

# EFFICIENT ESTIMATION OF AVERAGE TREATMENT EFFECTS USING THE ESTIMATED PROPENSITY SCORE

BY KEISUKE HIRANO, GUIDO W. IMBENS, AND GEERT RIDDER<sup>1</sup>

We are interested in estimating the average effect of a binary treatment on a scalar outcome. If assignment to the treatment is exogenous or unconfounded, that is, independent of the potential outcomes given covariates, biases associated with simple treatment-control average comparisons can be removed by adjusting for differences in the covariates. Rosenbaum and Rubin (1983) show that adjusting solely for differences between treated and control units in the propensity score removes all biases associated with differences in covariates. Although adjusting for differences in the propensity score removes all the bias, this can come at the expense of efficiency, as shown by Hahn (1998), Heckman, Ichimura, and Todd (1998), and Robins, Mark, and Newey (1992). We show that weighting by the inverse of a nonparametric estimate of the propensity score, rather than the true propensity score, leads to an efficient estimate of the average treatment effect. We provide intuition for this result by showing that this estimator can be interpreted as an empirical likelihood estimator that efficiently incorporates the information about the propensity score.

**KEYWORDS:** Propensity score, treatment effects, semiparametric efficiency, sieve estimator.

## 1. INTRODUCTION

ESTIMATING THE AVERAGE EFFECT of a binary treatment or policy on a scalar outcome is a basic goal of many empirical studies in economics. If assignment to the treatment is exogenous or unconfounded (i.e., independent of potential outcomes conditional on covariates or pre-treatment variables, an assumption also known as selection on observables), the average treatment effect can be estimated by matching<sup>2</sup> or by averaging within-subpopulation differences of treatment and control averages. If there are many covariates, such strategies may not be desirable or even feasible. An alternative approach is based on the propensity score, the conditional probability of receiving treatment given covariates. Rosenbaum and Rubin (1983, 1985) show that, under the assumption of unconfoundedness, adjusting solely for differences in the propensity score between treated and control units removes all biases. Recent applications of propensity score methods

<sup>1</sup> We thank Gary Chamberlain, Jinyong Hahn, James Robins, Donald Rubin, Jeffrey Wooldridge, four anonymous referees, seminar participants at the University of Chicago, UC Davis, the University of Michigan, Michigan State University, UC Irvine, the University of Miami, Johns Hopkins, and Harvard-MIT, and especially Whitney Newey for comments. Financial support for this research was generously provided through NSF Grants SBR-9818644 and SES-0136789 (Imbens) and SES-9985257 (Hirano).

<sup>2</sup> See Abadie and Imbens (2002) for a formal discussion of matching estimators in this context.

in economics include Dehejia and Wahba (1999), Heckman, Ichimura, and Todd (1997), and Lechner (1999).

Although adjusting for differences in the propensity score removes all bias, it need not be as efficient as adjusting for differences in all covariates, as shown by Hahn (1998), Heckman, Ichimura, and Todd (1998), and Robins, Mark, and Newey (1992). However, Rosenbaum (1987), Rubin and Thomas (1996), and Robins, Rotnitzky, and Zhao (1995) show that using parametric estimates of the propensity score, rather than the true propensity score, can avoid some of these efficiency losses.

In this paper we propose estimators that are based on adjusting for nonparametric estimates of the propensity score. The proposed estimators weight observations by the inverse of nonparametric estimates of the propensity score, rather than the true propensity score. Extending results from Newey (1994) to derive the large sample properties of these semiparametric estimators, we show that they achieve the semiparametric efficiency bound. We also show that for the case in which the propensity score is known the proposed estimators can be interpreted as empirical likelihood estimators (e.g., Qin and Lawless (1994), Imbens, Spady, and Johnson (1998)) that efficiently incorporate the information about the propensity score.

Our proposed estimators are relevant whether the propensity score is known or not. In randomized experiments, for example, the propensity score is known by design. In that case the proposed estimators can be used to improve efficiency over simply differencing treatment and control averages. With the propensity score known, an attractive choice for the nonparametric series estimator for the propensity score is to use the true propensity score as the leading term in the series. The proposed estimators can also be used in the case where the propensity score is unknown. In that case they are alternatives to the previously proposed efficient estimators that require nonparametric estimation of functions in addition to the propensity score.

In the next section we lay out the problem and discuss the prior literature. In Section 3 we provide some intuition for our efficiency results by examining a simplified version of the problem. In Section 4 we give the formal conditions under which weighting by the estimated propensity score results in an efficient estimator. Section 5 concludes.

## 2. THE BASIC SETUP AND PREVIOUS RESULTS

### 2.1. *The Model*

We have a random sample of size  $N$  from a large population. For each unit  $i$  in the sample, for  $i = 1, \dots, N$ , let  $T_i$  indicate whether the treatment of interest was received, with  $T_i = 1$  if unit  $i$  receives the active treatment, and  $T_i = 0$  if unit  $i$  receives the control treatment. Using the potential outcome notation popularized by Rubin (1974), let  $Y_i(0)$  denote the outcome for each unit  $i$  under control

and  $Y_i(1)$  the outcome under treatment.<sup>3</sup> We observe  $T_i$  and  $Y_i$ , where  $Y_i \equiv T_i \cdot Y_i(1) + (1 - T_i) \cdot Y_i(0)$ . In addition, we observe a vector of covariates denoted by  $X_i$ .<sup>4</sup> Initially we focus on the population average treatment effect:

$$(1) \quad \tau \equiv \mathbf{E}[Y(1) - Y(0)].$$

We shall also discuss estimation of weighted average treatment effects,

$$(2) \quad \tau_{\text{wate}} \equiv \frac{\int \mathbf{E}[Y(1) - Y(0)|X = x]g(x)dF(x)}{\int g(x)dF(x)},$$

where  $g(\cdot)$  is a known function of the covariates.<sup>5</sup> In the special case where the weight function  $g(x)$  is equal to the propensity score  $p(x) = \Pr(T = 1|X = x)$ , this leads under the unconfoundedness assumption to the average effect for the treated:

$$(3) \quad \tau_{\text{treated}} \equiv \mathbf{E}[Y(1) - Y(0)|T = 1].$$

The central problem of evaluation research is that for unit  $i$  we observe either  $Y_i(0)$  or  $Y_i(1)$ , but never both. To solve the identification problem, we maintain throughout the paper the unconfoundedness assumption (Rubin (1978), Rosenbaum and Rubin (1983)), related to the selection-on-observables assumption (Barnow, Cain, and Goldberger (1980)), which asserts that conditional on the observed covariates, the treatment indicator is independent of the potential outcomes. Formally, we have the following assumption:

ASSUMPTION 1 (Unconfounded Treatment Assignment):

$$T \perp (Y(0), Y(1)) \mid X.$$

Heckman, Ichimura, and Todd (1998) point out that for identification of the average treatment effect  $\tau$  this assumption can be weakened to mean independence ( $\mathbf{E}[Y(t)|T, X] = \mathbf{E}[Y(t)|X]$  for  $t = 0, 1$ ). If one is interested in the average effect for the treated, the assumption can be further weakened to only require

<sup>3</sup> Implicit in this notation is the stability assumption or SUTVA (Rubin (1978)) that units are not affected by receipt of treatment by others, and that there is only one version of the treatment.

<sup>4</sup> These variables are assumed not to be affected by the treatment.

<sup>5</sup> An alternative estimand, which we do not consider here, is the direct weighted average treatment effect of the form

$$\tau_{\text{dwate}} = \frac{\int \mathbf{E}[Y(1) - Y(0)|X = x]g(x)dx}{\int g(x)dx},$$

where the weighting is only over the known function  $g(x)$ . Note that in general  $F(x)$  is unknown so that knowledge of  $g(x)$  does not imply knowledge of  $g(x)dF(x)$  and the other way around; estimation strategies for the two estimands  $\tau_{\text{wate}}$  and  $\tau_{\text{dwate}}$  are in general different. Estimands of the latter type can be fitted into the framework of Robins and Ritov (1997).

that  $\mathbf{E}[Y(0)|T, X] = \mathbf{E}[Y(0)|X]$ . In this paper we focus on the full independence assumption in order to be consistent with much of the literature.

Under unconfoundedness we can estimate the average treatment effect conditional on covariates,  $\tau(x) \equiv \mathbf{E}[Y(1) - Y(0)|X = x]$ , because

$$\begin{aligned}\tau(x) &= \mathbf{E}[Y(1) - Y(0)|X = x] \\ &= \mathbf{E}[Y(1)|T = 1, X = x] - \mathbf{E}[Y(0)|T = 0, X = x] \\ &= \mathbf{E}[Y|T = 1, X = x] - \mathbf{E}[Y|T = 0, X = x].\end{aligned}$$

The population average treatment effect can then be obtained by averaging the  $\tau(x)$  over the distribution of  $X$ :  $\tau = \mathbf{E}[\tau(X)]$ . In practice, the strategy of forming cells and comparing units with exactly the same value of  $X$  may fail if  $X$  takes on too many distinct values.<sup>6</sup> To avoid the need to match units on the values of all covariates, Rosenbaum and Rubin (1983, 1985) developed an approach based on the propensity score, the probability of selection into the treatment group:

$$(4) \quad p(x) \equiv \Pr(T = 1|X = x) = \mathbf{E}[T|X = x],$$

which is assumed to be bounded away from zero and one. Their key insight was that if treatment and potential outcomes are independent conditional on all covariates, they are also independent conditional on the conditional probability of receiving treatment given covariates. Formally, as shown by Rosenbaum and Rubin (1983), unconfoundedness implies

$$(5) \quad T \perp (Y(0), Y(1)) \mid p(X),$$

implying that adjustment for the propensity score suffices for removing all biases associated with differences in the covariates.

## 2.2. Previous Results

The model set out above, as well as related models, have been examined by a number of researchers. In an important paper Hahn (1998), studying the same model as we do here, calculates the semiparametric efficiency bounds, and proposes efficient estimators, for  $\tau$  and  $\tau_{treated}$ . Hahn's estimator for  $\tau$ , which is efficient irrespective of whether the propensity score is known, nonparametrically estimates the two conditional expectations  $\mathbf{E}[YT|X = x]$  and  $\mathbf{E}[Y(1 - T)|X = x]$  as well as the propensity score  $p(x)$ , and then imputes the missing potential outcomes as  $\hat{Y}_i(1) = \hat{\mathbf{E}}[YT|X_i]/\hat{p}(X_i)$  and  $\hat{Y}_i(0) = \hat{\mathbf{E}}[Y(1 - T)|X_i]/(1 - \hat{p}(X_i))$ . Hahn shows that conditioning only on the true propensity score rather than on the

<sup>6</sup> A related issue is whether standard asymptotic theory provides adequate approximations to the sampling distributions of estimators based on initial nonparametric estimates of conditional means, especially when the dimension of the conditioning variable is high. For discussions of these issues, see Robins and Ritov (1997) and Angrist and Hahn (1999) and references therein.

full set of covariates does not in general lead to an efficient estimator. In addition Hahn concludes that knowledge of the propensity score is informative for estimating  $\tau_{treated}$  and derives efficient estimators both with and without such knowledge. A difference between Hahn's estimators and our proposed estimators is that Hahn requires nonparametric estimation of the propensity score as well as the two conditional means  $E[YT|X=x]$  and  $E[Y(1-T)|X=x]$ , whereas our proposed estimator only requires nonparametric estimation of the propensity score.

Heckman, Ichimura, and Todd (1997, 1998) and Heckman, Ichimura, Smith, and Todd (1998) focus on  $\tau_{treated}$ , the average treatment effect for the treated. They consider estimators based on local linear regressions of the outcome on treatment status and either covariates or the propensity score. They conclude that in general there is no clear ranking of their estimators; under some conditions the estimator based on adjustment for all covariates is superior to the estimator based on adjustment for the propensity score, and under other conditions the second estimator is to be preferred. Lack of knowledge of the propensity score does not alter this conclusion.

Rosenbaum (1987) and Rubin and Thomas (1996) investigate the differences between using the estimated and the true propensity score when the propensity score belongs to a parametric family. They conclude that there can be efficiency gains from using the estimated propensity score. Our results show that by making the specification of the propensity score sufficiently flexible, this approach leads to a fully efficient estimator.

Robins, Mark, and Newey (1992), Robins and Rotnitzky (1995), Robins, Rotnitzky, and Zhao (1995), and Rotnitzky and Robins (1995) study the related problem of inference for parameters in regression models where some data are Missing At Random (MAR; Rubin (1976), Little and Rubin (1987)). Rotnitzky and Robins (1995) show that in parametric settings weighting using the estimated rather than true selection probability can improve efficiency. They suggest it may be possible to achieve full efficiency by allowing the dimension of the model for the selection probability to grow with the sample size. For this missing data case Robins and Rotnitzky (1995) also propose an efficient estimator that relies on an initial consistent, but not necessarily efficient, estimator of the full population parameters. The estimator we propose is efficient (as is the estimator proposed by Hahn), but does not require an initial consistent estimator.

### 3. A SIMPLE EXAMPLE WITH BINARY COVARIATES

To develop some intuition for the formal results that will be presented in Section 4, we consider the simpler problem of estimating the population average of a variable  $Y$ ,  $\beta_0 = E[Y]$ , given a random sample of size  $N$  of the triple  $(T_i, X_i, T_i \cdot Y_i)$ . In other words,  $T_i$  and  $X_i$  are observed for all units in the sample, but  $Y_i$  is only observed if  $T_i = 1$ . We provide a heuristic argument for efficiency of using estimated weights, deferring formal results to Section 4.

The analog to the unconfoundedness assumption here is the assumption that the  $Y_i$  are Missing At Random (MAR; Rubin (1976)), or

$$T \perp Y \mid X.$$

The role of the propensity score is played here by the selection probability:  $p(x) = \mathbf{E}[T|X = x] = \Pr(T = 1|X = x)$ . First, we restrict our attention in this section to the case with a single binary covariate.<sup>7</sup> Let  $N_{tx}$  denote the number of observations with  $T_i = t$  and  $X_i = x$ , for  $t, x \in \{0, 1\}$ . Furthermore, suppose the true selection probability is constant,  $p(x) = 1/2$  for all  $x \in \{0, 1\}$ .<sup>8</sup> The normalized variance bound for  $\beta_0$  is

$$(6) \quad V_{bound} = 2 \cdot \mathbf{E}[\mathbf{V}(Y|X)] + \mathbf{V}(\mathbf{E}[Y|X]).$$

The “true weights” estimator weights the complete observations by the inverse of the true selection probability:

$$(7) \quad \hat{\beta}_{tw} = \frac{1}{N} \sum_{i=1}^N \frac{Y_i \cdot T_i}{p(X_i)} = \frac{1}{N} \sum_{i=1}^N \frac{Y_i \cdot T_i}{1/2}.$$

Its large sample normalized variance is

$$V_{tw} = 2 \cdot \mathbf{E}[\mathbf{V}(Y|X)] + \mathbf{V}(\mathbf{E}[Y|X]) + \mathbf{E}[\mathbf{E}[Y|X]^2] = V_{bound} + \mathbf{E}[\mathbf{E}[Y|X]^2],$$

strictly larger than the variance bound (6) unless  $\mathbf{E}[Y|X] = 0$ .

The second estimator weights the complete observations by the inverse of a nonparametric estimate of the selection probability. This estimator is the main focus of the paper and it will be discussed in Section 4 in more general settings. In the current setting the estimated selection probability is simply the proportion of observed outcomes for a given value of the covariate. For units with  $X_i = 0$  the proportion of observed outcomes is  $N_{10}/(N_{00} + N_{10})$ , and for units with  $X_i = 1$  the proportion of observed outcomes is  $N_{11}/(N_{01} + N_{11})$ . Thus the estimated selection probability is

$$\hat{p}(x) = \begin{cases} N_{10}/(N_{00} + N_{10}) & \text{if } x = 0, \\ N_{11}/(N_{01} + N_{11}) & \text{if } x = 1. \end{cases}$$

The proposed “estimated weights” estimator is then

$$(8) \quad \hat{\beta}_{ew} = \frac{1}{N} \sum_{i=1}^N \frac{Y_i \cdot T_i}{\hat{p}(X_i)}.$$

<sup>7</sup> An efficient estimator is easily obtained by averaging the within-subsample difference of treatment/control averages. It can also be found by specializing the more general estimators in Robins and Rotnitzky (1995) and Hahn (1998) to this simple case. The discussion here is solely intended to convey intuition for the formal results that will be presented in Section 4.

<sup>8</sup> Thus the missing data are Missing Completely At Random (MCAR; Rubin (1976), Little and Rubin (1987)).

The normalized variance of this estimator is equal to the variance bound:

$$V_{ew} = 2 \cdot \mathbf{E}[\mathbf{V}(Y|X)] + \mathbf{V}(\mathbf{E}[Y|X]) = V_{bound}.$$

Not only does the weighting estimator with nonparametrically estimated weights have a lower variance than the estimator using the “true” weights in this simple case, but it is in fact fully efficient. In the remainder of this section we shall provide some intuition for this result. This will suggest why this efficiency property may carry over to the case with continuous and vector-valued covariates, as well as with general dependence of the selection probability or propensity score on the covariates.

An alternative interpretation of the estimated-weights estimator is based on a Generalized Method of Moments (GMM) representation (Hansen (1982)). Under the assumption that the selection probability is  $p(x) = 1/2$ , we can estimate  $\beta_0$  using the single moment restriction  $\mathbf{E}[\psi_1(Y, X, T, \beta_0)] = 0$ , with

$$\psi_1(y, t, x, \beta) = \frac{y \cdot t}{p(x)} - \beta = \frac{y \cdot t}{1/2} - \beta.$$

The GMM estimator based on the single moment restriction  $\psi_1(\cdot)$ , given knowledge of the selection probability, is the true-weights estimator  $\hat{\beta}_{tw}$  in (7). However, this estimator is not necessarily efficient, because it ignores the additional information that is available in the form of knowledge of the selection probability. This additional information can be written in moment condition form as  $\mathbf{E}[T - p(X)|X] = \mathbf{E}[T - 1/2|X] = 0$ . With a binary covariate this conditional moment restriction corresponds to two marginal moment restrictions,  $\mathbf{E}[\psi_2(Y, T, X, \beta_0)] = 0$ , with:

$$\psi_2(y, t, x, \beta) = \begin{pmatrix} x \cdot (t - 1/2) \\ (1 - x) \cdot (t - 1/2) \end{pmatrix}.$$

Estimating  $\beta_0$  in a generalized method of moments framework using the moments  $\psi_1(\cdot)$  and  $\psi_2(\cdot)$  leads to a fully efficient estimator.<sup>9</sup> Here it is of particular interest to consider the empirical likelihood estimator (e.g., Qin and Lawless (1994), Imbens (1997), Kitamura and Stutzer (1997), Imbens, Spady, and Johnson (1998)). Empirical likelihood estimation is based on maximization, both over a nuisance parameter  $\pi = (\pi_1, \dots, \pi_N)$  and over the parameter of interest  $\beta$ , of the logarithm of the empirical likelihood function:

$$(9) \quad L(\pi) = \sum_{i=1}^N \ln \pi_i,$$

<sup>9</sup> Although  $\psi_2(\cdot)$  does not depend on the parameter of interest,  $\psi_2(\cdot)$  is generally correlated with  $\psi_1(\cdot)$ . Thus there can be efficiency gains from using both sets of moment conditions as in seemingly unrelated regressions. See, e.g., Hellerstein and Imbens (1999) and Qian and Schmidt (1999).

subject to the adding-up restriction  $\sum_i \pi_i = 1$  and the moment conditions  $\sum_i \pi_i \psi(y_i, t_i, x_i, \beta) = 0$ . Solving for  $\hat{\pi}_i$  and  $\hat{\beta}_{el}$  by maximizing (9) subject to the restrictions leads, after some manipulation, to:

$$\hat{\pi}_i = \left( 1 + \frac{\frac{N_{11}}{N_{01}+N_{11}} - 1/2}{1/4} \cdot x_i \cdot (t_i - 1/2) + \frac{\frac{N_{10}}{N_{00}+N_{10}} - 1/2}{1/4} \cdot (1 - x_i) \cdot (t_i - 1/2) \right)^{-1},$$

which in turn implies

$$\hat{\beta}_{el} = \sum_{i=1}^N 2 \cdot \hat{\pi}_i \cdot Y_i \cdot T_i = \hat{\beta}_{ew},$$

equal to the estimated weights estimator.

The above discussion generalizes directly to the case with general discrete covariates. With continuous covariates knowledge of the propensity score implies a conditional moment restriction corresponding to an infinite number of unconditional moment restrictions (e.g., Chamberlain (1987)). Using a series estimator for the propensity score captures the information content of such a conditional moment restriction by a sequence of unconditional moment restrictions.

The empirical likelihood interpretation suggests that moving from the true-weights estimator to the estimated-weights estimator increases efficiency in the same way that adding moment restrictions in a generalized method of moments framework improves efficiency. A similar finding appears in Crepon, Kramarz, and Trognon (1998) who find that using a reduced set of moment conditions, in which nuisance parameters are replaced by solutions to the sample analogs of the remaining moment conditions, is asymptotically equivalent to using the full set of moment conditions, whereas using the true values of the nuisance parameters may lead to efficiency losses. These results are also linked to the literature on weighting in stratified sampling. Translated to this simple example, the results by Lancaster (1990) suggest studying the distribution of the various estimators conditional on the ancillary statistics  $\sum T_i$ ,  $\sum X_i$ , and  $\sum T_i \cdot X_i$ . Conditional on those three statistics the true-weights estimator is biased, while the estimated-weights estimator remains unbiased. Rosenbaum (1987) discusses this issue specifically in the context of estimated versus true propensity scores. In a general discussion of weighted M-estimators, Wooldridge (1999, 2002) shows that weighting by the inverse of estimated rather than population probabilities can lead to efficiency gains.

#### 4. EFFICIENT ESTIMATION USING ESTIMATED WEIGHTS

In this section we present the main results of the paper. We discuss three distinct cases. First, we consider the problem of estimating the population average treatment effect under the unconfoundedness assumption. This includes as



a special case the extension of the binary-covariate MAR example of the previous section to continuous covariates. Second, we consider estimation of weighted average treatment effects. Finally, we consider estimation of the effect of the treatment on the treated which, in the known propensity score case, will follow directly from the solution to the general weighted average treatment effect case. This discussion will shed additional light on Hahn's (1998) interesting result that for this parameter knowledge of the propensity score affects the efficiency bound, as well as on the findings in Heckman, Ichimura, and Todd (1998) that in the case of the average treatment effect for the treated, neither using the true nor using the estimated propensity score dominates the other.

#### 4.1. Estimating Population Average Treatment Effects

In this section we use the set up from Section 2 with a pair of potential outcomes  $(Y(0), Y(1))$  for each unit and focus on efficient estimation of the population average treatment effect,  $\tau^* = \mathbf{E}[Y(1) - Y(0)]$ .<sup>10</sup> As before,  $p(x) = \Pr(T = 1|X = x)$  is the propensity score, the probability of receiving the active treatment. We maintain the unconfoundedness assumption. Define  $\mu_t(x) \equiv \mathbf{E}[Y(t)|X = x]$  and  $\sigma_t^2(x) = \mathbf{V}(Y(t)|X = x)$  to be the conditional mean and variance of  $Y(t)$  respectively. Under unconfoundedness we have  $\mu_t(x) = \mathbf{E}[Y|T = t, X = x]$  and  $\sigma_t^2(x) = \mathbf{V}(Y|T = t, X = x)$ . We can characterize  $\tau^*$  through the moment equation:

$$\mathbf{E}[\psi(Y, T, X, \tau^*, p^*(X))] = 0,$$

where

$$(10) \quad \psi(y, t, x, \tau, p(x)) = \frac{y \cdot t}{p(x)} - \frac{y \cdot (1-t)}{1-p(x)} - \tau.$$

Given an estimator  $\hat{p}(x)$  for the propensity score, we estimate  $\tau^*$  by setting the average moment evaluated at the estimated selection probability equal to zero as a function of  $\tau$ :  $(1/N) \sum_{i=1}^N \psi(Y_i, T_i, X_i, \hat{\tau}, \hat{p}(X_i)) = 0$ , leading to the estimator

$$(11) \quad \hat{\tau} = \frac{1}{N} \sum_{i=1}^N \left( \frac{Y_i \cdot T_i}{\hat{p}(X_i)} - \frac{Y_i \cdot (1-T_i)}{1-\hat{p}(X_i)} \right).$$

Because  $p^*(x)$  is a conditional expectation, this semiparametric estimation problem directly fits into the framework of Newey (1994). So we could apply his results directly if we estimate  $p^*(x)$  by a series of least squares regressions of treatment on polynomials in the covariates. (See the working paper version, Hirano, Imbens, and Ridder (2000).) However, because  $p^*(x)$  is a probability, such an approach has the unattractive feature that it approximates a probability by a linear function. We therefore estimate  $p^*(x)$  in a sieve approach

<sup>10</sup> Whenever necessary to avoid confusion we will use a superscript  $*$  to denote true (population) values, so that  $\tau^*$  denotes the population average treatment effect and  $p^*(x)$  the true (population) propensity score.

(e.g., Geman and Hwang (1982)) by the Series Logit Estimator (SLE). For  $K = 1, 2, \dots$ , let  $R^K(x) = (r_{1K}(x), r_{2K}(x), \dots, r_{KK}(x))'$  be a  $K$ -vector of functions. Although the theory is derived for general sequences of approximating functions, the most common class of functions are power series. Let  $\lambda = (\lambda_1, \dots, \lambda_r)'$  be an  $r$ -dimensional vector of nonnegative integers (multi-indices), with norm  $|\lambda| = \sum_{j=1}^r \lambda_j$ , let  $(\lambda(k))_{k=1}^\infty$  be a sequence that includes all distinct multi-indices and satisfies  $|\lambda(k)| \leq |\lambda(k+1)|$ , and let  $x^\lambda = \prod_{j=1}^r x_j^{\lambda_j}$ . For a sequence  $\lambda(k)$  we consider the series  $r_{kK}(x) = x^{\lambda(k)}$ . If we denote the logistic cdf by  $L(a) = \exp(a)/(1 + \exp(a))$ , the SLE for  $p^*(x)$  is defined by  $\hat{p}(x) = L(R^K(x)' \hat{\pi}_K)$  with

$$\hat{\pi}_K = \arg \max_{\pi} \sum_{i=1}^N (T_i \cdot \ln L(R^K(X_i)' \pi) + (1 - T_i) \cdot \ln(1 - L(R^K(X_i)' \pi))).$$

In Appendix A we discuss the relevant asymptotic theory for  $\hat{p}(x)$ .

In addition to the unconfoundedness assumption the following assumptions are used to derive the properties of the estimator. First, we restrict the distribution of  $X$ ,  $Y(0)$ , and  $Y(1)$ :

ASSUMPTION 2 (Distribution of  $X$ ):

- (i) the support  $\mathbf{X}$  of the  $r$ -dimensional covariate  $X$  is a Cartesian product of compact intervals,  $\mathbf{X} = \prod_{j=1}^r [x_{lj}, x_{uj}]$ ;
- (ii) the density of  $X$  is bounded, and bounded away from 0, on  $\mathbf{X}$ .

ASSUMPTION 3 (Distribution of  $Y(0)$ ,  $Y(1)$ ):

- (i)  $\mathbf{E}[Y(0)^2] < \infty$  and  $\mathbf{E}[Y(1)^2] < \infty$ ;
- (ii)  $\mu_0(x)$  and  $\mu_1(x)$  are continuously differentiable for all  $x \in \mathbf{X}$ .

The next assumption requires sufficient smoothness of the propensity score.

ASSUMPTION 4 (Selection Probability): The propensity score  $p^*(x)$  satisfies the following conditions. For all  $x \in \mathbf{X}$ :

- (i)  $p^*(x)$  is continuously differentiable of order  $s \geq 7 \cdot r$  where  $r$  is the dimension of  $X$ ;
- (ii)  $p^*(x)$  is bounded away from zero and one:  $0 < \underline{p} \leq p^*(x) \leq \bar{p} < 1$ .

Finally, we restrict the rate at which additional terms are added to the series approximation to  $p^*(x)$ , depending on the dimension of  $X$  and the number of derivatives of  $p^*(x)$ .

ASSUMPTION 5 (Series Estimator): The series logit estimator of  $p^*(x)$  uses a power series with  $K = N^\nu$  for some  $1/(4(s/r - 1)) < \nu < \frac{1}{9}$ .

The restriction on the derivatives (Assumption 4(i)) guarantees the existence of a  $\nu$  that satisfies the conditions in Assumption 5. Under these conditions we can state the first result.

THEOREM 1: *Suppose Assumptions 1–5 hold. Then:*

- (i)  $\hat{\tau} \xrightarrow{p} \tau^*$ ;  
(ii)  $\sqrt{N}(\hat{\tau} - \tau^*) \xrightarrow{d} \mathcal{N}(0, V)$ , where

$$\begin{aligned} V &= \mathbf{E} \left[ \left( \left( \frac{YT}{p^*(X)} - \frac{Y(1-T)}{1-p^*(X)} - \tau^* \right) \right. \right. \\ &\quad \left. \left. - \left( \frac{\mu_1(X)}{p^*(x)} + \frac{\mu_0(X)}{1-p^*(X)} \right) (T - p^*(X)) \right)^2 \right] \\ &= \mathbf{E} \left[ (\tau(X) - \tau)^2 + \frac{\sigma_1^2(X)}{p^*(X)} + \frac{\sigma_0^2(X)}{1-p^*(X)} \right]; \quad \text{and} \end{aligned}$$

- (iii)  $\hat{\tau}$  reaches the semiparametric efficiency bound.

PROOF: See Appendix B.

REMARK 1: This result also covers the extension of the binary-covariate MAR example in Section 3 to the continuous covariate case. For this case set  $Y = 0$  if  $T = 0$  and set  $Y(0)$  identically equal to 0.

REMARK 2: Theorem 1 establishes the result for continuous  $X$ . If  $X$  has both continuous and discrete components, this can be dealt with in a conceptually straightforward manner by using the continuous covariate estimator within samples homogenous in the discrete covariates, at the expense of additional notation.

Derivations presented in Appendix B show that the estimator in Theorem 1 can be represented as asymptotically linear:

$$\hat{\tau} = \tau^* + \frac{1}{N} \sum_{i=1}^N (\psi(Y_i, T_i, X_i, \tau^*, p^*(X_i)) + \alpha(T_i, X_i)) + o_p(1/\sqrt{N}),$$

where  $\psi(\cdot)$  is defined in (10) and

$$(12) \quad \alpha(t, x) = - \left( \frac{\mu_1(x)}{p^*(x)} + \frac{\mu_0(x)}{1-p^*(x)} \right) \cdot (t - p^*(x)).$$

The known-weights estimator, (11) with  $\hat{p}(x)$  replaced by  $p^*(x)$ , is asymptotically linear with score function  $\psi(\cdot)$ . The function  $\alpha(t, x)$  represents the effect on the score function of estimating  $p^*(x)$ . Its first factor,  $-(\mu_1(x)/p^*(x) + \mu_0(x)/(1-p^*(x)))$ , is the conditional expectation of the derivative of the moment condition  $\psi(y, t, x, \tau^*, p^*(x))$  with respect to  $p^*(x)$ . Hence, the score linearizes the estimator with respect to  $\tau$  (which is trivial since the estimator is already linear in  $\tau$ ) and  $p(\cdot)$ .

The asymptotically linear representation of  $\hat{\tau}$  implies that its asymptotic variance equals

$$(13) \quad \mathbf{E}[(\psi(Y, T, X, \tau^*, p^*(X)) + \alpha(T, X))^2],$$

shown in Appendix B to be equal to the variance expression in Theorem 1. We estimate this variance by replacing the unknown quantities  $\tau$ ,  $p^*(\cdot)$ , and  $\alpha(\cdot)$  by estimates and replacing the expectation by a sample average:

$$(14) \quad \widehat{V} = \frac{1}{N} \sum_{i=1}^N (\psi(Y_i, T_i, X_i, \hat{\tau}, \hat{p}(X_i)) + \hat{\alpha}(T_i, X_i))^2.$$

The estimation of  $\alpha(t, x)$  requires some additional explanation. The second factor,  $t - p^*(x)$  is estimated as  $t - \hat{p}(x)$ . The first factor,  $-(\mu_1(x)/p^*(x) + \mu_0(x)/(1 - p^*(x)))$ , can be written as the conditional expectation of  $-(YT/p^*(X)^2 + Y(1-T)/(1 - p^*(X))^2)$  given  $X$ . We therefore estimate the first factor in  $\alpha(t, x)$  by nonparametric regression of  $-(YT/\hat{p}(X)^2 + Y(1-T)/(1 - \hat{p}(X))^2)$  on  $X$ , using the same series approach as we used for estimating  $p^*(x)$ . Thus

$$\begin{aligned} & - \left( \frac{1}{N} \sum_{i=1}^N \left( \frac{Y_i T_i}{\hat{p}(X_i)^2} + \frac{Y_i(1-T_i)}{(1-\hat{p}(X_i))^2} \right) R^K(X_i) \right)' \\ & \times \left( \frac{1}{N} \sum_{i=1}^N R^K(X_i) R^K(X_i)' \right)^{-1} R^K(x), \end{aligned}$$

with  $R^K(x)$  the same series of approximating functions as before, is used as an estimator for  $-(\mu_1(x)/p^*(x) + \mu_0(x)/(1 - p^*(x)))$ , and the function  $\alpha(t, x)$  is estimated by  $\hat{\alpha}(t, x)$ :

$$(15) \quad \begin{aligned} \hat{\alpha}(t, x) = & - \left( \frac{1}{N} \sum_{i=1}^N \left( \frac{Y_i T_i}{\hat{p}(X_i)^2} + \frac{Y_i(1-T_i)}{(1-\hat{p}(X_i))^2} \right) R^K(X_i) \right)' \\ & \times \left( \frac{1}{N} \sum_{i=1}^N R^K(X_i) R^K(X_i)' \right)^{-1} R^K(x) (t - \hat{p}(x)). \end{aligned}$$

The following theorem describes the formal result.

**THEOREM 2:** *Suppose Assumptions 1–5 hold. Then  $\widehat{V}$  is consistent for  $V$ .*

**PROOF:** See Appendix B.

In practice bootstrapping methods may be a valuable alternative to the above variance estimator.

#### 4.2. Estimating the Weighted Average Treatment Effect

In this section we generalize the previous result to  $\tau_{\text{wate}}^*$ , the weighted average treatment effect for a known weight function  $g(x)$ . One motivation for considering this estimand is that by choosing  $g(x)$  appropriately, we can obtain treatment effects for subpopulations defined by  $X$ . In addition, by choosing  $g(x)$  equal to

the propensity score  $p^*(x)$ , we can recover the average effect of the treatment on the treated, as will be discussed below.

To estimate  $\tau_{wate}$ , we use the following moment function:

$$(16) \quad \psi(y, t, x, \tau_{wate}, p(x)) = g(x) \cdot \left( \frac{y \cdot t}{p(x)} - \frac{y \cdot (1-t)}{1-p(x)} - \tau_{wate} \right),$$

leading to the estimator

$$\hat{\tau}_{wate} = \sum_i g(X_i) \left[ \frac{Y_i \cdot T_i}{\hat{p}(X_i)} - \frac{Y_i \cdot (1-T_i)}{1-\hat{p}(X_i)} \right] / \sum_i g(X_i).$$

This estimator is asymptotically linear:

$$\begin{aligned} \hat{\tau}_{wate} = \frac{1}{\mathbf{E}[g(X)]} \frac{1}{N} \sum_{i=1}^N (\psi(Y_i, T_i, X_i, \tau_{wate}^*, p^*(x)) \\ + \alpha(T_i, X_i)) + o_p(1/\sqrt{N}), \end{aligned}$$

where now

$$\alpha(t, x) = -g(x) \cdot \left( \frac{\mu_1(x)}{p^*(x)} + \frac{\mu_0(x)}{1-p^*(x)} \right) (t - p^*(x)).$$

The asymptotic variance can be estimated as

$$\hat{V} = \frac{1}{(\sum_i g(X_i)/N)^2} \frac{1}{N} \sum_{i=1}^N (\psi(Y_i, T_i, X_i, \hat{\tau}_{wate}, \hat{p}(X_i)) + \hat{\alpha}(T_i, X_i))^2,$$

with an estimator for  $\alpha(t, x)$  analogous to that for the average treatment effect:

$$\begin{aligned} \hat{\alpha}(t, x) = -g(x) \frac{1}{N} \sum_{i=1}^N \left( \left( \frac{Y_i T_i}{\hat{p}_K(X_i)^2} + \frac{Y_i (1-T_i)}{(1-\hat{p}_K(X_i))^2} \right) R^K(X_i) R^K(X_i)' \right)' \\ \times \left( \frac{1}{N} \sum_{i=1}^N R^K(X_i) R^K(X_i)' \right)^{-1} R^K(x) (t - \hat{p}_K(x)). \end{aligned}$$

Similar reasoning to the previous results gives the following theorem:

**THEOREM 3:** *Suppose Assumptions 1–5 hold, that  $|g(x)|$  is bounded from above and that  $\mathbf{E}[g(X)] > 0$ . Then:*

- (i)  $\hat{\tau}_{wate} \xrightarrow{p} \tau_{wate}^*$ ;
- (ii)  $\sqrt{N}(\hat{\tau}_{wate} - \tau_{wate}^*) \xrightarrow{d} \mathcal{N}(0, V)$ , with

$$\begin{aligned} V = \frac{1}{\mathbf{E}[g(X)]^2} \mathbf{E} \left[ g(X)^2 (\tau(X) - \tau_{wate}^*)^2 + \frac{g(X)^2}{p^*(X)} \sigma_1^2(X) \right. \\ \left. + \frac{g(X)^2}{1-p^*(X)} \sigma_0^2(X) \right]; \quad \text{and} \end{aligned}$$

- (iii)  $\hat{V}$  is consistent for  $V$ .

The proof for this theorem follows the same line of argument as that for Theorems 1 and 2 and is omitted.

REMARK: We could weaken Assumption 4(ii), the assumption that the propensity score is bounded away from 0 and 1, by the assumption that  $g(x)/p^*(x)$  and  $g(x)/(1-p^*(x))$  are bounded from above. Thus, if there is insufficient overlap in the distributions of the treated and untreated units, one may wish to choose  $g(\cdot)$  to restrict attention to a subpopulation for which there is sufficiently large probability of observing both treated and untreated units.

A semiparametric efficiency bound for  $\tau_{wate}$  has not been previously calculated in the literature. The next result shows that our estimator is efficient.

THEOREM 4: *The semiparametric efficiency bound for estimation of  $\tau_{wate}$  is*

$$V = \frac{1}{\mathbf{E}[g(X)]^2} \mathbf{E} \left[ g(X)^2 (\tau(X) - \tau_{wate})^2 + \frac{g(X)^2}{p^*(X)} \sigma_1^2(X) + \frac{g(X)^2}{1-p^*(X)} \sigma_0^2(X) \right].$$

PROOF: See Appendix B.

#### 4.3. Estimating the Average Treatment Effect for the Treated

Under unconfoundedness the average treatment effect for the treated (Rubin (1977), Heckman and Robb (1985), Heckman, Ichimura, and Todd (1997, 1998)) is a special case of the weighted average treatment effect, corresponding to the weighting function  $g(x) = p^*(x)$ . To see this, first note that under unconfoundedness

$$\begin{aligned} \tau_{treated}^* &= \mathbf{E}[Y(1) - Y(0)|T = 1] = \mathbf{E}[\mathbf{E}[Y(1) - Y(0)|X, T = 1]|T = 1] \\ &= \mathbf{E}[\mathbf{E}[Y(1) - Y(0)|X]|T = 1] = \mathbf{E}[\tau(X)|T = 1]. \end{aligned}$$

Second, the latter is equal to

$$\begin{aligned} \mathbf{E}[\tau(X)|T = 1] &= \int \tau(x) dF(x|T = 1) \\ &= \int \tau(x) p^*(x) dF(x) / \int p^*(x) dF(x), \end{aligned}$$

which is  $\tau_{wate}^*$  with  $g(x)$  equal to  $p^*(x)$ . Hence we can use the moment equation (16) with  $p^*(x)$  substituted for  $g(x)$ :

$$(17) \quad \psi(y, t, x, \tau_{treated}, p(x)) = p^*(x) \cdot \left( \frac{y \cdot t}{p(x)} - \frac{y \cdot (1-t)}{1-p(x)} - \tau_{treated} \right).$$

The estimator is the solution to

$$(18) \quad 0 = \sum_{i=1}^N p^*(X_i) \cdot \left( \frac{Y_i \cdot T_i}{\hat{p}(X_i)} - \frac{Y_i \cdot (1 - T_i)}{1 - \hat{p}(X_i)} - \tau_{treated} \right),$$

with the same nonparametric series estimator  $\hat{p}(x)$  as before.

The next result, which follows directly from Theorem 4, shows that this estimator achieves the efficiency bound calculated by Hahn (1998) for estimation of the effect of treatment on the treated, assuming that the propensity score is known.

**COROLLARY 1:** *Suppose that Assumptions 1–5 hold. Then:*

- (i)  $\hat{\tau}_{treated} \xrightarrow{P} \tau_{treated}^*$ ;
- (ii)  $\sqrt{N}(\hat{\tau}_{treated} - \tau_{treated}^*) \xrightarrow{d} \mathcal{N}(0, V)$ , with

$$V = \frac{1}{\mathbf{E}[p^*(X)]^2} \mathbf{E} \left[ p^*(X)^2 (\tau(X) - \tau_{treated})^2 + p^*(X) \sigma_1^2(X) + \frac{p^*(X)^2}{1 - p^*(X)} \sigma_0^2(X) \right], \quad \text{and}$$

- (iii)  $\hat{\tau}_{treated}$  achieves the semiparametric efficiency bound.

The proof for this corollary is omitted as the result directly follows from Theorem 4.

Note that in the moment function (17) the propensity score appears in two places, first as  $p^*(x)$  multiplying the remainder of the moment function where it replaces the general weight function  $g(x)$  in (16), and second as  $p(x)$  in the denominator of the two terms. We only use the estimated propensity score in the second part in the efficient estimator in (18). The result of the theorem above implies that this is more efficient than using the true propensity score everywhere and solving

$$(19) \quad 0 = \sum_{i=1}^N p^*(X_i) \cdot \left( \frac{Y_i \cdot T_i}{p^*(X_i)} - \frac{Y_i \cdot (1 - T_i)}{1 - p^*(X_i)} - \tau_{treated} \right),$$

or using the estimated propensity score everywhere, which amounts to solving

$$(20) \quad 0 = \sum_{i=1}^N \hat{p}(X_i) \cdot \left( \frac{Y_i \cdot T_i}{\hat{p}(X_i)} - \frac{Y_i \cdot (1 - T_i)}{1 - \hat{p}(X_i)} - \tau_{treated} \right).$$

A direct implication of this result is that the sample average of the outcomes for the treated  $\sum_i Y_i T_i / \sum_i T_i$  is less efficient for the population average  $\mathbf{E}[Y(1)|T=1]$  than the weighted average  $\sum_i Y_i T_i (p^*(X_i) / \hat{p}(X_i)) / \sum_i p^*(X_i)$  where the weights are the ratio of the true and estimated propensity score. Another implication is that the estimators characterized by (19) and (20) cannot in general be ranked in terms of efficiency as there are effects of opposite signs (e.g., Heckman, Ichimura, and Todd (1997)).

If the propensity score is not known, then Hahn (1998) shows that this affects the efficiency bound for the effect of treatment on the treated. Our previous estimator  $\hat{\tau}_{treated}$  cannot be used since it makes use of  $p^*(x)$ . However, we can use the estimated propensity score in place of  $p^*(x)$  in the weighting of observations as in (20). Call this estimator  $\hat{\tau}_{te}$ . The next theorem shows that this estimator is efficient if the propensity is not known.

THEOREM 5: *Suppose that Assumptions 1–5 hold. Then:*

- (i)  $\hat{\tau}_{te} \xrightarrow{p} \tau_{treated}^*$ ;
- (ii)  $\sqrt{N}(\hat{\tau}_{te} - \tau_{treated}^*) \xrightarrow{d} \mathcal{N}(0, V)$ , with

$$V = \frac{1}{\mathbf{E}[p^*(X)]^2} \mathbf{E} \left[ p^*(X)(\tau(X) - \tau_{treated}^*)^2 + p^*(X)\sigma_1^2(X) + \frac{p^*(X)^2}{1 - p^*(X)} \sigma_0^2(X) \right], \quad \text{and}$$

- (iii)  $\hat{\tau}_{te}$  achieves the semiparametric efficiency bound for estimation of  $\tau_{treated}$  when the propensity score is not known.

The proof goes along the same lines as that for Theorems 1 and 2 and is omitted.

## 5. CONCLUSION

In this paper we have studied efficient estimation of various average treatment effects under an unconfounded treatment assignment assumption. Although weighting observations by the inverse of the true propensity score does not lead to efficient estimators, weighting each observation by the inverse of a nonparametric estimate of the propensity score does lead to efficient estimators. We provide intuition for this result through connections to the literatures on empirical likelihood estimators and choice-based sampling.

The estimators proposed in this paper require fewer functions to be estimated nonparametrically than other efficient estimators previously proposed in the literature. Which estimators have more attractive finite sample properties, and which have more attractive computational properties, remain open questions. The results underline the important role played by the propensity score in estimation of average causal effects.

*Dept. of Economics, University of Miami, P.O. Box 248126, Coral Gables FL 33124-6550; khirano@miami.edu;*

*Dept. of Agricultural and Resource Economics, University of California at Berkeley, 330 Giannini Hall, Berkeley, CA 94720-3310; imbens@econ.berkeley.edu; <http://elsa.berkeley.edu/users/imbens/>; and Dept. of Economics, University of California at Berkeley, and NBER;*



and

*Dept. of Economics, University of Southern California, Kaprielian Hall, University Park Campus, Los Angeles, CA 90089; ridder@usc.edu.*

*Manuscript received January, 2000; final revision received August, 2002.*

#### APPENDIX A: LOGISTIC SERIES ESTIMATOR

In this appendix we derive the relevant properties of the logistic series estimator, which can be interpreted as a sieve estimator (e.g., Geman and Hwang (1982)). Let  $r^K(x) = (r_{1K}(x), \dots, r_{KK}(x))'$  be a  $K$ -vector of functions. The triangular array of functions  $r^K(x)$ ,  $K = 1, 2, \dots$ , is the basis for the approximation of the propensity score. In particular, we approximate a function  $f: \mathbf{R}^r \rightarrow \mathbf{R}$  by  $\gamma'_K r^K(x)$ . Because  $\gamma'_K r^K(x) = \gamma'_K A_K^{-1} A_K r^K(x)$  we can also use  $R^K(x) = A_K r^K(x)$  as the basis of approximation. By choosing  $A_K$  appropriately we obtain a system of orthogonal (with respect to some weight function) functions. Specifically we choose  $A_K$  so that  $\mathbf{E}[R^K(X)R^K(X)'] = I_K$ . The properties of the series logit estimator and the proof of Theorem 1 are mostly for a general system of approximating functions. We shall indicate where the properties of the specific approximating class of functions are used. We will use the matrix norm  $\|A\| = \sqrt{\text{tr}(A'A)}$ . Note that this is the usual Euclidean norm if  $A$  is a column vector.<sup>11</sup> If  $A$  is a scalar, we denote the norm by  $|A|$ . Define

$$(21) \quad \zeta(K) = \sup_{x \in X} \|R^K(x)\|.$$

In general, this bound depends on the array of approximating functions that is used. For orthonormal polynomials Newey (1994, 1997) shows  $\zeta(K) \leq CK$ . Here, and in the sequel,  $C$  denotes a generic positive constant.<sup>12</sup>

We consider approximation of the log odds ratio by a power series. One possible choice for a triangular sequence of powers of  $x$  is

$$(22) \quad r^1(x) = 1, \quad r^2(x) = \begin{bmatrix} 1 \\ x_1 \end{bmatrix}, \quad r^{r+1}(x) = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_r \end{bmatrix}, \quad r^{r+1}(x) = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_r \\ x_1^2 \end{bmatrix}, \dots$$

Linear combinations of the elements of the vectors  $r^K(x)$  are the approximating power series. A power series with  $(n+1)^r$  terms has  $x_1, \dots, x_r$  included at least up to power  $n$ . Hence, if we use the sequence in (22) and set  $K = (n+1)^r$ , then  $r^K(x)$  has powers in all variables at least up to  $n$ . If a function  $f$  is  $s$  times continuously differentiable and  $K = (n+1)^r$ , then by Theorem 8, p. 90, in Lorentz (1986) there is a  $K$ -vector  $\gamma_K$  such that for  $R^K(x) = A_K r^K(x)$ , and on the compact set  $\mathbf{X}$ ,

$$(23) \quad \sup_{x \in X} |f(x) - R^K(x)' \gamma_K| < C_1 n^{-s} \leq C_2 K^{-\frac{s}{r}}.$$

<sup>11</sup> It is useful to list some properties of this norm that will be used in the following discussion. Let  $A$  and  $B$  be  $K \times K$  matrices and  $c$  be a  $K$  vector. Then  $\|AB\|^2 = \sum_i \sum_j (\sum_k a_{ik} b_{kj})^2$ . Applying the vector Cauchy-Schwartz inequality to the inner sum, we find  $\|AB\| \leq \|A\| \|B\|$ . By the maximum inequality for quadratic forms  $\|Ac\| \leq \sqrt{\lambda_{\max}(A'A)} \|c\|$ , where  $\lambda_{\max}$  is the largest eigenvalue, which gives a sharp upper bound (the upper bound  $\|A\| \|c\|$  is not sharp in general). We also frequently use the Cauchy-Schwartz inequality for expectations, which implies that for nonnegative random variables  $X, Y$ ,  $\mathbf{E}(XY) \leq \sqrt{\mathbf{E}(X^2)\mathbf{E}(Y^2)}$ .

<sup>12</sup> If two constants are needed, we will use the generic notation  $C_1, C_2$ .

To ensure that the approximation of  $p^*(x)$  is between 0 and 1 we first approximate the log odds ratio, which is also  $s$  times continuously differentiable and which is bounded on  $\mathbf{X}$  if the propensity score is bounded away from zero and one. Hence by (23) there is a  $\pi_K$  such that

$$(24) \quad \sup_{x \in \mathbf{X}} \left| \ln \left( \frac{p^*(x)}{1 - p^*(x)} \right) - R^K(x)' \pi_K \right| < CK^{-\frac{s}{r}}.$$

Let  $L(z) = \exp(z)/(1 + \exp(z))$  be the logistic cdf and  $L'(z) = L(z) \cdot (1 - L(z))$ . The series logit estimator of the population propensity score  $p^*(x)$  is  $\hat{p}_K(x) = L(R^K(x)' \hat{\pi}_K)$ , where

$$(25) \quad \hat{\pi}_K = \arg \max_{\pi} L_N(\pi),$$

for

$$(26) \quad L_N(\pi) = \sum_{i=1}^N (T_i \cdot \ln L(R^K(X_i)' \pi) + (1 - T_i) \cdot \ln(1 - L(R^K(X_i)' \pi))).$$

For  $N \rightarrow \infty$  and  $K$  fixed we have  $\hat{\pi}_K \xrightarrow{p} \pi_K^*$ , with  $\pi_K^*$  the pseudo true value:

$$(27) \quad \pi_K^* = \arg \max_{\pi} \mathbb{E}[p^*(X) \ln L(R^K(X)' \pi) + (1 - p^*(X)) \ln(1 - L(R^K(X)' \pi))].$$

We also define the pseudo true propensity score:  $p_K^*(x) = L(R^K(x)' \pi_K^*)$ .

In the proofs for the theorems we need some properties of this series logit estimator. For these properties it is convenient to distinguish between (i) the deterministic difference between the true propensity score and the pseudo true propensity score and (ii) the stochastic difference between the estimated propensity score and the pseudo true value. In the remainder of this appendix we therefore derive (i) a uniform bound on the difference between  $p^*(x)$  and  $p_K^*(x)$  and (ii) a bound on the sampling variance in the form of the stochastic order of  $\|\hat{\pi}_K - \pi_K^*\|$ .

LEMMA 1 (Approximation of Propensity Score): *Suppose that:*

- (i) *the support  $\mathbf{X}$  of  $X$  is a compact subset of  $\mathbf{R}^r$ ;*
- (ii) *the propensity score  $p^*(x)$  is  $s$  times continuously differentiable, with  $s/r \geq 4$ ;*
- (iii) *the propensity score  $p^*(x)$  is bounded away from zero and one on  $\mathbf{X}$ ;*
- (iv) *the density of  $X$  is bounded away from zero on  $\mathbf{X}$ .*

*Then for  $\pi_K$  in (24),*

$$\|\pi_K - \pi_K^*\| = O(K^{-s/(2r)}),$$

and

$$\sup_{x \in \mathbf{X}} |p^*(x) - p_K^*(x)| = O(K^{-s/(2r)} \zeta(K)).$$

PROOF: From (24), and by monotonicity of  $L(\cdot)$ , for all  $x \in \mathbf{X}$ ,

$$(28) \quad \begin{aligned} L(R^K(x)' \pi_K - CK^{-s/r}) - L(R^K(x)' \pi_K) &< p^*(x) - L(R^K(x)' \pi_K) \\ &< L(R^K(x)' \pi_K + CK^{-s/r}) - L(R^K(x)' \pi_K). \end{aligned}$$

By the mean value theorem applied to the lower and upper bound and by  $L'(R^K(x)' \tilde{\pi}) = L(R^K(x)' \tilde{\pi})(1 - L(R^K(x)' \tilde{\pi})) < 1/4$  we find that the lower and upper bound are bounded by  $-\frac{1}{4}CK^{-s/r}$  and  $\frac{1}{4}CK^{-s/r}$ , respectively. Hence, for the  $\pi_K$  that satisfies (24), we have

$$(29) \quad \sup_{x \in \mathbf{X}} |p^*(x) - L(R^K(x)' \pi_K)| < CK^{-s/r}.$$

Define

$$Q^*(\pi) = \mathbf{E}[p^*(X) \ln L(R^K(X)' \pi) + (1 - p^*(X)) \ln(1 - L(R^K(X)' \pi))],$$

and

$$Q_K(\pi) = \mathbf{E}[L(R^K(X)' \pi_K) \ln L(R^K(X)' \pi) + (1 - L(R^K(X)' \pi_K)) \ln(1 - L(R^K(X)' \pi))].$$

Then, by definition we have  $\pi_K^* = \arg \max_{\pi} Q^*(\pi)$ , and by the information inequality we have

$$\pi_K = \arg \max_{\pi} Q_K(\pi).$$

Let  $\eta = \inf_{x \in X} p^*(x) \cdot (1 - p^*(x))$ , so that by assumption (iii)  $\eta > 0$ . Define

$$\Pi_K = \left\{ \pi \in \mathbf{R}^K \mid \inf_{x \in X} L(R^K(x)' \pi)(1 - L(R^K(x)' \pi)) \geq \eta/2 \right\}.$$

Because of (29), for  $K$  large enough, we have  $\pi_K \in \Pi_K$ . Also, because  $L(R^K(x)' \pi)$  is bounded away from zero and one for  $\pi \in \Pi_K$ , it follows that  $\ln L(R^K(x)' \pi)$  is bounded, and thus by (29) there is a  $C_1$  such that

$$(30) \quad \sup_{\pi \in \Pi_K} |Q^*(\pi) - Q_K(\pi)| \leq C_1 K^{-s/r}.$$

Next, define for fixed  $C_2$

$$(31) \quad \tilde{\Pi}_K = \{\pi \in \mathbf{R}^K \mid \|\pi - \pi_K\| \leq C_2 K^{-s/(2r)}\}.$$

Because

$$\begin{aligned} \sup_{x \in X, \pi \in \tilde{\Pi}_K} |L(R^K(x)' \pi) - L(R^K(x)' \pi_K)| &\leq \sup_{x \in X, \pi \in \tilde{\Pi}_K, \tilde{\pi}} |L'(R^K(x)' \tilde{\pi}) R^K(x)' (\pi - \pi_K)| \\ &\leq \sup_{x \in X} \|R^K(x)\| \cdot \sup_{\pi \in \tilde{\Pi}_K} \|\pi - \pi_K\| \leq C \zeta(K) K^{-s/(2r)}, \end{aligned}$$

it follows that for a polynomial series estimator with  $\zeta(K) \leq CK$ , and for large enough  $K$ ,  $\tilde{\Pi}_K \subset \Pi_K$ . Thus, for  $\pi \in \tilde{\Pi}_K$  and with  $\lambda_{\min}(A)$  the smallest eigenvalue of  $A$ , and for large enough  $K$  so that  $\tilde{\Pi}_K \subset \Pi_K$ , given that  $\mathbf{E}[R^K(X) R^K(X)'] = I_K$ , we have

$$(32) \quad \lambda_{\min} \left( -\frac{\partial^2 Q_K}{\partial \pi \partial \pi'}(\pi) \right) = \lambda_{\min}(\mathbf{E}[L'(R^K(x)' \pi) R^K(X) R^K(X)']) \geq \eta/2.$$

Now choose the  $C_2$  in (31) to satisfy  $C_2 > \sqrt{4C_1/\eta}$ , for the  $C_1$  in (30). Let  $K$  be large enough so that  $\tilde{\Pi}_K \subset \Pi_K$ . Then, for  $\pi$  such that  $\|\pi - \pi_K\| = C_2 K^{-s/(2r)}$ , i.e.,  $\pi \in \tilde{\Pi}_K$ , the difference  $Q^*(\pi_K) - Q^*(\pi)$  satisfies

$$\begin{aligned} Q^*(\pi_K) - Q^*(\pi) &\geq Q_K(\pi_K) - C_1 K^{-s/r} - Q_K(\pi) - C_1 K^{-s/r} \\ &\geq -\frac{\partial Q_K}{\partial \pi}(\pi_K)(\pi - \pi_K) - \frac{1}{2}(\pi - \pi_K)' \frac{\partial^2 Q_K}{\partial \pi \partial \pi'}(\tilde{\pi})(\pi - \pi_K) - 2C_1 K^{-s/r} \\ &\geq \frac{C_2^2 \eta}{2} K^{-s/r} - 2C_1 K^{-s/r} > 0. \end{aligned}$$

Hence, there is a local maximum of  $Q^*(\pi)$  in the interior of  $\tilde{\Pi}_K$ , so that with  $Q^*(\pi)$  concave,  $\pi_K^* = \arg \max Q^*(\pi)$  must satisfy  $\pi_K^* \in \tilde{\Pi}_K$ , proving the first assertion. Then

$$|L(R^K(X)' \pi_K) - L(R^K(X)' \pi_K^*)| \leq \|R^K(X)\| \cdot \|\pi_K^* - \pi_K\| = O(\zeta(K) K^{-s/(2r)}),$$

and thus by (29) and the triangle inequality we have

$$\sup_{x \in X} |p^*(x) - L(R^K(x)' \pi_K^*)| = O(K^{-s/(2r)} \zeta(K)). \quad Q.E.D.$$

Next we derive the stochastic order of  $\|\hat{\pi}_K - \pi_K^*\|$  as  $K$  increases with  $N$ .

LEMMA 2 (Convergence of  $\hat{\pi}_K - \pi_K^*$ ): Suppose the same four conditions as in Lemma 1 hold. In addition, suppose that:

(v)  $K(N)$  is a sequence of values of  $K$  satisfying  $K(N) \rightarrow \infty$ , and  $\zeta(K(N))^4/N \rightarrow 0$ .  
Then

$$\|\hat{\pi}_{K(N)} - \pi_{K(N)}^*\| = O_p\left(\sqrt{\frac{K(N)}{N}}\right).$$

PROOF: In the sequel we write  $K$  for  $K(N)$ . By definition of  $R^K(x)$ ,

$$(33) \quad \hat{S}_K = \frac{1}{N} \sum_{i=1}^N R^K(X_i) R^K(X_i)'$$

has expectation equal to  $I_K$ . By Newey (1997), it satisfies

$$\|\hat{S}_K - I_K\| = O_p\left(\zeta(K) \sqrt{\frac{K}{N}}\right),$$

which converges to zero in probability by condition (v). Hence the probability that the smallest eigenvalue of  $\hat{S}_K$  is larger than  $1/2$  goes to one.

Next, we show that

$$(34) \quad \frac{1}{N} \frac{\partial L_N}{\partial \pi}(\pi_K^*) = O_p\left(\sqrt{\frac{K}{N}}\right).$$

Consider

$$\begin{aligned} \mathbf{E} \left\| \frac{1}{N} \frac{\partial L_N}{\partial \pi}(\pi_K^*) \right\|^2 &= \frac{1}{N} \text{tr} \mathbf{E}[L(R^K(X)' \pi_K^*)(1 - L(R^K(X)' \pi_K^*)) R^K(X) R^K(X)'] \\ &\leq \frac{1}{N} \text{tr} \mathbf{E}[R^K(X) R^K(X)'] = K/N. \end{aligned}$$

Hence

$$\mathbf{E} \left\| \frac{1}{N} \frac{\partial L_N}{\partial \pi}(\pi_K^*) \right\|^2 = O(K/N),$$

and the Markov inequality implies (34).

Next, let  $\eta = \inf_{x \in \mathbf{X}, K} L(R^K(x)' \pi_K^*)(1 - L(R^K(x)' \pi_K^*)) / 8$ , which by conditions (i) and (iii) and Lemma 1 is positive. For any  $\varepsilon > 0$  choose  $C$  such that for  $N$  large enough

$$(35) \quad \Pr\left(\left\| \frac{1}{N} \frac{\partial L_N}{\partial \pi}(\pi_K^*) \right\| < \eta C \sqrt{\frac{K}{N}}\right) \geq 1 - \frac{\varepsilon}{2}.$$

Note that,

$$\begin{aligned} &\sup_{x \in \mathbf{X}, \|\pi - \pi_K^*\| \leq C \sqrt{\frac{K}{N}}} |L(R^K(x)' \pi) - L(R^K(x)' \pi_K^*)| \\ &\leq \sup_{x \in \mathbf{X}, \|\pi - \pi_K^*\| \leq C \sqrt{\frac{K}{N}}} |R^K(x)'(\pi - \pi_K^*)| \leq \zeta(K) C \sqrt{\frac{K}{N}}, \end{aligned}$$

which goes to zero, so that for large enough  $N$ ,

$$\inf_{x \in \mathbf{X}, \|\pi - \pi_K^*\| \leq C \sqrt{\frac{K}{N}}} (L(R^K(x)' \pi)(1 - L(R^K(x)' \pi))) \geq 4\eta.$$

Choose  $N$  large enough so that this inequality holds, that (35) holds with probability at least  $1 - \varepsilon/2$ , and that the probability that the smallest eigenvalue of  $\hat{S}_K$  is larger than  $1/2$  is at least  $1 - \varepsilon/2$ . Then the probability that both of these hold is at least  $1 - \varepsilon$ . Then for every  $\pi$  with  $\|\pi - \pi_K^*\| = C\sqrt{K/N}$ , a second-order expansion gives

$$(36) \quad \frac{1}{N} L_N(\pi) = \frac{1}{N} L_N(\pi_K^*) + \frac{1}{N} \frac{\partial L_N(\pi_K^*)}{\partial \pi} (\pi - \pi_K^*) + \frac{1}{2N} (\pi - \pi_K^*)' \frac{\partial^2 L_N(\bar{\pi})}{\partial \pi \partial \pi'} (\pi - \pi_K^*)$$

where  $\|\bar{\pi} - \pi_K^*\| \leq \|\pi - \pi_K^*\| = C\sqrt{K/N}$ . We have

$$\begin{aligned} \frac{1}{2N} \frac{\partial^2 L_N(\bar{\pi})}{\partial \pi \partial \pi'} &= -\frac{1}{2N} \sum_{i=1}^n (L(R^K(X_i)' \bar{\pi}) (1 - L(R^K(X_i)' \bar{\pi})) R^K(X_i) R^K(X_i)') \\ &\leq -2\eta \hat{S}_K, \end{aligned}$$

with its eigenvalues bounded away from zero in absolute value by  $\eta$ . Then, rearranging (36) and using the triangle inequality, with probability greater than  $1 - \varepsilon$ , for  $\|\pi - \pi_K^*\| = C\sqrt{K/N}$ ,

$$\begin{aligned} \frac{1}{N} L_N(\pi) - \frac{1}{N} L_N(\pi_K^*) &\leq \frac{1}{N} \frac{\partial L_N}{\partial \pi} (\pi_K^*)' (\pi - \pi_K^*) - \eta \|\pi - \pi_K^*\|^2 \\ &\leq \left\| \frac{1}{N} \frac{\partial L_N}{\partial \pi} (\pi_K^*) \right\| \cdot \|\pi - \pi_K^*\| - \eta \|\pi - \pi_K^*\|^2 \\ &= \left( \left\| \frac{1}{N} \frac{\partial L_N}{\partial \pi} (\pi_K^*) \right\| - \eta C \sqrt{\frac{K}{N}} \right) \cdot \|\pi - \pi_K^*\| < 0. \end{aligned}$$

That is, we have with probability greater than  $1 - \varepsilon$ ,

$$\frac{1}{N} L_N(\pi) < \frac{1}{N} L_N(\pi_K^*) \quad \text{for all } \pi \text{ with } \|\pi - \pi_K^*\| = C\sqrt{\frac{K}{N}}.$$

Since  $L_N(\pi)$  is continuous, it has a maximum on the compact set  $\{\pi : \|\pi - \pi_K^*\| \leq C\sqrt{K/N}\}$ . By the last inequality, this maximum must occur for some  $\hat{\pi}_K$  with  $\|\hat{\pi}_K - \pi_K^*\| < C\sqrt{K/N}$ . Hence the first-order conditions are satisfied at  $\hat{\pi}_K$  and by concavity of  $L_N(\pi)$ ,  $\hat{\pi}_K$  maximizes  $L_N(\pi)$  over all of  $\mathbf{R}^K$ . Because the probability of this is greater than  $1 - \varepsilon$  with  $\varepsilon$  arbitrary, we conclude that  $\hat{\pi}_K$  exists and satisfies the first order conditions with probability approaching one, and that  $\|\hat{\pi}_K - \pi_K^*\| = O_p(\sqrt{K/N})$ . *Q.E.D.*

## APPENDIX B: PROOFS OF THEOREMS

**PROOF OF THEOREM 1:** To ease the notational burden we present the proof for the special case with  $Y(0)$  identically equal to 0. This can be interpreted as the special case where one is interested in estimating the average outcome  $\beta = \mathbf{E}[Y(1)]$ , where  $Y(1)$  is missing at random conditional on the covariates  $X$ . Thus it is the direct extension of the binary-covariate example in Section 3. Since the average treatment effect case simply amounts to estimating two averages where in both cases the variables are missing at random, the argument for the general case is exactly analogous, only involving substantially longer equations. In the proof we therefore follow the missing at random set up with the parameter of interest equal to  $\beta = \mathbf{E}[Y(1)]$ , making the missing at random assumption  $Y(1) \perp T|X$ , and with a random sample of  $(T_i, X_i, Y_i)_{i=1}^N$ , where  $Y_i = Y_i(1) \cdot T_i$ .

The estimated weight estimator  $\hat{\beta}_{ew}$  is

$$(37) \quad \hat{\beta}_{ew} = \frac{1}{N} \sum_{i=1}^N \frac{T_i \cdot Y_i}{\hat{p}_K(X_i)}$$

with  $\hat{p}_K(X_i) = L(R^K(X_i)' \hat{\pi}_K)$ . The key part of the proof is to show that

$$(38) \quad \left\| \sqrt{N}(\hat{\beta}_{ew} - \beta^*) - \frac{1}{\sqrt{N}} \sum_{i=1}^N \left\{ \left( \frac{T_i \cdot Y_i}{p^*(X_i)} - \beta^* \right) - \frac{\mu_1(X_i)}{p^*(X_i)} (T_i - p^*(X_i)) \right\} \right\| = o_p(1).$$

This implies that  $\hat{\beta}_{ew}$  is asymptotically linear, i.e. behaves asymptotically as a sample average, with score function  $\psi(Y, T, X, \beta^*, p^*(\cdot)) + \alpha(T, X)$ , where

$$\psi(y, t, x, \beta, p(\cdot)) = \frac{t \cdot y}{p(x)} - \beta, \quad \text{and} \quad \alpha(t, x) = -\frac{\mu_1(x)}{p^*(x)} \cdot (t - p^*(x)).$$

The first term of the score function,  $\psi(\cdot)$ , is equal to the score that would obtain if we substitute the population probability  $p^*$  for the estimator  $\hat{p}_K$  in (37). The second term,  $\alpha(\cdot)$ , gives the contribution of estimating  $p^*$  to the asymptotic distribution of  $\hat{\beta}_{ew}$ . This contribution is linear in  $T - p^*(X)$ . Hence, the score linearizes the estimator with respect to  $\beta$  and  $p(\cdot)$ . The asymptotic variance of  $\hat{\beta}_{ew}$  is equal to the variance of  $\psi(Y, T, x, \beta^*, p^*(X)) + \alpha(T, X)$  (note that its mean is 0). The three components of this variance are

$$\begin{aligned} \mathbf{E}[\psi(Y, T, X, \beta^*, p^*(\cdot))^2] &= \mathbf{E}\left[\frac{\mu_1(X)^2}{p^*(X)}\right] + \mathbf{E}\left[\frac{\sigma_1^2(X)}{p^*(X)}\right] - (\beta^*)^2, \\ \mathbf{E}[\alpha(T, X)^2] &= \mathbf{E}\left[\frac{\mu_1(X)^2}{p^*(X)}\right] - \mathbf{E}[\mu_1(X)^2], \\ \mathbf{E}[\psi(Y, T, X, \beta^*, p^*(\cdot)) \cdot \alpha(T, X)] &= -\mathbf{E}\left[\frac{\mu_1(X)^2}{p^*(X)}\right] + \mathbf{E}[\mu_1(X)^2], \end{aligned}$$

so that

$$\begin{aligned} \mathbf{E}[(\psi(Y, T, X, \beta^*, p^*(\cdot)) + \alpha(T, X))^2] &= \mathbf{E}[\mu_1(X)^2] - (\beta^*)^2 + \mathbf{E}\left[\frac{\sigma_1^2(X)}{p^*(X)}\right] \\ &= \mathbf{V}(\mathbf{E}[Y(1)|X]) + \mathbf{E}[\mathbf{V}(Y(1)|X)/p^*(X)], \end{aligned}$$

which is the variance in Theorem 1, specialized to the case with  $\mu_0(x) = \sigma_0^2(x) = 0$ .

In the proof of (38) we rewrite the difference by adding and subtracting a number of terms, so that we can bound the differences. We give the asymptotic order of all differences, which makes it easier to understand the role of the assumptions. We have

$$(39) \quad \sqrt{N}(\hat{\beta}_{ew} - \beta^*) = \frac{1}{\sqrt{N}} \sum_{i=1}^N \left( \frac{T_i Y_i}{\hat{p}_K(X_i)} - \frac{T_i Y_i}{p^*(X_i)} + \frac{T_i Y_i}{p^*(X_i)^2} (\hat{p}_K(X_i) - p^*(X_i)) \right)$$

$$(40) \quad + \frac{1}{\sqrt{N}} \sum_{i=1}^N \left( -\frac{T_i Y_i}{p^*(X_i)^2} (\hat{p}_K(X_i) - p^*(X_i)) \right. \\ \left. + \int_{\mathbf{x}} \frac{\mu_1(x)}{p^*(x)} (\hat{p}_K(x) - p^*(x)) dF_0(x) \right)$$

$$(41) \quad - \sqrt{N} \int_{\mathbf{x}} \frac{\mu_1(x)}{p^*(x)} (\hat{p}_K(x) - p^*(x)) dF_0(x) \\ - \frac{1}{\sqrt{N}} \sum_{i=1}^N \tilde{\delta}_K(X_i) \frac{T_i - p_K(X_i)}{\sqrt{p_K(X_i)(1 - p_K(X_i))}}$$

$$(42) \quad + \frac{1}{\sqrt{N}} \sum_{i=1}^N (\tilde{\delta}_K(X_i) - \delta_K(X_i)) \frac{T_i - p_K(X_i)}{\sqrt{p_K(X_i)(1 - p_K(X_i))}}$$

$$(43) \quad + \frac{1}{\sqrt{N}} \sum_{i=1}^N \left( \delta_K(X_i) \frac{T_i - p_K(X_i)}{\sqrt{p_K(X_i)(1 - p_K(X_i))}} \right. \\ \left. - \delta_0(X_i) \frac{T_i - p^*(X_i)}{\sqrt{p^*(X_i)(1 - p^*(X_i))}} \right)$$

$$(44) \quad + \frac{1}{\sqrt{N}} \sum_{i=1}^N \left\{ \left( \frac{T_i \cdot Y_i}{p^*(X_i)} - \beta^* \right) + \delta_0(X_i) \frac{T_i - p^*(X_i)}{\sqrt{p^*(X_i)(1 - p^*(X_i))}} \right\}.$$

In this expression  $F_0$  is the population cdf of  $X$  and

$$(45) \quad \tilde{\delta}_K(x) = - \int_{\mathbf{x}} \frac{\mu_1(z)}{p^*(z)} L'(R^K(z)' \tilde{\pi}_K) R^K(z)' dF_0(z) \tilde{\Sigma}_K^{-1} \sqrt{L'(R^K(x)' \pi_K^*)} R^K(x),$$

$$(46) \quad \delta_K(x) = - \int_{\mathbf{x}} \frac{\mu_1(z)}{p^*(z)} L'(R^K(z)' \pi_K^*) R^K(z)' dF_0(z) \Sigma_K^{-1} \sqrt{L'(R^K(x)' \pi_K^*)} R^K(x),$$

$$(47) \quad \delta_0(x) = - \frac{\mu_1(x)}{p^*(x)} \sqrt{p^*(X_i)(1 - p^*(X_i))},$$

with

$$\Sigma_K = E[R^K(X) R^K(X)' L'(R^K(X)' \pi_K^*)] \quad \text{and}$$

$$\tilde{\Sigma}_K = \frac{1}{N} \sum_{i=1}^N R^K(X_i) R^K(X_i)' L(R^K(X_i)' \tilde{\pi}_K),$$

and  $\tilde{\pi}_K$  between  $\hat{\pi}_K$  and  $\pi_K^*$ .

Note that (44) is equal to the linearized expression for  $\sqrt{N}(\hat{\beta}_{ew} - \beta^*)$ . To show that the estimator is indeed asymptotically linear, we must derive bounds on the terms (39)–(43). If a bound depends on both  $K$  and  $N$ , we derive the bound for sequences  $K(N)$  that go to  $\infty$  with  $N$ . Because during the derivation some restrictions on these sequences are imposed, the resulting bounds are not uniform in  $K$ . We have seen this type of argument in the derivation of the order of  $\|\hat{\pi}_{K(N)} - \pi_{K(N)}\|$  where we imposed the large sample identification condition  $\zeta(K(N))^4/N \rightarrow 0$ .

Below we present the bounds on the terms (39)–(43). Details for the derivations for these bounds are available from the authors (Hirano, Imbens, and Ridder (2002)). The bound for (39) is

$$\begin{aligned} & \left| \frac{1}{\sqrt{N}} \sum_{i=1}^N \left( \frac{T_i Y_i}{\hat{p}_K(X_i)} - \frac{T_i Y_i}{p^*(X_i)} + \frac{T_i Y_i}{p^*(X_i)^2} (\hat{p}_K(X_i) - p^*(X_i)) \right) \right| \\ &= O_p \left( \frac{\zeta(K)^3}{\sqrt{N}} \right) + O_p(\sqrt{N} \zeta(K)^2 K^{-\frac{5}{2}}) \\ &+ O_p(\zeta(K)^{5/2} K^{-\frac{5}{2r}}). \end{aligned}$$

The bound for (40) is

$$\begin{aligned} & \frac{1}{\sqrt{N}} \sum_{i=1}^N \left( - \frac{T_i Y_i}{p^*(X_i)^2} (\hat{p}_K(X_i) - p^*(X_i)) + \int_{\mathbf{x}} \frac{\mu_1(x)}{p^*(x)} (\hat{p}_K(x) - p^*(x)) dF_0(x) \right) \\ &= O_p(\zeta(K) K^{-\frac{5}{2r}}) + O_p \left( \frac{\zeta(K)^2}{\sqrt{N}} \right). \end{aligned}$$

The bound for (41) is

$$\begin{aligned} & \left| - \sqrt{N} \int_{\mathbf{x}} \frac{\mu_1(x)}{p^*(x)} (\hat{p}_K(x) - p^*(x)) dF_0(x) - \frac{1}{\sqrt{N}} \sum_{i=1}^N \tilde{\delta}_K(X_i) \frac{T_i - p_K(X_i)}{\sqrt{p_K(X_i)(1 - p_K(X_i))}} \right| \\ &= O(\sqrt{N} \zeta(K) K^{-\frac{5}{2r}}). \end{aligned}$$

The bound for (42) is

$$\left| \frac{1}{\sqrt{N}} \sum_{i=1}^N (\hat{\delta}_K(X_i) - \delta_K(X_i)) \frac{T_i - p_K(X_i)}{\sqrt{p_K(X_i)(1 - p_K(X_i))}} \right| = O_p \left( \frac{\zeta(K)^{9/2}}{N^{1/2}} \right).$$

The bound for (43) is

$$\begin{aligned} & \left| \frac{1}{\sqrt{N}} \sum_{i=1}^N \left( \delta_K(X_i) \frac{T_i - p_K(X_i)}{\sqrt{p_K(X_i)(1-p_K(X_i))}} - \delta_0(X_i) \frac{T_i - p^*(X_i)}{\sqrt{p^*(X_i)(1-p^*(X_i))}} \right) \right| \\ &= O_p(\max(K^{-\frac{1}{2} \frac{s}{r}}, \zeta(K) K^{-\frac{s}{2r}})). \end{aligned}$$

From these five expressions we obtain

$$\begin{aligned} (48) \quad & \left| \sqrt{N}(\hat{\beta}_{ew} - \beta^*) - \frac{1}{\sqrt{N}} \sum_{i=1}^N \left\{ \left( \frac{T_i Y_i}{p^*(X_i)} - \beta^* \right) - \frac{\mu_1(X_i)}{p^*(X_i)} (T_i - p^*(X_i)) \right\} \right| \\ &= O_p\left(\frac{\zeta(K)^3}{\sqrt{N}}\right) + O_p(\sqrt{N} \zeta(K)^2 K^{-\frac{s}{r}}) + O_p(\zeta(K)^{5/2} K^{-\frac{s}{2r}}) \\ &\quad + O_p(\zeta(K) K^{-\frac{s}{2r}}) + O_p\left(\frac{\zeta(K)^2}{\sqrt{N}}\right) + O(\sqrt{N} \zeta(K) K^{-\frac{s}{2r}}) \\ &\quad + O_p\left(\frac{\zeta(K)^{9/2}}{\sqrt{N}}\right) + O_p(\max(K^{-\frac{s}{2r}}, \zeta(K) K^{-\frac{s}{2r}})) \\ &= O_p(\sqrt{N} \zeta(K)^2 K^{-\frac{s}{r}}) + O_p(\zeta(K)^{5/2} K^{-\frac{s}{2r}}) + O_p\left(\frac{\zeta(K)^{9/2}}{\sqrt{N}}\right). \end{aligned}$$

Note that the second term of the final expression is a bias term, the third a variance term, and the first a combination of a variance and bias term.

As noted  $\zeta(K)$  depends on the sequence of approximating functions. For power series we have  $\zeta(K) = O(K)$ . If we consider sequences  $K(N) = N^c$  we can find the range of  $c$  for which (48) is  $o_p(1)$ . Substitution in the right-hand side of (48) gives that the first term on the right-hand side requires that  $c > 1/(2(s/r - 2))$ , the second that  $s/r > 5$ , and the third that  $c < 1/9$ . These inequalities can be simultaneously satisfied if  $s/r \geq 7$ . Q.E.D.

PROOF OF THEOREM 2: Define

$$(49) \quad \Psi_K = -\frac{1}{N} \sum_{i=1}^N \frac{Y_i T_i}{p^*(X_i)^2} R^K(X_i),$$

$$(50) \quad \hat{\Psi}_K = -\frac{1}{N} \sum_{i=1}^N \frac{Y_i T_i}{\hat{p}_K(X_i)^2} R^K(X_i),$$

$$(51) \quad \hat{S}_K = \frac{1}{N} \sum_{i=1}^N R^K(X_i) R^K(X_i)'$$

Then  $\Psi'_K \hat{S}_K^{-1} R^K(x)$  is the predicted value in a least squares series regression of  $-Y_i T_i / (p^*(X_i)^2)$  on  $R^K(X_i)$ .<sup>13</sup> This predicted value estimates  $-E(Y|x)/p^*(x)$ , which is the conditional expectation (given  $X = x$ ) of the derivative of the moment condition with respect to  $p^*$ . The usual bound for series estimators applies:

$$(52) \quad \sup_{x \in \mathbf{X}} \left| \Psi'_K \hat{S}_K^{-1} R^K(x) + \frac{\mu_1(x)}{p^*(x)} \right| \leq C_1 \zeta(K(N)) O_p\left(\sqrt{\frac{\zeta(K)}{N}}\right) + C_2 K^{-\frac{s}{r}}$$

<sup>13</sup> The number of terms in this series estimator need not be equal to that in the series estimator of the propensity score. The notation can be changed to reflect this.



with  $s'$  the number of continuous derivatives of  $\mu_1(x)$ . Also

$$(53) \quad \begin{aligned} \|\hat{\Psi}_K - \Psi_K\| &= \left\| \frac{1}{N} \sum_{i=1}^N \frac{(\hat{p}_K(X_i) - p^*(X_i))(p^*(X_i) + \hat{p}_K(X_i))}{\hat{p}_K(X_i)^2 p^*(X_i)^2} Y_i T_i R^K(X_i) \right\| \\ &\leq \frac{1}{N} \sum_{i=1}^N \left| \frac{p^*(X_i) + \hat{p}_K(X_i)}{\hat{p}_K(X_i)^2 p^*(X_i)^2} \right| \|\hat{p}_K(X_i) - p^*(X_i)\| \cdot |Y_i| \cdot T_i \cdot \|R^K(X_i)\|. \end{aligned}$$

As in the proof of Theorem 1 we have that  $\hat{p}_K(x)$  is bounded from 0 and 1 on  $\mathbf{X}$  if  $N \rightarrow \infty$  and hence we have the following bound for (53):

$$(54) \quad \begin{aligned} C \sup_{x \in \mathbf{X}} |\hat{p}_K(x) - p^*(x)| \sup_{x \in \mathbf{X}} \|R^K(x)\| \frac{1}{N} \sum_{i=1}^N |Y_i| + o_p(1) \\ = C_1 \zeta(K)^2 O_p\left(\sqrt{\frac{\zeta(K)}{N}}\right) + C_2 \zeta(K)^2 K^{-\frac{s'}{r}}. \end{aligned}$$

We use the bounds (52) and (54) to obtain a bound on

$$(55) \quad \begin{aligned} \hat{\alpha}_K(t, x) - \alpha(t, x) &= (\hat{\Psi}_K - \Psi_K)' \hat{S}_K^{-1} R^K(x) (t - \hat{p}_K(x)) \\ &\quad + \left[ \Psi_K' \hat{S}_K^{-1} R^K(x) + \frac{\mu_1(x)}{p^*(x)} \right] (t - \hat{p}_K(x)) + \frac{\mu_1(x)}{p^*(x)} (\hat{p}_K(x) - p^*(x)). \end{aligned}$$

Under the asymptotic identification condition

$$(56) \quad \begin{aligned} \sup_{x \in \mathbf{X}} |\hat{\alpha}_K(t, x) - \alpha(t, x)| &\leq C_1 \|\hat{\Psi}_K - \Psi_K\| \sup_{x \in \mathbf{X}} \|R^K(x)\| + C_2 \sup_{x \in \mathbf{X}} \left| \Psi_K' \hat{S}_K^{-1} R^K(x) + \frac{\mu_1(x)}{p^*(x)} \right| \\ &\quad + C_3 \sup_{x \in \mathbf{X}} \mu_1(x) \sup_{x \in \mathbf{X}} |\hat{p}_K(x) - p^*(x)|. \end{aligned}$$

Because  $\mathbf{X}$  is compact and  $\mu_1(x)$  is continuous,  $\sup_{x \in \mathbf{X}} \mu_1(x) < \infty$ . Substitution of the bounds (52) and (54), collecting terms of the same order and omitting terms of lower order gives the bound

$$(57) \quad \sup_{x \in \mathbf{X}} |\hat{\alpha}_K(t, x) - \alpha(t, x)| \leq C_1 \zeta(K)^3 O_p\left(\sqrt{\frac{\zeta(K)}{N}}\right) + C_2 \zeta(K)^3 K^{-\frac{s'}{r}} + C_3 K^{-\frac{s'}{r}}.$$

It can be shown that the difference between (14) and (13) is bounded by (57) (details of these calculations are available from the authors). Under the rates specified in Theorem 1 this bound is  $o_p(1)$ . Hence (14) is a consistent estimator for (13). *Q.E.D.*

**PROOF OF THEOREM 4:** The derivation of the efficiency bound follows the proof in Hahn (1998). We consider the case where the propensity score is known. From Theorem 3, it will be evident that the bound can be achieved without knowledge of the propensity score.

The density of  $(Y(0), Y(1), T, X)$  with respect to some  $\sigma$ -finite measure is

$$q(y(0), y(1), t, x) = f(y(0), y(1)|x) p^*(x)^t (1 - p^*(x))^{1-t} f(x).$$

The density of the observed data  $(y, t, x)$ , using the unconfoundedness assumption, is

$$q(y, t, x) = [f_1(y|x) p^*(x)]^t [f_0(y|x) (1 - p^*(x))]^{1-t} f(x),$$

where  $f_1(\cdot|x) = \int f(y(0), \cdot|x) dy(0)$ , and  $f_0(\cdot|x) = \int f(\cdot, y(1)|x) dy(1)$ . Consider a regular parametric submodel indexed by  $\theta$ , with density

$$q(y, t, x|\theta) = [f_1(y|x, \theta) p^*(x)]^t [f_0(y|x, \theta) (1 - p^*(x))]^{1-t} f(x|\theta),$$

which equals  $q(y, t, x)$  for  $\theta = \theta_0$ . Note that  $\theta$  does not enter into the term  $p^*(x)$ , because we are assuming that the propensity score is known. The score is given by

$$\frac{d}{d\theta} \ln q(y, t, x|\theta) = s(y, t, x|\theta) = t \cdot s_1(y|x, \theta) + (1-t) \cdot s_0(y|x, \theta) + s_x(x|\theta),$$

where

$$s_1(y|x, \theta) = \frac{d}{d\theta} \ln f_1(y|x, \theta),$$

$$s_0(y|x, \theta) = \frac{d}{d\theta} \ln f_0(y|x, \theta),$$

$$s_x(x|\theta) = \frac{d}{d\theta} \ln f(x|\theta).$$

The tangent space of the model is the set of functions

$$\mathcal{S} = \{t \cdot s_1(y, x) + (1-t) \cdot s_0(y, x) + s_x(x)\}$$

for  $s_1, s_0$ , and  $s_x$  satisfying

$$\int s_1(y, x) f_1(y|x) dy = 0, \quad \forall x,$$

$$\int s_0(y, x) f_0(y|x) dy = 0, \quad \forall x,$$

$$\int s_x(x) f(x) dx = 0.$$

We are interested in estimating

$$\tau_{\text{wate}} \equiv \frac{\iint g(x) y f_1(y|x) f(x) dy dx - \iint g(x) y f_0(y|x) f(x) dy dx}{\int g(x) f(x) dx}.$$

So for the parametric submodel indexed by  $\theta$ ,

$$\tau_{\text{wate}}(\theta) \equiv \frac{\iint g(x) y f_1(y|x, \theta) f(x|\theta) dy dx - \iint g(x) y f_0(y|x, \theta) f(x|\theta) dy dx}{\int g(x) f(x|\theta) dx}.$$

We need to find a function  $F_\tau(y, t, x)$  such that for all regular parametric submodels,

$$\frac{\partial \tau_{\text{wate}}(\theta_0)}{\partial \theta} = \mathbf{E}[F_\tau(Y, T, X) s(Y, T, X|\theta_0)].$$

First we calculate  $\partial \tau_{\text{wate}}(\theta)/\partial \theta$ . Let  $\mu_g \equiv \int g(x) f(x) dx$ . Then

$$\begin{aligned} \frac{\partial \tau_{\text{wate}}(\theta_0)}{\partial \theta} &= \frac{1}{\mu_g} \left[ \iint g(x) y s_1(y|x, \theta_0) f_1(y|x, \theta_0) f(x|\theta_0) dy dx \right. \\ &\quad \left. - \iint g(x) y s_0(y|x, \theta_0) f_0(y|x, \theta_0) f(x|\theta_0) dy dx \right] \\ &\quad + \frac{1}{\mu_g} \left[ \int g(x) \{ \mathbf{E}[Y(1) - Y(0)|X=x] - \tau_{\text{wate}} \} \right. \\ &\quad \left. \times s_x(x|\theta_0) f(x|\theta_0) dx \right]. \end{aligned}$$

The following choice for  $F_\tau$  satisfies the condition:

$$\begin{aligned} F_\tau(Y, T, X) &= \frac{T \cdot g(X)}{\mu_g \cdot p^*(x)} (Y - \mathbf{E}[Y(1)|X]) - \frac{(1-T) \cdot g(X)}{\mu_g \cdot (1-p^*(x))} (Y - \mathbf{E}[Y(0)|X]) \\ &\quad + \frac{g(X)}{\mu_g} (\mathbf{E}[Y(1) - Y(0)|X] - \tau_{wate}). \end{aligned}$$

Hence  $\tau_{wate}$  is pathwise differentiable. By Theorem 2, in Section 3.3 of Bickel, Klaassen, Ritov, and Wellner (1993), the variance bound is the expected square of the projection of  $F_\tau(Y, T, X)$  on  $\mathcal{S}$ . Since  $F_\tau \in \mathcal{S}$ , the variance bound is

$$\begin{aligned} \mathbf{E}[F_\tau(Y, T, X)^2] &= \mathbf{E}\left[\frac{g(X)^2}{(\mu_g)^2 p^*(X)} V(Y(1)|X)\right] + \mathbf{E}\left[\frac{g(X)^2}{(\mu_g)^2 (1-p^*(X))} V(Y(0)|X)\right] \\ &\quad + \mathbf{E}\left[\frac{g(X)^2}{(\mu_g)^2} (E(Y(1)|X) - E(Y(0)|X) - \tau_{wate})^2\right]. \quad Q.E.D. \end{aligned}$$

## REFERENCES

- ABADIE, A., AND G. IMBENS (2002): "Simple and Bias-Corrected Matching Estimators for Average Treatment Effects," NBER Technical Working Paper 283.
- ANGRIST, J. D., AND J. HAHN (1999): "When to Control for Covariates? Panel-Asymptotic Results for Estimates of Treatment Effects," NBER Technical Working Paper 241.
- BARNOW, B., G. CAIN, AND A. GOLDBERGER (1980): "Issues in the Analysis of Selectivity Bias," in *Evaluation Studies*, Vol. 5, ed. by E. Stromsdorfer and G. Farkas. San Francisco: Sage.
- BICKEL, P. J., C. A. J. KLAASSEN, Y. RITOV, AND J. A. WELLNER (1993): *Efficient and Adaptive Estimation for Semiparametric Models*. Baltimore: Johns Hopkins University Press.
- CHAMBERLAIN, G. (1987): "Asymptotic Efficiency in Estimation with Conditional Moment Restrictions," *Journal of Econometrics*, 34, 305–334.
- CREPON, B., F. KRAMARZ, AND A. TROGNON (1998): "Parameters of Interest, Nuisance Parameters and Orthogonality Conditions: An Application to Autoregressive Error Component Models," *Journal of Econometrics*, 82, 135–156.
- DEHEJIA, R., AND S. WAHBA (1999): "Causal Effects in Non-Experimental Studies: Re-Evaluating the Evaluation of Training Programs," *Journal of the American Statistical Association*, 94, 1053–1062.
- GEMAN, S., AND C. HWANG (1982): "Nonparametric Maximum Likelihood Estimation by the Method of Sieves," *Annals of Statistics*, 10, 401–414.
- HAHN, J. (1998): "On the Role of the Propensity Score in Efficient Semiparametric Estimation of Average Treatment Effects," *Econometrica*, 66, 315–331.
- HANSEN, L. (1982): "Large Sample Properties of Generalized Method of Moments Estimators," *Econometrica*, 50, 1029–1054.
- HECKMAN, J., H. ICHIMURA, AND P. TODD (1997): "Matching as an Econometric Evaluation Estimator: Evidence from Evaluating a Job Training Program," *Review of Economic Studies*, 64, 605–654.
- (1998): "Matching As An Econometric Evaluations Estimator," *Review of Economic Studies*, 65, 261–294.
- HECKMAN, J., H. ICHIMURA, J. SMITH, AND P. TODD (1998): "Characterizing Selection Bias Using Experimental Data," *Econometrica*, 66, 1017–1098.
- HECKMAN, J., AND R. ROBB (1985): "Alternative Methods for Evaluating the Impact of Interventions," in *Longitudinal Analysis of Labor Market Data*, ed. by J. Heckman and B. Singer. New York: Cambridge University Press.

- HELLERSTEIN, J., AND G. IMBENS (1999): "Imposing Moment Restrictions from Auxiliary Data by Weighting," *Review of Economics and Statistics*, 81, 1–14.
- HIRANO, K., G. IMBENS, AND G. RIDDER (2000): "Efficient Estimation of Average Treatment Effects Using the Estimated Propensity Score," NBER Technical Working Paper 251.
- (2002): "Addendum to: Efficient Estimation of Average Treatment Effects Using the Estimated Propensity Score," <http://elsa.berkeley.edu/users/imbens/>.
- IMBENS, G. (1997): "One-step Estimators in Overidentified Generalized Method of Moments Estimator," *Review of Economic Studies*, 64, 359–383.
- IMBENS, G., R. SPADY, AND P. JOHNSON (1998): "Information Theoretic Approaches to Inference in Moment Condition Models," *Econometrica*, 66, 333–357.
- KITAMURA, Y., AND M. STUTZER (1997): "An Information-Theoretic Alternative to Generalized Method of Moments Estimation," *Econometrica*, 65, 861–874.
- LANCASTER, T. (1990): "A Paradox in Choice-based Sampling," Mimeo, Department of Economics, Brown University.
- LECHNER, M. (1999): "Earnings and Employment Effects of Continuous Off-the-job Training in East Germany after Unification," *Journal of Business and Economic Statistics*, 17, 74–90.
- LITTLE, R., AND D. RUBIN (1987): *Statistical Analysis with Missing Data*. New York: Wiley.
- LORENTZ, G. (1986): *Approximation of Functions*. New York: Chelsea Publishing Company.
- NEWBY, W. (1994): "The Asymptotic Variance of Semiparametric Estimators," *Econometrica*, 62, 1349–1382.
- (1997): "Convergence Rates and Asymptotic Normality for Series Estimators," *Journal of Econometrics*, 79, 147–168.
- QIAN, H., AND P. SCHMIDT (1999): "Improved Instrumental Variables and Generalized Method of Moments Estimators," *Journal of Econometrics*, 91, 145–169.
- QIN, J., AND J. LAWLESS (1994): "Generalized Estimating Equations," *Annals of Statistics*, 22, 300–325.
- ROBINS, J., S. MARK, AND W. NEWBY (1992): "Estimating Exposure Effects by Modeling the Expectation of Exposure Conditional on Confounders," *Biometrics*, 48, 479–495.
- ROBINS, J., AND Y. RITOV (1997): "Towards a Curse of Dimensionality Appropriate (CODA) Asymptotic Theory for Semi-parametric Models," *Statistics in Medicine*, 16, 285–319.
- ROBINS, J., AND A. ROTNITZKY (1995): "Semiparametric Efficiency in Multivariate Regression Models with Missing Data," *Journal of the American Statistical Association*, 90, 122–129.
- ROBINS, J., A. ROTNITZKY, AND L. ZHAO (1995): "Analysis of Semiparametric Regression Models for Repeated Outcomes in the Presence of Missing Data," *Journal of the American Statistical Association*, 90, 106–121.
- ROSENBAUM, P. (1987): "Model-Based Direct Adjustment," *Journal of the American Statistical Association*, 82, 387–394.
- ROSENBAUM, P., AND D. RUBIN (1983): "The Central Role of the Propensity Score in Observational Studies for Causal Effects," *Biometrika*, 70, 41–55.
- (1985): "Reducing Bias in Observational Studies Using Subclassification on the Propensity Score," *Journal of the American Statistical Association*, 79, 516–524.
- ROTNITZKY, A., AND J. ROBINS (1995): "Semiparametric Regression Estimation in the Presence of Dependent Censoring," *Biometrika*, 82, 805–820.
- RUBIN, D. (1974): "Estimating Causal Effects of Treatments in Randomized and Nonrandomized Studies," *Journal of Educational Psychology*, 66, 688–701.
- (1976): "Inference and Missing Data," *Biometrika*, 63, 581–592.
- (1977): "Assignment to Treatment Group on the Basis of a Covariate," *Journal of Educational Statistics*, 2, 1–26.
- (1978): "Bayesian Inference for Causal Effects: the Role of Randomization," *Annals of Statistics*, 6, 34–58.
- RUBIN, D., AND N. THOMAS (1996): "Matching Using Estimated Propensity Scores: Relating Theory to Practice," *Biometrics*, 52, 249–264.

- WOOLDRIDGE, J. (1999): "Asymptotic Properties of Weighted M-Estimators for Variable Probability Samples," *Econometrica*, 67, 1385–1406.
- (2002): "Inverse Probability Weighted M-Estimators for Sample Selection, Attrition and Stratification," Institute for Fiscal Studies, Cemmap Working Paper cwp11/02.