimation of the information matrix,  $I_{\theta^*}$ , is standard:

$$\hat{I}_{\theta^*} = \frac{1}{N} \sum_{i=1}^N \frac{f(X_i' \hat{\theta}_N)^2}{F(X_i' \hat{\theta}_N)(1 - F(X_i' \hat{\theta}_N))} X_i X_i'.$$

Consider next estimation of the variances corresponding to matching on the true propensity score,  $\sigma^2$  and  $\sigma_t^2$ . For these components, we use estimators

that are based on those in [Abadie and Imbens \(2006\)](#). Let  $K_{M,\theta}(i)$  be the number of times that observation  $i$  is used as a match (when matching on  $F(X'\theta)$ ):

$$(8) \quad K_{M,\theta}(i) = \sum_{j=1}^N 1_{[i \in \mathcal{J}_M(j, \theta)]},$$

and let  $\hat{\sigma}^2(W_i, F(X'_i\theta^*))$  be an asymptotically unbiased (but not necessarily consistent) estimator of  $\bar{\sigma}^2(W_i, F(X'_i\theta^*))$ . The [Abadie and Imbens \(2006\)](#) variance estimators are

$$\begin{aligned} \hat{\sigma}^2 &= \frac{1}{N} \sum_{i=1}^N \left( (2W_i - 1) \left( Y_i - \frac{1}{M} \sum_{j \in \mathcal{J}_M(i, \hat{\theta})} Y_j \right) - \hat{\tau}_N \right)^2 \\ &\quad + \frac{1}{N} \sum_{i=1}^N \left( \left( \frac{K_{M,\hat{\theta}}(i)}{M} \right)^2 + \frac{2M-1}{M} \left( \frac{K_{M,\hat{\theta}}(i)}{M} \right) \right) \hat{\sigma}^2(W_i, F(X'_i\theta^*)) \end{aligned}$$

and

$$\begin{aligned} \hat{\sigma}_t^2 &= \frac{N}{N_1^2} \sum_{i=1}^N W_i \left( Y_i - \frac{1}{M} \sum_{j \in \mathcal{J}_M(i, \hat{\theta})} Y_j - \hat{\tau}_{t,N} \right)^2 \\ &\quad + \frac{N}{N_1^2} \sum_{i=1}^N (1 - W_i) \left( \frac{K_{M,\hat{\theta}}(i)(K_{M,\hat{\theta}}(i) - 1)}{M^2} \right) \hat{\sigma}^2(W_i, F(X'_i\theta^*)). \end{aligned}$$

To obtain the estimator  $\hat{\sigma}^2(W_i, F(X'_i\theta^*))$ , let  $\mathcal{H}_L(i, \theta)$  be the set of units in the same treatment arm as unit  $i$  that have the closest  $L$  values of  $F(X'\theta)$  to  $F(X'_i\theta)$ ,

$$\begin{aligned} \mathcal{H}_L(i, \theta) &= \left\{ j = 1, \dots, N : W_j = W_i, \right. \\ &\quad \left. \left( \sum_{k: W_k = W_i} 1_{[|F(X'_i\theta) - F(X'_k\theta)| \leq |F(X'_i\theta) - F(X'_j\theta)|]} \right) \leq L \right\}, \end{aligned}$$

where  $L$  is generic notation for a (small) positive integer. Later, we will also use the set  $\mathcal{H}_L^{(-i)}(i, \theta)$ , which is similarly defined but excludes  $i$ :

$$\begin{aligned} \mathcal{H}_L^{(-i)}(i, \theta) &= \left\{ j = 1, \dots, N : i \neq j, W_j = W_i, \right. \\ &\quad \left. \left( \sum_{k: W_k = W_i, k \neq i} 1_{[|F(X'_i\theta) - F(X'_k\theta)| \leq |F(X'_i\theta) - F(X'_j\theta)|]} \right) \leq L \right\}. \end{aligned}$$

The sets  $\mathcal{H}_L(i, \theta)$ ,  $\mathcal{H}_L^{(-i)}(i, \theta)$ , and  $\mathcal{J}_L(i, \theta)$  (this last one defined in Section 2) will be used to estimate the different components of  $\sigma_{\text{adj}}^2$  and  $\sigma_{t, \text{adj}}^2$ . The value of  $L$  can vary for different components.

For  $L \geq 2$  (typically,  $L = 2$ ), consider the following matching estimator of  $\bar{\sigma}^2(W_i, F(X_i' \theta^*))$ :

$$\hat{\bar{\sigma}}^2(W_i, F(X_i' \theta^*)) = \frac{1}{L-1} \sum_{j \in \mathcal{H}_L(i, \hat{\theta}_N)} \left( Y_j - \frac{1}{L} \sum_{k \in \mathcal{H}_L(i, \hat{\theta}_N)} Y_k \right)^2.$$

That is,  $\hat{\bar{\sigma}}^2(W_i, F(X_i' \theta^*))$  is a local variance estimator that uses information only from units with the same value of  $W$  as unit  $i$  and with similar values of  $F(X' \hat{\theta}_N)$ . Because  $\hat{\theta}_N - \theta^*$  converges in probability to zero,  $\hat{\bar{\sigma}}^2(W_i, F(X_i' \theta^*))$  becomes asymptotically unbiased as  $N \rightarrow \infty$ . But because  $L$  is fixed, the variance of  $\hat{\bar{\sigma}}^2(W_i, F(X_i' \theta^*))$  does not converge to zero and  $\hat{\bar{\sigma}}^2(W_i, F(X_i' \theta^*))$  is not consistent for  $\bar{\sigma}^2(W_i, F(X_i' \theta^*))$ . The objects of interest, however, are  $\sigma^2$  and  $\sigma_t^2$ , rather than  $\bar{\sigma}^2(W_i, F(X_i' \theta^*))$ . The estimators  $\hat{\sigma}^2$  and  $\hat{\sigma}_t^2$  average terms that are bounded in probability and asymptotically unbiased. As a result, as shown in [Abadie and Imbens \(2006\)](#),  $\hat{\sigma}^2$  and  $\hat{\sigma}_t^2$  are consistent for  $\sigma^2$  and  $\sigma_t^2$ , respectively.

Next consider estimation of  $c$  and  $c_t$ . Notice first that

$$\begin{aligned} & \text{cov}(X, Y | F(X' \theta^*), W = 1) \\ &= \text{cov}(X, \mu(1, X) | F(X' \theta^*), W = 1) \\ & \quad + \text{cov}(X, Y - \mu(1, X) | F(X' \theta^*), W = 1) \\ &= \text{cov}(X, \mu(1, X) | F(X' \theta^*)) \\ & \quad + \text{cov}(X, Y - \mu(1, X) | F(X' \theta^*), W = 1). \end{aligned}$$

Using the Law of Iterated Expectations,

$$\begin{aligned} & \text{cov}(X, Y - \mu(1, X) | F(X' \theta^*), W = 1) \\ &= E[X(Y - \mu(1, X)) | F(X' \theta^*), W = 1] \\ & \quad - E[X | F(X' \theta^*), W = 1] \\ & \quad \times E[Y - \mu(1, X) | F(X' \theta^*), W = 1] \\ &= E[X(\mu(1, X) - \mu(1, X)) | F(X' \theta^*), W = 1] \\ & \quad - E[X | F(X' \theta^*), W = 1] \\ & \quad \times E[\mu(1, X) - \mu(1, X) | F(X' \theta^*), W = 1] \\ &= 0. \end{aligned}$$



Therefore,

$$\text{cov}(X, \mu(1, X)|F(X'\theta^*)) = \text{cov}(X, Y|F(X'\theta^*), W = 1),$$

and the analogous result is valid conditional on  $W_i = 0$ :

$$\text{cov}(X, \mu(0, X)|F(X'\theta^*)) = \text{cov}(X, Y|F(X'\theta^*), W = 0).$$

If  $W_i = w$ ,  $\text{cov}(X_i, \mu(w, X_i)|F(X_i'\theta^*))$  can be estimated as

$$\begin{aligned} & \widehat{\text{cov}}(X_i, \mu(w, X_i)|F(X_i'\theta^*)) \\ &= \frac{1}{L-1} \sum_{j \in \mathcal{H}_L(i, \hat{\theta}_N)} \left( X_j - \frac{1}{L} \sum_{k \in \mathcal{H}_L(i, \hat{\theta}_N)} X_k \right) \left( Y_j - \frac{1}{L} \sum_{k \in \mathcal{H}_L(i, \hat{\theta}_N)} Y_k \right), \end{aligned}$$

for  $L \geq 2$  (typically,  $L = 2$ ). If  $W_i \neq w$ , then

$$\begin{aligned} & \widehat{\text{cov}}(X_i, \mu(w, X_i)|F(X_i'\theta^*)) \\ &= \frac{1}{L-1} \sum_{j \in \mathcal{J}_L(i, \hat{\theta}_N)} \left( X_j - \frac{1}{L} \sum_{k \in \mathcal{J}_L(i, \hat{\theta}_N)} X_k \right) \left( Y_j - \frac{1}{L} \sum_{k \in \mathcal{J}_L(i, \hat{\theta}_N)} Y_k \right), \end{aligned}$$

also for  $L \geq 2$  (typically,  $L = 2$ ). Like  $\hat{\sigma}^2(W_i, F(X_i'\theta^*))$ , the estimators  $\widehat{\text{cov}}(X_i, \mu(w, X_i)|F(X_i'\theta^*))$  are asymptotically unbiased and bounded in probability. This allows us to construct a consistent analog estimator of  $c$  that averages  $\widehat{\text{cov}}(X_i, \mu(w, X_i)|F(X_i'\theta^*))$  over the sample:

$$\begin{aligned} \hat{c} &= \frac{1}{N} \sum_{i=1}^N \left( \frac{\widehat{\text{cov}}(X_i, \mu(1, X_i)|F(X_i'\theta^*))}{F(X_i'\hat{\theta}_N)} \right. \\ & \quad \left. + \frac{\widehat{\text{cov}}(X_i, \mu(0, X_i)|F(X_i'\theta^*))}{1 - F(X_i'\hat{\theta}_N)} \right) f(X_i'\hat{\theta}_N). \end{aligned}$$

For  $c_t$ , we propose separate estimators for the two components,  $c_{t,1}$  and  $c_{t,2}$ , where

$$\begin{aligned} c_{t,1} &= \frac{1}{E[F(X'\theta^*)]} E[Xf(X'\theta^*) \\ & \quad \times (\bar{\mu}(1, F(X'\theta^*)) - \bar{\mu}(0, F(X'\theta^*)) - \tau_t)], \\ c_{t,2} &= \frac{1}{E[F(X'\theta^*)]} E \left[ \left( \text{cov}(X, \mu(1, X)|F(X'\theta^*)) \right. \right. \\ & \quad \left. \left. + \frac{F(X'\theta^*)}{1 - F(X'\theta^*)} \text{cov}(X, \mu(0, X)|F(X'\theta^*)) \right) f(X'\theta^*) \right], \end{aligned}$$

and  $c_t = c_{t,1} + c_{t,2}$ . The second component,  $c_{t,2}$ , is similar to  $c_t$ , and our proposed estimator for  $c_{t,2}$  is correspondingly similar to the estimator for  $c$ :

$$(9) \quad \widehat{c}_{t,2} = \frac{1}{N_1} \sum_{i=1}^N \left( \widehat{\text{cov}}(X_i, \mu(1, X_i) | F(X'_i \theta^*)) \right. \\ \left. + \frac{F(X'_i \widehat{\theta}_N)}{1 - F(X'_i \widehat{\theta}_N)} \widehat{\text{cov}}(X_i, \mu(0, X_i) | F(X'_i \theta^*)) \right) f(X'_i \widehat{\theta}_N).$$

The first component,  $c_{t,1}$ , involves the regression functions  $\bar{\mu}(w, F(X' \theta^*))$ . We estimate these regression functions using matching:

$$\widehat{\bar{\mu}}(0, F(X'_i \theta^*)) = \begin{cases} \frac{1}{L} \sum_{j \in \mathcal{H}_L^{(-i)}(i, \widehat{\theta}_N)} Y_j & \text{if } W_i = 0, \\ \frac{1}{L} \sum_{j \in \mathcal{J}_L(i, \widehat{\theta}_N)} Y_j & \text{if } W_i = 1, \end{cases}$$

and

$$\widehat{\bar{\mu}}(1, F(X'_i \theta^*)) = \begin{cases} \frac{1}{L} \sum_{j \in \mathcal{J}_L(i, \widehat{\theta}_N)} Y_j & \text{if } W_i = 0, \\ \frac{1}{L} \sum_{j \in \mathcal{H}_L^{(-i)}(i, \widehat{\theta}_N)} Y_j & \text{if } W_i = 1, \end{cases}$$

for  $L \geq 1$  (typically,  $L = 1$ ). Our proposed estimator for  $c_{t,1}$  is

$$\widehat{c}_{t,1} = \frac{1}{N_1} \sum_{i=1}^N X_i f(X'_i \widehat{\theta}_N) (\widehat{\bar{\mu}}(1, F(X'_i \theta^*)) - \widehat{\bar{\mu}}(0, F(X'_i \theta^*)) - \widehat{\tau}_t).$$

Notice that if  $W_i = w$ , the estimator of  $\bar{\mu}(w, F(X'_i \theta^*))$  is an average over  $\mathcal{H}_L^{(-i)}(i, \widehat{\theta}_N)$ , rather than over  $\mathcal{H}_L(i, \widehat{\theta}_N)$ . If observation  $i$  was not excluded to estimate  $\bar{\mu}(W_i, F(X'_i \theta^*))$ , then  $Y_i$  would be one of the terms of the average  $\widehat{\bar{\mu}}(W_i, F(X'_i \theta^*))$ . Therefore, if observation  $i$  was not excluded, the estimator  $\widehat{c}_{t,1}$  would contain terms of the type  $X_i f(X'_i \widehat{\theta}_N) (Y_i - \widehat{\bar{\mu}}(0, F(X'_i \theta^*)) - \widehat{\tau}_t)$  when  $W_i = 1$  and terms of the type  $X_i f(X'_i \widehat{\theta}_N) (\widehat{\bar{\mu}}(1, F(X'_i \theta^*)) - Y_i - \widehat{\tau}_t)$  when  $W_i = 0$ . These terms estimate  $E[Xf(X' \theta^*)(\mu(1, X) - \bar{\mu}(0, F(X' \theta^*)) - \tau_t)]$  and  $E[Xf(X' \theta^*)(\bar{\mu}(1, F(X' \theta^*)) - \mu(0, X) - \tau_t)]$ , respectively, rather than  $E[Xf(X' \theta^*)(\bar{\mu}(1, F(X' \theta^*)) - \bar{\mu}(0, F(X' \theta^*)) - \tau_t)]$ . To avoid this problem, we exclude observation  $i$  for the estimation of  $\bar{\mu}(W_i, F(X'_i \theta^*))$ .

For the remaining variance component,  $\partial\tau_t(\theta^*)/\partial\theta$ , notice that

$$\begin{aligned}\frac{\partial\tau_t}{\partial\theta}(\theta^*) &= \frac{1}{E[F(X'\theta^*)]}E[Xf(X'\theta^*)(Y(1) - Y(0) - \tau_t)] \\ &= \frac{1}{E[F(X'\theta^*)]}E[Xf(X'\theta^*)(\mu(1, X) - \mu(0, X) - \tau_t)].\end{aligned}$$

To estimate this component, we need to estimate the regression functions  $\mu(1, X)$  and  $\mu(0, X)$ , which is done by matching on the covariates rather than on the propensity score. Define the matching set (on covariates):

$$\mathcal{J}_L^X(i) = \left\{ j = 1, \dots, N : W_j = 1 - W_i, \left( \sum_{k: W_k = 1 - W_i} 1_{[\|X_i - X_k\| \leq \|X_i - X_j\|]} \right) \leq L \right\}.$$

Our estimator of  $\partial\tau_t(\theta^*)/\partial\theta$  is

$$\frac{\widehat{\partial\tau_t}}{\partial\theta} = \frac{1}{N_1} \sum_{i=1}^N X_i f(X_i' \widehat{\theta}_N) \left( (2W_i - 1) \left( Y_i - \frac{1}{L} \sum_{j \in \mathcal{J}_L^X(i)} Y_j \right) - \widehat{\tau}_t \right).$$

Putting these results together, our estimator of the large sample variance of the propensity score matching estimator of the average treatment effect, adjusted for first step estimation of the propensity score, is

$$\widehat{\sigma}_{\text{adj}}^2 = \widehat{\sigma}^2 - \widehat{c}' \widehat{I}_{\theta^*}^{-1} \widehat{c}.$$

The corresponding estimator for the variance for the estimator for the average effect for the treated is

$$\widehat{\sigma}_{\text{adj},t}^2 = \widehat{\sigma}_t^2 - \widehat{c}_t' \widehat{I}_{\theta^*}^{-1} \widehat{c}_t + \frac{\widehat{\partial\tau_t}'}{\partial\theta} \widehat{I}_{\theta^*}^{-1} \frac{\partial\tau_t}{\partial\theta}.$$

Consistency of these estimators for fixed  $L$  can be shown using the results in [Abadie and Imbens \(2006\)](#) and the arguments employed in Section 3.

## 5. CONCLUSIONS AND EXTENSIONS

In this article, we derive the large sample distribution of propensity score matching estimators for the case where the propensity score is unknown and needs to be estimated in a first step prior to matching. We show that first step estimation of the propensity score generally affects the asymptotic variance of matching estimators, and derive adjustments for propensity score matching

estimators of ATE and ATET. These results allow, for the first time, valid large sample inference for estimators that use matching on the estimated propensity score.

For concreteness, we frame the article within the context of estimation of average treatment effects under the assumption that treatment assignment is independent of potential outcomes conditional on a set of covariates,  $X$  (Assumption 1(i)). Without this assumption, the results of this article apply still to the estimation of the “controlled comparison” parameters

$$E[E[Y|X, W = 1] - E[Y|X, W = 0]]$$

and

$$E[Y|W = 1] - E[E[Y|X, W = 0]|W = 1],$$

which have the same form as ATE and ATET parameters but lack a causal interpretation in the absence of Assumption 1(i). Controlled contrasts are the building blocks of Oaxaca–Blinder-type decompositions, commonly applied in economics (Oaxaca (1973), Blinder (1973), DiNardo, Fortin, and Lemieux (1996)). The ideas and results in this article can easily be applied to other contexts where it is required to adjust for differences in the distribution of covariates between two samples. An important example is estimation with missing data when missingness is random conditional on a set of covariates (see, e.g., Little and Rubin (2002), Wooldridge (2007)).

## APPENDIX

Before proving Proposition 2, we introduce some additional notation. Using the definition of  $K_{M,\theta}(i)$  in (8), the estimator  $\hat{\tau}_N(\theta)$  can be written as

$$\hat{\tau}_N(\theta) = \frac{1}{N} \sum_{i=1}^N (2W_i - 1) \left( 1 + \frac{K_{M,\theta}(i)}{M} \right) Y_i.$$

Define  $\bar{\mu}_\theta(w, p) = E_\theta[Y|W = w, F(X'\theta) = p]$ , where  $E_\theta$  is the expectation operator under  $P^\theta$ ,

$$\begin{aligned} D_N(\theta) &= \frac{1}{\sqrt{N}} \sum_{i=1}^N (\bar{\mu}_\theta(1, F(X'_i\theta)) - \bar{\mu}_\theta(0, F(X'_i\theta)) - \tau) \\ &\quad + \frac{1}{\sqrt{N}} \sum_{i=1}^N (2W_i - 1) \left( 1 + \frac{K_{M,\theta}(i)}{M} \right) (Y_i - \bar{\mu}_\theta(W_i, F(X'_i\theta))) \end{aligned}$$

and

$$R_N(\theta) = \frac{1}{\sqrt{N}} \sum_{i=1}^N (2W_i - 1) \\ \times \left( \bar{\mu}_\theta(1 - W_i, F(X'_i \theta)) - \frac{1}{M} \sum_{j \in \mathcal{J}_M(i)} \bar{\mu}_\theta(1 - W_j, F(X'_j \theta)) \right).$$

Now the normalized estimator can be written as

$$\sqrt{N}(\widehat{\tau}_N(\theta) - \tau) = D_N(\theta) + R_N(\theta).$$

PROOF OF PROPOSITION 2: It can be seen that the result of Lemma S.1 in the Supplemental Material ([Abadie and Imbens \(2016\)](#)) holds uniformly in  $\theta$  for  $\|\theta - \theta^*\| \leq \varepsilon$ . This implies  $R_N(\theta_N) \xrightarrow{p} 0$ . Therefore, in order to prove the result in the proposition, it suffices to prove that, under  $P^{\theta_N}$ ,

$$\begin{pmatrix} D_N(\theta_N) \\ \sqrt{N}(\widehat{\theta}_N - \theta_N) \\ \Lambda_N(\theta^* | \theta_N) \end{pmatrix} \xrightarrow{d} N \left( \begin{pmatrix} 0 \\ 0 \\ -h' I_{\theta^*} h / 2 \end{pmatrix}, \begin{pmatrix} \sigma^2 & c' I_{\theta^*}^{-1} & -c' h \\ I_{\theta^*}^{-1} c & I_{\theta^*}^{-1} & -h \\ -h' c & -h' & h' I_{\theta^*} h \end{pmatrix} \right).$$

By Lemma 1, under  $P^{\theta_N}$ ,

$$\Lambda_N(\theta^* | \theta_N) = -h' \Delta_N(\theta_N) - \frac{1}{2} h' I_{\theta^*} h + o_p(1)$$

and

$$\sqrt{N}(\widehat{\theta}_N - \theta_N) = I_{\theta^*}^{-1} \Delta_N(\theta_N) + o_p(1).$$

Therefore, it suffices to prove that, under  $P^{\theta_N}$ ,

$$(A.1) \quad \begin{pmatrix} D_N(\theta_N) \\ \Delta_N(\theta_N) \end{pmatrix} \xrightarrow{d} N \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma^2 & c' \\ c & I_{\theta^*} \end{pmatrix} \right).$$

To prove (A.1), we extend the martingale representation of matching estimators ([Abadie and Imbens \(2012\)](#)) to allow for estimation of the propensity

score. Consider the linear combination  $C_N = z_1 D_N(\theta_N) + z_2' \Delta_N(\theta_N)$ :

$$\begin{aligned} C_N = & z_1 \frac{1}{\sqrt{N}} \sum_{i=1}^N (\bar{\mu}_{\theta_N}(1, F(X'_{N,i} \theta_N)) - \bar{\mu}_{\theta_N}(0, F(X'_{N,i} \theta_N)) - \tau) \\ & + z_1 \frac{1}{\sqrt{N}} \sum_{i=1}^N (2W_{N,i} - 1) \left( 1 + \frac{K_{M, \theta_N}(i)}{M} \right) \\ & \times (Y_{N,i} - \bar{\mu}_{\theta_N}(W_{N,i}, F(X'_{N,i} \theta_N))) \\ & + z_2' \frac{1}{\sqrt{N}} \sum_{i=1}^N X_{N,i} \frac{W_{N,i} - F(X'_{N,i} \theta_N)}{F(X'_{N,i} \theta_N)(1 - F(X'_{N,i} \theta_N))} f(X'_{N,i} \theta_N). \end{aligned}$$

We analyze  $C_N$  using martingale methods. First, notice that

$$C_N = \sum_{k=1}^{3N} \xi_{N,k},$$

where

$$\begin{aligned} \xi_{N,k} = & z_1 \frac{1}{\sqrt{N}} (\bar{\mu}_{\theta_N}(1, F(X'_{N,k} \theta_N)) - \bar{\mu}_{\theta_N}(0, F(X'_{N,k} \theta_N)) - \tau) \\ & + z_2' \frac{1}{\sqrt{N}} E[X_{N,k} | F(X'_{N,k} \theta_N)] \\ & \times \frac{W_{N,k} - F(X'_{N,k} \theta_N)}{F(X'_{N,k} \theta_N)(1 - F(X'_{N,k} \theta_N))} f(X'_{N,k} \theta_N) \end{aligned}$$

for  $1 \leq k \leq N$ ,

$$\begin{aligned} \xi_{N,k} = & z_2' \frac{1}{\sqrt{N}} (X_{N,k-N} - E[X_{N,k-N} | F(X'_{N,k-N} \theta_N)]) \\ & \times \frac{(W_{N,k-N} - F(X'_{N,k-N} \theta_N)) f(X'_{N,k-N} \theta_N)}{F(X'_{N,k-N} \theta_N)(1 - F(X'_{N,k-N} \theta_N))} \\ & + z_1 \frac{1}{\sqrt{N}} (2W_{N,k-N} - 1) \left( 1 + \frac{K_{M, \theta_N}(k-N)}{M} \right) \\ & \times (\mu(W_{N,k-N}, X_{N,k-N}) - \bar{\mu}_{\theta_N}(W_{N,k-N}, F(X'_{N,k-N} \theta_N))) \end{aligned}$$

for  $N + 1 \leq k \leq 2N$ , and

$$\xi_{N,k} = z_1 \frac{1}{\sqrt{N}} (2W_{N,k-2N} - 1) \left( 1 + \frac{K_{M,\theta_N}(k-2N)}{M} \right) \\ \times (Y_{N,k-2N} - \mu(W_{N,k-2N}, X_{N,k-2N})),$$

for  $2N + 1 \leq k \leq 3N$ . Consider the  $\sigma$ -fields  $\mathcal{F}_{N,k} = \sigma\{W_{N,1}, \dots, W_{N,k}, X'_{N,1}\theta_N, \dots, X'_{N,k}\theta_N\}$  for  $1 \leq k \leq N$ ,  $\mathcal{F}_{N,k} = \sigma\{W_{N,1}, \dots, W_{N,N}, X'_{N,1}\theta_N, \dots, X'_{N,N}\theta_N, X_{N,1}, \dots, X_{N,k-N}\}$  for  $N + 1 \leq k \leq 2N$ , and  $\mathcal{F}_{N,k} = \sigma\{W_{N,1}, \dots, W_{N,N}, X_{N,1}, \dots, X_{N,N}, Y_{N,1}, \dots, Y_{N,k-N}\}$  for  $2N + 1 \leq k \leq 3N$ . Then,

$$\left\{ \sum_{j=1}^i \xi_{N,j}, \mathcal{F}_{N,i}, 1 \leq i \leq 3N \right\}$$

is a martingale for each  $N \geq 1$ . Therefore, the limiting distribution of  $C_N$  can be studied using a Martingale Central Limit Theorem (e.g., Theorem 35.12 in Billingsley (1995, p. 476); importantly, notice that this theorem allows that the probability space varies with  $N$ ). Because of Assumption 4, and because  $K_{M,\theta}(i)$  has uniformly bounded moments (see Abadie and Imbens (2016)), it follows that

$$\sum_{k=1}^{3N} E_{\theta_N} [|\xi_{N,k}|^{2+\delta}] \rightarrow 0 \quad \text{for some } \delta > 0.$$

Lindeberg's condition in Billingsley's theorem follows easily from the last equation (Lyapunov's condition). As a result, we obtain that, under  $P^{\theta_N}$ ,

$$C_N \xrightarrow{d} N(0, \sigma_1^2 + \sigma_2^2 + \sigma_3^2),$$

where

$$\sigma_1^2 = \text{plim} \sum_{k=1}^N E_{\theta_N} [\xi_{N,k}^2 | \mathcal{F}_{N,k-1}], \\ \sigma_2^2 = \text{plim} \sum_{k=N+1}^{2N} E_{\theta_N} [\xi_{N,k}^2 | \mathcal{F}_{N,k-1}], \\ \sigma_3^2 = \text{plim} \sum_{k=2N+1}^{3N} E_{\theta_N} [\xi_{N,k}^2 | \mathcal{F}_{N,k-1}].$$

Assumption 5 implies

$$\begin{aligned}\sigma_1^2 &= z_1^2 E[(\bar{\mu}(1, F(X'\theta^*)) - \bar{\mu}(0, F(X'\theta^*)) - \tau)^2] \\ &\quad + z_2^2 E\left[\frac{f^2(X'\theta^*)}{F(X'\theta^*)(1 - F(X'\theta^*))}\right. \\ &\quad \left. \times E[X|F(X'\theta^*)]E[X'|F(X'\theta^*)]\right] z_2.\end{aligned}$$

Expectations of the sums of terms involving  $(1 + K_{M, \theta_N}(i)/M)^2$  can be calculated as in [Abadie and Imbens \(2016\)](#). We obtain

$$\begin{aligned}\sigma_2^2 &= z_2^2 E\left[\frac{f^2(X'\theta^*)}{F(X'\theta^*)(1 - F(X'\theta^*))} \text{var}(X|F(X'\theta^*))\right] z_2 \\ &\quad + z_1^2 E\left[\frac{\text{var}(\mu(1, X)|F(X'\theta^*))}{F(X'\theta^*)} + \frac{\text{var}(\mu(0, X)|F(X'\theta^*))}{1 - F(X'\theta^*)}\right] \\ &\quad + z_1^2 \frac{1}{2M} E\left[\left(\frac{1}{F(X'\theta^*)} - F(X'\theta^*)\right) \text{var}(\mu(1, X)|F(X'\theta^*))\right] \\ &\quad + z_1^2 \frac{1}{2M} E\left[\left(\frac{1}{1 - F(X'\theta^*)} - (1 - F(X'\theta^*))\right)\right. \\ &\quad \left. \times \text{var}(\mu(0, X)|F(X'\theta^*))\right] \\ &\quad + 2z_2^2 E\left[\left(\frac{\text{cov}(X, \mu(1, X)|F(X'\theta^*))}{F(X'\theta^*)}\right.\right. \\ &\quad \left. \left. + \frac{\text{cov}(X, \mu(0, X)|F(X'\theta^*))}{1 - F(X'\theta^*)}\right) f(X'\theta^*)\right] z_1.\end{aligned}$$

Here we use the fact that, conditional on the propensity score,  $X$  is independent of  $W$ . Finally, notice that

$$\begin{aligned}&\frac{1}{N} \sum_{i=1}^N \left(1 + \frac{K_{M, \theta_N}(i)}{M}\right)^2 \\ &\quad \times (\text{var}(Y_i|W_i, X_i) - E[\text{var}(Y_i|W_i, X_i)|W_i, F(X'_i\theta_N)])\end{aligned}$$



is a sum of martingale differences

$$\begin{aligned}\zeta_{N,i} &= \frac{1}{N} \left( 1 + \frac{K_{M,\theta_N}(i)}{M} \right)^2 \\ &\quad \times (\text{var}(Y_i|W_i, X_i) - E[\text{var}(Y_i|W_i, X_i)|W_i, F(X'_i\theta_N)])\end{aligned}$$

with respect to the filtration  $\mathcal{F}_{N,i} = \sigma\{W_1, \dots, W_N, X'_1\theta_N, \dots, X'_i\theta_N, X_1, \dots, X_i\}$ . As a result, we obtain that, for  $i > j$ ,  $E[\zeta_{N,i}\zeta_{N,j}|\mathcal{F}_{N,i-1}] = E[\zeta_{N,i}|\mathcal{F}_{N,i-1}]\zeta_{N,j} = 0$ . Therefore, using the Law of Iterated Expectations to eliminate the cross-products, we obtain

$$E\left[\left(\frac{1}{N} \sum_{i=1}^N \zeta_{N,i}\right)^2\right] = \frac{1}{N^2} E\left[\sum_{i=1}^N \zeta_{N,i}^2\right] \rightarrow 0.$$

Therefore,

$$\begin{aligned}\sigma_3^2 &= z_1^2 \text{plim} \frac{1}{N} \sum_{i=1}^N \left( 1 + \frac{K_{M,\theta_N}(i)}{M} \right)^2 E[\text{var}(Y_i|W_i, X_i)|W_i, F(X'_i\theta_N)] \\ &= z_1^2 E\left[\frac{E[\text{var}(Y|X, W=1)|F(X'\theta^*)]}{F(X'\theta^*)} \right. \\ &\quad \left. + \frac{E[\text{var}(Y|X, W=0)|F(X'\theta^*)]}{1-F(X'\theta^*)} \right] \\ &\quad + z_1^2 \frac{1}{2M} E\left[\left(\frac{1}{F(X'\theta^*)} - F(X'\theta^*)\right) \right. \\ &\quad \left. \times E[\text{var}(Y|X, W=1)|F(X'\theta^*)] \right] \\ &\quad + z_1^2 \frac{1}{2M} E\left[\left(\frac{1}{1-F(X'\theta^*)} - (1-F(X'\theta^*))\right) \right. \\ &\quad \left. \times E[\text{var}(Y|X, W=0)|F(X'\theta^*)] \right].\end{aligned}$$

Collecting terms and applying the fact that  $W$  is independent of  $X$  given  $F(X'\theta)$ , we obtain

$$\sigma_1^2 + \sigma_2^2 + \sigma_3^2 = z_1^2 \sigma^2 + z_2' I_{\theta^*} z_2 + 2z_2' c z_1.$$

Hence, by the Martingale Central Limit Theorem and the Cramer–Wold device, under  $P^{\theta_N}$ ,

$$\begin{pmatrix} D_N(\theta_N) \\ \Delta_N(\theta_N) \end{pmatrix} \xrightarrow{d} N \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma^2 & c' \\ c & I_{\theta^*} \end{pmatrix} \right),$$

proving (A.1) and thus Proposition 2.

*Q.E.D.*

PROOF OF THEOREM 1: Given our preliminary results, Theorem 1 follows from Andreou and Werker (2012).

*Q.E.D.*

The proof of Theorem 2 can be found in the Supplemental Material.

## REFERENCES

- ABADIE, A. (2005): “Semiparametric Difference-in-Differences Estimators,” *Review of Economic Studies*, 72 (1), 1–19. [782]
- ABADIE, A., AND G. W. IMBENS (2006): “Large Sample Properties of Matching Estimators for Average Treatment Effects,” *Econometrica*, 74 (1), 235–267. [781,784,785,789,795,796,799]
- (2008): “On the Failure of the Bootstrap for Matching Estimators,” *Econometrica*, 76 (6), 1537–1557. [782]
- (2011): “Bias-Corrected Matching Estimators for Average Treatment Effects,” *Journal of Business and Economic Statistics*, 29 (1), 1–11. [784,787]
- (2012): “A Martingale Representation for Matching Estimators,” *Journal of the American Statistical Association*, 107 (498), 833–843. [788,791,793,801]
- (2016): “Supplement to ‘Matching on the Estimated Propensity Score’,” *Econometrica Supplemental Material*, 84, <http://dx.doi.org/10.3982/ECTA11293>. [786,790,801,803,804]
- ANDREOU, E., AND B. J. M. WERKER (2012): “An Alternative Asymptotic Analysis of Residual-Based Statistics,” *Review of Economics and Statistics*, 94 (1), 88–99. [788,791,792,806]
- ANGRIST, J. D. (1998): “Estimating the Labor Market Impact of Voluntary Military Service Using Social Security Data on Military Applicants,” *Econometrica*, 66 (2), 249–288. [789]
- ANGRIST, J. D., AND G. M. KUERSTEINER (2011): “Causal Effects of Monetary Shocks: Semiparametric Conditional Independence Tests With a Multinomial Propensity Score,” *Review of Economics and Statistics*, 93 (3), 725–747. [782]
- BICKEL, P. J., C. A. KLAASSEN, Y. RITOV, AND J. A. WELLNER (1998): *Efficient and Adaptive Estimation for Semiparametric Models*. New York: Springer. [790]
- BILLINGSLEY, P. (1995): *Probability and Measure*. New York: Wiley. [803]
- BLINDER, A. S. (1973): “Wage Discrimination: Reduced Form and Structural Estimates,” *Journal of Human Resources*, 8 (4), 436–455. [800]
- COCHRAN, W. G. (1968): “The Effectiveness of Adjustment by Subclassification in Removing Bias in Observational Studies,” *Biometrics*, 24 (2), 295–313. [789]
- DAWID, A. P. (1979): “Conditional Independence in Statistical Theory,” *Journal of the Royal Statistical Society, Series B*, 41 (1), 1–31. [783]
- DEHEJIA, R., AND S. WAHBA (1999): “Causal Effects in Nonexperimental Studies: Reevaluating the Evaluation of Training Programs,” *Journal of the American Statistical Association*, 94 (448), 1053–1062. [782,787]
- DI NARDO, J., N. M. FORTIN, AND T. LEMIEUX (1996): “Labor Market Institutions and the Distribution of Wages, 1973–1992: A Semiparametric Approach,” *Econometrica*, 64 (5), 1001–1044. [800]
- GOLDBERGER, A. S. (1991): *A Course in Econometrics*. Cambridge, MA: Harvard University Press. [792]

- HAHN, J. (1998): "On the Role of the Propensity Score in Efficient Semiparametric Estimation of Average Treatment Effects," *Econometrica*, 66 (2), 315–331. [781,783]
- HECKMAN, J., H. ICHIMURA, AND P. TODD (1997): "Matching as an Econometric Evaluation Estimator: Evidence From a Job Training Programme," *Review of Economic Studies*, 64 (4), 605–654. [782]
- (1998): "Matching as an Econometric Evaluation Estimator," *Review of Economic Studies*, 65 (2), 261–294. [781,782]
- HIRANO, K., G. W. IMBENS, AND G. RIDDER (2003): "Efficient Estimation of Average Treatment Effects Using the Estimated Propensity Score," *Econometrica*, 71 (4), 1161–1189. [782]
- IMBENS, G. (2004): "Nonparametric Estimation of Average Treatment Effects Under Exogeneity: A Review," *Review of Economics and Statistics*, 86 (1), 1–29. [781]
- IMBENS, G., AND J. WOOLDRIDGE (2009): "Recent Developments in the Econometrics of Program Evaluation," *Journal of Economic Literature*, 47 (1), 5–86. [781]
- KHAN, S., AND E. TAMER (2010): "Irregular Identification, Support Conditions, and Inverse Weight Estimation," *Econometrica*, 78 (6), 2021–2042. [783]
- KIEFER, N. M., AND T. J. VOGELSANG (2005): "A New Asymptotic Theory for Heteroskedasticity-Autocorrelation Robust Tests," *Econometric Theory*, 21, 1130–1164. [784]
- LEHMANN, E. L., AND J. P. ROMANO (2005): *Testing Statistical Hypothesis*. New York: Springer. [790]
- LITTLE, R. J., AND D. B. RUBIN (2002): *Statistical Analysis With Missing Data* (Second Ed.). Hoboken, NJ: Wiley-Interscience. [800]
- NEWWEY, W. K., AND D. MCFADDEN (1994): "Large Sample Estimation and Hypothesis Testing," in *Handbook of Econometrics*, Vol. 4, ed. by R. F. Engle and D. McFadden. Amsterdam: Elsevier Science. [782]
- OAXACA, R. (1973): "Male–Female Wage Differentials in Urban Labor Markets," *International Economic Review*, 14 (3), 693–709. [800]
- ROSENBAUM, P., AND D. B. RUBIN (1983): "The Central Role of the Propensity Score in Observational Studies for Causal Effects," *Biometrika*, 70 (1), 41–55. [781,783,784,786]
- RUBIN, D. B. (1974): "Estimating Causal Effects of Treatments in Randomized and Non-Randomized Studies," *Journal of Educational Psychology*, 66 (5), 688–701. [783,784]
- SMITH, J., AND P. TODD (2005): "Does Matching Overcome LaLonde's Critique of Nonexperimental Estimators?" *Journal of Econometrics*, 125 (1–2), 305–353. [782]
- VAN DER VAART, A. (1998): *Asymptotic Statistics*. New York: Cambridge University Press. [790, 791]
- WOOLDRIDGE, J. M. (2007): "Inverse Probability Weighted Estimation for General Missing Data Problems," *Journal of Econometrics*, 141 (2), 1281–1301. [782,800]
- YATCHEW, A. (1997): "An Elementary Estimator of the Partial Linear Model," *Economics Letters*, 75, 135–143. [784]

*John F. Kennedy School of Government, 79 John F. Kennedy Street, Cambridge, MA 02138, U.S.A. and NBER; alberto\_abadie@harvard.edu*

*and*

*Stanford Graduate School of Business, 655 Knight Way, Stanford, CA 94305-7298, U.S.A. and NBER; imbens@stanford.edu.*

*Co-editor Elie Tamer handled this manuscript.*

*Manuscript received December, 2012; final revision received August, 2015.*