

# How Much Weak Overlap Can Doubly Robust T-Statistics Handle?

Jacob Dorn\*

April 25, 2025

## Abstract

In the presence of sufficiently weak overlap, it is known that no regular root-n-consistent estimators exist and standard estimators may fail to be asymptotically normal. This paper shows that a thresholded version of the standard doubly robust estimator is asymptotically normal with well-calibrated Wald confidence intervals even when constructed using nonparametric estimates of the propensity score and conditional mean outcome. The analysis implies a cost of weak overlap in terms of black-box nuisance rates, borne when the semiparametric bound is infinite, and the contribution of outcome smoothness to the outcome regression rate, which is incurred even when the semiparametric bound is finite. As a byproduct of this analysis, I show that under weak overlap, the optimal global regression rate is the same as the optimal pointwise regression rate, without the usual polylogarithmic penalty. The high-level conditions yield new rules of thumb for thresholding in practice. In simulations, thresholded AIPW can exhibit moderate overrejection in small samples, but I am unable to reject a null hypothesis of exact coverage in large samples. In an empirical application, the clipped AIPW estimator that targets the standard average treatment effect yields similar precision to a heuristic 10% fixed-trimming approach that changes the target sample.

## 1 Introduction

In observational data, it is common to find at least some treated observations with small propensity scores, in which case overlap is weak. Under weak overlap, outcome regression is more difficult and inverse propensity estimators can divide by small numbers.

This paper studies the extent to which standard asymptotic results for doubly robust estimators continue to apply when strict overlap does not hold, but there is a tail bound on inverse propensity scores. The paper

---

\*I am grateful for suggestions from Xiaohong Chen, Rebecca Dorn, Kevin Guo, Edward Kennedy, Samir Khan, Michal Kolesár, Lihua Lei, Xinwei Ma, Ulrich Müller, Yuya Sasaki, Yulong Wang, and Larry Wasserman, and participants at seminars at the University of Pennsylvania. Artificial intelligence was used to suggest changes in the editing process. This material is based upon work supported by the National Science Foundation Graduate Research Fellowship Program under Grant No. DGE-2039656 and by grant T32 HS026116 from the Agency for Healthcare Research and Quality. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author and do not necessarily reflect the views of the National Science Foundation or the Agency for Healthcare Research and Quality.

aims to answer two questions:

- (1) *Black-box rates.* Under what conditions on nuisance errors do Wald confidence intervals constructed using doubly robust estimators have well-calibrated coverage?
- (2) *Feasibility.* When are the black-box rates in (1) feasible under weak overlap?

Overlap weakness is measured with a tail-bound parameter  $\gamma_0$ . It is known that there is a phase transition when  $\gamma_0$  crosses two, which corresponds to a uniform propensity density.

Values of the tail parameter  $\gamma_0$  above two correspond to stronger overlap assurances, a case which I call *somewhat weak overlap*. Under somewhat weak overlap, it is known that the semiparametric bound for estimation of  $\psi_0$  is finite and can be achieved by the doubly robust Augmented Inverse Propensity Weighted (AIPW) estimator with known conditional mean outcome and propensity functions. I show that semiparametric efficiency holds even when using nonparametric estimates of the two AIPW nuisance functions, provided the nuisance estimates are constructed by cross-fitting and achieve error rates  $r_{\mu,n}$  and  $r_{e,n}$  that satisfy the product condition that  $n^{1/2}r_{\mu,n}r_{e,n}$  tends to zero in probability. The one modification I make relative to the usual analysis is that I require the nuisance rates to apply uniformly in regions with small propensity scores, which is immediate under  $L_\infty$  bounds.

Values of  $\gamma_0$  below two allow the density of propensity scores to be unbounded at zero, a case which I call *very weak overlap*. Under very weak overlap, it is known that the semiparametric efficiency bound is infinite and Inverse Propensity Weighting (IPW) estimators fail to be asymptotically normal, even if the propensity function is known. [Ma and Wang \(2020\)](#) show that thresholding strategies that clip (Winsorize) or trim (discard) observations with small propensity scores at a sequence of  $b_n$  tending to zero can restore asymptotic normality, but at a slower-than- $\sqrt{n}$  rate and with first-order bias. Previous work has proposed constructing confidence intervals for this case using debiasing strategies and self-normalized subsampling. I show that if thresholding is applied to AIPW, then there is no first-order bias and simple Wald confidence have exact asymptotic coverage. The analysis provides sufficient black-box conditions for AIPW with estimated nuisance functions to achieve these asymptotic results: first, that the threshold  $b_n$  tends to zero more slowly than the propensity error  $r_{e,n}$ ; and second, that the product  $n^{1/2}r_{\mu,n}b_n^{\gamma_0/2}$  tends to zero in probability.

That analysis characterizes the feasibility of AIPW Wald confidence intervals under black-box conditions on the propensity estimates and outcome regression. An added difficulty is that in regions with weak overlap, there cannot be many treated observations to use in outcome regression. I quantify the added difficulty for the case of regression within a Hölder smoothness class of order  $\beta_\mu > 0$ . I show that weak overlap scales the effective outcome smoothness by  $1 - 1/\gamma_0$ , exhibiting a cost of weak overlap even in the somewhat weak

overlap regime in which the other asymptotic results carry through nearly unchanged. I show that in this setting, the optimal global rate is equal to the optimal pointwise rate, without the usual polylogarithmic factor separating the optimal pointwise and global rates (Stone, 1982). The argument leverages a novel construction partitioning the covariate space into regions with increasingly strong overlap assurances in such a way as the contribution of each region to the worst-case error is half as large as the previous contribution. In this way, the contribution of infinitely many regions to the global error is controlled uniformly.

Taken together, these results provide a precise answer to the question posed by this work's title: doubly-robust t-statistics can handle weak overlap of tail bound  $\gamma_0$ , provided the outcome and propensity nuisance functions are in Hölder smoothness classes of order  $\beta_\mu$  and  $\beta_e$  and

$$\frac{\beta_\mu(1 - 1/\gamma_0)}{2\beta_\mu(1 - 1/\gamma_0) + d} + \frac{\beta_e \min\{\gamma_0/2, 1\}}{2\beta_e + d} > \frac{1}{2}.$$

In this case, the threshold  $b_n = n^{-\beta_e/(2\beta_e+d)} \log(n)^{(3\beta_e+d)/(2\beta_e+d)}$  suffices, regardless of the weak overlap parameter  $\gamma_0$ . When the outcome and propensity smoothness orders are the same  $\beta > 0$ , then thresholded AIPW can handle weak overlap of order  $\gamma_0$ , so long as:

$$\gamma_0 > \max \left\{ \frac{2\beta^2 + 2\beta d + d^2}{\beta(2\beta + d)}, \frac{4\beta^2}{4\beta^2 - d^2} \right\}.$$

Under Lipschitz continuity of both nuisance functions in one dimension, doubly robust t-statistics can handle weak overlap of order  $\gamma_0 > \frac{5}{3}$ . In higher dimensions, there is always some sufficient smoothness order that yields valid t-statistics for any fixed tail bound.

The conditions here yield new rules of thumb for thresholding in applied work. In my favored regime, the econometrician is willing to posit a minimal consistency rate for one of the two nuisance function estimates. Given such a minimal rate, a simple plug-in procedure predicts the threshold with the laxest restriction on the other nuisance function needed to achieve well-calibrated Wald confidence intervals. In the absence of any such information, a third rule of thumb derives a threshold that imposes the laxest equal minimal consistency rate on both nuisance estimates. None of these rules of thumb depend directly on knowledge of the tail bound parameter  $\gamma_0$ .

In simulations, I find that clipped AIPW achieves the promised properties asymptotically. I consider a setting of very weak overlap with nonparametric outcome regression and propensity estimates. Unthresholded IPW and AIPW estimators perform poorly, with large errors and nonnormal asymptotic distributions. In this setting, clipped IPW displays its known first-order bias, and clipped AIPW displays the second-order bias justified by the theoretical analysis. With access to 1,000 or 10,000 observations, I find that p-values based

on clipped AIPW t-statistics exhibit moderate overrejection. In large samples with 100,000 observations, a Kolmogorov-Smirnov test based on 5,000 simulations is unable to reject a null hypothesis that clipped AIPW p-values on the true causal effect are exactly uniformly distributed.

I apply the clipped AIPW estimator to data on right heart catheterization. I consider the setting of [Connors et al. \(1996\)](#), which has become a canonical setting with weak overlap, including providing the empirical application for [Crump et al. \(2009\)](#)’s paper proposing a 10% fixed-trimming rule of thumb. I compare the clipped AIPW estimator that targets the full-population effect to estimators that apply AIPW to a sample trimmed based on a fixed rule. I find that by including observations with small estimated propensities, the clipped AIPW strategy increases the estimated harm of the procedure by 0.17 standard errors relative to the 10 percent fixed-trimming rule, while increasing the estimated standard error by only 5.1%. These results show that targeting the full-population treatment effect does not need to introduce a major efficiency loss, and show that thresholded AIPW can easily be added as a robustness test when practitioners apply a fixed-trimming rule.

Weak overlap is a common phenomenon in practice and in theory. The dominant response to weak overlap in inverse propensity score practice is trimming: dropping samples with small propensity estimates in order to estimate average effects within a more precise population ([Currie and Walker, 2011](#); [Bailey and Goodman-Bacon, 2015](#); [Galiani et al., 2005](#)), typically following the 10% rule of thumb from [Crump et al. \(2009\)](#). Other proposals targeting new samples include reweighting towards higher-precision populations ([Yang and Ding, 2018](#); [Li et al., 2018](#)) or clipping strategies that Winsorize weights above ([Lee et al., 2011](#); [Ionides, 2008](#)).<sup>1</sup> [D’Amour et al. \(2021\)](#) argue that weak overlap is likely to be prevalent in modern settings with high-dimensional covariates. [Imbens \(2004\)](#) argues that changing the target estimand may be necessary in the absence of sufficient precision.

There is theoretical work characterizing estimation of nonstandard estimands under weak overlap. [Khan and Tamer \(2010\)](#) show that very weak overlap yields an irregularly identified parameter, an infinite semi-parametric efficiency bound, slower-than- $\sqrt{n}$  estimation rates for the traditional average causal effects, and no clear notion of best estimator. An important theoretical literature has proposed novel point and confidence interval estimators with desirable properties under weak overlap ([Rothe, 2017](#); [Armstrong and Kolesár, 2017, 2021](#); [Sasaki and Ura, 2022](#); [Ma et al., 2023](#); [Chaudhuri and Hill, 2024](#)). Many of these procedures have favorable properties relative to the simpler thresholded estimator I consider here, but to my knowledge there has been little take-up by practitioners. [Ma and Wang \(2020\)](#) and [Khan and Ugander \(2022\)](#) show that sufficiently trimmed AIPW and IPW can remain asymptotically normal, but at the cost of introducing

---

<sup>1</sup>Awkwardly, the epidemiological literature sometimes refers to the Winsorization strategy as “trimming.” My results hold for both dropping or Winsorizing extreme propensities, so the confused reader can view this as a work deriving simple asymptotics for trimmed AIPW regardless of their preferred meaning of “trim.”

first-order bias for the standard causal effects that often calls for a nonstandard debiasing strategy that can enable laxer regression conditions than the standard smoothness conditions I explore. [Crump et al. \(2009\)](#), [Yang and Ding \(2018\)](#), [Li et al. \(2018\)](#), and [Goldsmith-Pinkham et al. \(2024\)](#) propose changing estimands in response to weak overlap, which introduces a discontinuous estimand based on whether or not the econometrician detects meaningful overlap weakness.

Other theoretical literature so far has either proposed a nonstandard causal estimator, required non-standard techniques to construct confidence intervals, or done both. [Ma and Wang \(2020\)](#) and [Heiler and Kazak \(2021\)](#) propose using self-normalized subsampling methods that enable valid statistical inference for standard estimands without clipping or trimming, but empirical practice has favored simple t-tests. [Ma and Wang](#) also propose a debiasing procedure. [Heiler and Kazak](#) also find that estimated untrimmed AIPW is first-order equivalent to the oracle AIPW estimator with known nuisance functions, and the associated asymptotic distribution is alpha-stable, if the product of nuisance estimation rates is of a lower order than the oracle standard deviation; I find this result does not extend to the thresholded AIPW estimator that I show is asymptotically normal. [Ma et al. \(2024\)](#) and [Lei et al. \(2021\)](#) propose statistical tests under a null of sufficient or strict overlap, respectively, presumably in the hopes of avoiding these complications. My analysis of nonparametric regression rates is also relevant to the literature on regression with degenerate designs, although to my knowledge the possibility of a global rate with no polylogarithmic penalty is new ([Stone, 1982](#); [Hall et al., 1997](#); [Gaïffas, 2005](#); [Mou et al., 2023](#)).

The plan of the paper is as follows. Section 2 presents the setting and main theoretical results. Section 3 interprets these results as minimal black-box consistency rates and as minimal smoothness rates. Section 4 considers implications for parametric estimators and derives some rules of thumb for empirical use. Section 5 presents numerical results for simulations and the empirical application to right-heart catheterization. Section 6 concludes.

**Notation.** I follow [Heiler and Kazak \(2021\)](#) and use “strict overlap” to refer to the case in which the propensity score is bounded away from zero almost surely; I use “weak overlap” to refer to the case in which the infimum of the support of the propensity score is zero, which is sometimes called “limited overlap” ([Khan and Tamer, 2010](#); [Chaudhuri and Hill, 2024](#)). I focus my attention on distributions with weak overlap that may possess subexponential tails. I use “very weak overlap” to refer to case in which the associated heavy tails can fail to generate inverse propensity moments, a class which is sometimes called “heavy tailed” ([Chaudhuri and Hill, 2024](#)). I use “somewhat weak overlap” to refer to the case in which I allow only subexponential tails that are sufficiently light, a class which is sometimes said to satisfy “strict overlap” or “overlap” ([Heiler and Kazak, 2021](#); [Bruns-Smith et al., 2024](#)). I write  $\tilde{\psi}_{(Oracle)}^{AIPW}(b_n) = \frac{1}{n} \sum \phi(Z \mid b_n, \eta)$  for the oracle AIPW estimate with pseudo-outcome  $\phi(Z \mid b_n, \{\bar{e}(X), \bar{\mu}(X)\}) = \bar{\mu}(X) + \frac{D(Y - \bar{\mu}(X))}{\max\{\bar{e}(X), b\}}$  and clipping threshold  $b_n$  and

$\sigma_n = n^{-1/2} \sqrt{\frac{1}{n} \sum \phi(Z | b_n, \eta)^2 - \tilde{\psi}_{(Oracle)}^{AIPW}(b_n)^2}$  and  $\hat{\sigma}_n = n^{-1/2} \sqrt{\frac{1}{n} \sum \phi(Z | b_n, \hat{\eta})^2 - (\frac{1}{n} \sum \phi(Z | b_n, \hat{\eta}))^2}$   
 for the associated oracle and estimated sample standard deviation, respectively. I refer to regions of the  
 covariate space in which the propensity can be arbitrarily close to zero as singularities. I use the notation  
 $E_P[\cdot]$  and  $E[\cdot]$  to refer to the expectation under the maintained distribution  $P$ , and I use the notation  $\psi(P)$   
 to refer to the statistical average potential outcome  $E_P[E_P[Y | X, D = 1]]$  where the right-hand side is  
 well-defined under  $P$ . I abuse notation and write  $\psi = \psi(P)$  and use  $\sup_{P \in A} B$  to refer to the supremum of  
 $B$  over distributions  $P$  in  $A$  under any maintained restrictions on the distribution and nuisance functions.  
 I write that a set of nuisance functions are cross-fit if the data is partitioned into  $K$  folds and the nuisance  
 functions in fold  $k$  are independent of the data in fold  $k$ . I write  $A_n \leq_P B_n$  to refer to the case that for all  
 $\epsilon > 0$ ,  $P(A_n > B_n + \epsilon) \rightarrow 0$ . I write  $P(E_n)$  for the probability of event  $E_n$  occurring under the distribution  
 $P$ , with the number of draws  $n$  sometimes left implicit. I use the notation  $c_n \ll d_n$  for nonnegative sequences  
 $c_n, d_n$  to indicate that  $d_n > 0$  for all  $n$  large enough and  $c_n/d_n \rightarrow 0$ . I use the notation  $c_n \lesssim d_n$  and  $d_n \gtrsim c_n$   
 to indicate that there is some  $\delta > 0$  such that  $d_n \geq \delta c_n$  for all  $n$  large enough. I write  $c_n = o_P(d_n)$  for  
 sequence of  $d_n > 0$  to indicate that for all  $\delta > 0$ ,  $P(|c_n|/d_n > \delta) \rightarrow 0$ ; if there is only one distribution in a  
 statement,  $c_n = o(d_n)$  should be understood to mean  $c_n = o_P(d_n)$ . I use  $\log$  to refer to the natural logarithm  
 and  $a \vee b$  to indicate  $\max\{a, b\}$ . I define Hölder smoothness using a multivariate version of the notation  
 of [Tsybakov \(2009\)](#): a function  $f$  is in the Hölder smoothness class  $\Sigma(\beta, L)$  if the  $\lfloor \beta \rfloor$ -order multivariate  
 derivatives  $D^\alpha f = \frac{\partial \|\alpha\|}{\partial x_1^{\alpha_1} \partial x_2^{\alpha_2} \dots} f$  satisfy  $\|D^\alpha f(x) - D^\alpha f(x')\| \leq L \|x - x'\|^{\beta - \lfloor \beta \rfloor}$ , where I write  $D^\alpha f(x)$  for  
 $D^\alpha f$  evaluated at  $x$ . For simplicity, I use local polynomial regression to refer to specifically kernel regression  
 with uniform bandwidth:  $\hat{\mu}^{(NW)}(x | h) = \frac{\sum D \mathbf{1}\{\|X - x\| \leq h\} Y}{\sum D \mathbf{1}\{\|X - x\| \leq h\}}$  when feasible and  $\hat{\mu}^{(NW)}(x | h) = 0$  when no  
 nearby treated observations are available. For an estimator  $\hat{\eta}$  of  $\eta$ , I use  $\|\hat{\eta} - \eta\|_\infty$  to refer to the sup norm  
 $\sup_{x \in \text{Support}(P)} |\hat{\eta}(x) - \eta(x)|$ .

## 2 Setting, Consistency, and Asymptotic Normality

This section presents asymptotic results under black-box nuisance conditions.

### 2.1 Setting

I derive uniform convergence rates under lower bounds on overlap weakness. I follow [Ma and Wang \(2020\)](#), who provide important building blocks in my analysis, and parameterize overlap weakness through a  
 tail parameter  $\gamma_0$ . Unlike their analysis, the results will be uniform over a model family  $\mathcal{P}$  satisfying certain  
 restrictions, including some basic regularity conditions.

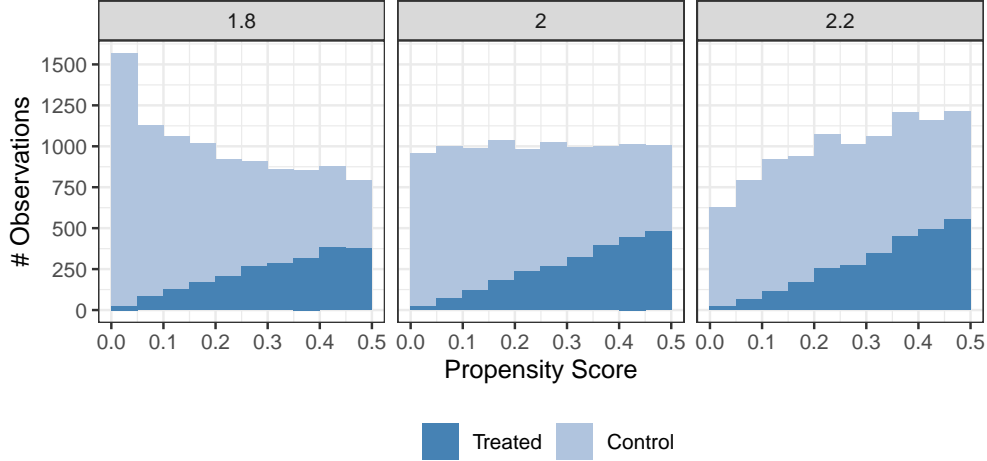


Figure 1: Simulations of 10,000 observations of  $e(X)$  with  $P(e(X) \leq \pi) = \pi^{\gamma_0-1}$  for increasing values of  $\gamma_0$ .

**Assumption 1.** Let  $\mathcal{P}$  be a nonempty family of distributions, and write  $e(X) = P(D = 1 | X)$  and  $\mu(X) = E_P[Y | X, D = 1]$ . Then every  $P \in \mathcal{P}$  satisfies the following conditions for some  $q > 3, M, \sigma_{\min}, C > 0$  and  $\gamma_0 > 1$ :

- (a) *Conditional moments.*  $\mathbb{E}[|Y - \mu(X)|^q | X, D = 1] \leq M^q < \infty$  almost surely.
- (b) *Unconditional moments.*  $\text{Var}(\mu(X)) \leq M$ .
- (c) *Residuals.*  $\text{Var}(Y | X, D) \geq \sigma_{\min}^2$ .
- (d) *Propensity tail.*  $P(e(X) \leq \pi) \leq C\pi^{\gamma_0-1}$  for all  $\pi \in [0, 1]$ .

Definition 1 generalizes [Ma and Wang \(2020\)](#)’s slowly varying tails assumption. Assumptions 1(a) through 1(c) are regularity conditions that rule out cases like perfectly predictable outcomes. Assumption 1(d) provides the substantial restriction on  $\mathcal{P}$ : overlap may be weak in the sense that  $\gamma_0$  is finite, but there is some minimal  $\gamma_0$  and  $C$  that provides a lower bound on the propensity’s tail behavior. Under strict overlap, Assumption 1(d) holds for any finite  $\gamma_0 > 1$ , and most results here hold after replacing  $\gamma_0$  with infinity. Under weak overlap, Assumption 1(d) may only hold for some values of  $\gamma_0$ , in which case the inverse propensity distribution may be heavy-tailed. I refer to the case  $\gamma_0 > 2$  as “somewhat weak overlap” and refer to the case of  $\gamma_0 < 2$  as “very weak overlap.” As  $\gamma_0$  shrinks below 2, overlap is permitted to be increasingly weak.  $\gamma_0 \leq 1$  corresponds to no bound on the propensity distribution.

Figure 1 illustrates behavior for simulated data with various values of  $\gamma_0$ . When the propensity score  $e(X)$  has a well-defined density,  $\gamma_0 = 2$  corresponds to a roughly uniform distribution of propensity scores ([Ma and Wang, 2020](#)). When  $\gamma_0$  is above two, the density of propensity scores tends to zero at zero; when  $\gamma_0$  is below two, the density of propensity scores can tend to infinity at zero. Heuristically, there are never

too many treated observations with very small propensity scores, but  $\gamma_0$  governs the degree to which there can be many untreated observations with very small propensity scores.

A phase transition occurs when  $\gamma_0$  crosses two.

**Proposition 1.** (i) Suppose Assumption 1 holds for some  $\gamma_0 > 2$ . Then the semiparametric bound is finite for all  $P \in \mathcal{P}$ . (ii) Suppose Assumption 1 holds for some  $\gamma_0 \in (1, 2)$ , and there is a  $P \in \mathcal{P}$  and  $C' > 0$  such that  $P(e(X) \leq \pi) \geq C'\pi^{\gamma_0-1}$  for all  $\pi \in (0, 1]$ . Then the semiparametric bound is infinite for  $P$ .

I will require certain rates on the nuisance functions  $e(X)$  and  $\mu(X)$ . I write the worst-case rates as  $r_{e,n}$  and  $r_{\mu,n}$ .

**Assumption 2** (Cross-fitting and  $L_\infty$  rates). The nuisances  $\hat{\mu}$  and  $\hat{e}$  are estimated with cross-fitting with a fixed number of folds  $K$ . If  $n_k$  is the number of observations per fold, then  $\inf_k n_k / \sup_k n_k \rightarrow 1$ . Further, for all  $k \in 1, \dots, K$  and all  $P \in \mathcal{P}$ , the cross-fit nuisances satisfy the uniform consistency rates  $\mathbb{E}_P[\|\hat{\mu}_n^{(-k)} - \mu\|_\infty] \leq r_{\mu,n}$  and  $\mathbb{E}_P[\|\hat{e}_n^{(-k)} - e\|_\infty] \leq r_{e,n}$  where  $r_{\mu,n}, r_{e,n}$  are uniformly bounded above.

Cross-fitting is a common strategy for simplifying the analysis of Neyman-orthogonal estimators like AIPW (Chernozhukov et al., 2018). In practice, nuisances satisfying Assumption 2 may only be achieved with arbitrarily high probability. I impose uniformity to ensure these rates hold in regions of  $x$  with weak overlap, but can be bypassed with  $L_2$  error conditions in other regions. Such uniformity assumptions are standard in studying semiparametric estimators under irregular identification (Semenova, 2024).

## 2.2 Estimator and Consistency

My formal analysis considers the clipped AIPW estimator with cross-fit nuisance function estimates. Results for the other standard thresholding procedure, trimming, generally follow by the same arguments. I begin by providing sufficient conditions for consistency.

For simplicity, I focus the theoretical analysis on estimating the average potential outcome  $\psi = E[DY/e(X)]$ . The average treatment effect follows as a corollary. The clipped AIPW estimator of  $\psi$  is:

$$\hat{\psi}_{clip}^{AIPW}(b_n) = \frac{1}{n} \sum_{k=1}^K \sum_{i \in \mathcal{F}^k} \phi(Z_i | b_n, \hat{\eta}^{(-k)}), \text{ where } \phi(Z | b, \hat{\eta}) = \hat{\mu}(X) + \frac{D(Y - \hat{\mu}(X))}{\max\{\hat{e}(X), b\}}. \quad (1)$$

In that equation,  $\mathcal{F}^k$  is the set of observations  $i$  randomly partitioned in fold  $k$ ,  $\hat{\eta}^{(-k)}$  is the nuisance function estimates constructed only on observations in folds other than  $k$ . The unthresholded AIPW estimator is the special case of  $b_n = 0$ . I analyze the clipped AIPW estimator because results for the trimmed AIPW estimator follow somewhat more easily.



A standard result for the unthresholded AIPW estimator is double robustness: when  $e(X)$  is bounded away from zero, unthresholded AIPW is consistent for  $\psi$  if either  $r_{e,n}$  or  $r_{\mu,n}$  tends to zero. The existence of weak overlap introduces a subtlety to double robustness.

**Proposition 2** (Consistency). *Suppose  $b_n$  satisfies  $n^{-1/2} \ll b_n \ll 1$ , the conditions of Assumption 2 hold, and either (i)  $r_{e,n} b_n^{\min\{\gamma_0-2,0\}} \rightarrow 0$  or (ii)  $r_{\mu,n} \frac{r_{e,n}+b_n}{b_n} \rightarrow 0$ . Then for all  $\epsilon > 0$ ,*

$$\sup_{P \in \mathcal{P}} P \left( \left| \hat{\psi}_{clip}^{AIPW}(b_n) - \psi(P) \right| > \epsilon \right) \rightarrow 0.$$

Under very weak overlap, condition (i) is stronger than the classic strict overlap condition that  $r_{e,n}$  consistency implies estimator consistency. It requires that  $r_{e,n}$  go to zero faster than  $b_n^{2-\gamma_0}$ , so that as overlap is allowed to be weaker, the propensity consistency rate may need to be as fast as  $b_n$  itself. Under even somewhat weak overlap, condition (ii) is stronger than the classic  $r_{\mu,n} \rightarrow 0$  condition. With an inconsistent propensity score, a meaningful fraction of the data may be clipped even asymptotically, in which case the outcome regression error rate must offset the positive probability of incorrectly assigning an inverse propensity weight of  $b_n^{-1}$ .

## 2.3 Asymptotic Normality and Confidence Intervals

This subsection presents the main theoretical claims of the paper. It shows that under suitable rate restrictions, the clipped AIPW estimator is first-order equivalent to an oracle clipped AIPW estimator, both estimators are consistent and asymptotically normal, and simple Wald confidence intervals are well-calibrated.

A common strategy for deriving confidence intervals for unthresholded AIPW under strict overlap is Neyman orthogonality. In those classic settings, the difference between the feasible AIPW estimator with estimated nuisance functions and the hypothetical oracle AIPW estimator with known nuisance functions is

$$\frac{1}{n} \sum \phi(Z | 0, \hat{\eta}) - \phi(Z | 0, \eta) = \frac{1}{n} \sum (\hat{\mu} - \mu) \left( \frac{D}{\hat{e}} - 1 \right) + (Y - \mu) \left( \frac{D}{\hat{e}} - \frac{D}{e} \right).$$

Intuitively, the regression errors  $\hat{\mu} - \mu$  are debiased by the inverse propensity estimates  $\frac{D}{\hat{e}}$  of the number one. As a result, in classical settings, slowly consistent nuisance estimates can yield quickly consistent causal estimates. When all nuisances are consistent at  $o(n^{-1/4})$  rates and  $e(X)$  is bounded away from zero, inverse propensity errors are of the same order as propensity errors, classical AIPW estimates are first-order equivalent to oracle estimates with known nuisances, and simple Wald confidence intervals cover the true causal effect by appeal to the asymptotically normal oracle AIPW estimator.

Under very weak overlap, thresholded AIPW does not obtain the standard debiasing benefit. The analogous decomposition for clipped AIPW is

$$\frac{1}{n} \sum \phi(Z | b_n, \hat{\eta}) - \phi(Z | b_n, \eta) = \frac{1}{n} \sum (\hat{\mu} - \mu) \left( \frac{D}{\max\{\hat{e}, b_n\}} - 1 \right) + (Y - \mu) \left( \frac{D}{\max\{\hat{e}, b_n\}} - \frac{D}{\max\{e, b_n\}} \right).$$

Above the clipping threshold  $b_n$ , thresholded AIPW's nuisance estimation error enjoys a product-of-errors character that is similar to classical settings, albeit with multiplication by weights as large as  $b_n^{-1}$ . Below the clipping threshold, the regression errors are not debiased by any inverse propensity estimate.

Thresholded AIPW enjoys a subtly different form of debiasing: as the threshold tends to zero, increasingly little mass remains. If the threshold  $b_n$  tends to zero quickly enough, the bias in the thresholded region is debiased by the threshold itself. If the threshold  $b_n$  tends to zero slowly enough, the product of errors condition can remain feasible. My formal contribution is to show that there can be a Goldilocks range where  $b_n$  tends to zero neither too slowly nor too quickly, so that thresholded AIPW is first-order equivalent to oracle AIPW and is asymptotically normal by appeal to the asymptotically normal thresholded oracle AIPW estimator.

My results for asymptotic normality and statistical inference will proceed under the following rate requirements.

**Assumption 3** (Minimal rates). Assumption 2 holds, with the following rates on the regression error  $r_{\mu,n}$  and the propensity error  $r_{e,n}$ :

- (a) *Consistency.*  $r_{\mu,n}, r_{e,n} \rightarrow 0$ .
- (b) *Product of errors.*  $r_{\mu,n} r_{e,n} \left( 1 + b_n^{(\gamma_0-2)/2} \log(1/b_n) \mathbf{1}_{\{\gamma_0=2\}} \right) \ll n^{-1/2}$ .
- (c) *Regression error near singularities.*  $r_{\mu,n} b_n^{\gamma_0/2} \ll n^{-1/2}$ .
- (d) *Asymptotically known thresholding.*  $r_{e,n} \ll b_n$ .

Under very weak overlap, conditions (b) and (c) are stronger than the standard product-of-errors condition  $r_{\mu,n} r_{e,n} \ll n^{-1/2}$ . For example, when  $\gamma_0 \geq 1.5$  and  $r_{\mu,n} = n^{-1/4}$ , then  $r_{e,n} \ll n^{-1/3}$  will suffice, provided  $r_{e,n} \ll b_n \ll n^{-1/3}$ . However, these conditions never require parametric  $n^{-1/2}$  consistency rates: shared regression rates of  $n^{-1/3}$  will always suffice for these conditions, provided the clipping threshold  $b_n$  goes to zero at a rate sufficiently close to  $n^{-1/3}$ . I discuss these conditions further in Section 3.2.

Certain technical possibilities call for one of two alternative further assumptions: a distributional smoothness assumption or a stronger rate assumption.

**Assumption 4** (Nongeneracy or faster rates). One of the following two conditions hold:

- (i) *Nondegenerate overlap.* There exists some  $\rho > 0$  such that for all  $P \in \mathcal{P}$  and  $\pi \in [0, 1]$ ,  $P(e(X) \leq \pi/2) \leq (1 - \rho)P(e(X) \leq \pi)$ .
- (ii) *Faster rates.*  $r_{\mu,n} b_n^{(\gamma_0-1)2/\gamma_0} \ll n^{-1/2}$ .

Assumption 4(i) is a uniform version of the requirement that  $P(e(X) \leq x) = c(x)x^{\gamma_0-1}$  for  $c(x)$  tending to a constant at zero. The definition formalizes the notion that a distribution may place some propensity mass near zero, but it may not place mass adversarially within the region of the origin. When  $\gamma_0 < 2$ , Assumption 4(ii) is stronger than Assumption 3(c). As  $\gamma_0$  tends to one, the condition approaches the parametric requirement  $r_{\mu,n} = O(n^{-1/2})$ .

I now provide the main theoretical result.

**Theorem 1** ((Slow) Asymptotic Normality). *Suppose  $b_n$  satisfies  $n^{-1/2} \ll b_n \ll 1$ , and Assumptions 1, 2, 3, and 4 hold. Then the clipped AIPW estimator is oracle-equivalent:*

$$\lim_{n \rightarrow \infty} \sup_{P \in \mathcal{P}} \sigma_n^{-2} E_P \left[ \left( \hat{\psi}_{clip}^{AIPW}(b_n) - \tilde{\psi}_{(Oracle)}^{AIPW}(b_n) \right)^2 \right] = 0.$$

*Further, clipped AIPW is asymptotically normal:*

$$\lim_{n \rightarrow \infty} \sup_{P \in \mathcal{P}} \sup_{t \in \mathbb{R}} \left| P \left( \frac{\hat{\psi}_{clip}^{AIPW}(b_n) - \psi(P)}{\hat{\sigma}_n} \leq t \right) - \Phi(t) \right| = 0.$$

Theorem 1 is the core theoretical claim of this paper. The first result shows that thresholded AIPW is first-order equivalent to an oracle estimator with known nuisances: the effect of nuisance estimation error on the treatment effect estimate tends to zero faster than the standard deviation of the oracle estimator. The second result leverages this first-order equivalence to characterize the asymptotic distribution of the clipped AIPW estimates and t-statistics: the estimator is asymptotically normal, and estimated t-statistics are asymptotically standard normal.

Both results are standard for AIPW under strict overlap, but substantial care is required to handle unbounded inverse propensities under weak overlap. The argument for normality builds on Ma and Wang (2020)'s proof that aggressively-trimmed oracle IPW with known propensities achieves asymptotic normality with first-order bias. I extend their argument to a uniform family of distributions using the Berry-Esseen Theorem and note that oracle AIPW must have zero finite-sample bias.

The main task of Theorem 1 is to show that replacing the true nuisances with estimated nuisances has a second-order effect on clipped AIPW estimates under appropriate conditions. This is nontrivial, because under weak overlap, there is an asymptotically unbounded number of observations with arbitrarily large

inverse propensities with even known nuisance functions. Nevertheless, by taking appropriate care and leveraging that clipping introduces bias by reducing inverse propensities, I am able to show that the effect of nuisance estimation is second-order even under the very weak overlap case in which unthresholded AIPW fails to be asymptotically normal and no regular root-n estimators exist.

Next, I show that Theorem 1 yields the natural result for inference: simple t-tests based on Wald confidence intervals are well-calibrated.

**Corollary 1** (T-tests are well-calibrated). *Suppose the conditions of Theorem 1 hold. Consider the Wald confidence interval  $\hat{C}_n(\alpha) = [\hat{\psi}_{clip}^{AIPW}(b_n) + z_{\alpha/2}\hat{\sigma}_n, \hat{\psi}_{clip}^{AIPW}(b_n) + z_{1-\alpha/2}\hat{\sigma}_n]$ . Then for all  $\alpha \in (0, 1/2)$ ,*

$$\limsup_{n \rightarrow \infty} \sup_{P \in \mathcal{P}} \left| P(\psi(P) \in \hat{C}_n(\alpha)) - (1 - \alpha) \right| = 0.$$

Under somewhat weak overlap, unthresholded AIPW is semiparametrically efficient. I now show thresholding is also unnecessary.

**Corollary 2** (Thresholding is second-order under somewhat weak overlap). *Suppose Assumption 2 holds for some  $\gamma_0 > 2$ ,  $r_{e,n}$  and  $r_{\mu,n} \rightarrow 0$ , and  $r_{e,n}r_{\mu,n} \ll n^{-1/2}$ . Then the feasible AIPW estimator with  $b_n = 0$  is semiparametrically efficient, and the associated Wald confidence interval  $\hat{C}_n(\alpha)$  satisfies*

$$\limsup_{n \rightarrow \infty} \sup_{P \in \mathcal{P}} \left| P(\psi(P) \in \hat{C}_n(\alpha)) - (1 - \alpha) \right| = 0.$$

The logic of Corollary 2 is to show that any sequence of  $b_n \rightarrow 0$  has a second-order effect on estimation, so that unthresholded and thresholded AIPW are first-order equivalent, and there is some  $b_n \rightarrow 0$  slowly enough satisfying the conditions of Theorem 1, so that there is a thresholded AIPW estimator achieving Corollary 1.

Taken together, this subsection yields a remarkable result for practice. The distribution  $P$  may place so much propensity mass near the origin that the semiparametric efficiency bound is infinite, the lower bound on the density of propensity mass near the origin can be so weak that identification nearly fails, and the nuisance estimator may be so poorly designed that it pushes all observations' estimated propensities towards the origin at a slower-than-parametric rate. Nevertheless, Neyman orthogonality is sufficiently powerful to ensure the validity of the simple t-test.

The next section interprets the rate requirements of Assumption 3.

Table 1: Summary of degradation of asymptotic behavior and requirements as overlap is permitted to be increasingly weak.

Overlap Phase	Strict ( $\inf e(x) > 0$ )	Somewhat Weak ( $\gamma_0 > 2$ )	Very Weak ( $\gamma_0 < 2$ )
<b>Double Robustness</b> Consistency Conditions	$r_{e,n} \rightarrow 0$ or $r_{\mu,n} \rightarrow 0$	$r_{e,n} \rightarrow 0$ or $r_{\mu,n} b_n^{-1} \rightarrow 0$	$r_{e,n} b_n^{\gamma_0-2} \rightarrow 0$ or $r_{\mu,n} r_{e,n} b_n^{-1} \rightarrow 0$
<b>Oracle AIPW Asymptotics</b> Unthresholded Distribution	Normal	Normal	Nonnormal
Thresholded Convergence Rate	$n^{-1/2}$	$n^{-1/2}$	$n^{-1/2} b_n^{(\gamma_0-2)/2}$
Semiparametrically Efficient?	Yes	Yes	No
<b>Nuisance Requirements</b> Black Box:	$n^{1/2} r_{\mu,n} r_{e,n} \rightarrow 0$ ( $L_2$ )	$n^{1/2} r_{\mu,n} r_{e,n} \rightarrow 0$ ( $L_\infty$ )	$n^{1/2} r_{\mu,n} r_{e,n}^{\gamma_0/2} \rightarrow 0$ ( $L_\infty$ )
Smoothness:	$\frac{\beta_\mu}{2\beta_\mu+d} + \frac{\beta_e}{2\beta_e+d} > \frac{1}{2}$	$\frac{\beta_\mu}{2\beta_\mu+d-\frac{\gamma_0}{\gamma_0-1}} + \frac{\beta_e}{2\beta_e+d} > \frac{1}{2}$	$\frac{\beta_\mu}{2\beta_\mu+d-\frac{\gamma_0}{\gamma_0-1}} + \frac{\frac{\gamma_0}{2}\beta_e}{2\beta_e+d} > \frac{1}{2}$
<b>Regression Rates</b> Pointwise optimum:	$n^{\frac{-\beta_\mu}{2\beta_\mu+d}}$	$n^{\frac{-\beta_\mu(1-1/\gamma_0)}{2\beta_\mu(1-1/\gamma_0)+d}}$	$n^{\frac{-\beta_\mu(1-1/\gamma_0)}{2\beta_\mu(1-1/\gamma_0)+d}}$
Global optimum:	$(n/\log(n))^{\frac{-\beta_\mu}{2\beta_\mu+d}}$	$n^{\frac{-\beta_\mu(1-1/\gamma_0)}{2\beta_\mu(1-1/\gamma_0)+d}}$	$n^{\frac{-\beta_\mu(1-1/\gamma_0)}{2\beta_\mu(1-1/\gamma_0)+d}}$

### 3 Interpretation of Nuisance Requirements

This section interprets the rate requirements for Wald confidence intervals to cover asymptotically. I summarize the results in Table 1. Under somewhat weak overlap, thresholded and unthresholded AIPW remain semiparametrically efficient and  $\sqrt{n}$ -consistent, and the traditional product of errors nuisance condition remains in place with a modification to an  $L_\infty$  on errors. Under very weak overlap, clipped AIPW achieves a slower consistency rate and the required black box nuisance rates are more stringent. Both cases make outcome regression more difficult, but never so difficult as to require parametric assumptions. As a byproduct of this analysis, I show that the optimal pointwise and global regression rates under weak overlap are the same, without the usual polylogarithmic factor in the global rate.

#### 3.1 Degradation of Consistency Rate

The previous analysis suggests that smaller values of  $b_n$  are preferable because they admit weaker black-box requirements. However, under very weak overlap, larger values of  $b_n$  correspond to faster AIPW rates.

I characterize the consistency rate of any oracle-equivalent estimator as follows.

**Proposition 3** (Consistency rate). *Suppose the assumptions of Proposition 2 hold. Then there exist positive constants  $c_{\min}$  and  $c_{\max}$  such that  $c_{\min} n^{-1} \mathbb{E}_P \left[ \frac{D}{\max\{e(X), b_n\}^2} \right] \leq \sigma_n^2 \leq c_{\max} n^{-1} \mathbb{E}_P \left[ \frac{D}{\max\{e(X), b_n\}^2} \right]$  for all  $P \in \mathcal{P}$ , where  $\sigma_n^2 = n^{-1} \left( \frac{1}{n} \sum \phi(Z | b_n, \eta)^2 - \tilde{\psi}_{(Oracle)}^{AIPW}(b_n)^2 \right)$  is the oracle sample variance.*

If the estimator were trimmed instead of clipped,  $\mathbb{E}_P \left[ \frac{D}{\max\{e(X), b_n\}^2} \right]$  would be replaced by a term like  $\mathbb{E}_P \left[ \frac{D \mathbf{1}\{e(X) \geq b_n\}}{e(X)} \right]$ . Weaker overlap corresponds to larger values of  $\mathbb{E}_P[D/\max\{e(X), b_n\}^2]$  and slower consistency rates. Conditional on  $P$ , larger values of  $b_n$  correspond to a smaller value of  $\mathbb{E}_P[D/\max\{e(X), b_n\}^2]$ , faster oracle consistency, and greater asymptotic power.

Proposition 3 implies a worst-case consistency rate over distributions in  $\mathcal{P}$ . I focus on the case of very weak overlap, because Corollary 2 shows that under somewhat weak overlap, clipped and traditional AIPW achieve a traditional  $\sqrt{n}$  consistency rate.

**Corollary 3** (Worst-case consistency rate). *Suppose  $\gamma_0 < 2$  and let  $b_n$  be a fixed sequence of  $b_n$  satisfying  $1 \gg b_n \gg n^{-1/2}$ . There exists a  $C' > 0$  such that for  $\mathcal{P}$  satisfying Assumption 1,  $C' n^{-1} b_n^{\gamma_0-2} \geq \sup_{P \in \mathcal{P}} \sigma_n^2$  for all  $n$  large enough. Further, there exists a family  $\mathcal{P}$  satisfying Assumption 1 and a  $C'' \in (0, C')$  such that  $\sup_{P \in \mathcal{P}} \sigma_n^2 \geq C'' n^{-1} b_n^{\gamma_0-2}$  for all  $n$  large enough.*

The rate  $n^{-1} b_n^{\gamma_0-2}$  is a worst-case consistency rate in  $b_n$ : every distribution in  $\mathcal{P}$  achieves a consistency at least as fast as  $n^{-1} b_n^{\gamma_0-2}$ , and it is possible to find a distribution for which the consistency rate is no faster. The combination of Corollary 3 and Theorem 1 yields a trade-off under very weak overlap: smaller values of  $b_n$  yield laxer requirements on regression estimation near singularities, but lead to larger variance and slower consistency.

### 3.2 Degradation of Black-Box Nuisance Requirements

Under very weak overlap, the black-box rates of Assumption 3 are more stringent than the usual  $r_{\mu,n} r_{e,n} \ll n^{-1/2}$  condition. The main requirement is that that  $r_{\mu,n} r_{e,n}^{\min\{\gamma_0/2, 1\}}$  goes to zero faster than  $n^{-1/2}$ . As a result, outcome regression rates are more valuable than nominally equivalent propensity rates under very weak overlap.

The usual product-of-errors condition under strict overlap often takes a form like  $r_{\mu,n} r_{e,n} \ll \sigma_n$ , where  $\sigma_n$  is the standard deviation of the Oracle estimator. For example, Heiler and Kazak (2021) argue that this condition is sufficient for estimated unthresholded AIPW to be first-order equivalent to an oracle estimator. An alternative characterization of the usual product-of-errors condition that  $r_{\mu,n} r_{e,n} \ll n^{-1/2}$ ; this requirement is equivalent under somewhat weak overlap, but is more stringent under very weak overlap. Even this stronger product-of-errors requirement is insufficient for Wald confidence intervals constructed with thresholded AIPW to be valid.

**Corollary 4** (Under very weak overlap, faster rates are necessary). *For any  $\gamma_0 \in (1, 2)$  and sequence  $b_n$  satisfying  $n^{-1/2} \ll b_n \ll 1$ , there is a family  $\mathcal{P}$  and cross-fit nuisance estimators  $\hat{\mu}$  and  $\hat{e}$  such that:*

1. Overlap is not too weak.  $\mathcal{P}$  satisfies Assumption 1 for this  $\gamma_0$ .
2. Fast nuisance rates.  $E_{P \in \mathcal{P}} [\|\hat{e} - e\|_\infty] \ll b_n$  and  $\sup_{P \in \mathcal{P}} E_P [\|\hat{\mu} - \mu\|_\infty] \times E_P [\|\hat{e} - e\|_\infty] \ll n^{-1/2}$ .
3. Wald inference fails. For any fixed target coverage level  $\alpha \in (0, 1)$ , the Wald confidence intervals  $\hat{\mathcal{C}}_n(\alpha)$  based on AIPW with clipping at  $b_n$  have the zero-coverage property  $\sup_{P \in \mathcal{P}} P(\psi(P) \in \hat{\mathcal{C}}_n(\alpha)) \rightarrow 0$ .

A heuristic sufficient condition for Wald confidence interval validity is  $r_{\mu,n} r_{e,n}^{\min\{\gamma_0, 2\}/2} \ll n^{-1/2}$ . When  $\gamma_0 < 2$ , there is a range of nuisance estimates such that  $r_{\mu,n} r_{e,n} \ll n^{-1/2} \ll r_{\mu,n} r_{e,n}^{\min\{\gamma_0, 2\}/2}$  and estimation bias can be of a higher order than the oracle variance.

I calculate the black-box rate requirements under Assumption 4(i) in a few special cases. I omit an analysis of the stronger rate requirement in Assumption 4(ii) that would be needed to handle degenerate distributions.

**Assumption 5.** Assumptions 1, 2, and 4(i) hold, and  $r_{e,n}, r_{\mu,n} \rightarrow 0$ .

I characterize the following special cases.

**Example 1** (Somewhat weak overlap). Suppose Assumption 5 holds,  $\gamma_0 > 2$ , and  $r_{\mu,n} r_{e,n} \ll n^{-1/2}$ . Then there exists a  $b_n \rightarrow 0$  such that clipped AIPW t-statistics are asymptotically well-calibrated.

**Example 2** (Second moments barely fail to exist). Suppose Assumption 5 holds for  $\gamma_0 = 2$  and there is some  $\eta > 0$  such that  $r_{\mu,n} r_{e,n} \log(1/r_{e,n}) \ll n^{-1/2}$ . Then there exists a  $b_n \rightarrow 0$  such that clipped AIPW t-statistics are asymptotically well-calibrated.

**Example 3** (Shared rates, very weak overlap). Suppose Assumption 5 holds for some  $\gamma_0 > 1$  and  $r_{\mu,n}, r_{e,n} \ll n^{-1/3}$ . Then there exists a  $b_n \rightarrow 0$  such that clipped AIPW t-statistics are asymptotically well-calibrated.

**Example 4** (Parametric rates). Suppose Assumption 5 holds for some  $\gamma_0 > 1$  and either (i)  $r_{\mu,n} = O(n^{-1/2})$  and  $r_{e,n} = o(1)$  or (ii)  $r_{e,n} = O(n^{-1/2})$  and  $r_{\mu,n} = o(n^{(\gamma_0-2)/4})$ . Then there exists a  $b_n \rightarrow 0$  such that clipped AIPW t-statistics are asymptotically well-calibrated.

I now unpack the black box and quantify the degree to which weak overlap makes a given outcome regression rate more difficult to achieve.

### 3.3 Necessary Smoothness Conditions

Weak overlap makes outcome regression more difficult: there may be few treated observations in precisely the regions in which thresholded AIPW depends most acutely on outcome regression. In this section, I show that for pointwise rates, even somewhat weak overlap can be viewed as degrading effective outcome

smoothness for optimal local polynomial estimators. I also derive a blessing of weak overlap: the optimal global rate is equal to the optimal pointwise rate, without the usual polylogarithmic penalty.

I characterize optimal nonparametric regression rates under Hölder continuity. For convenience, I fix the covariates to be uniform over a specific hypercube in  $\mathbb{R}^d$  with constant variance.

**Assumption 6** (Hölder smoothness and fixed domain).  $\mathcal{P}$  satisfies Assumptions 1 and 4(i), and for all  $P \in \mathcal{P}$ ,  $X \sim \text{Unif}([-1, 1]^d)$  and  $Y \mid X, D \sim \mathcal{N}(D\mu_P(X) + (1 - D)\mu'_P(X), \sigma_{\min}^2)$ , with  $\mu, \mu'$  in the Hölder smoothness class  $\Sigma(\beta_\mu, L)$  for some fixed  $\beta_\mu, L > 0$ .

Most of these assumptions are standard assumptions for studying local polynomial regression under strict overlap (Stone, 1982). I assume normal outcomes in order to simplify the characterization of the optimal rate. It is known that in this case but under strict overlap, the optimal pointwise rate is  $n^{\frac{-\beta_\mu}{2\beta_\mu + d}}$ , and the optimal global rate  $(n/\log(n))^{\frac{-\beta_\mu}{2\beta_\mu + d}}$  has a polylogarithmic penalty, in the sense that the optimal global rate is worse by some polynomial factor of  $\log(n)$ .

I require that treated observations cannot concentrate in small regions.

**Assumption 7** (Non-trivial concentration). There are parameters  $\rho, \gamma > 0$  such that for all  $h > 0$  small enough and all  $P \in \mathcal{P}$  and  $x_0 \in [-1, 1]^d$   $P(e(X) \geq \rho \sup_{\|x - x_0\| \leq h} e(x) \mid D = 1, \|X - x_0\| \leq h) > \gamma$ .

I show in Appendix Proposition 5 that this condition holds if the propensity function is sufficiently smooth. When the propensity function is nonsmooth, it is possible for nature to concentrate treated observations within a given bandwidth in a region that is too small, introducing local polynomial degeneracy issues that are outside the scope of this work.

In the worst case, weak overlap of order  $\gamma_0 > 1$  plays a role equivalent to scaling the effective outcome smoothness downward by  $(1 - 1/\gamma_0)$ , but also removes the polylogarithmic factor in the optimal global rate.

**Theorem 2** (Weak overlap reduces effective outcome smoothness). Define  $\psi_n = n^{\frac{-\beta^*}{2\beta^* + d}}$ , where  $\beta^* = \beta_\mu(1 - 1/\gamma_0)$ . Then

- (i)  $\psi_n$  is a pointwise (and global) rate upper bound. There exists a  $c > 0$  and a  $\mathcal{P}$  satisfying Assumptions 6 and 7 such that

$$\liminf_{n \rightarrow \infty} \inf_{\hat{\mu}} \sup_{P \in \mathcal{P}} P(|\hat{\mu}(0) - \mu(0)| > c\psi_n) > 0.$$

- (ii)  $\psi_n$  is an achievable global (and pointwise) rate. Suppose  $\mathcal{P}$  satisfies Assumptions 6 and 7. Then there exists an estimator  $\hat{\mu}(x)$  such that for all  $\epsilon > 0$ , there is a finite  $c(\epsilon)$  such that

$$\limsup_{n \rightarrow \infty} \sup_{P \in \mathcal{P}} P(\|\hat{\mu} - \mu\|_\infty > c(\epsilon)\psi_n) \leq \epsilon.$$



Further, the estimator can be computed without knowledge of the overlap bound  $\gamma_0$ .

Recall that under strict overlap, the optimal pointwise convergence rate for a Hölder-smooth regression function is  $n^{\frac{-\beta_\mu}{2\beta_\mu+d}}$ . Theorem 2 shows that weak overlap has the effect of degrading the effective smoothness rate from  $\beta_\mu$  to  $\beta_\mu(1 - 1/\gamma_0)$ . As overlap is allowed to become increasingly weak and other parameters are held constant, there can be regions of the covariate space with increasingly few treated observations so that the optimal pointwise regression rate is slower. This penalty occurs even under somewhat weak overlap. In the limit in which  $\gamma_0$  tends to one, the rate in Theorem 2 can become arbitrarily poor. For example, when  $\gamma_0 = 2$ , the difficulty of estimating a twice continuously differentiable function is comparable to the difficulty of estimating a Lipschitz-continuous function under strict overlap.

Theorem 2 also presents a blessing of weak overlap: the optimal global rate is equal to the pointwise rate, with no polylogarithmic penalty. Usually, optimal global rates are worse than optimal pointwise rates by a polylogarithmic factor due to the need to have accuracy at a number of gridpoints that grows with  $n$ . Under weak overlap, the global rate still must be consistent at all gridpoints, but most gridpoints must satisfy strict overlap and have a negligible contribution to the global consistency rate. The challenge is to partition the remaining observations in a way which does not impose a polylogarithmic penalty. Naively partitioning the remaining observations into regions of increasing overlap with a constant penalty yields  $\log(\log(n))$  partitions and a penalty potentially as large as  $\log(\log(\log(n)))$ . The proof of Theorem 2 instead uses a more subtle construction, partitioning  $[-1, 1]^d$  into observations with increasing overlap in such a way as each partition's worst-case global penalty is half as large as the previous worst-case penalty, even after accounting for the increasing number of gridpoints. Appendix Proposition 6 shows that if smallest partition's penalty is sufficiently *large*, then this construction eventually includes all relevant observations. Then, summing over the penalties yields a global penalty that is large but bounded, so that the optimal pointwise rate is an achievable (and therefore optimal) global rate.

Theorem 2 yields minimal smoothness assumptions for Wald confidence interval validity.

**Corollary 5** (Minimal smoothness conditions). *Suppose Assumption 6 holds and there is a  $\beta_e > 0$  such that  $e(X) \in \Sigma(\beta_e, L)$  and*

$$\frac{\beta_\mu}{2\beta_\mu + d\gamma_0/(\gamma_0 - 1)} + \frac{\min\{\gamma_0/2, 1\}\beta_e}{2\beta_e + d} > 1/2. \quad (2)$$

*Then there is a sequence of nuisance estimators and thresholds that are independent of  $\gamma_0$  such that for all*

$\gamma_0 > 1$ , the associated Wald confidence interval  $\hat{\mathcal{C}}_n(\alpha)$  constructed using  $\hat{\psi}_{clip}^{AIPW}(b_n)$  satisfies

$$\limsup_{n \rightarrow \infty} \sup_{P \in \mathcal{P}} \left| P(\psi(P) \in \hat{\mathcal{C}}_n(\alpha)) - (1 - \alpha) \right| = 0.$$

In one dimension, thresholded AIPW with Lipschitz-continuous conditional outcome mean and propensity function can handle weak overlap of order  $\gamma_0 > \frac{4+1/\beta_e}{3}$ . In multiple dimensions, the econometrician must assume stronger smoothness restrictions than Lipschitz continuity in order to achieve the necessary nuisance rate guarantees under even strict overlap. When the propensity function is infinitely-differentiable, thresholded AIPW with a Lipschitz-continuous conditional outcome mean can handle weak overlap of order  $\gamma_0 > \frac{2(d+1)}{d+2}$ . More generally, under very weak overlap, if  $\beta_\mu, \beta_e > \frac{d(\sqrt{\gamma_0^2 + 4\gamma_0 - 4} + 2 - \gamma_0)}{4(\gamma_0 - 1)}$ , then it is feasible to achieve standard inference with thresholded AIPW.

This concludes the substantive theoretical analysis. In the next section, I use these results to infer lessons for empirical practice.

## 4 Lessons for Empirical Practice

This section leverages the theoretical analysis to consider misspecified parametric estimators and some rules of thumb for empirical use.

### 4.1 Parametric Estimators and Misspecification

When both nuisance functions are estimated nonparametrically, then consistency is achievable and AIPW is generally preferable under strict overlap. When both nuisance functions are estimated parametrically, then it is possible for one or both nuisance function to be inconsistent and the choice of estimator may be ambiguous. I now provide some intuition on the two estimators when nuisance functions are estimated parametrically and through cross-fitting. I will consider IPW and AIPW with the same sequence of thresholds  $b_n$  satisfying  $1 \gg b_n \gg n^{-1/2}$ . I write that a nuisance estimate  $\hat{\eta}$  is consistent if it tends to the correct limit  $\eta$ , and I write that  $\hat{\eta}$  is inconsistent otherwise.

In this subsection, I will assume that parametric nuisance estimators  $\hat{\eta}$  achieve an  $L_\infty$  error relative to a limiting nuisance function  $\bar{\eta}$  that is the order of  $n^{-1/2}$ . For example, consider logit estimation of a propensity model of the form  $\bar{e}(X) = \frac{\exp(X'\beta)}{1 + \exp(X'\beta)}$  for a pseudo-true parameter  $\beta$ . If the support of  $X$  is bounded, then  $n^{-1/2}$ -consistent estimate of  $\beta$  is sufficient to achieve  $n^{-1/2}$ -consistent estimation of  $\bar{e}(X)$  everywhere. However, weak overlap may emerge from unbounded tails, in which case the  $L_\infty$  rate may not go to zero. Unbounded covariates are an important case in general. For example, [Ma and Wang \(2020\)](#) motivate

weak overlap tails through the distribution of covariates under a logistic propensity model. Nevertheless, a careful treatment of parametric estimation of nuisances with unbounded covariates is outside the scope of this work.

The analysis above is easiest to extend when either both or neither nuisance function is consistent. If both nuisance estimates are consistent, then the AIPW and IPW estimators will be consistent and will have variance on the same order, but the IPW estimator may have higher-order bias than the AIPW estimator. This higher-order bias follows because IPW can be viewed as a particular case of AIPW with an inconsistent outcome regression estimator. If both the propensity and outcome regression estimates are inconsistent, then both the IPW and AIPW estimators fail to be consistent, and as in the case of inconsistent nuisance functions with strict overlap, there is no general reason to prefer one or the other.

When the outcome regression estimate is inconsistent, there is no general reason to prefer IPW or AIPW, but both estimators may have bias that is of a higher-order than the estimator's standard deviation. When  $\hat{\mu}$  is inconsistent, both IPW and AIPW can be viewed as instances of AIPW with an inconsistent outcome regression estimate. Suppose  $P$  is a distribution from the second half of Corollary 3, which has  $P(e(X) \leq \pi) \sim \pi^{\gamma_0-1}$  for all  $\pi$  small enough. The bias in the thresholded region with an inconsistent outcome regression estimate is generally on the order of  $P(e(X) \leq b_n) \sim b_n^{\gamma_0-1}$ . However, by Corollary 3, the oracle AIPW (and oracle IPW) standard deviation is on the order of  $n^{-1/2}b_n^{\gamma_0/2-1} \ll b_n^{\gamma_0-1}$ . This heuristic analysis suggests that in many cases, IPW or AIPW-with-inconsistent-outcome-regression will have bias that is of a higher order than the estimator's standard error. That intuition is similar to Ma et al. (2023)'s analysis of trimmed AIPW with a tailored debiasing procedure.

The case of a consistent outcome regression estimate with inconsistent propensity estimates is more interesting. In this case, AIPW should have lower-order bias than IPW, because IPW will be inconsistent. The Berry-Esseen argument for AIPW asymptotic normality with known nuisance functions only requires cross-fitting and  $b_n$  to go to zero slower than  $n^{-1/2}$ , so that thresholded AIPW should also be asymptotically normal under appropriate error product conditions. However, it is unclear how the bias compares to sampling error. In any event, this robustness intuition is useful, because I apply parametric nuisance estimators in the application to right heart catheterization. Careful treatment of the parametric case is left for future work.

Before proceeding to apply clipped AIPW, I derive some rules of thumb for choosing a threshold.

## 4.2 Choice of Threshold

The theoretical analysis above provides conditions under which there is some sequence of thresholds for which AIPW is asymptotically normal and centered around the true causal estimand. I now provide guidance

for how to choose the threshold.

**Input:** Rate upper bounds  $\bar{r}_{\mu,n}$  and  $\bar{r}_{e,n}$  (maybe null), data  $\{D_i, \hat{e}_i\}_{i=1}^n$

**Output:** Rule-of-thumb threshold  $b_n$

```

Function calculateRuleofThumb( $\bar{r}_{\mu,n}, \bar{r}_{e,n}, \{D_i, \hat{e}_i\}_{i=1}^n$ ):
  if  $is.null(\bar{r}_{\mu,n})$  and  $!is.null(\bar{r}_{e,n})$  then
     $b_n \leftarrow \bar{r}_{e,n}$ 
  else
     $A_\mu \leftarrow !is.null(\bar{r}_{\mu,n})$ 
     $A_e \leftarrow !is.null(\bar{r}_{e,n})$ 
     $b_n \leftarrow \sup b : 0 \geq \text{errorBoundDiff}(b, \bar{r}_{\mu,n}, \bar{r}_{e,n}, \{D_i, \hat{e}_i\}_{i=1}^n, A_\mu, A_e)$ 
  end
  return  $b_n$ 
Function errorBoundDiff( $b, \bar{r}_{\mu,n}, \bar{r}_{e,n}, \{D_i, \hat{e}_i\}_{i=1}^n, A_\mu, A_e$ ):
   $\tilde{r}_{\mu,n} \leftarrow A_\mu \bar{r}_{\mu,n} + (1 - A_\mu)b$ 
   $\tilde{r}_{e,n} \leftarrow A_e \bar{r}_{e,n} + (1 - A_e)b$ 
   $\text{second\_moment} \leftarrow \frac{1}{n} \sum_{i=1}^n \frac{D_i}{\max\{\hat{e}_i, b\}^2}$ 
   $\text{error\_bound} \leftarrow \tilde{r}_{\mu,n} \frac{\frac{1}{n} \sum \mathbf{1}\{\hat{e}_i \leq b\}}{\sqrt{\text{second\_moment}}} + \tilde{r}_{\mu,n} \tilde{r}_{e,n} \sqrt{\text{second\_moment}}$ 
  return  $\text{error\_bound} - n^{-1/2}$ 

```

**Algorithm 1:** Rule-of-thumb for choice of threshold  $b_n$  given (possibly null) rate upper bounds  $\bar{r}_{\mu,n}$  and  $\bar{r}_{e,n}$ , to ensure that  $r_{\mu,n} \ll \bar{r}_{\mu,n}$  (or  $r_{\mu,n} \ll b_n$  if no bound is given) and  $r_{e,n} \ll \bar{r}_{e,n}$  (or  $r_{e,n} \ll b_n$  if no bound is given) is sufficient to achieve well-calibrated Wald confidence intervals.

I propose different rules of thumb based on whether the econometrician is willing to provide an upper bound on the rate of convergence for the propensity estimate, the outcome regression estimate, or both. The combined proposal is presented in Algorithm 1.

In practice, it is often relatively easy to identify an upper bound on the propensity rate of convergence. For instance, if the propensity score is estimated with local polynomial regression of order  $\ell_e$ , then the econometrician is implicitly asserting a Hölder smoothness of some order  $\beta_e > \ell_e$ , and a feasible consistency rate of  $(n/\log(n))^{-\beta_e/(2\beta_e+d)}$  (Stone, 1982). In this case, the econometrician can safely conjecture that if their estimator is well-founded, then it will achieve a global consistency rate  $r_{e,n} \ll n^{-\ell_e/(2\ell_e+d)}$ , and a threshold of  $b_n = n^{-\ell_e/(2\ell_e+d)} = “\bar{r}_{e,n}”$ . This rule of thumb is practical and the safest rule with respect to outcome regression that does not impose stronger requirements on the propensity estimate, but also often corresponds to a relatively slow consistency rate with relatively little power.

Three alternative rules of thumb target faster consistency rates and laxer propensity requirements. While the main text focuses on black-box requirements on consistency rates in terms  $\gamma_0$ , the proof goes through showing  $b_n$  tends to zero slower than  $n^{-1/2}$ ,  $r_{e,n} \ll b_n$ , and

$$r_{\mu,n} \frac{P(e(X) \leq b_n)}{\sqrt{E_P \left[ \frac{D}{\max\{e(X), b_n\}^2} \right]}} + r_{\mu,n} r_{e,n} \sqrt{E_P \left[ \frac{D}{\max\{e(X), b_n\}^2} \right]} \ll n^{-1/2}. \quad (\text{Appendix Assumption 3'})$$

In these alternative cases, I propose replacing  $b_n$  in the right-hand side with the putative threshold  $b$  and  $r_{\mu,n}$  and  $r_{e,n}$  in the left-hand side with predicted upper bounds  $\bar{r}_{\mu,n}$  and  $\bar{r}_{e,n}$  where feasible and with the putative threshold  $b$  otherwise. This yields an empirical function `error_bound(b)`. I then solve for the putative threshold  $b$  where `error_bound(b)` crosses  $n^{-1/2}$ . The result is a threshold  $b_n$  aimed to target maximal efficiency and minimal propensity consistency requirements, calculated without use of the outcome data or direct knowledge of  $\gamma_0$ , and with the property that if the nuisance errors go to zero faster than the specified upper bound (or estimated threshold), then the resulting Wald confidence intervals will be well-calibrated. An interesting avenue for future work is whether there is a convenient choice of context-dependent constant multiples in the `error_bound` function.

This rule of thumb is also always feasible, for example in the case of nonspecified upper bounds.

**Lemma 1** (Well-defined rule of thumb). *Suppose  $\hat{e} \in (0, 1]$  and  $\sum D/\hat{e} > 0$ . Let  $f_n(b)$  be the `error_bound` function with no upper bound nuisance rates given. Then there is exactly one  $b_n$  such that  $\limsup_{b \rightarrow b_n^-} f_n(b) \leq 0 \leq \liminf_{b \rightarrow b_n^+} f_n(b)$ .*

A rule of thumb with nonspecified rates seems particularly attractive, since it is an entirely data-driven way to choose the AIPW threshold. However, in practice, a given outcome regression rate is more difficult to achieve than a given propensity rate under weak overlap, so that rules of thumb with specified nuisance rates is more appropriate for nonparametric outcome regression estimates.

## 5 Applications

In this section, I present simulated results for the clipped AIPW estimator as well as empirical results from an application to right heart catheterization. I find that clipped AIPW performs well asymptotically, producing near-perfect calibration of p-values with 100,000 observations, but exhibits some undercoverage in small samples. When studying the right heart catheterization data, I find that the rule of thumb approach increases the estimated harm of the procedure by 0.17 standard errors relative to the usual 10% trimming rule, while increasing the estimated standard error by 5.1%.

### 5.1 Simulation Evidence

I now study the performance of the clipped AIPW estimator in simulations.

My simulation design is based on the design in [Ma and Wang \(2020\)](#). As in their work, I simulate data with  $P(e(X) \leq \pi) = \pi^{\gamma-1}$  and  $DY = \kappa D(1 - e(X)) + D(\varepsilon - 4)/\sqrt{8}$ , where  $\varepsilon \mid X, D \sim \xi_4^2$  is scaled to achieve zero mean and unit variance. However, I increase  $\gamma$  from 1.5 to 1.8 to ensure feasible outcome regression

rates, set  $\kappa = 2$  rather than  $\kappa = 1$  to avoid coincidental offsetting bias of IPW lower and upper tails in small samples, and reduce  $DY$  by  $\kappa E[D(1 - e(X))]$  so that the true average potential outcome is zero. I achieve this propensity distribution by taking  $X \sim Unif([0, 1])$  i.i.d. and setting  $e(X) = X^{1/(\gamma_0 - 1)}$ . I present results for 5,000 simulations of increasingly large samples.

I estimate both the propensity and outcome regressions with five-fold cross-fitting. I use shrinkage cubic splines and REML estimation, as implemented by the `mgcv` package in R. In this setting, Theorem 2 establishes that kernel regression can achieve a pointwise rate of  $n^{-1/(3+1/(\gamma_0-1))}$ . I conjecture that  $r_{\mu,n} \ll n^{-1/5}$ , which is feasible if  $\gamma_0 > 1.5$ , and choose the clipping threshold  $b_n$  based on Algorithm 1.

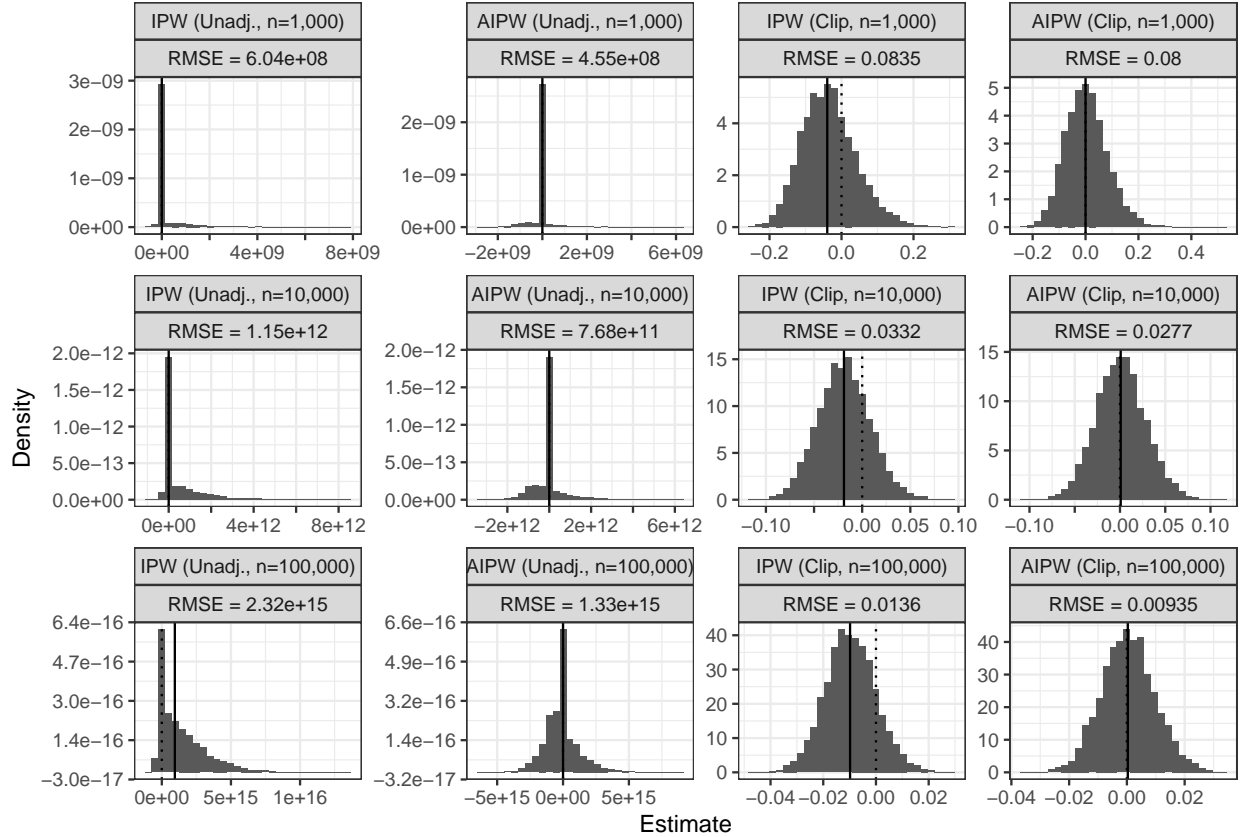


Figure 2: Histograms of point estimates in simulations for the various methods considered in the simulations. Vertical dotted and solid lines indicate true causal effect and median estimate, respectively. Clipped estimators achieve much better performance than unthresholded estimators, and clipped AIPW’s debiasing property is also apparent.

I begin by summarizing point estimates in Figure 2. The unthresholded estimators are approximately median-unbiased, but possess sufficiently heavy inverse propensity tails that the mean performance degrades with increasing sample size. The clipped estimators perform much better, but the clipped IPW estimator exhibits its known first-order bias. The clipped AIPW estimator exhibits less bias than the clipped IPW

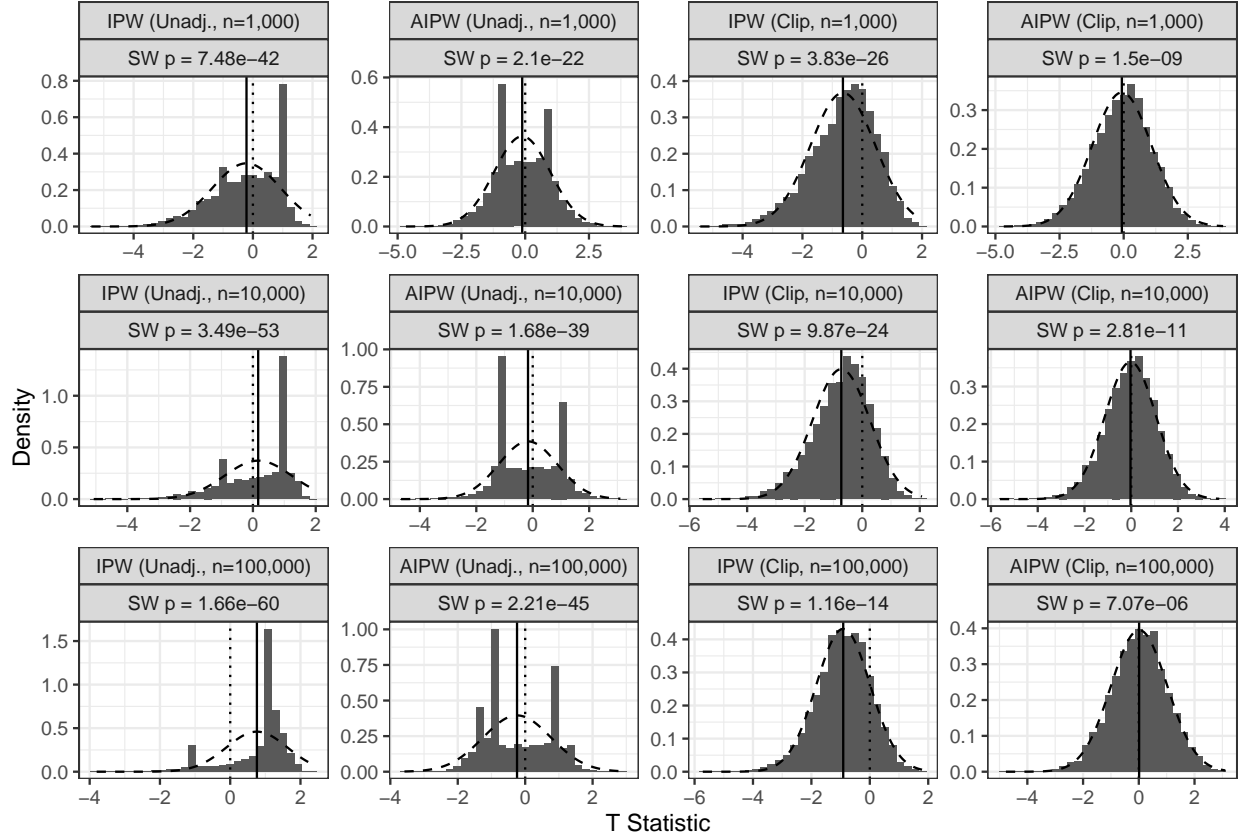


Figure 3: Histograms of simulated t-statistics on the true null hypothesis for various sample sizes. Vertical solid and dotted lines indicate mean t-statistic and target mean t-statistic of zero, respectively. Dashed line corresponds to the calibrated Gaussian density targeted in the Shapiro-Wilk test for normality.

estimator, and has slightly better performance in terms of mean squared error.

I find in Figure 3 that the clipped AIPW estimator’s t-statistics are reasonably well-calibrated. The plot presents t-statistics on the true average potential outcome. The t-statistics of unthresholded IPW and AIPW estimators are visibly non-Gaussian, and often exhibit a multimodal distribution. This poor performance is unsurprising: unthresholded estimators are known to fail to be asymptotically normal in this setting. Both thresholded estimators are known to be asymptotically normal in this setting when the propensity score is known, and both the asymptotic normality and the clipped IPW estimator’s first-order bias are visible to the naked eye, although the clipped IPW estimator also exhibits visible skew in small samples. I test for t-statistic normality using a Shapiro-Wilk test. The test rejects normality for both clipped estimates. Still, the clipped AIPW estimator’s violations are less severe by this criterion.

I find in Figure 4 that the clipped AIPW estimator’s p-values are well-calibrated in large samples. I use Wald confidence intervals to calculate two-sided p-values on the null of the true average potential outcome. If Wald confidence intervals are well-calibrated, then the simulated p-values on the true average potential

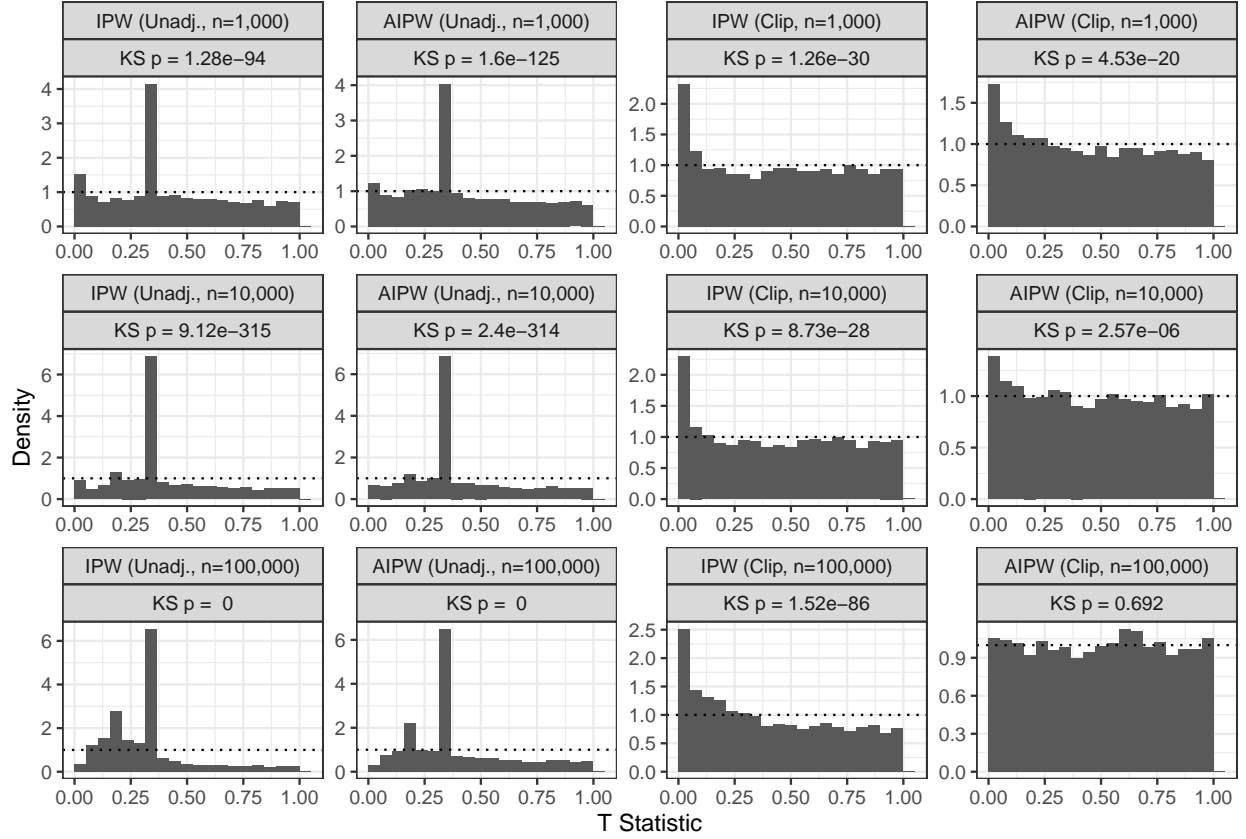


Figure 4: Histograms of simulation p-values on null hypothesis of true average potential outcome for various sample sizes. Dotted lines correspond to the target  $\text{Uniform}(0, 1)$  density. P-values in labels correspond to Kolmogorov-Smirnov tests for the  $\text{Uniform}(0, 1)$  distribution.

outcome will be exactly uniformly distributed. The unthresholded IPW and AIPW estimators exhibit known poor performance. The clipped IPW estimator exhibits overrejection even with large samples, as even oracle clipped IPW would provide well-calibrated inference for a biased estimand. The clipped AIPW estimator also overrejects in small samples, but the bias is less severe: with 1,000 observations, clipped IPW rejects the true null in 12.0% of simulations, while clipped AIPW rejects in 8.8% of simulations. As the sample size increases, the asymptotic calibration of Corollary 1 becomes apparent. With 100,000 observations, clipped IPW rejects the true null hypothesis in 12.8% of simulations, while clipped AIPW rejects in 5.3% of simulations. The Kolmogorov-Smirnov p-value on exact calibration of the two-sided test statistics for clipped AIPW with 100,000 observations is 0.692. This is a remarkable result: despite the known extreme difficulty of statistical inference in this setting, 5,000 simulated draws are insufficient to detect a meaningful failure of Wald confidence intervals based on the clipped AIPW estimator.

In moderate samples, clipped AIPW can undercover due to the difficulty of outcome regression in this setting. Figure 5 presents an example with 1,000 observations. It is rare to have treated observations with



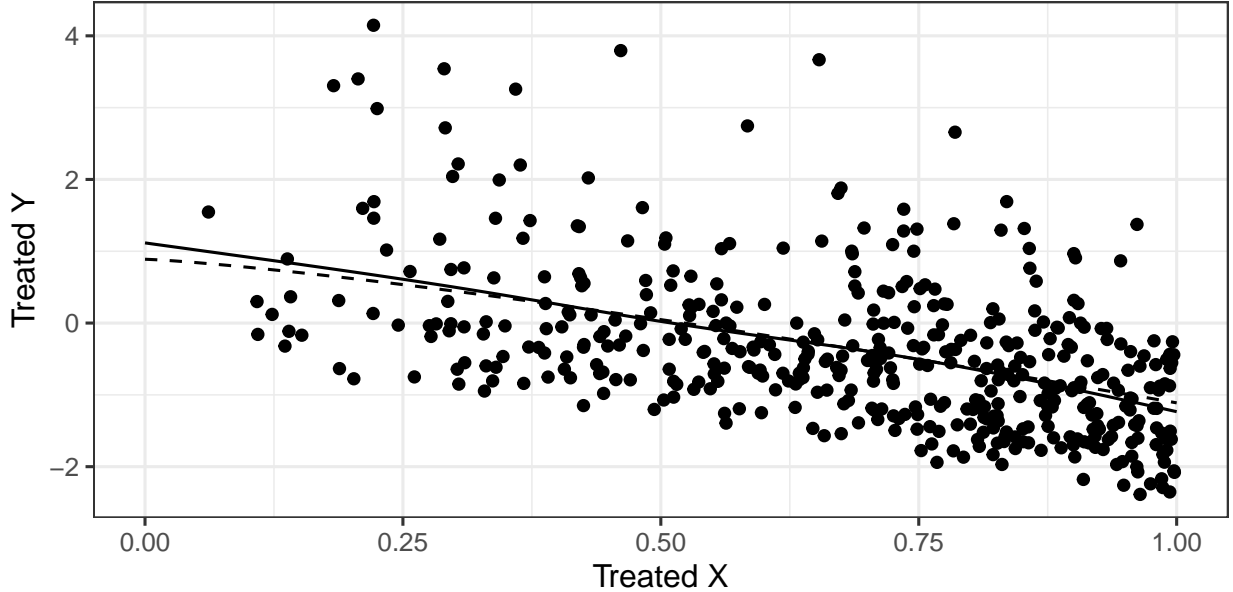


Figure 5: Simulated treated observations for one simulation of 1,000 observations. It is rare to see treated observations with small  $X$ , which corresponds to small values of  $e(X) = X^{1/(\gamma_0-1)}$ . As a result, such observations can have high leverage when predicting  $E[Y | X = 0, D = 1]$ , and can yield to important errors between the true (dashed) and predicted (solid) regression lines.

small values of  $e(X)$ . As a result, when such observations are treated, a small number of observations can receive substantial leverage in outcome regression, and the predictions of  $E[Y | X = 0, D = 1]$  can be driven by a small number of observations. In Appendix B (Figures 9 through 11), I conduct the same experiments, but with the estimated outcome regression function replaced by the true outcome regression function. The root-mean-squared error and failures of normality are comparable, suggesting these non-inferential patterns are driven by propensity estimation and clipping. However, the two-sided p-values exhibit better performance in small samples, and if anything slightly underreject with 100,000 observations.

In Appendix B (Figures 12 through 14), I show that these conclusions would largely carry through if clipping were replaced by trimming. The notable differences are that trimmed AIPW exhibits slightly better estimation performance in small samples, while if anything trimmed IPW is slightly worse; trimmed t-statistics exhibit less severe violations of normality; and p-values based on trimmed propensities exhibit more severe undercoverage for both IPW and AIPW.

## 5.2 Application to Right Heart Catheterization

I apply the clipped AIPW estimator to study the effect of right-heart catheterization (RHC) on survival. This dataset was first analyzed by Connors et al. (1996), and is a common benchmark in the weak overlap literature (Crump et al., 2009; Armstrong and Kolesár, 2017).

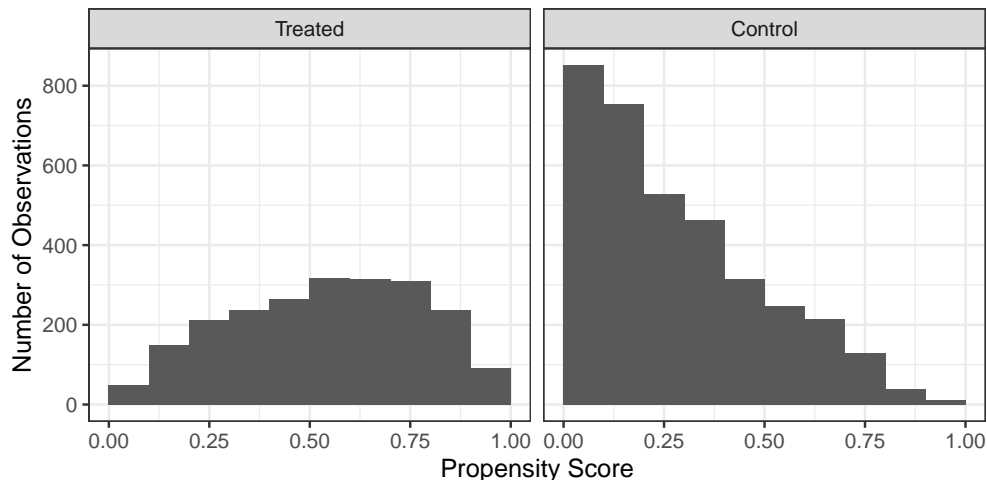


Figure 6: Histogram of estimated propensity scores for treated (left) and control (right) observations in right heart catheterization data. The plot is designed to parallel Figure 1 in [Crump et al. \(2009\)](#). Slight differences reflect the use of cross-fitting.

I analyze a version of the dataset from [Armstrong and Kolesár \(2017\)](#). The dataset is comprised of 5,735 adult patients, and the treatment  $D$  corresponds to receiving RHC within 24 hours of admission. The target causal effect is the average treatment effect of RHC on 30-day survival. The data includes 52 covariates  $X$  (72 covariates if counting factor levels separately). I estimate the nuisance functions  $e(X)$  and  $\mu(X)$  using five-fold cross-fitting. I estimate nuisance functions with logistic regression to align with [Crump et al. \(2009\)](#)’s empirical application. I estimate standard errors by bootstrapping the procedure. I keep fold assignment fixed in bootstraps to minimize the risk of over-fitting.

[Crump et al.](#) propose a weak overlap rule of thumb that estimates the treatment effect for the subpopulation with propensity scores between 10% and 90%. This rule-of-thumb trimming rule is chosen to approximately minimize asymptotic variance. This strategy ensures asymptotic normality, but changes the target estimand even asymptotically. By comparison, the clipped and trimmed AIPW estimators I analyze have thresholds  $b_n$  that tend to zero asymptotically. As a result, the estimators proposed here are able to target full population average treatment effect, potentially at the cost of increased variance. I compare these procedures to the 10% rule and other fixed trimming rules using the same nuisance estimates.

I present the distribution of estimated propensity scores for treated and control units in Figure 6. The figure is an analog of [Crump et al. \(2009\)](#)’s Figure 1. There is a meaningful density of units with estimated propensities near zero, suggesting weak overlap. This pattern is similar to the findings of [Crump et al.](#), although there are slight differences, presumably due to my use of cross-fitting.

I compare AIPW estimators for various trimmed subsamples to the clipped AIPW estimator. I choose the clipping threshold  $b_n$  through the no-specified-upper-bound version of Algorithm 1 because I estimate

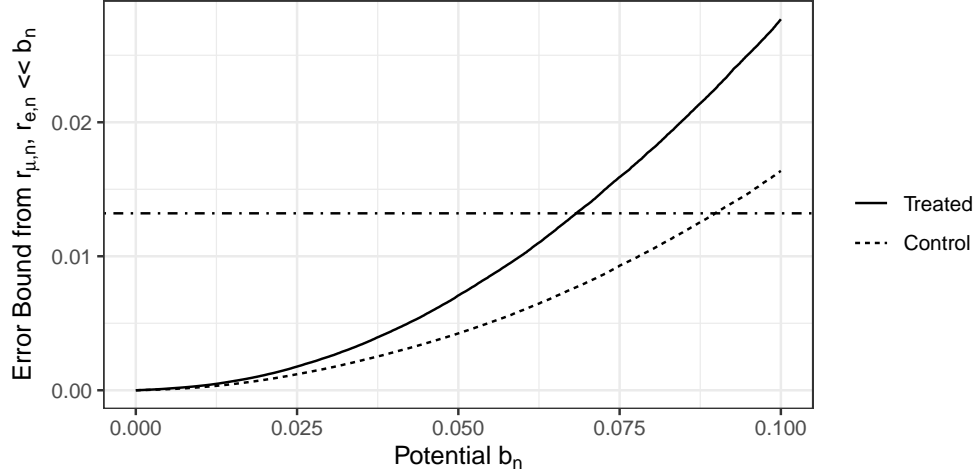


Figure 7: The value of `error_bound(b)` from Algorithm 1 for  $e(x)$  and  $1 - e(x)$  thresholding with no specified nuisance upper bounds, which corresponds to an error bound implied by  $r_{\mu,n}, r_{e,n} \ll b_n$ . The procedure chooses  $b_n$  to set this function equal to  $n^{-1/2}$ , which corresponds to the horizontal line.  $n^{-1/2}$  is indicated by horizontal dashed line. The more favorable distribution of estimated treatment propensities allows for a more aggressive clipping threshold.

both nuisance functions parametrically. I plot the functions used in choosing  $b_n$  in Figure 7. The estimated lower clipping threshold is 0.068 and affects 10.5% of observations. The Crump et al. 10% rule of thumb would exclude 16.3% of observations below. The estimated upper clipping threshold is 0.09 below one: there are few observations with large estimated propensities, so the rule of thumb concludes there is no need to trim observations with large estimated propensities. This upper threshold affects 1.4% of observations, comparable to the 1.8% of observations excluded above by the 10% rule of thumb.

I present estimated effects and confidence intervals for various potential fixed trimming rules in Figure 8. The 10% trimming rule yields an estimated reduction in survival rates of 5.79 percentage points among the trimmed sample, with an estimated 95% Wald confidence interval of  $[-9.14, -2.43]$ . Other trimming rules would yield larger confidence intervals, as expected because the 10% rule is chosen to roughly minimize asymptotic variance over target populations.

I compare the fixed-trimmed-sample AIPW estimates to a clipped AIPW estimator that targets the full population treatment effect. The estimated harm increases to -6.07 percentage points, a change of 0.168 standard errors under the 10% rule of thumb estimator. The clipped AIPW confidence interval of  $[-9.6, -2.55]$ , has a 5.14% larger width than the 10% trimmed sample interval. The clipped AIPW point estimates are similar to the point estimates under a 1% or 5% trimming rule, but the associated confidence interval is narrower under the full-population estimator. Part of the added width is driven by inverse propensities among clipped observations: if I used a trimmed, rather than clipped, AIPW estimator, the estimated effect would move by 0.256 standard errors, and the standard error would only increase by 0.54%. However, the

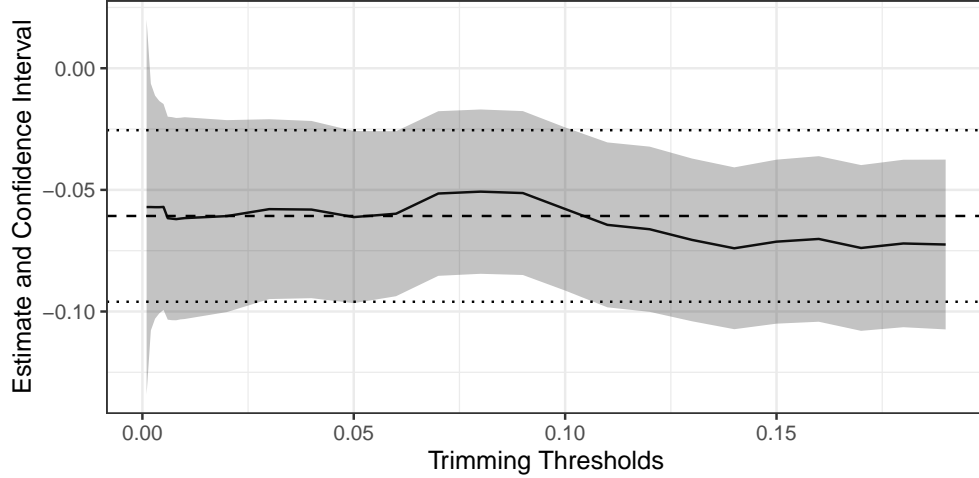


Figure 8: Estimated effects (solid line) and 95% confidence interval (shaded region) for AIPW applied to various trimmed subsamples. Estimate and confidence for clipped AIPW is represented by the dashed and dotted horizontal lines, respectively. Clipped AIPW produces similar estimates and standard errors as the clipping procedure while targeting a more interpretable estimand. A threshold of zero is omitted from the graph because the resulting confidence interval of  $[-568.8, 581.3]$  would make the graph difficult to read.

simulation results of Section 5.1 suggest that trimmed AIPW may slightly undercover.

Taken together, these results illustrate that under weak overlap, targeting the causal effect within the full population need not come at a large precision cost. In this application, clipped AIPW with a rule-of-thumb clipping rate yields similar estimates to estimators that target a fixed trimmed sample, while targeting a population that is often more relevant and adding only a small precision cost.

## 6 Conclusion

This work shows that standard Wald confidence intervals for clipped AIPW can achieve target coverage for standard causal effects under plausible conditions. I provide sufficient conditions on nuisance regression rates for clipped (or trimmed) AIPW to be uniformly valid over distributions with even very weak overlap. I use these theoretical results to derive new rules of thumb for choosing a threshold. I find that Wald confidence intervals perform well in simulations, especially in large samples, and can achieve comparable precision to a fixed 10% trimming rule in practice.

These results can be extended in many interesting directions. This work exploits Neyman orthogonality to achieve standard statistical inference in the presence of a small region of irregular identification. [Sasaki and Ura \(2022\)](#) and [Ma et al. \(2023\)](#) propose estimators for ratio estimands beyond IPW; the arguments here are likely to extend to their more general framework. Issues of weak overlap hold for inverse propensity and other importance sampling estimators in settings like difference-in-difference estimation ([Callaway and Sant’Anna,](#)

2021) or statistical inference for parameters that are identified at infinity (Andrews and Schafgans, 1998; Khan and Nekipelov, 2024); the results and rules of thumb here can likely be adapted to those settings. Semenova (2024) applies thresholding strategies to intersection bounds, where at a high level a margin condition plays the role of the minimal overlap bound here. Perhaps similar ideas could apply to other forms of irregular identification. The regression analysis in Section 3.3 may also extend to estimating the effects of continuous treatments.

The results here suggest that thresholded AIPW is a viable alternative to fixed-trimming rules. I provide rules of thumb that enable practitioners to easily report results that target the population average effect. When, as in my empirical application, the fixed-trimming and sequence-of-thresholds approaches yield similar causal conclusions, then there is strong evidence that causal conclusions are driven by causal effects, and not how the researcher treats observations with extreme propensity scores.

## References

- Andrews, D. W. K. and Schafgans, M. M. A. (1998). Semiparametric estimation of the intercept of a sample selection model. *The Review of Economic Studies*, 65(3):497–517.
- Armstrong, T. B. and Kolesár, M. (2017). A simple adjustment for bandwidth snooping. *The Review of Economic Studies*, 85(2):732–765.
- Armstrong, T. B. and Kolesár, M. (2021). Finite-sample optimal estimation and inference on average treatment effects under unconfoundedness. *Econometrica*, 89(3):1141–1177.
- Bailey, M. J. and Goodman-Bacon, A. (2015). The war on poverty’s experiment in public medicine: Community health centers and the mortality of older americans. *American Economic Review*, 105(3):1067–1104.
- Bruns-Smith, D., Dukes, O., Feller, A., and Ogburn, E. L. (2024). Augmented balancing weights as linear regression.
- Callaway, B. and Sant’Anna, P. H. (2021). Difference-in-differences with multiple time periods. *Journal of Econometrics*, 225(2):200–230. Themed Issue: Treatment Effect 1.
- Chaudhuri, S. and Hill, J. B. (2024). Heavy tail robust estimation and inference for average treatment effects. *Econometric Reviews*.
- Chen, X., Hong, H., and Tarozi, A. (2008). Semiparametric efficiency in gmm models of nonclassical measurement errors, missing data and treatment effects. Technical Report 1644, Cowles Foundation.

- Chernozhukov, V., Chetverikov, D., Demirer, M., Duflo, E., Hansen, C., Newey, W., and Robins, J. (2018). Double/debiased machine learning for treatment and structural parameters. *Econometrics Journal*, 21(1):C1–C68.
- Connors, Alfred F., J., Speroff, T., Dawson, N. V., Thomas, C., Harrell, Frank E., J., Wagner, D., Desbiens, N., Goldman, L., Wu, A. W., Califf, R. M., Fulkerson, William J., J., Vidaillet, H., Broste, S., Bellamy, P., Lynn, J., and Knaus, W. A. (1996). The effectiveness of right heart catheterization in the initial care of critically ill patients. *JAMA*, 276(11):889–897.
- Crump, R. K., Hotz, V. J., Imbens, G. W., and Mitnik, O. A. (2009). Dealing with limited overlap in estimation of average treatment effects. *Biometrika*, 96(1):187–199.
- Currie, J. and Walker, R. (2011). Traffic congestion and infant health: Evidence from E-ZPass. *American Economic Journal: Applied Economics*, 3(1):65–90.
- D’Amour, A., Ding, P., Feller, A., Lei, L., and Sekhon, J. (2021). Overlap in observational studies with high-dimensional covariates. *Journal of Econometrics*, 221(2):644–654.
- Galiani, S., Gertler, P., and Schargrodsky, E. (2005). Water for life: The impact of the privatization of water services on child mortality. *Journal of Political Economy*, 113(1):83–120.
- Gaïffas, S. (2005). Convergence rates for pointwise curve estimation with a degenerate design. *Mathematical Methods of Statistics*, 14(1).
- Goldsmith-Pinkham, P., Hull, P., and Kolesár, M. (2024). Contamination bias in linear regressions. *American Economic Review*, 114(12):4015–51.
- Hahn, J. (1998). On the role of the propensity score in efficient semiparametric estimation of average treatment effects. *Econometrica*, 66(2):315–331.
- Hall, P., Marron, J. S., Neumann, M. H., and Titterton, D. M. (1997). Curve estimation when the design density is low. *The Annals of Statistics*, 25(2):756 – 770.
- Heiler, P. and Kazak, E. (2021). Valid inference for treatment effect parameters under irregular identification and many extreme propensity scores. *Journal of Econometrics*, 222(2):1083–1108.
- Horn, R. A. and Johnson, C. R. (2013). *Matrix Analysis: Second Edition*. Cambridge University Press.
- Imbens, G. W. (2004). Nonparametric estimation of average treatment effects under exogeneity: A review. *Review of Economics and Statistics*, 86(1):4–29.

- Ionides, E. L. (2008). Truncated importance sampling. *Journal of Computational and Graphical Statistics*, 17(2):295–311.
- Khan, S. and Nekipelov, D. (2024). On uniform inference in nonlinear models with endogeneity. *Journal of Econometrics*, 240(2):105261.
- Khan, S. and Tamer, E. (2010). Irregular identification, support conditions, and inverse weight estimation. *Econometrica*, 78(6):2021–2042.
- Khan, S. and Ugander, J. (2022). Doubly-robust and heteroscedasticity-aware sample trimming for causal inference.
- Lee, B. K., Lessler, J., and Stuart, E. A. (2011). Weight trimming and propensity score weighting. *PLoS ONE*, 6(3).
- Lei, L., D’Amour, A., Ding, P., Feller, A., and Sekhon, J. (2021). Distribution-free assessment of population overlap in observational studies.
- Li, F., Morgan, K. L., and Zaslavsky, A. M. (2018). Balancing covariates via propensity score weighting. *Journal of the American Statistical Association*, 113(521):390–400.
- Ma, X., Sasaki, Y., and Wang, Y. (2024). Testing limited overlap. *Econometric Theory*.
- Ma, X. and Wang, J. (2020). Robust inference using inverse probability weighting. *Journal of the American Statistical Association*, 115(532):1851–1860.
- Ma, Y., Sant’Anna, P. H. C., Sasaki, Y., and Ura, T. (2023). Doubly robust estimators with weak overlap.
- Mou, W., Ding, P., Wainwright, M. J., and Bartlett, P. L. (2023). Kernel-based off-policy estimation without overlap: Instance optimality beyond semiparametric efficiency.
- Newey, W. K. (1994). The asymptotic variance of semiparametric estimators. *Econometrica*, 62(6):1349–1382.
- Rothe, C. (2017). Robust confidence intervals for average treatment effects under limited overlap. *Econometrica*, 85(2):645–660.
- Sasaki, Y. and Ura, T. (2022). Estimation and inference for moments of ratios with robustness against large trimming bias. *Econometric Theory*, 38(1):66–112.
- Semenova, V. (2024). Aggregated intersection bounds and aggregated minimax values.

- Stone, C. J. (1982). Optimal global rates of convergence for nonparametric estimators. *The Annals of Statistics*, 10(4):1040–1053.
- Tsybakov, A. B. (2009). *Introduction to Nonparametric Estimation*. Springer.
- Yang, S. and Ding, P. (2018). Asymptotic inference of causal effects with observational studies trimmed by the estimated propensity scores. *Biometrika*, 105(2):487–493.



## A Key Technical Assumptions and Claims

I make use of the following assumptions.

**Assumption 3’.** Assumption 2 holds, with the following rates on the regression error  $r_{\mu,n}$  and the propensity error  $r_{e,n}$  for any sequence of  $P(n) \in \mathcal{P}$ :

- (a) *Consistency.*  $r_{\mu,n}, r_{e,n} \rightarrow 0$ .
- (b) *Product of errors.*  $r_{\mu,n} r_{e,n} \sqrt{\mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e(X), b_n\}^2} \right]} \ll n^{-1/2}$ .
- (c) *Regression error near singularities.*  $r_{\mu,n} \frac{P(n)(e(X) \leq b_n)}{\sqrt{\mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e, b_n\}^2} \right]}} \ll n^{-1/2}$ .
- (d) *Asymptotically known thresholding.*  $r_{e,n} \ll b_n$ .

These conditions adapt to the distributions in the sequence  $P(n)$ .

The conditions of Assumption 3’ are weaker than the more interpretable conditions in the main text.

**Corollary 6** (Sufficiency of Assumption 3’). *Suppose Assumptions 3, 4, and Assumption 4(i) hold and let  $\rho > 0$  be given. Then Assumption 3’ holds.*

I will show that the feasible clipped estimator  $\hat{\psi}_{clip}^{AIPW}(b_n)$  is first-order equivalent to the oracle clipped estimator  $\tilde{\psi}_{(Oracle)}^{AIPW}(b_n)$ . The oracle clipped AIPW estimator is asymptotically normal by the trimmed IPW arguments in Ma and Wang (2020). By construction, the oracle clipped AIPW estimator is finite-sample unbiased. The following asymptotic normality follows as a result.

**Proposition 4** (Oracle asymptotic normality). *Suppose  $b_n$  satisfies  $n^{-1/2} \ll b_n \ll 1$ . Then the oracle clipped AIPW estimator has uniform convergence to a normal distribution in the sense that*

$$\limsup_{n \rightarrow \infty} \sup_{P \in \mathcal{P}} \sup_{t \in \mathbb{R}} \left| P \left( \frac{\tilde{\psi}_{(Oracle)}^{AIPW}(b_n) - \psi(P)}{\sigma_n} \leq t \right) - \Phi(t) \right| = 0.$$

Proposition 4 will be an extension of the following claim. In addition to this modified theorem, Theorem 1 replaces the oracle standard deviation  $\sigma_n$  with the estimated standard deviation  $\hat{\sigma}_n$  when constructing t-statistics.

**Theorem 1’** ((Slow) Asymptotic Normality). *Suppose the conditions of Theorem 1 hold, and  $P(n)$  is a sequence of distributions  $P \in \mathcal{P}$ . Then  $\sigma_n^{-1} \left( \hat{\psi}_{clip}^{AIPW}(b_n) - \psi_n \right) \overset{P(n)}{\rightsquigarrow} N(0, 1)$ , where  $\sigma_n$  is the oracle standard deviation defined in Proposition 4.*

Next, I describe the key new results for nonparametric regression. This result shows that in nonparametric regression, if the propensity function is sufficiently smooth, then nature cannot severely concentrate treated

observations within a given bandwidth of any point. The non-concentration ensures that the eigenvalues of the local polynomial regression matrix are nondegenerate.

**Proposition 5** (Non-trivial concentration). *Suppose Assumption 6 holds and there is an  $L > 0, \beta_e > \frac{d}{\gamma_0 - 1}$  such that  $P(D = 1 \mid X) \in \Sigma(\beta_e, L)$  for all  $P \in \mathcal{P}$ . Then Assumption 7 holds.*

Note that  $\frac{d}{\gamma_0 - 1}$  is also a key parameter in Mou et al. (2023). A broader connection is outside the scope of this work.

The next result presents the main construction involved in avoiding a polylogarithmic penalty under weak overlap. The claim iteratively constructs a sequence of minimal bandwidths and gridpoint counts that ensure that no more than  $m_n^{(k)}$  gridpoints  $x_{n,j}$  can have  $nE[D\mathbf{1}\{\|X - x_{n,j}\| \leq h\}] = h^{-2\beta_\mu}$  solved by  $h \leq h_n^{(k)}$ , and that each  $k$  can contribute at most half as much as the previous iteration to the polylogarithmic penalty. Here, I show that if the first set is expansive enough to incur what I later show is a polylogarithmic penalty bounded by  $\log(\delta)$ , then if  $\delta$  is at least some large, but finite, number, then this iterative process eventually covers all gridpoints.

**Proposition 6** (Inductive grouping). *Let  $\beta_\mu > 0, \gamma_0 > 1, d \geq 1$  be given. For any  $\delta > 1$ , inductively construct a sequence of  $h$  and  $m$  as follows. Take  $h_n^{(1)} = n^{\frac{-1}{2\beta_\mu + d\frac{\gamma_0}{\gamma_0 - 1}}}$ . Then for all  $k = 1, 2, \dots$ , inductively define  $m_n^{(k)} = \exp\left(2^{-k}\delta\left(\frac{h_n^{(k)}}{h_n^{(1)}}\right)^{-2\beta_\mu}\right)$  and  $h_n^{(k+1)} = \left(n\left(\frac{m_n^{(k)}}{\delta}\right)^{\gamma_0 - 1}\right)^{\frac{-1}{2\beta_\mu + d\frac{\gamma_0}{\gamma_0 - 1}}}$ . Then there is a  $\underline{\delta} > 0$  and a  $\pi \in (0, 1/4^{2\beta_\mu}]$  such that if  $\delta \geq \underline{\delta}$ , then  $\sum_{k=1}^{\infty} m_n^{(k)} = \infty$ , and for all  $k > 1, h_n^{(k)} \leq h_n^{(1)}\pi^{k-1}$ .*

## B Other Simulation Evidence

In this section, I presented simulated evidence for trimmed estimators.

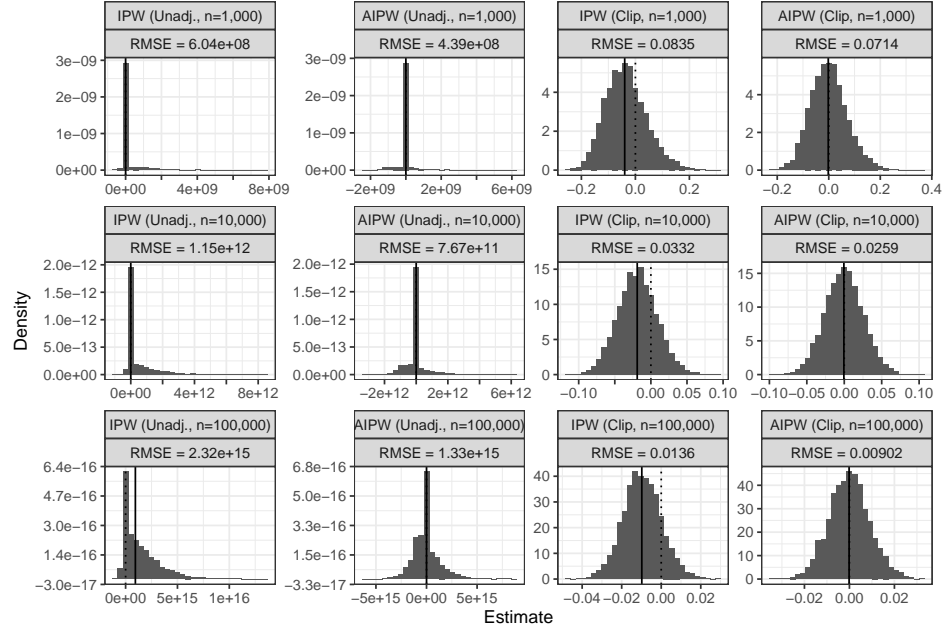


Figure 9: Histograms of point estimates in simulations for the various methods considered in the simulations, but using the oracle  $\mu$  regression function instead of the estimated  $\hat{\mu}$  regression function.

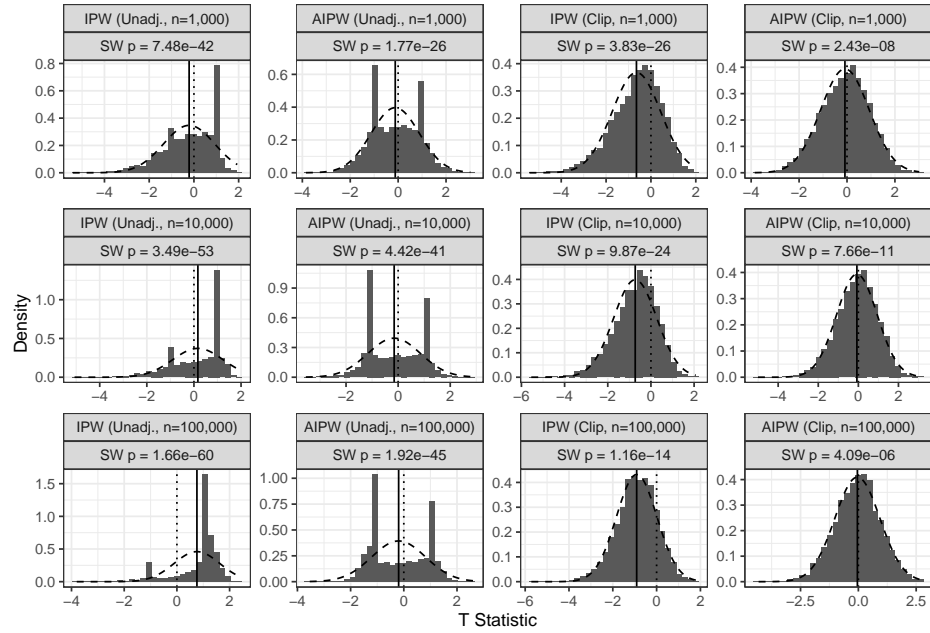


Figure 10: Histograms of simulation t-statistics for various sample sizes, but using the oracle  $\mu$  regression function instead of the estimated  $\hat{\mu}$  regression function.

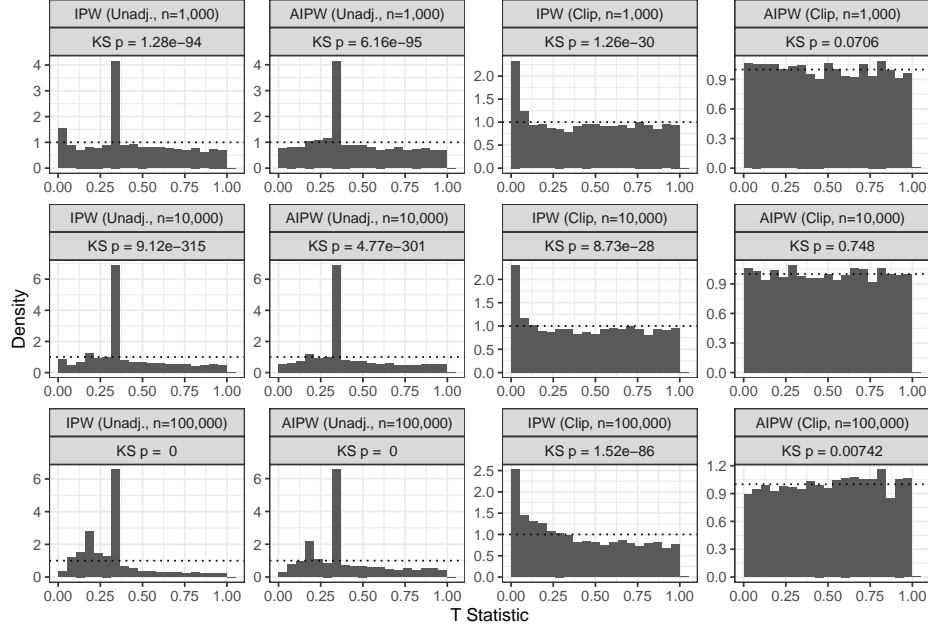


Figure 11: Histograms of simulation p-values on null hypothesis of true average potential outcome for various sample sizes, but using the oracle  $\mu$  regression function instead of the estimated  $\hat{\mu}$  regression function.

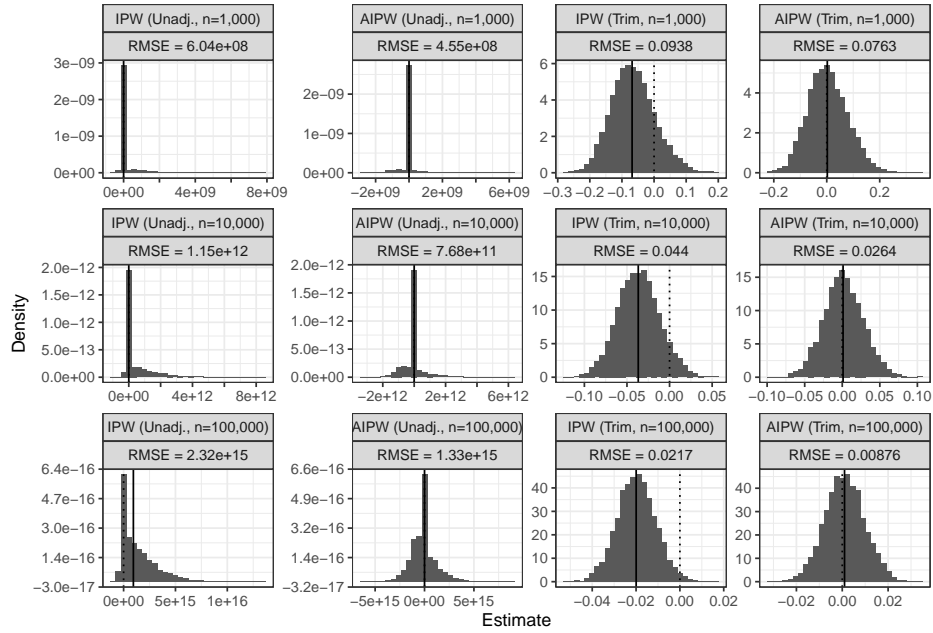


Figure 12: Histograms of point estimates in simulations for the various methods considered in the simulations, but with trimming instead of clipping.

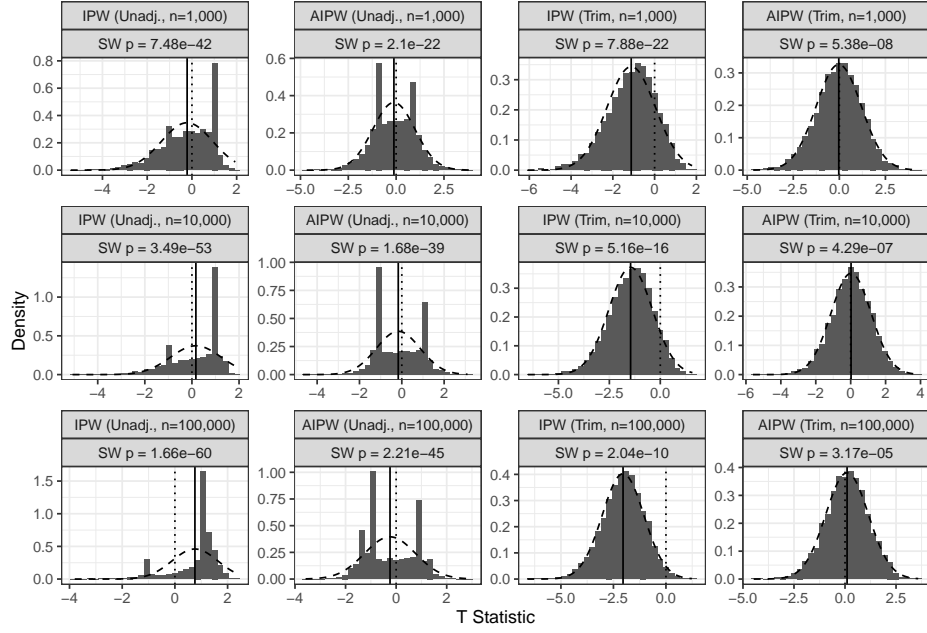


Figure 13: Histograms of simulation t-statistics for various sample sizes, but with trimming instead of clipping.

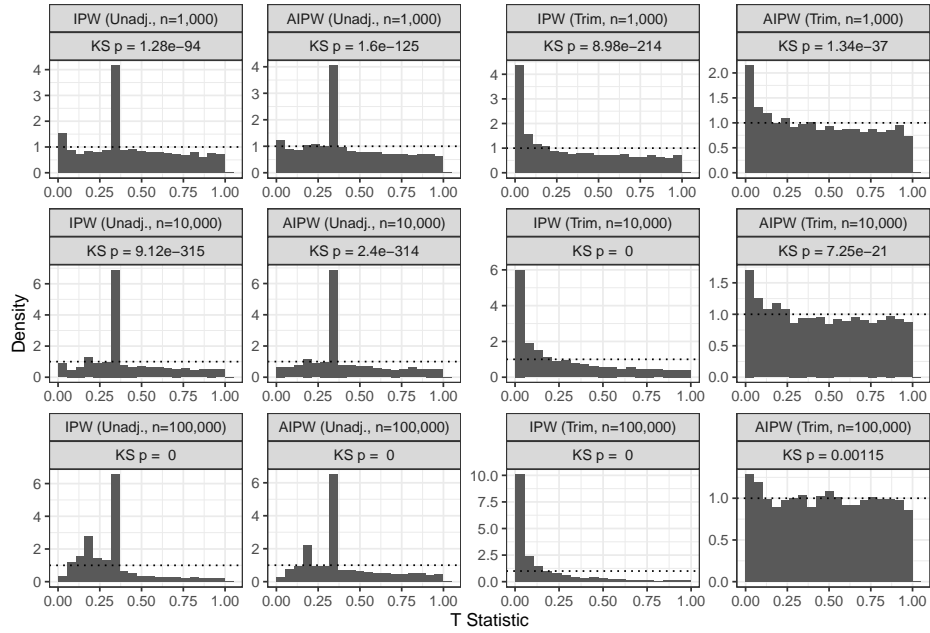


Figure 14: Histograms of simulation p-values on null hypothesis of true average potential outcome for various sample sizes, but with trimming instead of clipping.

## C Proofs

The proofs, including proofs of the claims in Appendix A, are split into sections showing asymptotic properties of oracle clipped AIPW (Appendix C.1), consistency of estimated clipped AIPW (Appendix C.2), and black-box consistency rates (Appendix C.3). Note that Appendix C.3 is out of order from the perspective of the main text and corresponds to claims in Section 3.1, but Proposition 3 is used to show claims that appear earlier in the main text. Appendix C.4 presents proofs for the main claims: the black-box asymptotic properties of estimated clipped AIPW. Finally, Appendix C.5, Appendix C.6, and Appendix C.7 prove claims about black-box nuisance rates, outcome regression rates, and rules of thumb, respectively.

**Additional Notation.** In these proofs, I use  $P(n)$  to refer to an arbitrary sequence of distributions for the purposes of computing suprema; for such sequences, I use  $\psi_n = \psi(P(n))$  to denote the sequence of average potential outcomes. I use  $P_n[c_n]$  to refer to the average of  $c_n$  over  $n$  draws from  $P$  (sometimes abusing notation and including nuisance functions in  $c_n$ ), and I use  $P[c_n]$  to refer to the expectation of  $c_n$  over  $P$ . This can occasionally lead to unfortunate notation like  $P(n)_n(E_n)$  for a sequence of event probabilities under a sequence of distributions. I write  $\lim_{x \rightarrow z+} f(x)$  and  $\lim_{x \rightarrow z-}$  for the right- and left-hand limits of  $f$  at  $z$ . I write  $c_n = o_{P(n)}(1)$  if for all  $\delta > 0$ ,  $P(n)(|c_n|/d_n > \delta) \rightarrow 0$ , and if no  $P(n)$  is defined, I use  $c_n = o_{P(n)}(d_n)$  to mean that for any sequence of  $P(n) \subset \mathcal{P}$ ,  $c_n = o_{P(n)}(d_n)$ . I write  $c_n = O_{P(n)}(1)$  if for all  $\epsilon > 0$ , there exists a  $\delta > 0$  such that  $P(n)(|c_n|/d_n > \delta) < \epsilon$ . If there is a sequence of distributions to be considered, then I use  $o(d_n)$  and  $O(d_n)$  to implicitly refer to  $o_{P(n)}(d_n)$  and  $O_{P(n)}(d_n)$ . I write that  $c_n \overset{P(n)}{\rightsquigarrow} N(0, 1)$  if  $\sup_{t \in \mathbb{R}} |P(n)(c_n \leq t) - \Phi(t)| \rightarrow 0$ , where  $\Phi$  is the standard normal cumulative distribution function; I write that  $c_n \rightarrow_{P(n)} c$  if  $c_n - c = o_{P(n)}(1)$ ; and I write that  $c_n \xrightarrow{\mathcal{P}} c$  if for all sequences of  $P(n) \in \mathcal{P}$ ,  $c_n \rightarrow_{P(n)} c$ . I write  $c_n = \Theta(d_n)$  if there exists a  $k_1, k_2 > 0$  such that  $P(n)[c_n \in [k_1 d_n, k_2 d_n]] \rightarrow 1$ , and I write  $c_n = \Omega(d_n)$  if there exists a  $k_1 > 0$  such that  $P(n)[c_n \geq k_1 d_n] \rightarrow 0$ .

### C.1 Oracle Normality

*Proof of Proposition 1.* For (i), it is known that the efficiency bound is  $n^{-1} \left( \text{Var}(\mu(X)) + E \left[ \frac{\text{Var}(Y|X, D=1)}{e(X)} \right] \right)$  (Hahn, 1998). By Assumptions 1(a) and 1(b),  $\text{Var}(\mu(X))$  and  $\text{Var}(Y | X, D = 1)$  are uniformly bounded above. It therefore only remains to show the standard condition (Newey, 1994; Hahn, 1998; Chen et al., 2008) that  $\sup_{P \in \mathcal{P}} E_P[1/e(X)]$  is bounded above:

$$\begin{aligned} E_P[1/e(X)] &= \int_0^\infty P\left(\frac{1}{e(X)} > t\right) dt \\ &= \int_0^\infty P(e(X) < 1/t) dt \end{aligned}$$

$$\begin{aligned}
&= 1 + \int_1^\infty P(e(X) < 1/t) dt \\
&\leq 1 + \int_1^\infty C(t^{-1})^{\gamma_0-1} dt \quad (\text{Lemma 2}) \\
&= 1 + C \int_1^\infty t^{1-\gamma_0} dt \\
&= 1 + \frac{C}{\gamma_0 - 2},
\end{aligned}$$

which is finite.

(ii) holds by extension of [Khan and Tamer \(2010\)](#) Theorem 4.1. In particular, note that by Assumption 1(c),  $\text{Var}(Y \mid X, D = 1) \geq \sigma_{\min}^2 > 0$ , so it only remains to show that  $E_P[1/e(X)]$  is infinite:

$$\begin{aligned}
E_P[1/e(X)] &= 1 + \int_1^\infty P(e(X) < 1/t) dt \\
&\geq 1 + C' \int_1^\infty t^{1-\gamma_0} dt = \infty.
\end{aligned}$$

□

**Lemma 2.** Define  $\pi_{\min} = \frac{C}{\gamma_0}$ . Then  $\inf_{P \in \mathcal{P}} P(D = 1) \geq \pi_{\min} > 0$ .

*Proof of Lemma 2.* For any  $P \in \mathcal{P}$ ,

$$P[D = 1] = E_P[e(X)] = \int_0^\infty P(e(X) > t) dt = \int_0^1 P(e(X) \leq t) dt \geq \int_0^1 C t^{\gamma_0-1} dt = \frac{C}{\gamma_0}.$$

□

**Lemma 3.** Assume  $b_n \rightarrow 0$ . Then for all large  $n$ , the following inequalities hold throughout  $\mathcal{P}$ :

- (i)  $P(e(X) > \pi_{\min}/2) \geq \pi_{\min}/2$
- (ii)  $\mathbb{E}[e(X)/\{e(X) \vee b_n\}^2] \geq \pi_{\min}/2$
- (iii)  $\mathbb{E}[|\phi_n - \mathbb{E}_{P(n)}[\phi_n]|^q] \leq (4M)^q \mathbb{E}[e(X)/\{e(X) \vee b_n\}^2]/b_n^{q-2}$
- (iv)  $\mathbb{E}[|\phi_n|^q] \leq (8M)^q \mathbb{E}[e(X)/\{e(X) \vee b_n\}^2]/b_n^{q-2}$

*Proof of Lemma 3.* I take these proofs one at a time.

(i) Start from the following inequalities:

$$\begin{aligned}
\pi_{\min} &\leq \mathbb{E}[e(X)] \\
&= \mathbb{E}[e(X)\mathbf{1}\{e(X) \leq \pi_{\min}/2\}] + \mathbb{E}[e(X)\mathbf{1}\{e(X) > \pi_{\min}/2\}] \\
&\leq (\pi_{\min}/2)[1 - P(e(X) > \pi_{\min}/2)] + P(e(X) > \pi_{\min}/2)
\end{aligned}$$

$$< \pi_{\min}/2 + P(e(X) > \pi_{\min}/2)$$

Subtracting  $\pi_{\min}/2$  from the far left- and right-hand sides of this inequality gives the desired conclusion.

(ii) If  $b_n \leq \pi_{\min}/2$  (which happens for all large  $n$ ), then:

$$\mathbb{E}[e(X)/\{e(X) \vee b_n\}^2] \geq \mathbb{E}[1/e(X)\mathbf{1}\{e(X) \geq b_n\}] \geq P(e(X) \geq b_n) \geq P(e(X) \geq \pi_{\min}/2) \geq \pi_{\min}/2.$$

(iii) By Jensen's inequality:

$$\begin{aligned} \mathbb{E}[|\phi_n - \mathbb{E}_{P(n)}[\phi_n]|^q] &\leq 2^{q-1}(\mathbb{E}[|\mu(X) - \mathbb{E}_{P(n)}[\mu(x)]|^q] + \mathbb{E}[|Y - \mu(X)|^q D / \{e(X) \vee b_n\}^q]) \\ &\leq 2^{q-1}(2^q \mathbb{E}[|\mu(X)|^q] + \mathbb{E}[\mathbb{E}[|Y - \mu(X)|^q \mid X, D = 1] e(X) / \{e(X) \vee b_n\}^q]) \\ &\leq 2^{q-1}(2^q M^q + 2^q \mathbb{E}[\mathbb{E}[|Y|^q \mid X, D = 1] e(X) / \{e(X) \vee b_n\}^q]) \\ &\leq 2^{q-1}(2^q M^q + 2^q M^q \mathbb{E}[e(X) / \{e(X) \vee b_n\}^2] \times 1 / \{e(X) \vee b_n\}^{q-2}) \\ &\leq 2^{q-1}(2^q M^q + 2^q M^q \mathbb{E}[e(X) / \{e(X) \vee b_n\}^2] / b_n^{q-2}). \end{aligned}$$

Since  $\mathbb{E}[e(X) / \{e(X) \vee b_n\}^2] / b_n^{q-2} \geq \pi_{\min} / 2 b_n^{q-2} \rightarrow \infty$  by Item (ii), I may further bound the above quantity by  $2^{2q} M^q \mathbb{E}[e(X) / \{e(X) \vee b_n\}^2] / b_n^{q-2}$  once  $n$  is large enough.

(iv) By Jensen's inequality:

$$\begin{aligned} \mathbb{E}[|\phi_n|^q] &= \mathbb{E}[|\phi_n - \mathbb{E}_{P(n)}[\phi_n] + \mathbb{E}_{P(n)}[\phi_n]|^q] \\ &\leq 2^{q-1}(\mathbb{E}[|\phi_n - \mathbb{E}_{P(n)}[\phi_n]|^q] + |\mathbb{E}_{P(n)}[\phi_n]|^q) \\ &\leq 2^{q-1}(4M)^q \mathbb{E}[e(X) / \{e(X) \vee b_n\}^2] / b_n^{q-2} + 2^{q-1} |\mathbb{E}_{P(n)}[\mu(x)]|^q \quad (\text{Item (iii)}) \\ &\leq 2^{q-1}(4M)^q \mathbb{E}[e(X) / \{e(X) \vee b_n\}^2] / b_n^{q-2} + 2^{q-1} \mathbb{E}[\mathbb{E}[|Y|^q \mid X, D = 1]] \quad (\text{Jensen}) \\ &\leq 2^{q-1}(4M)^q \mathbb{E}[e(X) / \{e(X) \vee b_n\}^2] / b_n^{q-2} + 2^{q-1} M^q. \end{aligned}$$

As before, since  $\mathbb{E}[e(X) / \{e(X) \vee b_n\}^2] / b_n^{q-2} \rightarrow \infty$ , the first term in the upper bound is eventually larger than the second and I may bound the whole expression by  $(8M)^q \mathbb{E}[e(X) / \{e(X) \vee b_n\}^2] / b_n^{q-2}$  once  $n$  is large enough.

□

**Lemma 4.** Let  $c(\gamma) = \frac{\gamma-1}{\gamma} C^{-1/(\gamma-1)} > 0$ . Then for any  $P \in \mathcal{P}$ , I have:

$$\mathbb{E}_P[e(X)\mathbf{1}\{e(X) \leq b_n\}] \geq c(\gamma)P(e(X) \leq b_n)^{\gamma/(\gamma-1)}. \quad (3)$$



This lower bound is attained when  $P(e(X) \leq t) = t^{\gamma-1}$ .

*Proof of Lemma 4.* Let  $p = P(e(X) \leq b_n)$ . If  $p = 0$ , then the bound holds trivially so I will assume throughout that  $p > 0$ . Then I may write:

$$\begin{aligned}\mathbb{E}_P[e(X)\mathbf{1}\{e(X) \leq b_n\}] &= \int_0^\infty P(e(X)\mathbf{1}\{e(X) \leq b_n\} > t)dt = \int_0^{b_n} P(t < e(X) \leq b_n)dt = \int_0^{b_n} p - P(e(X) \leq t)dt \\ &\geq b_np - \int_0^{b_n} \min\{p, Ct^{\gamma-1}\}dt = b_np - \left(\frac{C}{\gamma}\right) \left(\frac{p}{C}\right)^{\frac{\gamma}{\gamma-1}} - b_np + p \left(\frac{p}{C}\right)^{\frac{1}{\gamma-1}} = c(\gamma)p^{\gamma/(\gamma-1)}.\end{aligned}$$

This proves the lower bound. When  $P(e(X) \leq t) = t^{\gamma-1}$ , a direct calculation gives  $\mathbb{E}_P[e(X)\mathbf{1}\{e(X) \leq b_n\}] = [(\gamma-1)/\gamma]b_n^\gamma = [(\gamma-1)/\gamma]P(e(X) \leq b_n)^{\gamma/(\gamma-1)}$ . Therefore, the lower bound is also sharp.  $\square$

**Lemma 5.** For any  $P \in \mathcal{P}$ ,

$$\begin{aligned}\text{Var}_P(\phi(Z \mid b_n, \eta)) &\geq \sigma_{\min}^2 \mathbb{E}_P [e(X)/\max\{e(X), b_n\}^2] \\ &\geq \sigma_{\min}^2 [c(\gamma)P(e(X) \leq b_n)^{\gamma/(\gamma-1)}/b_n^2 + \pi_{\min}/2] \\ &\geq \sigma_{\min}^2 \pi_{\min}/2 > 0.\end{aligned}$$

*Proof of Lemma 5.* For the first line:

$$\begin{aligned}\text{Var}_P(\phi_n) &= \mathbb{E}[\text{Var}(\phi_n \mid X)] + \text{Var}(\mathbb{E}[\phi_n \mid X]) = \mathbb{E}[\text{Var}(\phi_n \mid X)] \\ &= \mathbb{E} \left[ \mathbb{E}[|Y - \mu(X)|^2 \mid X, D = 1] \frac{e(X)}{\{e(X) \vee b_n\}^2} \right] \geq \sigma_{\min}^2 \mathbb{E} \left[ \frac{e(X)}{\{e(X) \vee b_n\}^2} \right].\end{aligned}$$

Since Definition 1 implies  $e(X) > 0$ ,  $\text{Var}_P(\phi_n) > 0$ .

For the second line, I assume  $n$  is so large that  $b_n \leq \pi_{\min}/2$ . Then:

$$\begin{aligned}\mathbb{E}[e(X)/\max\{e(X), b_n\}^2] &= \mathbb{E}[e(X)/b_n^2 \mathbf{1}\{e(X) \leq b_n\}] + \mathbb{E}[1/e(X) \mathbf{1}\{e(X) > b_n\}] \\ &\geq \mathbb{E}[e(X)/b_n^2 \mathbf{1}\{e(X) \leq b_n\}] + P(e(X) > b_n) \\ &\geq \mathbb{E}[e(X)/b_n^2 \mathbf{1}\{e(X) \leq b_n\}] + P(e(X) > \pi_{\min}) \\ &\geq \mathbb{E}[e(X)/b_n^2 \mathbf{1}\{e(X) \leq b_n\}] + \pi_{\min}/2 && \text{(Lemma 3.(i))} \\ &\geq c(\gamma)(1/b_n)^2 P(e(X) \leq b_n)^{\gamma/(\gamma-1)} + \pi_{\min}/2. && \text{(Lemma 4)}\end{aligned}$$

The final line is immediate.  $\square$

**Lemma 6.**  $\frac{1}{\sqrt{\text{Var}(\phi(Z|b_n, \eta))}} \leq \frac{1}{\sqrt{\sigma_{\min}^2 \mathbb{E}_{P(n)}[D/\max\{e(X), b_n\}^2]}}$ .

*Proof of Lemma 6.* By Lemma 5, I have:

$$\sigma_n^{-1} \leq n^{1/2} / \sqrt{\sigma_{\min}^2 \mathbb{E}_{P(n)} [D / \max\{e(X), b_n\}^2]},$$

where  $\sigma_n^{-1} = n^{-1/2} / \sqrt{\text{Var}_{P(n)}(\phi_n)}$ . □

**Lemma 7.** Define  $\tilde{\phi}(Z | b, P) \equiv \phi(Z | b, \eta(P)) - \mathbb{E}_P[\mu(X)]$  for  $P \in \mathcal{P}$ . Further define  $\rho(b, P) \equiv \mathbb{E}_P[|\tilde{\phi}(Z | b, P)|^3]$  and  $\sigma(b, P) \equiv \sqrt{\text{Var}_P(\tilde{\phi}(Z | b, P))}$ .

Then the following hold:

1.  $\mathbb{E}_P[\tilde{\phi}(Z | b, P)] = 0$
2.  $\sigma(b, P) > 0$
3.  $\rho(b, P) < \infty$  (though it may be arbitrarily large)
4. If  $b_n$  be a sequence of positive real numbers such that  $n^{-1/2} \ll b_n$  and  $P(n)$  be a sequence of distributions in  $\mathcal{P}$ , then  $\frac{\rho(b_n, P(n))}{\sigma(b_n, P(n))^3 \sqrt{n}} = o(1)$ .

*Proof of Lemma 7.*  $\mathbb{E}_P[\tilde{\phi}(Z | b_n, P)] = 0$  is immediate.

$\text{Var}_P[\tilde{\phi}(Z | b, P)] > 0$  follows by Lemma 5.

For the third moment being finite:

$$\begin{aligned} \rho(b, P) &= \mathbb{E}_P[|\tilde{\phi}(Z | b, P)|^3] \leq 8 \mathbb{E}_P[|\mu(X) - \mathbb{E}_P[\mu(X)]|^3 + b^{-3} |Y - \mu(X)|^3] \\ &\leq O(M^q) + 16b^{-3} \mathbb{E}_P[|Y|^3]. \end{aligned}$$

This is finite (and  $O(b^{-3})O(M^q)$ ) by assumption.

Finally, I have the claim for sequences. Recall that by Lemmas 5 and 3,  $\frac{1}{\sigma(b_n, P(n))^3 \sqrt{n}} = o(1)$  and  $\mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e(X), b_n\}^2} \right] \geq \sigma_{\min}^2/2$ . As a result:

$$\begin{aligned} \frac{\rho(b_n, P(n))}{\sigma(b_n, P(n))^3 \sqrt{n}} &\leq 8 \frac{\mathbb{E}_{P(n)} \left[ \frac{D|Y - \mu(X)|^3}{\max\{e(X), b_n\}^3} + |\mu(X) - \mathbb{E}_{P(n)}[\mu(X)]|^3 \right]}{\sigma(b_n, P(n))^3 \sqrt{n}} \\ &\leq O(M^q) \frac{\mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e(X), b_n\}^2} \right]}{b_n \sigma(b_n, P(n))^3 \sqrt{n}} + \frac{O(M^q)}{\sigma(b_n, P(n))^3 \sqrt{n}} \\ &= O(M^q, \sigma_{\min}^2) \mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e(X), b_n\}^2} \right]^{-1/2} (b_n^2 n)^{-1/2} + o(1) = o(1). \end{aligned}$$

□

*Proof of Proposition 4.* Let  $P(n)$  be a sequence of distributions in  $\mathcal{P}$ . By Lemma 7 and the Berry Esseen Theorem, the difference between the CDF of oracle clipped AIPW t-statistic  $\frac{\tilde{\psi}_{clip}^{AIPW} - \psi_n}{\sigma_n} = \frac{\sum \tilde{\phi}(Z|b_n, P(n))}{\sqrt{\text{Var}(\phi(Z|b_n, \eta))}\sqrt{n}}$  and the standard normal CDF  $\Phi$  is uniformly bounded above by  $\frac{3\rho(b_n, P(n))}{\sigma(b_n, P(n))^3\sqrt{n}}$ . By Lemma 7.4, this difference tends to zero. Therefore:

$$\limsup_{n \rightarrow \infty} \sup_{P \in \mathcal{P}} \sup_{t \in \mathbb{R}} \left| P \left( \frac{\tilde{\psi}_{(Oracle)}^{AIPW}(b_n) - \psi(P)}{\sigma_n} \leq t \right) - \Phi(t) \right| = \limsup_{n \rightarrow \infty} o(1) = 0.$$

□

## C.2 Consistency

**Lemma 8.** *Under cross-fitting, I have the bias bound:*

$$\left| \mathbb{E} \left[ \hat{\psi}_{clip}^{AIPW}(b_n) - \psi_n \right] \mid \hat{\mu}, \hat{e} \right| \leq r_{\mu, n} P(n) (e(X) \leq b_n + r_{e, n}) + r_{\mu, n} r_{e, n} \mathbb{E} \left[ \frac{\mathbf{1}\{e > b_n + r_{e, n}\}}{e - r_{e, n}} \right].$$

*Proof.* Fix one fold and take  $\hat{\mu}^{(-k)}$  and  $\hat{e}^{(-k)}$  as given. I write the bias relative to oracle clipped AIPW as:

$$\mathbb{E} \left[ (\hat{\mu} - \mu) \left( 1 - \frac{D}{\max\{\hat{e}, b_n\}} \right) + (\mu - Y) \left( \frac{D}{\max\{e, b_n\}} - \frac{D}{\max\{\hat{e}, b_n\}} \right) \right] = \mathbb{E} \left[ (\hat{\mu} - \mu) \left( 1 - \frac{D}{\max\{\hat{e}, b_n\}} \right) \right],$$

with the equality following by cross-fitting.

For  $p = 1, 2$ , let  $c_{n, p}$  solve:

$$\frac{|\max\{c_{n, p}, b_n\}^p - \max\{c_{n, p} - r_{e, n}, b_n\}^p|}{\max\{c_{n, p} - r_{e, n}, b_n\}^p} = \frac{|\max\{c_{n, p}, b_n\}^p - \max\{c_{n, p} + r_{e, n}, b_n\}^p|}{\max\{c_{n, p} + r_{e, n}, b_n\}^p}. \quad (4)$$

$c_{n, p}$  is useful because it is the changeover point between whether the worst-case  $\hat{e}$  is above or below  $e$ . Note that  $c_{n, p} \in (b_n, b_n + r_{e, n})$  by the intermediate value theorem: when  $c_{n, p} = b_n$ , the left-hand side is zero, while when  $c_{n, p} = b_n + r_{e, n}$ , the left hand side has the smaller denominator but equal numerator.

$c_{n, 1}$  is useful, because when  $e(X) = c_{n, 1}$ ,  $\hat{e}(X) = e(X) - r_{e, n}$  and  $\hat{e}(X) = e(X) + r_{e, n}$  produce equal levels of observation-wise bias from clipped inverse propensities relative to unity.

Then I have the bound:

$$\begin{aligned} \left| \mathbb{E} \left[ \hat{\psi}_{clip}^{AIPW}(b_n) - \psi_n \right] \right| &\leq \left| \mathbb{E} \left[ (\hat{\mu} - \mu) \frac{\max\{\hat{e}, b_n\} - e}{\max\{\hat{e}, b_n\}} \mathbf{1}\{e \leq b_n - r_{e, n}\} \right] \mid \hat{e}, \hat{\mu} \right| \\ &\quad + \left| \mathbb{E} \left[ (\hat{\mu} - \mu) \frac{\max\{\hat{e}, b_n\} - e}{\max\{\hat{e}, b_n\}} \mathbf{1}\{e \in (b_n - r_{e, n}, c_{n, 1}]\} \right] \mid \hat{e}, \hat{\mu} \right| \\ &\quad + \left| \mathbb{E} \left[ (\hat{\mu} - \mu) \frac{\max\{\hat{e}, b_n\} - e}{\max\{\hat{e}, b_n\}} \mathbf{1}\{e \in (c_{n, 1}, b_n + r_{e, n}]\} \right] \mid \hat{e}, \hat{\mu} \right| \end{aligned}$$

$$\begin{aligned}
& + \left| \mathbb{E} \left[ (\hat{\mu} - \mu) \frac{\max\{\hat{e}, b_n\} - e}{\max\{\hat{e}, b_n\}} \mathbf{1}\{e > b_n + r_{e,n}\} \right] \mid \hat{e}, \hat{\mu} \right| \\
& \leq \left| \mathbb{E} \left[ r_{\mu,n} \frac{b_n - e}{b_n} \mathbf{1}\{e \leq b_n - r_{e,n}\} \right] \mid \hat{e}, \hat{\mu} \right| \\
& + \left| \mathbb{E} \left[ r_{\mu,n} \frac{r_{e,n}}{e + r_{e,n}} \mathbf{1}\{e \in (b_n - r_{e,n}, c_{n,1}]\} \right] \mid \hat{e}, \hat{\mu} \right| \\
& + \left| \mathbb{E} \left[ r_{\mu,n} \frac{e - b_n}{b_n} \mathbf{1}\{e \in (c_{n,1}, b_n + r_{e,n}]\} \right] \mid \hat{e}, \hat{\mu} \right| \\
& + \left| \mathbb{E} \left[ (\hat{\mu} - \mu) \frac{r_{e,n}}{e - r_{e,n}} \mathbf{1}\{e > b_n + r_{e,n}\} \right] \mid \hat{e}, \hat{\mu} \right| \\
& \leq r_{\mu,n} P(n)(e(X) \leq b_n + r_{e,n}) + r_{\mu,n} r_{e,n} \mathbb{E} \left[ \frac{1}{e - r_{e,n}} \mathbf{1}\{e > b_n + r_{e,n}\} \right],
\end{aligned}$$

where the final line follows because  $r_{\mu,n}$  is always multiplied by a term that is bounded above by one for all  $e(X) \leq b_n + r_{e,n}$ .  $\square$

*Proof of Proposition 2.* Let  $P(n)$  be a sequence of distributions in  $\mathcal{P}$ , and fix some  $k \in 1, \dots, K$ .

Write  $\hat{\psi}_{clip}^{AIPW}(b_n) = \frac{1}{K} \sum_k \hat{\psi}_{clip,k}^{AIPW}(b_n)$ , where  $\hat{\psi}_{clip,k}^{AIPW}(b_n)$  is the fold- $k$  average potential outcome estimate. I will show that  $\hat{\psi}_{clip,k}^{AIPW}(b_n) - \psi_n = o_{P(n)}(1)$ .

First I show that  $\mathbb{E} \left[ \hat{\psi}_{clip,k}^{AIPW}(b_n) - \psi_n \mid \hat{\mu}^{(-k)}, \hat{e}^{(-k)} \right] = o_{P(n)}(1)$ . This holds by the assumptions of Proposition 2 applied to the bias bound from Lemma 8:

$$\left| E \left[ \hat{\psi}_{clip,k}^{AIPW}(b_n) - \psi_n \mid \hat{\mu}^{(-k)}, \hat{e}^{(-k)} \right] \right| \leq C r_{\mu,n} (b_n + r_{e,n})^{\gamma_0 - 1} + r_{\mu,n} r_{e,n} E \left[ \frac{\mathbf{1}\{e(X) > b_n + r_{e,n}\}}{e(X) - r_{e,n}} \right].$$

If  $r_{\mu,n} \frac{r_{e,n} + b_n}{b_n} \rightarrow_{P(n)} 0$ , then this term is  $o_{P(n)}(1)$  because  $r_{\mu,n}$  and  $r_{e,n} + b_n$  are bounded above by Assumption 2, and at least one of the two terms must tend to zero because  $b_n \rightarrow_{P(n)} 0$  and  $r_{\mu,n} \frac{r_{e,n} + b_n}{b_n} \rightarrow_{P(n)} 0$ .

Suppose  $r_{e,n} b_n^{\min\{\gamma_0 - 2, 0\}} \rightarrow_{P(n)} 0$ . For the first term,  $r_{\mu,n} (b_n + r_{e,n})^{\gamma_0 - 1} \rightarrow_{P(n)} 0$  because  $r_{\mu,n}$  is bounded above and  $\gamma_0 > 1$ . For the final term, suppose that  $\gamma_0 < 2$ , so that the claim is not immediate. Then:

$$\begin{aligned}
E_{P(n)} \left[ \frac{\mathbf{1}\{e(X) > b_n + r_{e,n}\}}{e(X) - r_{e,n}} \right] &= C(b_n + r_{e,n})^{\gamma_0 - 1} b_n^{-1} + C(\gamma_0 - 1) \int_{b_n + r_{e,n}}^{C^{-1}/(\gamma_0 - 1)} (t - r_{e,n})^{-1} t^{\gamma_0 - 2} dt \\
&\leq C(b_n + r_{e,n})^{\gamma_0 - 1} b_n^{-1} + C(\gamma_0 - 1) \int_{b_n + r_{e,n}}^{C^{-1}/(\gamma_0 - 1)} (t - r_{e,n})^{\gamma_0 - 3} dt \\
&\leq C(b_n + r_{e,n})^{\gamma_0 - 1} b_n^{-1} + C(\gamma_0 - 1) \int_{b_n}^1 x^{\gamma_0 - 3} dx \\
&= C(b_n + r_{e,n})^{\gamma_0 - 1} b_n^{-1} + \frac{C(\gamma_0 - 1)}{2 - \gamma_0} (b_n^{\gamma_0 - 2} - 1) \\
&\leq \left( 2^{\gamma_0 - 1} C + \frac{\gamma_0 - 1}{2 - \gamma_0} C \right) b_n^{\gamma_0 - 2} + 2C r_{e,n}^{\gamma_0 - 1} b_n^{-1} \\
&= O(b_n^{\gamma_0 - 2} (1 + \delta_n^{1 - \gamma_0})) = O(b_n^{\gamma_0 - 2}),
\end{aligned}$$

so that  $r_{e,n}E_{P(n)} \left[ \frac{\mathbf{1}\{e(X) > b_n + r_{e,n}\}}{e(X) - r_{e,n}} \right] = o(1)$ . Note that this bound may be lax. Whether the propensity rate requirement could be weakened is an open question for future work. Regardless, in this remaining case under the propensity rate requirement (i),  $r_{\mu,n}r_{e,n}b_n^{\gamma_0-2}\delta_n \rightarrow_{P(n)} 0$  by construction of  $\delta_n$ , so that the final term of the bias bound tends to zero.

Next, I show that  $V(\hat{\psi}_{clip,k}^{AIPW}(b_n) \mid \hat{\mu}^{(-k)}, \hat{e}^{(-k)}) = o_{P(n)}(1)$ . I have:

$$\begin{aligned} V(\hat{\psi}_{clip,k}^{AIPW}(b_n) \mid \hat{\mu}^{(-k)}, \hat{e}^{(-k)}) &\leq n^{-1} \mathbb{E} \left[ \left( \hat{\mu} + \frac{D(Y - \hat{\mu})}{\max\{\hat{e}, b_n\}} \right)^2 \mid \hat{\mu}, \hat{e} \right] \\ &\leq 8n^{-1}b_n^{-2} \quad (\text{Lemma 3.(iv)}) \\ &= o(1). \end{aligned}$$

Therefore  $\mathbb{E} \left[ \left( \hat{\psi}_{clip,k}^{AIPW}(b_n) - \psi_n \right)^2 \mid \hat{\mu}^{(-k)}, \hat{e}^{(-k)} \right] = o_{P(n)}(1)$ . Therefore, for all  $\epsilon > 0$ , every  $\delta > 0$ , and every sequence of  $P(n)$ , there is an  $n$  large enough such that  $P(n) \left( \left( \hat{\psi}_{clip,k}^{AIPW}(b_n) - \psi_n \right)^2 > \epsilon^2 \right) \leq \delta$ . Thus, consistency holds.  $\square$

### C.3 Degradation of Consistency Rate

*Proof of Proposition 3.* Define  $\sigma_{\max}^2 = \sup_{P \in \mathcal{P}} \sup_{X,D} \text{Var}(Y \mid X, D)$ . By the presence of  $q > 2$  moments,  $\sigma_{\max}^2$  is finite.

By Lemma 5,  $\mathbb{E}_{P(n)} \left[ \frac{D\sigma_{\min}^2}{\max\{e(X), b_n\}^2} \right] \not\rightarrow 0$ .

By iid sampling and the oracle AIPW conditional mean being equal to  $\mu(X)$ , I obtain:

$$\begin{aligned} \text{Var}_{P(n)} \left( \tilde{\psi}_{(Oracle)}^{AIPW}(b_n) \right) - n^{-1} \text{Var}_{P(n)} (\mu(X)) &= \mathbb{E}_{P(n)} \left[ \text{Var}_{P(n)} \left( \tilde{\psi}_{(Oracle)}^{AIPW}(b_n) \mid \{X\} \right) \right] \\ &= n^{-1} \mathbb{E}_{P(n)} \left[ \text{Var}_{P(n)} \left( \frac{D(Y - \mu(X))}{\max\{e(X), b_n\}} \mid X \right) \right] \\ &= n^{-1} \mathbb{E}_{P(n)} \left[ e(X) \text{Var}_{P(n)} \left( \frac{(Y - \mu(X))}{\max\{e(X), b_n\}} \mid X, D = 1 \right) \right] \\ &= n^{-1} \mathbb{E}_{P(n)} \left[ e(X) \frac{\text{Var}(Y \mid X, D = 1)}{\max\{e(X), b_n\}^2} \right] \\ &= \Theta \left( n^{-1} \mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e(X), b_n\}^2} \right] \right). \end{aligned}$$

In addition,  $n^{-1} \text{Var}_{P(n)} (\mu(X)) \leq n^{-1}M = O \left( n^{-1} \mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e(X), b_n\}^2} \right] \right)$ , proving the claim.  $\square$

*Proof of Corollary 3.* Let  $n$  be large enough that  $b_n \leq 1$  and  $b_n^{\gamma_0-2} > 2$ . Recall the definition of  $\sigma_{\max}^2$  from the proof of Proposition 3.

For the upper bound, let  $\mathcal{P}$  be arbitrary:

$$\begin{aligned}
\sigma_n^2 - n^{-1}Var(\mu) &= n^{-1}Var_P\left(\mu(X) + \frac{D(Y - \mu(X))}{\max\{e(X), b_n\}}\right) - n^{-1}Var_P(\mu(X)) = n^{-1}E_P\left[\frac{D(Y - \mu(X))^2}{\max\{e(X), b_n\}^2}\right] \\
&\leq n^{-1}E_P\left[\frac{D\sigma_{\max}^2}{\max\{e(X), b_n\}^2}\right] = \sigma_{\max}^2 E_P\left[\frac{e(X)}{\max\{e(X), b_n\}^2}\right] = n^{-1}\sigma_{\max}^2 \int_0^\infty P\left(\frac{e(X)}{\max\{e(X), b_n\}^2} \geq t\right) dt \\
&= n^{-1}\sigma_{\max}^2 \int_0^\infty P(e(X) \leq b_n, e(X) \geq tb_n^2) dt + n^{-1}\sigma_{\max}^2 \int_0^\infty P(e(X) > b_n, e(X) \leq 1/t) dt \\
&= n^{-1}\sigma_{\max}^2 \int_0^{b_n^{-1}} P(e(X) \in [tb_n^2, b_n]) dt + n^{-1}\sigma_{\max}^2 \int_0^{b_n^{-1}} P(e(X) \in [b_n, 1/t]) dt \\
&= n^{-1}\sigma_{\max}^2 \int_0^{b_n^{-1}} P(e(X) \in [tb_n^2, 1/t]) dt \leq n^{-1}\sigma_{\max}^2 \int_0^{b_n^{-1}} P(e(X) \leq 1/t) dt \\
&\leq Cn^{-1}\sigma_{\max}^2 \int_0^{b_n^{-1}} t^{1-\gamma_0} dt = \underbrace{\frac{C\sigma_{\max}^2}{\gamma_0 - 2}}_{C'} n^{-1}b_n^{\gamma_0-2}.
\end{aligned}$$

For the remaining term,  $n^{-1}Var(\mu) = O(n^{-1}) = o(C'n^{-1}b_n^{\gamma_0-2})$ .

For the lower bound, define  $\mathcal{P} = \{P\}$ , where  $P$  is the distribution which draws  $e(X)$  from the CDF  $P(e(X) \leq \pi) = (1 - \pi_{\min}) \min\{C\pi^{\gamma_0-1}, 1\} + \pi_{\min}\mathbf{1}\{\pi \geq 1\}$  and  $Y | X, D \sim \mathcal{N}(0, \sigma_{\min}^2)$ . This distribution has valid conditional moments and residual variance by the minimal value of  $M$  and the choice of  $Var(Y | X < D)$ . The treated fraction is at least  $\pi_{\min} > 0$ . For all  $\pi < C^{-1/(\gamma_0-1)}$ ,  $P(e(X) \leq \pi) = (1 - \pi_{\min})C\pi^{\gamma_0-1} \leq C\pi^{\gamma_0-1}$ ; for all  $\pi > C^{-1/(\gamma_0-1)}$ ,  $P(e(X) \leq \pi) = (1 - \pi_{\min}) + \pi_{\min}\mathbf{1}\{\pi = 1\}$ , which must be below  $C\pi^{\gamma_0-1}$  for all such  $\pi$  in order for  $\mathcal{P}$  to be non-empty. Finally, note that:

$$\begin{aligned}
\sigma_n^2 - n^{-1}Var_P(\mu) &= n^{-1}E_P\left[\frac{D(Y - \mu(X))^2}{\max\{e, b_n\}^2}\right] \\
&= n^{-1}\sigma_{\min}^2 \left( \begin{aligned} &\int_0^{b_n} \frac{t}{b_n^2} (1 - \pi_{\min})(\gamma_0 - 1)Ct^{\gamma_0-2} dt \\ &+ \int_{b_n}^1 \frac{1}{t} (1 - \pi_{\min})(\gamma_0 - 1)Ct^{\gamma_0-2} dt \\ &+ \pi_{\min} \end{aligned} \right) \\
&= n^{-1}\sigma_{\min}^2 \left( \begin{aligned} &b_n^{-2}(1 - \pi_{\min})(\gamma_0 - 1)C \int_0^{b_n} t^{\gamma_0-1} dt \\ &+ (1 - \pi_{\min})(\gamma_0 - 1)C \int_{b_n}^1 t^{\gamma_0-3} dt \\ &+ \pi_{\min} \end{aligned} \right) \\
&= n^{-1}\sigma_{\min}^2 \left( \begin{aligned} &b_n^{\gamma_0-2}(1 - \pi_{\min})C \left( \frac{\gamma_0-1}{\gamma_0} + \frac{\gamma_0-1}{2-\gamma_0} \right) \\ &+ \pi_{\min} - (1 - \pi_{\min})C \frac{\gamma_0-1}{2-\gamma_0} \end{aligned} \right) \\
&\geq \underbrace{\sigma_{\min}^2 C(1 - \pi_{\min}) \left( \frac{\gamma_0-1}{\gamma_0} + \frac{\gamma_0-1}{2(2-\gamma_0)} \right)}_{C''} n^{-1}b_n^{\gamma_0-2}.
\end{aligned}$$

Note also that  $C'' > 0$ . Thus,  $C''n^{-1}b_n^{\gamma_0-2} \leq \sup_{P \in \mathcal{P}} \sigma_n^2 - \text{Var}_P(\mu(X)) \leq C'n^{-1}b_n^{\gamma_0-2}$ . Analogously to before,  $n^{-1}\text{Var}(\mu) = o(C''n^{-1}b_n^{\gamma_0-2})$ , completing the proof.  $\square$

## C.4 Asymptotic Normality and Rates

**Lemma 9.** *Suppose the conditions of Proposition 4 hold and let  $P(n)$  be a sequence of distributions in  $\mathcal{P}$ . Then,  $\mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e(X), b_n\}^2} \right] = O(1 + b_n^{\gamma_0-2} \log(1/b_n)^{1\{\gamma_0 \neq 2\}})$ , with a constant that only depends on  $C$  and  $\gamma_0$ .*

*Proof of Lemma 9.* Let  $P(n)$  be given. Then:

$$\begin{aligned} \mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e(X), b_n\}^2} \right] &= \mathbb{E}_{P(n)} \left[ \frac{e(X)}{\max\{e(X), b_n\}^2} \right] \leq \mathbb{E}_{P(n)} \left[ \frac{1}{\max\{e(X), b_n\}} \right] = \int_0^\infty P(n) \left( \frac{1}{\max\{e(X), b_n\}} > t \right) dt \\ &= \int_0^\infty P(n) (\max\{e(X), b_n\} < 1/t) dt = 1 + \int_1^\infty P(n) (\max\{e(X), b_n\} < 1/t) dt \\ &= 1 + \int_1^{1/b_n} P(n) (\max\{e(X), b_n\} < 1/t) dt \leq 1 + C \int_1^{b_n^{-1}} t^{1-\gamma_0} dt. \end{aligned}$$

First, suppose  $\gamma_0 = 2$ . Then  $\mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e(X), b_n\}^2} \right] \leq 1 + C(\log(1/b_n) - 1)$ . Alternatively, suppose  $\gamma_0 \neq 2$ . Then  $\mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e(X), b_n\}^2} \right] \leq 1 + \frac{C}{\gamma_0-2} (1 - b_n^{\gamma_0-2})$ .  $\square$

**Lemma 10.** *Suppose the conditions of Proposition 4 hold and  $P(n)$  is a sequence of distributions in  $\mathcal{P}$ . Then*

$$\frac{1}{\sqrt{\mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e, b_n\}^2} \right]}} = O \left( \frac{1}{\sqrt{1 + b_n^{-2} P(n)(e(X) \leq b_n)^{\gamma_0/(\gamma_0-1)}}} \right).$$

*Proof of Lemma 10.* For any  $m \geq 0$ , define  $\mathcal{P}_{m,n} = \{P \in \mathcal{P} \mid P(e(X) \leq b_n) \leq m\}$ . For each  $n$ , define  $m_n = P(n)(e(X) \leq b_n)^{\gamma_0/(\gamma_0-1)}$ .

I have:

$$\begin{aligned} \sup_{P \in \mathcal{P}_{m_n,n}} \mathbb{E}_P \left[ \frac{D}{\max\{e, b_n\}^2} \right] &= \sup_{P \in \mathcal{P}_{m_n,n}} \mathbb{E}_P \left[ \frac{D \mathbf{1}\{e(X) > b_n\}}{\max\{e, b_n\}^2} \right] + \mathbb{E}_P \left[ \frac{D \mathbf{1}\{e(X) \leq b_n\}}{b_n^2} \right] \\ &\geq 1 + \sup_{P \in \mathcal{P}_{m_n,n}} b_n^{-2} \mathbb{E}_P[e(X) \mathbf{1}\{e(X) \leq b_n\}] \\ &\geq 1 + c(\gamma_0) \sup_{P \in \mathcal{P}_{m_n,n}} P(e(X) \leq b_n)^{\gamma_0/(\gamma_0-1)} \quad (\text{Lemma 4}) \\ &= 1 + c(\gamma_0) b_n^{-2} m_n^{\gamma_0/(\gamma_0-1)}. \end{aligned}$$

Therefore  $(1 + b_n^{-2} P(n)(e(X) \leq b_n)^{\gamma_0/(\gamma_0-1)})^2 = O \left( \mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e, b_n\}^2} \right] \right)$ . Taking the square root and inverting both sides completes the proof.  $\square$

*Proof of Corollary 6.* The two changes are the product-of-errors term being stated as 3'(b) and the regression-error-near-singularities term being stated as 3'(c).

I begin with the product-of-errors term (b). Suppose Assumption 3(b) holds. I wish to show that if  $r_{\mu,n}r_{e,n} \left(1 + b_n^{(\gamma_0-2)/2} \log(1/b_n) \mathbf{1}_{\{\gamma_0=2\}/2}\right) = o(n^{-1/2})$ , then  $r_{\mu,n}r_{e,n} \sqrt{\mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e(X), b_n^2\}} \right]} = o_{P(n)}(n^{-1/2})$ . Let  $\alpha > 0$  from Assumption 3 be given. By Lemma 9, for all  $\alpha > 0$ ,

$$r_{\mu,n}r_{e,n} \sqrt{\mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e(X), b_n^2\}} \right]} = O_{P(n)} \left( r_{\mu,n}r_{e,n} \left(1 + b_n^{(\gamma_0-2)/2} \log(1/b_n) \mathbf{1}_{\{\gamma_0=2\}/2}\right) \right) = o(n^{-1/2}).$$

Next I consider the regression error near the singularities term (c). I first verify (c) for a sequence of  $P(n) \in \mathcal{P}$  under Assumption 4(i). Let a sequence of  $P(n) \in \mathcal{P}$  and  $b_n$  be given, and consider an arbitrary sub-sequence. If there is a further sub-sub-sequence with  $P(n)(e(X) \leq b_n) = 0$ , take this sub-sub-sequence and the claim holds. If not, take a sub-sub-sequence with  $P(n)(e(X) \leq b_n) > 0$ . On this sub-sub-sequence, I have the bound:

$$\begin{aligned} \mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e, b_n\}^2} \right] &\geq \mathbb{E}_{P(n)} \left[ \frac{D \mathbf{1}\{e \in (b_n/2, b_n]\}}{\max\{e, b_n\}^2} \right] \\ &\geq \frac{1}{2b_n} P(n)(e \in (b_n/2, b_n]) \\ &= \frac{1}{2b_n} (P(n)(e \leq b_n) - P(n)(e \leq b_n/2)) \\ &\geq \frac{\rho}{2b_n} P(n)(e \leq b_n). \end{aligned} \tag{Assumption 4(i)}$$

Then, applying Assumption 1(d),

$$r_{\mu,n} \frac{P(n)(e(X) \leq b_n)}{\sqrt{\mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e, b_n\}^2} \right]}} \leq r_{\mu,n} \left( \frac{2b_n P(n)(e \leq b_n)}{\rho} \right)^{1/2} \leq \sqrt{\frac{2C}{\rho}} r_{\mu,n} b_n^{\gamma_0/2} = o(n^{-1/2}).$$

Finally, I verify (c) assuming Assumption 4(ii) holds. I wish to show that if there is a sequence of  $P(n) \in \mathcal{P}$  and associated constants such that  $r_{\mu,n} b_n^{\gamma_0/2} = o(n^{-1/2})$ , then  $r_{\mu,n} \frac{P(n)(e(X) \leq b_n)}{\sqrt{\mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e, b_n\}^2} \right]}} = o(n^{-1/2})$ .

Under somewhat weak overlap ( $\gamma_0 > 2$ ), then

$$r_{\mu,n} \frac{P(n)(e(X) \leq b_n)}{\sqrt{\mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e, b_n\}^2} \right]}} = O(r_{\mu,n} b_n^{\gamma_0-1}) = O\left(r_{\mu,n} b_n^{2(\gamma_0-1)/\gamma_0 + (\gamma_0^2+2-3\gamma_0)/\gamma_0}\right) = o(n^{-1/2}).$$



I therefore proceed for  $\gamma_0 \leq 2$ . I will use the bound from Lemma 10:

$$r_{\mu,n} \frac{P(n)(e(X) \leq b_n)}{\sqrt{\mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e, b_n\}^2} \right]}} \leq r_{\mu,n} \frac{P(n)(e(X) \leq b_n)}{\sqrt{1 + b_n^{-2} P(n)(e(X) \leq b_n)^{\gamma_0/(\gamma_0-1)}}} = \text{“}RHS\text{”}$$

Consider a sub-sequence of  $P(n) \in \mathcal{P}$  and constants. I will show that there is a further sub-sub-sequence for which this right-hand side  $RHS$  is  $o(n^{-1/2})$ . Suppose there is a sub-sub-sequence such that  $b_n^{-2} P(n)(e(X) \leq b_n)^{\gamma_0/(\gamma_0-1)} \rightarrow 0$ . Then for that sub-sub-sequence, I have:

$$RHS \leq r_{\mu,n} P(n)(e(X) \leq b_n) = r_{\mu,n} o\left(b_n^{2(\gamma_0-1)/\gamma_0}\right) = o(n^{-1/2}).$$

If not, then  $b_n^2 \lesssim P(n)(e(X) \leq b_n)^{\gamma_0/(\gamma_0-1)}$  and I have the bound:

$$\begin{aligned} RHS &\leq O\left(r_{\mu,n} b_n P(n)(e(X) \leq b_n)^{1-\gamma_0/(2(\gamma_0-1))}\right) \\ &= O\left(r_{\mu,n} b_n \left(P(n)(e(X) \leq b_n)^{(1/2-1/\gamma_0)*\gamma_0/(\gamma_0-1)}\right)\right) \\ &= O\left(r_{\mu,n} b_n \left(P(n)(e(X) \leq b_n)^{\gamma_0/(\gamma_0-1)}\right)^{(\gamma_0-2)/(2\gamma_0)}\right) \\ &= O\left(r_{\mu,n} b_n (b_n^2)^{(\gamma_0-2)/(2\gamma_0)}\right) = O\left(r_{\mu,n} b_n^{2(\gamma_0-1)/\gamma_0}\right) = o(n^{-1/2}). \end{aligned}$$

Therefore Assumption 3'(c) holds. □

**Lemma 11.** Suppose the requirements of Theorem 1' hold. Then by implication,  $r_{\mu,n} P(n)(e(X) \leq b_n) \ll n^{-1/2} \sqrt{\mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e, b_n\}^2} \right]}$ .

*Proof of Lemma 11.*

$$\begin{aligned} r_{\mu,n} P(n)(e(X) \leq b_n) &= \left( r_{\mu,n} \frac{P(n)(e(X) \leq b_n)}{\sqrt{\mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e, b_n\}^2} \right]}} \right) \sqrt{\mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e, b_n\}^2} \right]} \\ &\ll n^{-1/2} \sqrt{\mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e, b_n\}^2} \right]}. \end{aligned} \tag{A3'(c)}$$

□

**Lemma 12** (Oracle consistency). If  $n^{-1/2} \ll b_n \ll 1$ , then  $|P_n[\phi_n] - \mathbb{E}_{P(n)}[\mu(x)]| \xrightarrow{\mathcal{P}} 0$ .

*Proof of Lemma 12.* Let  $P(n)$  be a sequence of distributions in  $\mathcal{P}$ . For any  $t > 0$ , I have:

$$P(n) (|P(n)_n[\phi_n] - \mathbb{E}_{P(n)}[\mu(x)]| > t) \leq \frac{\mathbb{E}[|\phi_n - \mathbb{E}_{P(n)}[\mu(x)]|^2]}{nt^2} \quad (\text{Chebyshev's inequality})$$

$$\begin{aligned}
&\leq \frac{\mathbb{E}[|\phi_n - \mathbb{E}_{P(n)}[\phi_n]|^q]^{2/q}}{nt^2} && \text{(Jensen's inequality)} \\
&\leq \frac{[(4M)^q \mathbb{E}[e(X)/\{e(X) \vee b_n\}^2]]^{2/q}}{t^2 nb_n^{2(q-2)/q}} && \text{(Lemma 3.(iii))} \\
&\leq \frac{(4M)^2}{t^2} \frac{1}{nb_n^2}.
\end{aligned}$$

This upper bound tends to zero and holds simultaneously for all  $P \in \mathcal{P}$ . Hence,  $|P(n)_n[\phi_n] - \mathbb{E}_{P(n)}[\mu(x)]| = o_{P(n)}(1)$ .  $\square$

**Lemma 13** (Oracle variance consistency). *Let  $\sigma_n^2 = n^{-1}(P_n[\phi_n^2] - P_n[\phi_n]^2)$  be the oracle sample variance. If  $n^{-1/2} \ll b_n \ll 1$ , then  $n\sigma_n^2/\text{Var}_{P(n)}(\phi_n) \xrightarrow{\mathcal{P}} 1$ .*

*Proof of Lemma 13.* Let  $P(n)$  be a sequence of distributions in  $\mathcal{P}$ .

First, I argue that for any  $q > 2$ :

$$\begin{aligned}
P\left(\left|\frac{P(n)_n[\phi_n^2] - P[\phi_n^2]}{\text{Var}_{P(n)}(\phi_n)}\right| > t\right) &\leq \frac{\mathbb{E}\{|P(n)_n[\phi_n^2] - P[\phi_n^2]|^{q/2}\}}{t^{q/2} \text{Var}_{P(n)}(\phi_n)^{q/2}} && \text{(Markov inequality)} \\
&\leq \frac{2}{t^{q/2} n^{q/2-1}} \frac{\mathbb{E}\{|\phi_n^2 - P[\phi_n^2]|^{q/2}\}}{\text{Var}_{P(n)}(\phi_n)^{q/2}} && \text{(von Bahr-Esseen inequality)} \\
&\leq \frac{2^{q/2+1}}{t^{q/2} n^{q/2-1}} \frac{\mathbb{E}[|\phi_n|^q]}{\text{Var}_{P(n)}(\phi_n)^{q/2}} && \text{(Jensen's inequality)} \\
&\leq \frac{2^{q/2+1}}{t^{q/2} n^{q/2-1}} \frac{(8M)^q \mathbb{E}[e(X)/\{e(X) \vee b_n\}^2]}{b_n^{q-2} (\text{Var}_{P(n)}(\phi_n))^{q/2}} && \text{(Lemma 3.(iv))} \\
&\leq \frac{2^{q/2+1}}{t^{q/2} n^{q/2-1}} \frac{(8M)^q \mathbb{E}[e(X)/\{e(X) \vee b_n\}^2]}{b_n^{q-2} \sigma_{\min}^q \mathbb{E}[e(X)/\{e(X) \vee b_n\}^2]^{q/2}} && \text{(Lemma 5)} \\
&\leq \frac{(8M)^{q/2+1}}{t^{q/2} \sigma_{\min}^q (\pi_{\min}/2)^{q/2-1}} \frac{1}{n^{q/2-1} b_n^{q-2}}. && \text{(Lemma 5)}
\end{aligned}$$

Since  $b_n \gg n^{-1/2}$ ,  $n^{q/2-1} b_n^{q-2} \rightarrow \infty$ , so that  $\left|\frac{P(n)_n[\phi_n^2] - P[\phi_n^2]}{\text{Var}_{P(n)}(\phi_n)}\right| = o_{P(n)}(1)$ .

Then, by the triangle inequality:

$$\begin{aligned}
|n\sigma_n^2/\text{Var}_{P(n)}(\phi_n) - 1| &\leq \left|\frac{P(n)_n[\phi_n]^2 - P[\phi_n]^2}{\text{Var}_{P(n)}(\phi_n)}\right| + \left|\frac{P(n)_n[\phi_n^2] - P[\phi_n^2]}{\text{Var}_{P(n)}(\phi_n)}\right| \\
&\leq |P(n)_n[\phi_n] + P[\phi_n]| \times |P(n)_n[\phi_n] - P[\phi_n]|/(\sigma_{\min}^2 \pi_{\min}/2) + o_{P(n)}(1) \\
&\quad \text{(Lemma 5 + above)} \\
&\leq (2M + o_{P(n)}(1)) o_{P(n)}(1) O_{P(n)}(1) + o_{P(n)}(1) \quad (P[\phi_n] \leq M + \text{Lemma 12}) \\
&= o_{P(n)}(1),
\end{aligned}$$

where  $\sigma_n^2$  is the oracle sample variance. Therefore, this upper bound tends to zero uniformly over  $\mathcal{P}$ .  $\square$

**Lemma 14** (Orthogonalized inverse propensities). *Suppose the conditions of Proposition 2 hold,  $b_n$  is a sequence such that  $b_n \geq 2r_{e,n}$  for all  $n$  large enough, and  $P(n)$  is a sequence of distributions in  $\mathcal{P}$ . Then*

$$\mathbb{E}_{P(n)} \left[ \left( \frac{D}{\max\{\hat{e}, b_n\}} - \frac{D}{\max\{e, b_n\}} \right)^2 \right] = o_{P(n)} \left( \mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e, b_n\}^2} \right] \right).$$

*Proof of Lemma 14.* Write  $(I) = \mathbb{E}_{P(n)} \left[ \left( \frac{D}{\max\{\hat{e}, b_n\}} - \frac{D}{\max\{e, b_n\}} \right)^2 \right]$ . Since  $r_{e,n} \ll b_n$ , let  $n$  be large enough that  $b_n \geq 2r_{e,n}$ , so that  $e \geq b_n + r_{e,n}$  implies  $e - r_{e,n} \geq e/2$ .

Then the squared error has the following decomposition

$$\begin{aligned} (I) &= \mathbb{E}_{P(n)} \left[ \frac{D}{\max\{\hat{e}, b_n\}^2} - \frac{D}{\max\{e, b_n\}^2} \right] - 2\mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e, b_n\}} \left( \frac{D}{\max\{\hat{e}, b_n\}} - \frac{D}{\max\{e, b_n\}} \right) \right] \\ &= \mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e, b_n\}^2} \frac{\max\{e, b_n\}^2 - \max\{\hat{e}, b_n\}^2}{\max\{\hat{e}, b_n\}^2} \right] - 2\mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e, b_n\}^2} \left( \frac{\max\{e, b_n\} - \max\{\hat{e}, b_n\}}{\max\{\hat{e}, b_n\}} \right) \right] \\ &= \mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e, b_n\}^2} \frac{\max\{e, b_n\}^2 - \max\{\hat{e}, b_n\}^2}{\max\{\hat{e}, b_n\}^2} \right] + 2\mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e, b_n\}^2} \left( \frac{\max\{\hat{e}, b_n\} - \max\{e, b_n\}}{\max\{\hat{e}, b_n\}} \right) \right] \\ &\leq \mathbb{E}_{P(n)} \left[ \frac{D\mathbf{1}\{\hat{e} \leq e\}}{\max\{e, b_n\}^2} \frac{\max\{e, b_n\}^2 - \max\{\hat{e}, b_n\}^2}{\max\{\hat{e}, b_n\}^2} \right] + 2\mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e, b_n\}^2} \left( \frac{r_{e,n}}{\max\{e + r_{e,n}, b_n\}} \right) \right] \\ &\leq \mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e, b_n\}^2} \frac{\max\{e, b_n\}^2 - \max\{e - r_{e,n}, b_n\}^2}{\max\{e - r_{e,n}, b_n\}^2} \right] + 2\mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e, b_n\}^2} \left( \frac{r_{e,n}}{b_n} \right) \right] \\ &\leq \mathbb{E}_{P(n)} \left[ \frac{D\mathbf{1}\{e \in [b_n, b_n + r_{e,n}]\}}{\max\{e, b_n\}^2} \frac{(b_n + r_{e,n})^2 - b_n^2}{b_n^2} \right] + \mathbb{E}_{P(n)} \left[ \frac{D\mathbf{1}\{e \in [b_n + r_{e,n}, 1]\}}{\max\{e, b_n\}^2} \frac{e^2 - (e - r_{e,n})^2}{(e - r_{e,n})^2} \right] \\ &\quad + 2\mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e, b_n\}^2} \left( \frac{r_{e,n}}{b_n} \right) \right] \\ &\leq \mathbb{E}_{P(n)} \left[ \frac{D\mathbf{1}\{e \in [b_n, b_n + r_{e,n}]\}}{\max\{e, b_n\}^2} \frac{2b_n r_{e,n} + r_{e,n}^2}{b_n^2} \right] + \mathbb{E}_{P(n)} \left[ \frac{D\mathbf{1}\{e \in [b_n + r_{e,n}, 1]\}}{\max\{e, b_n\}^2} \frac{2er_{e,n}}{(e - r_{e,n})^2} \right] \\ &\quad + 2\mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e, b_n\}^2} \left( \frac{r_{e,n}}{b_n} \right) \right] \\ &\leq \mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e, b_n\}^2} \left( \frac{4b_n r_{e,n} + r_{e,n}^2}{b_n^2} \right) \right] + 4\mathbb{E}_{P(n)} \left[ \frac{D\mathbf{1}\{e \in [b_n + r_{e,n}, 1]\}}{\max\{e, b_n\}^2} \frac{2r_{e,n}}{e} \right] \\ &\leq \mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e, b_n\}^2} \left( \frac{12b_n r_{e,n} + r_{e,n}^2}{b_n^2} \right) \right]. \end{aligned}$$

Since  $r_{e,n} = o(b_n)$  by assumption, this upper bound is  $o \left( \mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e, b_n\}^2} \right] \right)$ .  $\square$

**Lemma 15.** *For all  $P \in \mathcal{P}$ ,  $\mathbb{E}_P[D / \max\{e(X), b_n\}^2] \geq 1 - \left( b_n^2 \frac{\gamma_0 - 1}{c(\gamma_0)} \right)^{\gamma_0 - 1} \gamma_0^{-\gamma_0}$ .*

*Proof of Lemma 15.*

$$\begin{aligned} \mathbb{E}_P[D / \max\{e(X), b_n\}^2] &= \mathbb{E}_P[D\mathbf{1}\{e(X) \leq b_n\} / \max\{e(X), b_n\}^2] + \mathbb{E}_P[D\mathbf{1}\{e(X) > b_n\} / \max\{e(X), b_n\}^2] \\ &\geq b_n^{-2} \mathbb{E}_P[D\mathbf{1}\{e(X) \leq b_n\}] + (1 - P(e(X) \leq b_n)) \end{aligned}$$

$$\begin{aligned}
&\geq b_n^{-2} c(\gamma_0) P(e(X) \leq b_n)^{\gamma_0/(\gamma_0-1)} + (1 - P(e(X) \leq b_n)) \quad (\text{Lemma 4}) \\
&= 1 + P(e(X) \leq b_n) \left( c(\gamma_0) b_n^{-2} P(e(X) \leq b_n)^{1/(\gamma_0-1)} - 1 \right).
\end{aligned}$$

This term is minimized over  $P(e(X) \leq b_n)$  at  $P(e(X) \leq b_n) = \left( \frac{\gamma_0-1}{c(\gamma_0) b_n^{-2} \gamma_0} \right)^{\gamma_0-1}$ , which produces

$$\mathbb{E}_P[D / \max\{e(X), b_n\}^2] = 1 - \frac{P(e(X) \leq b_n)}{\gamma_0} = 1 - \left( b_n^2 \frac{\gamma_0-1}{c(\gamma_0)} \right)^{\gamma_0-1} \gamma_0^{-\gamma_0}.$$

□

**Lemma 16.** *Suppose the conditions of Lemma 14 hold and  $r_{\mu,n} \rightarrow 0$ . Let  $P(n)$  be a sequence of distributions in  $\mathcal{P}$ . Recall the definitions of  $\phi_n$  as the oracle clipped influence function and  $\hat{\phi}_n$  the estimated influence function. Then  $\mathbb{E}_{P(n)}[\phi_n^2] = \text{Var}_{P(n)}(\phi_n) + O(1)$  and  $P(n)_n [\hat{\phi}_n^2 - \phi_n^2] = o_{P(n)}(\text{Var}_{P(n)}(\phi_n))$ .*

*Proof of Lemma 16.* First, note that:

$$\mathbb{E}_{P(n)}[\phi_n^2] = \text{Var}_{P(n)}(\phi_n) + \mathbb{E}_{P(n)}[\phi_n]^2 = \text{Var}_{P(n)}(\phi_n) + O(1)^2. \quad (\text{Assumption 1(a)})$$

Next, I show that  $P(n)_n [\hat{\phi}_n^2 - \phi_n^2] = o_{P(n)}(\mathbb{E}_{P(n)}[D / \max\{e(X), b_n\}^2])$ . I have:

$$\begin{aligned}
\left| P(n)_n [\hat{\phi}_n^2 - \phi_n^2] \right| &= P(n)_n [(\hat{\phi}_n - \phi_n)^2] + \left| P(n)_n [2(\hat{\phi}_n - \phi_n)\phi_n] \right| \\
&\leq P(n)_n [(\hat{\phi}_n - \phi_n)^2] + 2\sqrt{P(n)_n [(\hat{\phi}_n - \phi_n)^2] P(n)_n [\phi_n^2]} \quad (\text{Cauchy-Schwarz}) \\
&= P(n)_n [(\hat{\phi}_n - \phi_n)^2] + \sqrt{P(n)_n [(\hat{\phi}_n - \phi_n)^2] O_{P(n)}(\mathbb{E}_{P(n)}[\phi_n^2])} \quad (\text{Lemma 13}) \\
&= \mathbb{E}_{P(n)}[(\hat{\phi}_n - \phi_n)^2] + \sqrt{\mathbb{E}_{P(n)}[(\hat{\phi}_n - \phi_n)^2] O_{P(n)}(\mathbb{E}_{P(n)}[\phi_n^2])} \\
&\quad + o_{P(n)}\left(\sqrt{\mathbb{E}_{P(n)}[\phi_n^2]}\right). \quad (\text{SLLN})
\end{aligned}$$

I decompose the nuisance error as:

$$\begin{aligned}
(\hat{\phi}_n - \phi_n)^2 &\leq \left( |\hat{\mu} - \mu| \left| 1 - \frac{D}{\max\{e, b_n\}} \right| + |Y - \hat{\mu}| \left| \frac{D}{\max\{\hat{e}, b_n\}} - \frac{D}{\max\{e, b_n\}} \right| \right)^2 \quad (\text{Tri. Ineq.}) \\
&\leq 4 \left( r_{\mu,n}^2 \left( 1 + \frac{D}{\max\{e, b_n\}^2} \right) + ((Y - \mu)^2 + o_{P(n)}(r_{\mu,n})) \left( \frac{D}{\max\{\hat{e}, b_n\}} - \frac{D}{\max\{e, b_n\}} \right)^2 \right) \\
&\lesssim r_{\mu,n}^2 \frac{D}{\max\{e, b_n\}^2} + \left( \frac{D}{\max\{\hat{e}, b_n\}} - \frac{D}{\max\{e, b_n\}} \right)^2 + o_{P(n)}(1)
\end{aligned}$$

$$= o_{P(n)}(1) \frac{D}{\max\{e, b_n\}^2} \mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e, b_n\}^2} \right] + \left( \frac{D}{\max\{\hat{e}, b_n\}} - \frac{D}{\max\{e, b_n\}} \right)^2 + o_{P(n)}(1)$$

( $r_{\mu, n} \rightarrow 0$ , Lemma 15)

$$\begin{aligned} \mathbb{E}_{P(n)} \left[ (\hat{\phi}_n - \phi_n)^2 \right] &= o_{P(n)} \left( \mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e, b_n\}^2} \right] \right) \\ &\quad + \mathbb{E}_{P(n)} \left[ \left( \frac{D}{\max\{\hat{e}, b_n\}} - \frac{D}{\max\{e, b_n\}} \right)^2 \right] + o_{P(n)}(1) \\ &= o_{P(n)} \left( \mathbb{E}_{P(n)} \left[ \frac{D}{\max\{e, b_n\}^2} \right] \right) + o_{P(n)}(1). \quad (\text{Lemma 14}) \\ &= o_{P(n)} \left( \text{Var}_{P(n)}(\phi_n) \right) + o_{P(n)}(1). \quad (\text{Lemma 5}) \end{aligned}$$

As a result:

$$\begin{aligned} \left| P(n)_n \left[ \hat{\phi}_n^2 - \phi_n^2 \right] \right| &= \sqrt{o_{P(n)} \left( \text{Var}_{P(n)}(\phi_n) \right) O_{P(n)} \left( \mathbb{E}_{P(n)}[\phi_n^2] \right) + o_{P(n)} \left( \text{Var}_{P(n)}(\phi_n) \right) + o_{P(n)} \left( \sqrt{\mathbb{E}_{P(n)}[\phi_n^2]} \right)} \\ &= o_{P(n)} \left( \text{Var}_{P(n)}(\phi_n) + \sqrt{\text{Var}_{P(n)}(\phi_n)} \right) = o_{P(n)} \left( \text{Var}_{P(n)}(\phi_n) \right). \quad (\text{Lemma 15}) \end{aligned}$$

□

**Lemma 17** (Estimated variance consistency). *Suppose the conditions of Proposition 2 and Lemma 14 hold. Let  $P(n)$  be a sequence of distributions in  $\mathcal{P}$ . Then  $\hat{\sigma}_n^2 / \sigma_n^2 \xrightarrow{\mathcal{P}} 1$ .*

*Proof of Lemma 17.* The conditions of Proposition 2 hold, so the conditions of Proposition 3 hold as well. The conditions of Lemma 14 hold, so the conditions of Lemma 16 hold as well.

Recall the definition  $\bar{\sigma}_n^2 = n^{-1} \text{Var}_{P(n)}(\phi_n)$ .

Let  $\sigma_n^2 = n^{-1} (P(n)_n[\phi_n^2] - P(n)_n[\phi_n]^2)$  be the oracle sample variance. By Lemma 13,  $\sigma_n^2 / \bar{\sigma}_n^2 \xrightarrow{\mathcal{P}} 1$ . Therefore it suffices to show that  $(\hat{\sigma}_n^2 - \sigma_n^2) / \bar{\sigma}_n^2 = \frac{P(n)_n[\hat{\phi}_n^2] - P(n)_n[\phi_n^2] - \hat{\psi}_{clip}^{AIPW}(b_n)^2 + \tilde{\psi}_{(Oracle)}^{AIPW}(b_n)^2}{\text{Var}_{P(n)}(\phi_n)} \xrightarrow{\mathcal{P}} 0$ .

Note that by Proposition 3,  $\mathbb{E}_{P(n)} [D / \max\{e(X), b_n\}^2] = \Theta(n\bar{\sigma}_n^2) = \Theta_{P(n)} (\text{Var}_{P(n)}(\phi_n))$ .

By the triangle inequality:

$$\begin{aligned} \left| \frac{\hat{\sigma}_n^2 - \sigma_n^2}{\bar{\sigma}_n^2} \right| &\leq \left| \frac{P(n)_n [\hat{\phi}_n^2 - \phi_n^2]}{\text{Var}_{P(n)}(\phi_n)} \right| + \left| \frac{\hat{\psi}_{clip}^{AIPW}(b_n)^2 - \tilde{\psi}_{(Oracle)}^{AIPW}(b_n)^2}{\text{Var}_{P(n)}(\phi_n)} \right| \\ &\lesssim \left| P(n)_n [\hat{\phi}_n^2 - \phi_n^2] \right| O \left( \frac{1}{\mathbb{E}_{P(n)} [D / \max\{e(X), b_n\}^2]} \right) \quad (\text{Lemma 6}) \\ &\quad + O_{P(n)} \left( \left| \hat{\psi}_{clip}^{AIPW}(b_n) - \tilde{\psi}_{(Oracle)}^{AIPW}(b_n) \right| \right) \quad (\text{Lemma 5}) \\ &= \left| P(n)_n [\hat{\phi}_n^2 - \phi_n^2] \right| O_{P(n)} \left( \frac{1}{\mathbb{E}_{P(n)} [D / \max\{e(X), b_n\}^2]} \right) + o_{P(n)}(1) \quad (\text{Proposition 2}) \\ &= \left| P(n)_n [\hat{\phi}_n^2 - \phi_n^2] \right| O_{P(n)} (1 / \text{Var}_{P(n)}(\phi_n)) + o_{P(n)}(1) \quad (\text{Proposition 3}) \end{aligned}$$

$$\begin{aligned}
&= o_{P(n)}(\text{Var}_{P(n)}(\phi_n)) O_{P(n)}(1/\text{Var}_{P(n)}(\phi_n)) + o_{P(n)}(1) \\
&= o_{P(n)}(1).
\end{aligned} \tag{Lemma 16}$$

Therefore  $(\hat{\sigma}_n^2 - \sigma_n^2)/\bar{\sigma}_n^2 \rightarrow_{P(n)} 0$ . By Lemma 13,  $\sigma_n^2/\bar{\sigma}_n^2 \rightarrow_{P(n)} 1$ . As a result,  $(\hat{\sigma}_n^2 - \bar{\sigma}_n^2)/\bar{\sigma}_n^2 \rightarrow_{P(n)} 0$  and  $\hat{\sigma}_n^2/\bar{\sigma}_n^2 \xrightarrow{\mathcal{P}} 1$ .  $\square$

**Lemma 18.** *Suppose the conditions of Theorem 1' hold. Then  $\sigma_n^{-1} \left( \hat{\psi}_{clip}^{AIPW}(b_n) - \tilde{\psi}_{(Oracle)}^{AIPW}(b_n) \right) = o_{P(n)}(1)$ .*

*Proof of Lemma 18.* I write  $k(i)$  for observation  $i$ 's fold and  $n_k$  for the number of observations in fold  $k$ . Then the oracle and clipped AIPW estimators are:

$$\begin{aligned}
\tilde{\psi}_{(Oracle)}^{AIPW}(b_n) &= \frac{1}{n} \sum_{i=1}^n \phi(Z_i | b_n, \eta) = \sum_k \frac{n_k}{n} \frac{1}{n_k} \overbrace{\sum_{i:k(i)=k} \phi(Z_i | b_n, \eta)}^{“\tilde{\psi}_{clip}^{AIPW,(k)}(b_n)”} \\
\hat{\psi}_{clip}^{AIPW}(b_n) &= \frac{1}{n} \sum_{i=1}^n \phi(Z_i | b_n, \hat{\eta}^{(-k)}) = \sum_k \frac{n_k}{n} \frac{1}{n_k} \underbrace{\sum_{i:k(i)=k} \phi(Z_i | b_n, \hat{\eta}^{(-k)})}_{“\hat{\psi}_{clip}^{AIPW,(k)}(b_n)”}.
\end{aligned}$$

I write  $\hat{r}_k \equiv \sigma_n^{-1} \left( \tilde{\psi}_{clip}^{AIPW,(k)}(b_n) - \hat{\psi}_{clip}^{AIPW,(k)}(b_n) \right)$ . I wish to show that  $\sum_k \frac{n_k}{n} \hat{r}_k = o_{P(n)}(1)$ . I consider an arbitrary  $k$  and quantify the bias and variance of  $\hat{r}_k$  given the data and nuisance estimates from the other folds  $-k$ .

I write the standard decomposition:

$$\hat{r}_k = \sigma_n^{-1} P_n \left[ (\hat{\mu} - \mu) \left( 1 - \frac{D}{\max\{\hat{e}, b_n\}} \right) + (\mu - Y) \left( \frac{D}{\max\{e, b_n\}} - \frac{D}{\max\{\hat{e}, b_n\}} \right) \right].$$

By cross-fitting, the bias satisfies  $\mathbb{E}[\hat{r}_k | \hat{\eta}^{(-k)}] = \sigma_n^{-1} \mathbb{E} \left[ (\hat{\mu} - \mu) \frac{\max\{\hat{e}, b_n\} - e}{\max\{\hat{e}, b_n\}} \right]$ . I now bound this term.

$$\begin{aligned}
\left| \mathbb{E}[\hat{r}_k | \hat{e}^{(-k)}, \hat{\mu}^{(-k)}] \right| &\leq \sigma_n^{-1} r_{\mu,n} P(n)(e(X) \leq b_n + r_{e,n}) \\
&\quad + \sigma_n^{-1} r_{\mu,n} r_{e,n} \mathbb{E}_{P(n)} \left[ \frac{1}{e - r_{e,n}} \mathbf{1}\{e > b_n + r_{e,n}\} \right] \tag{Lemma 8} \\
&\leq o_{P(n)}(1) + n^{1/2} r_{\mu,n} r_{e,n} \frac{\mathbb{E}[D/\max\{e, b_n\}^2]}{\sqrt{\mathbb{E}[D/\max\{e, b_n\}^2]}} \tag{Lemma 11 + Proposition 3} \\
&= o_{P(n)}(1). \tag{Assumption 3'(b)}
\end{aligned}$$

Next, I show that  $\text{Var}_{P(n)} \left( \hat{\psi}_{clip}^{AIPW}(b_n) - \tilde{\psi}_{(Oracle)}^{AIPW}(b_n) \right) = o_{P(n)}(\bar{\sigma}_n^2)$ , where  $\bar{\sigma}_n^2 = n^{-1} \text{Var}_{P(n)}(\phi_n)$ . . Consider the estimates in one fold  $k$ , with the nuisances from other folds fixed. For the estimates in that

fold:

$$\begin{aligned}
Var_{P(n)} \left( \hat{\psi}_{clip}^{AIPW}(b_n) - \tilde{\psi}_{(Oracle)}^{AIPW}(b_n) \right) &= Var \left( P(n)_n \hat{\phi} - \phi \right) = n^{-1} \mathbb{E}_{P(n)} \left[ (\hat{\phi} - \phi)^2 - \mathbb{E}_{P(n)} E[\hat{\phi} - \phi]^2 \right] \\
&= n^{-1} \mathbb{E}_{P(n)} \left[ (\hat{\phi} - \phi)^2 - o_{P(n)}(1) \right] \quad (\text{Bias}) \\
&= n^{-1} o_{P(n)} \left( Var_{P(n)}(\phi_n) \right) + o_{P(n)}(1) \quad (\text{Lemma 16}) \\
&= n^{-1} o_{P(n)} \left( \mathbb{E}_{P(n)} [D / \max\{e(X), b_n\}^2]^2 \right) + o_{P(n)}(1) \\
&\quad (\text{Lemma 3.(iii)}) \\
&= o_{P(n)}(\bar{\sigma}_n^2 + 1) \quad (\text{Lemma 6}) \\
&= o_{P(n)} \left( \frac{E_{P(n)}[D / \max\{e(X), b_n\}^2]}{n} + 1 \right) \quad (\text{Proposition 3}) \\
&= o_{P(n)}(\bar{\sigma}_n^2). \quad (b_n \gg n^{-1/2})
\end{aligned}$$

As a result:

$$\begin{aligned}
E \left[ \hat{r}_k^2 \mid \hat{\eta}^{(-k)} \right] &= E \left[ \hat{r}_k \mid \hat{\eta}^{(-k)} \right]^2 + Var \left( \hat{r}_k \mid \hat{\eta}^{(-k)} \right) = o_{P(n)}(1) \\
\hat{r}_k &= o_{P(n)}(1) \\
\left| \sigma_n^{-1} \left( \hat{\psi}_{clip}^{AIPW}(b_n) - \tilde{\psi}_{(Oracle)}^{AIPW}(b_n) \right) \right| &= o_{P(n)}(1) \leq \sum_k \frac{n_k}{n} |\hat{r}_k| = \sum_k \frac{n_k}{n} o_{P(n)}(1) = o_{P(n)}(1).
\end{aligned}$$

Finally, I am ready to prove the central claims of this work.  $\square$

*Proof of Theorem 1'.* By Lemma 18,  $\sigma_n^{-1} \left( \hat{\psi}_{clip}^{AIPW}(b_n) - \tilde{\psi}_{(Oracle)}^{AIPW}(b_n) \right) = o_{P(n)}(1)$ . Therefore, by Proposition 4,  $\sigma_n^{-1} \left( \hat{\psi}_{clip}^{AIPW}(b_n) - \psi_n \right) = \sigma_n^{-1} \left( \hat{\psi}_{clip}^{AIPW}(b_n) - \tilde{\psi}_{(Oracle)}^{AIPW}(b_n) \right) + \sigma_n^{-1} \left( \tilde{\psi}_{(Oracle)}^{AIPW}(b_n) - \psi_n \right) \overset{P(n)}{\rightsquigarrow} N(0, 1)$ .  $\square$

*Proof of Theorem 1.* For either claim, let  $P(n)$  be a sequence of distributions  $P$  in the relevant set. Note that in either case, the assumptions of Theorem 1' hold by Corollary 6. Therefore, by Theorem 1',

$$\sup_{t \in \mathbb{R}} \left| P(n)_n \left( \frac{\hat{\psi}_{clip}^{AIPW}(b_n) - \psi_n}{\sigma_n} \leq t \right) - \Phi(t) \right| \rightarrow 0,$$

where  $P(n)_n$  denotes the empirical average under distribution  $P(n)$  and  $\sigma_n$  is defined in Proposition 4.

Now I expand the empirical t-statistics for any fixed  $t$ :

$$P(n)_n \left( \frac{\hat{\psi}_{clip}^{AIPW}(b_n) - \psi_n}{\hat{\sigma}_n} \leq t \right) = P(n)_n \left( \frac{\hat{\psi}_{clip}^{AIPW}(b_n) - \psi_n}{\sigma_n} \left( \frac{\sigma_n}{\hat{\sigma}_n} \right) \leq t \right)$$

$$\begin{aligned}
&= P(n)_n \left( \frac{\hat{\psi}_{clip}^{AIPW}(b_n) - \psi_n}{\sigma_n} \leq t \frac{\hat{\sigma}_n}{\sigma_n} \right) \\
&= P(n)_n \left( \frac{\hat{\psi}_{clip}^{AIPW}(b_n) - \psi_n}{\sigma_n} - t \frac{\hat{\sigma}_n}{\sigma_n} \leq 0 \right) \\
&\rightarrow \Phi(t),
\end{aligned}$$

with the final result holding by Slutsky's theorem.

Therefore,

$$\sup_{t \in \mathbb{R}} \left| P(n)_n \left( \frac{\hat{\psi}_{clip}^{AIPW}(b_n) - \psi_n}{\hat{\sigma}_n} \leq t \right) - \Phi(t) \right| \rightarrow 0,$$

by properties of a cumulative distribution function.  $\square$

*Proof of Corollary 1.* For simplicity of exposition, I prove the result for the class  $\mathcal{P}$  under Assumption 4(ii):

$$\begin{aligned}
\limsup_{n \rightarrow \infty} \sup_{P \in \mathcal{P}} \left| P(\psi(P) \in \hat{\mathcal{C}}_n) - (1 - \alpha) \right| &= \limsup_{n \rightarrow \infty} \sup_{P \in \mathcal{P}} \left| P^n \left( \frac{\psi(P) - \hat{\psi}_{clip}^{AIPW}(b_n)}{\hat{\sigma}_n} \in [z_{\frac{\alpha}{2}}, z_{1-\frac{\alpha}{2}}] \right) - (1 - \alpha) \right| \\
&= \limsup_{n \rightarrow \infty} \sup_{P \in \mathcal{P}} \left| P^n \left( \frac{\psi(P) - \hat{\psi}_{clip}^{AIPW}(b_n)}{\hat{\sigma}_n} > z_{1-\frac{\alpha}{2}} \right) - \frac{\alpha}{2} \right| \\
&\quad - \left| P^n \left( \frac{\psi(P) - \hat{\psi}_{clip}^{AIPW}(b_n)}{\hat{\sigma}_n} > z_{\frac{\alpha}{2}} \right) - (1 - \frac{\alpha}{2}) \right| \\
&= \limsup_{n \rightarrow \infty} \sup_{P \in \mathcal{P}} \left| P^n \left( \frac{\psi(P) - \hat{\psi}_{clip}^{AIPW}(b_n)}{\hat{\sigma}_n} < z_{1-\frac{\alpha}{2}} \right) - (1 - \frac{\alpha}{2}) \right| \\
&\quad - \left| P^n \left( \frac{\psi(P) - \hat{\psi}_{clip}^{AIPW}(b_n)}{\hat{\sigma}_n} < z_{\frac{\alpha}{2}} \right) - \frac{\alpha}{2} \right| \\
&\leq \limsup_{n \rightarrow \infty} \sup_{P \in \mathcal{P}} \left| P^n \left( \frac{\hat{\psi}_{clip}^{AIPW}(b_n) - \psi(P)}{\hat{\sigma}_n} < z_{1-\frac{\alpha}{2}} \right) - \Phi(z_{1-\frac{\alpha}{2}}) \right| \\
&\quad + \limsup_{n \rightarrow \infty} \sup_{P \in \mathcal{P}} \left| P^n \left( \frac{\hat{\psi}_{clip}^{AIPW}(b_n) - \psi(P)}{\hat{\sigma}_n} < z_{\frac{\alpha}{2}} \right) - \Phi(z_{\frac{\alpha}{2}}) \right| \\
&= 0. \tag{Theorem 1}
\end{aligned}$$

$\square$

*Proof of Corollary 2.* By construction, there is a sequence of  $b_n \rightarrow 0$  such that  $r_{e,n} \ll b_n$ ,  $r_{e,n} \ll b_n \ll (n^{-1/2}/r_{\mu,n})^{2 \cdot \min\{1/\gamma_0, \gamma_0/(4(\gamma_0-1))\}}$ . For such a  $b_n$ , Assumption 2(b) holds by  $r_{\mu,n}r_{e,n} \ll n^{-1/2}$  and  $\gamma_0 > 2$ . Further, Assumption 4(ii) holds because  $(\gamma_0 - 1)2/\gamma_0 > 1$ .

Recall the definition of the oracle clipped AIPW estimator  $\tilde{\psi}_{(Oracle)}^{AIPW}(b_n)$  and the feasible estimator  $\hat{\psi}_{clip}^{AIPW}(b_n)$ .

It is clear by the proof of Proposition 1(i) that  $\tilde{\psi}_{(Oracle)}^{AIPW}(0)$  achieves the finite semiparametric efficiency



bound. I claim that  $\hat{\psi}_{clip}^{AIPW}(b_n) - \tilde{\psi}_{(Oracle)}^{AIPW}(0) = o_{P(n)}(n^{-1/2})$ , so that the remaining claims will hold by the proofs of Theorem 1 and Corollary 1. By the proof of Theorem 1,  $\hat{\psi}_{clip}^{AIPW}(b_n) - \tilde{\psi}_{(Oracle)}^{AIPW}(b_n) = o_{P(n)}(n^{-1/2})$ , so it only remains to show that  $\tilde{\psi}_{(Oracle)}^{AIPW}(b_n) - \tilde{\psi}_{(Oracle)}^{AIPW}(0) = o_{P(n)}(n^{-1/2})$ .

By construction,  $E_{P(n)} \left[ \tilde{\psi}_{(Oracle)}^{AIPW}(b_n) - \tilde{\psi}_{(Oracle)}^{AIPW}(0) \right] = 0$ , so it suffices to show that:

$$\begin{aligned}
Var_{P(n)} \left( n^{1/2} \left( \tilde{\psi}_{(Oracle)}^{AIPW}(b_n) - \tilde{\psi}_{(Oracle)}^{AIPW}(0) \right) \right) &= Var_{P(n)} \left( n^{-1/2} \sum \frac{D(Y - \mu(X))}{\max\{e(X), b_n\}} - \frac{D(Y - \mu(X))}{e(X)} \right) \\
&= Var_{P(n)} \left( n^{-1/2} \sum \frac{D(Y - \mu(X))(e(X) - b_n) \mathbf{1}\{e(X) \leq b_n\}}{b_n e(X)} \right) \\
&= Var_{P(n)} \left( \frac{D(Y - \mu(X))(e(X) - b_n) \mathbf{1}\{e(X) \leq b_n\}}{b_n e(X)} \right) \\
&= E_{P(n)} \left[ \left( \frac{D(Y - \mu(X))(e(X) - b_n) \mathbf{1}\{e(X) \leq b_n\}}{b_n e(X)} \right)^2 \right] \\
&= E_{P(n)} \left[ \frac{(Y - \mu(X))^2 (b_n - \max\{e(X), b_n\})^2}{b_n^2 e(X)} \right] \\
&\leq \sigma_{\max}^2 E_{P(n)} \left[ \frac{\mathbf{1}\{e(X) \leq b_n\}}{e(X)} \right] \\
&= \sigma_{\max}^2 \int_0^\infty P(n) \left( \frac{\mathbf{1}\{e(X) \leq b_n\}}{e(X)} > t \right) dt \\
&= \sigma_{\max}^2 \int_0^\infty P(n) (e(X) \leq \min\{b_n, 1/t\}) dt \\
&= \sigma_{\max}^2 \left( \int_0^{1/b_n} P(n) (e(X) \leq \min\{b_n, 1/t\}) dt + \int_{1/b_n}^\infty P(n) (e(X) \leq \min\{b_n, 1/t\}) dt \right) \\
&= \sigma_{\max}^2 \left( P(n) (e(X) \leq b_n) (1/b_n) + \int_{1/b_n}^\infty P(n) (e(X) \leq 1/t) dt \right) \\
&\leq C \sigma_{\max}^2 \left( b_n^{\gamma_0 - 2} + \int_{1/b_n}^\infty t^{1 - \gamma_0} dt \right) \\
&= C \sigma_{\max}^2 \left( \frac{\gamma_0 - 1}{\gamma_0 - 2} \right) \\
&= o_{P(n)}(1).
\end{aligned}$$

□

## C.5 Degradation of Black-Box Nuisance Requirements

*Proof of Corollary 4.* Take  $P$  to be the distribution that draws  $e(X)$  from the CDF  $P(e(X) \leq \pi) = \pi^{\gamma_0 - 1}$  and draws  $Y \mid X, D \sim \mathcal{N}(0, 1)$ . let  $\mathcal{P} = \{P\}$ . Such a family confirms to the requirements of Assumption 1 by construction.

Because  $\gamma_0 < 2$ , there exists  $r_{\mu,n}, r_{e,n} \rightarrow 0$  such that  $r_{\mu,n}r_{e,n} \ll r_{\mu,n}b_n \ll n^{-1/2} \ll r_{\mu,n}r_{e,n}^{\gamma_0/2} \ll r_{\mu,n}b_n^{\gamma_0/2}$ . Take these sequences. Note by construction that  $r_{\mu,n} \gg n^{-1/2}$ .

Take the nuisance estimates as  $\hat{\mu}_n(X) = \mu(X) - r_{\mu,n}$  and  $\hat{e}_n(X) = e(X) + r_{e,n}$ . These estimators do not depend on the data, so that they are cross-fit estimators. Therefore these nuisance estimates satisfy the second condition.

It remains to show that for any fixed  $\alpha$  and sequence of  $b_n \gg n^{-1/2}$ , Wald confidence intervals have the zero-coverage property. I show that every subsequence has a subsubsequence with the zero-coverage property.

Pick some subsequence of  $n$ . I first show that the estimator's bias is at least of at least the order  $r_{\mu,n}r_{e,n}^{\gamma_0-1}$ .

Write  $m_n = \max\{b_n - r_{e,n}, 0\}$ . The bias of the clipped estimator  $\hat{\psi}_{clip}^{AIPW}(b_n)$  with  $n$  observations is:

$$\begin{aligned} E_P \left[ \hat{\psi}_{clip}^{AIPW}(b_n) - \psi(P) \right] &= E \left[ (\hat{\mu}(X) - \mu(X)) \left( \frac{D}{\max\{\hat{e}(X), b_n\}} - 1 \right) \right] = r_{\mu,n} E \left[ 1 - \frac{D}{\max\{e(X) + r_{e,n}, b_n\}} \right] \\ &= \frac{r_{\mu,n}}{b_n} E[(b_n - e(X)) \mathbf{1}\{e(X) \leq b_n - r_{e,n}\}] + r_{\mu,n} E \left[ \frac{r_{e,n} \mathbf{1}\{e(X) \geq b_n - r_{e,n}\}}{e(X) + r_{e,n}} \right] \\ &= \frac{r_{\mu,n}(\gamma_0 - 1)}{b_n} \int_0^{m_n} (b_n - \pi) \pi^{\gamma_0-2} d\pi + r_{\mu,n} r_{e,n} (\gamma_0 - 1) \int_{m_n}^1 (\pi + r_{e,n})^{-1} \pi^{\gamma_0-2} d\pi. \end{aligned}$$

By assumption,  $b_n \gtrsim r_{e,n}$ , so that

$$\begin{aligned} E_P \left[ \hat{\psi}_{clip}^{AIPW}(b_n) - \psi(P) \right] &\geq \frac{r_{\mu,n}r_{e,n}(\gamma_0 - 1)}{2} \int_{b_n+r_{e,n}}^1 \pi^{\gamma_0-3} d\pi \\ &\geq \frac{r_{\mu,n}(b_n - r_{e,n})^{\gamma_0}}{b_n \gamma_0} + \frac{r_{\mu,n}r_{e,n}(\gamma_0 - 1)}{2(2 - \gamma_0)} (b_n + r_{e,n})^{\gamma_0-2} - \overbrace{\frac{r_{\mu,n}r_{e,n}(\gamma_0 - 1)}{2(2 - \gamma_0)}}^{o(1)} \\ &\geq \frac{(\gamma_0 - 1)2^{\gamma_0-3}}{2 - \gamma_0} r_{\mu,n}r_{e,n}b_n^{\gamma_0-2} + \frac{r_{\mu,n}(b_n - r_{e,n})^{\gamma_0}}{b_n \gamma_0} - o(1) \gtrsim r_{\mu,n}r_{e,n}^{\gamma_0-1}. \end{aligned}$$

Either way, there is a subsubsequence with the desired bias property. Therefore the original subsequence has the desired bias property.

It is convenient to write  $B_n = \frac{E_P[\hat{\psi}_{clip}^{AIPW}(b_n) - \psi(P)]}{\sigma_n}$  for this proof.

By the proof of Corollary 3, there is a  $C^{-1} > 0$  such that  $\sigma_n \geq Cn^{-1/2}b_n^{\gamma_0/2-1}$  for all  $n$  large enough.

For such  $n$ :

$$B_n = \frac{E_P \left[ \hat{\psi}_{clip}^{AIPW}(b_n) - \psi(P) \right]}{\sigma_n} \geq C2^{\gamma_0-2}n^{1/2}r_{\mu,n}b_n^{\gamma_0-1}b_n^{1-\gamma_0/2} = C2^{\gamma_0-2}n^{1/2}r_{\mu,n}b_n^{\gamma_0/2} \geq C2^{\gamma_0-2}n^{1/2}r_{\mu,n}r_{e,n}^{\gamma_0/2} \rightarrow \infty.$$

Note also that:

$$\begin{aligned}
\text{Var}\left(\hat{\psi}_{clip}^{AIPW}(b_n)\right) &= \frac{1}{n} \text{Var}\left(\mu + \frac{D}{\max\{\hat{e}, b_n\}}(Y - \mu) + r_{\mu,n} \left(1 - \frac{D}{\max\{b_n, e + r_{e,n}\}}\right)\right) \\
&= \frac{1}{n} \text{Var}\left(\mu + \frac{D}{\max\{e, b_n\}}(Y - \mu)\right) + \frac{1}{n} E\left[\text{Var}\left(\left(\frac{D}{\max\{\hat{e}, b_n\}} - \frac{D}{\max\{e, b_n\}}\right)(Y - \mu) \mid X, D\right)\right] \\
&\quad + \frac{r_{\mu,n}^2}{n} E\left[\text{Var}\left(1 - \frac{D}{\max\{b_n, e + r_{e,n}\}} \mid X\right)\right] + \frac{r_{\mu,n}^2}{n} \text{Var}\left(E\left[1 - \frac{D}{\max\{b_n, e + r_{e,n}\}} \mid X\right]\right) \\
&= \sigma_n^2 + \frac{1}{n} E\left[e(X) \text{Var}(Y \mid X, D = 1) \frac{\mathbf{1}\{e \geq b_n - r_{e,n}\}(e + r_{e,n} - \max\{e, b_n\})^2}{\max\{e, b_n\} \max\{e + r_{e,n}, b_n\}}\right] \\
&\quad + \frac{r_{\mu,n}^2}{n} E\left[\frac{e(1-e)}{\max\{b_n, e + r_{e,n}\}^2}\right] + \frac{r_{\mu,n}^2}{n} \text{Var}\left(\frac{\max\{b_n, e + r_{e,n}\} - e}{\max\{b_n, e + r_{e,n}\}}\right) \\
&\leq \sigma_n^2 + \frac{\sigma_{\max}^2 r_{e,n}^2}{n} E\left[e \frac{\mathbf{1}\{b_n - r_{e,n} \leq e \leq b_n\}}{b_n^2} + e \frac{\mathbf{1}\{b_n \leq e\}}{e^2}\right] \\
&\quad + \frac{r_{\mu,n}^2}{n} E\left[\frac{e(1-e) + b_n^2 \mathbf{1}\{e \leq b_n - r_{e,n}\} + r_{e,n}^2 \mathbf{1}\{e \geq b_n - r_{e,n}\}}{\max\{b_n, e + r_{e,n}\}^2}\right] \\
&\leq \sigma_n^2 + \frac{\sigma_{\max}^2 r_{e,n}^2}{n} E\left[\frac{\mathbf{1}\{e \leq b_n\}}{b_n} + \frac{\mathbf{1}\{e \geq b_n\}}{e}\right] \\
&\quad + \frac{r_{\mu,n}^2}{n} E\left[\frac{(e + b_n^2) \mathbf{1}\{e \leq b_n - r_{e,n}\}}{b_n^2}\right] + \frac{r_{\mu,n}^2}{n} E\left[\frac{(e + r_{e,n}^2) \mathbf{1}\{e \geq b_n - r_{e,n}\}}{\max\{b_n, e + r_{e,n}\}^2}\right] \\
&\leq \sigma_n^2 + \frac{\sigma_{\max}^2 r_{e,n}^2}{n} + \frac{r_{\mu,n}^2}{n} b_n^{-1} P(e \leq b_n - r_{e,n}) + \frac{r_{\mu,n}^2}{n} E\left[\frac{\mathbf{1}\{b_n - r_{e,n} \leq e \leq b_n\}}{b_n}\right] + \frac{r_{\mu,n}^2}{n} E\left[\frac{\mathbf{1}\{e \geq b_n\}}{e}\right] \\
&\leq \sigma_n^2 + \frac{\sigma_{\max}^2 r_{e,n}^2 + r_{\mu,n}^2 \left(2b_n^{\gamma_0-2} + \int_{b_n}^1 \frac{\pi^{\gamma_0-3}}{\gamma_0-1} d\pi\right)}{n} \leq \sigma_n^2 + \frac{\sigma_{\max}^2 r_{e,n}^2 + \frac{2(2-\gamma_0)(\gamma_0-1)+1}{(2-\gamma_0)(\gamma_0-1)} r_{\mu,n}^2 b_n^{\gamma_0-2}}{n}.
\end{aligned}$$

By the proof of Corollary 3 applied to the distribution  $P$ ,  $\sigma_n^2$  is on the order of  $n^{-1}b_n^{\gamma_0-2}$ . Therefore  $\text{Var}\left(\hat{\psi}_{clip}^{AIPW}(b_n)\right) \leq \sigma_n^2 + o(\sigma_n^2)$ .

Note that since  $r_{\mu,n} \rightarrow 0$  and  $b_n \geq 2r_{e,n}$  for all  $n$  large enough, the conditions of Lemma 17 hold. Therefore  $\hat{\sigma}_n^2 / \text{Var}(\hat{\psi}_{clip}^{AIPW}(b_n)) \rightarrow^P 1$  and  $E\left[\hat{\psi}_{clip}^{AIPW}(b_n) - \psi(P)\right] \rightarrow^P \infty$ . As a result,  $\frac{E_P[\hat{\psi}_{clip}^{AIPW}(b_n) - \psi(P)]}{\hat{\sigma}_n} \rightarrow^P \infty$  and for any fixed  $\alpha \in (0, 1/2)$ ,

$$\begin{aligned}
P(\psi(P) \in \hat{\mathcal{C}}_n) &= P\left(\frac{\hat{\psi}_{clip}^{AIPW}(b_n) - \psi(P)}{\hat{\sigma}_n} \in [z_{\alpha/2}, z_{1-\alpha/2}]\right) \\
&= P\left(\frac{\hat{\psi}_{clip}^{AIPW}(b_n) - E\left[\hat{\psi}_{clip}^{AIPW}(b_n)\right]}{\sigma_n + o_P(\sigma_n)} \in [B_n + z_{\alpha/2} + o_P(1), B_n + z_{1-\alpha/2} + o_P(1)]\right) \\
&= P\left(\frac{O_P(\sigma_n)}{\sigma_n + o_P(\sigma_n)} \in [B_n + z_{\alpha/2} + o_P(1), B_n + z_{1-\alpha/2} + o_P(1)]\right) \rightarrow^P 0,
\end{aligned}$$

with the limit holding because  $B_n$  tends to infinity.  $\square$

*Proof of Example 1.* For simplicity, I proceed assuming  $r_{e,n} \gg n^{-1/2}$ . By Theorem 1, it remains to show

that there is a  $b_n \rightarrow 0$  such that 3(c) ( $r_{\mu,n}b_n^{\gamma_0/2} \ll n^{-1/2}$ ) and 3(d) ( $r_{e,n} \ll b_n$ ).

Because  $r_{\mu,n}r_{e,n} \ll n^{-1/2}$ , there exists some  $\delta_n \rightarrow \infty$  such that  $r_{\mu,n}r_{e,n}\delta_n \ll n^{-1/2}$ . Choose some  $b_n$  such that  $1 \gg b_n$  and  $r_{e,n} \ll b_n \ll (r_{e,n}\delta_n)^{2/\gamma_0}$ . This is feasible because  $\gamma_0 > 2$ , so that  $r_{e,n}^{2/\gamma_0} \gg r_{e,n}$  and  $\delta_n^{2/\gamma_0} \rightarrow \infty$ . Then  $r_{e,n} \ll b_n$  and  $r_{\mu,n}b_n^{\gamma_0/2} \ll r_{\mu,n}r_{e,n}\delta_n \ll n^{-1/2}$ . Thus, both conditions hold.  $\square$

*Proof of Example 2.* For simplicity, I proceed assuming  $r_{e,n} \gg n^{-1/2}$ . By Theorem 1, it remains to show that there is a  $b_n \rightarrow 0$  such that  $r_{\mu,n}r_{e,n} \log(1/b_n) \ll n^{-1/2}$ ,  $r_{\mu,n}b_n \ll n^{-1/2}$ , and  $r_{e,n} \ll b_n$ .

Because  $r_{\mu,n}r_{e,n} \log(1/r_{e,n}) \ll n^{-1/2}$ , there exists a  $b_n$  such that  $r_{e,n} \ll b_n \ll 1$  and  $r_{\mu,n}b_n \log(1/b_n) \ll n^{-1/2}$ . For this  $b_n$ , all three conditions hold by inspection.  $\square$

*Proof of Example 3.* For simplicity, I proceed assuming  $r_{e,n}, r_{\mu,n} \gg n^{-1/2}$ .

If  $\gamma_0 \geq 2$ , the claim holds by Example 1 and Example 2.

Now suppose that  $\gamma_0 < 2$ . By Theorem 1, it remains to show that there is a  $b_n \rightarrow 0$  such that  $r_{\mu,n}r_{e,n}b_n^{(\gamma_0-2)/2} \ll n^{-1/2}$ ,  $r_{\mu,n}b_n^{\gamma_0/2} \ll n^{-1/2}$ , and  $r_{e,n} \ll b_n$ . Because  $r_{e,n} \ll n^{-1/3}$ , there exists a  $b_n$  such that  $n^{-1/3} \gg b_n \gg \max\{r_{e,n}, n^{-1}r_{\mu,n}^{-2}\}$ . For this  $b_n$ :

$$\begin{aligned} r_{\mu,n}r_{e,n}b_n^{(\gamma_0-2)/2} &\ll r_{\mu,n}r_{e,n}b_n^{-1/2} \ll r_{\mu,n}r_{e,n} (n^{-1}r_{\mu,n}^{-2})^{-1/2} = n^{1/2}r_{\mu,n}^2r_{e,n} \ll n^{-1/2} \\ r_{\mu,n}b_n^{\gamma_0/2} &\ll r_{\mu,n}b_n^{1/2} \ll n^{-1/2}, \end{aligned}$$

verifying all conditions.  $\square$

*Proof of Example 4.* If  $\gamma_0 \geq 2$ , the claim holds by Example 1 and Example 2.

Now suppose that  $\gamma_0 < 2$ . By Theorem 1, it remains to show that there is a  $b_n \rightarrow 0$  such that  $r_{\mu,n}r_{e,n}b_n^{(\gamma_0-2)/2} \ll n^{-1/2}$ ,  $r_{\mu,n}b_n^{\gamma_0/2} \ll n^{-1/2}$ , and  $r_{e,n} \ll b_n$ .

(i)  $r_{\mu,n} = O(n^{-1/2})$ . Take  $b_n \rightarrow 0$  such that  $b_n \gg r_{e,n}$ . The first condition holds because  $b_n^{(\gamma_0-2)/2} \ll r_{e,n}^{(\gamma_0-2)/2} \ll r_{e,n}^{-1/2}$ . The second condition holds because  $r_{\mu,n}b_n^{\gamma_0/2} \ll n^{-1/2}$ .

(ii)  $r_{e,n} = O(n^{-1/2})$ . There is a  $\delta_n \rightarrow \infty$  such that  $\delta_n^{1/2} \ll n^{-1/4}/r_{\mu,n}$ . Take  $b_n$  such that  $n^{-1/2} \gg b_n \gg n^{-1/2}\delta_n$ , so that  $r_{e,n} \ll b_n$ . Then:

$$\begin{aligned} r_{\mu,n}r_{e,n}b_n^{(\gamma_0-2)/2} &\ll r_{\mu,n}r_{e,n}b_n^{-1/2} \ll r_{\mu,n}r_{e,n}n^{1/4}\delta_n^{-1/2} \ll n^{-1/2} \\ r_{\mu,n}b_n^{\gamma_0/2} &\ll r_{\mu,n}b_n^{1/2} \ll r_{\mu,n}n^{-1/4}\delta_n^{1/2} \ll n^{-1/2}, \end{aligned}$$

verifying all conditions.  $\square$

## C.6 Necessary Smoothness Conditions

**Lemma 19** (Minimal expected nearby observations). *Suppose the conditions of Theorem 2 hold. Define  $A(x_0 | h) = \{x : x \in [-1, 1]^d, \|x - x_0\| \leq h\}$ . Then*

$$\inf_{x_0 \in [-1, 1]^d, P \in \mathcal{P}, h > 0} E_P [D | X \in A(x_0 | h)] \geq 2^{-(\gamma_0+1)/(\gamma_0-1)} C^{-1/(\gamma_0-1)} h^{\frac{d}{\gamma_0-1}}.$$

*Proof of Lemma 19.* Proof by contradiction. Suppose not, and there is a  $P \in \mathcal{P}$ ,  $x_0 \in [-1, 1]^d$ ,  $h > 0$  with

$$E_P [D | X \in A(x_0 | h)] < 2^{-(\gamma_0+1)/(\gamma_0-1)} C^{-1/(\gamma_0-1)} h^{\frac{d}{\gamma_0-1}}.$$

Define  $\pi = 2^{-(\gamma_0+1)/(\gamma_0-1)} C^{-1/(\gamma_0-1)} h^{\frac{d}{\gamma_0-1}}$  and  $B = \{X \in A(x_0 | h) : e(X) \leq \pi\}$ . Then  $P(X \in B) < P(X \in A(x_0 | h))/2 \geq h^d/2$ . As a result:

$$P(e(X) \leq \pi) > \frac{h^d}{2} = C\pi^{\gamma_0-1} \frac{C^{-1}h^d\pi^{1-\gamma_0}}{2} = C\pi^{\gamma_0-1} C^{-1} 2^{-1-\gamma_0} h^d \left(2^{\frac{\gamma_0+1}{1-\gamma_0}} C^{\frac{-1}{\gamma_0-1}} h^{\frac{d}{\gamma_0-1}}\right)^{1-\gamma_0} = C\pi^{\gamma_0-1}.$$

Contradiction. □

**Lemma 20** (Minimal coefficient). *Suppose Assumption 6 holds and  $e(X) \in \Sigma(\beta_e, L)$  for some  $\beta_e > \frac{d}{\gamma_0-1}$ . For  $0 \leq j < \beta_e$ , Let  $c_j(v | x_0)$  be the coefficient in the  $j^{\text{th}}$ -order Taylor expansion of  $e(x)$  around  $x_0$  applied in the direction of  $v$ . Then there is a  $c^* > 0$  such that for all  $x_0$ , there exists an  $x \neq x_0$  and a  $0 \leq j \leq \alpha^{(Mou)}$  such that  $x_0 + \frac{x-x_0}{\|x-x_0\|} \in [-1, 1]^d$  and  $c_j\left(\frac{x-x_0}{\|x-x_0\|} | x_0\right) \geq c^*$ .*

*Proof of Lemma 20.* If  $\beta_e \leq 1$ , the claim is immediate by Lemma 19. I therefore proceed assuming  $\beta_e > 1$ .

Proof by contradiction. Suppose for all  $c > 0$ , there exists a  $P$  and an  $x_0$  such that  $c_j\left(\frac{x-x_0}{\|x-x_0\|} | x_0\right) < c$  for all  $j \leq \alpha^{(Mou)}$ .

Define:

$$h^* = \operatorname{argsup}_{h \in (0, 1]} \sum_{\alpha^{(Mou)} < j \leq \ell_e} \bar{c}_j(h)^j + L_e(h)^{\beta_e} \leq \frac{C^{1/d}}{\lceil \alpha^{(Mou)} + 2 \rceil} (h)^{\frac{d}{\gamma_0-1}}.$$

This  $h^*$  is well-defined because as  $h$  tends to zero from above, the left-hand side is of a lower order than the right-hand side. By continuity,  $h_n^* \in (0, 1]$  and satisfies this weak inequality.

Take  $c^* = \min_{0 \leq j \leq \alpha^{(Mou)}} \frac{C^{\frac{1}{1-\gamma_0}}}{\lceil \alpha^{(Mou)} + 2 \rceil} (h^*)^{\frac{d}{\gamma_0-1}-j}$ . I claim that for this  $c^*$ , the the content of Lemma 20 holds.

Proof by contradiction. Suppose, not, and there is an  $x_0 \in [-1, 1]^d$  and an  $x$  as above such that

$c_j \left( \frac{x-x_0}{\|x-x_0\|} \mid x_0 \right) < c^*$  for all  $j \leq \alpha^{(Mou)}$ . Then for all  $x \in [-1, 1]^d/x_0$  such that  $\|x - x_0\| \leq h^*$  and  $x_0 + \frac{x-x_0}{\|x-x_0\|} \in [-1, 1]^d$ :

$$\begin{aligned} e(x) &= f(x \mid x_0) + g(x \mid x_0) \\ &= \sum_{0 \leq j \leq \alpha^{(Mou)}} c_j \left( \frac{x-x_0}{\|x-x_0\|} \mid x_0 \right) \|x-x_0\|^j + \sum_{\alpha^{(Mou)} < j \leq \ell_e} c_j \left( \frac{x-x_0}{\|x-x_0\|} \mid x_0 \right) \|x-x_0\|^j + L \|x-x_0\|^{\beta_e} \\ &\leq \sum_{0 \leq j \leq \alpha^{(Mou)}} \frac{C^{\frac{1}{1-\gamma_0}}}{[\alpha^{(Mou)} + 2]} (h^*)^{\frac{d}{\gamma_0-1}-j} \|x-x_0\|^j + \frac{C^{\frac{1}{1-\gamma_0}}}{[\alpha^{(Mou)} + 2]} (h^*)^{\frac{d}{\gamma_0-1}} \leq C^{\frac{1}{1-\gamma_0}} \frac{[\alpha^{(Mou)} + 1]}{[\alpha^{(Mou)} + 2]} (h^*)^{\frac{d}{\gamma_0-1}}. \end{aligned}$$

So that:

$$\begin{aligned} P \left( e(X) \leq C^{\frac{1}{1-\gamma_0}} \frac{[\alpha^{(Mou)} + 1]}{[\alpha^{(Mou)} + 2]} (h^*)^{\frac{d}{\gamma_0-1}} \right) &\geq P \left( x_0 + \frac{X-x_0}{\|X-x_0\|} \in [-1, 1]^d, \|X-x_0\| \geq h^* \right) \\ &\geq (h^*)^d = C \left( C^{\frac{1}{1-\gamma_0}} (h^*)^{d/(\gamma_0-1)} \right)^{\gamma_0-1} \\ &= C \left( \frac{[\alpha^{(Mou)} + 2]}{[\alpha^{(Mou)} + 1]} \right)^{\gamma_0-1} \left( C^{\frac{1}{1-\gamma_0}} \frac{[\alpha^{(Mou)} + 1]}{[\alpha^{(Mou)} + 2]} (h^*)^{\frac{d}{\gamma_0-1}} \right)^{\gamma_0-1} \\ &> C \left( C^{\frac{1}{1-\gamma_0}} \frac{[\alpha^{(Mou)} + 1]}{[\alpha^{(Mou)} + 2]} (h^*)^{\frac{d}{\gamma_0-1}} \right)^{\gamma_0-1}. \end{aligned}$$

But by Assumption 1(d),  $P \left( e(X) \leq C^{\frac{1}{1-\gamma_0}} \frac{[\alpha^{(Mou)} + 1]}{[\alpha^{(Mou)} + 2]} (h^*)^{\frac{d}{\gamma_0-1}} \right) \leq C \left( C^{\frac{1}{1-\gamma_0}} \frac{[\alpha^{(Mou)} + 1]}{[\alpha^{(Mou)} + 2]} (h^*)^{\frac{d}{\gamma_0-1}} \right)^{\gamma_0-1}$ . Contradiction.

Therefore, for this  $c^* > 0$  and all  $x_0$ , there exists an  $x \neq x_0$  and a  $0 \leq j \leq \alpha^{(Mou)}$  such that  $c_j \left( \frac{x-x_0}{\|x-x_0\|} \mid x_0 \right) \geq c^*$ .  $\square$

*Proof of Proposition 5.* Proof by contradiction. Suppose not, and for all  $\rho, \gamma > 0$ , there is an

$$h \in \left( 0, (C^{1/d}/(2L))^{\frac{1}{\beta_e} - \frac{d}{\gamma_0-1}} \right],$$

$x_0 \in [-1, 1]^d$ , and  $P \in \mathcal{P}$  such that  $P(e(X) \geq \rho \sup_{\|x-x_0\| \leq h} e(x) \mid D=1, \|X-x_0\| \leq h) \leq \gamma$ .

Take some  $\rho > 0$  and some sequence of  $\gamma_n \rightarrow 0^+$ , with associated bandwidths  $h_n$ , such that

$$P(n) \left( e(X) \geq \rho \left( \sup_{\|x-x_0\| \leq h_n} e(x) \right) \mid D=1, \|X-x_0\| \leq h_n \right) \leq \gamma_n.$$

Because  $[-1, 1]^d$  is compact and the coefficients are in a compact space, there is a subsequence of  $n$  for which  $x_0$  and the local polynomial coefficients of  $e(x)$  around  $x_0$  are convergent. Without loss of generality, suppose this is the sequence. Write  $g_n(x \mid x_0) = P(n)(D=1 \mid X=x) - f_n(x \mid x_0)$  be the local polynomial propensity residuals. Also write  $f^*(\cdot \mid x_0)$  for the local polynomial coefficients at the limiting coefficients.

First, I show that  $h_n \rightarrow 0$ . Suppose not, and  $\limsup_{n \rightarrow \infty} h_n(x) = h^* > 0$ . Without loss of generality, suppose  $\lim_{n \rightarrow \infty} h_n(x) = h^*$ . Write  $A = \{x \mid \|x - x_0\| \leq h^*\}$ . Then by continuity of densities and bounds on derivatives, for all  $n$  large enough,

$$\begin{aligned}
\gamma_n &\geq P(n) \left( e(X) \geq \frac{\rho}{2} \left( \sup_{\|x-x_0\| \leq h^*} e(x) \right) \mid D=1, \|X-x_0\| \leq h^* \right) \\
&\geq P(n) \left( e(X) \geq \frac{3\rho}{8} \left( \sup_{\|x-x_0\| \leq h^*} f^*(x \mid x_0) \right) \mid D=1, \|X-x_0\| \leq h^* \right) \\
&\quad - 1 \left\{ \sup_{\|x-x_0\| \leq h^*} |e(x) - f^*(x \mid x_0)| \geq \frac{\rho}{8} \sup_{\|x-x_0\| \leq h^*} f^*(x \mid x_0) \right\} \\
&\geq P(n) \left( e(X) \geq \frac{3\rho}{8} \left( \sup_{\|x-x_0\| \leq h^*} f^*(x \mid x_0) \right) \mid D=1, \|X-x_0\| \leq h^* \right) \\
&\quad - 1 \left\{ \sup_{x \in A} \|f_n(x \mid x_0) - f^*(x \mid x_0)\| + |g_n(x \mid x_0)| \geq \frac{\rho}{16} \sup_{\|x-x_0\| \leq h^*} f^*(x \mid x_0) \right\} \\
&= P(n) \left( e(X) \geq \frac{3\rho}{8} \left( \sup_{\|x-x_0\| \leq h^*} f^*(x \mid x_0) \right) \mid D=1, \|X-x_0\| \leq h^* \right) - 1 \left\{ O((h^*)^{\beta_e}) \geq \Omega((h^*)^{\frac{-d}{\gamma_0-1}}) \right\} \\
&\geq P(n) \left( e(X) \geq \frac{\rho}{4} \left( \sup_{\|x-x_0\| \leq h^*} f^*(x \mid x_0) \right) \mid D=1, \|X-x_0\| \leq h^* \right),
\end{aligned}$$

which is a positive constant. Therefore  $\gamma_n \not\rightarrow 0$ . Contradiction.

I therefore proceed assuming  $h_n \rightarrow 0$ . Let the lowest-order nonzero coefficient in  $f^*$  be of order  $j^*$ .  $j^*$  is defined and finite by Lemma 20. Define

$$G_n = P(n) \left( e(X) \geq \frac{\rho}{2} \left( \sup_{\|x-x_0\| \leq h^*} e(x) \right) \mid D=1, \|X-x_0\| \leq h^* \right) - \gamma_n.$$

I wish to show that  $G_n$  does not converge to zero. By the Bolzano-Weierstrass Theorem, it suffices to show that there is a nonconvergent subsequence of  $n$ .

Write  $m_{n,j} \equiv \sup_{\|v\|=1} |\tilde{c}_j(v)| h_n^{j-j^*}$  and  $m_n \equiv \max_{0 \leq j \leq j^*} m_{n,j}$ . I consider two cases: (i)  $m_n \rightarrow 0$  or (ii)  $\liminf_{n \rightarrow \infty} m_n > 0$  (and potentially infinite).

Case (i): suppose  $m_n \rightarrow 0$ . Then for all  $x \in A_n$ ,

$$\begin{aligned}
e(x) &= f^*(x \mid x_0) + \tilde{f}_n(x \mid x_0) + g_n(x \mid x_0) \\
&= \sup_{\|\alpha\|=j^*} D^\alpha e(x_0) (x - x_0)^\alpha + O(h_n^{\min\{j^*+1, \beta_e\}}) + o(h_n^{j^*}) + O(h_n^{\beta_e}) \\
&= h_n^{j^*} c_{j^*} \left( \frac{x - x_0}{\|x - x_0\|} \right) \left( \frac{\|x - x_0\|}{h_n} \right)^{j^*} + o(h_n^{j^*}).
\end{aligned}$$

Let  $n \geq n'$  imply that the  $o(h_n^{j^*})$  term is at most half as large as  $\sup_{x \in A_n} h_n^{j^*} c_{j^*} \left( \frac{x-x_0}{\|x-x_0\|} \right) \left( \frac{\|x-x_0\|}{h_n} \right)^{j^*}$ , as well as to imply that  $x_0 + h_n v \in [-1, 1]^d$  if and only if there is an  $h > 0$  such that  $x_0 + hv \in [-1, 1]^d$ . Then:

$$\begin{aligned} & P(n) \left( e(X) \geq \sup_{x \in A_n} e(x) \mid D = 1, X \in A_n \right) \\ & \geq P(n) \left( h_n^{j^*} c_{j^*} \left( \frac{X-x_0}{\|X-x_0\|} \right) \left( \frac{\|X-x_0\|}{h_n} \right)^{j^*} \geq \frac{\rho}{2} \sup_{x \in A_n} h_n^{j^*} c_{j^*} \left( \frac{x-x_0}{\|x-x_0\|} \right) \left( \frac{\|x-x_0\|}{h_n} \right)^{j^*} \mid D = 1, X \in A_n \right) \\ & = P(n') \left( c_{j^*} \left( \frac{X-x_0}{\|X-x_0\|} \right) \left( \frac{\|X-x_0\|}{h_{n'}} \right)^{j^*} \geq \frac{\rho}{2} \sup_{x \in A_{n'}} c_{j^*} \left( \frac{x-x_0}{\|x-x_0\|} \right) \left( \frac{\|x-x_0\|}{h_{n'}} \right)^{j^*} \mid D = 1, X \in A_{n'} \right) > 0. \end{aligned}$$

Contradiction.

Case (ii): suppose  $\liminf_{n \rightarrow \infty} m_n > 0$ . Write

$$\tilde{f}_n(x \mid x_0) = h_n^{j^*} m_n \sum_{j=0}^{\ell_e} \underbrace{\frac{\sum_{\|\alpha\|=j} D^\alpha e(x_0) \left( \frac{x-x_0}{\|x-x_0\|} \right)^\alpha h_n^{j-j^*}}{m_n}}_{\tilde{d}_{j,n} \left( \frac{x-x_0}{\|x-x_0\|} \right)} \left( \frac{\|x-x_0\|}{h_n} \right)^j.$$

By construction,  $\tilde{d}_{j,n} \left( \frac{x-x_0}{\|x-x_0\|} \right)$  is bounded between  $-1$  and  $1$ , so that there is a convergent subsequence for all  $\alpha$ . Without loss of generality I proceed assuming this is the full sequence.

Write  $\tilde{d}_j^*(v) = \lim_{n \rightarrow \infty} \tilde{d}_{j,n}(v)$  for all  $v$ . Write  $\tilde{d}_{j^*}^*(v) = c_{j^*}(v)/m_n$ , where  $c_{j^*}(v)$  is the  $j^*$ -order coefficient in  $f^*$ . By construction, for all  $\|x - x_0\| \leq h_n$ ,

$$\begin{aligned} e(x \mid x_0) &= f^*(x \mid x_0) + \tilde{f}_n(x \mid x_0) + g_n(x \mid x_0) \\ &= m_n h_n^{j^*} \left( \sum_{j=0}^{j^*} \tilde{d}_j^* \left( \frac{x-x_0}{\|x-x_0\|} \right) \left( \frac{\|x-x_0\|}{h_n} \right)^j + o(1) + O \left( h_n^{\min\{1, \beta_e - j^*\}} \right) \right) \\ &= m_n h_n^{j^*} \sum_{j=0}^{j^*} \tilde{d}_j^* \left( \frac{x-x_0}{\|x-x_0\|} \right) \left( \frac{\|x-x_0\|}{h_n} \right)^j + o(m_n h_n^{j^*}). \end{aligned}$$

Let  $n \geq n'$  imply that the  $o(m_n h_n^{j^*})$  term is at most half as large as the largest value of the first term over  $x \in A_n$ , as well as that  $x_0 + h_n v \in [-1, 1]^d$  if and only if there is an  $h > 0$  such that  $x_0 + hv \in [-1, 1]^d$ .

Then:

$$\begin{aligned} & P(n) \left( e(X) \geq \sup_{x \in A_n} e(x) \mid D = 1, X \in A_n \right) \\ & \geq P(n) \left( m_n h_n^{j^*} \sum_{j=0}^{j^*} \tilde{d}_j^* \left( \frac{X-x_0}{\|X-x_0\|} \right) \left( \frac{\|X-x_0\|}{h_n} \right)^j \geq \frac{\rho}{2} \sup_{x \in A_n} m_n h_n^{j^*} \sum_{j=0}^{j^*} \tilde{d}_j^* \left( \frac{x-x_0}{\|x-x_0\|} \right) \left( \frac{\|x-x_0\|}{h_n} \right)^j \mid D = 1, X \in A_n \right) \end{aligned}$$



$$= P(n') \left( \sum_{j=0}^{j^*} \tilde{d}_j^* \left( \frac{X - x_0}{\|X - x_0\|} \right) \left( \frac{\|X - x_0\|}{h_{n'}} \right)^j \geq \frac{\rho}{2} \sup_{x \in A_{n'}} \sum_{j=0}^{j^*} \tilde{d}_j^* \left( \frac{x - x_0}{\|x - x_0\|} \right) \left( \frac{\|x - x_0\|}{h_{n'}} \right)^j \mid D = 1, X \in A_{n'} \right) > 0.$$

Contradiction.

Therefore there is a  $\rho, \gamma > 0$  such that for all  $h > 0$  small enough, for all  $P \in \mathcal{P}$  and  $x_0 \in [-1, 1]^d$   $P(e(X) \geq \rho \sup_{\|x - x_0\| \leq h} e(x) \mid D = 1, \|X - x_0\| \leq h) > \gamma$ .  $\square$

**Lemma 21** (Minimal eigenvalue technical result). *Suppose the conditions of Theorem 2 hold. There is an  $h' > 0$  such that*

$$\lim_{\varepsilon \rightarrow +0} \sup_{P \in \mathcal{P}, x_0 \in [-1, 1]^d, h \in (0, h'], \|v\|=1} P(\|v^T U\| \leq \varepsilon \mid D = 1, \|X - x_0\| \leq h) = 0.$$

*Proof of Lemma 21.* Take some sequence of  $\varepsilon_n \rightarrow^+ 0$ . Take  $h', \rho, \gamma$  from Assumption 7.

Let  $P(n)$  be a sequence of distributions in  $\mathcal{P}$ , let  $v_n$  be a sequence of vectors with  $\|v\| = 1$ , let  $h_n$  be a sequence in  $(0, h']$ , and write  $A_n = \{x : \|x - x_0\| \leq h_n\}$ . Then:

$$\begin{aligned} P(n) (\|v_n^T U\| \leq \sqrt{\varepsilon_n} \mid D = 1, X \in A_n) &= \frac{E_{P(n)} [e(X) 1 \{ \|v_n^T U\| \leq \sqrt{\varepsilon_n} \} \mid X \in A_n]}{E_{P(n)} [e(X) \mid X \in A_n]} \\ &\leq \frac{(\sup_{x \in A_n} e(x)) P(n) (\|v_n^T U\| \leq \sqrt{\varepsilon_n} \mid X \in A_n)}{P(n) (e(X) \geq \rho (\sup_{x \in A_n} e(x))) \rho (\sup_{x \in A_n} e(x))} \\ &\leq \frac{(\sup_{x \in A_n} e(x)) P(n) (\|v_n^T U\| \leq \sqrt{\varepsilon_n} \mid X \in A_n)}{\gamma \rho \sup_{x \in A_n} e(x)} \\ &= \frac{P(n) (\|v_n^T U\| \leq \sqrt{\varepsilon_n} \mid X \in A_n)}{\rho \gamma} = O(\sqrt{\varepsilon_n}) = o(1). \end{aligned}$$

$\square$

**Lemma 22** (Uniform minimal expected eigenvalue). *Suppose the conditions of Theorem 2 hold, and let  $U(v)$  be the vector of zero-through- $\lfloor \beta_e \rfloor$ -order interactions of  $v$ . Define:*

$$\begin{aligned} \lambda(\bar{h}) &\equiv \inf_{h \in (0, \bar{h}], \|v\|=1, x_0 \in [-1, 1]^d, P \in \mathcal{P}} v^T E_P \left[ U \left( \frac{X - x_0}{h} \right) U \left( \frac{X - x_0}{h} \right)^T \mid D = 1, \|X - x_0\| \leq h \right] v, \\ \lambda^* &\equiv \liminf_{\bar{h} \rightarrow +0} \lambda(\bar{h}), \end{aligned}$$

where  $U(\cdot)$  is the  $\lfloor \beta_\mu \rfloor$ -order local polynomial interaction matrix. Then  $\lambda^* > 0$ .

*Proof of Lemma 22.* Let  $h'$  be from Lemma 21. Fix some  $\bar{h} \in (0, h']$ . Let  $h_n$  be a sequence in  $(0, \bar{h}]$ , let  $v_n$  be a sequence of vectors with  $\|v\| = 1$ , let  $x_{0,n}$  be a sequence of points in  $[-1, 1]^d$ , and let  $P(n)$  be a sequence of distributions in  $\mathcal{P}$ . Let  $A_n = \{x : \|x - x_{0,n}\| \leq h_n\}$ .

By Lemma 21, there is a  $\varepsilon > 0$  such that:

$$\inf_{P \in \mathcal{P}, x_0 \in [-1, 1]^d, h \in (0, h_n], \|v\|=1} P \left( \|v^T U\| \geq \varepsilon \mid D = 1, \|X - x_0\| \leq h \right) > 0.$$

Call this infimum  $\delta > 0$ .  $\|v^T U\|$  is bounded, so  $\delta$  is finite.

Define:

$$\begin{aligned} \lambda_n &\equiv v_n^T E_{P(n)} \left[ U \left( \frac{X - x_{0,n}}{h_n} \right) U \left( \frac{X - x_{0,n}}{h_n} \right)^T \mid D = 1, X \in A_n \right] v_n \\ &= E_{P(n)} \left[ v_n^T U \left( \frac{X - x_{0,n}}{h_n} \right) \left( v_n^T U \left( \frac{X - x_{0,n}}{h_n} \right) \right)^T \mid D = 1, X \in A_n \right] \\ &= E_{P(n)} \left[ \left( v_n^T U \left( \frac{X - x_{0,n}}{h_n} \right) \right)^2 \mid D = 1, X \in A_n \right] \\ &\geq \varepsilon^2 P(n) \left( \left| v_n^T U \left( \frac{X - x_{0,n}}{h_n} \right) \right| \geq \varepsilon \mid D = 1, X \in A_n \right) \geq \varepsilon^2 \delta > 0. \end{aligned}$$

Therefore, for all  $\bar{h} \leq h'$ ,  $\lambda(\bar{h}) \geq \varepsilon^2 \delta$ . Therefore  $\lambda^* \geq \varepsilon^2 \delta > 0$ .  $\square$

**Lemma 23** (Uniform functional approximation across multiple gridpoints). *Suppose the conditions of Theorem 2 hold. Let  $f_{k,n}(v) : [-1, 1]^d \rightarrow \mathbb{R}$  be a set of  $k_n$  functions that are uniformly bounded. Let  $\mathcal{S}(m, \underline{h})$  be the set of sets of  $m$  tuples  $(x, h)$  with  $x \in [-1, 1]^d$  and  $h \in [\underline{h}_n, h']$ , where  $h'$  comes from Lemma 22's notion of  $h$  small enough and where for all  $(x_1, h_1), (x_2, h_2) \in S \in \mathcal{S}(m, \underline{h})$  with  $x_2 \neq x_1$ , there are no points  $x \in [-1, 1]^d$  with  $\|x - x_1\| \leq h_1$  and  $\|x - x_2\| \leq h_2$ . Write  $\tilde{c} = 2^{\frac{2(1+\gamma_0)}{1-\gamma_0}} C^{\frac{2}{1-\gamma_0}} (\sup_v |f(v)|)^{-2}$ . Suppose (i)  $\underline{h}_n^{\frac{d\gamma_0}{\gamma_0-1}} \gg n^{-1}$ ; (ii)  $m_n$  be a sequence tending to infinity such for all fixed  $a > 0$ ,  $m_n \ll \exp\left(-an\underline{h}_n^{\frac{d\gamma_0}{\gamma_0-1}}\right)$ ; and (iii) and  $k_n \left(1 - \left(\frac{e-1}{e}\right)^{\frac{m_n}{\exp\left(\log(1/2) + \tilde{c}^{-1} n^{-1} \underline{h}_n^{\frac{d\gamma_0}{\gamma_0-1}}\right)}}\right) \rightarrow 0$ . Then there is a sequence of  $\varepsilon_n \rightarrow^+ 0$  such that for all  $k = 1, \dots, k_n$ :*

$$\sup_{P \in \mathcal{P}, S \in \mathcal{S}(m_n, \underline{h}_n)} P \left( \max_j \left| \frac{\sum_i D_i K \left( \frac{\|X_i - x_j\|}{h_{n,j}} \right) f_{k,n} \left( \frac{X_i - x_j}{h_{n,j}} \right) - E \left[ DK \left( \frac{\|X - x_j\|}{h_{n,j}} \right) f_{k,n} \left( \frac{X - x_j}{h_{n,j}} \right) \right]}{n E_{P(n)} \left[ DK \left( \frac{\|X - x_j\|}{h_{n,j}} \right) \right]} \right| \geq \varepsilon_n \right) = o(1).$$

*Proof of Lemma 23.* This proof will get hairy. For simplicity, I proceed assuming  $K$  is the uniform bandwidth scaled downwards by one-half; the proof would require far more care if the kernel took on nonzero values for an unbounded set. Also, for simplicity, assume that  $e(X)$  is bounded above by one-half — the difficulty here comes from small propensity scores, and is already substantial. Let the upper bound of  $|f|$  be  $b \geq 0$ .

Define  $R_n \equiv \tilde{c} n h_n^{\frac{d\gamma_0}{\gamma_0-1}}$  and

$$V_{n,j,k} \equiv \sum_i D_i K \left( \frac{\|X_i - x_j\|}{h_{n,j}} \right) \frac{f_{k,n} \left( \frac{X_i - x_j}{h_{n,j}} \right) - E \left[ f_{k,n} \left( \frac{X - x_j}{h_{n,j}} \right) \mid DK \left( \frac{\|X - x_j\|}{h_{n,j}} \right) = 1 \right]}{n E_{P(n)} \left[ DK \left( \frac{\|X - x_j\|}{h_{n,j}} \right) \right]}.$$

By construction, there is a sequence of  $\varepsilon_n^+ \rightarrow 0$  such that  $k_n \left( 1 - (1 - 1/e)^{\frac{m_n}{\exp(\log(1/2) + R_n \varepsilon_n^2)}} \right) \rightarrow 0$ .

Note that the events  $D_i K \left( \frac{\|X_i - x_j\|}{h_{n,j}} \right)$  are mutually disjoint for a given  $i$  across  $j$ ; therefore write  $j(i)$  for the  $j$  such that  $D_i K \left( \frac{\|X_i - x_{j(i)}\|}{h_{n,j(i)}} \right) = 1$  if feasible, and write  $j(i) = 0$  if no such  $j$  exists.

Let  $P(n)$  be a sequence of distributions in  $\mathcal{P}$ , let  $S_n$  be a sequence of sets in  $\mathcal{S}(m_n, h_n)$ , and let  $\{(x_{n,j}, h_{n,j})\}$  be a sequence of associated points and bandwidths. By Lemma 19, for all  $j, n$ ,

$$E_{P(n)} \left[ DK \left( \frac{\|X - x_j\|}{h_{n,j}} \right) \right] \geq C^{1/(1-\gamma_0)} 2^{\frac{1+\gamma_0}{1-\gamma_0}} h_n^{\frac{d\gamma_0}{\gamma_0-1}}.$$

Note that I use a laxer bound for  $h_n$  because the polynomial order here is found elsewhere, and the polynomial order in Lemma 19 is not found elsewhere. Thus, the argument will continue to hold in the presence of certain typos.

Consider the event  $A$  of  $\{i, j(i)\}$ . Note that for any given  $j$ , by the Chernoff bound for binomial random variables,

$$P(n) \left( \frac{\sum D_i K \left( \frac{\|X_i - x_j\|}{h_{n,j}} \right)}{n E \left[ DK \left( \frac{\|X - x_j\|}{h_{n,j}} \right) \right]} \geq 2 \right) \leq \exp \left( \frac{-n E \left[ DK \left( \frac{\|X - x_j\|}{h_{n,j}} \right) \right]}{3} \right) \leq \exp \left( -\tilde{c} \frac{b^2}{3} n h_n^{\frac{d\gamma_0}{\gamma_0-1}} \right),$$

$$\begin{aligned} \text{So that } P(n) \left( \max_j \frac{\sum D_i K \left( \frac{\|X_i - x_j\|}{h_{n,j}} \right)}{n E \left[ DK \left( \frac{\|X - x_j\|}{h_{n,j}} \right) \right]} \leq 2 \right) &\leq 1 - \sum_j P(n) \left( \frac{\sum D_i K \left( \frac{\|X_i - x_j\|}{h_{n,j}} \right)}{n E \left[ DK \left( \frac{\|X - x_j\|}{h_{n,j}} \right) \right]} \geq 2 \right) \\ &\leq 1 - O \left( m_n \exp \left( -\tilde{c} \frac{b^2}{3} n h_n^{\frac{d\gamma_0}{\gamma_0-1}} \right) \right) = 1 - o(1). \end{aligned}$$

I therefore proceed under the high probability event that  $A$  is such that  $\max_j \frac{\sum D_i K \left( \frac{\|X_i - x_j\|}{h_{n,j}} \right)}{n E \left[ DK \left( \frac{\|X - x_j\|}{h_{n,j}} \right) \right]} \leq 2$ .

I now apply the Hoeffding inequality to the  $\sum DK \left( \frac{\|X - x_j\|}{h_{n,j}} \right) \leq 2 E \left[ DK \left( \frac{\|X - x_j\|}{h_{n,j}} \right) \right]$  elements of  $V_{n,j}$  conditional on  $A$ , for all  $k, n, j$

$$P(n) (|V_{n,j,k}| \geq \varepsilon_n) \leq 2 \exp \left( \frac{-2 \left( n E_{P(n)} \left[ DK \left( \frac{\|X - x_j\|}{h_{n,j}} \right) \right] \right)^2 \varepsilon_n^2}{\left( \sum D_i K \left( \frac{\|X_i - x_j\|}{h_{n,j}} \right) \right) b^2} \right) \leq 2 \exp (-R_n \varepsilon_n^2).$$

Then:

$$\begin{aligned}
P(n) \left( \max_{k=1}^{k_n} \max_{j=1}^{m_n} |V_{n,j,k}| \geq \varepsilon_n \mid A \right) &\leq \sum_{k=1}^{k_n} \left( 1 - \prod_{j=1}^{m_n} (1 - P(n) (|V_{n,j,k}| \geq \varepsilon_n)) \right) \\
&\leq k_n \left( 1 - \prod_{j=1}^{m_n} (1 - 2 \exp(-R_n \varepsilon_n^2)) \right) = k_n \left( 1 - \left( 1 - 2 \exp\left(\frac{-1}{16b^2} n h_n^{\frac{d}{\gamma_0-1}} \varepsilon_n^2\right) \right)^{m_n} \right) \\
&= k_n \left( 1 - \left( \left( 1 - \frac{1}{\exp(\log(1/2) + R_n \varepsilon_n^2)} \right)^{\exp(\log(1/2) + R_n \varepsilon_n^2)} \right)^{\frac{m_n}{\exp(\log(1/2) + R_n \varepsilon_n^2)}} \right) \\
&= k_n \left( 1 - (1 - 1/e - o(1))^{\frac{m_n}{\exp(\log(1/2) + R_n \varepsilon_n^2)}} \right) = o(1).
\end{aligned}$$

Therefore  $P \left( \max_{k=1}^{k_n} \max_{j=1}^{m_n} |V_{n,j,k}| \geq \varepsilon_n \right) = o(1)$ .  $\square$

**Lemma 24** (Nondegeneracy of local polynomial eigenvalues at estimated bandwidths over gridpoints). *Suppose the conditions of Theorem 2 hold. Fix some  $k > 0$  and let  $\mathcal{S}_n$  be a set of  $g_n$  points  $x_{n,j} \in [-1, 1]^d$ , with  $g_n \leq k'n$  for some fixed  $k' > 0$ . For each  $x_{n,j}$ , let  $h_{n,j} = \sup h : n \sum_i D_i \mathbf{1} \{ \|X_i - x_{n,j}\| \} \leq kh^{-2\beta_\mu}$  and let  $h_{n,j}^*$  solve  $nE[D\mathbf{1}\{\|X - x_{n,j}\|\}] = kh^{-2\beta_\mu}$ . Then there is a  $\delta_n \rightarrow^+ 0$  such that*

$$\limsup_{n \rightarrow \infty} \sup_{P \in \mathcal{P}} P \left( \max_j \left| \frac{\lambda_{\min} \left( \frac{\sum_i D_i K \left( \frac{\|X_i - x_{n,j}\|}{h_{n,j}} \right) U \left( \frac{X_i - x_{n,j}}{h_{n,j}} \right) U \left( \frac{X_i - x_{n,j}}{h_{n,j}} \right)^T}{\sum_i D_i K \left( \frac{\|X_i - x_{n,j}\|}{h_{n,j}} \right)} \right)}{\lambda_{\min} \left( \frac{E_P \left[ DK \left( \frac{\|X - x_{n,j}\|}{h_{n,j}^*} \right) U \left( \frac{X - x_{n,j}}{h_{n,j}^*} \right) U \left( \frac{X - x_{n,j}}{h_{n,j}^*} \right)^T \right]}{E_P \left[ DK \left( \frac{\|X - x_{n,j}\|}{h_{n,j}^*} \right) \right]} \right)} - 1 \right| \geq \delta_n \right) = o(1).$$

*Proof of Lemma 24.* For convenience, I proceed assuming  $K$  is the uniform bandwidth, scaled downwards by one-half. Let  $P(n)$  be a sequence of distributions in  $\mathcal{P}$ .

Apply Lemma 23 to the sequence  $h_n = n^{\frac{-1}{2\beta_\mu + d} \frac{\gamma_0}{\gamma_0 - 1}} / \log(n)$  and  $m_n = 3g_n \ll \exp\left(ah_n^{\frac{d}{\gamma_0-1}}\right)$  for all  $a > 0$ , to yield a sequence of  $\varepsilon_n^{(a)} \rightarrow^+ 0$ .

Let  $h_{n,j}$  solve  $\min_h |N_n(h \mid x_{n,j}) - kh^{-2\beta_\mu}|$ , and let  $h_{n,j}^*$  solve  $\min_h |N_n^*(h \mid x_{n,j}) - kh^{-2\beta_\mu}|$ , where  $N_n^*(h \mid x_{n,j}) = nE[e(X)\mathbf{1}\{\|X - x_{n,j}\| \leq h\}]$ . Further, let  $[h_{n,j}, \bar{h}_{n,j}]$  be the convex hull of the set of  $h$  that solve  $N_n^*(h \mid x_{n,j})(1 + \varepsilon_n^{(a)})^{-1} = kh^{-2\beta_\mu}$  or  $N_n^*(h \mid x_{n,j})(1 + \varepsilon_n^{(a)}) = kh^{-2\beta_\mu}$ . By Lemma 23, with probability tending to one,  $h_{n,j} \in [h_{n,j}, \bar{h}_{n,j}]$ .

Write:

$$A_{n,j}(h, h') = \frac{E_{P(n)} \left[ DK \left( \frac{\|X - x_{n,j}\|}{h} \right) U \left( \frac{X - x_{n,j}}{h} \right) U \left( \frac{X - x_{n,j}}{h} \right)^T \right]}{E_{P(n)} \left[ DK \left( \frac{\|X - x_{n,j}\|}{h'} \right) \right]}$$

$$B_{n,j}(h, h') = \frac{\sum_{i=1}^n D_i K\left(\frac{\|X_i - x_{n,j}\|}{h}\right) U\left(\frac{X_i - x_{n,j}}{h}\right) U\left(\frac{X_i - x_{n,j}}{h}\right)^T}{\sum_{i=1}^n D_i K\left(\frac{\|X_i - x_{n,j}\|}{h'}\right)}$$

The claim is that there is a  $\delta_n \rightarrow^+ 0$  such that

$$P(n) \left( \left| \frac{\lambda_{\min}(B_{n,j}(h_{n,j}, h_{n,j}))}{\lambda_{\min}(A_{n,j}(h_{n,j}^*, h_{n,j}^*))} - 1 \right| \geq \delta_n \right) = o(1).$$

$A_{n,j}^{-1}$  is symmetric, so that  $\|A_{n,j}^{-1}\|_{(op)}^2$  is equal to the squared largest eigenvalue of  $A_{n,j}^{-1}$ , where  $\|\cdot\|_{(op)}$  is equal to the operator norm. By Lemma 22,  $\|A_{n,j}^{-1}\|_{(op)}^2 = \lambda_{\min}(A_{n,j})^{-2}$  is bounded above.

Take  $m_n = 3g_n \leq 3k'n$ ,  $k_n = (n+1)$ , and  $h_n = n^{\frac{-1}{2\beta\mu + d\frac{\gamma_0}{\gamma_0-1}}} / \log(n)^{\frac{\gamma_0-1}{d\gamma_0}}$ , so that the first condition of Lemma 23 holds by Lemma 22. The second condition holds because  $m_n = O(n)$ . The third condition holds by L'Hopital's Rule applied to  $n \left(1 - \left(\frac{e-1}{e}\right)^{\frac{an}{\exp(b+cn^d/\log(n))}}\right) \rightarrow 0$  for any fixed  $a, b, c, d > 0$ . Thus, I may apply Lemma 23 to the  $k_n = n+1$  bounded functions  $f(v) = U_k(v)U_\ell(v)$  and  $1 \left\{ \frac{h_{n,j}}{h_{n,j}^*} \leq v \leq \frac{\bar{h}_{n,j}}{h_{n,j}^*} \right\}$  at the bandwidths  $h_n$ . Let the associated  $\varepsilon_n$  terms be  $\varepsilon_n^{(b)}$ . Then:

$$\begin{aligned} & \max_j |A_{n,j}(h_{n,j}^*, h_{n,j}^*)_{k,\ell} - B_{n,j}(h_{n,j}, h_{n,j})_{k,\ell}| \\ & \leq \max_j \left| \begin{array}{c} A_{n,j}(h_{n,j}^*, h_{n,j}^*)_{k,\ell} \\ - B_{n,j}(h_{n,j}^*, h_{n,j}^*)_{k,\ell} \end{array} \right| + \max_j \left| \begin{array}{c} B_{n,j}(h_{n,j}^*, h_{n,j}^*)_{k,\ell} \\ - B_{n,j}(h_{n,j}^*, h_{n,j}^*)_{k,\ell} \end{array} \right| + \max_j \left| \begin{array}{c} B_{n,j}(h_{n,j}^*, h_{n,j}^*)_{k,\ell} \\ - B_{n,j}(h_{n,j}, h_{n,j})_{k,\ell} \end{array} \right| \\ & \leq o_{P(n)}(1) + O_{P(n)} \left( \max_j |A_{n,j}(h_{n,j}^*, h_{n,j}^*)_{k,\ell}| \right) \max_j \left| \left( 1 - \frac{\sum D_i K\left(\frac{\|X_i - x_{n,j}\|}{h_{n,j}^*}\right)}{\sum D_i K\left(\frac{\|X_i - x_{n,j}\|}{h_{n,j}}\right)} \right) \right| + O \left( \max_j \frac{\sum D_i 1 \{ \|X_i - x_{n,j}\| \in [h_{n,j}, \bar{h}_{n,j}] \}}{\sum D_i 1 \{ \|X_i - x_{n,j}\| \leq h_{n,j} \}} \right) \\ & = o_{P(n)}(1) + O_{P(n)} \left( \varepsilon_n^{(a)} \right) + o \left( \left( 1 + \varepsilon_n^{(a)} \right)^2 (1 + \varepsilon_n^{(b)}) \right) = o_{P(n)}(1). \end{aligned}$$

Therefore  $\max_j \|A_{n,j} - B_{n,j}\|_{(op)} = o(1)$ .

Then, by well-known arguments (Horn and Johnson, 2013, p. 381):

$$\begin{aligned} \max_j |\lambda_{\min}(B_{n,j}) - \lambda_{\min}(A_{n,j})| &= \max_i |\lambda_{\max}(B_{n,j}^{-1}) - \lambda_{\max}(A_{n,j}^{-1})| \leq \max_j \|B_{n,j}^{-1} - A_{n,j}^{-1}\|_{(op)} \\ &\leq \max_j \frac{\|A_{n,j}^{-1}\|_{(op)}^2 \|A_{n,j} - B_{n,j}\|_{(op)}}{1 - \|A_{n,j}^{-1}(A_{n,j} - B_{n,j})\|_{(op)}} = \max_j \frac{O_{P(n)}(1) o_{P(n)}(1)}{1 - o_{P(n)}(1)} = o_{P(n)}(1), \end{aligned}$$

so that

$$\max_j \left| \frac{\lambda_{\min}(B_{n,j})}{\lambda_{\min}(A_{n,j})} - 1 \right| = \max_j \left| \frac{\lambda_{\min}(B_{n,j}) - \lambda_{\min}(A_{n,j})}{\lambda_{\min}(A_{n,j})} \right| \leq \max_j \frac{|\lambda_{\min}(B_{n,j}) - \lambda_{\min}(A_{n,j})|}{\min_j \lambda_{\min}(A_{n,j})} = o_{P(n)}(1),$$

so that the full claim holds.  $\square$

**Lemma 25** (Minimal small-propensity points). *Suppose the conditions of Theorem 2 hold, and there are  $m_n$  points  $x_j \in [-1, 1]^d$  such that  $\inf_{j \neq j'} \|x_j - x_{j'}\| \geq h$  and  $\max_j nE[D\mathbf{1}\{\|X - x_j\| \leq h/d\}] \leq (h/d)^{-2\beta_\mu}$ . Then there is a universal constant  $B \geq 2$  that depends on  $C, \gamma_0, \beta_\mu, d$  such that  $m_n \leq B \left( nh^{2\beta_\mu + d \frac{\gamma_0}{\gamma_0 - 1}} \right)^{1 - \gamma_0}$ .*

*Proof of Lemma 25.* Note that the set of points with  $\mathbf{1}\{\|X - x_j\| \leq h/d\}$  are mutually exclusive across  $x$ . As a result, for all such  $j$ :

$$\begin{aligned} E \left[ e(X) \mid \|X - x_j\| \leq \frac{h}{d} \right] &\leq \frac{nE[e(X)\mathbf{1}\{\|X - x_j\| \leq h/d\}]}{nP(\|X - x_j\| \leq h/d)} \leq \frac{d^{d+2\beta_\mu}}{nh^{d+2\beta_\mu}} \\ m_n \left( \frac{h}{d} \right)^d / 2 &\leq \sum_j \frac{P(\|X - x_j\| \leq \frac{h}{d})}{2} \leq P \left( e(X) \leq \frac{2d^{d+2\beta_\mu}}{nh^{d+2\beta_\mu}} \right) \leq C \left( \frac{2d^{d+2\beta_\mu}}{nh^{d+2\beta_\mu}} \right)^{\gamma_0 - 1} \\ m_n &\leq C 2^{\gamma_0} d^{2\beta_\mu(\gamma_0 - 1) + d\gamma_0} \left( nh^{2\beta_\mu + d \frac{\gamma_0}{\gamma_0 - 1}} \right)^{1 - \gamma_0}. \end{aligned}$$

Finally, if the constant is below 2, without loss of generality set  $B = 2$ .  $\square$

**Lemma 26** (KL divergence). *For any given  $L' > 0$  and  $x_0 \in [-1, 1]^d$ , construct distributions  $P_{n,m}$  for  $m = 1, 2$  as follows. Draw  $X \sim U([-1, 1]^d)$ . Draw  $D \mid X \sim \text{Bern} \left( C^{-(\gamma_0 - 1)} P_{n,m}(\|X - x_0\| \leq \|x - x_0\|)^{1/(\gamma_0 - 1)} \right)$ . Finally, draw  $Y \mid X, D \sim \mathcal{N}(D\mu_{n,m}(X), \sigma_{\min}^2)$  where  $\mu_{n,1}(X) = 0$  and*

$$\mu_{n,2}(X) = \frac{L'}{\exp\left(\frac{-4}{3}\right)} n^{\frac{-\beta_\mu}{2\beta_\mu + d \frac{\gamma_0}{\gamma_0 - 1}}} \exp \left( \frac{-1}{1 - \left( \frac{2\|X - x_0\|}{n^{\frac{-1}{2\beta_\mu + d \frac{\gamma_0}{\gamma_0 - 1}}}} \right)^2} \right) \mathbf{1} \left\{ \|X - x_0\| \leq \frac{n^{\frac{-1}{2\beta_\mu + d \frac{\gamma_0}{\gamma_0 - 1}}}}{2} \right\}.$$

Finally, define  $\mathcal{P} = \{P_{n,1}\}_{n=1^\infty, m=1,2}$ . Then (i) if there exists a  $\mathcal{P}_0$  satisfying Assumptions 6 and 7, then there exists a fixed  $L'$  such that  $\mathcal{P}$  satisfies Assumptions 6 and 7. (ii) there is an  $\alpha > 0$  finite such that  $KL(P_{n,1}, P_{n,2}) \leq \alpha$ .

*Proof of Lemma 26.* First, I show (i) that such an  $L'$  exists. Because there exists a  $P \in \mathcal{P}_0$  in this set and every  $P_{n,m}$  has the smallest possible range of  $Y - \mu(X) \mid X, D = 1$ , it must be that for every  $P_{n,m} \in \mathcal{P}$ , Assumption 1(a) and Assumption 1(c) hold. Also, if I define  $V(x) = P_{n,m}(\|X - x_0\| \leq \|x - x_0\|)$  which is distributed  $Unif([0, 1])$ , then for all  $P_{n,m}$ :

$$P_{n,m}(e(X) \leq \pi) = P_{n,m} \left( C^{-(\gamma_0 - 1)} V(X)^{\frac{1}{\gamma_0 - 1}} \leq \pi \right) = P_{n,m}(V(X) \leq C\pi^{\gamma_0 - 1}) = C\pi^{\gamma_0 - 1}.$$

Therefore, Assumption 1(d) holds. Also,

$$P_{n,m} \left( e(X) \leq \frac{\pi}{2} \right) = C \left( \frac{\pi}{2} \right)^{\gamma_0 - 1} = 2^{1 - \gamma_0} P_{n,m}(e(X) \leq \pi),$$

so that Assumption 4(i) holds with  $\rho = 2^{1-\gamma_0}$ . It is also clear that  $E_{P_{n,m}}[Y \mid X, D = 0] = E_{P_{n,0}}[Y \mid X, D = 1] \in \Sigma(\beta_\mu, L)$ , since these functions are a constant zero.

It remains to show that there is an  $L' > 0$  such that for all  $n, m$ ,  $\text{Var}_{P_{n,2}}(\mu_{n,2}(X)) \leq M$  (Assumption 1(b) and completing Assumption 6) and  $\mu_{n,2}(X) \in \Sigma(\beta_\mu, L)$  (Assumption 7). For the variance upper bound:

$$\text{Var}_{P_{n,2}}(\mu_{n,2}(X)) \leq (L')^2 \left( n^{\frac{-2\beta_\mu}{2\beta_\mu+d\frac{\gamma_0}{\gamma_0-1}}} \right) \exp\left(\frac{8}{3} - 2\right) \leq (L')^2 \exp(2/3).$$

so that it suffices to take  $L' \leq \sqrt{M \exp(-2/3)}$ . For Hölder continuity, write  $\mu_{n,2}(X) = \frac{L'}{a} n^{\frac{-\beta_\mu}{2\beta_\mu+d\frac{\gamma_0}{\gamma_0-1}}} g_a \left( \frac{\|X-x_0\|}{n^{\frac{-1}{2\beta_\mu+d\frac{\gamma_0}{\gamma_0-1}}}} \right)$ . If  $a > 0$  is small enough, then  $g_a$  is infinitely differentiable and in  $\Sigma(\beta_\mu, 1/2)$ . Thus, by standard arguments (Tsybakov, 2009), if  $L'$  is small enough,  $\mu_{n,2}(X) \in \Sigma(\beta_\mu, L)$ . Thus, there is an  $L' > 0$  such that  $\mathcal{P}$  satisfies Assumptions 6 and 7.

Second, I show the main claim (ii). It is useful to write  $h_n = n^{\frac{-1}{2\beta_\mu+d\frac{\gamma_0}{\gamma_0-1}}}$ . Then:

$$\begin{aligned} KL(P_{n,1}, P_{n,2}) &= nP_{n,1} \left( \log \frac{dP_{n,1}}{dP_{n,2}} \right) \\ &\leq nP_{n,1} \left( \|X - x_0\| \leq \frac{h_n}{4} \right) P_{n,1} \left( D = 1 \mid \|X - x_0\| \leq \frac{h_n}{4} \right) \frac{\left( \frac{L'}{\exp(-\frac{4}{3})} h_n^{\beta_\mu} \exp\left(\frac{-1}{1-1/4}\right) \right)^2}{2} \\ &= nk2^{-2d-1} (L')^2 \left( \frac{k}{C} \right)^{\frac{1}{\gamma_0-1}} h_n^{d+\frac{d}{\gamma_0-1}+2\beta_\mu} = k2^{-2d-1} (L')^2 \left( \frac{k}{C} \right)^{\frac{1}{\gamma_0-1}} = \alpha. \end{aligned}$$

□

*Proof of Proposition 6.* Let  $\underline{\delta} \geq 1$  solve:

$$\left( \frac{\underline{\delta}}{\log(\underline{\delta})} \right)^{\frac{2\beta_\mu+d\frac{\gamma_0}{\gamma_0-1}}{2\beta_\mu}} = (\underline{\delta})^{\frac{\left( \frac{\underline{\delta}}{\log(\underline{\delta})} \right)^{2\beta_\mu+d\frac{\gamma_0}{\gamma_0-1}}}{4} - 1}.$$

This is defined, because as  $\delta \rightarrow^+ 1$ , the left-hand side tends to infinity while the right-hand side tends to one; but as  $\delta$  tends to infinity, the left-hand side  $(\delta^{\frac{2\beta_\mu+d\frac{\gamma_0}{\gamma_0-1}}{2\beta_\mu}})$  grows more slowly than the right-hand side  $(\delta^{\delta/\log(\delta)})$ . I will assume that  $\left( \frac{\underline{\delta}}{\log(\underline{\delta})} \right)^{\frac{1}{2\beta_\mu}} \geq 4$ . If not, increase  $\underline{\delta}$  to have  $\left( \frac{\underline{\delta}}{\log(\underline{\delta})} \right)^{\frac{1}{2\beta_\mu}} = 4^{1/(2\beta_\mu)}$ . In either case,  $\left( \frac{\underline{\delta}}{\log(\underline{\delta})} \right)^{\frac{2\beta_\mu+d\frac{\gamma_0}{\gamma_0-1}}{2\beta_\mu}} \leq (\underline{\delta})^{\frac{\left( \frac{\underline{\delta}}{\log(\underline{\delta})} \right)^{2\beta_\mu+d\frac{\gamma_0}{\gamma_0-1}}}{4} - 1}$ . Also define  $r = \left( \frac{\underline{\delta}}{\log(\underline{\delta})} \right)^{\frac{2\beta_\mu+d\frac{\gamma_0}{\gamma_0-1}}{2\beta_\mu}}$  and  $\pi = \left( \frac{\underline{\delta}}{\log(\underline{\delta})} \right)^{\frac{-1}{2\beta_\mu}}$ . Note that by construction,  $\pi \leq 1/4^{1/(2\beta_\mu)}$ .

First, I claim that if  $\delta \geq \underline{\delta}_n^{(0)}$ , then  $h_n^{(k+1)} \leq r^{\frac{-1}{2\beta_\mu+d\frac{\gamma_0}{\gamma_0-1}}} h_n^{(k)}$  and  $m_n^{(k+1)} \geq r m_n^{(k)}$  for all  $k = 0, 1, \dots$ . I claim this inductively. It is convenient to write  $m_n^{(0)} = \delta$  and  $h_n^{(0)}$  to solve  $\delta = \exp\left(\delta \left( \left( \frac{h_n^{(0)}}{h_n^{(1)}} \right)^{-2\beta_\mu} \right)\right)$ , so

that for all  $k \geq 0$ ,

$$m_n^{(k+1)} = \left(m_n^{(k)}\right)^{\left(\frac{h_n^{(k)}}{2^{1/\beta_\mu} h_n^{(k+1)}}\right)^{2\beta_\mu}} = \left(m_n^{(0)}\right)^{\frac{(h_n^{(0)})^{2\beta_\mu}}{2^{2(k+1)} (h_n^{(k+1)})^{2\beta_\mu}}}$$

$$h_n^{(k+1)} = h_n^{(0)} \left(\frac{m_n^{(k)}}{m_n^{(0)}}\right)^{\frac{-1}{2\beta_\mu + d \frac{\gamma_0}{\gamma_0 - 1}}} = h_n^{(0)} \left(m_n^{(0)}\right)^{1 - \frac{(h_n^{(0)})^{2\beta_\mu}}{2^{2k} (h_n^{(k)})^{2\beta_\mu}} \frac{1}{2\beta_\mu + d \frac{\gamma_0}{\gamma_0 - 1}}}.$$

Then, for  $k = 0$ :

$$\delta = \exp\left(\delta \left(\left(\frac{h_n^{(1)}}{h_n^{(0)}}\right)^{2\beta_\mu}\right)\right)$$

$$\frac{h_n^{(1)}}{h_n^{(0)}} = \left(\frac{\log(\delta)}{\delta}\right)^{1/2\beta_\mu} \leq r^{\frac{-1}{2\beta_\mu + d \frac{\gamma_0}{\gamma_0 - 1}}}.$$

$$\frac{m_n^{(1)}}{m_n^{(0)}} = \left(m_n^{(0)}\right)^{\frac{1}{4} \left(\frac{h_n^{(0)}}{h_n^{(1)}}\right)^{2\beta_\mu} - 1} \geq (\delta)^{\frac{r^{\frac{2\beta_\mu}{2\beta_\mu + d \frac{\gamma_0}{\gamma_0 - 1}}}}{4} - 1} = (\delta)^{\frac{\left(\left(\frac{\delta}{\log(\delta)}\right)^{\frac{2\beta_\mu + d \frac{\gamma_0}{\gamma_0 - 1}}{2\beta_\mu}}\right)^{2\beta_\mu}}{4} - 1} \geq r.$$

I now proceed assuming the claim holds for all  $k' = 0, \dots, k-1$ . Then:

$$\frac{h_n^{(k+1)}}{h_n^{(k)}} = \left(\frac{m_n^{(k)}}{m_n^{(k-1)}}\right)^{\frac{-1}{2\beta_\mu + d \frac{\gamma_0}{\gamma_0 - 1}}} \leq r^{\frac{-1}{2\beta_\mu + d \frac{\gamma_0}{\gamma_0 - 1}}}$$

$$\frac{m_n^{(k+1)}}{m_n^{(k)}} = \left(m_n^{(0)}\right)^{\frac{(h_n^{(0)})^{2\beta_\mu}}{4^k (h_n^{(k)})^{2\beta_\mu}} \left(\frac{1}{4h_n^{(k+1)}} - 1\right)} \geq (\delta)^{\frac{(h_n^{(0)})^{2\beta_\mu}}{(h_n^{(0)})^{2\beta_\mu}} \left(\frac{1}{4h_n^{(1)}} - 1\right)} \geq r.$$

Thus, if  $\delta \geq \underline{\delta}$ ,  $\sum_{k=1}^{\infty} m_n^{(k)} \geq \sum_{k=1}^{\infty} \delta = \infty$ . □

*Proof of Theorem 2.* Recall that  $\psi_n = n^{\frac{-\beta_\mu}{2\beta_\mu + d \frac{\gamma_0}{\gamma_0 - 1}}}$ . There are two main directions to show.

**Lower bound pointwise rate.** Define  $x_0 = (0, \dots, 0)$ . Let  $\mathcal{P}$  be as in Lemma 26, with associated distributions  $P_{n,m}$  for  $m = 1, 2$ . By Lemma 26,  $K(P_{n,1}, P_{n,2}) \leq \alpha$ . Define the seminorm  $d(P, Q) = |E_P[Y \mid X = x_0, D = 1] - E_Q[Y \mid X = x_0, D = 1]|$ . By construction,

$$d(P_{n,1}, Q_{n,2}) = \overbrace{L' \exp(1/3)}^{\text{"A/2"}} n^{\frac{-\beta_\mu}{2\beta_\mu + d \frac{\gamma_0}{\gamma_0 - 1}}} n^{\frac{-\beta_\mu}{2\beta_\mu + d \frac{\gamma_0}{\gamma_0 - 1}}},$$

where  $L' > 0$  is fixed. Therefore  $d(P_{n,1}, P_{n,d}) \geq 2s_n$ , where  $s_n = An^{\frac{-\beta_\mu}{2\beta_\mu + d \frac{\gamma_0}{\gamma_0 - 1}}}$ . Thus, standard arguments



(Tsybakov, 2009) show that for any fixed estimator  $\hat{\mu}$  of  $E[Y \mid X = x_0, D = 1]$  and all  $n$  large enough,

$$\sup_{P \in \mathcal{P}} P(|\hat{\mu}(x_0) - \mu(x_0)| \geq s_n) \geq \max_{j=1,2} P_j(|\hat{\mu}(x_0) - \mu_{n,j}(x_0)| \geq s_n) \geq \max\left(\frac{\exp(-\alpha)}{4}, \frac{1 - \sqrt{\alpha/2}}{2}\right) > 0.$$

Thus, for all  $t > 1$  (see Tsybakov Theorem 2.3),

$$\liminf_{n \rightarrow \infty} \inf_{\hat{\mu}} \sup_{P \in \mathcal{P}} P\left(n^{\frac{\beta_\mu}{2\beta_\mu + d} \frac{\gamma_0}{\gamma_0 - 1}} |\hat{\mu}(x_0) - \mu(x_0)| \geq t^{\frac{\beta_\mu}{2\beta_\mu + d} \frac{\gamma_0}{\gamma_0 - 1}}\right) \geq \max\left(\frac{\exp(-ct)}{4}, \frac{1 - \sqrt{ct/2}}{2}\right),$$

where  $c$  is a constant that only depends on the parameters of the problem. Thus,  $n^{\frac{-\beta_\mu}{2\beta_\mu + d} \frac{\gamma_0}{\gamma_0 - 1}}$  is a lower bound on the *pointwise* rate of convergence.

**Achievable global rate.** This is the more difficult and interesting direction. let  $\bar{\delta}$  and  $\pi$  be taken from Proposition 6. Also let  $k_n^*$  be the smallest  $k$  for which  $\pi^{(k-1)2\beta_\mu} \leq 1/\log(n)^2$ . Note that  $k_n^* = O(\log(\log(n)))$ .

Choose some  $\delta \geq \bar{\delta}$ , where  $\bar{\delta}$  is constructed in Proposition 6. Define  $h_n^{(k)}$  and  $m_n^{(k)}$  as in Proposition 6.

I define the estimator through a series of steps. At a high level, a first step finds the grid width through a minimum level of implied overlap within grid regions. A second step chooses the local polynomial bandwidth for each gridpoint based on the number of treated observations nearby, and then conducts the associated regression. A third step interpolates between gridpoints.

Split  $[-1, 1]^d$  into a first-pass grid of hypercubes of edge length  $1/\lfloor 1/n^{\frac{-1}{2\beta_\mu + d}} \rfloor$ . Let the associated first-pass gridpoints be  $\{\tilde{x}_{n,k}\}$ . Define  $N_n(h \mid x) = \sum D_i 1\{\|X_i - x\| \leq h\}$ . Let  $\tilde{h}_{n,j} = \sup_{h \leq 1} h : N_n(h \mid \tilde{x}_{n,j}) \leq 2h^{-2\beta_\mu}$ , which is well-defined because  $N_n(0 \mid x) \leq n$ . Define  $\bar{h}_n = \max_k \tilde{h}_{n,j}$ .

Now construct the true grid using the upper-bound gridpoint distance  $\bar{h}_n$ . Split  $[-1, 1]^d$  into a grid of hypercubes of edge length  $1/\lfloor 1/\bar{h}_n \rfloor$ . For simplicity, I proceed assuming  $1/\bar{h}_n$  is an integer, so that the edges are of length  $\bar{h}_n$ ; the rounding error is second-order. Write  $g_n = (2/\bar{h}_n)^d$  for the number of gridpoints and call the points on the grid  $x_{n,1}, x_{n,2}, \dots, x_{n,g_n}$ .

For each gridpoint  $x_{n,j}$ , choose the bandwidth  $h_{n,j}$  as follows. Choose  $h_{n,j} = \sup h : N_n(h \mid x_{n,j}) \leq 2h^{-2\beta_\mu}$ . By Lemma 23 applied to  $O(n)$  points and  $f = 1$  at the fixed bandwidths  $h_{n,0}/d$ , with probability tending to one,  $h_{n,j} \geq h_{n,0}/d$  for all  $j$ .

Construct the regression estimate at gridpoints as  $\hat{\mu}(x_{n,j})$  from local polynomial regression of  $Y$  on  $U\left(\frac{X - x_{n,j}}{h_{n,j}}\right)$  with weights  $DK\left(\frac{\|X - x_{n,j}\|}{h_{n,j}}\right)$  with uniform kernel. (If this regression is degenerate, take  $\hat{\mu}(x_{n,j}) = 0$ .) For all other points  $x \in [-1, 1]^d$ , construct  $\hat{\mu}(x)$  through linear regression of  $Y$  on  $U\left(\frac{x_{n,j} - x}{h_n}\right)$  among gridpoints  $j$  with  $\|x_{n,j} - x\| \leq \lceil \beta_\mu \rceil \bar{h}_n$ . By inspection, this estimator only requires knowledge of  $\beta_\mu$ .

It remains to show that this estimator achieves a global consistency rate of at least  $\psi_n$ . Define  $h_n^{(k(j))}$  to be the largest  $h_n^{(k)}$  for  $k = 1, \dots, k_n^*$  satisfying  $nE[D1\{\|X - x_{n,j}\| \leq h/d\}] \leq (h/d)^{-2\beta_\mu}$ . As a result of this

conditioning,  $h_{n,j} \leq h_n^{(k(j))}$ . Write  $k(j) = 1$  if no such  $k$  exists; this will turn out to be an ignorable event.

I first show that there is a  $k > 0$  such that  $\bar{h}_n \leq kn^{-1/(2\beta_\mu + d\gamma_0/(\gamma_0-1))}$  almost surely. By Lemma 19, there is a  $k' > 0$  such that  $E[D\mathbf{1}\{\|X - x\| \leq h\}] \geq k'h^{d\gamma_0/(\gamma_0-1)}$  for all  $h \leq 1$  and all  $x \in [-1, 1]^d$ . Let  $h_{(init)}$  solve  $k'n h_{(init)}^{d\gamma_0/(\gamma_0-1)} = 4h_{(init)}^{-2\beta_\mu}$ , so that  $h_{(init)} = (k'/4n)^{-1/(2\beta_\mu + d\gamma_0/(\gamma_0-1))}$ . Split up  $[-1, 1]^d$  into a second-step set of hypercubes of edge length  $1/\lfloor d/h_{(init)} \rfloor$ , with gridpoints  $\tilde{x}_{n,i}^*$ . For convenience, assume that  $d/h_{(init)}$  is an integer; the difference is a rounding error. Note that the number of second-step gridpoints is  $O(n)$ . By Lemma 24, with probability tending to one,  $N_n(h_{(init)} \mid \tilde{x}_{n,i}^*) \geq 2h_{(init)}^{-2\beta_\mu}$ . I proceed under this event. Now fix an  $x \in [-1, 1]^d$  and consider a bandwidth of area  $h \geq (1 + 1/d)h_{(init)}$ . This hypercube includes at least one  $\tilde{x}_{n,i}^*$  (within a distance of  $h_{(init)}/d$ ) and all points within  $h_{(init)}$  of  $\tilde{x}_{n,i}^*$ . Call such a point  $i(x)$  arbitrarily. Thus, for such an  $h \geq (1 + 1/d)h_{(init)}$   $N_n(h \mid x) \geq N_n(h \mid \tilde{x}_{n,i(x)}^*) \geq N_n(h_{(init)} \mid \tilde{x}_{n,i(x)}^*) \geq \min_i N_n(h_{(init)} \mid \tilde{x}_{n,i}^*) \geq 2h_{(init)}^{-2\beta_\mu}$ . Thus, on the high-probability event that I am proceeding on,  $\underline{h}_n \leq h_{(init)} = kn^{-1/(2\beta_\mu + d\gamma_0/(\gamma_0-1))}$ , where  $k = (4/k')^{1/(2\beta_\mu + d\gamma_0/(\gamma_0-1))}$  for some  $k' > 0$ .

I now apply Lemma 19, Lemma 22, Lemma 23, and Lemma 24 to the gridpoints. There are  $g_n = O(n)$  gridpoints, so that Lemma 19, Lemma 22, and Lemma 24 apply immediately, so that the eigenvalues at  $h_{n,j}$  are nondegenerate. Now consider the tuples  $(x_{n,j}, h_n^{(k_n^*)})$  and the  $k_n = k_n^* + 1$  functions: first,  $f_k(v) = \mathbf{1}\{h_n^{(k_n^*)}v \leq h_n^{(k)}\}$ , and  $f_{k_n^*+1}(v) = 1$ . Define  $\underline{h}_n^* = \underline{h}_n \log(n)^{2(\gamma_0-1)/d}$ . By Lemma 19 and construction,  $h_n^{(k_n^*)} \geq \underline{h}_n^*$  for large enough  $n$ . It remains to show that (i)  $(\underline{h}_n^*)^{d\frac{\gamma_0}{\gamma_0-1}} \gg n^{-1}$ , (ii)  $g_n \ll \exp\left(\frac{1}{16}n(\underline{h}_n^*)^{d\frac{\gamma_0}{\gamma_0-1}}\right)$ , and (iii)  $k_n \left(1 - \left(\frac{e-1}{e}\right)^{\frac{g_n}{\exp\left(\log(1/2) + \frac{1}{16}n(\underline{h}_n^*)^{d\frac{\gamma_0}{\gamma_0-1}}\right)}}\right) \rightarrow 0$ . (i) and (ii) hold by inspection. (iii) holds by bounding  $k_n = O(n \log(\log(n)))$  and then applying L'Hopital's Rule.

As a result, I proceed under the high probability event that  $N_n(h_n^{(k)}/d \mid x_{n,j}) \geq \frac{n}{2}E[D\mathbf{1}\{\|X - x_{n,j}\| \leq h_n^{(k)}/d\}]$ , so that the  $k(1)$  additional assignment is never necessary, and the smallest eigenvalue of the  $\sum DKUU^T / \sum DK$  matrices at the gridpoint bandwidths  $h_{n,j}$  is at least  $\lambda^*/2 > 0$ . These conditions also ensure that there are at least  $(h_n(1)/d)^{-2\beta_\mu} = (\underline{h}_n/d)^{-2\beta_\mu} \rightarrow \infty$  treated observations within the bandwidth  $h_{n,j}$  of each gridpoint  $j$ .

Next, I claim that there are at most  $m_n^{(k)}$  gridpoints with  $h_{n,j} \geq h_n^{(k)}$ . Recall by the above that  $h_{n,j} \leq h_n^{(k(j))}$ . Thus, the number of gridpoints with  $h_{n,j} \geq h_n^{(k)}$  is bounded by the number of  $j$  with  $h_n^{(k(j))} \geq h_n^{(k)}$ . By Lemma 25, this is bounded above by:

$$B\left(n(h_n^{(k)})^{2\beta_\mu + d\frac{\gamma_0}{\gamma_0-1}}\right)^{1-\gamma_0} = B\left(\frac{h_n^{(k)}}{h_n^{(1)}}\right)^{-(\gamma_0-1)(2\beta_\mu + d\frac{\gamma_0}{\gamma_0-1})} \leq \exp\left(2^{-k}\delta\left(\frac{h_n^{(k)}}{h_n^{(1)}}\right)^{-2\beta_\mu}\right) = m_n^{(k)}.$$

Next, I bound the largest conditional expected squared gridpoint error above. Let  $Z$  be the data of

$\{X, D\}$ . The local polynomial estimator is:

$$\hat{\mu}(x_{n,j}) = \mathbf{1}^T \left( \underbrace{\frac{\sum D_i K \left( \frac{\|X_i - x_{n,j}\|}{h_n} \right) U \left( \frac{X_i - x_{n,j}}{h_n} \right) U \left( \frac{X_i - x_{n,j}}{h_n} \right)^T}{\sum D_i K \left( \frac{\|X_i - x_{n,j}\|}{h_n} \right)}}_{\hat{\Sigma}_n} \right)^{-1} \frac{\sum D_i K \left( \frac{\|X_i - x_{n,j}\|}{h_n} \right) U \left( \frac{X_i - x_{n,j}}{h_n} \right) Y}{\sum D_i K \left( \frac{\|X_i - x_{n,j}\|}{h_n} \right)}.$$

Also write  $\tilde{\mu}_n(x)$  for the prediction of  $E[Y \mid X = x, D = 1]$  based on the  $\ell_e$ -order Taylor expansion of  $\mu(X)$  at  $x_{n,j}$ . On the above event, the conditional bias of the local polynomial estimator for a given gridpoint is:

$$\begin{aligned} E[\hat{\mu}(x_{n,j}) \mid Z] - \mu(x_{n,j}) &= \mathbf{1}^T \hat{\Sigma}_n^{-1} \frac{\sum D_i K \left( \frac{\|X_i - x_{n,j}\|}{h_n} \right) U \left( \frac{X_i - x_{n,j}}{h_n} \right) (\mu(X_i) - \tilde{\mu}_n(X_i))}{\sum D_i K \left( \frac{\|X_i - x_{n,j}\|}{h_n} \right)} \\ &= O \left( \lambda_{\min}(\hat{\Sigma}_n)^{-1} \right) O(h_n^{\beta_\mu}) = O(h_n^{\beta_\mu}), \end{aligned}$$

with a constant that is independent of  $P(n)$  and  $x_{0,h}$ . Thus, on the above event, the largest gridpoint conditional bias is bounded. On the other hand, one gridpoint's conditional variance is:

$$\begin{aligned} \text{Var}(\hat{\mu}(x_{n,j}) \mid Z) &= \mathbf{1}^T \hat{\Sigma}_n^{-1} \frac{\sum D_i K \left( \frac{\|X_i - x_{n,j}\|}{h_n} \right)^2 U \left( \frac{X_i - x_{n,j}}{h_n} \right) U \left( \frac{X_i - x_{n,j}}{h_n} \right)^T \text{Var}(Y \mid X = X_i, D = 1)}{\left( \sum D_i K \left( \frac{\|X_i - x_{n,j}\|}{h_n} \right) \right)^2} \hat{\Sigma}_n^{-1} \mathbf{1} \\ &\leq O(1) \mathbf{1}^T \hat{\Sigma}_n^{-1} \frac{\sum D_i K \left( \frac{\|X_i - x_{n,j}\|}{h_n} \right) U \left( \frac{X_i - x_{n,j}}{h_n} \right) U \left( \frac{X_i - x_{n,j}}{h_n} \right)^T \sigma_{\max}^2}{\left( \sum D_i K \left( \frac{\|X_i - x_{n,j}\|}{h_n} \right) \right)^2} \hat{\Sigma}_n^{-1} \mathbf{1} \quad (K \leq 1) \\ &= O \left( \left( \sum D_i K \left( \frac{\|X_i - x_{n,j}\|}{h_n} \right) \right)^{-1} \right) \mathbf{1}^T \hat{\Sigma}_n^{-1} \mathbf{1} \sigma_{\max}^2 \\ &= O \left( \left( \sum D_i K \left( \frac{\|X_i - x_{n,j}\|}{h_n} \right) \right)^{-1} \right) O \left( \lambda_{\min}(\hat{\Sigma}_n)^{-1} \right) \\ &= O \left( n^{-1} h_n^{d+d/(\gamma_0-1)} \right) O \left( \lambda_{\min}(E_{P(n)}[\hat{\Sigma}_n])^{-1} \right) \quad (\text{Lemma 19, Lemma 24}) \\ &= O \left( n^{-1} n^{\frac{d\gamma_0(\gamma_0-1)}{2\beta_\mu+d\gamma_0/(\gamma_0-1)}} \right) O(1) \quad (\text{Lemma 22}) \\ &= O \left( n^{\frac{-2\beta_\mu}{2\beta_\mu+d\gamma_0/(\gamma_0-1)}} \right), \end{aligned}$$

once again with a constant that is independent of  $P(n)$  and  $x_{0,h}$ .

The argument for characterizing the largest gridpoint error is not quite standard. By standard arguments (Tsybakov, 2009), if there are at most  $m_n^{(j)}$  gridpoints with  $h_{n,j} \geq h_n^{(k)}$  so that  $N_n(h_{n,j} \mid x_{n,j}) \geq (h_{n,j})^{-2\beta_\mu}$  and the  $h_{n,j}$  bandwidths around  $x_{n,j}$  are by construction non-overlapping subsets of  $[-1, 1]^d$ ,

then  $E_P \left[ \max_{j: h_{n,j} \geq h_n^{(k)}} (\hat{\mu}(x_{n,j}) - E[\hat{\mu}(x_{n,j}) | Z])^2 | Z \right] = O_P \left( \log(m_n^{(j)}) (h_{n,j})^{2\beta_\mu} \right)$ . Then the conditional expected largest gridpoint squared difference overall is bounded above by

$$\begin{aligned}
E_P \left[ \max_{j=1}^{g_n} (\hat{\mu}(x_{n,j}) - E[\hat{\mu}(x_{n,j}) | Z])^2 | Z \right] &\leq \sum_{k=1}^{k_n^*} E_P \left[ \max_{j=1}^{g_n} (\hat{\mu}(x_{n,j}) - E[\hat{\mu}(x_{n,j}) | Z])^2 | Z \right] \\
&\leq \sum_{k=1}^{k_n^*} E_P \left[ \max_{h_{n,j}=h_n^{(k)}} (\hat{\mu}(x_{n,j}) - E[\hat{\mu}(x_{n,j}) | Z])^2 | Z \right] \\
&= O_p \left( \sum_{k=1}^{k_n^*-1} \log(m_n^{(j)}) (h_n^{(k)})^{2\beta_\mu} \right) + O_p \left( \log(g_n) (h_n^{(k^*)})^{2\beta_\mu} \right) \\
&= O_p \left( \delta \sum_{k=1}^{k_n^*} 2^{-k} \left( \frac{h_n^{(k)}}{h_n^{(1)}} \right)^{-2\beta_\mu} (h_n^{(k)})^{2\beta_\mu} \right) \quad (\text{Proposition 6}) \\
&\quad + O_p \left( \log(g_n) \pi^{(k-1)2\beta_\mu} (h_n^{(1)})^{2\beta_\mu} \right) \\
&= O_p \left( \sum_{k=1}^{k_n^*} 2^{-k} (h_n^{(1)})^{2\beta_\mu} \right) + O_p \left( \frac{\log(n)}{\log(n)^2} (h_n^{(1)})^{2\beta_\mu} \right) \quad (\delta \text{ fixed}) \\
&= O_p \left( \sum_{k=1}^{\infty} 2^{-k} (h_n^{(1)})^{2\beta_\mu} \right) O_p \left( \frac{\log(n)}{\log(n)^2} (h_n^{(1)})^{2\beta_\mu} \right) \quad (\delta \text{ fixed}) \\
&\leq 2 (h_n^{(1)})^{2\beta_\mu} = 2n^{\frac{-2\beta_\mu}{2\beta_\mu+d}\frac{\gamma_0}{\gamma_0-1}}.
\end{aligned}$$

Thus, returning to more standard arguments,  $\max_{j=1}^{g_n} |\hat{\mu}(x_{n,j}) - \mu(x_{n,j})| = O_{P(n)} \left( n^{\frac{-\beta_\mu}{2\beta_\mu+d}\frac{\gamma_0}{\gamma_0-1}} \right)$ . But the prediction for non-gridpoints has nondegenerate eigenvalues, so by standard arguments,  $\max_{x \in [-1,1]^d} |\hat{\mu}(x) - \mu(x)| = O \left( \max_{j=1}^{g_n} |\hat{\mu}(x_{n,j}) - \mu(x_{n,j})| + L \left( \frac{1}{\lfloor h_n \rfloor} \right)^{\beta_\mu} \right) = O_{P(n)}(\psi_n)$ . Therefore  $\hat{\mu}$  achieves the global rate  $\psi_n$ , and with no polylogarithmic penalty.

**Completing the proof.** An achievable global rate of convergence is also an achievable pointwise rate of convergence. Thus,  $n^{\frac{-\beta_\mu}{2\beta_\mu+d}\frac{\gamma_0}{\gamma_0-1}}$  is the optimal pointwise rate of convergence. A lower bound on the pointwise rate of convergence is also a lower bound on the global rate of convergence. Thus,  $n^{\frac{-\beta_\mu}{2\beta_\mu+d}\frac{\gamma_0}{\gamma_0-1}}$  is also the optimal global rate of convergence.  $\square$

*Proof of Corollary 5.* Write  $\alpha_\mu = \frac{\beta_\mu}{2\beta_\mu+d\gamma_0/(\gamma_0-1)}$  and  $\alpha_e = \frac{\beta_e}{2\beta_e+d}$ .

By standard arguments (Stone, 1982) and Theorem 2, there are cross-fit estimators  $\hat{\mu}(X)$  and  $\hat{e}(X)$  such that  $r_{\mu,n} \lesssim (n/\log(n))^{-\alpha_\mu}$  and  $r_{e,n} \lesssim (n/\log(n))^{-\alpha_e}$ , none of which depend on  $\gamma_0$ . As a result, Assumption 2 holds.

Take  $b_n = r_{e,n} \log(n)$ . Because  $\alpha_e > 0$ ,  $1 \gg b_n \gg r_{e,n}$ . By Corollary 1, it only remains to show that the conditions of Assumption 3 hold.

(a) *Consistency.* It is clear that  $r_{\mu,n}, r_{e,n} \rightarrow 0$ .

(b) *Product of errors.* If  $\gamma_0 \geq 2$ , the claim holds by inspection. If not:

$$\begin{aligned} r_{\mu,n} r_{e,n} b_n^{(\gamma_0-2)/2} &\ll r_{\mu,n} b_n^{\gamma_0/2} \ll r_{\mu,n} r_{e,n}^{\gamma_0/2} \log(n)^{\gamma_0/2} \lesssim (n/\log(n))^{-\alpha_\mu - \alpha_e \gamma_0/2} \log(n)^{\gamma_0/2} \\ &= \log(n)^{\alpha_\mu + \alpha_e \gamma_0/2 + \gamma_0/2} n^{-\alpha_\mu - \alpha_e \gamma_0/2} \\ &= n^{-1/2} \log(n)^{\alpha_\mu + \alpha_e \gamma_0/2 + \gamma_0/2} n^{1/2 - \alpha_\mu - \alpha_e \gamma_0/2} \ll n^{-1/2}. \end{aligned} \quad (\text{Equation (2)})$$

(c) *Regression error near singularities.* By inspection of the previous argument,  $r_{\mu,n} b_n^{\gamma_0/2} \ll n^{-1/2}$ .

(d) *Asymptotically known thresholding.* By construction,  $r_{e,n} \ll b_n$ .

As a result, by Corollary 1, the result for Wald confidence interval validity holds.  $\square$

## C.7 Choice of Threshold

*Proof of Lemma 1.* First, I show that there is at least one such solution.

Recall the equation:

$$f_n(b) = \frac{b \frac{1}{n} \sum \mathbf{1}\{\hat{e}(X) \leq b\}}{\sqrt{\frac{1}{n} \sum \frac{D}{\max\{\hat{e}, b\}^2}}} + b^2 \sqrt{\frac{1}{n} \sum \frac{D}{\max\{\hat{e}, b\}^2}} - n^{-1/2}.$$

When  $b = 0$ ,  $f_n(b)$  is well-defined:  $\sum D/\bar{e}$  is finite, so  $\sup D/\bar{e}^2$  is finite. Because the first two terms of  $f_n(b)$  include multiplication by  $b$ ,  $f_n(0) = 0$ .

When  $b = 1$ :

$$f_n(1) = \left( \sqrt{\frac{1}{n} \sum D} \right)^{-1} + \sqrt{\frac{1}{n} \sum D} - n^{-1/2} > \left( \sqrt{\frac{1}{n} \sum D} \right)^{-1} - 1 \geq 0.$$

The final line holds because  $\frac{1}{n} \sum D \in (0, 1]$  by assumption.

Define  $b_n^- = \sup b \leq 1 \mid f_n(b) \leq 0$ . Define  $b_n^+ = \inf b \geq b_n^- \mid f_n(b) \geq 0$ . Because  $f_n(0) \leq 0 \leq f_n(1)$ , both of these values are well-defined. Therefore, for every  $b$  satisfying  $b_n^- < b < b_n^+$ , it is the case that  $f_n(b)$  is a well-defined real number that satisfies both  $f_n(b) > 0$  and  $f_n(b) < 0$ . No such number exists, so it must be that  $b_n^- = b_n^+$ . Define  $b_n$  to be that value.

Next, I show that there is a unique solution. In particular, I show that  $\hat{g}_n(b) \equiv \frac{b \frac{1}{n} \sum \mathbf{1}\{\hat{e}(X) \leq b\}}{\sqrt{\frac{1}{n} \sum \frac{D}{\max\{\hat{e}, b\}^2}}} + b^2 \sqrt{\frac{1}{n} \sum \frac{D}{\max\{\hat{e}, b\}^2}}$  is a strictly increasing function of  $b$  for  $b \geq \min_i \hat{e}_i$ . As  $b$  increases, the first term's numerator strictly increases and the denominator weakly decreases. As a result, the first term strictly increases in that range. For  $b < \min_i \hat{e}_i$ , the first term is zero and as a result is weakly increasing. The

second term can be rewritten as

$$\sqrt{\frac{1}{n} \sum D \min\{\hat{e}^{-2}b^4, b^2\}},$$

which is a strictly increasing function. As a result,  $f_n(b)$  is a strictly increasing function in the desired range, so that there can be at most one solution. □