# Optimal Doubly Robust Estimation of Heterogeneous Causal Effects

Edward H. Kennedy

Department of Statistics & Data Science Carnegie Mellon University

edward@stat.cmu.edu

#### Abstract

Heterogeneous effect estimation plays a crucial role in causal inference, with applications across medicine and social science. Many methods for estimating conditional average treatment effects (CATEs) have been proposed in recent years, but there are important theoretical gaps in understanding if and when such methods are optimal. This is especially true when the CATE has nontrivial structure (e.g., smoothness or sparsity). Our work contributes in several main ways. First, we study a two-stage doubly robust CATE estimator and give a generic model-free error bound, which, despite its generality, yields sharper results than those in the current literature. We apply the bound to derive error rates in nonparametric models with smoothness or sparsity, and give sufficient conditions for oracle efficiency. Underlying our error bound is a general oracle inequality for regression with estimated or imputed outcomes, which is of independent interest; this is the second main contribution. The third contribution is aimed at understanding the fundamental statistical limits of CATE estimation. To that end, we propose and study a local polynomial adaptation of double-residual regression. We show that this estimator can be oracle efficient under even weaker conditions, if used with a specialized form of sample splitting and careful choices of tuning parameters. These are the weakest conditions currently found in the literature, and we conjecture that they are minimal in a minimax sense. We go on to give error bounds in the non-trivial regime where oracle rates cannot be achieved. Some finite-sample properties are explored with simulations.

Keywords: conditional effects, influence function, minimax rate, nonparametric regression.

## 1 Introduction

Heterogeneous effect estimation plays a crucial role in causal inference, with applications across medicine and social science, e.g., improving understanding of variation, and informing policy or optimizing treatment decisions. The most common target parameter in this setup is the conditional average treatment effect (CATE) function,  $\mathbb{E}(Y^1 - Y^0 \mid X = x)$ , which measures the expected difference in outcomes had those with covariates X = x been treated versus not. The CATE is typically identified under standard causal assumptions (including no unmeasured confounding) as the difference between two regression functions,  $\tau(x) \equiv \mathbb{E}(Y \mid X = x, A = 1) - \mathbb{E}(Y \mid X = x, A = 0)$ .

Important early methods for estimating the CATE often employed semiparametric models, for example partially linear models assuming  $\tau(x)$  to be constant [Robins et al., 1992, Robinson, 1988], or structural nested models in which  $\tau(x)$  followed some known parametric form [Robins, 1994, van der Laan, 2006, van der Laan and Robins, 2003, Vansteelandt and Joffe, 2014]. These approaches reflect a commonly held conviction that the CATE may be more structured and simple than the rest of the data-generating process (with the zero treatment effect case an obvious example). This is similar in spirit to the common belief that interaction terms are more often zero than "main effects".

Recent years have seen a move towards more flexible estimators of  $\tau(x)$ . The first nonparametric model where the CATE has its own complexity separate from regression functions seems to be Example 4 of Robins et al. [2008]. van der Laan [2013] (Sections 3.1–3.2) and Luedtke and van der Laan [2016] proposed an important model-free "meta-algorithm" for estimating the CATE, a variant of which we study in this paper. The last 5–10 years has seen even more emphasis on nonparametrics and incorporating machine learning [Athey and Imbens, 2016, Foster and Syrgkanis, 2019, Foster et al., 2011, Hahn et al., 2020, Imai and Ratkovic, 2013, Künzel et al., 2019, Nie and Wager, 2017, Semenova and Chernozhukov, 2017, Shalit et al., 2017, Wager and Athey, 2018]. The present work is in a similar vein, focusing on (i) providing more flexible CATE estimators with stronger theoretical guarantees, and (ii) pushing forward our understanding of optimality and the fundamental limits of CATE estimation. At the start of each of the Sections 3, 4, and 5, we detail related work and describe how our results fit.

After describing the setup and presenting a simple motivating illustration in Section 2, we go on to present our three main contributions: (i) a model-free oracle inequality for regression with estimated pseudo-outcomes (given in Section 3); (ii) general conditions for oracle efficiency of a doubly robust estimator we term the DR-Learner, with applications to specific regression methods under smoothness and sparsity conditions (in Section 4); and (iii) a more refined analysis of a specialized estimator we call the lp-R-Learner, showing that faster rates can be achieved in the non-oracle regime, and giving a partial answer towards understanding the fundamental limits of CATE estimation (in Section 5).

## 2 Setup & Illustration

We assume access to an iid sample of observations of  $Z_i = (X_i, A_i, Y_i)$ , where  $X \in \mathbb{R}^d$  are covariates,  $A \in \{0, 1\}$  is a binary treatment or exposure, and  $Y \in \mathbb{R}^d$  an outcome of interest. The distribution of Z is indexed by the covariate distribution and nuisance functions:

$$\pi(x) = \mathbb{P}(A = 1 \mid X = x)$$

$$\mu_a(x) = \mathbb{E}(Y \mid X = x, A = a)$$

$$\eta(x) = \mathbb{E}(Y \mid X = x)$$

Our goal is estimation of the difference in regression functions under treatment versus control  $\tau(x) \equiv \mu_1(x) - \mu_0(x)$ . Under standard causal assumptions of no unmeasured confounding, consistency, and positivity or overlap ( $\epsilon \leq \pi \leq 1 - \epsilon$  wp1, which we assume throughout), the function  $\tau(x)$  also equals  $\mathbb{E}(Y^1 - Y^0 \mid X = x)$ , where  $Y^a$  is the counterfactual outcome under A = a. We refer to  $\mu_1(x) - \mu_0(x)$  as the CATE, noting that our results hold regardless of whether the causal assumptions do. The average treatment effect (ATE) is given by  $\mathbb{E}\{\tau(X)\}$ .

### 2.1 Notation

We use  $\mathbb{P}_n(f) = \mathbb{P}_n\{f(Z)\} = \frac{1}{n} \sum_i f(Z_i)$  as shorthand for sample averages. When  $x \in \mathbb{R}^d$  we let  $||x||^2 = \sum_j x_j^2$  denote the usual (squared) Euclidean norm, and for generic (possibly random) functions f we let  $||f||^2 = \int f(z)^2 d\mathbb{P}(z)$  denote the (squared)  $L_2(\mathbb{P})$  norm. We use the notation  $a \lesssim b$  to mean  $a \leq Cb$  for some universal constant C, and  $a \approx b$  to mean  $cb \leq a \leq Cb$  so that  $a \lesssim b$  and  $b \lesssim a$ . We let  $a_n \sim b_n$  mean  $a_n/b_n \to 1$  as  $n \to \infty$ .

At various points we refer to s-smooth functions, which we define as those contained in the Hölder class  $\mathcal{H}(s)$ . Intuitively, these are smooth functions that are close to their  $\lfloor s \rfloor$ -order Taylor approximations. More precisely,  $\mathcal{H}(s)$  is the set of functions  $f: \mathcal{X} \to \mathbb{R}$  that are  $\lfloor s \rfloor$ -times continuously differentiable with partial derivatives bounded, and for which

$$|D^m f(x) - D^m f(x')| \lesssim ||x - x'||^{s - \lfloor s \rfloor}$$

for all x, x' and  $m = (m_1, ..., m_d)$  such that  $\sum_j m_j = \lfloor s \rfloor$ , where  $D^m = \frac{\partial^{\lfloor s \rfloor}}{\partial_{x_1}^{m_1} ... \partial_{x_d}^{m_d}}$  is the multivariate partial derivative operator.

## 2.2 Simple Motivating Illustration

Consider a simple data-generating process where the covariates X are uniform on [-1,1],

$$\pi(x) = 0.5 + 0.4 \times sign(x)$$

and  $\mu_1(x) = \mu_0(x)$  are equal to the piecewise polynomial function defined on page 10 of Györfi et al. [2002], which is illustrated in Figure 1. Figure 1 also shows n = 1000 simulated data points from this data-generating process, approximately half of which are treated (shown on the left panel) and the other half untreated (shown on the right). Also shown are estimates of the corresponding  $\mu_1$  and  $\mu_0$  functions, using default tuning parameters with the smoothing.spline function in base R.

An interesting but probably common phenomenon occurs in this simple example. The individual regression functions are very non-smooth, and difficult to estimate well on their own; this is especially true in the region where there are fewer treated individuals. Thus the estimate  $\hat{\mu}_1$  tends to oversmooth on the left, where there are more untreated individuals; in contrast, the estimate  $\hat{\mu}_0$  tends to undersmooth on the right, where there are more treated individuals. This means a naive plug-in estimator of the CATE that simply takes the difference  $\hat{\mu}_1 - \hat{\mu}_0$  will be a poor and overly complex estimator of the true difference, which is not only a constant but zero.

In contrast, suppose for simplicity that the propensity scores  $\pi$  were known. Then a regression of the inverse-probability-weighted (IPW) pseudo-outcome  $\xi \equiv \frac{(A-\pi)Y}{\pi(1-\pi)}$  would, up to constants, behave just as an oracle estimator that had access to the actual counterfactual difference  $Y^1-Y^0$ , since  $\mathbb{E}(\xi\mid X=x)=\tau(x)$  exactly. Figure 2 shows results from this procedure, as well as two other more efficient and doubly robust versions described in subsequent sections, again all using default tuning parameter choices from smoothing.spline. For these simulated data, the doubly robust estimators are much more efficient than the IPW estimator, and do a much better job of adapting to the correct underlying simplicity of the true  $\tau$ .

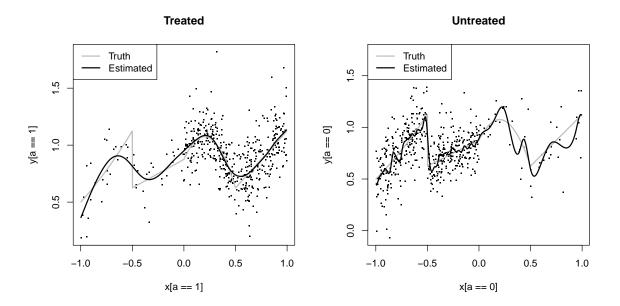


Figure 1: Plot of simulated data where the regression functions  $\mu_1$  and  $\mu_0$  individually are complex and difficult to estimate, but their difference is simply constant and equal to zero. Thus a naive plug-in estimator of the CATE will be overly complex, yielding large errors.

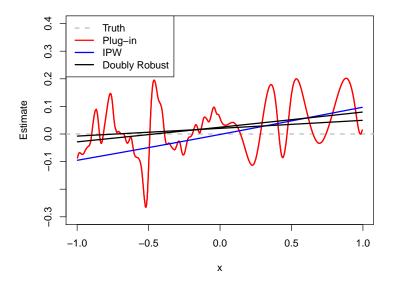


Figure 2: Estimated CATE curves in a simple simulated example. The plug-in method inherits the large errors from estimating the individual regression functions, which are complex and non-smooth. The IPW and doubly robust methods adapt to the smoothness of the CATE, which is constant in this example, with the doubly robust methods more efficient.

The results shown in Figure 2 are typical for this data-generating process: across 500 simulations, the IPW and doubly robust estimators gave smaller integrated squared bias across X by a factor of 10–100, respectively, and the doubly robust estimators improved on the integrated variance of the IPW estimator by nearly a factor of 20.

In the following sections we study the error of procedures like those illustrated above, giving model-free guarantees for practical use, as well as a more theoretical study of the fundamental limits of CATE estimation in a large nonparametric model.

## 3 Oracle Inequality for Pseudo-Outcome Regression

In this section we give a general oracle inequality for two-stage regression estimators that regress estimated "pseudo-outcomes" on a covariate vector. The inequality relates the error of such a procedure to that of an oracle with access to the true, unknown outcome. We give a finite-sample bound on the corresponding gap in performance, which we show in the next section can be informative in practice, despite making minimal assumptions about the regression procedure itself. In addition to laying a foundation for the analysis of a doubly robust CATE estimator in the next section, the result should also be of independent interest in other problems involving regression with estimated or imputed outcomes [Ai and Chen, 2003, Fan and Gijbels, 1994, Kennedy et al., 2017, Rubin and van der Laan, 2005, 2006].

With a few exceptions, previous work on regression with estimated outcomes has not appeared to exploit sample splitting, and has largely focused on particular pseudo-outcomes, and particular regression estimators (in both stages); in contrast we lean on sample splitting in order to obtain a more general result that is agnostic about the methods used, beyond some basic stability conditions. Prominent examples of previous work include Ai and Chen [2003], Rubin and van der Laan [2005], and Foster and Syrgkanis [2019], all of which gave results for pseudo-outcomes of a general form. Ai and Chen [2003] and Rubin and van der Laan [2005] did not use sample splitting and so limited their attention to particular estimators. Ai and Chen [2003] used sieves, and focused more on a finite-dimensional component appearing in the pseudo-outcome. Rubin and van der Laan [2005] considered least squares, penalized methods, and linear smoothers in the second stage, while being agnostic about the first-stage regression. However, their error bounds do not allow one to exploit double robustness in the pseudo-outcome, which is a crucial advantage in practice and which our result does allow.

Our results are closest to Foster and Syrgkanis [2019], who give a similar oracle inequality for generic empirical risk minimization when the loss involves complex nuisance functions. However there are some important distinctions to be made. Most importantly, Foster and Syrgkanis [2019] assume the loss satisfies a Neyman orthogonality property in order to obtain squared error rates; in contrast, our bound does not rely on this structure, but can exploit it when it holds, as shown in the following section. Crucially, our bound will also be seen to yield doubly robust errors, whereas the orthogonality-based results of Foster and Syrgkanis [2019] yield errors that are second-order but not doubly robust. Double robustness is essential when different nuisance components are estimated with different errors. Our somewhat sharper results are likely due to the fact that we focus on regression, and so can exploit more particular structure. These and some other differences are discussed further in the next section.

**Theorem 1.** Suppose  $Z_0^n = (Z_{01}, ..., Z_{0n})$  and  $Z^n = (Z_1, ..., Z_n)$  are independent training and test samples, respectively. Let  $\hat{f}(z) = \hat{f}(z; Z_0^n)$  be an estimate of a function f(z) using only the training data  $Z_0^n$ , and define  $m(x) \equiv \mathbb{E}\{f(z) \mid X = x\}$ .

Let  $\widehat{\mathbb{E}}_n(Y \mid X = x)$  denote a generic estimator of the regression function  $\mathbb{E}(Y \mid X = x)$ , using the test data  $(X_i, Y_i) \subseteq Z_i$ , i = 1, ..., n. Assume the estimator  $\widehat{\mathbb{E}}_n$  satisfies

1. 
$$\widehat{\mathbb{E}}_n(Y \mid X = x) + c = \widehat{\mathbb{E}}_n(Y + c \mid X = x)$$
 for any constant c.

2. If 
$$\mathbb{E}(W \mid X = x) = \mathbb{E}(Y \mid X = x)$$
 then

$$\mathbb{E}\left[\left\{\widehat{\mathbb{E}}_n(W\mid X=x) - \mathbb{E}(W\mid X=x)\right\}^2\right] \approx \mathbb{E}\left[\left\{\widehat{\mathbb{E}}_n(Y\mid X=x) - \mathbb{E}(Y\mid X=x)\right\}^2\right].$$

Let  $\widehat{m}(x) = \widehat{\mathbb{E}}_n\{\widehat{f}(Z) \mid X = x\}$  denote the regression of  $\widehat{f}(Z)$  on X in the test samples, and let  $\widetilde{m}(x) = \widehat{\mathbb{E}}_n\{f(Z) \mid X = x\}$  denote the corresponding oracle regression of f(Z) on X.

Define the error function  $\widehat{r}(x) = \widehat{r}(x; Z_0^n) \equiv \mathbb{E}\{\widehat{f}(Z) \mid X = x, Z_0^n\} - m(x)$ . Then

$$\mathbb{E}\left[\left\{\widehat{m}(x)-m(x)\right\}^2\right]\lesssim \mathbb{E}\left[\left\{\widetilde{m}(x)-m(x)\right\}^2\right]+\mathbb{E}\left\{\widehat{r}(x)^2\right\}.$$

Remark 1. The pointwise bound of Theorem 1 implies a corresponding result for the integrated mean squared error, i.e., that

$$\mathbb{E}\|\widehat{m} - m\|^2 \lesssim \mathbb{E}\|\widetilde{m} - m\|^2 + \int \mathbb{E}\left\{\widehat{r}(x)^2\right\} d\mathbb{P}(x)$$

under mild boundedness conditions.

Along with proofs of other main results, a proof of Theorem 1 can be found in Section 7. Sample splitting plays an important role here: note that the regression procedure as defined estimates the pseudo-outcome on a separate sample, independent from the one used in the second-stage regression via  $\hat{\mathbb{E}}_n$ . Examples of such procedures can be found in subsequent sections, e.g., as illustrated in Figures 3 and 5. The main role of sample splitting is that it allows for informative error analysis while being agnostic about the first- and second-stage methods.

Remark 2. With iid data, one can always obtain separate independent samples by randomly splitting the data in half (or in folds); further, to regain full sample size efficiency one can always swap the samples, repeat the procedure, and average the results, popularly called cross-fitting and used for example by Bickel and Ritov [1988], Robins et al. [2008], Zheng and van der Laan [2010], and Chernozhukov et al. [2018]. In this paper, to simplify notation we always analyze a single split procedure, with the understanding that extending to an analysis of an average across independent splits is straightforward.

The assumptions of Theorem 1 are both quite weak, and only reflect a mild form of stability of the second stage regression estimator  $\widehat{\mathbb{E}}_n$ . Assumption 1 says that adding a constant

to an outcome before doing a regression gives the same result as adding a constant post-regression. In fact, this assumption can be substantially weakened. For example, as long as  $\mathbb{E}[\{\widehat{\mathbb{E}}_n(Y\mid X=x)+c-\widehat{\mathbb{E}}_n(Y+c\mid X=x)\}^2] \lesssim \mathbb{E}[\{\widehat{\mathbb{E}}_n(Y\mid X=x)-\mathbb{E}(Y\mid X=x)\}^2]$ , i.e., the difference between adding a constant pre- versus post-regression is no more than the regression error itself, then the same oracle inequality holds. The stronger version of Assumption 1 is used in the theorem statement for simplicity and because it often holds in practice; for example, it holds for any linear smoother with so-called normal weights, which sum to one [Stone, 1977].

Assumption 2 of Theorem 1 says that the second-stage regression method gives the same error, up to constants, when applied to random variables with the same conditional means. This is also a very mild form of stability, since typical error bounds have rates that only depend on the smoothness or sparsity of the regression function, with other distributional features (e.g., conditional variance bounds) affecting only constants. For example, this appears to be true for all the methods analyzed in Györfi et al. [2002].

Thus under mild conditions, Theorem 1 indicates that the error of a sample-split pseudo-outcome regression procedure is at most that of the oracle, plus a particular error gap. This gap is the expected square of the error function  $\hat{r}$ , which is defined as the conditional bias of the pseudo-outcome estimate, given X and the training sample. In many important examples, including the CATE setup studied in detail in subsequent sections, the error function  $\hat{r}$  will itself take the form of a second-order product of errors. This will allow the pseudo-outcome regression to attain the oracle error rate, even when the nuisance quantities appearing in the estimated pseudo-outcome are estimated at slower rates.

### 4 DR-Learner

In this section we analyze a two-stage doubly robust estimator we call the DR-Learner, following the naming scheme from Nie and Wager [2017] and Künzel et al. [2019]. We first describe the algorithm in detail, then give model-agnostic error bounds which apply for arbitrary first-stage estimators, and as long as the second-stage estimators are mildly stable in a certain sense. We go on to apply the error bound in multiple nonparametric models incorporating smoothness or sparsity structure, and then explore the performance of the DR-Learner in simulations.

#### 4.1 Previous Work

Variants of the DR-Learner have been used before, though often tied to particular estimators and not incorporating double sample splitting, which allows for model-agnostic error bounds and reduced bias. van der Laan [2013] (Section 3.1–3.2) appears to be the first to propose the general DR-Learner approach, i.e., flexible regressions of the pseudo-outcome (1) below on covariates. Specifically, van der Laan [2013] and Luedtke and van der Laan [2016] advocate regressing the pseudo-outcome on  $V \subseteq X$  to construct candidate CATE estimators, and then selecting among them with a tailored cross-validation approach [van der Laan and Dudoit, 2003]. The main distinction with our work is they did not give specific error guarantees for the CATE (and did not use the same sample splitting scheme, though this is less important).

To the best of our knowledge, previous papers on the DR-Learner that do give specific error rates either employ stronger nuisance estimation conditions than we show are required, or else do not allow the CATE to be smoother than the regression functions. Lee et al. [2017] studied a local-linear version of the DR-Learner, but assumed the first-stage nuisance error was negligible. Semenova and Chernozhukov [2017] and Zimmert and Lechner [2019] studied series and local-constant variants, respectively, but required conditions on nuisance estimation that are as restrictive as for the ATE. Fan et al. [2019] also studied a local-constant variant, but did not consider the case where the CATE is smoother than the regression functions. As mentioned in the previous section, Foster and Syrgkanis [2019] considered a model-agnostic version of the DR-Learner, and gave results which are conceptually similar in the sense of giving a secondorder error bound relative to an oracle. A crucial difference, however, is that our error rates are doubly robust, meaning that our resulting conditions for oracle efficiency are weaker when the propensity score and outcome regressions are estimated at different rates. Künzel et al. [2019] considered various method-agnostic "meta-learners", but not the DR-Learner; further, none of their methods are doubly robust, and so in general would be expected to inherit larger error rates from individual  $\hat{\mu}_a$  estimators.

Remark 3. The DR-Learner approach is motivated by the fact that (1) is the (uncentered) efficient influence function for the ATE [Hahn, 1998, Robins and Rotnitzky, 1995]; this drives many of its favorable properties. For a review of influence functions and semiparametric theory we refer to van der Laan and Robins [2003], Tsiatis [2006], and Kennedy [2016].

## 4.2 Construction & Analysis

The algorithm below describes our proposed construction of the DR-Learner.

**Algorithm 1** (DR-Learner). Let  $(D_{1a}^n, D_{1b}^n, D_2^n)$  denote three independent samples of n observations of  $Z_i = (X_i, A_i, Y_i)$ .

Step 1. Nuisance training:

- (a) Construct estimates  $\hat{\pi}$  of the propensity scores  $\pi$  using  $D_{1a}^n$ .
- (b) Construct estimates  $(\widehat{\mu}_0, \widehat{\mu}_1)$  of the regression functions  $(\mu_0, \mu_1)$  using  $D_{1b}^n$ .

Step 2. Pseudo-outcome regression: Construct the pseudo-outcome

$$\widehat{\varphi}(Z) = \frac{A - \widehat{\pi}(X)}{\widehat{\pi}(X)\{1 - \widehat{\pi}(X)\}} \left\{ Y - \widehat{\mu}_A(X) \right\} + \widehat{\mu}_1(X) - \widehat{\mu}_0(X) \tag{1}$$

and regress it on covariates X in the test sample  $D_2^n$ , yielding

$$\widehat{\tau}_{dr}(x) = \widehat{\mathbb{E}}_n \{ \widehat{\varphi}(Z) \mid X = x \}. \tag{2}$$

Step 3. Cross-fitting (optional): Repeat Step 1–2 twice, first using  $(D_{1b}^n, D_2^n)$  for nuisance training and  $D_{1a}^n$  as the test sample, and then using  $(D_{1a}^n, D_2^n)$  for training and  $D_{1b}^n$  as the test sample. Use the average of the resulting three estimators as a final estimate of  $\tau$ .

Figure 3 gives a schematic illustrating the DR-Learner construction.

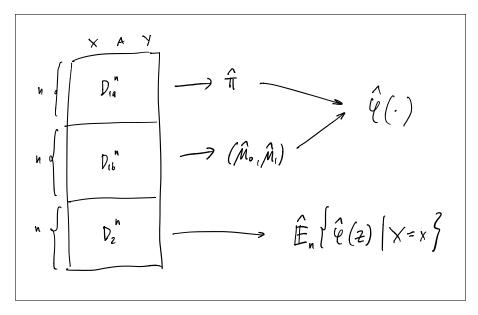


Figure 3: Schematic illustrating the DR-Learner approach. In the first stage, the nuisance functions  $\widehat{\pi}$  and  $(\widehat{\mu}_0, \widehat{\mu}_1)$  are estimated from training samples  $D_{1a}^n$  and  $D_{1b}^n$ , respectively. In the second stage, these estimates are used to construct an estimate of the pseudo-outcome  $\widehat{\varphi}$ , which is then regressed on X using the sample  $D_2^n$ .

Remark 4. The double sample splitting in Algorithm 1 (with  $\hat{\pi}$  and  $\hat{\mu}_a$  fit on separate samples) is similar to that used by Newey and Robins [2018], though there it is used with undersmoothing to reduce bias. Here we only use it to obtain a convenient product of mean squared errors, and do not require any undersmoothing. In Section 5, however, double sample splitting is used in the same spirit as Newey and Robins [2018].

Next we present the main result of this section, which gives error bounds for the DR-Learner procedure for any arbitrary first- and second-stage estimators, as long as the latter obey the mild stability conditions of Theorem 1.

**Theorem 2.** Let  $\hat{\tau}_{dr}(x)$  denote the DR-Learner estimator detailed in Algorithm 1. Assume:

- 1. The propensity score estimates are bounded as  $\epsilon \leq \widehat{\pi}(x) \leq 1 \epsilon$  for some  $\epsilon > 0$  wp1.
- 2. The second-stage regression  $\widehat{\mathbb{E}}_n(\cdot \mid X = x)$  satisfies Assumptions 1–2 of Theorem 1.

Let  $\widetilde{\tau}(x) = \widehat{\mathbb{E}}_n(Y^1 - Y^0 \mid X = x)$  denote an oracle estimator that regresses the difference  $(Y^1 - Y^0)$  on X, where  $Y^a$  is such that  $\mathbb{P}(Y^a \leq y \mid X = x) = \mathbb{P}(Y \leq y \mid X = x, A = a)$ .

Then

$$\mathbb{E}\left[\left\{\widehat{\tau}_{dr}(x) - \tau(x)\right\}^{2}\right] \lesssim R^{*}(x) + \mathbb{E}\left[\left\{\widehat{\pi}(x) - \pi(x)\right\}^{2}\right] \sum_{a=0}^{1} \mathbb{E}\left[\left\{\widehat{\mu}_{a}(x) - \mu_{a}(x)\right\}^{2}\right]$$

where  $R^*(x) = \mathbb{E}[\{\widetilde{\tau}(x) - \tau(x)\}^2]$  is the oracle risk.

The finite-sample bound on the DR-Learner error given in Theorem 2 shows that it can only deviate from the oracle error by at most the product of the mean squared errors of the propensity score and regression estimators, up to constants, thus allowing faster rates for estimating the CATE even when the nuisance estimates converge at slower rates. Importantly the result is agnostic about the methods used, and requires no special tuning or undersmoothing.

Theorem 2 gives a smaller risk bound compared to earlier work. Semenova and Chernozhukov [2017] and Zimmert and Lechner [2019] gave a remainder error larger than the product of the nuisance mean squared errors, with oracle efficiency requiring this product to shrink faster than 1/nd and 1/(n/h) for series and kernel-based second stage regressions, respectively (for h a shrinking bandwidth). Fan et al. [2019] assumed the CATE was only as smooth as the individual regression functions, giving a larger oracle risk. The results from Foster and Syrgkanis [2019] are closest to ours, but their error bounds are not doubly robust, and instead involve  $L_4$  errors of all nuisance components.

Remark 5. Several important oracle inequalities for cross-validated selection of estimators exist in the literature [Díaz et al., 2018, Luedtke and van der Laan, 2016, van der Laan and Dudoit, 2003]. These are relevant for cross-validated CATE estimation, but are conceptually different from our Theorem 2; for example, they use a different oracle. Namely, in cross-validated selection the oracle is the best performing among a group of learners, whereas in our setup the oracle is the specified learner  $\widehat{\mathbb{E}}_n$  when it is given access to the true difference  $Y^1-Y^0$ .

## 4.3 Examples & Illustrations

An important feature of the oracle inequality in Theorem 2 is that it is essentially model-free: it only requires that the second-stage regression estimator satisfies the mild stability assumptions of Theorem 1. In the following corollaries, we illustrate the flexibility of this result by applying it in settings where the nuisance functions and CATE are smooth or sparse (i.e., in settings where local polynomial, series, lasso, or random forest estimators would work well). Similar results could be immediately obtained in any model where bounds on mean squared error rates are known.

Corollary 1. Suppose the assumptions of Theorem 2 hold. Further assume:

- 1. The propensity score  $\pi$  is  $\alpha$ -smooth, and is estimated with squared error  $n^{-2\alpha/(2\alpha+d)}$ .
- 2. The regression functions  $\mu_a$  are  $\beta$ -smooth, and are estimated with squared error  $n^{-2\beta/(2\beta+d)}$ .
- 3. The CATE  $\tau$  is  $\gamma$ -smooth.

If the second-stage estimator  $\widehat{\mathbb{E}}_n$  yields the minimax optimal squared error rate  $n^{-2\gamma/(2\gamma+d)}$  for  $\gamma$ -smooth functions, then

$$\mathbb{E}\left[\left\{\widehat{\tau}_{dr}(x) - \tau(x)\right\}^{2}\right] \lesssim n^{\frac{-2\gamma}{2\gamma+d}} + n^{\frac{-2\alpha}{2\alpha+d} + \frac{-2\beta}{2\beta+d}}$$

and thus the DR-Learner is oracle efficient if

$$\alpha\beta \ge \frac{d^2}{4} - \frac{\left(\alpha + \frac{d}{2}\right)\left(\beta + \frac{d}{2}\right)}{\left(1 + \frac{2\gamma}{d}\right)},\tag{3}$$

which reduces to  $s \ge \frac{d/2}{1+d/\gamma}$  when  $\alpha = \beta = s$ .

Corollary (1) illustrates how the DR-Learner can adapt to smoothness in the CATE even when the propensity score and regression functions may be less smooth, and gives sufficient conditions for achieving the oracle rate  $n^{-2\gamma/(2\gamma+d)}$  depending on the nuisance smoothness and dimension d.

It is instructive to compare the sufficient condition in (3) to the analogous condition for root-n consistency of a standard doubly robust estimator of the average treatment effect, which is  $\sqrt{\alpha\beta} \geq d/2$  (cf. Equation 25 of Robins et al. [2009a]). First, as the CATE smoothness gets larger, the sufficient condition (3) for the CATE approaches that for the ATE, as should be expected (intuitively, an infinitely smooth CATE should be nearly as easy to estimate as the ATE). Second, the positive term subtracted from  $d^2/4$  in (3) can be interpreted as a "lowered bar" for optimal estimation, due to the fact that the oracle rate  $n^{\frac{-1}{2+d/\gamma}}$  is slower than root-n. This phenomenon was also noted in the dose-response estimation problem by Kennedy et al. [2017], and is in contrast with the error bounds given in Nie and Wager [2017] and Zimmert and Lechner [2019], which required ATE-like conditions for oracle efficiency of CATE estimators, via  $n^{-1/4}$  or faster rates on the nuisance estimators (note that if  $\alpha = \beta$  then  $\sqrt{\alpha\beta} \geq d/2$  means the nuisance error rate is faster than  $n^{-1/4}$ ). In the next section we will show how the sufficient condition (3) can even be improved upon.

The following theorem gives a result analogous to that in Corollary 1, but adapted to error rates found in sparse rather than smooth models.

Corollary 2. Suppose the assumptions of Theorem 2 hold. Further assume:

- 1. The propensity score  $\pi$  is  $\alpha$ -sparse and is estimated with squared error  $\frac{\alpha \log d}{n}$ .
- 2. The regression functions  $\mu_a$  are  $\beta$ -sparse and are estimated with squared error  $\frac{\beta \log d}{n}$ .
- 3. The CATE  $\tau$  is  $\gamma$ -sparse.

If the second-stage estimator  $\widehat{\mathbb{E}}_n$  yields squared error  $\frac{\gamma \log d}{n}$  for  $\gamma$ -sparse functions, then

$$\mathbb{E}\left[\left\{\widehat{\tau}_{dr}(x) - \tau(x)\right\}^{2}\right] \lesssim \frac{\gamma \log d}{n} + \frac{\alpha\beta \log^{2} d}{n^{2}}$$

and thus the DR-Learner is oracle efficient if

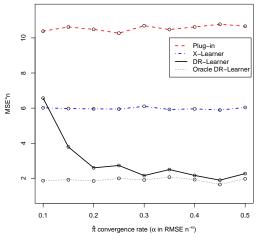
$$\alpha \beta \lesssim \frac{\gamma n}{\log d}.$$
 (4)

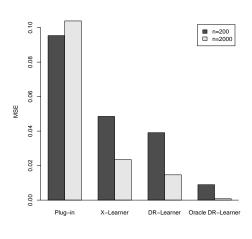
Sparsity conditions for root-n consistent estimation of the ATE with a standard doubly robust estimator are given in Farrell [2015], and are roughly of the form  $\alpha\beta \lesssim n/\log^2 d$ . In comparison, the sufficient condition (4) is less stringent: it allows the product of the nuisance sparsities to be larger by a factor  $\gamma \log d$ . For example if  $\gamma \sim \sqrt{n}/\log d$  so that the oracle rate is  $n^{-1/4}$ , then oracle efficiency for the CATE would be guaranteed as long as  $\alpha\beta \lesssim n\sqrt{n}/\log^2 d$ , which is  $\sqrt{n}$  times larger than the  $n/\log^2 d$  threshold required for the ATE. This is another illustration of the lower bar on nuisance errors for oracle efficient CATE versus ATE estimation.

### 4.4 Simulation Experiments

In this section we study some finite-sample properties via simulations (R code is available in Section 8). We use four methods for CATE estimation: a plug-in that estimates the regression functions  $\mu_0$  and  $\mu_1$  and takes the difference (called the T-Learner by Künzel et al. [2019]), the X-Learner from Künzel et al. [2019], the DR-Learner from Section 4, and an oracle DR-Learner that uses the true pseudo-outcome in the second-stage regression.

First we use the piecewise polynomial model from the motivating example in Section 2.2, with outcome and second-stage regressions fit using smoothing.spline in R. Figure 4a shows the mean squared error for the four CATE methods at n=2000 (based on 500 simulations with MSE averaged over 500 independent test samples), across a range of convergence rates for the propensity score estimator  $\hat{\pi}$ . To control the convergence rate we constructed this estimator as  $\hat{\pi} = \exp{it\{\log{it}(\pi) + \epsilon_n\}}$ , where  $\epsilon_n \sim N(n^{-\alpha}, n^{-2\alpha})$  so that  $RMSE(\hat{\pi}) \sim n^{-\alpha}$ . The results show that the plug-in estimator inherits the large error in estimating the individual regression functions, while the DR-Learner achieves much smaller errors and adapts to the smoothness of the CATE. The X-Learner has MSE in between the two. Consistent with Theorem 2, the MSE of the DR-Learner approaches that of the oracle as the propensity score estimation error decreases (i.e., as the convergence rate gets faster).





(a) Piecewise polynomial model from 2.2

(b) High-dimensional model

Figure 4: Simulation results.

Next we consider a high-dimensional logistic model where  $X=(X_1,...,X_d)\sim N(0,I_d)$ , and  $\operatorname{logit}\{\pi(x)\}=\frac{1}{2\sqrt{\alpha}}\sum_{j=1}^{\alpha}x_j$ , and  $\operatorname{logit}\{\mu_a(x)\}=\frac{1}{\sqrt{\beta}}\sum_{j=1}^{\beta}x_j$ , so the propensity score and individual regressions are  $\alpha$  and  $\beta$  sparse, respectively, while the CATE is zero. The normalization of the coefficients ensures  $\pi(X)\in[0.2,0.8]$  with high probability and similarly for  $\mu_a$ . We use standard cross-validation-tuned lasso for all model-fitting, i.e., to estimate the propensity scores, regression functions, and second-stage fits. Figure 4b shows mean squared errors for the four CATE methods when d=500 and  $\alpha=\beta=50$ , for n=200 and n=2000, across 100 simulations (median of mean square errors is reported due to high skewness). As expected from Theorem 2, the DR-Learner is closer to the oracle than the plug-in or X-Learner, though this high-dimensional setup yields larger estimation error than the previous simulation model. The relative performance of the DR-Learner seems to improve with sample size.

## 5 Local Polynomial R-Learner

The previous section gave general sufficient conditions under which the DR-Learner attains the oracle error rate of an estimator with direct access to the difference  $Y^1 - Y^0$ , showing that this rate can in fact be achieved whenever the product of nuisance errors is of smaller order. This raises the crucial question of what happens when this product is not sufficiently small: in such regimes, is there any hope at still attaining the oracle error rate?

The current section provides a first answer to this question, with a more refined analysis of a different estimator. Specifically, we analyze a double-sample-split local polynomial adaptation of the R-Learner [Nie and Wager, 2017, Robinson, 1988], which we call the lp-R-Learner for short. The R-Learner of Nie and Wager [2017] is a nonparametric RKHS regression-based extension of the double-residual regression method of Robinson [1988]. A nonparametric series-based version of the R-Learner was also proposed in Example 4 of Robins et al. [2008], though assuming known propensity scores and not incorporating outcome regression. Chernozhukov et al. [2017] studied a lasso version of the R-learner. Some important distinctions between our results and those in previous work include the following: (i) our estimator is built from local polynomials, and incorporates a specialized form of sample splitting inspired by Newey and Robins [2018] for bias reduction, (ii) our sufficient conditions for attaining oracle efficiency are substantially weaker than the  $n^{-1/4}$  rates in Nie and Wager [2017] and Chernozhukov et al. [2017], and (iii) we give specific rates of convergence outside the oracle regime.

We first describe the lp-R-Learner in detail, then give the main error bound result, which holds under a Hölder-smooth model, and is valid for a wide variety of tuning parameter choices. Following that, we optimize the bound with specific tuning parameter choices, under different sets of conditions, and discuss the resulting rates.

#### 5.1 Construction

The algorithm below describes the lp-R-Learner construction.

**Algorithm 2** (lp-R-Learner). Let  $(D_{1a}^n, D_{1b}^n, D_2^n)$  denote three independent samples of n observations of  $Z_i = (X_i, A_i, Y_i)$ .

Let  $b: \mathbb{R}^d \to \mathbb{R}^p$  denote the vector of basis functions consisting of all powers of each covariate, up to order  $\lfloor \gamma \rfloor$ , and all interactions up to degree  $\lfloor \gamma \rfloor$  polynomials (cf. Masry [1996]). Let  $K_{hx}(X) = \frac{1}{h^d}K\left(\frac{X-x}{h}\right)$  for  $K: \mathbb{R}^d \to \mathbb{R}$  a bounded kernel function with support  $[-1,1]^d$ , and h a bandwidth parameter.

#### Step 1. Nuisance training:

- (a) Using  $D_{1a}^n$ , construct estimates  $\widehat{\pi}_a$  of the propensity scores  $\pi$ .
- (b) Using  $D_{1b}^n$ , construct estimates  $\widehat{\eta}$  of the regression function  $\eta = \pi \mu_1 + (1 \pi)\mu_0$ , and estimates  $\widehat{\pi}_b$  of the propensity scores  $\pi$ .

Step 2. Localized double-residual regression:

Define  $\hat{\tau}_r(x)$  as the fitted value from a kernel-weighted least-squares regression (in the test sample  $D_2^n$ ) of outcome residual  $(Y - \hat{\eta})$  on basis terms b scaled by the treatment residual  $(A - \hat{\pi}_b)$ , with weights  $\left(\frac{A - \hat{\pi}_a}{A - \hat{\pi}_b}\right) K_{hx}$ . Thus  $\hat{\tau}_r(x) = b(0)^T \hat{\theta}$  for

$$\widehat{\theta} = \underset{\theta \in \mathbb{R}^p}{\operatorname{arg\,min}} \ \mathbb{P}_n \left( K_{hx}(X) \left\{ \frac{A - \widehat{\pi}_a(X)}{A - \widehat{\pi}_b(X)} \right\} \left[ \left\{ Y - \widehat{\eta}(X) \right\} - \theta^{\mathrm{T}} b(X - x) \left\{ A - \widehat{\pi}_b(X) \right\} \right]^2 \right).$$
(5)

Step 3. Cross-fitting (optional): Repeat Step 1–2 twice, first using  $(D_{1b}^n, D_2^n)$  for nuisance training and  $D_{1a}^n$  as the test sample, and then using  $(D_{1a}^n, D_2^n)$  for training and  $D_{1b}^n$  as the test sample. Use the average of the resulting three estimators of  $\tau$  as the final estimator  $\widehat{\tau}_r$ .

Remark 6. The kernel weights in the second step regression need to be multiplied by the ratio  $(A - \hat{\pi}_a)/(A - \hat{\pi}_b)$  in order to ensure the independence of relevant products of nuisance estimators (i.e., that they are built from separate samples  $D_{1a}^n$  and  $D_{1b}^n$ ). This allows for multiplicative biases and thus faster rates due to undersmoothing, first introduced by Newey and Robins [2018] but for  $\sqrt{n}$ -estimable functionals.

Figure 5 gives a schematic illustrating the lp-R-Learner construction.

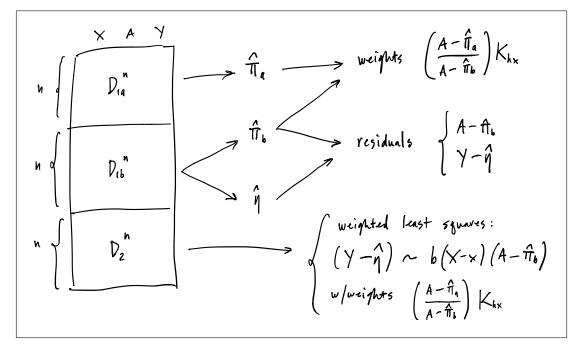


Figure 5: Schematic illustrating the lp-R-Learner approach. In the first stage, the nuisance functions  $\widehat{\pi}_a$  and  $(\widehat{\pi}_b, \widehat{\eta})$  are estimated from training samples  $D_{1a}^n$  and  $D_{1b}^n$ , respectively. In the second stage, these estimates are used in a kernel-weighted least squares regression of residuals  $(Y - \widehat{\eta})$  on residual-scaled basis terms  $(A - \widehat{\pi}_b)b$ , with weights  $\left(\frac{A - \widehat{\pi}_a}{A - \widehat{\pi}_b}\right)K_{hx}$ .

### 5.2 Main Error Bound & Oracle Results

Before giving the main error bound in this section, we first present the following condition on the nuisance estimators that we use in the analysis. At a high level, this condition requires the nuisance estimators to be linear smoothers with particular bias and variance bounds.

Condition 1. The nuisance estimators  $(\widehat{\pi}_a, \widehat{\pi}_b, \widehat{\eta})$  are linear smoothers of the form

$$\widehat{\pi}_{j}(x) = \sum_{i \in D_{1j}^{n}} w_{i\alpha}(x; X_{1j}^{n}) A_{i}$$

$$\widehat{\eta}(x) = \sum_{i \in D_{1b}^{n}} w_{i\beta}(x; X_{1b}^{n}) Y_{i}$$
(1a)

with weights  $w_{i}(x; X_{1}^{n})$  depending on tuning parameter k, satisfying

$$\left(\sum_{i=1}^{n} w_{i\alpha}(x; X_{1j}^{n})^{2}\right) \vee \left(\sum_{i=1}^{n} w_{i\beta}(x; X_{1j}^{n})^{2}\right) \lesssim \frac{k}{n}$$
(1b)

and yielding pointwise conditional bias bounds

$$\left| \mathbb{E}\{\widehat{\pi}_j(x) \mid X_{1j}^n\} - \pi(x) \right| \lesssim k^{-\alpha/d} \quad \text{and} \quad \left| \mathbb{E}\{\widehat{\eta}(x) \mid X_{1b}^n\} - \eta(x) \right| \lesssim k^{-\beta/d}. \tag{1c}$$

Conditions (1a)–(1c) are relatively standard. Many popular estimators take the form given in (1a), including kernel estimators and local polynomials, but also series estimators and smoothing splines, nearest neighbor matching, RKHS regression, Gaussian process regression, some versions of random forests (e.g., Mondrian and kernel forests), as well as weighted combinations or ensembles of such methods. However, greedy versions of random forests and locally adaptive methods would generally be excluded.

Condition (1b) holds for several prominent linear smoothers. For example, it holds for series estimators built from k basis terms, whenever the basis b(x) satisfies  $\sup_x \|b(x)\| \lesssim \sqrt{k}$ ; this includes for example the Fourier basis, spline series, CDV wavelets, and local polynomial partitioning series [Belloni et al., 2015]. Condition (1b) also holds for standard kernel or local polynomial estimators when taking the bandwidth parameter as  $h \sim k^{-d}$ , as shown for example in Proposition 1.13 of Tsybakov [2009].

Condition (1c) also has been shown to hold for series and local polynomial estimators, for example, when the underlying regression function is appropriately smooth. In particular, under standard conditions, (1c) would hold for these methods when the propensity score  $\pi$  is  $\alpha$ -smooth and the regression function  $\eta$  is  $\beta$ -smooth; we again refer to Belloni et al. [2015] and Tsybakov [2009] for a review of related results. For (1c) to hold uniformly over all  $x \in \mathcal{X}$  would typically mean these bounds would only hold up to log factors; however our result will only require the bias to be controlled locally, near the point at which the CATE is to be estimated.

The next result gives error bounds on the lp-R-Learner from Algorithm 2, under Hölder smoothness conditions.

**Theorem 3.** Let  $\hat{\tau}_r(x)$  denote the lp-R-Learner estimator detailed in Algorithm 2. Assume:

- 1. The estimator  $\hat{\eta}$  and observations Z are bounded, and X has bounded density.
- 2. The propensity score estimates satisfy  $\epsilon \leq \widehat{\pi}_i(x) \leq 1 \epsilon$  for some  $\epsilon > 0$ .
- 3. The eigenvalues of the sample design matrices  $\widehat{Q}$  and  $\widetilde{Q}$  defined in (12) are bounded away from zero in probability.
- 4. The nuisance estimators  $(\widehat{\pi}_a, \widehat{\pi}_b, \widehat{\eta})$  satisfy Condition 1, with the bias bounds holding for all x' such that  $||x' x|| \le h$ .

Let  $s = \frac{\alpha + \beta}{2}$  denote the average smoothness of the propensity score and regression function. Then, if the CATE  $\tau(x)$  is  $\gamma$ -smooth,

$$\widehat{\tau}_r(x) - \tau(x) = O_{\mathbb{P}}\left(h^{\gamma} + k^{-2s/d} + k^{-2\alpha/d} + \frac{1}{\sqrt{nh^d}}\left(1 + \frac{k}{n}\right)\right)$$

as long as  $\frac{k/n}{\sqrt{nh^d}} \to 0$ .

Before detailing the result and implications of Theorem 3, we first discuss the assumptions. The first part of Assumption 1 is mostly to simplify presentation, and could be weakened at the cost of added complexity; the second part ensuring X has bounded density is more crucial, but still mild. Assumption 2 is standard in the causal literature, and in theory could be guaranteed by simply thresholding propensity score estimates; however, extreme propensity values (i.e., positivity violations) are an important issue in practice, especially in the nonparametric and high-dimensional setup [D'Amour et al., 2017]. Assumption 3 is relatively standard (see, e.g., Assumption (LP1) of Tsybakov [2009]) but would restrict how n and h scale; when d is fixed, standard bandwidth choices should suffice. Assumption 4 is arguably most crucial, and does the most work in the proof (together with the specialized sample splitting); however, as detailed in the discussion of Condition 1 prior to the theorem statement, Assumption 4 uses standard conditions commonly found in the nonparametric regression literature.

Now we give some discussion and interpretation of the (in-probability) error bound of Theorem 3. The first three terms are the bias, and the last two are the variance (on the standard deviation scale). The bias has three components. The first  $h^{\gamma}$  bias term comes from the bias of an oracle estimator with access to the true propensity score and regression function, and matches the bias of an oracle with direct access to the difference  $Y^1 - Y^0$ . The other two bias terms come from nuisance estimation: the first  $k^{-2s/d}$  term is the product of the biases of the propensity score and regression estimators, whereas the second  $k^{-2\alpha/d}$  term is the squared bias of the propensity score estimator. If the propensity score is at least as smooth as the regression function, then the first  $k^{-2s/d}$  term will dominate. Some heuristic intuition about why these specific bias terms arise is as follows: by virtue of its least squares construction, the lp-R-Learner can be viewed as a product of the inverse of an " $X^TX$ -like" term involving products of  $\widehat{\pi}_a$  and  $\widehat{\pi}_b$ , and an " $X^TY$ -like" term involving products of  $\widehat{\pi}_a$  and  $\widehat{\eta}$ .

The variance has two components. As with the bias, the first  $1/\sqrt{nh^d}$  term is the standard deviation of an oracle estimator with access to the true nuisance functions. The second term

is a product of this oracle standard error with k/n, which is the additional variance coming from nuisance estimation (in fact k/n is the product of the standard deviations of the nuisance estimators). In standard regimes the tuning parameter k would have to be chosen of smaller order than n (e.g.,  $k \log k/n \to 0$  as in Belloni et al. [2015] and elsewhere) in order for Condition 1 to hold, making the variance contribution from nuisance estimation asymptotically negligible. This last point will be discussed further shortly.

Now we give conditions under which the oracle rate is achieved by the lp-R-Learner, which we conjecture are not only sufficient but also necessary conditions.

Corollary 3. Suppose the assumptions of Theorem 3 hold. Further assume  $\alpha \geq \beta$  so the propensity score is smoother than the regression function, and take

- 1.  $h \sim n^{-1/(2\gamma+d)}$ , and
- 2.  $k \sim n/\log^2 n$  so that  $k \log k/n \to 0$ .

Then, up to log factors,

$$\widehat{\tau}(x) - \tau(x) = O_{\mathbb{P}}\left(n^{-\gamma/(2\gamma+d)} + n^{-2s/d}\right)$$

and the oracle rate is achieved if the average nuisance smoothness satisfies  $s \geq \frac{d/4}{1+d/2\gamma}$ .

Corollary 3 shows that an undersmoothed lp-R-Learner can be oracle efficient under much weaker conditions than the error bound given in Theorem 2 would indicate for the DR-Learner. Note taking  $k \sim n/\log^2 n$  amounts to undersmoothing as much as possible: it drives down nuisance bias, letting the variance k/n go to zero only very slowly, since the contribution of the latter is asymptotically negligible for CATE estimation as long as  $k \lesssim n$ .

The result in Corollary 3 is remniscient of a similar phenomenon in marginal ATE estimation, where Robins et al. [2009b] showed that the condition  $s \geq d/4$  is necessary and sufficient for the existence of root-n consistent estimators of the ATE functional  $\mathbb{E}\{\mathbb{E}(Y\mid X,A=1)\}$ . Our result thus shows that  $s\geq d/4$  is also sufficient for oracle efficient estimation of the CATE, but now the oracle rate is  $n^{-\gamma/(2\gamma+d)}$  rather than root-n, and so there is in fact a weaker bar for oracle efficiency (namely  $s\geq \frac{d/4}{1+d/2\gamma}$ ), depending on the smoothness  $\gamma$ . At one extreme, when the CATE is infinitely smooth, the condition we give for oracle efficiency of the lp-R-Learner recovers the usual  $s\geq d/4$  condition for the ATE. At the other extreme, when the CATE is non-smooth, oracle efficiency can be much easier to achieve (e.g., if the CATE is only  $\gamma=1$ -smooth, it is only required that  $s\geq 1/2$  even for arbitrarily large dimension d). We conjecture that the condition  $s\geq \frac{d/4}{1+d/2\gamma}$  is not only sufficient but may also be necessary for oracle efficiency in the above Hölder model, making the proposed lp-R-Learner minimax optimal in this regime when  $\alpha\geq\beta$ ; however, we leave a proof of this to future work.

When  $s < \frac{d/4}{1+d/2\gamma}$  and the oracle rate is not achieved, the rate  $n^{-2s/d}$  is slower than the usual functional estimation rate  $n^{-4s/(4s+d)}$ , which is minimax optimal for the ATE if the covariate density is smooth enough [Robins et al., 2009b], and also occurs for simpler functionals like the expected density [Bickel and Ritov, 1988, Birgé and Massart, 1995]. To illustrate this gap, if s = d/8 then the rate in Corollary 3 is  $n^{-1/4}$  while the usual functional minimax rate is  $n^{-1/3}$ .

## 5.3 Faster Rates with Known Covariate Density

Here we briefly consider how the lp-R-Learner rates from Corollary 3 can be improved by exploiting structure in the covariate density, as in Robins et al. [2008, 2017]. This is somewhat paradoxical since the CATE itself does not depend on the covariate density.

Some intuition for a possible rate improvement is as follows. If the density of the covariates X was known, then one could construct nuisance estimators  $(\widehat{\pi}_a, \widehat{\pi}_b, \widehat{\eta})$  satisfying Condition 1 without any restriction on the tuning parameter k, such as the  $k \log k/n \to 0$  restriction employed in Corollary 3. For example, a series estimator  $\widehat{\eta}$  could satisfy Condition 1 for any choice of k, if it took the form

$$\widehat{\eta}(x) = b(x)^{\mathrm{T}} \Big[ \mathbb{E} \left\{ b(X)b(X)^{\mathrm{T}} \right\} \Big]^{-1} \mathbb{P}_n \{ b(X)Y \}$$
(6)

where b is a vector of appropriate basis functions (different from those used in the lp-R-Learner construction), and similarly for  $\widehat{\pi}$ . Intuitively this is because restrictions like  $k \log k/n \to 0$  are only required to ensure the inverse of the "sample" design matrix  $\mathbb{P}_n(bb^T)$  converges to a limit (in operator norm), whereas if the density of X was known, then one could simply compute the true design matrix  $\mathbb{E}(bb^T)$  exactly.

We conjecture however that the  $n^{-2s/d}$  rate from Corollary 3 may not be improvable in the non-random fixed X case. A fixed design setup was recently considered by Gao and Han [2020], though in a somewhat specialized model where the propensity scores have zero smoothness. In fact when  $\alpha = 0$  our rate matches their lower bound in the non-oracle regime.

The next result gives rates for the lp-R-Learner when there are no restrictions on the choice of nuisance tuning parameter k.

Corollary 4. Suppose the assumptions of Theorem 3 hold, and in particular suppose Assumption 4 holds without any restrictions on k (e.g., as if the density of X is known).

Then if  $\alpha \geq \beta$  the convergence rate in Theorem 3 is optimized by taking

$$h \sim n^{\frac{-3s}{2s\gamma+(s+\gamma)d}}$$
 and  $k \sim n^{\frac{3\gamma d/2}{2s\gamma+(s+\gamma)d}}$ 

which gives

$$\widehat{\tau}(x) - \tau(x) = O_{\mathbb{P}}\left(n^{-\gamma/(2\gamma+d)} + n^{\frac{-3s}{2s+d(1+s/\gamma)}}\right),\,$$

and so the oracle rate is achieved if the average nuisance smoothness satisfies  $s \ge \frac{d/4}{1+d/2\gamma}$ .

Remark 7. When the density of the covariates X is unknown but smooth and estimable at fast enough rates, we expect results similar to those in Corollary 4 to still hold. This phenomenon also occurs for the ATE functional [Robins et al., 2008, 2017]. However, here the analysis of the lp-R-Learner becomes much more complicated, so we leave this to future work.

When the nuisance tuning parameter k is unrestricted, one needs to balance all three dominant terms from Theorem 3: the two bias terms  $h^{\gamma}$  and  $k^{-2s/d}$ , as well as the variance term  $\frac{k/n}{\sqrt{nh^d}}$ . Notice the "elbow" at  $s \geq \frac{d/4}{1+d/2\gamma}$  occurs here even when the density is known, perhaps again suggesting that this condition for oracle efficiency may be both sufficient and necessary. Whether the rate  $n^{\frac{-3s}{2s+d(1+s/\gamma)}}$  is minimax optimal or not in the non-oracle regime is unknown; we will pursue this in future work.

Interestingly, the rate  $n^{\frac{-3s}{2s+d(1+s/\gamma)}}$  is slightly slower than the usual functional estimation rate  $n^{-4s/(4s+d)}$ . For example, when  $\gamma \to \infty$  and s=d/8, the CATE rate from Corollary 4 is  $n^{-3/10}$  whereas the classic functional estimation rate is  $n^{-1/3}$ . Therefore there do appear to be benefits of exploiting structure in the covariate density, but whether the gap between the improved rate of Corollary 4 and the usual functional rate can be closed is unclear. An interesting more philosophical question is whether structure in the covariate density should be exploited for CATE estimation in practical settings, since the CATE itself does not depend on the covariate density.

## 6 Discussion

In this paper we studied the problem of estimating conditional treatment effects, giving new results that apply to a wider variety of methods and that better exploit when the CATE itself is structured, compared to the current state-of-the-art. Sections 3 and 4 were more practically oriented, giving model-free yet informative error bounds for general regression with estimated outcomes, and the DR-Learner method for CATE estimation, along with examples from smooth and sparse models, illustrating the flexibility of Theorems 1 and 2. In contrast, Section 5 was more theoretically oriented, aiming instead at understanding the fundamental statistical limits of CATE estimation, i.e., the smallest possible achievable error. We derived upper bounds on this error with a specially constructed (and tuned) estimator called the lp-R-Learner. Namely we showed that in a Hölder model, oracle efficiency is possible under weaker conditions on the nuisance smoothness than indicated from the DR-Learner error bound given in Theorem 2 – which itself is an improvement over conditions given in previous work.

Figure 6 summarizes the various error rates in this paper graphically, in an illustrative setup where the dimension is d=20, CATE smoothness is  $\gamma=2d$ , and nuisance smoothness  $\alpha=\beta=s$  varies from 0–10. This shows the various gaps between methods/rates, in an interesting regime where the CATE is relatively smooth compared to the nuisance functions.

Our work raises some interesting open questions, some of the more immediate of which we list here for reference:

- 1. Can the error bounds in Theorem 1 and 2 be improved without committing to particular first- or second-stage methods?
- 2. What rates can be achieved by specialized sample splitting and tuning of a DR-Learner, rather than an R-Learner?
- 3. What are the analogous results of Corollaries 1 and 2 in non-smooth/sparse models?

#### Dimension d=20, CATE smoothness $\gamma$ =2d

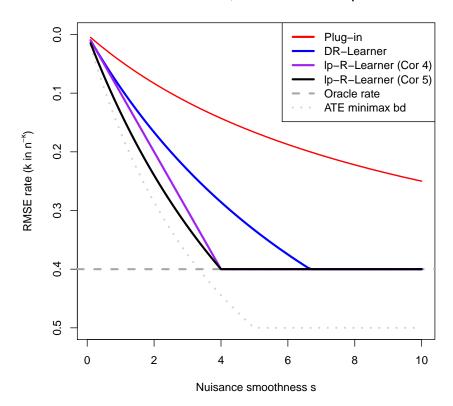


Figure 6: Illustration of error bounds in this work, as a function of nuisance smoothness  $s=\alpha=\beta$ . The dotted gray line represents the classical minimax lower bound for functional estimation, which is  $\sqrt{n}$  when  $s\geq d/4$ . The dashed gray line is the oracle rate that would be achieved by a minimax optimal regression using  $Y^1-Y^0$ . The red line is the bound for a plug-in estimator, which is just the rate for estimating the individual regression functions. The blue line is the rate for the DR-Learner given in Theorem 2, which matches the oracle when  $s\geq \frac{d/2}{1+d/\gamma}$ . The purple line is the rate achieved by the lp-R-Learner when tuned as in Corollary 3, which matches the oracle when  $s\geq \frac{d/4}{1+d/2\gamma}$ . The black line is the improved rate achieved by the lp-R-Learner when tuned as in Corollary 4, e.g., with known covariate density.

4. Can the error bounds in Theorem 3 be obtained with other methods? A natural alternative is a higher-order influence function approach [Robins et al., 2008, 2017], which could avoid the  $\alpha \geq \beta$  condition, perhaps at the cost of some extra complexity.

We have obtained partial answers to some of these questions, but most remain unanswered and left for future work. One of the deepest open questions is whether the rates given in Corollaries 3 and 4 are minimax optimal or not (in the fixed and random design setups, respectively).

## Acknowledgements

This research was supported by NSF Grant DMS1810979. The author thanks Sivaraman Balakrishnan, Jamie Robins, and Larry Wasserman for very helpful discussions.

## References

- C. Ai and X. Chen. Efficient estimation of models with conditional moment restrictions containing unknown functions. *Econometrica*, 71(6):1795–1843, 2003.
- S. Athey and G. Imbens. Recursive partitioning for heterogeneous causal effects. *Proceedings* of the National Academy of Sciences, 113(27):7353–7360, 2016.
- A. Belloni, V. Chernozhukov, D. Chetverikov, and K. Kato. Some new asymptotic theory for least squares series: pointwise and uniform results. *Journal of Econometrics*, 186(2): 345–366, 2015.
- P. J. Bickel and Y. Ritov. Estimating integrated squared density derivatives: sharp best order of convergence estimates. *Sankhyā*, pages 381–393, 1988.
- L. Birgé and P. Massart. Estimation of integral functionals of a density. *The Annals of Statistics*, 23(1):11–29, 1995.
- V. Chernozhukov, M. Goldman, V. Semenova, and M. Taddy. Orthogonal machine learning for demand estimation: High dimensional causal inference in dynamic panels. arXiv preprint arXiv:1712.09988, 2017.
- V. Chernozhukov, D. Chetverikov, M. Demirer, E. Duflo, C. Hansen, W. Newey, and J. M. Robins. Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal*, 21(1):C1–C68, 2018.
- A. D'Amour, P. Ding, A. Feller, L. Lei, and J. Sekhon. Overlap in observational studies with high-dimensional covariates. arXiv preprint arXiv:1711.02582, 2017.
- I. Díaz, O. Savenkov, and K. Ballman. Targeted learning ensembles for optimal individualized treatment rules with time-to-event outcomes. *Biometrika*, 105(3):723–738, 2018.
- J. Fan and I. Gijbels. Censored regression: local linear approximations and their applications. Journal of the American Statistical Association, 89(426):560–570, 1994.
- Q. Fan, Y.-C. Hsu, R. P. Lieli, and Y. Zhang. Estimation of conditional average treatment effects with high-dimensional data. arXiv preprint arXiv:1908.02399, 2019.
- M. H. Farrell. Robust inference on average treatment effects with possibly more covariates than observations. *Journal of Econometrics*, 189(1):1–23, 2015.
- D. J. Foster and V. Syrgkanis. Orthogonal statistical learning. arXiv preprint arXiv:1901.09036, 2019.
- J. C. Foster, J. M. Taylor, and S. J. Ruberg. Subgroup identification from randomized clinical trial data. *Statistics in Medicine*, 30(24):2867–2880, 2011.

- Z. Gao and Y. Han. Minimax optimal nonparametric estimation of heterogeneous treatment effects. arXiv preprint arXiv:2002.06471, 2020.
- L. Györfi, M. Kohler, A. Krzykaz, and H. Walk. A Distribution-Free Theory of Nonparametric Regression. Springer, 2002.
- J. Hahn. On the role of the propensity score in efficient semiparametric estimation of average treatment effects. *Econometrica*, 66(2):315–331, 1998.
- P. R. Hahn, J. S. Murray, C. M. Carvalho, et al. Bayesian regression tree models for causal inference: regularization, confounding, and heterogeneous effects. *Bayesian Analysis*, 2020.
- K. Imai and M. Ratkovic. Estimating treatment effect heterogeneity in randomized program evaluation. *The Annals of Applied Statistics*, 7(1):443–470, 2013.
- E. H. Kennedy. Semiparametric theory and empirical processes in causal inference. *In: Statistical Causal Inferences and Their Applications in Public Health Research*, pages 141–167, 2016.
- E. H. Kennedy, Z. Ma, M. D. McHugh, and D. S. Small. Nonparametric methods for doubly robust estimation of continuous treatment effects. *Journal of the Royal Statistical Society:* Series B, 79(4):1229–1245, 2017.
- S. R. Künzel, J. S. Sekhon, P. J. Bickel, and B. Yu. Metalearners for estimating heterogeneous treatment effects using machine learning. *Proceedings of the National Academy of Sciences*, 116(10):4156–4165, 2019.
- S. Lee, R. Okui, and Y.-J. Whang. Doubly robust uniform confidence band for the conditional average treatment effect function. *Journal of Applied Econometrics*, 32(7):1207–1225, 2017.
- A. R. Luedtke and M. J. van der Laan. Super-learning of an optimal dynamic treatment rule. *The International Journal of Biostatistics*, 12(1):305–332, 2016.
- E. Masry. Multivariate regression estimation local polynomial fitting for time series. *Stochastic Processes and their Applications*, 65(1):81–101, 1996.
- W. K. Newey and J. M. Robins. Cross-fitting and fast remainder rates for semiparametric estimation. arXiv preprint arXiv:1801.09138, 2018.
- X. Nie and S. Wager. Quasi-oracle estimation of heterogeneous treatment effects. arXiv preprint arXiv:1712.04912, 2017.
- J. M. Robins. Correcting for non-compliance in randomized trials using structural nested mean models. Communications in Statistics Theory and Methods, 23(8):2379–2412, 1994.
- J. M. Robins and A. Rotnitzky. Semiparametric efficiency in multivariate regression models with missing data. *Journal of the American Statistical Association*, 90(429):122–129, 1995.
- J. M. Robins, S. D. Mark, and W. K. Newey. Estimating exposure effects by modelling the expectation of exposure conditional on confounders. *Biometrics*, 48(2):479–495, 1992.
- J. M. Robins, L. Li, E. J. Tchetgen Tchetgen, and A. W. van der Vaart. Higher order influence functions and minimax estimation of nonlinear functionals. *Probability and Statistics: Essays in Honor of David A. Freedman*, pages 335–421, 2008.

- J. M. Robins, L. Li, E. Tchetgen Tchetgen, and A. W. van der Vaart. Quadratic semiparametric von mises calculus. *Metrika*, 69(2-3):227–247, 2009a.
- J. M. Robins, E. J. Tchetgen Tchetgen, L. Li, and A. W. van der Vaart. Semiparametric minimax rates. *Electronic Journal of Statistics*, 3:1305–1321, 2009b.
- J. M. Robins, L. Li, R. Mukherjee, E. Tchetgen Tchetgen, and A. W. van der Vaart. Minimax estimation of a functional on a structured high dimensional model. *The Annals of Statistics*, 45(5):1951–1987, 2017.
- P. M. Robinson. Root-n-consistent semiparametric regression. *Econometrica*, 56(4):931–954, 1988.
- D. B. Rubin and M. J. van der Laan. A general imputation methodology for nonparametric regression with censored data. *UC Berkeley Division of Biostatistics Working Paper Series*, 194, 2005.
- D. B. Rubin and M. J. van der Laan. Doubly robust censoring unbiased transformations. *UC Berkeley Division of Biostatistics Working Paper Series*, 208, 2006.
- V. Semenova and V. Chernozhukov. Estimation and inference about conditional average treatment effect and other structural functions. arXiv, pages arXiv-1702, 2017.
- U. Shalit, F. D. Johansson, and D. Sontag. Estimating individual treatment effect: generalization bounds and algorithms. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 3076–3085. JMLR. org, 2017.
- C. J. Stone. Consistent nonparametric regression. The Annals of Statistics, pages 595–620, 1977.
- A. A. Tsiatis. Semiparametric Theory and Missing Data. New York: Springer, 2006.
- A. B. Tsybakov. Introduction to Nonparametric Estimation. New York: Springer, 2009.
- M. J. van der Laan. Statistical inference for variable importance. The International Journal of Biostatistics, 2(1), 2006.
- M. J. van der Laan. Targeted learning of an optimal dynamic treatment, and statistical inference for its mean outcome. UC Berkeley Division of Biostatistics Working Paper Series, 317:1–90, 2013.
- M. J. van der Laan and S. Dudoit. Unified cross-validation methodology for selection among estimators and a general cross-validated adaptive epsilon-net estimator: Finite sample oracle inequalities and examples. UC Berkeley Division of Biostatistics Working Paper Series, Paper 130, 2003.
- M. J. van der Laan and J. M. Robins. *Unified Methods for Censored Longitudinal Data and Causality*. New York: Springer, 2003.
- S. Vansteelandt and M. M. Joffe. Structural nested models and g-estimation: the partially realized promise. *Statistical Science*, 29(4):707–731, 2014.
- S. Wager and S. Athey. Estimation and inference of heterogeneous treatment effects using random forests. *Journal of the American Statistical Association*, 113(523):1228–1242, 2018.

- W. Zheng and M. J. van der Laan. Asymptotic theory for cross-validated targeted maximum likelihood estimation. *UC Berkeley Division of Biostatistics Working Paper Series*, Paper 273:1–58, 2010.
- M. Zimmert and M. Lechner. Nonparametric estimation of causal heterogeneity under high-dimensional confounding. arXiv preprint arXiv:1908.08779, 2019.

## 7 Proofs

### 7.1 Proof of Theorem 1

First note that a debiased regression estimate  $\widehat{m}(x) - \widehat{r}(x)$  is equal to a single regression of  $\widehat{f}(Z) - \widehat{r}(x)$  on X, since by Assumption 1 it follows that

$$\widehat{m}(x) - \widehat{r}(x) = \widehat{\mathbb{E}}_n \{ \widehat{f}(Z) \mid X = x \} - \widehat{r}(x)$$

$$= \widehat{\mathbb{E}}_n \{ \widehat{f}(Z) - \widehat{r}(x) \mid X = x \}$$
(7)

because  $\hat{r}(x)$  is constant across the test data.

Next note that a regression of  $\widehat{f}(Z) - \widehat{r}(x)$  on X matches the oracle mean squared error up to constants. This follows since  $\mathbb{E}\{\widehat{f}(Z) - \widehat{r}(x) \mid X = x, Z_0^n\} = \mathbb{E}\{f(Z) \mid X = x\} = m(x)$  by definition of the error  $\widehat{r}(x)$ , so that

$$\mathbb{E}\{\widehat{f}(Z) - \widehat{r}(x) \mid X = x\} = \mathbb{E}[\mathbb{E}\{\widehat{f}(Z) - \widehat{r}(x) \mid X = x, Z_0^n\} \mid X = x] = m(x)$$

by iterated expectation.

Therefore since  $\widehat{f}(Z) - \widehat{r}(x)$  has the same conditional mean as f(Z), Assumption 2 implies

$$\mathbb{E}\left(\left[\widehat{\mathbb{E}}_n\{\widehat{f}(Z) - \widehat{r}(x) \mid X = x\} - m(x)\right]^2\right) \lesssim \mathbb{E}\left(\left[\widehat{\mathbb{E}}_n\{f(Z) \mid X = x\} - m(x)\right]^2\right). \tag{8}$$

Hence we have

$$\mathbb{E}\left(\left[\widehat{\mathbb{E}}_{n}\left\{\widehat{f}(Z)\mid X=x\right\} - \mathbb{E}\left\{f(Z)\mid X=x\right\}\right]^{2}\right)$$

$$= \mathbb{E}\left(\left[\widehat{\mathbb{E}}_{n}\left\{\widehat{f}(Z) - \widehat{r}(x)\mid X=x\right\} - \mathbb{E}\left\{f(Z)\mid X=x\right\} + \widehat{r}(x)\right]^{2}\right)$$

$$\leq 2\mathbb{E}\left(\left[\widehat{\mathbb{E}}_{n}\left\{\widehat{f}(Z) - \widehat{r}(x)\mid X=x\right\} - \mathbb{E}\left\{f(Z)\mid X=x\right\}\right]^{2}\right) + 2\mathbb{E}\left\{\widehat{r}(x)^{2}\right\}$$

$$\lesssim \mathbb{E}\left(\left[\widehat{\mathbb{E}}_{n}\left\{f(Z)\mid X=x\right\} - \mathbb{E}\left\{f(Z)\mid X=x\right\}\right]^{2}\right) + \mathbb{E}\left\{\widehat{r}(x)^{2}\right\}$$

where the first equality follows from (7) by Assumption 1, the second since  $(a+b)^2 \le 2(a^2+b^2)$ , and the third from (8) by Assumption 2.

## 7.2 Proof of Theorem 2

This result follows by an application of Theorem 1. The conditions all hold by assumption or design, so the main step is bounding the expected square of the error function  $\widehat{r}(x)$ .

By definition  $\widehat{r}(x) = \mathbb{E}\{\widehat{\varphi}(Z) \mid X = x, Z_0^n\} - \tau(x)$ . By iterated expectation we have

$$\widehat{r}(x) = \left\{ \frac{\pi(x)}{\widehat{\pi}(x)} - 1 \right\} \left\{ \mu_1(x) - \widehat{\mu}_1(x) \right\} - \left\{ \frac{1 - \pi(x)}{1 - \widehat{\pi}(x)} - 1 \right\} \left\{ \mu_0(x) - \widehat{\mu}_0(x) \right\}$$

so that

$$\widehat{r}(x)^{2} \leq 2 \left\{ \frac{\pi(x)}{\widehat{\pi}(x)} - 1 \right\}^{2} \left\{ \mu_{1}(x) - \widehat{\mu}_{1}(x) \right\}^{2} + 2 \left\{ \frac{1 - \pi(x)}{1 - \widehat{\pi}(x)} - 1 \right\}^{2} \left\{ \mu_{0}(x) - \widehat{\mu}_{0}(x) \right\}^{2}$$

$$\leq \left( \frac{2}{\epsilon^{2}} \right) \left\{ \pi(x) - \widehat{\pi}(x) \right\}^{2} \left[ \left\{ \mu_{1}(x) - \widehat{\mu}_{1}(x) \right\}^{2} + \left\{ \mu_{0}(x) - \widehat{\mu}_{0}(x) \right\}^{2} \right]$$

where the first inequality follows since  $(a+b)^2 \leq 2(a^2+b^2)$ , and the second by the bound on the propensity scores. Therefore the result follows by the independence of  $\widehat{\pi}$  and  $(\widehat{\mu}_0, \widehat{\mu}_1)$ .

### 7.3 Proof of Theorem 3

To simplify notation, in this subsection we largely omit function arguments, for example writing  $\hat{\tau} = \hat{\tau}(x)$ , b = b(X - x),  $K_{hx} = K_{hx}(X)$ , etc. We also define  $b_0 = b(0)$  and

$$\widehat{\phi} = \widehat{\phi}(Z) = (A - \widehat{\pi}_a)(Y - \widehat{\eta}) \qquad \qquad \phi = \phi(Z) = (A - \pi)(Y - \eta)$$

$$\widehat{\varphi}_a = \widehat{\varphi}_a(Z) = (A - \widehat{\pi}_a)(A - \widehat{\pi}_b) \qquad \qquad \nu = \nu(X) = \pi(1 - \pi).$$

We let  $X^n$  denote all the covariates across the training and test samples  $(D_{1a}^n, D_{1b}^n, D_2^n)$ , and we let  $D_1^n = (D_{1a}^n, D_{1b}^n)$  denote the training data.

First note that by definition

$$\widehat{\tau} = b_0^{\mathrm{T}} \mathbb{P}_n (bK_{hx} \widehat{\varphi}_a b^{\mathrm{T}})^{-1} \mathbb{P}_n (bK_{hx} \widehat{\phi}).$$

Thus we begin with the central decomposition

$$\widehat{\tau} - \tau = b_0^{\mathrm{T}} \mathbb{P}_n (bK_{hx} \nu b^{\mathrm{T}})^{-1} \mathbb{P}_n (bK_{hx} \phi) - \tau$$

$$+ b_0^{\mathrm{T}} \mathbb{P}_n (bK_{hx} \nu b^{\mathrm{T}})^{-1} \mathbb{P}_n \{bK_{hx} (\widehat{\phi} - \phi)\}$$

$$+ b_0^{\mathrm{T}} \left\{ \mathbb{P}_n (bK_{hx} \widehat{\varphi}_a b^{\mathrm{T}})^{-1} - \mathbb{P}_n (bK_{hx} \nu b^{\mathrm{T}})^{-1} \right\} \mathbb{P}_n (bK_{hx} \widehat{\phi})$$

$$\equiv b_0^{\mathrm{T}} \widetilde{Q}^{-1} \mathbb{P}_n (bK_{hx} \phi) - \tau$$

$$+ b_0^{\mathrm{T}} \widetilde{Q}^{-1} \mathbb{P}_n \{bK_{hx} (\widehat{\phi} - \phi)\}$$

$$+ b_0^{\mathrm{T}} \left( \widehat{Q}^{-1} - \widetilde{Q}^{-1} \right) \mathbb{P}_n (bK_{hx} \widehat{\phi})$$

$$(10)$$

where we define

$$\widetilde{Q} \equiv \mathbb{P}_n \Big\{ b(X - x) K_{hx}(X) \nu(X) b(X - x)^{\mathrm{T}} \Big\}$$

$$\widehat{Q} \equiv \mathbb{P}_n \Big\{ b(X - x) K_{hx}(X) \widehat{\varphi}_a(Z) b(X - x)^{\mathrm{T}} \Big\}.$$
(12)

The general approach in this proof is to use conditional error bounds for each term in the decomposition above, from which one arrives at bounds in probability (cf. Lemma 6.1 of Chernozhukov et al. [2018]).

The term on the right-hand side of (9) is the difference between  $\tau$  and an oracle version of the lp-R-Learner that has access to the true nuisance functions and so is built from  $\nu$  and  $\phi$ ; we

will show that it attains the same order as the oracle rate  $n^{-\gamma/(2\gamma+d)}$ . The term (10) captures the error from estimating  $(\widehat{\pi}_b, \widehat{\eta})$  in  $\widehat{\phi}$ ; we will show its order is the product of the biases from estimating  $(\widehat{\pi}_b, \widehat{\eta})$  plus a typically smaller variance term. The term (11) captures the error from estimating  $(\widehat{\pi}_a, \widehat{\pi}_b)$  in  $\widehat{\phi}_a$ ; we will show it behaves similarly to (10), except involving the product of the biases of  $(\widehat{\pi}_a, \widehat{\pi}_b)$ . In regimes where the oracle rate is not achievable and the propensity score  $\pi$  is smoother than the regression function  $\mu$ , we will show that the term (10) dominates.

#### 7.3.1 Term (9)

The analysis of the term in (9) follows that of a standard local polynomial estimator of pseudooutcome  $\phi/\nu$  on X with a special choice of kernel. This is because we can write

$$b_0^{\mathrm{T}} \mathbb{P}_n(bK_{hx}\nu b^{\mathrm{T}})^{-1} \mathbb{P}_n(bK_{hx}\phi) = \sum_{i=1}^n w_i(x; X^n) \frac{\phi(Z_i)}{\nu(X_i)}$$

where the weights are

$$w_i(x; X^n) \equiv \frac{1}{n} b(0)^{\mathrm{T}} \widetilde{Q}^{-1} b(X_i - x) K_{hx}(X_i) \nu(X_i).$$
 (13)

Thus the oracle estimator in (9) is a local polynomial estimator of the regression of  $\phi/\nu$  on X using scaled kernel function  $K_{hx}\nu$ , which has the same support as  $K_{hx}$  and has a smaller upper bound since  $\nu \leq 1/4$ .

The above implies that the weights  $w_i$  satisfy analogues of Proposition 1.12 and Lemma 1.3 in Tsybakov [2009], i.e., that they reproduce polynomials in X up to degree  $\lfloor \gamma \rfloor$ , and have the localizing properties given in the following lemma.

**Lemma 1.** Assume X is bounded, the kernel K satisfies  $K(u) \leq K_{max}\mathbb{1}(\|u\| \leq 1)$ , the eigenvalues of  $\widetilde{Q}$  are bounded below by  $\lambda_n$ , and  $\mathbb{P}_n\{\mathbb{1}(\|X-x\| \leq h)\} \lesssim v_n h^d$ . Then the weights in (13) satisfy:

- 1.  $\sup_{i,x} |w_i(x;X^n)| \lesssim \frac{1/\lambda_n}{nh^d}$ .
- 2.  $\sum_{i=1}^{n} |w_i(x; X^n)| \lesssim v_n/\lambda_n$ .
- 3.  $w_i(x; X^n) = 0$  when  $||X_i x|| \le h$ .

*Proof.* For (1) note

$$\sup_{i,x} |w_i(x;X^n)| \le \frac{1}{4nh^d} \left| \left| \widetilde{Q}^{-1}b(X_i - x)K\left(\frac{X_i - x}{h}\right) \right| \right|$$

$$\le \frac{K_{max}/\lambda_n}{4nh^d} \left| \left| b(X_i - x) \right| \left| \mathbb{1}\left( \left| \left| X_i - x \right| \right| \le h \right) \lesssim \frac{1/\lambda_n}{nh^d}$$

by the fact that ||b(0)|| = 1 and  $\nu \le 1/4$ , and from Assumptions 1–2. Similarly for (2)

$$\sum_{i=1}^{n} |w_i(x; X^n)| \le \frac{K_{max}/\lambda_n}{4nh^d} \sum_{i=1}^{n} ||b(X_i - x)|| \, \mathbb{1}\Big( \, ||X_i - x|| \le h \Big) \lesssim v_n/\lambda_n$$

where the first equality used the same logic as in (1) and the second Assumption 3. Property (3) follows since  $K(u) \lesssim \mathbb{1}(\|u\| \leq 1)$ .

Note also that  $\mathbb{E}(\phi/\nu \mid X) = \tau$  since

$$\mathbb{E}\{\phi(Z) \mid X\} = \mathbb{E}\Big[\{A - \pi(X)\}\{Y - \eta(X)\} \mid X\Big] = \pi(X)\{1 - \pi(X)\}\tau(X)$$

by iterated expectation. Therefore by the same logic as in Proposition 1.13 of Tsybakov [2009], using Lemma 1 with the Hölder condition on  $\tau$ , we have

$$\mathbb{E}\Big\{b_0^{\mathsf{T}} \mathbb{P}_n(bK_{hx}\nu b^{\mathsf{T}})^{-1} \mathbb{P}_n(bK_{hx}\phi) - \tau \mid X^n\Big\} = \sum_{i=1}^n w_i(x; X^n) \Big\{\tau(X_i) - \tau(x)\Big\}$$

$$= \sum_{i=1}^n w_i(x; X^n) \Big[\Big\{\tau_{\gamma, x}(X_i) - \tau(x)\Big\} + \Big\{\tau_{\gamma, x}(X_i) - \tau(X_i)\Big\}\Big]$$

$$\lesssim 0 + \sum_{i=1}^n |w_i(x; X^n)| ||X_i - x||^{\gamma} \mathbb{1}(||X_i - x|| \le h) \lesssim h^{\gamma} \left(\frac{v_n}{\lambda_n}\right)$$

where in the second line  $\tau_{\gamma,x}(X_i) = \sum_{|j| \leq \lfloor \gamma \rfloor} \frac{(X_i - x)^j}{j!} D^j \tau(x)$  is the  $\lfloor \gamma \rfloor$ -order multivariate Taylor approximation of  $\tau$  at x evaluated at  $X_i$ , the third line follows by the polynomial-reproducing property of  $w_i$ , the Hölder assumption on  $\tau$ , and Property 1 of Lemma 1, and the last line by Property 2 of Lemma 1.

For the variance we have

$$\operatorname{var}\left\{b_0^{\mathrm{T}} \mathbb{P}_n(bK_{hx}\nu b^{\mathrm{T}})^{-1} \mathbb{P}_n(bK_{hx}\phi) \mid X^n\right\} = \sum_{i=1}^n w_i(x;X^n)^2 \operatorname{var}(\phi \mid X = X_i)$$

$$\lesssim \sup_{i,x} |w_i(x;X^n)| \sum_{i=1}^n |w_i(x;X^n)| \lesssim \frac{1}{nh^d} \left(\frac{v_n}{\lambda_n^2}\right)$$

since the variance of  $\phi$  is bounded, and using Properties 1–2 of Lemma 1. Therefore term (9) satisfies

$$\mathbb{E}\left[\left\{b_0^{\mathrm{T}}\widetilde{Q}^{-1}\mathbb{P}_n(bK_{hx}\phi) - \tau\right\}^2 \mid X^n\right] \lesssim \left(h^{2\gamma} + \frac{1}{nh^d}\right) \left(\frac{v_n}{\lambda_n}\right)^2. \tag{14}$$

### $7.3.2 \quad \text{Term} \ (10)$

Now we bound the conditional mean and variance of term (10). First note we have

$$\mathbb{E}(\widehat{\phi} - \phi \mid D^n, X^n) = \mathbb{E}\left\{ (A - \widehat{\pi}_a)(Y - \widehat{\eta}) - (A - \pi)(Y - \eta) \mid D^n, X^n \right\}$$

$$= \mathbb{E}\left\{ (A - \pi + \pi - \widehat{\pi}_a)(Y - \eta + \eta - \widehat{\eta}) - (A - \pi)(Y - \eta) \mid D^n, X^n \right\}$$

$$= \{\widehat{\pi}_a(X_i) - \pi(X_i)\}\{\widehat{\eta}(X_i) - \eta(X_i)\} \equiv \widehat{R}_2(X_i)$$
(15)

by iterated expectation. Let

$$\mathcal{B}_n(x;\widehat{f}) = \mathbb{E}\{\widehat{f}(x) \mid X^n\} - f(x)$$

denote the pointwise conditional bias of a generic estimator  $\hat{f}$  of f at x.

Then the conditional mean of term (10) is

$$\mathbb{E}\Big[b_0^{\mathsf{T}}\mathbb{P}_n(bK_{hx}\nu b^{\mathsf{T}})^{-1}\mathbb{P}_n\{bK_{hx}(\widehat{\phi}-\phi)\} \mid X^n\Big] = \sum_{i=1}^n w_i(x;X^n)\nu(X_i)^{-1}\mathbb{E}(\widehat{\phi}_i-\phi_i\mid X^n)$$

$$= \sum_{i=1}^n w_i(x;X^n)\nu(X_i)^{-1}\mathbb{E}\{\widehat{R}_2(X_i)\mid X^n\}$$

$$= \sum_{i=1}^n w_i(x;X^n)\nu(X_i)^{-1}\mathcal{B}_n(X_i;\widehat{\pi}_a)\mathcal{B}_n(X_i;\widehat{\eta})$$

$$\lesssim \sum_{i=1}^n \frac{1/\lambda_n}{nh^d} |\mathcal{B}_n(X_i;\widehat{\pi}_a)\mathcal{B}_n(X_i;\widehat{\eta})| \mathbb{1}(||X_i-x|| \leq h)$$

$$\lesssim \sup_{||x'-x|| \leq h} |\mathcal{B}_n(x';\widehat{\pi}_a)| |\mathcal{B}_n(x';\widehat{\eta})| \left(\frac{v_n}{\lambda_n}\right) \lesssim k^{-\alpha/d}k^{-\beta/d}\left(\frac{v_n}{\lambda_n}\right)$$

where the second line follows by iterated expectation, the third since  $\widehat{\pi}_a \perp \!\!\! \perp \widehat{\eta}$ , the fourth and fifth by Properties 1–3 of Lemma 1 and since  $\nu \geq \epsilon(1-\epsilon)$ , and the last by Condition 1c via Assumption 4. For the conditional variance we have the decomposition

$$\operatorname{var}\left[b_{0}^{\mathsf{T}}\mathbb{P}_{n}(bK_{hx}\nu b^{\mathsf{T}})^{-1}\mathbb{P}_{n}\left\{bK_{hx}(\widehat{\phi}-\phi)\right\} \mid X^{n}\right] = \operatorname{var}\left[\sum_{i=1}^{n}w_{i}(x;X^{n})\left\{\frac{\widehat{\phi}(Z_{i})-\phi(Z_{i})}{\nu(X_{i})}\right\} \mid X^{n}\right]$$

$$= \mathbb{E}\left(\operatorname{var}\left[\sum_{i=1}^{n}w_{i}(x;X^{n})\left\{\frac{\widehat{\phi}(Z_{i})-\phi(Z_{i})}{\nu(X_{i})}\right\} \mid D^{n},X^{n}\right] \mid X^{n}\right)$$

$$+\operatorname{var}\left(\mathbb{E}\left[\sum_{i=1}^{n}w_{i}(x;X^{n})\left\{\frac{\widehat{\phi}(Z_{i})-\phi(Z_{i})}{\nu(X_{i})}\right\} \mid D^{n},X^{n}\right] \mid X^{n}\right)$$

$$(16)$$

For the term in (16) note that

$$\operatorname{var}\left[\sum_{i=1}^{n} \left\{\frac{w_{i}(x;X^{n})}{\nu(X_{i})}\right\} \left(\widehat{\phi}_{i} - \phi_{i}\right) \mid D^{n}, X^{n}\right] = \sum_{i=1}^{n} \left\{\frac{w_{i}(x;X^{n})}{\nu(X_{i})}\right\}^{2} \operatorname{var}\left(\widehat{\phi}_{i} - \phi_{i} \mid D^{n}, X^{n}\right)$$

$$\lesssim \frac{1}{\epsilon(1-\epsilon)} \sup_{i,x} |w_{i}(x;X^{n})| \sum_{i=1}^{n} |w_{i}(x;X^{n})| \lesssim \frac{1}{nh^{d}} \left(\frac{v_{n}}{\lambda_{n}}\right)$$

where the second line used that  $var(\widehat{\phi} - \phi \mid D^n, X^n)$  and  $1/\nu$  are bounded, and the last Properties 1–2 of Lemma 1. The second term (17) equals

$$\operatorname{var}\left[\sum_{i=1}^{n} \left\{\frac{w_{i}(x; X^{n})}{\nu(X_{i})}\right\} \widehat{R}_{2}(X_{i}) \mid X^{n}\right] \\
= \operatorname{var}\left[\sum_{i=1}^{n} \frac{w_{i}(x; X^{n})}{\nu(X_{i})} \sum_{j=1}^{n} w_{j\alpha}(X_{i}; X^{n}) \left\{A_{j} - \pi(X_{i})\right\} \sum_{j'=1}^{n} w_{j'\beta}(X_{i}; X^{n}) \left\{Y_{j'} - \eta(X_{i})\right\} \mid X^{n}\right] \\
= \operatorname{var}\left[\sum_{i,j,j'} \frac{w_{i}(x; X^{n})}{\nu(X_{i})} w_{j\alpha}(X_{i}; X^{n}) w_{j'\beta}(X_{i}; X^{n}) \left\{A_{j} - \pi(X_{i})\right\} \left\{Y_{j'} - \eta(X_{i})\right\} \mid X^{n}\right] \\
\lesssim \sum_{i=1}^{n} w_{i}(x; X^{n})^{2} \sum_{j=1}^{n} w_{j\alpha}(X_{i}; X^{n})^{2} \sum_{j'=1}^{n} w_{j'\beta}(X_{i}; X^{n})^{2} \lesssim \frac{1/\lambda_{n}}{nh^{d}} \left(\frac{k}{n}\right)^{2}$$

where the first equality used the definition  $\widehat{R}_2$  in (15) and Condition 1a (via Assumption 4), the third that  $\operatorname{var}(V_1V_2) \leq \mathbb{E}(V_1^2)\mathbb{E}(V_2^2)$  if  $V_1 \perp \!\!\! \perp V_2$  and boundedness of (A,Y) and the propensity scores, and the last Properties 1–2 of Lemma 1 and Condition 1b via Assumption 4. Therefore term (10) satisfies

$$\mathbb{E}\left(\left[b_0^{\mathrm{T}}\widetilde{Q}^{-1}\mathbb{P}_n\{bK_{hx}(\widehat{\phi}-\phi)\}\right]^2 \mid X^n\right) \lesssim \left(k^{-4s/d} + \frac{1}{nh^d}\left(1 + \left(\frac{k}{n}\right)^2\right)\right) \left(\frac{v_n}{\lambda_n}\right)^2. \tag{18}$$

## 7.3.3 Term (11)

Note term (11) equals

$$b_0^{\mathrm{T}} \Big\{ \mathbb{P}_n (bK_{hx} \widehat{\varphi}_a b^{\mathrm{T}})^{-1} - \mathbb{P}_n (bK_{hx} \nu b^{\mathrm{T}})^{-1} \Big\} \mathbb{P}_n (bK_{hx} \widehat{\phi}) = b_0^{\mathrm{T}} \left( \widehat{Q}^{-1} - \widetilde{Q}^{-1} \right) \mathbb{P}_n (bK_{hx} \widehat{\phi})$$

$$= \left\{ b_0^{\mathrm{T}} \widetilde{Q}^{-1} \left( \widetilde{Q} - \widehat{Q} \right) \right\} \left\{ \widehat{Q}^{-1} \mathbb{P}_n (bK_{hx} \widehat{\phi}) \right\}$$

$$(19)$$

For the first term in the product in (19) we have

$$b_0^{\mathrm{T}} \widetilde{Q}^{-1} \left( \widehat{Q} - \widetilde{Q} \right) = \sum_{i=1}^n \left\{ \frac{w_i(x; X^n)}{\nu(X_i)} \right\} \left\{ \widehat{\varphi}_a(Z_i) - \nu(X_i) \right\} b(X_i - x)^{\mathrm{T}}$$
$$= \sum_{i=1}^n \left\{ \frac{b(X_i - x)}{\nu(X_i)} \right\} w_i(x; X^n) \left\{ \widehat{\varphi}_a(Z_i) - \nu(X_i) \right\}$$

which we can tackle with similar logic as for term (10). First note

$$\mathbb{E}(\widehat{\varphi}_{a} - \nu \mid D^{n}, X^{n}) = \mathbb{E}\left\{ (A - \widehat{\pi}_{a})(A - \widehat{\pi}_{b}) - (A - \pi)^{2} \mid D^{n}, X^{n} \right\}$$

$$= \mathbb{E}\left\{ (A - \pi + \pi - \widehat{\pi}_{a})(A - \pi + \pi - \widehat{\pi}_{b}) - (A - \pi)^{2} \mid D^{n}, X^{n} \right\}$$

$$= \{\widehat{\pi}_{a}(X_{i}) - \pi(X_{i})\}\{\widehat{\pi}_{b}(X_{i}) - \pi(X_{i})\} \equiv \widehat{R}_{2\pi}(X_{i})$$
(20)

by iterated expectation. Defining  $\mathcal{B}_n(x; \hat{f})$  as in the previous subsection, the conditional mean of the first term of (19) is

$$\mathbb{E}\left\{b_0^{\mathsf{T}}\widetilde{Q}^{-1}\left(\widehat{Q}-\widetilde{Q}\right) \mid X^n\right\} = \sum_{i=1}^n \left\{\frac{b(X_i-x)}{\nu(X_i)}\right\} w_i(x;X^n) \mathbb{E}\left\{\widehat{\varphi}_a(Z_i) - \nu(X_i) \mid X^n\right\}$$

$$= \sum_{i=1}^n \left\{\frac{b(X_i-x)}{\nu(X_i)}\right\} w_i(x;X^n) \mathbb{E}\left\{\widehat{R}_{2\pi}(X_i) \mid X^n\right\}$$

$$= \sum_{i=1}^n \left\{\frac{b(X_i-x)}{\nu(X_i)}\right\} w_i(x;X^n) \ \mathcal{B}_n(X_i;\widehat{\pi}_a) \mathcal{B}_n(X_i;\widehat{\pi}_b)$$

$$\lesssim \sum_{i=1}^n \frac{1/\lambda_n}{nh^d} \left|\mathcal{B}_n(X_i;\widehat{\pi}_a)\mathcal{B}_n(X_i;\widehat{\pi}_b)\right| \mathbb{1}(\|X_i-x\| \leq h)$$

$$\lesssim \sup_{\|x'-x\| \leq h} \left|\mathcal{B}_n(x';\widehat{\pi}_a)\right| \left|\mathcal{B}_n(x';\widehat{\pi}_b)\right| \left(\frac{v_n}{\lambda_n}\right) \lesssim k^{-2\alpha/d} \left(\frac{v_n}{\lambda_n}\right)$$

where the second line follows by iterated expectation, the third since  $\widehat{\pi}_a \perp \!\!\! \perp \widehat{\pi}_b$ , the fourth and fifth by Properties 1–3 of Lemma 1 and since  $\nu \geq \epsilon(1-\epsilon)$ , and the last by Condition 1c via Assumption 4.

The analysis of the conditional variance follows exactly the same logic as for term (10), and is of the same order. For the second term in the product in (19) we have

$$\left| \left| \widehat{Q}^{-1} \mathbb{P}_n(bK_{hx}\widehat{\phi}) \right| \right| \lesssim \lambda_n^{-1} \left| \left| \mathbb{P}_n(bK_{hx}\widehat{\phi}) \right| \right|$$

and note

$$\mathbb{E}\{\|\mathbb{P}_n(bK_{hx}\widehat{\phi})\|^2 \mid X^n\} = \sum_j \mathbb{E}\{\mathbb{P}_n(b_jK_{hx}\widehat{\phi})^2 \mid X^n\}$$
$$= \sum_j \left[\mathbb{E}\{\mathbb{P}_n(b_jK_{hx}\widehat{\phi}) \mid X^n\}^2 + \operatorname{var}\{\mathbb{P}_n(b_jK_{hx}\widehat{\phi}) \mid X^n\}\right].$$

Therefore for  $b_{\ell}$  the  $\ell$ -th component of  $b(X_i - x)$ 

$$\mathbb{E}\Big\{\mathbb{P}_n(b_{\ell}K_{hx}\widehat{\phi})\mid X^n\Big\} = \mathbb{P}_n\Big\{b_{\ell}K_{hx}\mathbb{E}(\widehat{\phi}_i\mid X^n)\Big\} \lesssim \frac{1}{nh^d}\sum_{i=1}^n \mathbb{1}(\|X_i - x\| \le h) \lesssim v_n$$

using boundedness of the kernel K and observations Z. For the variance we have

$$\operatorname{var}\left\{\mathbb{P}_{n}(b_{\ell}K_{hx}\widehat{\phi})\mid X^{n}\right\} = \mathbb{E}\left[\operatorname{var}\left\{\mathbb{P}_{n}(b_{\ell}K_{hx}\widehat{\phi})\mid D^{n}, X^{n}\right\}\mid X^{n}\right] + \operatorname{var}\left[\mathbb{E}\left\{\mathbb{P}_{n}(b_{\ell}K_{hx}\widehat{\phi})\mid D^{n}, X^{n}\right\}\mid X^{n}\right]$$
(21)

For the term in (21) note

$$\operatorname{var}\left\{\mathbb{P}_{n}(b_{\ell}K_{hx}\widehat{\phi})\mid D^{n}, X^{n}\right\} = \left(\frac{1}{nh^{d}}\right)^{2} \sum_{i=1}^{n} b_{\ell}(X_{i} - x)^{2} K\left(\frac{X_{i} - x}{h}\right)^{2} \operatorname{var}(\widehat{\phi}_{i} \mid D^{n}, X^{n})$$

$$\lesssim \left(\frac{1}{nh^{d}}\right)^{2} \sum_{i=1}^{n} \mathbb{1}(\|X_{i} - x\| \leq h) \lesssim \frac{v_{n}}{nh^{d}}$$

using the boundedness of the kernel, X, and  $\widehat{\phi}$ . For the term in (22)

$$\operatorname{var}\left[\mathbb{E}\left\{\mathbb{P}_{n}(b_{\ell}K_{hx}\widehat{\phi})\mid D^{n}, X^{n}\right\} \mid X^{n}\right] = \operatorname{var}\left[\mathbb{P}_{n}\left\{b_{\ell}K_{hx}\mathbb{E}(\widehat{\phi}\mid D^{n}, X^{n})\right\} \mid X^{n}\right]$$

$$= \operatorname{var}\left\{\mathbb{P}_{n}\left\{b_{\ell}K_{hx}\widehat{R}_{2}(X_{i})\mid X^{n}\right\}\right\}$$

$$= \operatorname{var}\left[\frac{1}{n}\sum_{i=1}^{n}b_{\ell}K_{hx}\sum_{j=1}^{n}w_{j\alpha}(X_{i}; X^{n})\{A_{j} - \pi(X_{i})\}\sum_{j'=1}^{n}w_{j'\beta}(X_{i}; X^{n})\{Y_{j'} - \eta(X_{i})\}\mid X^{n}\right]$$

$$= \operatorname{var}\left[\frac{1}{n}\sum_{i,j,j'}b_{\ell}K_{hx}w_{j\alpha}(X_{i}; X^{n})w_{j'\beta}(X_{i}; X^{n})\{A_{j} - \pi(X_{i})\}\{Y_{j'} - \eta(X_{i})\}\mid X^{n}\right]$$

$$\lesssim \frac{1}{n^{2}}\sum_{i=1}^{n}K_{hx}(X_{i})^{2}\sum_{j=1}^{n}w_{j\alpha}(X_{i}; X^{n})^{2}\sum_{j'=1}^{n}w_{j'\beta}(X_{i}; X^{n})^{2}\lesssim \frac{v_{n}}{nh^{d}}\left(\frac{k}{n}\right)^{2}$$

where the second line follows from the definition of  $\widehat{R}_2(X_i)$  and since  $\mathbb{E}(\phi \mid X_i)$  is constant given  $X^n$ , and the rest follows the same logic as for the term in (17). Therefore given  $X^n$  we have that

$$\widehat{Q}^{-1}\mathbb{P}_n(bK_{hx}\widehat{\phi}) = O_{\mathbb{P}}\left(\left(\frac{v_n}{\lambda_n}\right)\left(1 + \frac{1}{\sqrt{v_n nh^d}}\left(1 + \frac{k}{n}\right)\right)\right)$$

and so conditional on  $X^n$ , the term (11) is of order

$$O_{\mathbb{P}}\left(\left(k^{-2\alpha/d} + \frac{1}{\sqrt{nh^d}}\left(1 + \frac{k}{n}\right)\right)\left(\frac{v_n}{\lambda_n}\right)^2 \left(1 + \frac{1}{\sqrt{v_n nh^d}}\left(1 + \frac{k}{n}\right)\right)\right). \tag{23}$$

### 7.3.4 Combining Bounds

Now we use Lemma 6.1 of Chernozhukov et al. [2018] or equivalently the following lemma, to deduce unconditional convergence from the previously derived conditional bounds.

**Lemma 2.** Suppose a random variable  $Z_n$  satisfies

$$|\mathbb{E}(Z_n \mid X^n)| \lesssim b_n$$
 and  $var(Z_n \mid X^n) \lesssim s_n^2$ .

Then  $Z_n = O_{\mathbb{P}}(b_n + s_n)$ 

*Proof.* We have

$$|\mathbb{E}(Z_n)| = |\mathbb{E}\{\mathbb{E}(Z_n \mid X^n)\}| \le \mathbb{E}|\mathbb{E}(Z_n \mid X^n)| \lesssim b_n$$

and

$$\operatorname{var}(Z_n) = \mathbb{E}\{\operatorname{var}(Z_n \mid X^n)\} + \operatorname{var}\{\mathbb{E}(Z_n \mid X^n)\}$$

$$\lesssim s_n^2 + \mathbb{E}[\{\mathbb{E}(Z_n \mid X^n) - \mathbb{E}(Z_n)\}^2]$$

$$\leq s_n^2 + 2\mathbb{E}\{\mathbb{E}(Z_n \mid X^n)^2 + \mathbb{E}(Z_n)^2\} \lesssim s_n^2 + b_n^2$$

where the first equality follows by the law of total variance, the second by the variance bound and by definition (for the first and second terms, respectively), the third since  $(a + b)^2 \le 2(a^2 + b^2)$ , and the fourth by the mean bound. Therefore

$$\mathbb{P}\left(\frac{|Z_n|}{b_n + s_n} \ge M_{\epsilon}\right) \le \frac{\mathbb{E}(Z_n^2)/(b_n + s_n)^2}{M_{\epsilon}^2} \le \frac{C(b_n^2 + s_n^2)}{(b_n + s_n)^2} \frac{1}{M_{\epsilon}^2}$$

where the first equality follows from Chebyshev's inequality, and the second since  $(b_n + s_n)^2 \ge b_n^2 + s_n^2$ . Therefore the result follows since one can always take  $M_{\epsilon} = 1/\sqrt{\epsilon/C}$ .

Therefore combining the conditional bounds in (14), (18), and (23), together with the facts that  $v_n = O_{\mathbb{P}}(1)$  and  $1/\lambda_n = O_{\mathbb{P}}(1)$  from Assumptions 1 and 3, respectively, we have the unconditional convergence results

$$b_0^{\mathrm{T}} \widehat{Q}^{-1} \mathbb{P}_n(bK_{hx}\phi) - \tau = O_{\mathbb{P}} \left( h^{\gamma} + \frac{1}{\sqrt{nh^d}} \right)$$
 (24)

$$b_0^{\mathrm{T}} \mathbb{P}_n (bK_{hx} \nu b^{\mathrm{T}})^{-1} \mathbb{P}_n \{ bK_{hx} (\widehat{\phi} - \phi) \} = O_{\mathbb{P}} \left( k^{-2s/d} + \frac{1}{\sqrt{nh^d}} \left( 1 + \frac{k}{n} \right) \right)$$
 (25)

$$b_0^{\mathrm{T}} \widetilde{Q}^{-1} \left( \widetilde{Q} - \widehat{Q} \right) \widehat{Q}^{-1} \mathbb{P}_n(bK_{hx} \widehat{\phi}) = O_{\mathbb{P}} \left( k^{-2\alpha/d} + \frac{1}{\sqrt{nh^d}} \left( 1 + \frac{k}{n} \right) \right)$$
 (26)

where for the last result we also used the condition that  $(k/n)/\sqrt{nh^d} \to 0$  to discard lower order terms.

## 8 R Code

Piecewise polynomial model from motivating example in Section 2.2:

```
set.seed(1234)
expit <- function(x){ \exp(x)/(1+\exp(x)) }; logit <- function(x){ \log(x/(1-x)) }
n < 4*2000; nsim < 500; rateseq < seq(0.1,0.5,by=0.05); res2 < NULL
for (rate in rateseq){
  res <- data.frame(matrix(nrow=nsim,ncol=4))</pre>
  colnames(res) <- c("plugin", "xl", "drl", "oracle.drl")</pre>
for (i in 1:nsim){
## simulate data
s \leftarrow sort(rep(1:4,n/4)); x \leftarrow (runif(n,-1,1)); ps \leftarrow 0.1 + 0.8*(x>0)
mu0 < (x <= -.5)*0.5*(x+2)^2 + (x/2+0.875)*(x>-1/2 & x<0) +
  (x>0 & x<.5)*(-5*(x-0.2)^2 +1.075) + (x>.5)*(x+0.125); mu1 <- mu0; tau <- 0
a \leftarrow rbinom(n,1,ps); y \leftarrow a*mu1 + (1-a)*mu0 + rnorm(n,sd=(.2-.1*cos(2*pi*x)))
## estimate nuisance functions
pihat <- expit( logit(ps) + rnorm(n,mean=1/(n/4)^rate,sd=1/(n/4)^rate))
mu1hat <- predict(smooth.spline(x[a==1 \& s==2],y[a==1 \& s==2]),x)$y
mu0hat <- predict(smooth.spline(x[a==0 \& s==2],y[a==0 \& s==2]),x)$y
## construct estimators
plugin <- mu1hat-mu0hat</pre>
x1 \leftarrow predict(smooth.spline(x[a==1 \& s==3],(y-mu0hat)[a==1 \& s==3]),x)$y
x0 \leftarrow predict(smooth.spline(x[a==0 \& s==3],(mu1hat-y)[a==0 \& s==3]),x)$y
xl \leftarrow pihat*x0 + (1-pihat)*x1
pseudo <- ((a-pihat)/(pihat*(1-pihat)))*(y-a*mu1hat-(1-a)*mu0hat) + mu1hat-mu0hat
drl <- predict(smooth.spline(x[s==3],pseudo[s==3]),x)$y</pre>
pseudo.or <- ((a-ps)/(ps*(1-ps)))*(y-a*mu1-(1-a)*mu0) + mu1-mu0
oracle.drl <- predict(smooth.spline(x[s==3],pseudo.or[s==3]),x)$y</pre>
## save MSEs
res$plugin[i] <- (n/4)*mean((plugin-tau)[s==4]^2)
res$xl[i] <- (n/4)*mean((xl-tau)[s==4]^2)
res$drl[i] <- (n/4)*mean((drl-tau)[s==4]^2)
ressoracle.drl[i] <- (n/4)*mean((oracle.drl-tau)[s==4]^2)
}
res2 <- rbind(res2, c(rate, apply(res,2,mean)))
}
```

High-dimensional model example:

```
set.seed(1234)
library(glmnet)
expit <- function(x){ exp(x)/(1+exp(x)) }; logit <- function(x){ log(x/(1-x)) }
nsim <- 100; res <- data.frame(matrix(nrow=nsim,ncol=4))</pre>
colnames(res) <- c("plugin","xl","drl","oracle.drl")</pre>
n < 4*2000; d < 500; alpha d/10; beta d/10; beta
for (i in 1:nsim){
## simulate data
s \leftarrow sort(rep(1:4,n/4)); x \leftarrow matrix(rnorm(n*d), n, d)
mu0 <- expit(as.numeric(x %*% rep(c(1, 0), c(beta,d-beta)))/sqrt(beta/1))</pre>
ps <- expit(as.numeric(x %*% rep(c(1, 0), c(alpha,d-alpha)))/sqrt(alpha/0.25))</pre>
tau <- 0; mu1 <- mu0 + tau
a \leftarrow rbinom(n,1,ps); y \leftarrow rbinom(n,1,a*mu1+(1-a)*mu0)
## estimate nuisance functions
pihat <- predict(cv.glmnet(x[s==1,],a[s==1], family="binomial"),newx=x,</pre>
         type="response", s="lambda.min")
mu0hat <- predict(cv.glmnet(x[a==0 & s==2,],y[a==0 & s==2], family="binomial"), newx=x,
          type="response", s="lambda.min")
mu1hat <- predict(cv.glmnet(x[a==1 & s==2,],y[a==1 & s==2], family="binomial"),newx=x,
          type="response", s="lambda.min")
## construct estimators
plugin <- mu1hat-mu0hat</pre>
x1 <- predict(cv.glmnet(x[a==1 & s==3,],(y-mu0hat)[a==1 & s==3]),newx=x)
x0 \leftarrow predict(cv.glmnet(x[a==0 \& s==3,],(mu1hat-y)[a==0 \& s==3]),newx=x)
xl \leftarrow pihat*x0 + (1-pihat)*x1
pseudo <- ((a-pihat)/(pihat*(1-pihat)))*(y-a*mu1hat-(1-a)*mu0hat) + mu1hat-mu0hat
drl <- predict(cv.glmnet(x[s==3,],pseudo[s==3]),newx=x)</pre>
pseudo.or <- ((a-ps)/(ps*(1-ps)))*(y-a*mu1-(1-a)*mu0) + mu1-mu0
oracle.drl <- predict(cv.glmnet(x[s==3,],pseudo.or[s==3]),newx=x)</pre>
## save MSEs
res$plugin[i] <- (n/4)*mean((plugin-tau)[s==4]^2)
res$xl[i] <- (n/4)*mean((xl-tau)[s==4]^2)
res$drl[i] <- (n/4)*mean((drl-tau)[s==4]^2)
res$oracle.drl[i] <- (n/4)*mean((oracle.drl-tau)[s==4]^2)
}
```