ESTIMATION BASED ON NEAREST NEIGHBOR MATCHING: FROM DENSITY RATIO TO AVERAGE TREATMENT EFFECT

ZHEXIAO LIN

Department of Statistics, University of California, Berkeley

PENG DING

Department of Statistics, University of California, Berkeley

FANG HAN

Department of Statistics, University of Washington

Nearest neighbor (NN) matching is widely used in observational studies for causal effects. Abadie and Imbens (2006) provided the first large-sample analysis of NN matching. Their theory focuses on the case with the number of NNs, M fixed. We reveal something new out of their study and show that once allowing M to diverge with the sample size an intrinsic statistic in their analysis constitutes a consistent estimator of the density ratio with regard to covariates across the treated and control groups. Consequently, with a diverging M, the NN matching with Abadie and Imbens' (2011) bias correction yields a doubly robust estimator of the average treatment effect and is semiparametrically efficient if the density functions are sufficiently smooth and the outcome model is consistently estimated. It can thus be viewed as a precursor of the double machine learning estimators.

KEYWORDS: Graph-based statistics, stochastic geometry, double robustness, double machine learning, propensity score.

1. INTRODUCTION

MATCHING METHODS (Greenwood (1945), Chapin (1947), Cochran and Rubin (1973), Rubin (2006), Rosenbaum (2010)) aim to balance observations from different groups through minimizing group differences in observed covariates. Such methods have proven their usefulness for causal inference in various disciplines, including economics (Imbens (2004)), epidemiology (Brookhart, Schneeweiss, Rothman, Glynn, Avorn, and Stürmer (2006)), political science (Ho, Imai, King, and Stuart (2007), Sekhon (2008)), and sociology (Morgan and Harding (2006)).

Among all the matching methods, nearest neighbor (NN) matching (Rubin (1973)) is likely the most frequently used and easiest to implement approach. In the simplest treatment-control study, NN matching assigns each treatment (control) individual to M control (treatment) individuals with the smallest distance to it. In this regard, two questions arise. First, how do we select the number of matches, M? This is referred to in the literature as ratio matching, and is both important and delicate, well known to be related to the bias-variance trade-off in nonparametic statistics (Smith (1997), Rubin and Thomas (2000), Imbens and Rubin (2015)). Second, how do we perform large-sample statistical inference for NN matching estimators? Such an analysis is usually nonstandard and technically challenging. Indeed, it was long-lacking in the literature until Abadie and Imbens (2006).

Zhexiao Lin: zhexiaolin@berkeley.edu Peng Ding: pengdingpku@berkeley.edu

Fang Han: fanghan@uw.edu

To answer the above two questions, a series of papers (Abadie and Imbens (2006, 2008, 2011, 2012)) established large-sample properties of M-NN matching for estimating the average treatment effect (ATE). These results are, however, only valid when in ratio matching, M is fixed. The according message is then mixed. As a matter of fact, the ATE estimator based on M-NN matching with a fixed M is asymptotically biased and inefficient. While bias correction is now feasible to alleviate the first issue (Abadie and Imbens (2011)), the lack of efficiency seems fundamental.

This manuscript revisits the study of Abadie and Imbens (2006) from a new perspective, bridging M-NN matching to density ratio estimation (Nguyen, Wainwright, and Jordan (2010), Sugiyama, Suzuki, and Kanamori (2012)) as well as double robustness (Scharfstein, Rotnitzky, and Robins (1999), Bang and Robins (2005)). To this end, our analysis stresses, in ratio matching, the benefits of forcing M to diverge with the sample size n in order to achieve statistical efficiency. Our claim is thus aligned with observations in the random graph-based inference literature (Wald and Wolfowitz (1940), Friedman and Rafsky (1979), Henze (1988), Liu and Singh (1993), Henze and Penrose (1999), Berrett, Samworth, and Yuan (2019), Bhattacharya (2019), Shi, Drton, and Han (2023, 2022), Lin and Han (2023)).

The contributions of this manuscript are two-fold. First, we show that a statistic that plays a pivotal role in the analysis of Abadie and Imbens (2006), $K_M(x)$ (Abadie and Imbens (2006, p. 240); to be defined in (2.2) of Section 2), which measures the number of matched times of the covariate value x, actually gives rise to a consistent density ratio estimator in the two-sample setting. Furthermore, from the angle of density ratio estimation, this NN matching-based estimator is to our knowledge the first one that simultaneously satisfies being conceptually one step, computationally efficient, and statistically rate-optimal. This estimator itself is thus an appealing alternative to existing density ratio estimators.

Getting back to the original ATE estimation problem, our second contribution is to use the above insights to bridge the bias-corrected matching estimator (Abadie and Imbens (2011)), doubly robust estimators (Scharfstein, Rotnitzky, and Robins (1999), Bang and Robins (2005), Farrell (2015)), and double machine learning estimators (Chernozhukov et al. (2018)). In fact, Abadie and Imbens' (2011) bias-corrected estimator can be formulated as

$$\widehat{\tau}_{M}^{\text{bc}} = \widehat{\tau}^{\text{reg}} + \frac{1}{n} \left[\sum_{i=1,D_{i}=1}^{n} \left(1 + \frac{K_{M}(i)}{M} \right) \widehat{R}_{i} - \sum_{i=1,D_{i}=0}^{n} \left(1 + \frac{K_{M}(i)}{M} \right) \widehat{R}_{i} \right]$$

(see Lemma 3.1, with notation introduced in Section 3 and $K_M(i)$ representing the number of times the unit i is matched), and then $1 + K_M(i)/M$ converges to the inverses of the propensity scores $1 - e(X_i)$ and $e(X_i)$ for units with $D_i = 0$ and 1, respectively. One could then leverage the general double robustness and double machine learning theory to validate the following two properties of $\widehat{\tau}_M^{\text{bc}}$:

- to validate the following two properties of $\widehat{\tau}_M^{\text{bc}}$:

 (1) *Consistency*: $\widehat{\tau}_M^{\text{bc}}$ converges in probability to the population ATE, if either the density (propensity score) functions satisfy certain conditions or the outcome (regression) model is consistently estimated, with $M \log n/n \to 0$ and $M \to \infty$ as $n \to \infty$.
 - (2) Semiparametric efficiency: $\hat{\tau}_{M}^{\text{bc}}$ is an asymptotically Normal estimator of τ with the asymptotic variance attaining the semiparametric efficiency lower bound (Hahn (1998)), if the density functions are sufficiently smooth, the outcome model is consistently estimated, and M scales with n at a proper rate. Furthermore, a simple consistent estimator of the asymptotic variance is available.

Although Abadie and Imbens (2006, Theorem 5) hints at the necessity of allowing M to diverge for gaining efficiency, we provide rigorous theory for their conjecture. Our results thus complement those made in Abadie and Imbens (2006, 2011) and provide additional theoretical justifications for practitioners to use NN matching for inferring the ATE.

Technically, our analysis hinges on a diverging M that grows with n. In contrast, existing results on NN matching for causal effects (Abadie and Imbens (2006, 2008, 2011, 2012)) all focused on a fixed M. Instead, we take a different route to establish nonasymptotic moment bounds on $K_M(x)$ with more flexibility in specifying the rate of M with respect to n (see Lin and Han (2023) for a similar idea in analyzing rank-based statistics).

Paper Organization. The rest of this manuscript proceeds as follows. Section 2 gives a brief overview of the NN matching-based density ratio estimator. Section 3 revisits Abadie and Imbens' (2011) bias-corrected NN matching-based estimator of the ATE, $\hat{\tau}_{M}^{bc}$. Section 4 elaborates on the double robustness and semiparametric efficiency of $\hat{\tau}_{M}^{bc}$ as well as its double machine learning version. Section 5 presents simulation studies to complement the theory. Section 6 includes some final remarks. We relegate technical details to the Appendix as well as an Online Appendix in the Supplementary Material (Lin, Ding, and Han (2023)). Appendices A and B introduce the algorithms and theory for the NN matching-based density ratio estimator. Appendix C and the Online Appendix present the proofs of results in the paper and in the Appendix, respectively.

Notation. For any integers $n, d \ge 1$, let $[n] = \{1, 2, ..., n\}$, n! be the factorial of n, and \mathbb{R}^d be the d-dimensional real space. A set consisting of distinct elements $x_1, ..., x_n$ is written as either $\{x_1, ..., x_n\}$ or $\{x_i\}_{i=1}^n$, and its cardinality is written by $|\{x_i\}_{i=1}^n|$. The corresponding sequence is denoted by $[x_1, ..., x_n]$ or $[x_i]_{i=1}^n$. Let $\mathbb{1}(\cdot)$ denote the indicator function. For any $a, b \in \mathbb{R}$, write $a \lor b = \max\{a, b\}$ and $a \land b = \min\{a, b\}$. We use $\stackrel{d}{\longrightarrow}$ and $\stackrel{p}{\longrightarrow}$ to denote convergence in distribution and in probability, respectively. For any sequence of random variables $\{X_n\}$, write $X_n = o_P(1)$ if $X_n \stackrel{p}{\longrightarrow} 0$ and $X_n = O_P(1)$ if X_n is bounded in probability. Let P_Z represent the law of a random variable Z.

2. DENSITY RATIO ESTIMATION VIA NN-MATCHING

Consider two general random vectors X, Z in \mathbb{R}^d that are defined on the same probability space, with d to be a fixed positive integer. Let ν_0 and ν_1 represent the probability measures of X and Z, respectively. Assume ν_0 and ν_1 are absolutely continuous with respect to the Lebesgue measure λ on \mathbb{R}^d equipped with the Euclidean norm $\|\cdot\|$; denote the corresponding densities (Radon–Nikodym derivatives) by f_0 and f_1 . Assume further that ν_1 is absolutely continuous with respect to ν_0 and write the corresponding density ratio, f_1/f_0 , as r; set 0/0=0 by default.

Assume X_1, \ldots, X_{N_0} are N_0 independent copies of X, Z_1, \ldots, Z_{N_1} are N_1 independent copies of Z, and $[X_i]_{i=1}^{N_0}$ and $[Z_j]_{j=1}^{N_1}$ are mutually independent. The problem of estimating the density ratio r based on $\{X_1, \ldots, X_{N_0}, Z_1, \ldots, Z_{N_1}\}$ is fundamental in economics (Cunningham (2021)), information theory (Cover and Thomas (2006)), machine learning (Sugiyama, Suzuki, and Kanamori (2012)), statistics (Imbens and Rubin (2015)), and other fields.

In density ratio estimation, NN-based estimators are advocated before due to its computational efficiency; cf. Lima, Cunha, Oyaizu, Frieman, Lin, and Sheldon (2008), Póczos and Schneider (2011), Kremer, Gieseke, Pedersen, and Igel (2015), Noshad, Moon,

Sekeh, and Hero (2017), Berrett, Samworth, and Yuan (2019), Zhao and Lai (2020), among many others. Based on Abadie and Imbens' (2006, 2008, 2011, 2012) NN matching framework, we propose a new density ratio estimator based on NN matching. To this end, some necessary notation is introduced first.

DEFINITION 2.1—NN Matching: For any $x, z \in \mathbb{R}^d$ and $M \in [N_0]$:

(i) let $\mathcal{X}_{(M)}(\cdot): \mathbb{R}^d \to \{X_i\}_{i=1}^{N_0}$ be the mapping that returns the value of the input z's Mth NN in $\{X_i\}_{i=1}^{N_0}$, that is, the value of $x \in \{X_i\}_{i=1}^{N_0}$ such that

$$\sum_{i=1}^{N_0} \mathbb{1}(\|X_i - z\| \le \|x - z\|) = M; \tag{2.1}$$

(ii) let $K_M(\cdot): \mathbb{R}^d \to \{0\} \cup [N_1]$ be the mapping that returns the number of matched times of x, that is,

$$K_M(x) = K_M(x; \{X_i\}_{i=1}^{N_0}, \{Z_j\}_{j=1}^{N_1}) = \sum_{j=1}^{N_1} \mathbb{1}(\|x - Z_j\| \le \|\mathcal{X}_{(M)}(Z_j) - Z_j\|); \quad (2.2)$$

(iii) let $A_M(\cdot): \mathbb{R}^d \to \mathcal{B}(\mathbb{R}^d)$ be the corresponding mapping from \mathbb{R}^d to the class of all Borel sets in \mathbb{R}^d so that

$$A_M(x) = A_M(x; \{X_i\}_{i=1}^{N_0}) = \{z \in \mathbb{R}^d : ||x - z|| \le ||\mathcal{X}_{(M)}(z) - z||\}$$
 (2.3)

returns the catchment area of x in the setting of (ii).

Because ν_0 is absolutely continuous with respect to the Lebesgue measure, (2.1) has a unique solution. Abadie and Imbens (2006, pp. 240 and 260) introduced the terms $K_M(\cdot)$ and $A_M(\cdot)$ to analyze the asymptotic behavior of their NN matching-based ATE estimator. We also adopt their terminology "catchment area" in Definition 2.1(iii). Proposition 2.1 below formally links $K_M(\cdot)$ to $A_M(\cdot)$. It was established in the proof of Abadie and Imbens (2006, Lemma 3), and is stated here to aid understanding.

PROPOSITION 2.1: For any
$$x \in \mathbb{R}^d$$
, we have $K_M(x) = \sum_{j=1}^{N_1} \mathbb{1}(Z_j \in A_M(x))$.

REMARK 2.1—Relation Between $A_M(X_i)$'s and Voronoi Tessellation When M=1: We can verify that, due to the absolute continuity of ν_0 , $[A_1(X_i)]_{i=1}^{N_0}$ are almost surely disjoint except for a Lebesgue measure zero area, and partition \mathbb{R}^d into N_0 polygons. Furthermore, we can also verify that $\{A_1(X_i)\}_{i=1}^{N_0}$ are exactly the Voronoi tessellation defined in Voronoi (1908), which plays a vital role in stochastic and computational geometry. In this case, each element $A_1(X_i)$ is a Voronoi cell from the definition of (2.3).

With these notation and concepts, we are now ready to introduce the following density ratio estimator based on NN matching.

DEFINITION 2.2—NN Matching-Based Density Ratio Estimator: For any $M \in [N_0]$ and $x \in \mathbb{R}^d$, we define the following estimator for r(x):

$$\widehat{r}_{M}(x) = \widehat{r}_{M}(x; \{X_{i}\}_{i=1}^{N_{0}}, \{Z_{j}\}_{j=1}^{N_{1}}) = \frac{N_{0}}{N_{1}} \frac{K_{M}(x)}{M}.$$
(2.4)

The estimator $\hat{r}_M(\cdot)$ is by construction a one-step estimator, and satisfies the following two properties simultaneously:

- (P1) Computationally of low complexity: it is of a subquadratic (and nearly linear when M is small) time complexity via a careful algorithmic formulation based on k-d trees (see Algorithms 1–2 and Theorem A.1 in Appendix A), and thus in many scientific applications is computationally more attractive than its optimization-based alternatives (Lima et al. (2008), Kremer et al. (2015), Borgeaud et al. (2021)).
- (P2) Statistically *rate-optimal*: it is information-theoretically efficient in terms of achieving an upper bound of estimation accuracy that matches the corresponding minimax lower bound over a class of Lipschitz density functions (see Appendix B).

3. REVISITING THE BIAS-CORRECTED MATCHING ESTIMATOR OF THE ATE

This section studies the bias-corrected NN matching-based estimator of the ATE, proposed in Abadie and Imbens (2011) to correct the asymptotic bias of the original matching-based estimator derived by Abadie and Imbens (2006). To this end, we leverage the new insights in Section 2 as well as the technical results in Appendices A–B, and bridge the study to both the classic double robustness and the modern double machine learning frameworks.

We first review the setup for the NN matching-based estimator and its bias-corrected version. Following Abadie and Imbens (2006), let $[(X_i, D_i, Y_i)]_{i=1}^n$ be n independent copies of (X, D, Y), where $D \in \{0, 1\}$ is a binary treatment variable, let $X \in \mathbb{R}^d$ represent the individual covariates, assumed to be absolute continuous admitting a density f_X , and let $Y \in \mathbb{R}$ stand for the outcome variable.

For each unit $i \in [n]$, we observe $D_i = 1$ if in the treated group and $D_i = 0$ if in the control group. Let $n_0 = \sum_{i=1}^n (1 - D_i)$ and $n_1 = \sum_{i=1}^n D_i$ be the numbers of control and treated units, respectively. Under the potential outcomes framework (Rubin (1974)), the unit i has two potential outcomes, $Y_i(1)$ and $Y_i(0)$, but we observe only one of them: $Y_i = D_i Y_i(1) + (1 - D_i) Y_i(0)$. The goal is to estimate the population ATE, $\tau = \mathrm{E}[Y_i(1) - Y_i(0)]$, based on the observations $\{(X_i, D_i, Y_i)\}_{i=1}^n$. To estimate ATE, we consider its empirical counterpart $\widehat{\tau}_M = n^{-1} \sum_{i=1}^n [\widehat{Y}_i(1) - \widehat{Y}_i(0)]$, where $\widehat{Y}_i(0)$ and $\widehat{Y}_i(1)$ are the imputed outcomes of $Y_i(0)$ and $Y_i(1)$. Following Abadie and Imbens (2006), we focus on the matching-based estimator by imputing missing potential outcomes as

$$\widehat{Y}_i(0) = \begin{cases} Y_i & \text{if } D_i = 0, \\ \frac{1}{M} \sum_{j \in \mathcal{J}_M(i)} Y_j & \text{if } D_i = 1 \end{cases} \quad \text{and} \quad \widehat{Y}_i(1) = \begin{cases} \frac{1}{M} \sum_{j \in \mathcal{J}_M(i)} Y_j & \text{if } D_i = 0, \\ Y_i & \text{if } D_i = 1. \end{cases}$$

Here, $\mathcal{J}_M(i)$ represents the index set of M-NNs of X_i in $\{X_j: D_j = 1 - D_i\}_{j=1}^n$, that is, the set of all indices $j \in [n]$ such that $D_j = 1 - D_i$ and $\sum_{\ell=1,D_\ell=1-D_i}^n \mathbb{1}(\|X_\ell - X_i\| \le \|X_j - X_i\|) \le M$. With a slight abuse of notation, let $K_M(i)$ represent the number of matched times for unit i, that is, $K_M(i) = \sum_{j=1,D_j=1-D_i}^n \mathbb{1}(i \in \mathcal{J}_M(j))$. We can then rewrite the above matching-based estimator as

$$\widehat{\tau}_{M} = \frac{1}{n} \left[\sum_{i=1, D_{i}=1}^{n} \left(1 + \frac{K_{M}(i)}{M} \right) Y_{i} - \sum_{i=1, D_{i}=0}^{n} \left(1 + \frac{K_{M}(i)}{M} \right) Y_{i} \right].$$
(3.1)

However, when d > 1, the bias of $\hat{\tau}_M$ is asymptotically nonnegligible (Abadie and Imbens (2006)). To fix this, Abadie and Imbens (2011) proposed the following bias-corrected

version for $\widehat{\tau}_M$. In detail, let $\widehat{\mu}_0(x)$ and $\widehat{\mu}_1(x)$ be mappings from \mathbb{R}^d to \mathbb{R} that estimate the conditional means of the outcomes $\mu_0(x) = \mathrm{E}[Y \mid X = x, D = 0]$ and $\mu_1(x) = \mathrm{E}[Y \mid X = x, D = 1]$, respectively, with the corresponding residuals $\widehat{R}_i = Y_i - \widehat{\mu}_{D_i}(X_i)$, $i \in [n]$. Define the estimator based on outcome regressions as $\widehat{\tau}^{\mathrm{reg}} = n^{-1} \sum_{i=1}^{n} [\widehat{\mu}_1(X_i) - \widehat{\mu}_0(X_i)]$. Consider the bias-corrected matching-based estimator in Abadie and Imbens (2011):

$$\widehat{\tau}_{M}^{bc} = \frac{1}{n} \sum_{i=1}^{n} [\widehat{Y}_{i}^{bc}(1) - \widehat{Y}_{i}^{bc}(0)], \qquad (3.2)$$

with

$$\widehat{Y}_i^{ ext{bc}}(0) = egin{cases} Y_i & ext{if } D_i = 0, \ rac{1}{M} \sum_{j \in \mathcal{J}_M(i)} \left(Y_j + \widehat{\mu}_0(X_i) - \widehat{\mu}_0(X_j)
ight) & ext{if } D_i = 1, \end{cases}$$

and

$$\widehat{Y}_i^{\text{bc}}(1) = \begin{cases} \frac{1}{M} \sum_{j \in \mathcal{J}_M(i)} (Y_j + \widehat{\mu}_1(X_i) - \widehat{\mu}_1(X_j)) & \text{if } D_i = 0, \\ Y_i & \text{if } D_i = 1. \end{cases}$$

Lemma 3.1 below shows an equivalent form of $\hat{\tau}_{M}^{bc}$.

LEMMA 3.1: The bias-corrected matching-based estimator in (3.2) can be rewritten in terms of $\hat{\tau}^{reg}$ and the residuals \hat{R}_i 's as

$$\widehat{\tau}_{M}^{\text{bc}} = \widehat{\tau}^{\text{reg}} + \frac{1}{n} \left[\sum_{i=1, D_{i}=1}^{n} \left(1 + \frac{K_{M}(i)}{M} \right) \widehat{R}_{i} - \sum_{i=1, D_{i}=0}^{n} \left(1 + \frac{K_{M}(i)}{M} \right) \widehat{R}_{i} \right]. \tag{3.3}$$

Otsu and Rai (2017) derived another linear form of $\widehat{\tau}_M^{bc}$ to motivate a bootstrap procedure for variance estimation. The form in (3.3) is related to doubly robust estimators reviewed shortly. In detail, we first have some outcome models and residuals defined in the same way as above, and then let $\widehat{e}(x) : \mathbb{R}^d \to \mathbb{R}$ be a generic estimator of the propensity score (Rosenbaum and Rubin (1983)), e(x) = P(D = 1 | X = x). The doubly robust estimator in Scharfstein, Rotnitzky, and Robins (1999) and Bang and Robins (2005) could then be formulated as

$$\widehat{\tau}^{\text{dr}} = \widehat{\tau}^{\text{reg}} + \frac{1}{n} \left[\sum_{i=1, D_i=1}^n \frac{1}{\widehat{e}(X_i)} \widehat{R}_i - \sum_{i=1, D_i=0}^n \frac{1}{1 - \widehat{e}(X_i)} \widehat{R}_i \right]. \tag{3.4}$$

Conditional on (D_1, \ldots, D_n) , $[X_i : D_i = \omega]_{i=1}^n$ are n_ω i.i.d. random variables sampled from the distribution of $X \mid D = \omega$, and the two groups of sample points, $[X_i : D_i = 0]_{i=1}^n$ and $[X_i : D_i = 1]_{i=1}^n$, are mutually independent. Let $f_{X\mid D=\omega}$ denote the density of $X\mid D=\omega$. From the construction of $K_M(i)$ and results in Appendix B, once allowing M to diverge to infinity, conditional on (D_1, \ldots, D_n) , $n_0/n_1 \cdot K_M(i)/M$ and $n_1/n_0 \cdot K_M(i)/M$ are consistent estimators of $f_{X\mid D=1}(X_i)/f_{X\mid D=0}(X_i)$ and $f_{X\mid D=0}(X_i)/f_{X\mid D=1}(X_i)$ for units with $D_i=0$ and $D_i=1$, respectively. Because n_1/n_0 converges almost surely to P(D=1)

1)/P(D=0) by the law of large numbers, the statistic $1+K_M(i)/M$ is then a consistent estimator of $1/(1-e(X_i))$ and $1/e(X_i)$ for units with $D_i=0$ and $D_i=1$, respectively. Thus, in view of (3.4), the bias-corrected matching-based estimator $\hat{\tau}_M^{\rm bc}$ in (3.3) is actually a doubly robust estimator of τ , and accordingly, should also enjoy all the corresponding desirable properties. This novel insight into $\hat{\tau}_M^{\rm bc}$ allows us to establish its asymptotic properties with a diverging M.

4. ASYMPTOTIC ANALYSIS WITH DIVERGING M

The theory for matching with a diverging M has been an important gap in the literature. With a univariate covariate, Abadie and Imbens (2006) provided a heuristic argument about the additional efficiency gain for $\hat{\tau}_M$ with larger M. With a general covariate, Abadie and Imbens (2011) used simulation to evaluate the finite-sample properties of $\hat{\tau}_M^{\rm bc}$ and highlighted the importance of bias correction with large M. Nevertheless, existing theoretical results for NN matching estimators all focused on fixed M (Abadie and Imbens (2006, 2008, 2011, 2016), Kallus (2020), Armstrong and Kolesár (2021), Ferman (2021)). In this section, we will present the corresponding theory with a diverging M and also make connections between $\hat{\tau}_M^{\rm bc}$ and double robustness/double machine learning estimators.

4.1. The Original Matching-Based Estimator

We first analyze the original bias-corrected matching-based estimator $\widehat{\tau}_M^{\text{bc}}$. Let $U_{\omega} = Y(\omega) - \mu_{\omega}(X)$ for $\omega \in \{0, 1\}$ and \mathbb{X} be the support of X. Let $\|\cdot\|_{\infty}$ denote the L_{∞} norm of a function.

We need following assumptions to prove the consistency of $\widehat{\tau}_{M}^{\text{bc}}$.

ASSUMPTION 4.1: (i) For almost all $x \in \mathbb{X}$, D is independent of (Y(0), Y(1)) conditional on X = x, and there exists a constant $\eta > 0$ such that $\eta < P(D = 1 | X = x) < 1 - \eta$.

- (ii) $[(X_i, D_i, Y_i)]_{i=1}^n$ are i.i.d. following the joint distribution of (X, D, Y).
- (iii) $E[U_{\omega}^2 | X = x]$ is uniformly bounded for almost all $x \in \mathbb{X}$ and $\omega \in \{0, 1\}$.
- (iv) $E[\mu_{\omega}^{2}(X)]$ is bounded for $\omega \in \{0, 1\}$.

ASSUMPTION 4.2: For $\omega \in \{0, 1\}$, there exists a deterministic function $\bar{\mu}_{\omega}(\cdot) : \mathbb{R}^d \to \mathbb{R}$ such that $\mathrm{E}[\bar{\mu}_{\omega}^2(X)]$ is bounded and the estimator $\widehat{\mu}_{\omega}(x)$ satisfies $\|\widehat{\mu}_{\omega} - \bar{\mu}_{\omega}\|_{\infty} = o_{\mathrm{P}}(1)$.

ASSUMPTION 4.3: For $\omega \in \{0, 1\}$, the estimator $\widehat{\mu}_{\omega}(x)$ satisfies $\|\widehat{\mu}_{\omega} - \mu_{\omega}\|_{\infty} = o_{\mathbb{P}}(1)$.

Assumption 4.1(i) is the unconfoundedness and overlap assumptions, and is often referred to as the strong ignorability condition (Rosenbaum and Rubin (1983)). Assumption 4.2 allows for outcome model misspecification; for example, if $\hat{\mu}_{\omega} = \bar{\mu}_{\omega} = 0$, $\hat{\tau}_{M}^{bc}$ then reduces to $\hat{\tau}_{M}$. Assumption 4.3 assumes that the outcome models are consistently estimated

We need the following assumptions to prove the efficiency of $\widehat{ au}_{M}^{\mathrm{bc}}$.

Assumption 4.4: (i) $E[U^2_{\omega} | X = x]$ is uniformly bounded away from zero for almost all $x \in \mathbb{X}$ and $\omega \in \{0, 1\}$.

(ii) There exists a constant $\kappa > 0$ such that $\mathbb{E}[|U_{\omega}|^{2+\kappa} | X = x]$ is uniformly bounded for almost all $x \in \mathbb{X}$ and $\omega \in \{0, 1\}$.

(iii) $\max_{t \in \Lambda_{\lfloor d/2 \rfloor + 1}} \|\partial^t \mu_{\omega}\|_{\infty}$ is bounded, where for any positive integer k, Λ_k is the set of all *d*-dimensional vectors of nonnegative integers $t = (t_1, \ldots, t_d)$ such that $\sum_{i=1}^d t_i = k$ and $|\cdot|$ stands for the floor function.

ASSUMPTION 4.5: For $\omega \in \{0, 1\}$, the estimator $\widehat{\mu}_{\omega}(x)$ satisfies

$$\max_{t \in \Lambda_{\lfloor d/2 \rfloor + 1}} \left\| \partial^t \widehat{\mu}_{\omega} \right\|_{\infty} = O_{\mathbb{P}}(1) \quad and \quad \max_{t \in \Lambda_{\ell}} \left\| \partial^t \widehat{\mu}_{\omega} - \partial^t \mu_{\omega} \right\|_{\infty} = O_{\mathbb{P}}(n^{-\gamma_{\ell}}) \quad for \ all \ \ell \in \left[\lfloor d/2 \rfloor \right],$$

with some constants γ_{ℓ} 's satisfying $\gamma_{\ell} > \frac{1}{2} - \frac{\ell}{d}$ for $\ell = 1, 2, ..., \lfloor d/2 \rfloor$.

Assumption 4.4 is comparable to Assumption A.4 and the assumptions in Abadie and Imbens (2011, Theorem 2). Compared with the assumptions in Abadie and Imbens (2011, Theorem 2), Assumption 4.4(iii) is weaker in the sense that it only requires a finite order of smoothness. Assumption 4.5 again assumes the approximation accuracy of the outcome models, with lower convergence rates required for higher-order derivatives of the outcome models. Under some smoothness conditions on the outcome model as made in Abadie and Imbens (2011), Assumption 4.5 holds using power series approximation (Abadie and Imbens (2011, Lemma A.1)). Lastly, compared with Chernozhukov et al. (2018), we need approximation accuracy concerning derivatives of the outcome model estimator, which is not required in Chernozhukov et al. (2018); see Section 4.2 for more discussions.

Theorem 4.1 below presents the double robustness and semiparametric efficiency properties of $\hat{\tau}_{M}^{bc}$. Recall the semiparametric efficiency lower bound for estimating ATE (see Hahn (1998)):

$$\sigma^{2} = \mathbb{E}\left[\mu_{1}(X) - \mu_{0}(X) + \frac{D(Y - \mu_{1}(X))}{e(X)} - \frac{(1 - D)(Y - \mu_{0}(X))}{1 - e(X)} - \tau\right]^{2},$$

and introduce an estimator for σ^2 based on NN matching:

$$\widehat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \left[\widehat{\mu}_1(X_i) - \widehat{\mu}_0(X_i) + (2D_i - 1) \left(1 + \frac{K_M(i)}{M} \right) \widehat{R}_i - \widehat{\tau}_M^{\text{bc}} \right]^2.$$

THEOREM 4.1: (i) (Double robustness of $\widehat{\tau}_{M}^{bc}$) On one hand, if the distribution of (X, D, Y) satisfies Assumptions 4.1, 4.2, either $(P_{X|D=0}, P_{X|D=1})$ or $(P_{X|D=1}, P_{X|D=0})$ satisfies Assumption B.1 in the Appendix, and $M \log n/n \to 0$ and $M \to \infty$ as $n \to \infty$, then $\widehat{\tau}_M^{\text{bc}} - \tau \stackrel{\text{p}}{\longrightarrow} 0$.

On the other hand, if the distribution of (X, D, Y) satisfies Assumptions 4.1 and

4.3, then $\widehat{\tau}_M^{\rm bc} - \tau \stackrel{\mathsf{p}}{\longrightarrow} 0$.

(ii) (Semiparametric efficiency of $\widehat{ au_M}^{bc}$) Assume the distribution of (X,D,Y) satisfies Assumptions 4.1, 4.4, 4.5, and either $(P_{X|D=0}, P_{X|D=1})$ or $(P_{X|D=1}, P_{X|D=0})$ satisfies Assumption B.1 in the Appendix. Define

$$\gamma = \left\{ \min_{\ell \in \llbracket \lfloor d/2 \rfloor \rrbracket} \left[1 - \left(\frac{1}{2} - \gamma_{\ell} \right) \frac{d}{\ell} \right] \right\} \wedge \left[1 - \frac{1}{2} \frac{d}{\lfloor d/2 \rfloor + 1} \right],$$

recalling that γ_ℓ 's were introduced in Assumption 4.5. Then, if $M \to \infty$ and $M/n^{\gamma} \to$ $0 \text{ as } n \to \infty, \text{ we have } \sqrt{n}(\widehat{\tau}_M^{\text{bc}} - \tau) \stackrel{d}{\longrightarrow} N(0, \sigma^2).$

If in addition Assumption 4.3 *holds, then* $\widehat{\sigma}^2 \stackrel{p}{\longrightarrow} \sigma^2$.

REMARK 4.1: To be in line with the double robustness terminology, we can call Assumptions B.1 used in Theorem 4.1 the "density (or propensity) model assumptions" and Assumptions 4.3–4.5 the "outcome (or regression) model assumptions."

REMARK 4.2: The first part of Theorem 4.1(i) requires $M \to \infty$ for achieving the consistency of the propensity score model. When M is fixed and the outcome model is misspecified, $\widehat{\tau}_M^{\rm bc}$ is no longer doubly robust in the sense of Theorem 4.1. However, it does not imply that $\widehat{\tau}_M^{\rm bc}$ is inconsistent for estimating τ . In fact, Abadie and Imbens (2006, Theorem 3) showed that $\widehat{\tau}_M^{\rm bc}$ with a fixed M can still be consistent even if we choose $\widehat{\mu}_w = 0$ for w = 0, 1. They showed that $\widehat{\tau}_M^{\rm bc}$ with a fixed M is consistent as long as the outcome models are smooth but misspecified since the matching discrepancy then converges to zero.

REMARK 4.3: Theorem 4.1 has implications for practical data analysis. We discuss two. First, it highlights the importance of allowing M to diverge in asymptotic analysis. Nevertheless, it is a challenging problem to choose M in finite samples. We use simulation to illustrate the choice of M. Second, it gives an alternative variance estimator $\hat{\sigma}^2$ for the bias-corrected matching estimator when M diverges. Abadie and Imbens (2006) gave another variance estimator for fixed M. While it is challenging to compare the two variance estimators in theory, we use simulation to compare them in finite samples. See Section 5 for the details of simulation.

If d=1 and we pick $\widehat{\mu}_{\omega}=0$ for $\omega\in\{0,1\}$, then Assumption 4.5 automatically holds and the bias-corrected estimator $\widehat{\tau}_{M}^{bc}$ reduces to the original estimator $\widehat{\tau}_{M}$ studied in Abadie and Imbens (2006). Theorem 4.1(ii) then directly implies the following corollary that corresponds to Abadie and Imbens (2006, Corollary 1) with one key difference that M goes to infinity here.

COROLLARY 4.1—Semiparametric efficiency of $\widehat{\tau}_M$ when d=1: Assume d=1, the distribution of (X,D,Y) satisfies Assumptions 4.1, 4.4, and either $(P_{X|D=0},P_{X|D=1})$ or $(P_{X|D=1},P_{X|D=0})$ satisfies Assumption B.1 in the Appendix. If $M \to \infty$ and $M/n^{\frac{1}{2}} \to 0$ as $n \to \infty$, then $\sqrt{n}(\widehat{\tau}_M - \tau) \stackrel{d}{\longrightarrow} N(0,\sigma^2)$.

REMARK 4.4: By picking $\widehat{\mu}_{\omega} = 0$ for $\widehat{\tau}_{M}$, Assumption 4.3 is in general no longer satisfied. Accordingly, in Corollary 4.1, $\widehat{\sigma}^{2}$ may not be a consistent estimator of σ^{2} without additional assumptions. However, by decomposing σ^{2} into the form of Theorem 1 in Hahn (1998), one could still estimate σ^{2} via a similar and direct way as what is outlined in Section 4 in Abadie and Imbens (2006). We do not pursue this track in detail here as the case of d=1 without Assumption 4.3 is beyond the main scope of this manuscript.

4.2. A Double Machine Learning Version of the Matching

Assumptions 4.4 and 4.5 enforce arguably strong requirements on the smoothness of the outcome model. To weaken such assumptions, Chernozhukov et al. (2018) introduced the idea of double machine learning. In this section, we consider the option to combine NN matching with double machine learning.

Assume *n* is divisible by *K* for simplicity. Let $[I_k]_{k=1}^K$ be a *K*-fold random partition of [n], with each of size equal to n' = n/K. For each $k \in [K]$ and $\omega \in \{0, 1\}$, construct $\widehat{\mu}_{\omega, k}(\cdot)$

using data $[(X_j, D_j, Y_j)]_{j=1, j \notin I_k}^n$, and let $K_{M,k}(i)$ be the number of matched times for unit *i* by adding (X_i, D_i, Y_i) into $[(X_i, D_j, Y_j)]_{i=1, j \notin I_i}^n$. Define

$$\widetilde{\tau}_{M,k}^{\text{bc}} = \frac{1}{n'} \sum_{i=1,i \in I_k}^{n} \left[\widehat{\mu}_{1,k}(X_i) - \widehat{\mu}_{0,k}(X_i) \right] + \frac{1}{n'} \left[\sum_{i=1,i \in I_k, D_i = 1}^{n} \left(1 + \frac{K_{M,k}(i)}{M} \right) (Y_i - \widehat{\mu}_{1,k}(X_i)) \right] \\
- \sum_{i=1,i \in I_k, D_i = 0}^{n} \left(1 + \frac{K_{M,k}(i)}{M} \right) (Y_i - \widehat{\mu}_{0,k}(X_i)) \right]$$

for k = 1, ..., K, and then define $\tilde{\tau}_{M,K}^{bc} = K^{-1} \sum_{k=1}^{K} \tilde{\tau}_{M,k}^{bc}$. We can use the same variance estimator $\widehat{\sigma}^2$ for $\widetilde{\tau}_{M,K}^{bc}$.

To analyze $\widetilde{\tau}_{M,K}^{bc}$ instead of $\widehat{\tau}_{M}^{bc}$, we replace Assumptions 4.4 and 4.5 with the following two assumptions.

ASSUMPTION 4.6: (i) $E[U_{\omega}^2]$ is bounded away from zero for $\omega \in \{0, 1\}$.

(ii) There exists a constant $\kappa > 0$ such that $E[|Y|^{2+\kappa}]$ is bounded.

ASSUMPTION 4.7: For $\omega \in \{0,1\}$, the estimator $\widehat{\mu}_{\omega}(x)$ satisfies $\|\widehat{\mu}_{\omega} - \mu_{\omega}\|_{\infty} =$ $o_{P}(n^{-d/(4+2d)}).$

REMARK 4.5: Assumption 4.6 corresponds to Assumption 5.1 in Chernozhukov et al. (2018), and is similar to Assumption 4 in Abadie and Imbens (2006). Assumption 4.7 assumes approximation accuracy of the outcome model under the L_{∞} norm. Abadie and Imbens (2011) used the power series approximation (Newey (1997)) to estimate the outcome model, which under some classic nonparametric statistics assumptions automatically satisfies Assumption 4.7 (cf. Lemma A.1 in Abadie and Imbens (2011)). The same conclusion also applies to spline and wavelet regression estimators; cf. Chen and Christensen (2015).

REMARK 4.6: Assumption 4.7 assumes an approximation rate under L_{∞} norm. This is different from the L_2 norm put in Chernozhukov et al. (2018, Assumption 5.1), but can be handled with some trivial modifications to the proof of Chernozhukov et al. (2018, Theorem 5.1) since one can replace the Cauchy–Schwarz inequality by the L_1 – L_{∞} Hölder's inequality. An L_1 -norm bound on $K_M(i)/M$, to be established in Theorem B.4 in the Appendix, can then be applied directly.

THEOREM 4.2: (i) (Double robustness of $\widetilde{\tau}_{M,K}^{bc}$) On one hand, if the distribution of (X,D,Y) satisfies Assumptions 4.1, 4.2, either $(P_{X|D=0},P_{X|D=1})$ or $(P_{X|D=1},P_{X|D=0})$ satisfies Assumption B.1 in the Appendix, and $M \log n/n \to 0$ and $M \to \infty$ as $n \to \infty$, then $\widetilde{\tau}_{M,K}^{bc} - \tau \stackrel{p}{\longrightarrow} 0$. On the other hand, if the distribution of (X,D,Y) satisfies Assumptions 4.1 and

4.3, then $\widetilde{\tau}_{MK}^{bc} - \tau \stackrel{p}{\longrightarrow} 0$.

(ii) (Semiparametric efficiency of $\widetilde{\tau}_{M,K}^{bc}$) Assume the distribution of (X, D, Y) satisfies Assumptions 4.1, 4.6, 4.7 and either $(P_{X|D=0}, P_{X|D=1})$ or $(P_{X|D=1}, P_{X|D=0})$ satisfies Assumetric efficiency of $\widetilde{\tau}_{M,K}^{bc}$) sumption B.3 in the Appendix. Then if we pick $M = \alpha n^{\frac{2}{2+d}}$ for some constant $\alpha > 0$, then $\sqrt{n}(\widetilde{\tau}_{MK}^{\text{bc}} - \tau) \stackrel{\text{d}}{\longrightarrow} N(0, \sigma^2).$

In addition, we have $\widehat{\sigma}^2 \stackrel{p}{\longrightarrow} \sigma^2$.

REMARK 4.7: There are two parts where Theorem 4.2(ii) requires stronger conditions than Theorem 4.1(ii). First, Theorem 4.2(ii) requires M to grow polynomially fast with n, whereas Theorem 4.1(ii) only requires M to (i) diverge not so fast for controlling the difference of matching units and (ii) diverge to infinity (no matter how slowly it is) for achieving semiparametric efficiency. The assumptions in Theorems 4.1(ii) and 4.2(ii) both ensure semiparametric efficiency for bias-corrected matching-based estimators. Second, Theorem 4.1(ii) only requires Assumption B.1 for the density model. This is again weaker than the Lipschitz-type conditions (Assumption B.3) assumed in Theorem 4.2(ii) but is in line with the observations made in Abadie and Imbens (2006) and Abadie and Imbens (2011). Of note, these relaxations are possible due to adding more smoothness assumptions on the outcome model (Assumptions 4.4–4.5 versus Assumptions 4.6–4.7).

REMARK 4.8: Technically, to use Chernozhukov et al.'s (2018) Theorem 5.1 to establish Theorem 4.2, we need some modifications due to a reparametrization of the nuisance parameters. This is because Chernozhukov et al. (2018) considered estimating 1/e(X) and 1/(1-e(X)) via plugging in an estimate of e(X), whereas $\widetilde{\tau}_{M,K}^{bc}$ directly uses $1+K_M(X)/M$ to estimate 1/e(X) and 1/(1-e(X)) for units with D=1 and D=0, respectively. We elaborate the modifications in the proof of Theorem 4.2(ii).

5. SIMULATION

This section uses simulation to complement the theory. We consider bias-corrected matching estimators with either a fixed or diverging M, with the asymptotic variance estimated by either $\hat{\sigma}^2$ or the estimator introduced in Abadie and Imbens (2006, Section 4).

The first data are from the National Supported Work (LaLonde (1986)). We use the specific sample studied in Dehejia and Wahba (1999). The data contain 185 treated and 260 control units. To simulate data from this study, we follow the Monte Carlo simulation design of Athey, Imbens, Metzger, and Munro (2023), and use the same pretreatment variables that include "age," "education," "black," 'Hispanic," "married," "nodegree," "re74," and "re75." By using the conditional Wasserstein Generative Adversarial Networks (WGAN), one could then create a large population of observations similar to the real data, and have access to both potential outcomes for evaluating the treatment effect. Specifically, we directly use the conditional WGAN generated data available on the repository of Athey et al. (2023). There the population size is 1,000,000. For a given sample size n, we set $n_1 = n * 185/(185 + 260)$ and $n_0 = n * 260/(185 + 260)$, and draw samples from the generated data separately for treated and control groups. We consider $n \in \{600, 1200, 4800, 9600\}$.

The second data are from Shadish, Clark, and Steiner (2008), which evaluated the effects of mathematical training on mathematics test performance. We use the data from the nonrandomized arm. The data contain 79 treated and 131 control units. We use nine pretreatment covariates including "vocabulary pretest,", "mathematics pretest,", "number of prior mathematics courses," 'Caucasian," "age," "male," "mother education," "father education," and "high school GPA." We follow the Monte Carlo simulation design of Athey et al. (2023) to generate new data with population size 1,000,000. For a given sample size n, we set $n_1 = n * 79/(79 + 131)$ and $n_0 = n * 131/(79 + 131)$. Other settings are the same as those of the first data.

We consider the estimator $\widehat{\tau}_M^{\text{bc}}$ with both fixed $M \in \{1, 4, 16\}$ and diverging $M = \lfloor \alpha n^{2/(2+d)} \rfloor$ of d = 4 for the first data and d = 7 for the second data; here, the diverg-

ing rate is suggested by Theorem B.4. In this study, we pick $\alpha \in \{0.5, 1, 2, 5, 10\}$. Notice that here we choose d=4 for the first data since in the eight pretreatment variables there are only four continuous variables; it is straightforward to check that the rest four binary variables will not affect the asymptotic properties established in this manuscript as well as those in Abadie and Imbens (2006). We choose d=7 for the second data based on the same reason. For the outcome models, we consider the second-order power series. The estimator's asymptotic variance is estimated using either $\hat{\sigma}$ in Theorem 4.1(ii) (SE) or Abadie and Imbens's (2006) (AISE). We implement 2000 repetitions and Tables I and II report the calculated root-mean-squared-error (RMSE), bias, standard deviation (SD), mean-absolute-error (MAE), and the empirical coverage rate for nominal 95% and 90% confidence intervals. Tables I and II also provide inside the parentheses the root-n scaled RMSE, bias, SD, and MAE divided by σ^* , with σ^* computed as the sample size-scaled standard deviation of $\widehat{\tau}_M^{\text{bc}}$ with the sample size chosen to be 100,000, $\alpha = 1$, and 2000

TABLE I SIMULATION RESULTS, LALONDE (1986), $\sigma^* = 9.55$.

						95% Coverage		90% Coverage	
n	M	RMSE	Bias	SD	MAE	SE	AISE	SE	AISE
600	M = 1	1.055 (2.71)	-0.039 (-0.10)	1.054 (2.70)	0.469 (1.20)	0.930	0.913	0.868	0.858
	M = 4	1.043 (2.68)	-0.038(-0.10)	1.042 (2.67)	0.442(1.13)	0.926	0.911	0.862	0.847
	M = 16	1.037 (2.66)	-0.027(-0.07)	1.036 (2.66)	0.435 (1.12)	0.931	0.913	0.873	0.858
	$\alpha = 0.5$	1.043 (2.68)	-0.038(-0.10)	1.042 (2.67)	0.442(1.13)	0.926	0.911	0.862	0.847
	$\alpha = 1$	1.039 (2.67)	-0.034(-0.09)	1.039 (2.67)	0.437 (1.12)	0.928	0.911	0.864	0.845
	$\alpha = 2$	1.037 (2.66)	-0.027(-0.07)	1.036 (2.66)	0.435 (1.12)	0.931	0.913	0.873	0.858
	$\alpha = 5$	1.037 (2.66)	-0.022(-0.06)	1.037 (2.66)	0.434 (1.11)	0.948	0.927	0.891	0.869
	$\alpha = 10$	1.037 (2.66)	$-0.058\ (-0.15)$	1.036 (2.66)	0.433 (1.11)	0.947	0.926	0.901	0.882
1200	M = 1	0.341 (1.24)	-0.001 (-0.00)	0.341 (1.24)	0.272 (0.99)	0.939	0.941	0.886	0.890
	M = 4	0.310 (1.12)	0.002 (0.01)	0.310 (1.12)	0.248(0.90)	0.934	0.936	0.884	0.887
	M = 16	0.305 (1.11)	0.014(0.05)	0.305 (1.11)	0.244 (0.88)	0.939	0.940	0.887	0.889
	$\alpha = 0.5$	0.309 (1.12)	0.003(0.01)	0.309 (1.12)	0.247 (0.90)	0.940	0.942	0.882	0.883
	$\alpha = 1$	0.305 (1.11)	0.008(0.03)	0.305 (1.11)	0.244 (0.89)	0.941	0.943	0.882	0.884
	$\alpha = 2$	0.306 (1.11)	$0.018\ (0.06)$	0.305 (1.11)	0.244 (0.89)	0.939	0.939	0.886	0.887
	$\alpha = 5$	0.307 (1.11)	0.029(0.10)	0.306 (1.11)	0.246(0.89)	0.950	0.950	0.898	0.895
	$\alpha = 10$	0.307 (1.11)	0.020 (0.07)	0.306 (1.11)	0.245 (0.89)	0.955	0.956	0.908	0.907
4800	M = 1	0.163 (1.18)	-0.001 (-0.01)	0.163 (1.18)	0.129 (0.94)	0.949	0.948	0.900	0.901
	M = 4	0.149 (1.08)	-0.000(-0.00)	0.149 (1.08)	0.118(0.86)	0.951	0.951	0.897	0.897
	M = 16	0.145 (1.05)	0.001 (0.01)	0.145 (1.05)	0.116(0.84)	0.952	0.950	0.903	0.903
	$\alpha = 0.5$	0.146 (1.06)	-0.000(-0.00)	0.146 (1.06)	0.117(0.85)	0.949	0.948	0.899	0.899
	$\alpha = 1$	0.145 (1.05)	0.001 (0.01)	0.145 (1.05)	0.116(0.84)	0.952	0.950	0.903	0.903
	$\alpha = 2$	0.145 (1.05)	0.006(0.04)	0.144 (1.05)	0.116(0.84)	0.953	0.953	0.906	0.907
	$\alpha = 5$	0.145 (1.05)	0.017(0.12)	0.144 (1.05)	0.116(0.84)	0.957	0.957	0.906	0.903
	$\alpha = 10$	0.147 (1.06)	0.027 (0.19)	0.144 (1.05)	0.117 (0.85)	0.958	0.958	0.909	0.910
9600	M = 1	0.115 (1.18)	-0.003(-0.03)	0.115 (1.18)	0.092 (0.94)	0.951	0.952	0.897	0.896
	M = 4	0.106(1.08)	-0.002(-0.02)	0.105 (1.08)	0.084(0.86)	0.950	0.950	0.901	0.901
	M = 16	0.103 (1.06)	-0.001 (-0.01)	0.103 (1.06)	0.082 (0.84)	0.948	0.948	0.902	0.902
	$\alpha = 0.5$	0.104 (1.07)	-0.001(-0.01)	0.104 (1.07)	0.082 (0.85)	0.953	0.951	0.904	0.905
	$\alpha = 1$	0.103 (1.06)	-0.001(-0.01)	0.103 (1.06)	0.082 (0.84)	0.950	0.950	0.904	0.903
	$\alpha = 2$	0.103 (1.05)	0.001 (0.02)	0.103(1.05)	0.082 (0.84)	0.954	0.953	0.906	0.906
	$\alpha = 5$	0.104(1.07)	$0.010\ (0.11)$	0.103(1.06)	0.083(0.85)	0.951	0.950	0.899	0.898
	$\alpha = 10$	0.106 (1.09)	$0.020\ (0.20)$	0.104(1.07)	0.084(0.87)	0.951	0.951	0.899	0.900

TABLE II Simulation results, Shadish, Clark, and Steiner (2008), $\sigma^* = 3.85$.

						95% Coverag		90% Coverage	
n	M	RMSE	Bias	SD	MAE	SE	AISE	SE	AISE
600	M = 1	0.190 (1.21)	0.005 (0.03)	0.190 (1.21)	0.152 (0.97)	0.879	0.942	0.806	0.888
	M = 4	0.182 (1.16)	0.007(0.04)	0.182 (1.16)	0.144 (0.92)	0.879	0.926	0.799	0.875
		0.180 (1.14)	0.011(0.07)	0.179 (1.14)	0.143(0.91)	0.865	0.921	0.785	0.866
	$\alpha = 0.5$	0.185 (1.18)	0.006(0.04)	0.185 (1.18)	0.147(0.94)	0.875	0.932	0.799	0.880
	$\alpha = 1$	0.182 (1.16)	0.007(0.04)	0.182 (1.16)	0.144(0.92)	0.879	0.926	0.799	0.875
	$\alpha = 2$	0.180(1.15)	0.009(0.06)	0.180(1.14)	0.143(0.91)	0.870	0.922	0.798	0.866
	$\alpha = 5$	0.179 (1.14)	0.012(0.07)	0.179 (1.14)	0.143 (0.91)	0.863	0.922	0.780	0.863
	$\alpha = 10$	0.179 (1.14)	0.012 (0.08)	0.179 (1.14)	0.142 (0.91)	0.856	0.924	0.780	0.864
1200	M = 1	0.130 (1.17)	0.006 (0.05)	0.129 (1.16)	0.103 (0.93)	0.905	0.948	0.840	0.893
	M = 4	0.123 (1.11)	0.008(0.07)	0.123(1.11)	0.098(0.88)	0.898	0.934	0.822	0.879
	M = 16	0.122(1.09)	0.013(0.12)	0.121(0.19)	0.097(0.87)	0.892	0.931	0.818	0.878
	$\alpha = 0.5$	0.125 (1.13)	0.008(0.07)	0.125 (1.13)	0.100(0.90)	0.898	0.943	0.836	0.885
	$\alpha = 1$	0.123 (1.11)	0.008(0.07)	0.123(1.11)	0.098(0.88)	0.898	0.934	0.822	0.879
	$\alpha = 2$	0.122(1.10)	0.011(0.10)	0.121(1.09)	0.097(0.87)	0.896	0.934	0.819	0.877
	$\alpha = 5$	0.122 (1.10)	0.015(0.13)	0.121(1.09)	0.097(0.87)	0.889	0.932	0.816	0.875
	$\alpha = 10$	0.121 (1.09)	0.017 (0.16)	0.120 (1.08)	0.097 (0.87)	0.880	0.930	0.799	0.876
4800	M = 1	0.064 (1.15)	0.006 (0.11)	0.063 (1.14)	0.051 (0.91)	0.918	0.943	0.858	0.890
	M = 4	0.060(1.09)	0.007(0.12)	0.060(1.08)	0.048(0.87)	0.912	0.939	0.839	0.877
		0.060(1.08)	0.009(0.17)	0.059(1.07)	0.048(0.86)	0.902	0.926	0.825	0.865
	$\alpha = 0.5$	0.061(1.09)	0.006(0.12)	0.060(1.09)	0.048(0.87)	0.918	0.941	0.844	0.876
	$\alpha = 1$	0.060(1.08)	0.007(0.13)	0.060(1.07)	0.048(0.86)	0.908	0.933	0.838	0.870
	$\alpha = 2$	0.060(1.08)	0.009(0.16)	0.059(1.07)	0.048(0.86)	0.902	0.930	0.829	0.865
	$\alpha = 5$	0.060(1.08)	0.012(0.21)	0.059(1.06)	0.048(0.86)	0.895	0.920	0.824	0.861
	$\alpha = 10$	0.060 (1.09)	0.015 (0.26)	0.059 (1.05)	0.048 (0.87)	0.891	0.916	0.819	0.858
9600	M = 1	0.045 (1.14)	0.005 (0.14)	0.044 (1.13)	0.036 (0.91)	0.923	0.940	0.864	0.886
	M = 4	0.042(1.07)	0.006(0.15)	0.042(1.06)	0.034(0.86)	0.920	0.933	0.853	0.881
		0.042 (1.06)	0.008(0.20)	0.041(1.04)	0.033(0.85)	0.910	0.928	0.847	0.869
		0.042 (1.08)	0.006(0.15)	0.042 (1.06)	0.034 (0.86)	0.922	0.938	0.856	0.882
	$\alpha = 1$	0.042 (1.06)	0.006(0.16)	0.041(1.05)	0.033(0.85)	0.916	0.934	0.851	0.872
	$\alpha = 2$	0.042 (1.06)	0.008(0.20)	0.041 (1.04)	0.033(0.85)	0.912	0.929	0.847	0.869
	$\alpha = 5$	0.042 (1.07)	0.010(0.25)	0.041 (1.04)	0.034 (0.86)	0.902	0.922	0.842	0.862
	$\alpha = 10$	0.042 (1.08)	0.012 (0.31)	0.041 (1.03)	0.034 (0.86)	0.897	0.916	0.829	0.860

Monte Carlo repetitions. Here, we use the value $(\sigma^*)^2$ to approximate the semiparametric efficiency lower bound if the assumptions in Theorem 4.1 hold.

For the first data, two observations are in line:

- 1. Regardless of which n is chosen, picking $M = \lfloor \alpha n^{2/(2+d)} \rfloor$ with α set to be 1 nearly always achieves the smallest SD, RMSE, and MAE. The simulation results thus support our recommendation to increase M for achieving better statistical performance.
- 2. Although consistency is established under different requirements for M, the two considered asymptotic variance estimators (SE and AISE) both yield good empirical coverage rates. The coverage rates are both very close to the nominal ones when n is large and there is not much difference between the two.

Some similar observations can be found for the second data. Notably, although picking $M = \lfloor \alpha n^{2/(2+d)} \rfloor$ with $\alpha = 1$ is not achieving the smallest RMSE this time, its RMSE is very

close to the smallest. However, for the second data, AISE yields generally better coverage rates than SE, although SE's coverage rates are improving as *n* increases.

To conclude, the simulation results generally support (a) increasing M with the sample size n for minimizing the RMSE and (b) exploiting Abadie and Imbens's (2006) approach to estimating the asymptotic variance of $\hat{\tau}_{M}^{bc}$. For choosing the M, the simulation results favor $M = \lfloor \alpha n^{2/(2+d)} \rfloor$ with α selected to be 1, while calculating the theoretically optimal α is believed to be difficult and also is beyond the scope of this manuscript.

6. SOME FINAL REMARKS

Some alternative matching estimators can also achieve double robustness or semiparametric efficiency. Yang and Zhang (2023) proposed to use the NN matching based on the propensity score (Rosenbaum and Rubin (1983), Abadie and Imbens (2016)) and the prognostic score (Hansen (2008)) simultaneously, and established the double robustness of the resulting matching estimator. They focused on fixed M, and consequently, their estimator did not achieve semiparametric efficiency. Wang and Zubizarreta (2023) proposed a matching method based on integer programming to ensure global balance of the covariates, and established the efficiency of the resulting difference-in-means estimator. They focused on fixed M, and even with fixed M, their integer programming problem was computationally challenging compared with NN matching.

There are three additional questions addressed in Abadie and Imbens (2006, 2012). First, estimation of the average treatment effect on the treated (ATT) can be incorporated in the double robustness and double machine learning framework (Theorem 4.2) and matching framework (Theorem 4.1(ii)) in a similar way. Second, asymptotic Normality (with an additional asymptotic bias term) of $\hat{\tau}_M$ in general d can be established as Theorem 4.1(ii). Third, unbalanced designs with n_0 much larger than n_1 cannot be incorporated in the double robustness and double machine learning framework, but can be studied in the same way as Theorem 4.1(ii).

APPENDIX A: DENSITY RATIO ESTIMATION I: COMPUTATION

Additional Notation. For any two real sequences $\{a_n\}_{n=1}^{\infty}$ and $\{b_n\}_{n=1}^{\infty}$, write $a_n \lesssim b_n$ (or equivalently, $b_n \gtrsim a_n$) if there exists a universal constant C > 0 such that $a_n/b_n \leq C$ for all sufficiently large n, and write $a_n \prec b_n$ (or equivalently, $b_n \succ a_n$) if $a_n/b_n \to 0$ as n goes to infinity. We write $a_n \asymp b_n$ if both $a_n \lesssim b_n$ and $b_n \lesssim a_n$ hold. We write $a_n = O(b_n)$ if $|a_n| \lesssim b_n$ and $a_n = o(b_n)$ if $|a_n| \prec b_n$. Denote the closed ball in \mathbb{R}^d centered at x with radius δ by $a_n = b_n$. In the sequel, let $a_n = b_n$ constants whose actual values may change at different locations.

This section discusses implementation and establishes Property (P1) for the proposed estimator $\hat{r}_M(\cdot)$. To this end, we separately discuss two cases:

Case I: estimating only the values of $\widehat{r}_M(\cdot)$ at the observed data points X_1, \ldots, X_{N_0} . **Case II**: estimating the values of $\widehat{r}_M(\cdot)$ at both the observed data points X_1, \ldots, X_{N_0} and n new points $x_1, \ldots, x_n \in \mathbb{R}^d$.

Case I. In many applications, we are only interested in a functional of density ratios at observed sample points, that is, the values of $\Phi(\{r(X_i)\}_{i=1}^{N_0})$ for some given functions Φ defined on \mathbb{R}^{N_0} . Check, for example, in a slightly different but symmetric form—(3.3) for such an example on ATE estimation. To this end, it is natural to consider the plug-in estimator $\Phi(\{\widehat{r}_M(X_i)\}_{i=1}^{N_0})$, for which it suffices to compute the values of $\{\widehat{r}_M(X_i)\}_{i=1}^{N_0}$.

Algorithm 1: Density ratio estimators at sample points.

```
Input: \{X_i\}_{i=1}^{N_0}, \{Z_j\}_{j=1}^{N_1}, and M.

Output: \{\widehat{r}_M(X_i)\}_{i=1}^{N_0}.

Build a k-d tree using \{X_i\}_{i=1}^{N_0};

for j=1:N_1 do

Search the M-NNs of Z_j in \{X_i\}_{i=1}^{N_0} using the k-d tree;

Store the indices of the M-NNs of Z_j as S_j;

Count and store the number of occurrence in \bigcup_{j=1}^{N_1} S_j for each element in [\![N_0]\!], which is then \{K_M(X_i)\}_{i=1}^{N_0};

Obtain \{\widehat{r}_M(X_i)\}_{i=1}^{N_0} based on (2.4).
```

Built on the k-d tree structure (Bentley (1975)) for tracking NNs, Algorithm 1 outlines an easy to implement algorithm to simultaneously compute all the values of $\{\widehat{r}_M(X_i)\}_{i=1}^{N_0}$. This algorithm could be regarded as a direct extension of the celebrated Friedman–Bentley–Finkel algorithm (Friedman, Bentley, and Finkel (1977)) to the NN matching setting.

Case II. Suppose we are interested in estimating density ratios at both the observed and n new points in \mathbb{R}^d . A naive algorithm is then to insert each new point into observed points and perform Algorithm 1 in order. However, this algorithm is not ideal as the corresponding time complexity would be n times the complexity of Algorithm 1, which could be computationally heavy with a large number of new points.

Instead, we develop a more sophisticated implementation. Let the new points be $\{x_i\}_{i=1}^n$. Algorithm 2 computes all the values of $\{\widehat{r}_M(x_i)\}_{i=1}^n$ as well as $\{\widehat{r}_M(X_i)\}_{i=1}^{N_0}$. The key message delivered here is that, compared with the aforementioned naive implementation, in Algorithm 2 we only need to construct one single k-d tree; the matching elements are then categorized to two different sets, corresponding to those with regard to X_i 's and X_i 's, separately. Such an implementation is thus intuitively much more efficient.

Theorem A.1 below elaborates on the computational advantage of the proposed estimator.

THEOREM A.1: (1) The average time complexity of Algorithm 1 to compute all the values of $\{\widehat{r}_M(X_i)\}_{i=1}^{N_0}$ is $O((d+N_1M/N_0)N_0\log N_0)$.

(2) Assume $[x_i]_{i=1}^n$ are independent and identically distributed (i.i.d.) following ν_0 and are

(2) Assume $[x_i]_{i=1}^n$ are independent and identically distributed (i.i.d.) following v_0 and are independent of $[X_i]_{i=1}^{N_0}$. Then the average time complexity of Algorithm 2 to compute all the values of $\{\widehat{r}_M(x_i)\}_{i=1}^n$ and $\{\widehat{r}_M(X_i)\}_{i=1}^{N_0}$ is $O((d+N_1M/N_0)(N_0+n)\log(N_0+n))$.

REMARK A.1—Comparison to Non-NN-Based Estimators: Assuming $N_0 \approx N_1 \approx N$, it is worth noting that optimization-based methods are commonly of a time complexity $O(N^2)$ if not worse (Noshad et al. (2017)). They are thus less appealing in terms of handling gigantic data as was argued in, for example, astronomy (Lima et al. (2008), Kremer et al. (2015)) and big text analysis (Borgeaud et al. (2021)) applications.

REMARK A.2—Comparison to the Two-Step NN-Based Density Ratio Estimator: Regarding Case I, a direct calculation yields that the time complexity of the simple two-step NN-based method, which separately estimates f_1 and f_0 based on individual M-NN

Algorithm 2: Density ratio estimators at both sample and new points.

```
Input: \{X_i\}_{i=1}^{N_0}, \{Z_j\}_{j=1}^{N_1}, M, and new points \{x_i\}_{i=1}^n.
Output: \{\widehat{r}_M(X_i)\}_{i=1}^{N_0} and \{\widehat{r}_M(x_i)\}_{i=1}^n.
Build a k-d tree using \{X_i\}_{i=1}^{N_0} \cup \{x_i\}_{i=1}^n;
for j = 1 : N_1 do
      Set S_j and S'_j be two empty sets;
      while |S_i| < M do
       Search the mth NN of Z_j in \{X_i\}_{i=1}^{N_0} \cup \{x_i\}_{i=1}^n; if the mth NN of Z_j is in \{X_i\}_{i=1}^{N_0} then add the index into S_j;
      else
m \leftarrow m + 1;
Store the indices sets S_j and S'_j;
Count and store the number of occurrence in \bigcup_{i=1}^{N_1} S_i for each element in [N_0], which
```

is then $\{K_M(X_i)\}_{i=1}^{N_0}$. Count and store the number of occurrence in $\bigcup_{j=1}^{N_1} S_j'$ for each element in [n], which is then $\{K_M(x_i)\}_{i=1}^n$;

Obtain $\{\widehat{r}_M(X_i)\}_{i=1}^{N_0}$ and $\{\widehat{r}_M(x_i)\}_{i=1}^n$ based on (2.4).

density estimators, is $O(dN_0 \log N_0 + dN_1 \log N_1 + N_0 M \log N_0 + N_0 M \log N_1)$. It is thus of the same order as Algorithm 1 when $N_1 \times N_0$, while computationally heavier when $N_1 \prec N_0$. Regarding Case II, the time complexity of the simple two-step NN-based method is $O(dN_0 \log N_0 + dN_1 \log N_1 + (N_0 + n)M \log N_0 + (N_0 + n)M \log N_1)$. Thus, if n is of less or equal order of N_0 , it is of the same order when $N_1 \times N_0$, while computationally heavier than Algorithm 2 when $N_1 \prec N_0$.

REMARK A.3—Comparison to the one-step NN-based density ratio estimator in Noshad et al. (2017): To estimate f-divergence measures, Noshad et al. (2017) constructed another one-step NN-based estimator admitting the simple form: $\widehat{r}_{M}(x) =$ $(N_0/N_1)(\mathcal{M}_i/(\mathcal{N}_i+1))$, where \mathcal{N}_i and \mathcal{M}_i are the numbers of points in $\{X_i\}_{i=1}^{N_0}$ and $\{Z_i\}_{i=1}^{N_1}$ among the M NNs of x; cf. Noshad et al. (2017, equation (20)). For Case I, its time complexity is $O(d(N_0 + N_1)\log(N_0 + N_1) + N_0M\log(N_0 + N_1))$; while for Case II, it is $O(d(N_0 + N_1) \log(N_0 + N_1) + (N_0 + n)M \log(N_0 + N_1))$. Both are at the same order as the naive NN-based one, but unlike the naive approach, this estimator is indeed onestep. However, it is still theoretically unclear if this estimator is statistically efficient; see Remark B.4 ahead for more details.

APPENDIX B: DENSITY RATIO ESTIMATION II: THEORY

This section introduces the theory for density ratio estimation based on NN matching. To this end, before establishing detailed theoretical properties (e.g., consistency and the rate of convergence) for $\widehat{r}_M(\cdot)$, we first exhibit a lemma elaborating on the asymptotic L^p moments of $\nu_1(A_M(x))$, the ν_1 -measure of the catchment area. This novel result did not appear in Abadie and Imbens's analysis. It is also of independent interest in stochastic and computational geometry in light of Remark 2.1.

LEMMA B.1—Asymptotic L^p Moments of Catchment Areas's ν_1 -Measure: Assuming $M \log N_0/N_0 \to 0$ as $N_0 \to \infty$, we have $\lim_{N_0 \to \infty} (N_0/M) \mathbb{E}[\nu_1(A_M(x))] = r(x)$ holds for ν_0 -almost all x. If we further assume $M \to \infty$, then for any positive integer p, we have $\lim_{N_0 \to \infty} (N_0/M)^p \mathbb{E}[\nu_1^p(A_M(x))] = [r(x)]^p$ holds for ν_0 -almost all x.

REMARK B.1—Relation to the Measure of Voronoi Cells: When M=1 and $\nu_0=\nu_1$, the measure of catchment areas reduces to the measure of Voronoi cells as pointed out in Remark 2.1. Interestingly, in the stochastic geometry literature, Devroye, Györfi, Lugosi, and Walk (2017) studied a related problem of bounding the moments of the measure of Voronoi cells (cf. Theorem 2.1 therein). Setting M=1 and $\nu_0=\nu_1$ in the first part of Lemma B.1 and recalling Remark 2.1, we can derive their Theorem 2.1(i). On the other hand, Devroye et al. (2017, Theorem 2.1(ii)) showed that when $\nu_0=\nu_1$, p=2, and $d\leq 3$, $(M^{-1}N_0)^2\mathrm{E}[\nu_1^2(A_M(x))]$ converges to 1 whereas $N_0^2\mathrm{E}[\nu_1^2(A_1(x))]$ does not; cf. Devroye et al. (2017, Section 4.2). This supports the necessity of forcing $M\to\infty$ for stabilizing the moments of $\widehat{r}_M(\cdot)$.

B.1. Consistency

We first establish the pointwise consistency of the estimator $\hat{r}_M(x)$ for r(x). This requires nearly no assumption on ν_0 , ν_1 except for those made at the beginning of Section 2, in line with similar observations made in NN-based density estimation (Biau and Devroye (2015, Theorem 3.1)).

THEOREM B.1—Pointwise Consistency: Assume $M \log N_0/N_0 \to 0$ as $N_0 \to \infty$.

- (i) (Asymptotic unbiasedness) For ν_0 -almost all x, we have $\lim_{N_0 \to \infty} \mathbb{E}[\widehat{r}_M(x)] = r(x)$.
- (ii) (Pointwise L_p consistency) Let p be any positive integer and assume further that $MN_1/N_0 \to \infty$ and $M \to \infty$ as $N_0 \to \infty$. Then for ν_0 -almost all x, we have $\lim_{N_0 \to \infty} \mathbb{E}[|\widehat{r}_M(x) r(x)|^p] = 0$.

For evaluating the global consistency of the estimator, it is necessary to introduce the following (global) L_p risk:

$$L_p \text{ risk} = \mathbb{E}[|\widehat{r}_M(X) - r(X)|^p | X_1, \dots, X_{N_0}, Z_1, \dots, Z_{N_1}] = \int_{\mathbb{R}^d} |\widehat{r}_M(x) - r(x)|^p f_0(x) dx,$$

where X is a copy drawn from ν_0 that is independent of the data. For the L_p risk consistency of the estimator, we impose conditions on ν_0 and ν_1 further as follows.

Define the supports of ν_0 and ν_1 as S_0 and S_1 , respectively. For any set $S \subset \mathbb{R}^d$, define the diameter of S as diam(S) = $\sup_{x,z \in S} ||x - z||$.

Assumption B.1: (i) ν_0 , ν_1 are two probability measures on \mathbb{R}^d , both are absolutely continuous with respect to λ , and ν_1 is absolutely continuous with respect to ν_0 .

- (ii) There exists a constant R > 0 such that $diam(S_0) \le R$.
- (iii) There exist two constants f_L , $f_U > 0$ such that for any $x \in S_0$ and $z \in S_1$, $f_L \le f_0(x) \le f_U$ and $f_1(z) \le f_U$.
- (iv) There exists a constant $a \in (0, 1)$ such that for any $\delta \in (0, \text{diam}(S_0)]$ and $z \in S_1$, $\lambda(B_{z,\delta} \cap S_0) \ge a\lambda(B_{z,\delta})$, recalling that $B_{z,\delta}$ represents the closed ball in \mathbb{R}^d with center at z and radius δ .

REMARK B.2: Assumption B.1 is standard in the literature for establishing the global consistency of density ratio estimators. The regularity conditions on the support ensure that the angle of the support is not too sharp, which trivially hold for any *d*-dimensional cube. These conditions were also enforced in Nguyen, Wainwright, and Jordan (2010, Theorem 1), Sugiyama, Suzuki, Nakajima, Kashima, von Bünau, and Kawanabe (2008, Assumption 1), Kpotufe (2017, Definition 1), among many others.

We then establish the L_p risk consistency of the estimator via the Hardy–Littlewood maximal inequality (Stein (2016)); cf. Lemma S3.2 in the Online Appendix. Of note, this inequality was used in Han, Jiao, Weissman, and Wu (2020) in a relative manner in order to study the information-theoretic limit of entropy estimation.

THEOREM B.2— L_p Risk Consistency: Assume the pair of ν_0 , ν_1 satisfies Assumption B.1. Let p be any positive integer. Assume further that $M \log N_0/N_0 \to 0$, $MN_1/N_0 \to \infty$, and $M \to \infty$ as $N_0 \to \infty$. We then have

$$\lim_{N_0 \to \infty} \mathbb{E} \left[\int_{\mathbb{R}^d} \left| \widehat{r}_M(x) - r(x) \right|^p f_0(x) \, \mathrm{d}x \right] = 0.$$

As a direct corollary of Theorem B.2, one can obtain the limit of any finite moment of $\nu_1(A_M(\cdot))$ with a random center. This can be regarded as a global extension to Lemma B.1.

COROLLARY B.1: Assume the same conditions as in Theorem B.2. We then have $\lim_{N_0\to\infty}(N_0/M)^p\mathrm{E}[\nu_1^p(A_M(W))]=\mathrm{E}([r(W)]^p)$, where W follows an arbitrary distribution that is absolutely continuous with respect to ν_0 and has density bounded above and below by two positive constants. In particular, it holds when W is drawn from ν_0 .

B.2. Rates of Convergence

In this section, we establish the rates of convergence for $\hat{r}(x)$ under both pointwise and global measures. We first consider the pointwise mean square error (MSE) convergence rate and show that $\hat{r}_M(\cdot)$ is minimax optimal in that regard. In the sequel, we fix an $x \in \mathbb{R}^d$ and consider the following local assumption on (ν_0, ν_1) .

- ASSUMPTION B.2—Local Assumption: (i) ν_0 , ν_1 are two probability measures on \mathbb{R}^d , both are absolutely continuous with respect to λ , and ν_1 is absolutely continuous with respect to ν_0 .
- (ii) There exist two constants f_L , $f_U > 0$ such that $f_0(x) \ge f_L$ and $f_1(x) \le f_U$.
- (iii) There exists a constant $\delta > 0$ such that for any $z \in B_{x,\delta}$, $|f_0(x) f_0(z)| \vee |f_1(x) f_1(z)| \le L||x z||$ for some constants L > 0.

Define the following probability class:

$$\mathcal{P}_{x,p}(f_L, f_U, L, d, \delta) = \{(\nu_0, \nu_1) : \text{Assumption B.2 holds}\}.$$

The following theorem establishes the uniform pointwise convergence rate of $\widehat{r}_M(\cdot)$.

THEOREM B.3—Pointwise Rates of Convergence: Assume $M \log N_0/N_0 \to 0$ and $M/\log N_0 \to \infty$ as $N_0 \to \infty$. Consider a sufficiently large N_0 .

(i) Asymptotic bias:

$$\sup_{(\nu_0,\nu_1)\in\mathcal{P}_{x,\mathrm{p}}(f_L,f_U,L,d,\delta)} \left|\mathrm{E}\big[\widehat{r}_M(x)\big] - r(x)\right| \leq C \bigg(\frac{M}{N_0}\bigg)^{1/d},$$

where C > 0 is a constant only depending on f_L , f_U , L, d. Further assume $MN_1/N_0 \to \infty$ as $N_0 \to \infty$.

(ii) Asymptotic variance:

$$\sup_{(\nu_0,\nu_1)\in\mathcal{P}_{x,p}(f_L,f_U,L,d,\delta)} \operatorname{Var}\left[\widehat{r}_M(x)\right] \leq C' \left[\left(\frac{1}{M}\right) + \left(\frac{N_0}{MN_1}\right)\right],$$

where C' > 0 is a constant only depending on f_L , f_U .

(iii) Asymptotic MSE:

$$\sup_{(\nu_0,\nu_1)\in\mathcal{P}_{x,p}(f_L,f_U,L,d,\delta)} \mathbf{E}\big[\widehat{r}_M(x)-r(x)\big]^2 \leq C''\bigg[\bigg(\frac{M}{N_0}\bigg)^{2/d}+\bigg(\frac{1}{M}\bigg)+\bigg(\frac{N_0}{MN_1}\bigg)\bigg],$$

where C'' > 0 is a constant only depending on f_L , f_U , L, d.

Further assume $N_1^{-\frac{d}{2+d}} \log N_0 \to 0$ as $N_0 \to \infty$.

(iv) Fix $\alpha > 0$ and take $M = \alpha \cdot \{N_0^{\frac{2}{2+d}} \vee (N_0 N_1^{-\frac{d}{2+d}})\}$. We have

$$\sup_{(\nu_0,\nu_1)\in\mathcal{P}_{x,p}(f_L,f_U,L,d,\delta)} \mathbb{E}\big[\widehat{r}_M(x) - r(x)\big]^2 \le C'''(N_0 \wedge N_1)^{-\frac{2}{2+d}}, \tag{B.1}$$

where C''' > 0 is a constant only depending on f_L , f_U , L, d, α .

The rate of convergence in (B.1) matches the established minimax lower bound in Lipschitz density function estimation (Tsybakov (2009, Section 2)). By some simple manipulation, the argument in Tsybakov (2009, Exercise 2.8) directly extends to density ratio as the latter is a harder statistical problem (Kpotufe (2017, Remark 3)). This is formally stated in the following proposition.

PROPOSITION B.1—Pointwise MSE minimax lower bound: For sufficiently large N_0 and N_1 ,

$$\inf_{\widetilde{r}} \sup_{(\nu_0,\nu_1)\in\mathcal{P}_{X,\mathbb{D}}(f_L,f_U,L,d,\delta)} \mathbb{E}\big[\widetilde{r}(x)-r(x)\big]^2 \geq c(N_0\wedge N_1)^{-\frac{2}{2+d}},$$

where c > 0 is a constant only depending on f_L , f_U , L, d, and the infimum is taken over all measurable functions.

We then move on to the global risk and study the rates of convergence. To this end, a global assumption on (ν_0, ν_1) is given below.

ASSUMPTION B.3—Global Assumption: (i) ν_0 , ν_1 are two probability measures on \mathbb{R}^d , both are absolutely continuous with respect to λ , and ν_1 is absolutely continuous with respect to ν_0 .

- (ii) There exists a constant R > 0 such that $diam(S_0) \le R$.
- (iii) There exist two constants f_L , $f_U > 0$ such that for any $x \in S_0$ and $z \in S_1$, $f_L \le f_0(x) \le 1$ f_U and $f_1(z) \leq f_U$.
- (iv) There exists a constant $a \in (0, 1)$ such that for any $\delta \in (0, \text{diam}(S_0)]$ and any $z \in S_1$, $\lambda(B_{z,\delta}\cap S_0)\geq a\lambda(B_{z,\delta}).$
- (v) There exists a constant H > 0 such that the surface area (Hausdorff measure, Evans and Garzepy (2018, Section 3.3)) of S_1 is bounded by H.
- (vi) There exists a constant L > 0 such that for any $x, z \in S_1$, $|f_0(x) f_0(z)| \vee |f_1(x) f_1(x)| = 0$ $|f_1(z)| \le L||x-z||.$

REMARK B.3: Assumption B.3 is standard in the literature for establishing the global risk of density ratio estimators; similar assumptions were made in Zhao and Lai (2022, Assumption 1) and Zhao and Lai (2020, Assumption 1). Note that the regularity conditions on the support automatically hold for d-dimensional cubes, and the restriction on the surface area is added to control the boundary effect on NN-based methods.

Define the following probability class:

$$\mathcal{P}_{g}(f_{L}, f_{U}, L, d, a, H, R) = \{(\nu_{0}, \nu_{1}) : \text{Assumption B.3 holds}\}. \tag{B.2}$$

The next theorem establishes the uniform rate of convergence of $\hat{r}(\cdot)$ within the above probability class under the L_1 risk. This rate is further matched by a minimax lower bound derived in Theorem 1 of Zhao and Lai (2022) using similar arguments as in the pointwise case.

THEOREM B.4—Global rates of convergence under the L_1 risk: Assume $M \log N_0$ $N_0 \to 0$, $M/\log N_0 \to \infty$, $MN_1/N_0 \to \infty$ as $N_0 \to \infty$. Consider a sufficiently large N_0 . (i) We have the following uniform upper bound:

$$\sup_{(\nu_0,\nu_1)\in\mathcal{P}_{g}(f_L,f_U,L,d,a,H,R)} \mathbb{E}\left[\int_{\mathbb{R}^d} |\widehat{r}_{M}(x) - r(x)| f_0(x) \, \mathrm{d}x\right]$$

$$\leq C \left[\left(\frac{M}{N_0}\right)^{1/d} + \left(\frac{1}{M}\right)^{1/2} + \left(\frac{N_0}{MN_1}\right)^{1/2} \right],$$

where C > 0 is a constant only depending on f_L , f_U , a, H, L, d.

(ii) Further assume $N_1^{-\frac{d}{2+d}} \log N_0 \to 0$ as $N_0 \to \infty$, fix $\alpha > 0$, and take $M = \alpha \cdot \{N_0^{\frac{2}{2+d}} \lor \{N_0^{\frac{2}{2+$ $(N_0 N_1^{-\frac{d}{2+d}})$ }. We then have

$$\sup_{(\nu_0,\nu_1)\in\mathcal{P}_{\mathbf{g}}(f_L,f_U,L,d,a,H,R)} \mathbf{E} \left[\int_{\mathbb{R}^d} \left| \widehat{r}_M(x) - r(x) \right| f_0(x) \, \mathrm{d}x \right] \leq C' (N_0 \wedge N_1)^{-\frac{1}{2+d}},$$

where C' > 0 is a constant only depending on f_L , f_U , a, H, L, d, α .

PROPOSITION B.2—Global Minimax Lower Bound Under the L_1 Risk: If a is sufficiently small and H, R are sufficiently large, then for sufficiently large N_0 and N_1 ,

$$\inf_{\widetilde{r}} \sup_{(\nu_0,\nu_1)\in\mathcal{P}_{\mathbf{Z}}(f_L,f_U,L,d,a,H,R)} \mathrm{E}\bigg[\int_{\mathbb{R}^d} \big|\widetilde{r}(x)-r(x)\big|f_0(x)\,\mathrm{d}x\bigg] \geq c(N_0\wedge N_1)^{-\frac{1}{2+d}},$$

where c > 0 is a constant only depending on f_L , f_U , L, d, and the infimum is taken over all measurable functions.

REMARK B.4—Comparison to the One-Step Estimator in Noshad et al. (2017): The estimator introduced in Remark A.3 by Noshad et al. (2017) is to our knowledge the only alternative density ratio estimator in the literature that is able to attain both the property (P1) and being one step. However, the arguments in Noshad et al. (2017, Section III) can only yield the bound $\mathrm{E}[\widehat{r}_M'(x)-r(x)]^2\lesssim (M/N_0)^{1/d}+M^{-1}$ for $(\nu_0,\nu_1)\in\mathcal{P}_{x,p}(f_L,f_U,L,d,\delta)$. This is via equation (21) therein, de-Poissonizing the estimator, and further assuming N_1/N_0 converges to a positive constant. The above bound is strictly looser than the bound $(M/N_0)^{2/d}+M^{-1}$ for $\widehat{r}_M(\cdot)$ shown in Theorem B.3. However, it seems mathematically challenging to improve their analysis and accordingly, unlike $\widehat{r}_M(\cdot)$, it is still theoretically unclear if the estimator $\widehat{r}_M'(x)$ is a statistically efficient density ratio estimator.

APPENDIX C: PROOFS OF THE RESULTS IN SECTIONS 3 AND 4

C.1. Proof of Lemma 3.1

PROOF OF LEMMA 3.1: By simple algebra, we have

$$\begin{split} \widehat{\tau}_{M}^{\text{bc}} &= \frac{1}{n} \sum_{i=1}^{n} \left[\widehat{Y}_{i}^{\text{bc}}(1) - \widehat{Y}_{i}^{\text{bc}}(0) \right] \\ &= \frac{1}{n} \sum_{i=1,D_{i}=1}^{n} \left[Y_{i} - \frac{1}{M} \sum_{j \in \mathcal{J}_{M}(i)} \left(Y_{j} + \widehat{\mu}_{0}(X_{i}) - \widehat{\mu}_{0}(X_{j}) \right) \right] \\ &+ \frac{1}{n} \sum_{i=1,D_{i}=0}^{n} \left[\frac{1}{M} \sum_{j \in \mathcal{J}_{M}(i)} \left(Y_{j} + \widehat{\mu}_{1}(X_{i}) - \widehat{\mu}_{1}(X_{j}) \right) - Y_{i} \right] \\ &= \frac{1}{n} \sum_{i=1,D_{i}=1}^{n} \left[\widehat{R}_{i} + \widehat{\mu}_{1}(X_{i}) - \widehat{\mu}_{0}(X_{i}) - \frac{1}{M} \sum_{j \in \mathcal{J}_{M}(i)} \widehat{R}_{j} \right] \\ &+ \frac{1}{n} \sum_{i=1,D_{i}=0}^{n} \left[\frac{1}{M} \sum_{j \in \mathcal{J}_{M}(i)} \widehat{R}_{j} - \widehat{R}_{i} + \widehat{\mu}_{1}(X_{i}) - \widehat{\mu}_{0}(X_{i}) \right] \\ &= \frac{1}{n} \sum_{i=1}^{n} \left[\widehat{\mu}_{1}(X_{i}) - \widehat{\mu}_{0}(X_{i}) \right] + \frac{1}{n} \left[\sum_{i=1,D_{i}=1}^{n} \left(1 + \frac{K_{M}(i)}{M} \right) \widehat{R}_{i} - \sum_{i=1,D_{i}=0}^{n} \left(1 + \frac{K_{M}(i)}{M} \right) \widehat{R}_{i} \right]. \end{split}$$

This completes the proof.

C.2. Proof of Theorem 4.1

PROOF OF THEOREM 4.1(i): **Part I.** Suppose the density function is sufficiently smooth. For any $i \in [n]$, let $\bar{R}_i = Y_i - \bar{\mu}_{D_i}(X_i)$. From (3.3),

$$\widehat{\tau}_{M}^{bc} = \widehat{\tau}^{reg} + \frac{1}{n} \left[\sum_{i=1}^{n} D_{i} \left(1 + \frac{K_{M}(i)}{M} \right) \widehat{R}_{i} - \sum_{i=1}^{n} (1 - D_{i}) \left(1 + \frac{K_{M}(i)}{M} \right) \widehat{R}_{i} \right] \\
= \frac{1}{n} \sum_{i=1}^{n} \left[\widehat{\mu}_{1}(X_{i}) - \overline{\mu}_{1}(X_{i}) \right] - \frac{1}{n} \sum_{i=1}^{n} \left[\widehat{\mu}_{0}(X_{i}) - \overline{\mu}_{0}(X_{i}) \right] \\
+ \frac{1}{n} \left[\sum_{i=1}^{n} (2D_{i} - 1) \left(1 + \frac{K_{M}(i)}{M} \right) \left(\overline{\mu}_{D_{i}}(X_{i}) - \widehat{\mu}_{D_{i}}(X_{i}) \right) \right] \\
+ \frac{1}{n} \left[\sum_{i=1}^{n} D_{i} \left(1 + \frac{K_{M}(i)}{M} - \frac{1}{e(X_{i})} \right) \overline{R}_{i} \right] \\
- \sum_{i=1}^{n} (1 - D_{i}) \left(1 + \frac{K_{M}(i)}{M} - \frac{1}{1 - e(X_{i})} \right) \overline{R}_{i} \right] \\
+ \frac{1}{n} \left[\sum_{i=1}^{n} \left(1 - \frac{D_{i}}{e(X_{i})} \right) \overline{\mu}_{1}(X_{i}) - \sum_{i=1}^{n} \left(1 - \frac{1 - D_{i}}{1 - e(X_{i})} \right) \overline{\mu}_{0}(X_{i}) \right] \\
+ \frac{1}{n} \left[\sum_{i=1}^{n} \frac{D_{i}}{e(X_{i})} Y_{i} - \sum_{i=1}^{n} \frac{1 - D_{i}}{1 - e(X_{i})} Y_{i} \right]. \tag{C.1}$$

For each pair of terms, we only establish the first half part under treatment, and the second half under control can be established in the same way.

For the first term in (C.1),

$$\left| \frac{1}{n} \sum_{i=1}^{n} \left[\widehat{\mu}_{1}(X_{i}) - \bar{\mu}_{1}(X_{i}) \right] \right| \leq \|\widehat{\mu}_{1} - \bar{\mu}_{1}\|_{\infty} = o_{P}(1). \tag{C.2}$$

For the second term in (C.1),

$$\left| \frac{1}{n} \sum_{i=1}^{n} D_{i} \left(1 + \frac{K_{M}(i)}{M} \right) \left(\bar{\mu}_{1}(X_{i}) - \widehat{\mu}_{1}(X_{i}) \right) \right|$$

$$\leq \|\widehat{\mu}_{1} - \bar{\mu}_{1}\|_{\infty} \frac{1}{n} \sum_{i=1}^{n} D_{i} \left(1 + \frac{K_{M}(i)}{M} \right) = \|\widehat{\mu}_{1} - \bar{\mu}_{1}\|_{\infty} = o_{P}(1), \tag{C.3}$$

where the last step is due to $\sum_{i=1}^{n} D_i K_M(i) = n_0 M$.

Notice that from Assumption 4.1(i), $P_{X|D=0}$ and $P_{X|D=1}$ share the same support, and their densities are both bounded and bounded away from zero as long as one is. Then $(P_{X|D=0}, P_{X|D=1})$ and $(P_{X|D=1}, P_{X|D=0})$ both satisfy Assumption B.1 as long as one satisfies.

For the third term in (C.1), by Theorem B.2,

$$\left\{ E \left[\left| \frac{1}{n} \sum_{i=1}^{n} D_{i} \left(1 + \frac{K_{M}(i)}{M} - \frac{1}{e(X_{i})} \right) \bar{R}_{i} \right| \right] \right\}^{2} \\
\leq \left\{ E \left[\left| D_{i} \left(1 + \frac{K_{M}(i)}{M} - \frac{1}{e(X_{i})} \right) \bar{R}_{i} \right| \right] \right\}^{2} \\
\leq E \left[1 + \frac{K_{M}(i)}{M} - \frac{1}{e(X_{i})} \right]^{2} E[D_{i} \bar{R}_{i}]^{2} \\
= E \left[1 + \frac{K_{M}(i)}{M} - \frac{1}{e(X_{i})} \right]^{2} E[D_{i} (Y_{i}(1) - \bar{\mu}_{1}(X_{i}))^{2}] \\
\leq E \left[1 + \frac{K_{M}(i)}{M} - \frac{1}{e(X_{i})} \right]^{2} E[\sigma_{1}^{2}(X_{i}) + (\mu_{1}(X_{i}) - \bar{\mu}_{1}(X_{i}))^{2}] = o(1), \quad (C.4)$$

where $\sigma_1^2(x) = \mathbb{E}[U_1^2 | X = x]$ for $x \in \mathbb{X}$. For the fourth term in (C.1), notice that

$$E\left[\frac{1}{n}\sum_{i=1}^{n}\left(1-\frac{D_i}{e(X_i)}\right)\bar{\mu}_1(X_i)\,\bigg|\,X_1,\ldots,X_n\right]=0,$$

and

$$\operatorname{Var}\left[\frac{1}{n}\sum_{i=1}^{n}\left(1-\frac{D_{i}}{e(X_{i})}\right)\bar{\mu}_{1}(X_{i})\right]$$

$$=\operatorname{E}\left[\operatorname{Var}\left[\frac{1}{n}\sum_{i=1}^{n}\left(1-\frac{D_{i}}{e(X_{i})}\right)\bar{\mu}_{1}(X_{i})\,\middle|\,X_{1},\ldots,X_{n}\right]\right]$$

$$=\frac{1}{n}\operatorname{E}\left[\bar{\mu}_{1}^{2}(X_{i})\left(\frac{1}{e(X_{i})}-1\right)\right]=O(n^{-1}).$$

Then

$$\frac{1}{n} \sum_{i=1}^{n} \left(1 - \frac{D_i}{e(X_i)} \right) \bar{\mu}_1(X_i) = o_{\mathbb{P}}(1). \tag{C.5}$$

For the fifth term in (C.1), notice that $E[Y^2]$ are bounded and $[(X_i, D_i, Y_i)]_{i=1}^n$ are i.i.d. Using the weak law of large numbers yields

$$\frac{1}{n} \left[\sum_{i=1}^{n} \frac{D_i}{e(X_i)} Y_i - \sum_{i=1}^{n} \frac{1 - D_i}{1 - e(X_i)} Y_i \right] \xrightarrow{p} E[Y_i(1) - Y_i(0)] = \tau.$$
 (C.6)

Plugging (C.2), (C.3), (C.4), (C.5), (C.6) into (C.1) completes the proof.

Part II. Suppose the outcome model is correct. By (3.3),

$$\widehat{\tau}_{M}^{bc} = \widehat{\tau}^{reg} + \frac{1}{n} \left[\sum_{i=1}^{n} D_{i} \left(1 + \frac{K_{M}(i)}{M} \right) \widehat{R}_{i} - \sum_{i=1}^{n} (1 - D_{i}) \left(1 + \frac{K_{M}(i)}{M} \right) \widehat{R}_{i} \right] \\
= \frac{1}{n} \sum_{i=1}^{n} \left[\widehat{\mu}_{1}(X_{i}) - \mu_{1}(X_{i}) \right] - \frac{1}{n} \sum_{i=1}^{n} \left[\widehat{\mu}_{0}(X_{i}) - \mu_{0}(X_{i}) \right] \\
+ \frac{1}{n} \left[\sum_{i=1}^{n} (2D_{i} - 1) \left(1 + \frac{K_{M}(i)}{M} \right) \left(\mu_{D_{i}}(X_{i}) - \widehat{\mu}_{D_{i}}(X_{i}) \right) \right] \\
+ \frac{1}{n} \left[\sum_{i=1}^{n} D_{i} \left(1 + \frac{K_{M}(i)}{M} \right) \left(Y_{i} - \mu_{1}(X_{i}) \right) \right] \\
- \sum_{i=1}^{n} (1 - D_{i}) \left(1 + \frac{K_{M}(i)}{M} \right) \left(Y_{i} - \mu_{0}(X_{i}) \right) \right] \\
+ \frac{1}{n} \sum_{i=1}^{n} \left[\mu_{1}(X_{i}) - \mu_{0}(X_{i}) \right]. \tag{C.7}$$

For the first term in (C.7),

$$\left| \frac{1}{n} \sum_{i=1}^{n} \left[\widehat{\mu}_{1}(X_{i}) - \mu_{1}(X_{i}) \right] \right| \leq \|\widehat{\mu}_{1} - \mu_{1}\|_{\infty} = o_{P}(1).$$
 (C.8)

For the second term in (C.7),

$$\left| \frac{1}{n} \sum_{i=1}^{n} D_{i} \left(1 + \frac{K_{M}(i)}{M} \right) \left(\mu_{1}(X_{i}) - \widehat{\mu}_{1}(X_{i}) \right) \right|
\leq \|\widehat{\mu}_{1} - \mu_{1}\|_{\infty} \frac{1}{n} \sum_{i=1}^{n} D_{i} \left(1 + \frac{K_{M}(i)}{M} \right) = \|\widehat{\mu}_{1} - \mu_{1}\|_{\infty} = o_{P}(1).$$
(C.9)

For the third term in (C.7), noticing that $K_M(\cdot)$ is a function of (X_1, \ldots, X_n) and (D_1, \ldots, D_n) , we can obtain

$$E\left[\frac{1}{n}\sum_{i=1}^{n}D_{i}\left(1+\frac{K_{M}(i)}{M}\right)\left(Y_{i}-\mu_{1}(X_{i})\right)\,\middle|\,X_{1},\ldots,X_{n},D_{1},\ldots,D_{n}\right]=0.$$

By a martingale representation (Abadie and Imbens (2012)) and then the martingale convergence theorem (which holds for both fixed and diverging M), we obtain

$$\frac{1}{n} \sum_{i=1}^{n} D_i \left(1 + \frac{K_M(i)}{M} \right) (Y_i - \mu_1(X_i)) = o_P(1).$$
 (C.10)

For the fourth term in (C.7), notice that $E[\mu_{\omega}^2(X)]$ is bounded for $\omega \in \{0, 1\}$. Using the weak law of large numbers, we obtain

$$\frac{1}{n} \sum_{i=1}^{n} \left[\mu_1(X_i) - \mu_0(X_i) \right] \xrightarrow{p} \mathrm{E} \left[\mu_1(X_i) - \mu_0(X_i) \right] = \tau. \tag{C.11}$$

Plugging
$$(C.8)$$
, $(C.9)$, $(C.10)$, $(C.11)$ into $(C.7)$ completes the proof. Q.E.D.

PROOF OF THEOREM 4.1(ii): For $\omega \in \{0, 1\}$ and $m \in \llbracket M \rrbracket$, let $j_m(i)$ represent the index of mth-NN of X_i in $\{X_j : D_j = 1 - D_i\}_{j=1}^n$, that is, the index $j \in \llbracket n \rrbracket$ such that $D_j = 1 - D_i$ and $\sum_{\ell=1,D_\ell=1-D_i}^n \mathbb{1}(\lVert X_\ell - X_i\rVert \leq \lVert X_j - X_i\rVert) = m$. With a little abuse of notation, let $\epsilon_i = Y_i - \mu_{D_i}(X_i)$ for any $i \in \llbracket n \rrbracket$. By the definition of $\widehat{\tau}_M^{\text{bc}}$ in (3.3), we can verify the decomposition $\widehat{\tau}_M^{\text{bc}} = \overline{\tau}(X) + E_M + B_M - \widehat{B}_M$, where

$$\begin{split} \bar{\tau}(X) &= \frac{1}{n} \sum_{i=1}^{n} \left[\mu_{1}(X_{i}) - \mu_{0}(X_{i}) \right], \\ E_{M} &= \frac{1}{n} \sum_{i=1}^{n} (2D_{i} - 1) \left(1 + \frac{K_{M}(i)}{M} \right) \epsilon_{i}, \\ B_{M} &= \frac{1}{n} \sum_{i=1}^{n} (2D_{i} - 1) \left[\frac{1}{M} \sum_{m=1}^{M} \left(\mu_{1-D_{i}}(X_{i}) - \mu_{1-D_{i}}(X_{j_{m}(i)}) \right) \right], \\ \widehat{B}_{M} &= \frac{1}{n} \sum_{i=1}^{n} (2D_{i} - 1) \left[\frac{1}{M} \sum_{m=1}^{M} \left(\widehat{\mu}_{1-D_{i}}(X_{i}) - \widehat{\mu}_{1-D_{i}}(X_{j_{m}(i)}) \right) \right]. \end{split}$$

We have the following central limit theorem on $\bar{\tau}(X) + E_M$.

LEMMA C.1:
$$\sqrt{n}\sigma^{-1/2}(\bar{\tau}(X) + E_M - \tau) \stackrel{d}{\longrightarrow} N(0, 1).$$

For the bias term $B_M - \widehat{B}_M$, define $U_{m,i} = X_{j_m(i)} - X_i$ for any $i \in [n]$ and $m \in [M]$. We then have the following lemma bounding the moments of $U_{M,i}$.

LEMMA C.2: Let p be any positive integer. Then there exists a constant $C_p > 0$ only depending on p such that for any $i \in [n]$ and $M \in [n_{1-D_i}]$,

$$E[\|U_{M,i}\|^p | D_1, \ldots, D_n] \le C_p (M/n_{1-D_i})^{p/d}.$$

In light of the smoothness conditions on μ_{ω} and approximation conditions on $\widehat{\mu}_{\omega}$ for $\omega \in \{0, 1\}$, we can establish the following lemma using Lemma C.2.

LEMMA C.3:
$$\sqrt{n}(B_M - \widehat{B}_M) \stackrel{p}{\longrightarrow} 0$$
.

Combining Lemma C.1 and Lemma C.3 completes the proof.

PROOF OF THEOREM 4.1(ii), CONSISTENCY OF $\hat{\sigma}^2$: By definition, we can verify the decomposition $\hat{\sigma}^2 - \sigma^2 = T_1 + T_2 + T_3 + T_4$, where

$$\begin{split} T_1 &= \frac{1}{n} \sum_{i=1}^n \left[\widehat{\mu}_1(X_i) - \widehat{\mu}_0(X_i) + D_i \left(1 + \frac{K_M(i)}{M} \right) \widehat{R}_i - (1 - D_i) \left(1 + \frac{K_M(i)}{M} \right) \widehat{R}_i - \widehat{\tau}_M^{\text{bc}} \right]^2 \\ &- \frac{1}{n} \sum_{i=1}^n \left[\mu_1(X_i) - \mu_0(X_i) + D_i \left(1 + \frac{K_M(i)}{M} \right) (Y_i - \mu_1(X_i)) \right. \\ &- (1 - D_i) \left(1 + \frac{K_M(i)}{M} \right) (Y_i - \mu_0(X_i)) - \widehat{\tau}_M^{\text{bc}} \right]^2, \\ T_2 &= \frac{1}{n} \sum_{i=1}^n \left[\mu_1(X_i) - \mu_0(X_i) + D_i \left(1 + \frac{K_M(i)}{M} \right) (Y_i - \mu_1(X_i)) \right. \\ &- (1 - D_i) \left(1 + \frac{K_M(i)}{M} \right) (Y_i - \mu_0(X_i)) - \widehat{\tau}_M^{\text{bc}} \right]^2 \\ &- \frac{1}{n} \sum_{i=1}^n \left[\mu_1(X_i) - \mu_0(X_i) + \frac{D_i}{e(X_i)} (Y_i - \mu_1(X_i)) \right. \\ &- \frac{1 - D_i}{1 - e(X_i)} (Y_i - \mu_0(X_i)) - \widehat{\tau}_M^{\text{bc}} \right]^2, \\ T_3 &= \frac{1}{n} \sum_{i=1}^n \left[\mu_1(X_i) - \mu_0(X_i) + \frac{D_i}{e(X_i)} (Y_i - \mu_1(X_i)) - \frac{1 - D_i}{1 - e(X_i)} (Y_i - \mu_0(X_i)) - \widehat{\tau}_M^{\text{bc}} \right]^2 \\ &- \frac{1}{n} \sum_{i=1}^n \left[\mu_1(X_i) - \mu_0(X_i) + \frac{D_i}{e(X_i)} (Y_i - \mu_1(X_i)) - \frac{1 - D_i}{1 - e(X_i)} (Y_i - \mu_0(X_i)) - \tau \right]^2, \\ T_4 &= \frac{1}{n} \sum_{i=1}^n \left[\mu_1(X_i) - \mu_0(X_i) + \frac{D_i}{e(X_i)} (Y_i - \mu_1(X_i)) - \frac{1 - D_i}{1 - e(X_i)} (Y_i - \mu_0(X_i)) - \tau \right]^2 \\ &- \sigma^2. \end{split}$$

By Assumption 4.3, Assumption 4.1, Theorem B.2, and the fact that $\widehat{\tau}_M^{\rm bc} = O_{\rm P}(1)$, we have $T_1 = o_{\rm P}(1)$. By Assumption 4.1, Theorem B.2, and $\widehat{\tau}_M^{\rm bc} = O_{\rm P}(1)$, we have $T_2 = o_{\rm P}(1)$. By Assumption 4.1 and $\widehat{\tau}_M^{\rm bc} - \tau = o_{\rm P}(1)$, we have $T_3 = o_{\rm P}(1)$. By the fact that $[(X_i, D_i, Y_i)]_{i=1}^n$ are i.i.d., Assumption 4.1 and the weak law of large numbers, we have $T_4 = o_{\rm P}(1)$. Combining the above four facts together then completes the proof. *Q.E.D.*

C.3. Proof of Theorem 4.2

PROOF OF THEOREM 4.2: For Theorem 4.2(i), analysis analogous to the proof of Theorem 4.1(i) can be performed on $\check{\tau}_{M,k}^{\text{bc}}$ for any $k \in [\![K]\!]$. Then the results apply to $\widetilde{\tau}_{M,K}^{\text{bc}}$ automatically since K is fixed.

For Theorem 4.2(ii), from Definition 3.1 in Chernozhukov et al. (2018), $\tilde{\tau}_{M,K}^{bc}$ is the double machine learning estimator. We then follow the proof of Theorem 5.1 (recalling Remark 4.8) and use the notation in Chernozhukov et al. (2018), essentially checking

Assumptions 3.1 and 3.2 therein. In the following, we adopt the notation in Chernozhukov et al. (2018).

For estimating the ATE, from equation (5.3) in Chernozhukov et al. (2018), the score (or the efficient influence function (Tsiatis (2006, Section 3.4))) is

$$\psi(X, D, Y; \widetilde{\tau}, \widetilde{\zeta}) = \widetilde{\mu}_1(X) - \widetilde{\mu}_0(X) + \frac{D(Y - \widetilde{\mu}_1(X))}{\widetilde{e}(X)} - \frac{(1 - D)(Y - \widetilde{\mu}_0(X))}{1 - \widetilde{e}(X)} - \widetilde{\tau},$$

where $\widetilde{\zeta}(x) = (\widetilde{\mu}_0(x), \widetilde{\mu}_1(x), \widetilde{\rho}_0(x), \widetilde{\rho}_1(x))$ are the nuisance parameters by letting $\widetilde{\rho}_0(x) = 1/(1-\widetilde{e}(x))$ and $\widetilde{\rho}_1(x) = 1/\widetilde{e}(x)$. Let $\rho_0(x) = 1/(1-e(x))$ and $\rho_1(x) = 1/e(x)$. Then the true value is $\zeta(x) = (\mu_0(x), \mu_1(x), \rho_0(x), \rho_1(x))$.

We can then write the score as

$$\psi(X, D, Y; \widetilde{\tau}, \widetilde{\zeta}) = \widetilde{\mu}_1(X) - \widetilde{\mu}_0(X) + D(Y - \widetilde{\mu}_1(X))\widetilde{\rho}_1(X) - (1 - D)(Y - \widetilde{\mu}_0(X))\widetilde{\rho}_0(X) - \widetilde{\tau}.$$

For any p > 0, let $||f||_p = ||f(X, D, Y)||_p = (\int |f(\omega)|^p dP_{(X,D,Y)}(\omega))^{1/p}$. For the κ in Assumption 4.1, let $q = 2 + \kappa/2$, $q_1 = 2 + \kappa$, and q_2 such that $q^{-1} = q_1^{-1} + q_2^{-1}$. Let \mathcal{T}_n be the set consisting of all $\tilde{\zeta}$ such that for $\omega \in \{0, 1\}$,

$$\|\widetilde{\mu}_{\omega} - \mu_{\omega}\|_{\infty} = o(n^{-d/(4+2d)}), \qquad \|\widetilde{\rho}_{\omega} - \rho_{\omega}\|_{1} = O(n^{-1/(d+2)}), \qquad \|\widetilde{\rho}_{\omega} - \rho_{\omega}\|_{q_{2}} = o(1).$$

Then the selection of \mathcal{T}_n satisfies Assumption 3.2(a) in Chernozhukov et al. (2018) from Assumption 4.7, Theorem B.4, and Theorem B.2, respectively.

For step 1 in the proof of Theorem 5.1 in Chernozhukov et al. (2018), we verify the Neyman orthogonality. We can show that $\mathrm{E}\psi(X,D,Y;\tau,\zeta)=0$. For any $\widetilde{\zeta}\in\mathcal{T}_n$, the Gateaux derivative in the direction $\widetilde{\zeta}-\zeta$ is

$$\begin{split} & \partial_{\widetilde{\zeta}} \mathrm{E} \psi(X, D, Y; \tau, \zeta) [\widetilde{\zeta} - \zeta] \\ & = \mathrm{E} \big[\widetilde{\mu}_{1}(X) - \mu_{1}(X) \big] - \mathrm{E} \big[\widetilde{\mu}_{0}(X) - \mu_{0}(X) \big] \\ & - \mathrm{E} \big[D \big(\widetilde{\mu}_{1}(X) - \mu_{1}(X) \big) \rho_{1}(X) \big] + \mathrm{E} \big[(1 - D) \big(\widetilde{\mu}_{0}(X) - \mu_{0}(X) \big) \rho_{0}(X) \big] \\ & + \mathrm{E} \big[D \big(Y - \mu_{1}(X) \big) \big(\widetilde{\rho}_{1}(X) - \rho_{1}(X) \big) \big] - \mathrm{E} \big[(1 - D) \big(Y - \mu_{0}(X) \big) \big(\widetilde{\rho}_{0}(X) - \rho_{0}(X) \big) \big]. \end{split}$$

We can check that the above quantity is zero, which completes this step.

Step 2 and step 3 therein can be directly applied.

For step 4 therein, we can establish in the same way that for $\omega \in \{0, 1\}$, $\|\mu_{\omega}\|_{q_1} = O(1)$ from $\|Y\|_{q_1} = O(1)$, and $\tau = O(1)$. Then from Hölder's inequality and $\|\rho_{\omega}\|_{\infty}$ is bounded for $\omega \in \{0, 1\}$, for any $\widetilde{\zeta} \in \mathcal{T}_n$,

$$\begin{split} & \left\| \psi(X, D, Y; \tau, \widetilde{\zeta}) \right\|_{q} \\ & = \left\| \widetilde{\mu}_{1}(X) - \widetilde{\mu}_{0}(X) + (2D - 1) \left(Y - \widetilde{\mu}_{D}(X) \right) \widetilde{\rho}_{D}(X) - \tau \right\|_{q} \\ & \leq \left\| \widetilde{\mu}_{1}(X) \right\|_{q} + \left\| \widetilde{\mu}_{0}(X) \right\|_{q} + \left\| \left(Y - \widetilde{\mu}_{1}(X) \right) \widetilde{\rho}_{1}(X) \right\|_{q} + \left\| \left(Y - \widetilde{\mu}_{0}(X) \right) \widetilde{\rho}_{0}(X) \right\|_{q} + \tau \\ & \leq \left\| \mu_{1} \right\|_{q} + \left\| \widetilde{\mu}_{1} - \mu_{1} \right\|_{\infty} + \left\| \mu_{0} \right\|_{q} + \left\| \widetilde{\mu}_{0} - \mu_{0} \right\|_{\infty} + \left(\left\| Y \right\|_{q_{1}} + \left\| \mu_{1} \right\|_{q_{1}} + \left\| \widetilde{\mu}_{1} - \mu_{1} \right\|_{\infty} \right) \left\| \widetilde{\rho}_{1} \right\|_{q_{2}} \\ & + \left(\left\| Y \right\|_{q_{1}} + \left\| \mu_{0} \right\|_{q_{1}} + \left\| \widetilde{\mu}_{0} - \mu_{0} \right\|_{\infty} \right) \left\| \widetilde{\rho}_{0} \right\|_{q_{2}} + \tau = O(1). \end{split}$$

The last step is from the definition of \mathcal{T}_n . Then we complete this step.

For step 5 therein, by Hölder's inequality, for any $\widetilde{\zeta} \in \mathcal{T}_n$,

$$\begin{split} \left\| \psi(X, D, Y; \tau, \widetilde{\zeta}) - \psi(X, D, Y; \tau, \zeta) \right\|_{2} \\ & \leq \|\widetilde{\mu}_{1} - \mu_{1}\|_{2} + \|\widetilde{\mu}_{0} - \mu_{0}\|_{2} + \left\| D(Y - \widetilde{\mu}_{1}(X))\widetilde{\rho}_{1}(X) - D(Y - \mu_{1}(X))\rho_{1}(X) \right\|_{2} \\ & + \left\| (1 - D)(Y - \widetilde{\mu}_{0}(X))\widetilde{\rho}_{0}(X) - (1 - D)(Y - \mu_{0}(X))\rho_{0}(X) \right\|_{2} \\ & \leq \|\widetilde{\mu}_{1} - \mu_{1}\|_{\infty} + \|\widetilde{\mu}_{0} - \mu_{0}\|_{\infty} + (\|Y\|_{q_{1}} + \|\mu_{1}(X)\|_{q_{1}})\|\widetilde{\rho}_{1} - \rho_{1}\|_{q_{2}} + \|\widetilde{\mu}_{1} - \mu_{1}\|_{\infty}\|\widetilde{\rho}_{1}\|_{2} \\ & + (\|Y\|_{q_{1}} + \|\mu_{0}(X)\|_{q_{1}})\|\widetilde{\rho}_{0} - \rho_{0}\|_{q_{2}} + \|\widetilde{\mu}_{0} - \mu_{0}\|_{\infty}\|\widetilde{\rho}_{0}\|_{2} = o(1). \end{split}$$

The last step is due to the definition of \mathcal{T}_n .

Notice that for any $t \in (0, 1)$,

$$\partial_t^2 \mathrm{E}\psi(X,D,Y;\tau,\zeta+t(\widetilde{\zeta}-\zeta)) = -2\mathrm{E}\big[(2D-1)\big(\widetilde{\mu}_D(X)-\mu_D(X)\big)\big(\widetilde{\rho}_D(X)-\rho_D(X)\big)\big].$$

Then by the definition of \mathcal{T}_n , for any $\widetilde{\zeta} \in \mathcal{T}_n$,

$$\left|\partial_t^2 \mathrm{E} \psi\big(X,D,Y;\tau,\zeta+t(\widetilde{\zeta}-\zeta)\big)\right| \leq 2 \sum_{w \in \{0,1\}} \|\widetilde{\mu}_w - \mu_w\|_\infty \|\widetilde{\rho}_w - \rho_w\|_1 = o\big(n^{-1/2}\big).$$

We then complete this step, and thus complete the proof.

The consistency of the variance estimator can be established in the same way as Theorem 4.1(ii). Q.E.D.

REFERENCES

ABADIE, ALBERTO, AND GUIDO W. IMBENS (2006): "Large Sample Properties of Matching Estimators for Average Treatment Effects," *Econometrica*, 74, 235–267. [2187-2191,2193,2195-2198,2200]

(2008): "On the Failure of the Bootstrap for Matching Estimators," *Econometrica*, 76 (6), 1537–1557. [2188-2190,2193]

(2011): "Bias-Corrected Matching Estimators for Average Treatment Effects," *Journal of Business and Economic Statistics*, 29, 1–11. [2187-2194,2196,2197]

— (2012): "A Martingale Representation for Matching Estimators," *Journal of the American Statistical Association*, 107 (498), 833–843. [2188-2190,2200,2210]

(2016): "Matching on the Estimated Propensity Score," *Econometrica*, 84, 781–807. [2193,2200]

ARMSTRONG, TIMOTHY B., AND MICHAL KOLESÁR (2021): "Finite-Sample Optimal Estimation and Inference on Average Treatment Effects Under Unconfoundedness," *Econometrica*, 89, 1141–1177. [2193]

ATHEY, SUSAN, GUIDO W. IMBENS, JONAS METZGER, AND EVAN MUNRO (2023): "Using Wasserstein Generative Adversarial Networks for the Design of Monte Carlo Simulations," *Journal of Econometrics* (forthcoming). [2197]

BANG, HEEJUNG, AND JAMES M. ROBINS (2005): "Doubly Robust Estimation in Missing Data and Causal Inference Models," *Biometrics*, 61 (4), 962–973. [2188,2192]

BENTLEY, JON L. (1975): "Multidimensional Binary Search Trees Used for Associative Searching," *Communications of the ACM*, 18 (9), 509–517. [2201]

BERRETT, THOMAS B., RICHARD J. SAMWORTH, AND MING YUAN (2019): "Efficient Multivariate Entropy Estimation via k-Nearest Neighbour Distances," *The Annals of Statistics*, 47 (1), 288–318. [2188,2190]

BHATTACHARYA, BHASWAR B. (2019): "A General Asymptotic Framework for Distribution-Free Graph-Based Two-Sample Tests," *Journal of the Royal Statistical Society. Series B*, 81 (3), 575–602. [2188]

BIAU, GÉRARD, AND LUC DEVROYE (2015): Lectures on the Nearest Neighbor Method. Springer. [2203]

BORGEAUD, SEBASTIAN, ARTHUR MENSCH, JORDAN HOFFMANN, TREVOR CAI, ELIZA RUTHERFORD, KATIE MILLICAN, GEORGE VAN DEN DRIESSCHE, JEAN-BAPTISTE LESPIAU, BOGDAN DAMOC, AIDAN CLARK et al. (2021): "Improving Language Models by Retrieving From Trillions of Tokens," in *Proceedings of the 39th International Conference on Machine Learning*, Vol. 162. Proceedings of Machine Learning Research, 2206–2240. [2191,2201]

- BROOKHART, M. ALAN, SEBASTIAN SCHNEEWEISS, KENNETH J. ROTHMAN, ROBERT J. GLYNN, JERRY AVORN, AND TIL STÜRMER (2006): "Variable Selection for Propensity Score Models," *American Journal of Epidemiology*, 163 (12), 1149–1156. [2187]
- CHAPIN, F. STUART (1947): Experimental Designs in Sociological Research. Harper and Brothers. [2187]
- CHEN, XIAOHONG, AND TIMOTHY M. CHRISTENSEN (2015): "Optimal Uniform Convergence Rates and Asymptotic Normality for Series Estimators Under Weak Dependence and Weak Conditions," *Journal of Econometrics*, 188 (2), 447–465. [2196]
- CHERNOZHUKOV, VICTOR, DENIS CHETVERIKOV, MERT DEMIRER, ESTHER DUFLO, CHRISTIAN HANSEN, WHITNEY NEWEY, AND JAMES ROBINS (2018): "Double/Debiased Machine Learning for Treatment and Structural Parameters," *The Econometrics Journal*, 21 (1), C1–C68. [2188,2194-2197,2212,2213]
- COCHRAN, WILLIAM G., AND DONALD B. RUBIN (1973): "Controlling Bias in Observational Studies: A Review," Sankhyā, Series A, 35 (4), 417–446. [2187]
- COVER, THOMAS M., AND JOY THOMAS (2006): *Elements of Information Theory* (Second Ed.). John Wiley and Sons. [2189]
- CUNNINGHAM, SCOTT (2021): Causal Inference: The Mixtape. Yale University Press. [2189]
- DEHEJIA, RAJEEV H., AND SADEK WAHBA (1999): "Causal Effects in Nonexperimental Studies: Reevaluating the Evaluation of Training Programs," *Journal of the American Statistical Association*, 94 (448), 1053–1062. [2197]
- DEVROYE, LUC, LÁSZLÓ GYÖRFI, GÁBOR LUGOSI, AND HARRO WALK (2017): "On the Measure of Voronoi Cells," *Journal of Applied Probability*, 54 (2), 394–408. [2203]
- EVANS, LAWRENCE C., AND RONALD F. GARZEPY (2018): Measure Theory and Fine Properties of Functions. Routledge. [2206]
- FARRELL, MAX H. (2015): "Robust Inference on Average Treatment Effects With Possibly More Covariates Than Observations," *Journal of Econometrics*, 189 (1), 1–23. [2188]
- FERMAN, BRUNO (2021): "Matching Estimators With few Treated and Many Control Observations," *Journal of Econometrics*, 225, 295–307. [2193]
- FRIEDMAN, JEROME H., AND LAWRENCE C. RAFSKY (1979): "Multivariate Generalizations of the Wald-Wolfowitz and Smirnov Two-Sample Tests," *The Annals of Statistics*, 7 (4), 697–717. [2188]
- FRIEDMAN, JEROME H., JON L. BENTLEY, AND RAPHAEL A. FINKEL (1977): "An Algorithm for Finding Best Matches in Logarithmic Expected Time," *ACM Transactions on Mathematical Software*, 3 (3), 209–226. [2201] GREENWOOD, ERNEST (1945): *Experimental Sociology*. Columbia University Press. [2187]
- HAHN, JINYONG (1998): "On the Role of the Propensity Score in Efficient Semiparametric Estimation of Average Treatment Effects," *Econometrica*, 66 (2), 315–331. [2188,2194,2195]
- HAN, YANJUN, JIANTAO JIAO, TSACHY WEISSMAN, AND YIHONG WU (2020): "Optimal Rates of Entropy Estimation Over Lipschitz Balls," *The Annals of Statistics*, 48 (6), 3228–3250. [2204]
- HANSEN, BEN B. (2008): "The Prognostic Analogue of the Propensity Score," *Biometrika*, 95 (2), 481–488. [2200]
- HENZE, NORBERT (1988): "A Multivariate Two-Sample Test Based on the Number of Nearest Neighbor Type Coincidences," *The Annals of Statistics*, 16 (2), 772–783. [2188]
- HENZE, NORBERT, AND MATHEW D. PENROSE (1999): "On the Multivariate Runs Test," *The Annals of Statistics*, 27 (1), 290–298. [2188]
- Ho, Daniel E., Kosuke Imai, Gary King, and Elizabeth A. Stuart (2007): "Matching as Nonparametric Preprocessing for Reducing Model Dependence in Parametric Causal Inference," *Political Analysis*, 15 (3), 199–236. [2187]
- IMBENS, GUIDO W. (2004): "Nonparametric Estimation of Average Treatment Effects Under Exogeneity: A Review," *Review of Economics and Statistics*, 86 (1), 4–29. [2187]
- IMBENS, GUIDO W., AND DONALD B. RUBIN (2015): Causal Inference in Statistics, Social, and Biomedical Sciences. Cambridge University Press. [2187,2189]
- KALLUS, NATHAN (2020): "Generalized Optimal Matching Methods for Causal Inference," *Journal of Machine Learning Research*, 21, 1–54. [2193]
- KPOTUFE, SAMORY (2017): "Lipschitz Density-Ratios, Structured Data, and Data-Driven Tuning," in 2017 International Conference on Artificial Intelligence and Statistics. PMLR, 1320–1328. [2204,2205]
- KREMER, JAN, FABIAN GIESEKE, K. STEENSTRUP PEDERSEN, AND CHRISTIAN IGEL (2015): "Nearest Neighbor Density Ratio Estimation for Large-Scale Applications in Astronomy," *Astronomy and Computing*, 12, 67–72. [2189,2191,2201]
- LALONDE, ROBERT J. (1986): "Evaluating the Econometric Evaluations of Training Programs With Experimental Data," *The American Economic Review*, 76 (4), 604–620. [2197,2198]
- LIMA, MARCOS, CARLOS E. CUNHA, HIROAKI OYAIZU, JOSHUA FRIEMAN, HUAN LIN, AND ERIN S. SHELDON (2008): "Estimating the Redshift Distribution of Photometric Galaxy Samples," *Monthly Notices of the Royal Astronomical Society*, 390 (1), 118–130. [2189,2191,2201]

- LIN, ZHEXIAO, AND FANG HAN (2023): "On Boosting the Power of Chatterjee's Rank Correlation," *Biometrika*, 110 (2), 283–299. [2188,2189]
- LIN, ZHEXIAO, PENG DING, AND FANG HAN (2023): "Supplement to 'Estimation Based on Nearest Neighbor Matching: From Density Ratio to Average Treatment Effect'," *Econometrica Supplemental Material*, 91, https://doi.org/10.3982/ECTA20598. [2189]
- LIU, REGINA Y., AND KESAR SINGH (1993): "A Quality Index Based on Data Depth and Multivariate Rank Tests," *Journal of the American Statistical Association*, 88 (421), 252–260. [2188]
- MORGAN, STEPHEN L., AND DAVID J. HARDING (2006): "Matching Estimators of Causal Effects: Prospects and Pitfalls in Theory and Practice," *Sociological Methods and Research*, 35 (1), 3–60. [2187]
- NEWEY, WHITNEY K. (1997): "Convergence Rates and Asymptotic Normality for Series Estimators," *Journal of Econometrics*, 79 (1), 147–168. [2196]
- NGUYEN, XUANLONG, MARTIN J. WAINWRIGHT, AND MICHAEL I. JORDAN (2010): "Estimating Divergence Functionals and the Likelihood Ratio by Convex Risk Minimization," *IEEE Transactions on Information Theory*, 56 (11), 5847–5861. [2188,2204]
- NOSHAD, MORTEZA, KEVIN R. MOON, SALIMEH Y. SEKEH, AND ALFRED O. HERO (2017): "Direct Estimation of Information Divergence Using Nearest Neighbor Ratios," in 2017 IEEE International Symposium on Information Theory (ISIT), 903–907. [2190,2201,2202,2207]
- OTSU, TAISUKE, AND YOSHIYASU RAI (2017): "Bootstrap Inference of Matching Estimators for Average Treatment Effects," *Journal of the American Statistical Association*, 112 (520), 1720–1732. [2192]
- Póczos, Barnabás, and Jeff Schneider (2011): "On the Estimation of Alpha-Divergences," in 2011 International Conference on Artificial Intelligence and Statistics, 609–617. [2189]
- ROSENBAUM, PAUL R. (2010): Design of Observational Studies. Springer. [2187]
- ROSENBAUM, PAUL R., AND DONALD B. RUBIN (1983): "The Central Role of the Propensity Score in Observational Studies for Causal Effects," *Biometrika*, 70 (1), 41–55. [2192,2193,2200]
- RUBIN, DONALD B. (1973): "Matching to Remove Bias in Observational Studies," *Biometrics*, 29 (1), 159–183. [2187]
- ——— (1974): "Estimating Causal Effects of Treatments in Randomized and Nonrandomized Studies," Journal of Educational Psychology, 66 (5), 688–701. [2191]
- (2006): Matched Sampling for Causal Effects. Cambridge University Press. [2187]
- RUBIN, DONALD B., AND NEAL THOMAS (2000): "Combining Propensity Score Matching With Additional Adjustments for Prognostic Covariates," *Journal of the American Statistical Association*, 95 (450), 573–585. [2187]
- SCHARFSTEIN, DANIEL O., ANDREA ROTNITZKY, AND JAMES M. ROBINS (1999): "Adjusting for Nonignorable Drop-out Using Semiparametric Nonresponse Models," *Journal of the American Statistical Association*, 94 (448), 1096–1120. [2188,2192]
- SEKHON, JASJEET S. (2008): "Multivariate and Propensity Score Matching Software With Automated Balance Optimization: The Matching Package for R," *Journal of Statistical Software*, 42 (7), 1–52. [2187]
- SHADISH, WILLIAM R., MARGARET H. CLARK, AND PETER M. STEINER (2008): "Can Nonrandomized Experiments Yield Accurate Answers? A Randomized Experiment Comparing Random and Nonrandom Assignments," Journal of the American Statistical Association, 103 (484), 1334–1344. [2197,2199]
- SHI, HONGJIAN, MATHIAS DRTON, AND FANG HAN (2022): "On the Power of Chatterjee's Rank Correlation," *Biometrika*, 109 (2), 317–333. [2188]
- ——— (2023): "On Azadkia-Chatterjee's Conditional Dependence Coefficient," *Bernoulli* (forthcoming). [2188]
- SMITH, HERBERT L. (1997): "Matching With Multiple Controls to Estimate Treatment Effects in Observational Studies," *Sociological Methodology*, 27 (1), 325–353. [2187]
- STEIN, ELIAS M. (2016): Singular Integrals and Differentiability Properties of Functions. Princeton University Press. [2204]
- SUGIYAMA, MASASHI, TAIJI SUZUKI, AND TAKAFUMI KANAMORI (2012): Density Ratio Estimation in Machine Learning. Cambridge University Press. [2188,2189]
- SUGIYAMA, MASASHI, TAIJI SUZUKI, SHINICHI NAKAJIMA, HISASHI KASHIMA, PAUL VON BÜNAU, AND MOTOAKI KAWANABE (2008): "Direct Importance Estimation for Covariate Shift Adaptation," *Annals of the Institute of Statistical Mathematics*, 60 (4), 699–746. [2204]
- TSIATIS, ANASTASIOS A. (2006): Semiparametric Theory and Missing Data. Springer. [2213]
- TSYBAKOV, ALEXANDRE B. (2009): Introduction to Nonparametric Estimation. Springer. [2205]
- VORONOI, GEORGES (1908): "Nouvelles Applications des Paramètres Continus à la Théorie des Formes Quadratiques. Deuxième Mémoire. Recherches sur les Parallélloèdres Primitifs," *Journal für die reine und angewandte Mathematik (Crelles Journal)*, 1908 (134), 198–287. [2190]

WALD, ABRAHAM, AND JACOB WOLFOWITZ (1940): "On a Test Whether Two Samples Are From the Same Population," *Annals of Mathematical Statistics*, 11 (2), 147–162. [2188]

WANG, YIXIN, AND JOSÉ R. ZUBIZARRETA (2023): "Large Sample Properties of Matching for Balance," Statistica Sinica, 33, 1789–1808. [2200]

YANG, SHU, AND YUNSHU ZHANG (2023): "Multiply Robust Matching Estimators of Average and Quantile Treatment Effects," Scandinavian Journal of Statistics, 50, 235–265. [2200]

ZHAO, PUNING, AND LIFENG LAI (2020): "Minimax Optimal Estimation of KL Divergence for Continuous Distributions," *IEEE Transactions on Information Theory*, 66 (12), 7787–7811. [2190,2206]

——— (2022): "Analysis of KNN Density Estimation," *IEEE Transactions on Information Theory*, 68 (12), 7971–7995. [2206]

Co-editor Guido Imbens handled this manuscript.

Manuscript received 22 February, 2022; final version accepted 7 September, 2023; available online 7 September, 2023.

The replication package for this paper is available at https://doi.org/10.5281/zenodo.8322609. The authors were granted an exemption to publish parts of their data because either access to these data is restricted or the authors do not have the right to republish them. Therefore, the replication package only includes the codes and the parts of the data that are not subject to the exemption. However, the authors provided the Journal with (or assisted the Journal to obtain) temporary access to the restricted data. The Journal checked the provided and restricted data and the codes for their ability to reproduce the results in the paper and approved online appendices.

SUPPLEMENT TO "ESTIMATION BASED ON NEAREST NEIGHBOR MATCHING: FROM DENSITY RATIO TO AVERAGE TREATMENT EFFECT"

(Econometrica, Vol. 91, No. 6, November 2023, 2187–2217)

ZHEXIAO LIN

Department of Statistics, University of California, Berkeley

PENG DING

Department of Statistics, University of California, Berkeley

FANG HAN

Department of Statistics, University of Washington

S1. PROOFS OF THE RESULTS IN SECTIONS 3 AND 4

Additional Notation. WE USE X and Z to represent $(X_1, X_2, ..., X_{N_0})$ and $(Z_1, Z_2, ..., Z_{N_0})$ Z_{N_1}), respectively. Let U(0,1) denote the uniform distribution on [0, 1]. Let $U \sim U(0,1)$ and $U_{(M)}$ be the Mth order statistic of N_0 independent random variables from U(0,1), assumed to be mutually independent and both independent of (X, Z). It is well known that $U_{(M)} \sim \operatorname{Beta}(M, N_0 + 1 - M)$. Let $\operatorname{Bin}(\cdot, \cdot)$ denote the binomial distribution. Let $L_1(\mathbb{R}^d)$ denote the space of all functions $f : \mathbb{R}^d \to \mathbb{R}$ such that $\int |f(x)| dx < \infty$. For any $x \in \mathbb{R}^d$ and function $f: \mathbb{R}^d \to \mathbb{R}$, we say x is a Lebesgue point (Bogachev and Ruas (2007, Theorem 5.6.2)) of f if

$$\lim_{\delta \to 0^+} \frac{1}{\lambda(B_{x,\delta})} \int_{B_{x,\delta}} |f(x) - f(z)| \, \mathrm{d}z = 0.$$

S2. PROOFS OF THE RESULTS IN APPENDIX A

S2.1. Proof of Theorem A.1

PROOF OF THEOREM A.1: We consider the complexities of two algorithms separately. Algorithm 1.

The worst-case computation complexity of building a balanced k-d tree is $O(dN_0 \log N_0)$ (cf. Brown (2015)) since the size of the k-d tree is N_0 .

The average computation complexity of searching a NN is $O(\log N_0)$ from Friedman, Bentley, and Finkel (1977), and then the average computation complexity of search M-NNs in $\{X_i\}_{i=1}^{N_0}$ for all $\{Z_j\}_{j=1}^{N_1}$ is $O(MN_1 \log N_0)$.

Notice that $|S_j| = M$ for any $j \in [N_1]$ and then $|\bigcup_{i=1}^{N_1} S_j| \le N_1 M$. Since the elements of each S_j are in $[N_0]$, the largest integer in $\bigcup_{j=1}^{N_1} S_j$ is N_0 . Then the computation complexity of counting step is $O(N_1M + N_0)$ due to the counting sort algorithm (Cormen, Leiserson, Rivest, and Stein (2009, Section 8.2)).

Combining the above three steps completes the proof for Algorithm 1.

Zhexiao Lin: zhexiaolin@berkelev.edu Peng Ding: pengdingpku@berkeley.edu

Fang Han: fanghan@uw.edu

Algorithm 2.

The computation complexity of building a k-d tree is $O(d(N_0 + n) \log(N_0 + n))$ from Algorithm 1 since the size of the k-d tree is $N_0 + n$.

For the searching step, for each $j \in [N_1]$, the number of NNs to be searched is $M + \sum_{i=1}^n \mathbb{1}(\|x_i - Z_j\| \le \|\mathcal{X}_{(M)}(Z_j) - Z_j\|)$. Then from (2.2), the total number of NNs searched for all $j \in [N_1]$ is $\sum_{j=1}^{N_1} (M + \sum_{i=1}^n \mathbb{1}(\|x_i - Z_j\| \le \|\mathcal{X}_{(M)}(Z_j) - Z_j\|)) = N_1 M + \sum_{i=1}^n K_M(x_i)$. Let X, Z be two independent copies from ν_0 , ν_1 , respectively, and are independent of the data. Since $[Z_j]_{j=1}^{N_1}$ are i.i.d. and $[X_i]_{i=1}^{N_0} \cup [x_i]_{i=1}^n$ are i.i.d, we have $E[\sum_{i=1}^n K_M(x_i)] = nE[K_M(X)] = N_1 nE[\nu_1(A_M(X))] = N_1 n \frac{M}{N_0+1}$ since $E[\nu_1(A_M(X))] = P(\|X - Z\| \le \|\mathcal{X}_{(M)}(Z) - Z\|) = P(U \le U_{(M)}) = \frac{M}{N_0+1}$ by using the probability integral transform. Then the average computation complexity for the searching step is $O(N_0^{-1}N_1M(N_0+n)\log(N_0+n))$.

For the counting step, the computation complexity for counting $\bigcup_{j=1}^{N_1} S_j$ is $O(N_0 + N_1 M)$ since the cardinality of $\bigcup_{j=1}^{N_1} S_j$ is at most $N_1 M$ and the largest integer is N_0 . The average computation complexity for counting $\bigcup_{j=1}^{N_1} S_j'$ is $O(N_0^{-1} N_1 M n + n)$ since the average cardinality of $\bigcup_{j=1}^{N_1} S_j'$ is at most $N_0^{-1} N_1 M n$ and the largest integer is n.

Combining the above three steps completes the proof for Algorithm 2. Q.E.D.

S3. PROOFS OF THE RESULTS IN APPENDIX B

S3.1. Proof of Lemma B.1

PROOF OF LEMMA B.1: From the Lebesgue differentiation theorem, for any $f \in L_1(\mathbb{R}^d)$, x is a Lebesgue point of f for λ -almost all x. Then for ν_0 -almost all x, we have $f_0(x) > 0$ and x is a Lebesgue point of f_0 and f_1 from the absolute continuity of ν_0 and ν_1 . We then only need to consider those $x \in \mathbb{R}^d$ such that $f_0(x) > 0$ and x is a Lebesgue point of f_0 and f_1 .

We first introduce a lemma about the Lebesgue point.

LEMMA S3.1: Let ν be a probability measure on \mathbb{R}^d admitting a density f with respect to the Lebesgue measure. Let $x \in \mathbb{R}^d$ be a Lebesgue point of f. Then for any $\epsilon \in (0, 1)$, there exists $\delta = \delta_x > 0$ such that for any $z \in \mathbb{R}^d$ satisfying $||z - x|| \le \delta$, we have

$$\left|\frac{\nu(B_{x,\|z-x\|})}{\lambda(B_{x,\|z-x\|})} - f(x)\right| \le \epsilon, \qquad \left|\frac{\nu(B_{z,\|z-x\|})}{\lambda(B_{z,\|z-x\|})} - f(x)\right| \le \epsilon.$$

Part I. This part proves the first claim. We separate the proof of Part I into two cases based on the value of $f_1(x)$.

Case I.1. $f_1(x) > 0$. Since x is a Lebesgue point of ν_0 and ν_1 , by Lemma S3.1, for any $\epsilon \in (0, 1)$, there exists some $\delta = \delta_x > 0$ such that for any $z \in \mathbb{R}^d$ with $||z - x|| \le \delta$, we have for $w \in \{0, 1\}$,

$$\left|\frac{\nu_w(B_{x,\|z-x\|})}{\lambda(B_{x,\|z-x\|})} - f_w(x)\right| \le \epsilon f_w(x), \qquad \left|\frac{\nu_w(B_{z,\|z-x\|})}{\lambda(B_{z,\|z-x\|})} - f_w(x)\right| \le \epsilon f_w(x).$$

Accordingly, if $||z - x|| \le \delta$, by $\lambda(B_{z,||x-z||}) = \lambda(B_{x,||x-z||})$, we have

$$\frac{1 - \epsilon}{1 + \epsilon} \frac{f_0(x)}{f_1(x)} \le \frac{\nu_0(B_{z, \|x - z\|})}{\lambda(B_{z, \|x - z\|})} \frac{\lambda(B_{x, \|x - z\|})}{\nu_1(B_{x, \|x - z\|})} = \frac{\nu_0(B_{z, \|x - z\|})}{\nu_1(B_{x, \|x - z\|})} \le \frac{1 + \epsilon}{1 - \epsilon} \frac{f_0(x)}{f_1(x)}.$$
 (S3.1)

On the other hand, for any $z \in \mathbb{R}^d$ such that $||z - x|| > \delta$, $\nu_0(B_{z,||z-x||}) \ge \nu_0(B_{z^*,\delta}) \ge (1-\epsilon)f_0(x)\lambda(B_{z^*,\delta}) = (1-\epsilon)f_0(x)\lambda(B_{0,\delta})$, where z^* is the intersection point of the surface of $B_{x,\delta}$ and the line connecting z and x.

Let $\eta_N = 4\log(N_0/M)$. Since $M\log N_0/N_0 \to 0$, we can take N_0 large enough so that $\eta_N \frac{M}{N_0} = 4\frac{M}{N_0}\log(\frac{N_0}{M}) < (1-\epsilon)f_0(x)\lambda(B_{0,\delta})$. Then for any $z \in \mathbb{R}^d$ such that $\nu_0(B_{z,\|z-x\|}) \le \eta_N M/N_0$, we have $\|z-x\| \le \delta$ since otherwise it would contradict the selection of N_0 .

Let Z be a copy from ν_1 independent of the data. Then

$$E[\nu_1(A_M(x))] = P(Z \in A_M(x)) = P(\nu_0(B_{Z,\|x-Z\|}) \le \nu_0(B_{Z,\|X_{(M)}(Z)-Z\|})).$$
 (S3.2)

For any given $z \in \mathbb{R}^d$, $\left[\nu_0(B_{z,\|X_i-z\|})\right]_{i=1}^{N_0}$ are i.i.d. from U(0,1) since $[X_i]_{i=1}^{N_0}$ are i.i.d. from ν_0 and we use the probability integral transform. Then $\nu_0(B_{Z,\|\mathcal{X}_{(M)}(Z)-Z\|})$ has the same distribution as $U_{(M)}$ and is independent of Z.

Upper bound. With a slight abuse of notation, we define $W = \nu_0(B_{Z,||x-Z||})$. We then have, from (S3.1) and (S3.2),

$$\begin{split} & E\left[\nu_{1}\left(A_{M}(x)\right)\right] \\ & = P\left(W \leq \nu_{0}(B_{Z,\|\mathcal{X}_{(M)}(Z)-Z\|})\right) \\ & \leq P\left(W \leq \nu_{0}(B_{Z,\|\mathcal{X}_{(M)}(Z)-Z\|}) \leq \eta_{N}\frac{M}{N_{0}}\right) + P\left(\nu_{0}(B_{Z,\|\mathcal{X}_{(M)}(Z)-Z\|}) > \eta_{N}\frac{M}{N_{0}}\right) \\ & = P\left(W \leq \nu_{0}(B_{Z,\|\mathcal{X}_{(M)}(Z)-Z\|}) \leq \eta_{N}\frac{M}{N_{0}}, \|Z-x\| \leq \delta\right) + P\left(U_{(M)} > \eta_{N}\frac{M}{N_{0}}\right) \\ & \leq P\left(\nu_{0}(B_{Z,\|\mathcal{X}-Z\|}) \leq \nu_{0}(B_{Z,\|\mathcal{X}_{(M)}(Z)-Z\|}), \|Z-x\| \leq \delta\right) + P\left(U_{(M)} > \eta_{N}\frac{M}{N_{0}}\right) \\ & \leq P\left(\frac{1-\epsilon}{1+\epsilon}\frac{f_{0}(x)}{f_{1}(x)}\nu_{1}(B_{x,\|\mathcal{X}-Z\|}) \leq \nu_{0}(B_{Z,\|\mathcal{X}_{(M)}(Z)-Z\|}), \|Z-x\| \leq \delta\right) + P\left(U_{(M)} > \eta_{N}\frac{M}{N_{0}}\right) \\ & \leq P\left(\frac{1-\epsilon}{1+\epsilon}\frac{f_{0}(x)}{f_{1}(x)}\nu_{1}(B_{x,\|\mathcal{X}-Z\|}) \leq \nu_{0}(B_{Z,\|\mathcal{X}_{(M)}(Z)-Z\|})\right) + P\left(U_{(M)} > \eta_{N}\frac{M}{N_{0}}\right) \\ & = P\left(\frac{1-\epsilon}{1+\epsilon}\frac{f_{0}(x)}{f_{1}(x)}U \leq U_{(M)}\right) + P\left(U_{(M)} > \eta_{N}\frac{M}{N_{0}}\right). \end{split} \tag{S3.3}$$

For the second term in (S3.3), notice that $\eta_N \to \infty$ as $N_0 \to \infty$. Then from the Chernoff bound and for N_0 sufficiently large, we have

$$\begin{split} \frac{N_0}{M} \mathbf{P} \bigg(U_{(M)} > \eta_N \frac{M}{N_0} \bigg) &= \frac{N_0}{M} \mathbf{P} \bigg(\mathbf{Bin} \bigg(N_0, \eta_N \frac{M}{N_0} \bigg) < M \bigg) \\ &\leq \frac{N_0}{M} \exp \bigg((1 + \log \eta_N - \eta_N) M \bigg) \\ &\leq \frac{N_0}{M} \exp \bigg(-\frac{1}{2} \eta_N M \bigg) = \bigg(\frac{N_0}{M} \bigg)^{1-2M}. \end{split}$$

Since $M/N_0 \to 0$ and $M \ge 1$, we then obtain

$$\lim_{N_0 \to \infty} \frac{N_0}{M} P\left(U_{(M)} > \eta_N \frac{M}{N_0}\right) = 0.$$
 (S3.4)

For the first term in (S3.3), we have

$$\frac{N_0}{M} P\left(\frac{1-\epsilon}{1+\epsilon} \frac{f_0(x)}{f_1(x)} U \le U_{(M)}\right)$$

$$= \frac{N_0}{M} \int_0^1 P\left(U_{(M)} \ge \frac{1-\epsilon}{1+\epsilon} \frac{f_0(x)}{f_1(x)} t\right) dt$$

$$= \frac{1+\epsilon}{1-\epsilon} \frac{f_1(x)}{f_0(x)} \int_0^{\frac{1-\epsilon}{1+\epsilon} \frac{f_0(x)}{f_1(x)} \frac{N_0}{M}} P\left(U_{(M)} \ge \frac{M}{N_0} t\right) dt \le \frac{1+\epsilon}{1-\epsilon} \frac{f_1(x)}{f_0(x)} \int_0^\infty P\left(\frac{N_0}{M} U_{(M)} \ge t\right) dt$$

$$= \frac{1+\epsilon}{1-\epsilon} \frac{f_1(x)}{f_0(x)} \frac{N_0}{M} E[U_{(M)}] = \frac{1+\epsilon}{1-\epsilon} \frac{f_1(x)}{f_0(x)} \frac{N_0}{N_0+1}. \tag{S3.5}$$

We then obtain

$$\limsup_{N_0 \to \infty} \frac{N_0}{M} P\left(\frac{1 - \epsilon}{1 + \epsilon} \frac{f_0(x)}{f_1(x)} U \le U_{(M)}\right) \le \frac{1 + \epsilon}{1 - \epsilon} \frac{f_1(x)}{f_0(x)}. \tag{S3.6}$$

Plugging (S3.4) and (S3.6) to (S3.3) then yields

$$\limsup_{N_0 \to \infty} \frac{N_0}{M} \mathbb{E}\left[\nu_1(A_M(x))\right] \le \frac{1+\epsilon}{1-\epsilon} \frac{f_1(x)}{f_0(x)}.$$
 (S3.7)

Lower bound. We have, from (S3.1) and (S3.2),

$$E[\nu_{1}(A_{M}(x))] = P(W \leq \nu_{0}(B_{Z,\|\mathcal{X}_{(M)}(Z)-Z\|})) \geq P(W \leq \nu_{0}(B_{Z,\|\mathcal{X}_{(M)}(Z)-Z\|}) \leq \eta_{N} \frac{M}{N_{0}})$$

$$= P(W \leq \nu_{0}(B_{Z,\|\mathcal{X}_{(M)}(Z)-Z\|}) \leq \eta_{N} \frac{M}{N_{0}}, \|Z - x\| \leq \delta)$$

$$\geq P(\frac{1+\epsilon}{1-\epsilon} \frac{f_{0}(x)}{f_{1}(x)} \nu_{1}(B_{x,\|x-Z\|}) \leq \nu_{0}(B_{Z,\|\mathcal{X}_{(M)}(Z)-Z\|}) \leq \eta_{N} \frac{M}{N_{0}}, \|Z - x\| \leq \delta)$$

$$= P(\frac{1+\epsilon}{1-\epsilon} \frac{f_{0}(x)}{f_{1}(x)} \nu_{1}(B_{x,\|x-Z\|}) \leq \nu_{0}(B_{Z,\|\mathcal{X}_{(M)}(Z)-Z\|}) \leq \eta_{N} \frac{M}{N_{0}})$$

$$\geq P(\frac{1+\epsilon}{1-\epsilon} \frac{f_{0}(x)}{f_{1}(x)} \nu_{1}(B_{x,\|x-Z\|}) \leq \nu_{0}(B_{Z,\|\mathcal{X}_{(M)}(Z)-Z\|})$$

$$- P(\nu_{0}(B_{Z,\|\mathcal{X}_{(M)}(Z)-Z\|}) > \eta_{N} \frac{M}{N_{0}})$$

$$= P(\frac{1+\epsilon}{1-\epsilon} \frac{f_{0}(x)}{f_{1}(x)} U \leq U_{(M)}) - P(U_{(M)} > \eta_{N} \frac{M}{N_{0}}). \tag{S3.8}$$

The second last equality is from the fact that for $||Z - x|| > \delta$,

$$\frac{1+\epsilon}{1-\epsilon}\frac{f_0(x)}{f_1(x)}\nu_1(B_{x,\|x-Z\|}) \geq \frac{1+\epsilon}{1-\epsilon}\frac{f_0(x)}{f_1(x)}\nu_1(B_{x,\delta}) \geq \frac{1+\epsilon}{1-\epsilon}\frac{f_0(x)}{f_1(x)}f_1(x)(1-\epsilon)\lambda(B_{0,\delta}) > \eta_N\frac{M}{N_0},$$

and then that $\frac{1+\epsilon}{1-\epsilon} \frac{f_0(x)}{f_1(x)} \nu_1(B_{x,\|x-Z\|}) \le \eta_N \frac{M}{N_0}$ implies $\|Z-x\| \le \delta$.

For the first term in (S3.8), we have

$$\frac{N_0}{M} P\left(\frac{1+\epsilon}{1-\epsilon} \frac{f_0(x)}{f_1(x)} U \le U_{(M)}\right) = \frac{1-\epsilon}{1+\epsilon} \frac{f_1(x)}{f_0(x)} \int_0^{\frac{1+\epsilon}{1-\epsilon} \frac{f_0(x)}{f_1(x)} \frac{N_0}{M}} P\left(U_{(M)} \ge \frac{M}{N_0} t\right) dt.$$

If $\frac{1+\epsilon}{1-\epsilon} \frac{f_0(x)}{f_1(x)} \ge 1$, then by $U_{(M)} \in [0, 1]$, we have

$$\frac{N_0}{M} P\left(\frac{1+\epsilon}{1-\epsilon} \frac{f_0(x)}{f_1(x)} U \le U_{(M)}\right) = \frac{1-\epsilon}{1+\epsilon} \frac{f_1(x)}{f_0(x)} \frac{N_0}{M} E[U_{(M)}] = \frac{1-\epsilon}{1+\epsilon} \frac{f_1(x)}{f_0(x)} \frac{N_0}{N_0+1}.$$

If $\frac{1+\epsilon}{1-\epsilon} \frac{f_0(x)}{f_1(x)} < 1$, from the Chernoff bound,

$$\begin{split} &\int_{\frac{1+\epsilon}{1-\epsilon}\frac{f_0(x)}{f_1(x)}\frac{N_0}{M}}^{\frac{N_0}{M}} \mathbf{P}\bigg(U_{(M)} \geq \frac{M}{N_0}t\bigg) \, \mathrm{d}t \\ &\leq \bigg[1 - \frac{1+\epsilon}{1-\epsilon}\frac{f_0(x)}{f_1(x)}\bigg] \frac{N_0}{M} \mathbf{P}\bigg(U_{(M)} \geq \frac{1+\epsilon}{1-\epsilon}\frac{f_0(x)}{f_1(x)}\bigg) \\ &\leq \bigg[1 - \frac{1+\epsilon}{1-\epsilon}\frac{f_0(x)}{f_1(x)}\bigg] \frac{N_0}{M} \exp\bigg[M - \frac{1+\epsilon}{1-\epsilon}\frac{f_0(x)}{f_1(x)}N_0 \\ &- M\log M + M\log\bigg(\frac{1+\epsilon}{1-\epsilon}\frac{f_0(x)}{f_1(x)}N_0\bigg)\bigg]. \end{split}$$

Since $f_0(x) > 0$ and $M \log N_0/N_0 \to 0$, we obtain

$$\lim_{N_0\to\infty}\int_{\frac{1+\epsilon}{1-\epsilon}}^{\frac{N_0}{M}} \sum_{t=0}^{N_0} \Pr\left(U_{(M)}\geq \frac{M}{N_0}t\right) \mathrm{d}t = 0.$$

Then we always have

$$\lim_{N_0 \to \infty} \frac{N_0}{M} P\left(\frac{1+\epsilon}{1-\epsilon} \frac{f_0(x)}{f_1(x)} U \le U_{(M)}\right) = \frac{1-\epsilon}{1+\epsilon} \frac{f_1(x)}{f_0(x)}.$$

Using the above identity along with (S3.4) to (S3.8) yields

$$\liminf_{N_0 \to \infty} \frac{N_0}{M} \mathbb{E}[\nu_1(A_M(x))] \ge \frac{1 - \epsilon}{1 + \epsilon} \frac{f_1(x)}{f_0(x)}.$$
 (S3.9)

Lastly, combining (S3.7) with (S3.9) and noticing that ϵ is arbitrary, we obtain

$$\lim_{N_0 \to \infty} \frac{N_0}{M} E[\nu_1(A_M(x))] = \frac{f_1(x)}{f_0(x)} = r(x).$$
 (S3.10)

Case I.2. $f_1(x) = 0$. Again, for any $\epsilon \in (0, 1)$, by Lemma S3.1, there exists some $\delta = \delta_x > 0$ such that for any $z \in \mathbb{R}^d$ with $||z - x|| \le \delta$, we have

$$\left|\frac{\nu_0(B_{z,\|z-x\|})}{\lambda(B_{z,\|z-x\|})} - f_0(x)\right| \le \epsilon f_0(x), \qquad \left|\frac{\nu_1(B_{x,\|z-x\|})}{\lambda(B_{x,\|z-x\|})}\right| \le \epsilon.$$

Recall that $W = \nu_0(B_{Z,||x-Z||})$. Then if $||Z - x|| \le \delta$, we have

$$W \ge (1 - \epsilon) f_0(x) \lambda(B_{Z, \|x - Z\|}) = (1 - \epsilon) f_0(x) \lambda(B_{x, \|x - Z\|}) \ge \epsilon^{-1} (1 - \epsilon) f_0(x) \nu_1(B_{x, \|x - Z\|}).$$

Proceeding in the same way as (\$3.3), we obtain

$$\begin{split} \mathrm{E}\big[\nu_1\big(A_M(x)\big)\big] &\leq \mathrm{P}\Big(W \leq \nu_0\big(B_{Z,\|\mathcal{X}_{(M)}(Z) - Z\|}\big)\eta_N\frac{M}{N_0}, \|Z - x\| \leq \delta\Big) + \mathrm{P}\Big(U_{(M)} > \eta_N\frac{M}{N_0}\Big) \\ &\leq \mathrm{P}\Big(\frac{1 - \epsilon}{\epsilon}f_0(x)U \leq U_{(M)}\Big) + \mathrm{P}\Big(U_{(M)} > \eta_N\frac{M}{N_0}\Big). \end{split}$$

For the first term above,

$$\frac{N_0}{M} P\left(\frac{1-\epsilon}{\epsilon} f_0(x) U \le U_{(M)}\right) = \frac{\epsilon}{1-\epsilon} \frac{1}{f_0(x)} \int_0^{\frac{1-\epsilon}{\epsilon} f_0(x) \frac{N_0}{M}} P\left(U_{(M)} \ge \frac{M}{N_0} t\right) dt
\le \frac{\epsilon}{1-\epsilon} \frac{1}{f_0(x)} \int_0^{\infty} P\left(\frac{N_0}{M} U_{(M)} \ge t\right) dt
= \frac{\epsilon}{1-\epsilon} \frac{1}{f_0(x)} \frac{N_0}{M} E[U_{(M)}] = \frac{\epsilon}{1-\epsilon} \frac{1}{f_0(x)} \frac{N_0}{N_0 + 1}.$$

By (S3.4) and noticing ϵ is arbitrary, we have

$$\lim_{N_0 \to \infty} \frac{N_0}{M} E[\nu_1(A_M(x))] = 0 = r(x).$$
 (S3.11)

Combining (S3.10) and (S3.11) completes the proof of the first claim.

Part II. This part proves the second claim. We also separate the proof of Part II into two cases based on the value of $f_1(x)$.

Case II.1. $f_1(x) > 0$. Again, for any $\epsilon \in (0, 1)$, we take δ in the same way as in Case I.1. Let $\eta_N = \eta_{N,p} = 4p \log(N_0/M)$. We also take N_0 sufficiently large so that $\eta_N \frac{M}{N_0} = 4p \frac{M}{N_0} \log(\frac{N_0}{M}) < (1 - \epsilon) f_0(x) \lambda(B_{0,\delta})$.

Let $\widetilde{Z}_1, \ldots, \widetilde{Z}_p$ be p independent copies that are drawn from ν_1 independent of the data. Then

$$\begin{split} & \mathrm{E}\big[\nu_{1}^{p}\big(A_{M}(x)\big)\big] \\ & = \mathrm{P}\big(\widetilde{Z}_{1}, \ldots, \widetilde{Z}_{p} \in A_{M}(x)\big) \\ & = \mathrm{P}\big(\nu_{0}(B_{\widetilde{Z}_{1}, \|x - \widetilde{Z}_{1}\|}) \leq \nu_{0}(B_{\widetilde{Z}_{1}, \|\mathcal{X}_{(M)}(\widetilde{Z}_{1}) - \widetilde{Z}_{1}\|}), \ldots, \nu_{0}(B_{\widetilde{Z}_{p}, \|x - \widetilde{Z}_{p}\|}) \leq \nu_{0}(B_{\widetilde{Z}_{p}, \|\mathcal{X}_{(M)}(\widetilde{Z}_{p}) - \widetilde{Z}_{p}\|})\big). \end{split}$$

Let $W_k = \nu_0(B_{\widetilde{Z}_k, \|x - \widetilde{Z}_k\|})$ and $V_k = \nu_0(B_{\widetilde{Z}_k, \|\mathcal{X}_{(M)}(\widetilde{Z}_k) - \widetilde{Z}_k\|})$ for any $k \in [p]$. Then $[W_k]_{k=1}^p$ are i.i.d. since $[\widetilde{Z}_k]_{k=1}^p$ are i.i.d. For any $k \in [p]$ and $\widetilde{Z}_k \in \mathbb{R}^d$ given, $V_k \mid \widetilde{Z}_k$ has the same

distribution as $U_{(M)}$. Then for any $k \in [p]$, V_k has the same distribution as $U_{(M)}$, and V_k is independent of \widetilde{Z}_k .

Let $W_{\max} = \max_{k \in \llbracket p \rrbracket} W_k$ and $V_{\max} = \max_{k \in \llbracket p \rrbracket} V_k$. Then

$$E[\nu_1^p(A_M(x))] \le P(W_{\text{max}} \le V_{\text{max}})$$

$$\le P\left(W_{\text{max}} \le V_{\text{max}} \le \eta_N \frac{M}{N_0}\right) + P\left(V_{\text{max}} > \eta_N \frac{M}{N_0}\right). \tag{S3.12}$$

For the second term in (S3.12),

$$P\bigg(V_{\max} > \eta_N \frac{M}{N_0}\bigg) \leq \sum_{k=1}^p P\bigg(V_k > \eta_N \frac{M}{N_0}\bigg) = pP\bigg(U_{(M)} > \eta_N \frac{M}{N_0}\bigg).$$

Proceeding as (S3.4),

$$\left(\frac{N_0}{M}\right)^p P\left(U_{(M)} > \eta_N \frac{M}{N_0}\right) \le \left(\frac{N_0}{M}\right)^p \exp\left(-\frac{1}{2}\eta_N M\right) = \left(\frac{N_0}{M}\right)^{p(1-2M)}.$$

We then obtain

$$\lim_{N_0 \to \infty} \left(\frac{N_0}{M}\right)^p P\left(V_{\text{max}} > \eta_N \frac{M}{N_0}\right) = 0.$$
 (S3.13)

For the first term in (S3.12), notice that $[\nu_1(B_{x,\|\widetilde{Z}_k-x\|})]_{k=1}^p$ are i.i.d. from U(0,1) since $[\widetilde{Z}_k]_{k=1}^p$ are i.i.d. We then have

$$\begin{split} &\left(\frac{N_{0}}{M}\right)^{p} \mathbf{P}\left(W_{\max} \leq V_{\max} \leq \eta_{N} \frac{M}{N_{0}}\right) \\ &= \left(\frac{N_{0}}{M}\right)^{p} \mathbf{P}\left(W_{\max} \leq V_{\max} \leq \eta_{N} \frac{M}{N_{0}}, \max_{k \in \llbracket p \rrbracket} \lVert \widetilde{Z}_{k} - x \rVert \leq \delta\right) \\ &\leq \left(\frac{N_{0}}{M}\right)^{p} \mathbf{P}\left(\frac{1 - \epsilon}{1 + \epsilon} \frac{f_{0}(x)}{f_{1}(x)} \max_{k \in \llbracket p \rrbracket} \nu_{1}(B_{x, \parallel \widetilde{Z}_{k} - x \parallel}) \leq V_{\max} \leq \eta_{N} \frac{M}{N_{0}}, \max_{k \in \llbracket p \rrbracket} \lVert \widetilde{Z}_{k} - x \rVert \leq \delta\right) \\ &\leq \left(\frac{N_{0}}{M}\right)^{p} \mathbf{P}\left(\frac{1 - \epsilon}{1 + \epsilon} \frac{f_{0}(x)}{f_{1}(x)} \max_{k \in \llbracket p \rrbracket} \nu_{1}(B_{x, \parallel \widetilde{Z}_{k} - x \parallel}) \leq V_{\max}\right) \\ &= \left(\frac{N_{0}}{M}\right)^{p} \int_{0}^{1} pt^{p-1} \mathbf{P}\left(V_{\max} \geq \frac{1 - \epsilon}{1 + \epsilon} \frac{f_{0}(x)}{f_{1}(x)}t \mid \max_{k \in \llbracket p \rrbracket} \nu_{1}(B_{x, \parallel \widetilde{Z}_{k} - x \parallel}) = t\right) dt \\ &= p\left(\frac{1 + \epsilon}{1 - \epsilon} \frac{f_{1}(x)}{f_{0}(x)}\right)^{p} \\ &\times \int_{0}^{\frac{1 - \epsilon}{1 + \epsilon} \frac{f_{0}(x)}{f_{1}(x)} \frac{N_{0}}{M}} t^{p-1} \mathbf{P}\left(V_{\max} \geq \frac{M}{N_{0}}t \mid \max_{k \in \llbracket p \rrbracket} \nu_{1}(B_{x, \parallel \widetilde{Z}_{k} - x \parallel}) = \frac{1 + \epsilon}{1 - \epsilon} \frac{f_{1}(x)}{f_{0}(x)} \frac{M}{N_{0}}t\right) dt \end{split}$$

$$= p \left(\frac{1+\epsilon}{1-\epsilon} \frac{f_1(x)}{f_0(x)} \right)^p \left[\int_0^1 t^{p-1} P\left(V_{\text{max}} \ge \frac{M}{N_0} t \, \Big| \, \max_{k \in \llbracket p \rrbracket} \nu_1(B_{x, \parallel \widetilde{Z}_k - x \parallel}) = \frac{1+\epsilon}{1-\epsilon} \frac{f_1(x)}{f_0(x)} \frac{M}{N_0} t \right) dt \right. \\ + \int_1^{\frac{1-\epsilon}{1+\epsilon} \frac{f_0(x)}{f_1(x)} \frac{N_0}{M}} t^{p-1} P\left(V_{\text{max}} \ge \frac{M}{N_0} t \, \Big| \, \max_{k \in \llbracket p \rrbracket} \nu_1(B_{x, \parallel \widetilde{Z}_k - x \parallel}) = \frac{1+\epsilon}{1-\epsilon} \frac{f_1(x)}{f_0(x)} \frac{M}{N_0} t \right) dt \right].$$

For the first term,

$$\int_0^1 t^{p-1} P\bigg(V_{\max} \ge \frac{M}{N_0} t \, \Big| \, \max_{k \in [\![p]\!]} \nu_1(B_{x, \|\widetilde{Z}_k - x\|}) = \frac{1 + \epsilon}{1 - \epsilon} \frac{f_1(x)}{f_0(x)} \frac{M}{N_0} t \bigg) \, \mathrm{d}t \le \int_0^1 t^{p-1} \, \mathrm{d}t = \frac{1}{p}.$$

For the second term, using the Chernoff bound, conditional on $\widetilde{\mathbf{Z}} = (\widetilde{Z}_1, \dots, \widetilde{Z}_p)$,

$$\int_{1}^{\frac{1-\epsilon}{1+\epsilon} \frac{f_0(x)}{f_1(x)} \frac{N_0}{M}} t^{p-1} P\left(V_{\text{max}} \ge \frac{M}{N_0} t \,\middle|\, \widetilde{\mathbf{Z}}\right) dt
\le \int_{0}^{\infty} (1+t)^{p-1} P\left(V_{\text{max}} \ge \frac{M}{N_0} (1+t) \,\middle|\, \widetilde{\mathbf{Z}}\right) dt
\le \int_{0}^{\infty} (1+t)^{p-1} \left[\sum_{k=1}^{p} P\left(V_k \ge \frac{M}{N_0} (1+t) \,\middle|\, \widetilde{\mathbf{Z}}\right) \right] dt
= p \int_{0}^{\infty} (1+t)^{p-1} P\left(U_{(M)} \ge \frac{M}{N_0} (1+t)\right) dt
\le p \int_{0}^{\infty} (1+t)^{p-1} (1+t)^{M} \exp(-tM) dt \le \sqrt{2\pi} p M^{-1/2} \left(1 + \frac{1}{M}\right)^{p-1} (1+o(1)),$$

where the last step follows from Stirling's approximation with $M \to \infty$.

Then we obtain

$$\limsup_{N_0 \to \infty} \left(\frac{N_0}{M}\right)^p P\left(W_{\text{max}} \le V_{\text{max}}, V_{\text{max}} \le \eta_N \frac{M}{N_0}\right) \le \left(\frac{1+\epsilon}{1-\epsilon} \frac{f_1(x)}{f_0(x)}\right)^p. \tag{S3.14}$$

Plugging (S3.13) and (S3.14) into (S3.12) yields

$$\limsup_{N_0 \to \infty} \left(\frac{N_0}{M}\right)^p \mathbb{E}\left[\nu_1^p \left(A_M(x)\right)\right] \le \left(\frac{1+\epsilon}{1-\epsilon} \frac{f_1(x)}{f_0(x)}\right)^p = \left(\frac{1+\epsilon}{1-\epsilon} r(x)\right)^p. \tag{S3.15}$$

Lastly, using Hölder's inequality,

$$\left(\frac{N_0}{M}\right)^p \mathrm{E}\big[\nu_1^p\big(A_M(x)\big)\big] \ge \left[\frac{N_0}{M} \mathrm{E}\big[\nu_1\big(A_M(x)\big)\big]\right]^p.$$

Employing the first claim, we have

$$\liminf_{N_0 \to \infty} \left(\frac{N_0}{M} \right)^p \mathbb{E} \left[\nu_1^p \left(A_M(x) \right) \right] \ge \left[r(x) \right]^p.$$
(S3.16)

Combining (S3.15) with (S3.16) and noting that ϵ is arbitrary, we obtain

$$\lim_{N_0 \to \infty} \left(\frac{N_0}{M}\right)^p \mathrm{E}\left[\nu_1^p \left(A_M(x)\right)\right] = \left[r(x)\right]^p. \tag{S3.17}$$

Case II.2. $f_1(x) = 0$. For any $\epsilon \in (0, 1)$, we take δ in the same way as in the proof of Case I.2 and take η_N as in the proof of Case II.1. By (S3.12),

$$\left(\frac{N_0}{M}\right)^p \mathrm{E}\big[\nu_1^p\big(A_M(x)\big)\big] \leq \left(\frac{N_0}{M}\right)^p \mathrm{P}\bigg(W_{\max} \leq V_{\max} \leq \eta_N \frac{M}{N_0}\bigg) + \left(\frac{N_0}{M}\right)^p \mathrm{P}\bigg(V_{\max} > \eta_N \frac{M}{N_0}\bigg).$$

For the first term,

$$\begin{split} &\left(\frac{N_0}{M}\right)^p \mathbf{P}\bigg(W_{\max} \leq V_{\max} \leq \eta_N \frac{M}{N_0}\bigg) \\ &\leq \left(\frac{N_0}{M}\right)^p \mathbf{P}\bigg(\frac{1-\epsilon}{\epsilon} f_0(x) \max_{k \in \llbracket p \rrbracket} \nu_1(B_{x, \lVert \widetilde{Z}_k - x \rVert}) \leq V_{\max}\bigg) \\ &= \left(\frac{N_0}{M}\right)^p \int_0^1 pt^{p-1} \mathbf{P}\bigg(V_{\max} \geq \frac{1-\epsilon}{\epsilon} f_0(x)t \ \Big| \ \max_{k \in \llbracket p \rrbracket} \nu_1(B_{x, \lVert \widetilde{Z}_k - x \rVert}) = t\bigg) \, \mathrm{d}t \\ &= p\bigg(\frac{\epsilon}{1-\epsilon} \frac{1}{f_0(x)}\bigg)^p \int_0^{\frac{1-\epsilon}{\epsilon} f_0(x) \frac{N_0}{M}} t^{p-1} \mathbf{P}\bigg(V_{\max} \geq \frac{M}{N_0} t \ \Big| \ \max_{k \in \llbracket p \rrbracket} \nu_1(B_{x, \lVert \widetilde{Z}_k - x \rVert}) = t\bigg) \, \mathrm{d}t. \end{split}$$

Then proceeding in the same way as (S3.14), we have

$$\limsup_{N_0 \to \infty} \left(\frac{N_0}{M}\right)^p P\left(W_{\max} \le V_{\max} \le \eta_N \frac{M}{N_0}\right) \le \left(\frac{\epsilon}{1 - \epsilon} \frac{1}{f_0(x)}\right)^p.$$

Lastly, using (S3.13) and noting again that ϵ is arbitrary, we obtain

$$\lim_{N_0 \to \infty} \left(\frac{N_0}{M}\right)^p \mathrm{E}\left[\nu_1^p \left(A_M(x)\right)\right] = 0 = \left[r(x)\right]^p. \tag{S3.18}$$

Combining (S3.17) and (S3.18) then completes the proof of the second claim. Q.E.D

S3.2. Proof of Theorem B.1

PROOF OF THEOREM B.1(i): By (2.4) and that $[Z_j]_{j=1}^{N_1}$ are i.i.d,

$$E[\widehat{r}_{M}(x)] = E\left[\frac{N_{0}}{N_{1}}\frac{K_{M}(x)}{M}\right] = \frac{N_{0}}{N_{1}M}E\left[\sum_{i=1}^{N_{1}}\mathbb{1}(Z_{i} \in A_{M}(x))\right] = \frac{N_{0}}{M}E[\nu_{1}(A_{M}(x))].$$

Employing Lemma B.1 then completes the proof.

PROOF OF THEOREM B.1(ii): By Hölder's inequality, it suffices to consider the case when p is even. Because x^p is convex for p > 1 and x > 0, we have

$$E[|\widehat{r}_{M}(x) - r(x)|^{p}]$$

$$\leq 2^{p-1} \left(E[|\widehat{r}_{M}(x) - E[\widehat{r}_{M}(x)|X]|^{p}] + E[|E[\widehat{r}_{M}(x)|X] - r(x)|^{p}] \right). \quad (S3.19)$$

For the second term in (S3.19), Lemma B.1 implies

$$\lim_{N_0 \to \infty} \mathbb{E}[|\mathbb{E}[\widehat{r}_M(x) | X] - r(x)|^p] = \lim_{N_0 \to \infty} \mathbb{E}\left[\left|\frac{N_0}{M}\nu_1(A_M(x)) - r(x)\right|^p\right] = 0$$
 (S3.20)

by expanding the product term.

For the first term in (S3.19), noticing that $[Z_j]_{j=1}^{N_1}$ are i.i.d, we have $K_M(x)|X \sim \text{Bin}(N_1, \nu_1(A_M(x)))$. Using Lemma B.1 and $MN_1/N_0 \to \infty$, for any positive integers p and q, we have

$$\lim_{N_0 \to \infty} \left(\frac{N_0}{N_1 M} \right)^p \mathbf{E} \left[N_1^p \nu_1^p \left(A_M(x) \right) \right] = \left[r(x) \right]^p,$$

$$\lim_{N_0 \to \infty} \left(\frac{N_0}{N_1 M} \right)^p \left(\frac{N_0}{M} \right)^q \mathbf{E} \left[N_1^p \nu_1^{p+q} \left(A_M(x) \right) \right] = \left[r(x) \right]^{p+q},$$

and then $\mathrm{E}[N_1^p \nu_1^p(A_M(x))]$ is the dominated term among $[\mathrm{E}[N_1^k \nu_1^{k+q}(A_M(x))]]_{k \leq p, q \geq 0}$. To complete the proof, for any positive integer c and $Z \sim \mathrm{Bin}(n, p')$, let $\mu_c = \mathrm{E}[(Z - \mathrm{E}[Z])^c]$ be the cth central moment. By Romanovsky (1923), we have

$$\mu_{c+1} = p'(1-p')\left(nc\mu_{c-1} + \frac{\mathrm{d}\mu_c}{\mathrm{d}p'}\right).$$

Then for even p, we obtain

$$\mathrm{E}\big[\big(K_M(x)-N_1\nu_1\big(A_M(x)\big)\big)^p\big]\lesssim \mathrm{E}\big[N_1\nu_1\big(A_M(x)\big)\big]^{p/2}\lesssim \left(\frac{N_1M}{N_0}\right)^{p/2}.$$

The first term in (S3.19) then satisfies

$$\mathrm{E}\big[\big|\widehat{r}_{M}(x) - \mathrm{E}\big[\widehat{r}_{M}(x) \,|\, X\big]\big|^{p}\big] = \left(\frac{N_{0}}{N_{1}M}\right)^{p} \mathrm{E}\big[\big(K_{M}(x) - N_{1}\nu_{1}\big(A_{M}(x)\big)\big)^{p}\big] \lesssim \left(\frac{N_{0}}{N_{1}M}\right)^{p/2}.$$

Since $MN_1/N_0 \to \infty$, we obtain

$$\lim_{N_0 \to \infty} \mathbb{E}\left[\left|\widehat{r}_M(x) - \mathbb{E}\left[\widehat{r}_M(x) \mid X\right]\right|^p\right] = 0.$$
 (S3.21)

Plugging (S3.20) and (S3.21) into (S3.19) then completes the proof. Q.E.D.

PROOF OF THEOREM B.2: We first cite the Hardy–Littlewood maximal inequality.

LEMMA S3.2—Hardy-Littlewood Maximal Inequality (Stein (2016)): For any locally integrable function $f: \mathbb{R}^d \to \mathbb{R}$, define

$$\mathsf{M} f(x) = \sup_{\delta > 0} \frac{1}{\lambda(B_{x,\delta})} \int_{B_{x,\delta}} \left| f(z) \right| \mathrm{d} z.$$

Then for $d \ge 1$, there exists a constant $C_d > 0$ only depending on d such that for all t > 0 and $f \in L_1(\mathbb{R}^d)$, we have

$$\lambda\big(\big\{x:\mathsf{M}f(x)>t\big\}\big)<\frac{C_d}{t}\|f\|_{L_1},$$

where $\|\cdot\|_{L_1}$ stands for the function L_1 norm.

Let $\epsilon > 0$ be given. We assume $\epsilon \le f_L$. From Assumption B.1, S_0 and S_1 are bounded, then ν_0 and ν_1 are compactly supported. Since $f_0, f_1 \in L_1$, and the class of continuous functions are dense in the class of compactly supported L_1 functions from simple use of Lusin's theorem, we can find g_0 , g_1 such that g_0 , g_1 are continuous and $||f_0 - g_0||_{L_1} \le \epsilon^3$ and $||f_1 - g_1||_{L_1} \le \epsilon^3$.

Since g_0 , g_1 are continuous with compact supports, they are uniformly continuous, that is, there exists $\delta > 0$ such that for any $x, z \in \mathbb{R}^d$ and $||z - x|| \le \delta$, we have $|g_0(x) - g_0(z)| \le \delta$ $\frac{\epsilon^2}{3} \text{ and } |g_1(x) - g_1(z)| \le \frac{\epsilon^2}{3}.$ For any $x \in \mathbb{R}^d$, we have

$$\frac{1}{\lambda(B_{x,\delta})} \int_{B_{x,\delta}} |f_0(x) - f_0(z)| dz$$

$$\leq \frac{1}{\lambda(B_{x,\delta})} \int_{B_{x,\delta}} \left[|f_0(x) - g_0(x)| + |g_0(x) - g_0(z)| + |f_0(z) - g_0(z)| \right] dz$$

$$= |f_0(x) - g_0(x)| + \frac{1}{\lambda(B_{x,\delta})} \int_{B_{x,\delta}} |g_0(x) - g_0(z)| dz$$

$$+ \frac{1}{\lambda(B_{x,\delta})} \int_{B_{x,\delta}} |f_0(z) - g_0(z)| dz. \tag{S3.22}$$

For the first term in (S3.22), using Markov's inequality, we have

$$\lambda(\{x: |f_0(x) - g_0(x)| > \epsilon^2/3\}) \le 3\epsilon^{-2} ||f_0 - g_0||_{L_1} \le 3\epsilon.$$
 (S3.23)

For the second term in (S3.22), by the selection of δ ,

$$\frac{1}{\lambda(B_{x,\delta})} \int_{B_{x,\delta}} |g_0(x) - g_0(z)| \, \mathrm{d}z \le \max_{z \in B_{x,\delta}} |g_0(x) - g_0(z)| \le \frac{\epsilon^2}{3}. \tag{S3.24}$$

For the third term,

$$\frac{1}{\lambda(B_{x,\delta})} \int_{B_{x,\delta}} \left| f_0(z) - g_0(z) \right| dz \le \sup_{\delta > 0} \frac{1}{\lambda(B_{x,\delta})} \int_{B_{x,\delta}} \left| f_0(z) - g_0(z) \right| dz = \mathsf{M}(f_0 - g_0)(x).$$

Lemma S3.2 then yields

$$\lambda(\{x: \mathsf{M}(f_0 - g_0)(x) > \epsilon^2/3\}) < 3C_d \epsilon^{-2} \|f_0 - g_0\|_{L_1} \le 3C_d \epsilon. \tag{S3.25}$$

We can establish similar results for f_1 , g_1 .

Let

$$A_{1} = \left\{ x : \left| f_{0}(x) - g_{0}(x) \right| > \epsilon^{2}/3 \right\} \cup \left\{ x : \left| f_{1}(x) - g_{1}(x) \right| > \epsilon^{2}/3 \right\}$$
$$\cup \left\{ x : \mathsf{M}(f_{0} - g_{0})(x) > \epsilon^{2}/3 \right\} \cup \left\{ x : \mathsf{M}(f_{1} - g_{1})(x) > \epsilon^{2}/3 \right\}.$$

Plugging (S3.23), (S3.24), (S3.25) into (S3.22), for any $x \notin A_1$ and $||z - x|| \le \delta$, we have

$$\frac{1}{\lambda(B_{x,\delta})} \int_{B_{x,\delta}} \left| f_0(x) - f_0(z) \right| \mathrm{d}z \le \epsilon^2, \qquad \frac{1}{\lambda(B_{x,\delta})} \int_{B_{x,\delta}} \left| f_1(x) - f_1(z) \right| \mathrm{d}z \le \epsilon^2,$$

and $\lambda(A_1) \leq 6(C_d + 1)\epsilon$.

Let $A_2 = \{x : f_1(x) \le \epsilon\}$. We then separate the proof into three cases. In the following, it suffices to consider $f_0(x) > 0$ due to the definition of L_p risk.

Case I. $x \notin A_1 \cup A_2$. By $\epsilon \leq f_L$ and the definition of A_2 , for any $x \notin A_1 \cup A_2$ and $||z - x|| \leq \delta$,

$$\frac{1}{\lambda(B_{x,\delta})} \int_{B_{x,\delta}} \left| f_0(x) - f_0(z) \right| dz \le \epsilon^2 \le \epsilon f_L \le \epsilon f_0(x),$$

$$\frac{1}{\lambda(B_{x,\delta})} \int_{B_{x,\delta}} \left| f_1(x) - f_1(z) \right| dz \le \epsilon^2 \le \epsilon f_1(x).$$

We then obtain for $w \in \{0, 1\}$,

$$\left|\frac{\nu_w(B_{x,\|z-x\|})}{\lambda(B_{x,\|z-x\|})} - f_w(x)\right| \le \epsilon f_w(x), \qquad \left|\frac{\nu_w(B_{z,\|z-x\|})}{\lambda(B_{z,\|z-x\|})} - f_w(x)\right| \le \epsilon f_w(x).$$

Let $\eta_N = \eta_{N,p} = 4p \log(N_0/M)$. We also take N_0 large enough so that $\eta_N \frac{M}{N_0} = 4p \frac{M}{N_0} \log(\frac{N_0}{M}) < (1 - \epsilon) f_L \lambda(B_{0,\delta})$. Then for any $x \in \mathbb{R}^d$ such that $f_0(x) > 0$, we have $\eta_N \frac{M}{N_0} < (1 - \epsilon) f_0(x) \lambda(B_{0,\delta})$.

Proceeding as in the proof of Case II.1 of Lemma B.1 and also Theorem B.1 by using Fubini's theorem, since ϵ is arbitrary, we obtain

$$\lim_{N_0 \to \infty} \mathbb{E} \left[\int_{\mathbb{R}^d} |\widehat{r}_M(x) - r(x)|^p f_0(x) \mathbb{1}(x \notin A_1 \cup A_2) \, \mathrm{d}x \right] = 0.$$
 (S3.26)

Case II. $x \in A_2 \setminus A_1$. In this case, we have

$$\left|\frac{\nu_0(B_{x,\|z-x\|})}{\lambda(B_{x,\|z-x\|})} - f_0(x)\right| \le \epsilon f_0(x), \qquad \left|\frac{\nu_0(B_{z,\|z-x\|})}{\lambda(B_{z,\|z-x\|})} - f_0(x)\right| \le \epsilon f_0(x),$$

$$\left|\frac{\nu_1(B_{x,\|z-x\|})}{\lambda(B_{x,\|z-x\|})} - f_1(x)\right| \le \epsilon^2, \qquad \left|\frac{\nu_1(B_{z,\|z-x\|})}{\lambda(B_{z,\|z-x\|})} - f_1(x)\right| \le \epsilon^2.$$

O.E.D.

Take η_N and take N_0 sufficiently large as in Case I above. Proceeding as the proof of Case II.2 of Lemma B.1 and also Theorem B.1 by using Fubini's theorem, since ϵ is arbitrary, we obtain

$$\lim_{N_0 \to \infty} \mathbb{E} \left[\int_{\mathbb{R}^d} \left| \widehat{r}_M(x) - r(x) \right|^p f_0(x) \mathbb{1}(x \in A_2 \setminus A_1) \, \mathrm{d}x \right] = 0.$$
 (S3.27)

Case III. $x \in A_1$. In this case, for any $x \in A_1$ and $z \in S_1$, $\nu_0(B_{z,\|z-x\|}) \ge f_L \lambda(B_{z,\|z-x\|} \cap S_0) \ge af_L \lambda(B_{z,\|z-x\|}) \ge \frac{af_L}{f_U} \nu_1(B_{x,\|z-x\|})$. Then for any $x \in A_1$, from (S3.12) and in the same way as (S3.14),

$$\begin{split} \left(\frac{N_0}{M}\right)^p \mathrm{E}\left[\nu_1^p \left(A_M(x)\right)\right] &\leq \left(\frac{N_0}{M}\right)^p \mathrm{P}(W_{\max} \leq V_{\max}) \\ &\leq \left(\frac{N_0}{M}\right)^p \mathrm{P}\left(\frac{af_L}{f_U} \max_{k \in \llbracket p \rrbracket} \nu_1(B_{x, \parallel \widetilde{Z}_k - x \parallel}) \leq V_{\max}\right) \\ &\leq \left(\frac{f_U}{af_L}\right)^p \left(1 + o(1)\right) = O(1). \end{split}$$

Proceeding as in the proof of Theorem B.1, and due to the boundedness assumptions on f_0 and f_1 , for any $x \in A_1$ and p even,

$$\mathbb{E}[|\widehat{r}_{M}(x) - r(x)|^{p}] \lesssim \mathbb{E}[|\widehat{r}_{M}(x) - \mathbb{E}[\widehat{r}_{M}(x) | X]|^{p}] + \mathbb{E}[(\mathbb{E}[\widehat{r}_{M}(x) | X])^{p}] + |r(x)|^{p} \lesssim 1.$$

Then

$$\mathrm{E}\bigg[\int_{\mathbb{R}^d} \big|\widehat{r}_M(x) - r(x)\big|^p f_0(x) \mathbb{1}(x \in A_1) \,\mathrm{d}x\bigg] \lesssim f_U \lambda(A_1) \lesssim \epsilon.$$

Since ϵ is arbitrary, we have

$$\lim_{N_0 \to \infty} \mathbb{E} \left[\int_{\mathbb{R}^d} \left| \widehat{r}_M(x) - r(x) \right|^p f_0(x) \mathbb{1}(x \in A_1) \, \mathrm{d}x \right] = 0.$$
 (S3.28)

Combining (S3.26), (S3.27), and (S3.28) completes the proof.

PROOF OF COROLLARY B.1: Corollary B.1 can be established following the same way as that of Theorem B.2 but with less effort since we only have to show

$$\lim_{N_0 \to \infty} \mathbb{E} \left[\int_{\mathbb{R}^d} \left| \mathbb{E} \left[\widehat{r}_M(x) \mid X \right] - r(x) \right|^p f_0(x) \, \mathrm{d}x \right] = 0.$$

In detail, denote the Radon–Nikodym derivative of the probability measure of W with respect to ν_0 by r_W . We then have

$$\limsup_{N_0 \to \infty} \mathbb{E} \left[\left| \frac{N_0}{M} \nu_1 (A_M(W)) - r(W) \right|^p \right] \\
= \limsup_{N_0 \to \infty} \mathbb{E} \left[\int_{\mathbb{R}^d} \left| \frac{N_0}{M} \nu_1 (A_M(x)) - r(x) \right|^p r_W(x) f_0(x) \, \mathrm{d}x \right]$$

$$\begin{split} &\lesssim \limsup_{N_0 \to \infty} \mathrm{E} \bigg[\int_{\mathbb{R}^d} \bigg| \frac{N_0}{M} \nu_1 \big(A_M(x) \big) - r(x) \bigg|^p f_0(x) \, \mathrm{d}x \bigg] \\ &= \limsup_{N_0 \to \infty} \mathrm{E} \bigg[\int_{\mathbb{R}^d} \big| \mathrm{E} \big[\widehat{r}_M(x) \, \big| \, X \big] - r(x) \big|^p f_0(x) \, \mathrm{d}x \bigg] = 0, \end{split}$$

where the last line has been established in the proof of Theorem B.2. Noticing that $E[r(W)]^p$ is bounded under Assumption B.1, the proof is thus complete. *Q.E.D.*

We only have to prove the first two claims as the rest are trivial.

PROOF OF THEOREM B.3(i): For any $z \in \mathbb{R}^d$ such that $||z - x|| \le \delta/2$, since $B_{z,||z-x||} \subset B_{x,2||z-x||} \subset B_{x,\delta}$, we have

$$\begin{split} & \left| \frac{\nu_0(B_{z,\|z-x\|})}{\lambda(B_{z,\|z-x\|})} - f_0(x) \right| \leq \frac{1}{\lambda(B_{z,\|z-x\|})} \int_{B_{z,\|z-x\|}} \left| f_0(y) - f_0(x) \right| \, \mathrm{d}y \leq 2L\|z-x\|, \\ & \left| \frac{\nu_1(B_{x,\|z-x\|})}{\lambda(B_{x,\|z-x\|})} - f_1(x) \right| \leq \frac{1}{\lambda(B_{x,\|z-x\|})} \int_{B_{x,\|z-x\|}} \left| f_1(y) - f_1(x) \right| \, \mathrm{d}y \leq L\|z-x\|. \end{split}$$

Consider any $\delta_N > 0$ such that $\delta_N \le \delta/2$. If $||z - x|| \le \delta_N$ and $f_0(x) > 2L\delta_N$, then

$$\frac{f_0(x) - 2L\delta_N}{f_1(x) + L\delta_N} \le \frac{\nu_0(B_{z,\|x-z\|})}{\lambda(B_{z,\|x-z\|})} \frac{\lambda(B_{x,\|x-z\|})}{\nu_1(B_{x,\|x-z\|})}.$$

If further $f_1(x) > L\delta_N$, then

$$\frac{\nu_0(B_{z,\|x-z\|})}{\lambda(B_{z,\|x-z\|})} \frac{\lambda(B_{x,\|x-z\|})}{\nu_1(B_{x,\|x-z\|})} \leq \frac{f_0(x) + 2L\delta_N}{f_1(x) - L\delta_N}.$$

On the other hand, if $||z - x|| \ge \delta_N$ and $f_0(x) > 2L\delta_N$, $\nu_0(B_{z,||z-x||}) \ge (f_0(x) - 2L\delta_N) \times \lambda(B_{0,\delta_N}) = (f_0(x) - 2L\delta_N)V_d\delta_N^d$, where V_d is the Lebesgue measure of the unit ball on \mathbb{R}^d .

Let $\delta_N = (\frac{4}{f_L V_d})^{1/d} (\frac{M}{N_0})^{1/d}$. Since $M/N_0 \to 0$, we have $\delta_N \to 0$ as $N_0 \to \infty$. Taking N_0 large enough so that $\delta_N < f_L/(4L)$ and $\delta_N \le \delta/2$, then $2LV_d \delta_N^{d+1} = \frac{M}{N_0} \frac{8L}{f_L} \delta_N < 2\frac{M}{N_0}$. Then for any $(\nu_0, \nu_1) \in \mathcal{P}_{x,p}(f_L, f_U, L, d, \delta)$,

$$(f_0(x) - 2L\delta_N)V_d\delta_N^d > 4\frac{f_0(x)}{f_L}\frac{M}{N_0} - 2\frac{M}{N_0} \ge 2\frac{M}{N_0}.$$

With a slight abuse of notation, let $W = \nu_0(B_{Z,\|x-Z\|})$. Then $W \leq 2\frac{M}{N_0}$ implies that $\|Z - x\| \leq \delta_N$.

Depending on the value of $f_1(x)$, the proof is separated into two cases.

Case I. $f_1(x) > L\delta_N$.

Upper bound. Proceeding similar to (S3.3), we have

$$\begin{split} & E[\widehat{r}_{M}(x)] = \frac{N_{0}}{M} E[\nu_{1}(A_{M}(x))] = \frac{N_{0}}{M} P(W \leq \nu_{0}(B_{Z, \|\mathcal{X}_{(M)}(Z) - Z\|})) \\ & \leq \frac{N_{0}}{M} P\left(W \leq \nu_{0}(B_{Z, \|\mathcal{X}_{(M)}(Z) - Z\|}) \leq 2\frac{M}{N_{0}}\right) + \frac{N_{0}}{M} P\left(U_{(M)} > 2\frac{M}{N_{0}}\right) \\ & \leq \frac{N_{0}}{M} P\left(W \leq \nu_{0}(B_{Z, \|\mathcal{X}_{(M)}(Z) - Z\|}), \|Z - x\| \leq \delta_{N}\right) + \frac{N_{0}}{M} P\left(U_{(M)} > 2\frac{M}{N_{0}}\right) \\ & \leq \frac{N_{0}}{M} P\left(\frac{f_{0}(x) - 2L\delta_{N}}{f_{1}(x) + L\delta_{N}} \nu_{1}(B_{x, \|x - Z\|}) \leq \nu_{0}(B_{Z, \|\mathcal{X}_{(M)}(Z) - Z\|}), \|Z - x\| \leq \delta_{N}\right) \\ & + \frac{N_{0}}{M} P\left(U_{(M)} > 2\frac{M}{N_{0}}\right) \\ & \leq \frac{N_{0}}{M} P\left(\frac{f_{0}(x) - 2L\delta_{N}}{f_{1}(x) + L\delta_{N}} U \leq U_{(M)}\right) + \frac{N_{0}}{M} P\left(U_{(M)} > 2\frac{M}{N_{0}}\right). \end{split} \tag{S3.29}$$

For the second term in (S3.29), since $M/\log N_0 \to \infty$, for any $\gamma > 0$,

$$\frac{N_0}{M} P \left(U_{(M)} > 2 \frac{M}{N_0} \right) = \frac{N_0}{M} P \left(Bin \left(N_0, 2 \frac{M}{N_0} \right) \le M \right)
\le \frac{N_0}{M} N_0^{-(1 - \log 2)M/\log N_0} < N_0^{-\gamma}.$$
(S3.30)

For the first term in (S3.29), proceeding as (S3.5), we obtain

$$\frac{N_0}{M} P\left(\frac{f_0(x) - 2L\delta_N}{f_1(x) + L\delta_N} U \le U_{(M)}\right) \le \frac{f_1(x) + L\delta_N}{f_0(x) - 2L\delta_N} \frac{N_0}{N_0 + 1}.$$

Then we obtain

$$E[\widehat{r}_{M}(x)] \le \frac{f_{1}(x) + L\delta_{N}}{f_{0}(x) - 2L\delta_{N}} \frac{N_{0}}{N_{0} + 1} + o(N_{0}^{-\gamma}).$$
 (S3.31)

Lower bound. Proceeding similar to (S3.8), we have

$$\begin{split} & \mathrm{E} \big[\widehat{r}_{M}(x) \big] = \frac{N_{0}}{M} \mathrm{E} \big[\nu_{1} \big(A_{M}(x) \big) \big] = \frac{N_{0}}{M} \mathrm{P} \big(W \leq \nu_{0} \big(B_{Z, \| \mathcal{X}_{(M)}(Z) - Z \|} \big) \big) \\ & \geq \frac{N_{0}}{M} \mathrm{P} \bigg(W \leq \nu_{0} \big(B_{Z, \| \mathcal{X}_{(M)}(Z) - Z \|} \big) \leq 2 \frac{M}{N_{0}} \bigg) \\ & = \frac{N_{0}}{M} \mathrm{P} \bigg(W \leq \nu_{0} \big(B_{Z, \| \mathcal{X}_{(M)}(Z) - Z \|} \big) \leq 2 \frac{M}{N_{0}}, \, \| Z - x \| \leq \delta_{N} \bigg) \\ & \geq \frac{N_{0}}{M} \mathrm{P} \bigg(\frac{f_{0}(x) + 2L\delta_{N}}{f_{1}(x) - L\delta_{N}} \nu_{1} \big(B_{x, \| x - Z \|} \big) \leq \nu_{0} \big(B_{Z, \| \mathcal{X}_{(M)}(Z) - Z \|} \big) \leq 2 \frac{M}{N_{0}}, \, \| Z - x \| \leq \delta_{N} \bigg) \\ & = \frac{N_{0}}{M} \mathrm{P} \bigg(\frac{f_{0}(x) + 2L\delta_{N}}{f_{1}(x) - L\delta_{N}} \nu_{1} \big(B_{x, \| x - Z \|} \big) \leq \nu_{0} \big(B_{Z, \| \mathcal{X}_{(M)}(Z) - Z \|} \big) \leq 2 \frac{M}{N_{0}} \bigg) \end{split}$$

$$\geq \frac{N_0}{M} P\left(\frac{f_0(x) + 2L\delta_N}{f_1(x) - L\delta_N} U \leq U_{(M)}\right) - \frac{N_0}{M} P\left(U_{(M)} > 2\frac{M}{N_0}\right)$$

$$= \frac{f_1(x) - L\delta_N}{f_0(x) + 2L\delta_N} \int_0^{\frac{f_0(x) + 2L\delta_N}{f_1(x) - L\delta_N} \frac{N_0}{M}} P\left(U_{(M)} \geq \frac{M}{N_0}t\right) dt - \frac{N_0}{M} P\left(U_{(M)} > 2\frac{M}{N_0}\right).$$

Consider the first term. If $\frac{f_0(x)+2L\delta_N}{f_1(x)-L\delta_N} \geq 1$, then

$$\frac{f_1(x) - L\delta_N}{f_0(x) + 2L\delta_N} \int_0^{\frac{f_0(x) + 2L\delta_N}{f_1(x) - L\delta_N} \frac{N_0}{M}} P\left(U_{(M)} \ge \frac{M}{N_0}t\right) dt = \frac{f_1(x) - L\delta_N}{f_0(x) + 2L\delta_N} \frac{N_0}{N_0 + 1}.$$

If $\frac{f_0(x)+2L\delta_N}{f_1(x)-L\delta_N} < 1$, using the Chernoff bound, for any $\gamma > 0$,

$$\begin{split} & \int_{\frac{f_0(x)+2L\delta_N}{f_1(x)-L\delta_N}}^{\frac{N_0}{M}} \mathbf{P}\bigg(U_{(M)} \geq \frac{M}{N_0}t\bigg) \, \mathrm{d}t \\ & \leq \int_{\frac{f_L}{f_U}}^{\frac{N_0}{M}} \mathbf{P}\bigg(U_{(M)} \geq \frac{M}{N_0}t\bigg) \, \mathrm{d}t \leq \bigg[1-\frac{f_L}{f_U}\bigg] \frac{N_0}{M} \mathbf{P}\bigg(U_{(M)} \geq \frac{f_L}{f_U}\bigg) \\ & \leq \bigg[1-\frac{f_L}{f_U}\bigg] \frac{N_0}{M} \exp\bigg[M-\frac{f_L}{f_U}N_0-M\log M+M\log\bigg(\frac{f_L}{f_U}N_0\bigg)\bigg] \prec N_0^{-\gamma}. \end{split}$$

The last step is due to $M \log N_0/N_0 \rightarrow 0$. Recalling (S3.30), we then obtain

$$E[\hat{r}_{M}(x)] \ge \frac{f_{1}(x) - L\delta_{N}}{f_{0}(x) + 2L\delta_{N}} \frac{N_{0}}{N_{0} + 1} - o(N_{0}^{-\gamma}).$$
 (S3.32)

Combining (S3.31) and (S3.32), and taking N_0 large enough so that $L\delta_N \leq f_U \wedge (f_L/4)$, we obtain

$$\begin{split} &|\mathbf{E}[\widehat{r}_{M}(x)] - r(x)| \\ &\leq \left| \frac{f_{1}(x) + L\delta_{N}}{f_{0}(x) - 2L\delta_{N}} \frac{N_{0}}{N_{0} + 1} - \frac{f_{1}(x)}{f_{0}(x)} \right| \vee \left| \frac{f_{1}(x) - L\delta_{N}}{f_{0}(x) + 2L\delta_{N}} \frac{N_{0}}{N_{0} + 1} - \frac{f_{1}(x)}{f_{0}(x)} \right| \\ &+ o(N_{0}^{-\gamma}) \leq \frac{f_{0}(x)L\delta_{N} + 2f_{1}(x)L\delta_{N}}{f_{0}(x)(f_{0}(x) - 2L\delta_{N})} + \frac{1}{N_{0} + 1} \frac{f_{1}(x) + L\delta_{N}}{f_{0}(x) - 2L\delta_{N}} + o(N_{0}^{-\gamma}) \\ &\leq \left(\frac{2}{f_{L}} + \frac{4f_{U}}{f_{L}^{2}}\right)L\delta_{N} + \frac{4f_{U}}{f_{L}} \frac{1}{N_{0} + 1} + o(N_{0}^{-\gamma}). \end{split}$$

By the selection of δ_N and that the right-hand side does not depend on x, we complete the proof for this case.

Case II. $f_1(x) \le L\delta_N$. The upper bound (S3.31) in Case I still holds for this case. Accordingly, taking N_0 large enough so that $L\delta_N \le f_L/4$, we have

$$\begin{aligned} \left| \mathbf{E}[\widehat{r}_{M}(x)] - r(x) \right| &\leq \mathbf{E}[\widehat{r}_{M}(x)] + r(x) \\ &\leq \frac{f_{1}(x) + L\delta_{N}}{f_{0}(x) - 2L\delta_{N}} \frac{N_{0}}{N_{0} + 1} + \frac{f_{1}(x)}{f_{0}(x)} + o(N_{0}^{-\gamma}) \end{aligned}$$

$$\leq \frac{4}{f_L} L \delta_N + \frac{1}{f_L} L \delta_N + o(N_0^{-\gamma}).$$

We thus complete the whole proof.

Q.E.D.

PROOF OF THEOREM B.3(ii): By the law of total variance,

$$\operatorname{Var}[\widehat{r}_{M}(x)] = \operatorname{E}[\operatorname{Var}[\widehat{r}_{M}(x) | X]] + \operatorname{Var}[\operatorname{E}[\widehat{r}_{M}(x) | X]]. \tag{S3.33}$$

For the first term in (S3.33), let Z be a copy drawn from ν_1 independently of the data. Then, since $[Z_j]_{j=1}^{N_1}$ are i.i.d,

$$E\left[\operatorname{Var}\left[\widehat{r}_{M}(x) \mid X\right]\right] = E\left[\operatorname{Var}\left[\frac{N_{0}}{N_{1}M}K_{M}(x) \mid X\right]\right]$$

$$= \left(\frac{N_{0}}{N_{1}M}\right)^{2} E\left[\operatorname{Var}\left[\sum_{j=1}^{N_{1}} \mathbb{1}\left(Z_{j} \in A_{M}(x)\right) \mid X\right]\right]$$

$$= \frac{N_{0}^{2}}{N_{1}M^{2}} E\left[\operatorname{Var}\left[\mathbb{1}\left(Z \in A_{M}(x)\right) \mid X\right]\right]$$

$$= \frac{N_{0}^{2}}{N_{1}M^{2}} E\left[\nu_{1}\left(A_{M}(x)\right) - \nu_{1}^{2}\left(A_{M}(x)\right)\right] \leq \frac{N_{0}^{2}}{N_{1}M^{2}} E\left[\nu_{1}\left(A_{M}(x)\right)\right]$$

$$= \frac{N_{0}}{N_{1}M} E\left[\widehat{r}_{M}(x)\right] \lesssim C\frac{N_{0}}{N_{1}M}, \tag{S3.34}$$

where C > 0 is a constant only depending on f_L , f_U . The last step is due to (S3.31). For the second term in (S3.33), notice that

$$\operatorname{Var}\left[\operatorname{E}\left[\widehat{r}_{M}(x) \mid X\right]\right] = \operatorname{Var}\left[\operatorname{E}\left[\frac{N_{0}}{N_{1}M}K_{M}(x) \mid X\right]\right] = \left(\frac{N_{0}}{M}\right)^{2}\operatorname{Var}\left[\nu_{1}\left(A_{M}(x)\right)\right].$$

Recalling that $W = \nu_0(B_{Z,||x-Z||})$, we have the following lemma about the density of W near 0.

LEMMA S3.3: Denote the density of W by f_W . Then for any $(\nu_0, \nu_1) \in \mathcal{P}_{x,p}(f_L, f_U, L, d, \delta)$,

$$f_W(0) = r(x)$$
.

Furthermore, for any $\epsilon > 0$ and N_0 sufficiently large, we have for all $0 \le w \le 2M/N_0$,

$$\sup_{(\nu_0,\nu_1)\in\mathcal{P}_{x,p}(f_L,f_U,\delta,L,d)} f_W(w) \leq (1+\epsilon) \frac{f_U}{f_L}.$$

Due to Lemma S3.3, we can take N_0 sufficiently large so that for any $0 \le w \le 2M/N_0$,

$$\sup_{(\nu_0,\nu_1)\in\mathcal{P}_{x,p}(f_L,f_U,\delta,L,d)} f_W(w) \le 2\frac{f_U}{f_L}.$$

Let Z, \widetilde{Z} be two independent copies from ν_1 that are further independent of the data. Let $W=\nu_0(B_{Z,\|x-Z\|})$ and $\widetilde{W}=\nu_0(B_{\widetilde{Z},\|x-\widetilde{Z}\|})$. Let $V=\nu_0(B_{Z,\|\mathcal{X}_{(M)}(Z)-Z\|})$ and $\widetilde{V}=\nu_0(B_{\widetilde{Z},\|\mathcal{X}_{(M)}(\widetilde{Z})-\widetilde{Z}\|})$. We then have

$$\operatorname{Var}[\nu_{1}(A_{M}(x))] = \operatorname{E}[\nu_{1}^{2}(A_{M}(x))] - \left(\operatorname{E}[\nu_{1}(A_{M}(x))]\right)^{2}$$

$$= \operatorname{P}(Z \in A_{M}(x), \widetilde{Z} \in A_{M}(x)) - \operatorname{P}(Z \in A_{M}(x))\operatorname{P}(\widetilde{Z} \in A_{M}(x))$$

$$= \operatorname{P}(W \leq V, \widetilde{W} \leq \widetilde{V}) - \operatorname{P}(W \leq V)\operatorname{P}(\widetilde{W} \leq \widetilde{V}).$$

Due to the independence between Z and \widetilde{Z} , W and \widetilde{W} are independent. Notice that $V \mid Z$ have the same distribution as $U_{(M)}$ for any $Z \in \mathbb{R}^d$, then V and Z are independent, so are \widetilde{V} and \widetilde{Z} .

Let us expand the variance further as

$$\operatorname{Var}\left[\nu_{1}\left(A_{M}(x)\right)\right] = \left[P\left(W \leq V, \widetilde{W} \leq \widetilde{V}, W \leq 2\frac{M}{N_{0}}, \widetilde{W} \leq 2\frac{M}{N_{0}}\right) - P\left(W \leq V, W \leq 2\frac{M}{N_{0}}\right)P\left(\widetilde{W} \leq \widetilde{V}, \widetilde{W} \leq 2\frac{M}{N_{0}}\right)\right] + \left[P\left(W \leq V, \widetilde{W} \leq \widetilde{V}\right) - P\left(W \leq V, \widetilde{W} \leq \widetilde{V}, W \leq 2\frac{M}{N_{0}}, \widetilde{W} \leq 2\frac{M}{N_{0}}\right)\right] - \left[P\left(W \leq V\right)P\left(\widetilde{W} \leq \widetilde{V}\right) - P\left(\widetilde{W} \leq V, \widetilde{W} \leq 2\frac{M}{N_{0}}\right)\right] - P\left(\widetilde{W} \leq V, W \leq 2\frac{M}{N_{0}}\right)P\left(\widetilde{W} \leq \widetilde{V}, \widetilde{W} \leq 2\frac{M}{N_{0}}\right)\right]. \tag{S3.35}$$

For the first term in (S3.35), we have the following lemma.

LEMMA S3.4: We have

$$\begin{split} & \left(\frac{N_0}{M}\right)^2 \left[P\left(W \le V, \widetilde{W} \le \widetilde{V}, W \le 2\frac{M}{N_0}, \widetilde{W} \le 2\frac{M}{N_0}\right) \\ & - P\left(W \le V, W \le 2\frac{M}{N_0}\right) P\left(\widetilde{W} \le \widetilde{V}, \widetilde{W} \le 2\frac{M}{N_0}\right) \right] \le C\frac{1}{M}, \end{split}$$

where C > 0 is a constant only depending on f_L , f_U .

For the second term in (S3.35),

$$\begin{split} & P(W \leq V, \widetilde{W} \leq \widetilde{V}) - P\bigg(W \leq V, \widetilde{W} \leq \widetilde{V}, W \leq 2\frac{M}{N_0}, \widetilde{W} \leq 2\frac{M}{N_0}\bigg) \\ & \leq P\bigg(W \leq V, \widetilde{W} \leq \widetilde{V}, W > 2\frac{M}{N_0}\bigg) + P\bigg(W \leq V, \widetilde{W} \leq \widetilde{V}, \widetilde{W} > 2\frac{M}{N_0}\bigg) \end{split}$$

$$\leq \mathbf{P}\bigg(V>2\frac{M}{N_0}\bigg)+\mathbf{P}\bigg(\widetilde{V}>2\frac{M}{N_0}\bigg)=2\mathbf{P}\bigg(U_{(M)}>2\frac{M}{N_0}\bigg).$$

Using the Chernoff bound and $M/\log N_0 \to \infty$, for any $\gamma > 0$,

$$\left(\frac{N_0}{M}\right)^2 \mathbf{P}\left(U_{(M)} > 2\frac{M}{N_0}\right) \le \left(\frac{N_0}{M}\right)^2 \exp\left[-(1-\log 2)M\right] \prec N_0^{-\gamma}.$$

We then have

$$\left(\frac{N_0}{M}\right)^2 \left[P(W \le V, \widetilde{W} \le \widetilde{V}) - P\left(W \le V, \widetilde{W} \le \widetilde{V}, W \le 2\frac{M}{N_0}, \widetilde{W} \le 2\frac{M}{N_0} \right) \right]
\le 2\left(\frac{N_0}{M}\right)^2 P\left(U_{(M)} > 2\frac{M}{N_0}\right) \prec N_0^{-\gamma}.$$
(S3.36)

For the third term in (S3.35), we can check

$$\left\lceil P(W \le V) P(\widetilde{W} \le \widetilde{V}) - P\left(W \le V, W \le 2 \frac{M}{N_0}\right) P\left(\widetilde{W} \le \widetilde{V}, \widetilde{W} \le 2 \frac{M}{N_0}\right) \right\rceil \ge 0.$$

Plugging Lemma S3.4 and (S3.36) into (S3.35) by taking $\gamma > 1$, we obtain

$$\left(\frac{N_0}{M}\right)^2 \operatorname{Var}\left[\nu_1(A_M(x))\right] \lesssim C\frac{1}{M},\tag{S3.37}$$

where C > 0 is a constant only depending on f_L , f_U .

Plugging (\$3.34) and (\$3.37) into (\$3.33) completes the proof. Q.E.D.

S3.6. Proof of Proposition B.1

PROOF OF PROPOSITION B.1: We take ν_0 and ν_1 to share the same support, and assume x to be the origin of \mathbb{R}^d without loss of generality.

When $N_1 \lesssim N_0$, we take ν_0 to be the uniform distribution with density f_L on $[-f_L^{-1/d}/2, f_L^{-1/d}/2]^d$. Then the MSE is lower bounded by the density estimation over Lipchitz class with N_1 samples.

When $N_0 \lesssim N_1$, we take ν_1 to be the uniform distribution with density f_U on $[-f_U^{-1/d}/2, f_U^{-1/d}/2]^d$. Notice that $1/f_0$ is also local Lipchitz from the lower boundness condition and local Lipchitz condition on f_0 . Then the MSE is lower bounded by the density estimation over Lipchitz class with N_0 samples.

We then complete the proof by combining the above two lower bounds and then using the famous minimax lower bound in Lipschitz density estimation (Tsybakov (2009, Exercise 2.8)),

Q.E.D.

S3.7. Proof of Theorem B.4

PROOF OF THEOREM B.4: We only have to prove the first claim as the second is trivial. Take $\delta_N = (\frac{4}{f_L V_d})^{1/d} (\frac{M}{N_0})^{1/d}$ as in the proof of Theorem B.3(i). Take $\delta_N' = (\frac{2}{af_L V_d})^{1/d} \times (\frac{M}{N_0})^{1/d}$. For any $x \in \mathbb{R}^d$, denote the distance of x to the boundary of S_1 by $\Delta(x)$, that is, $\Delta(x) = \inf_{z \in \partial S_1} ||z - x||$.

Depending on the position of x and the value of $\Delta(x)$, we separate the proof into three cases.

Case I. $x \in S_1$ and $\Delta(x) > 2\delta_N$. In this case, since $\Delta(x) > 2\delta_N$, for any $||z - x|| \le \delta_N$, we have $B_{z,||z-x||} \subset S_1$. From the smoothness conditions on f_0 and f_1 , similar to the proof of Theorem B.3, we have

$$E\left[\int_{\mathbb{R}^{d}} \left| \widehat{r}_{M}(x) - r(x) \right| f_{0}(x) \mathbb{1}\left(x \in S_{1}, \Delta(x) > 2\delta_{N}\right) dx \right] \\
\leq \int_{\mathbb{R}^{d}} \left(E\left[\widehat{r}_{M}(x) - r(x)\right]^{2} \right)^{1/2} f_{0}(x) \mathbb{1}\left(x \in S_{1}, \Delta(x) > 2\delta_{N}\right) dx \\
\leq C\left[\left(\frac{M}{N_{0}}\right)^{1/d} + \left(\frac{1}{M}\right)^{1/2} + \left(\frac{N_{0}}{MN_{1}}\right)^{1/2}\right] \int_{\mathbb{R}^{d}} f_{0}(x) \mathbb{1}\left(x \in S_{1}, \Delta(x) > 2\delta_{N}\right) dx \\
\leq C\left[\left(\frac{M}{N_{0}}\right)^{1/d} + \left(\frac{1}{M}\right)^{1/2} + \left(\frac{N_{0}}{MN_{1}}\right)^{1/2}\right], \tag{S3.38}$$

where the constant C > 0 only depends on f_L , f_U , L, d.

Case II. $x \in S_0 \setminus S_1$ and $\Delta(x) > \delta'_N$. In this case, r(x) = 0 and for any $z \in S_1$,

$$u_0(B_{z,\|z-x\|}) \ge f_L \lambda(B_{z,\|z-x\|} \cap S_0) \ge af_L \lambda(B_{z,\|z-x\|}) > af_L V_d \delta_N'^d \ge 2\frac{M}{N_0}.$$

Then for any $\gamma > 0$,

$$\begin{split} \mathbf{E}\big[\big|\widehat{r}_{M}(x) - r(x)\big|\big] &= \mathbf{E}\big[\widehat{r}_{M}(x)\big] = \frac{N_{0}}{M} \mathbf{E}\big[\nu_{1}\big(A_{M}(x)\big)\big] \\ &= \frac{N_{0}}{M} \mathbf{P}\big(W \leq \nu_{0}(B_{Z, \|\mathcal{X}_{(M)}(Z) - Z\|})\big) \leq \frac{N_{0}}{M} \mathbf{P}\Big(U_{(M)} > 2\frac{M}{N_{0}}\Big) \prec N_{0}^{-\gamma}. \end{split}$$

We then obtain

Case III. $x \in S_0$ and $\Delta(x) \leq (2\delta_N) \vee \delta'_N$. In this case, for any $z \in S_1$,

$$\nu_0(B_{z,\|z-x\|}) \ge f_L \lambda(B_{z,\|z-x\|} \cap S_0) \ge af_L \lambda(B_{z,\|z-x\|}) \ge \frac{af_L}{f_U} \nu_1(B_{x,\|z-x\|}).$$

Accordingly,

$$E[|\widehat{r}_{M}(x) - r(x)|] \leq E[\widehat{r}_{M}(x)] + r(x) = \frac{N_{0}}{M} P(W \leq \nu_{0}(B_{Z, \|X_{(M)}(Z) - Z\|})) + r(x)
\leq \frac{N_{0}}{M} P\left(\frac{af_{L}}{f_{U}} \nu_{1}(B_{x, \|x - Z\|}) \leq \nu_{0}(B_{Z, \|X_{(M)}(Z) - Z\|})\right) + r(x)
\leq \frac{N_{0}}{M} P\left(\frac{af_{L}}{f_{U}} U \leq U_{(M)}\right) + r(x) = \frac{f_{U}}{af_{L}} (1 + o(1)) + \frac{f_{U}}{f_{L}}.$$

From the definition of δ_N , δ_N' , and $M/N_0 \to 0$, we have δ_N , $\delta_N' \to 0$ as $N_0 \to \infty$. Since the surface area of S_1 is bounded by H, we have $\lambda(\{x : \Delta(x) \le (2\delta_N) \lor \delta_N'\}) \lesssim H\{(2\delta_N) \lor \delta_N'\}$. Then we obtain

$$E\left[\int_{\mathbb{R}^{d}} \left| \widehat{r}_{M}(x) - r(x) \right| f_{0}(x) \mathbb{1}\left(\Delta(x) \leq (2\delta_{N}) \vee \delta_{N}'\right) dx \right] \\
\leq \left(\frac{f_{U}}{af_{L}} (1 + o(1)) + \frac{f_{U}}{f_{L}}\right) \int_{\mathbb{R}^{d}} f_{0}(x) \mathbb{1}\left(\Delta(x) \leq (2\delta_{N}) \vee \delta_{N}'\right) dx \\
\leq \left(\frac{f_{U}}{af_{L}} (1 + o(1)) + \frac{f_{U}}{f_{L}}\right) f_{U} \lambda \left(\left\{x : \Delta(x) \leq (2\delta_{N}) \vee \delta_{N}'\right\}\right) \\
\lesssim \left(\frac{f_{U}}{af_{L}} + \frac{f_{U}}{f_{L}}\right) f_{U} H\left(\delta_{N} + \delta_{N}'\right) \leq C\left(\frac{M}{N_{0}}\right)^{1/d}, \tag{S3.40}$$

where the constant C > 0 only depends on f_L , f_U , a, H, d. Combining (S3.38), (S3.39), (S3.40) completes the proof. Q.E.D.

S3.8. Proof of Proposition B.2

PROOF OF PROPOSITION B.2: We take ν_0 and ν_1 to be of the same support.

When $N_1 \lesssim N_0$, we take ν_0 to be the uniform distribution with density f_L on $[-f_L^{-1/d}/2, f_L^{-1/d}/2]^d$. Then the L_1 risk is lower bounded by the L_1 risk over support of density estimation over Lipchitz class with N_1 samples.

When $N_0 \lesssim N_1$, we take ν_1 to be the uniform distribution with density f_U on $[-f_U^{-1/d}/2, f_U^{-1/d}/2]^d$. Notice $1/f_0$ is also Lipchitz from the lower boundness condition and Lipchitz condition on f_0 . From the lower boundness condition on f_0 , the L_1 risk is lower bounded by the L_1 risk over support of density estimation over Lipchitz class with N_0 samples.

We then complete the proof by combining the above two lower bounds and then using then the minimax lower bound of L_1 risk for density estimation over Lipchitz class (Zhao and Lai (2022, Theorem 1)). Q.E.D.

S4. PROOFS OF THE RESULTS IN APPENDIX C

S4.1. Proof of Lemma C.1

PROOF OF LEMMA C.1: For any $x \in \mathbb{X}$, define $\sigma_{\omega}^2(x) = \mathrm{E}[U_{\omega}^2 \, | \, X = x] = \mathrm{E}[[Y(\omega) - \mu_{\omega}(X)]^2 \, | \, X = x]$ for $\omega \in \{0, 1\}$. Let

$$V^{\tau} = \mathbb{E}[\mu_1(X) - \mu_0(X) - \tau]^2$$
 and $V^E = \frac{1}{n} \sum_{i=1}^n \left(1 + \frac{K_M(i)}{M}\right)^2 \sigma_{D_i}^2(X_i).$

From the central limit theorem (Billingsley (2008, Theorem 27.1)), we have

$$\sqrt{n}(\bar{\tau}(X) - \tau) \stackrel{d}{\longrightarrow} N(0, V^{\tau}).$$
 (S4.1)

Let $E_{M,i} = (2D_i - 1)(1 + K_M(i)/M)\epsilon_i$ for any $i \in [n]$. Conditional on $X, D, [E_{M,i}]_{i=1}^n$ are independent. Notice that $E[E_{M,i} | X, D] = 0$ and $\sum_{i=1}^n \text{Var}[E_{M,i} | X, D] = nV^E$. To apply the

Lindeberg–Feller central limit theorem (Billingsley (2008, Theorem 27.2)), it suffices to verify that: for a given (X, D),

$$\frac{1}{nV^E}\sum_{i=1}^n\mathrm{E}\big[(E_{M,i})^2\mathbb{I}\big(|E_{M,i}|>\delta\sqrt{nV^E}\big)\,\big|\,X,\boldsymbol{D}\big]\to 0,$$

for all $\delta > 0$.

Let $C_{\sigma} = \sup_{x \in \mathbb{X}, \omega \in \{0,1\}} \{ \mathbb{E}[|U_{\omega}|^{2+\kappa} | X = x] \vee \mathbb{E}[U_{\omega}^{2} | X = x] \} < \infty$. Let $p_{1} = 1 + \kappa/2$ and p_{2} be the constant such that $p_{1}^{-1} + p_{2}^{-1} = 1$. By Hölder's inequality and Markov's inequality,

$$\frac{1}{nV^{E}} \sum_{i=1}^{n} \mathbb{E}[(E_{M,i})^{2} \mathbb{1}(|E_{M,i}| > \delta \sqrt{nV^{E}}) | X, D]$$

$$\leq \frac{1}{nV^{E}} \sum_{i=1}^{n} (\mathbb{E}[|E_{M,i}|^{2+\kappa} | X, D])^{1/p_{1}} (\mathbb{P}(|E_{M,i}| > \delta \sqrt{nV^{E}} | X, D))^{1/p_{2}}$$

$$\leq \frac{1}{nV^{E}} \sum_{i=1}^{n} (\mathbb{E}[|E_{M,i}|^{2+\kappa} | X, D])^{1/p_{1}} \left(\frac{1}{\delta^{2} nV^{E}} \mathbb{E}[(E_{M,i})^{2} | X, D]\right)^{1/p_{2}}$$

$$\leq \frac{C_{\sigma}}{nV^{E}} \left(\frac{1}{\delta^{2} nV^{E}}\right)^{1/p_{2}} \sum_{i=1}^{n} \left(1 + \frac{K_{M}(i)}{M}\right)^{2(1+1/p_{2})}.$$

Notice that $\mathrm{E}[1+K_M(i)/M]^{2(1+1/p_2)}<\infty$ from Theorem B.2. Let $c_\sigma=\inf_{x\in\mathbb{X},\,\omega\in\{0,1\}}\mathrm{E}[U_\omega^2\mid X=x]>0$. From the definition of V^E , we have $V^E\geq c_\sigma$ for almost all X, D. Then

$$E\left[\frac{1}{nV^{E}}\sum_{i=1}^{n}E\left[(E_{M,i})^{2}\mathbb{1}\left(|E_{M,i}|>\delta\sqrt{nV^{E}}\right)|X,\boldsymbol{D}\right]\right]=O(n^{-1/p_{2}})=o(1).$$

We thus obtain

$$\frac{1}{nV^{E}} \sum_{i=1}^{n} \mathrm{E} \big[(E_{M,i})^{2} \mathbb{1} \big(|E_{M,i}| > \delta \sqrt{nV^{E}} \big) \, | \, X, \, \boldsymbol{D} \big] = o_{P}(1).$$

Applying the Lindeberg-Feller central limit theorem then yields

$$\sqrt{n}(V^E)^{-1/2}E_M = (nV^E)^{-1/2}\sum_{i=1}^n E_{M,i} \stackrel{d}{\longrightarrow} N(0,1).$$
 (S4.2)

Noticing that $\sqrt{n}(\bar{\tau}(X) - \tau)$ and $\sqrt{n}(V^E)^{-1/2}E_M$ are asymptotically independent, leveraging the same argument as made in Abadie and Imbens (2006, Proof of Theorem 4, p. 267) and then combining (S4.1) and (S4.2) reaches

$$\sqrt{n}(V^{\tau} + V^{E})^{-1/2}(\tilde{\tau}(X) + E_{M} - \tau) \stackrel{d}{\longrightarrow} N(0, 1). \tag{S4.3}$$

We decompose V^E as

$$V^{E} = \frac{1}{n} \sum_{i=1,D_{i}=1}^{n} \left(1 + \frac{K_{M}(i)}{M} \right)^{2} \sigma_{1}^{2}(X_{i}) + \frac{1}{n} \sum_{i=1,D_{i}=0}^{n} \left(1 + \frac{K_{M}(i)}{M} \right)^{2} \sigma_{0}^{2}(X_{i})$$

$$= \left[\frac{1}{n} \sum_{i=1,D_{i}=1}^{n} \left(\frac{1}{e(X_{i})} \right)^{2} \sigma_{1}^{2}(X_{i}) + \frac{1}{n} \sum_{i=1,D_{i}=0}^{n} \left(\frac{1}{1 - e(X_{i})} \right)^{2} \sigma_{0}^{2}(X_{i}) \right]$$

$$+ \frac{1}{n} \sum_{i=1,D_{i}=1}^{n} \left[\left(1 + \frac{K_{M}(i)}{M} \right)^{2} - \left(\frac{1}{e(X_{i})} \right)^{2} \right] \sigma_{1}^{2}(X_{i})$$

$$+ \frac{1}{n} \sum_{i=1,D_{i}=1}^{n} \left[\left(\frac{1}{1 - e(X_{i})} \right)^{2} - \left(1 + \frac{K_{M}(i)}{M} \right)^{2} \right] \sigma_{0}^{2}(X_{i}). \tag{S4.4}$$

For the first term in (S4.4), notice that $[(X_i, D_i, Y_i)]_{i=1}^n$ are i.i.d. and $E[D_i(e(X_i))^{-2} \times \sigma_1^2(X_i)]$, $E[(1-D_i)(1-e(X_i))^{-2}\sigma_0^2(X_i)] < \infty$. Using the weak law of large numbers, we have

$$\frac{1}{n} \sum_{i=1,D_i=1}^{n} \left(\frac{1}{e(X_i)}\right)^2 \sigma_1^2(X_i) + \frac{1}{n} \sum_{i=1,D_i=0}^{n} \left(\frac{1}{1-e(X_i)}\right)^2 \sigma_0^2(X_i) \stackrel{\mathsf{p}}{\longrightarrow} \mathrm{E}\bigg[\frac{\sigma_1^2(X)}{e(X)} + \frac{\sigma_0^2(X)}{1-e(X)}\bigg].$$

For the second term in (S4.4), using the Cauchy–Schwarz inequality,

$$E \left| \frac{1}{n} \sum_{i=1,D_i=1}^{n} \left[\left(1 + \frac{K_M(i)}{M} \right)^2 - \left(\frac{1}{e(X_i)} \right)^2 \right] \sigma_1^2(X_i) \right| \\
\leq C_\sigma E \left[D_i \left| \left(1 + \frac{K_M(i)}{M} \right)^2 - \left(\frac{1}{e(X_i)} \right)^2 \right| \right] \\
= C_\sigma E \left[D_i E \left[\left| \left(1 + \frac{K_M(i)}{M} \right)^2 - \left(\frac{1}{e(X_i)} \right)^2 \right| | \mathbf{D} \right] \right] \\
\leq C_\sigma E \left[D_i \left(E \left[\left(\frac{K_M(i)}{M} - \frac{1 - e(X_i)}{e(X_i)} \right)^2 \right| \mathbf{D} \right] E \left[\left(2 + \frac{K_M(i)}{M} + \frac{1 - e(X_i)}{e(X_i)} \right)^2 \right| \mathbf{D} \right] \right] \\
= o(1),$$

where the last step is due to Theorem B.2. Then we obtain

$$\frac{1}{n} \sum_{i=1,D:=1}^{n} \left[\left(1 + \frac{K_M(i)}{M} \right)^2 - \left(\frac{1}{e(X_i)} \right)^2 \right] \sigma_1^2(X_i) \stackrel{\mathsf{p}}{\longrightarrow} 0.$$

For the third term in (\$4.4), we can establish in the same way that

$$\frac{1}{n} \sum_{i=1}^{n} \left[\left(\frac{1}{1 - e(X_i)} \right)^2 - \left(1 + \frac{K_M(i)}{M} \right)^2 \right] \sigma_0^2(X_i) \stackrel{\mathsf{p}}{\longrightarrow} 0.$$

Then from (S4.4),

$$V^E \stackrel{\mathsf{p}}{\longrightarrow} \mathrm{E} \bigg[\frac{\sigma_1^2(X)}{e(X)} + \frac{\sigma_0^2(X)}{1 - e(X)} \bigg].$$

By (S4.3), Slutsky's lemma (van der Vaart (1998, Theorem 2.8)), and the definition of σ^2 , we complete the proof. Q.E.D.

S4.2. Proof of Lemma C.2

PROOF OF LEMMA C.2: From Assumption B.1 and Assumption 4.1, let $R = \text{diam}(\mathbb{X}) < \infty$ and $f_L = \inf_{x \in \mathbb{X}, \omega \in \{0,1\}} f_{\omega}(x) > 0$. For any $x \in \mathbb{X}$, $\omega \in \{0,1\}$, and $u \leq R$, from Assumption B.1, $\nu_{\omega}(B_{x,u} \cap \mathbb{X}) \geq f_L \lambda(B_{x,u} \cap \mathbb{X}) \geq f_L a \lambda(B_{x,u}) = f_L a V_d u^d$, where V_d is the Lebesgue measure of the unit ball on \mathbb{R}^d .

Let $c_0 = f_L a V_d$. For any $i \in [n]$, $x \in X$, $M \le n_{1-D_i}$, if $0 \le u \le R n_{1-D_i}^{1/d}$, we have

$$P(\|X_{j} - X_{i}\| \geq u n_{1-D_{i}}^{-1/d} | \mathbf{D}, X_{i} = x, j = j_{M}(i))$$

$$\leq P(Bin(n_{1-D_{i}}, \nu_{1-D_{i}}(B_{x,un_{1-D_{i}}^{-1/d}} \cap \mathbb{X})) \leq M | \mathbf{D})$$

$$\leq P(Bin(n_{1-D_{i}}, c_{0}u^{d}n_{1-D_{i}}^{-1}) \leq M | \mathbf{D}).$$

Using the Chernoff bound, if $M < c_0 u^d$, then

$$P(Bin(n_{1-D_i}, c_0u^d n_{1-D_i}^{-1}) \le M \mid \mathbf{D}) \le \exp\left(M - c_0u^d + M\log\left(\frac{c_0u^d}{M}\right)\right).$$

Notice that the above upper bound does not depend on x. We then obtain

$$P(\|X_{j} - X_{i}\| \ge u n_{1-D_{i}}^{-1/d} | \mathbf{D}, j = j_{M}(i))$$

$$\le \mathbb{1}(M < c_{0}u^{d}) \exp\left(M - c_{0}u^{d} + M \log\left(\frac{c_{0}u^{d}}{M}\right)\right) + \mathbb{1}(M \ge c_{0}u^{d}).$$

On the other hand, if $u > Rn_{1-D_i}^{1/d}$, then the probability is zero from the definition of R. Accordingly, the above bound holds for any $u \ge 0$.

For any $i \in [n]$, we thus have

$$n_{1-D_{i}}^{p/d} \mathbb{E} \left[\| U_{M,i} \|^{p} \, | \, \mathbf{D} \right]$$

$$= p \int_{0}^{\infty} P(\| X_{j} - X_{i} \| \ge u n_{1-D_{i}}^{-1/d} \, | \, \mathbf{D}, j = j_{M}(i)) u^{p-1} \, \mathrm{d}u$$

$$\leq p \int_{0}^{\infty} \left[\mathbb{1} \left(M < c_{0} u^{d} \right) \exp \left(M - c_{0} u^{d} + M \log \left(\frac{c_{0} u^{d}}{M} \right) \right) + \mathbb{1} \left(M \ge c_{0} u^{d} \right) \right] u^{p-1} \, \mathrm{d}u$$

$$= p c_{0}^{-p/d} d^{-1} \left[\int_{M}^{\infty} \left(\frac{e}{M} \right)^{M} t^{M + \frac{p}{d} - 1} e^{-t} \, \mathrm{d}t + \int_{0}^{M} t^{\frac{p}{d} - 1} \, \mathrm{d}t \right], \tag{S4.5}$$

where the last step is through taking $t = c_0 u^d$.

For the first term in (S4.5), from Stirling's formula and $M \to \infty$,

$$\int_{M}^{\infty} \left(\frac{e}{M}\right)^{M} t^{M+\frac{p}{d}-1} e^{-t} dt \leq \int_{0}^{\infty} \left(\frac{e}{M}\right)^{M} t^{M+\frac{p}{d}-1} e^{-t} dt \sim \sqrt{2\pi} M^{\frac{p}{d}-\frac{1}{2}},$$

where \sim means asymptotic convergence. For the second term in (S4.5), $\int_0^M t^{\frac{p}{d}-1} dt = \frac{d}{p} M^{\frac{p}{d}}$. Combining the above two terms then completes the proof. Q.E.D.

S4.3. Proof of Lemma C.3

PROOF OF LEMMA C.3: We bound $B_M - \widehat{B}_M$ by

$$|B_{M} - \widehat{B}_{M}|$$

$$= \left| \frac{1}{n} \sum_{i=1}^{n} (2D_{i} - 1) \left[\frac{1}{M} \sum_{m=1}^{M} (\mu_{1-D_{i}}(X_{i}) - \mu_{1-D_{i}}(X_{j_{m}(i)}) - \widehat{\mu}_{1-D_{i}}(X_{i}) + \widehat{\mu}_{1-D_{i}}(X_{j_{m}(i)}) \right] \right|$$

$$\leq \frac{1}{n} \sum_{i=1}^{n} \max_{m \in \llbracket M \rrbracket} \left| \mu_{1-D_{i}}(X_{i}) - \mu_{1-D_{i}}(X_{j_{m}(i)}) - \widehat{\mu}_{1-D_{i}}(X_{i}) + \widehat{\mu}_{1-D_{i}}(X_{j_{m}(i)}) \right|$$

$$\leq \frac{1}{n} \sum_{i=1}^{n} \max_{m \in \llbracket M \rrbracket, \omega \in \{0,1\}} \left| \mu_{\omega}(X_{i}) - \mu_{\omega}(X_{j_{m}(i)}) - \widehat{\mu}_{\omega}(X_{i}) + \widehat{\mu}_{\omega}(X_{j_{m}(i)}) \right|. \tag{S4.6}$$

Let $k = \lfloor d/2 \rfloor + 1$. For any $\omega \in \{0, 1\}$, by Taylor expansion to kth order,

$$\left| \mu_{\omega}(X_{j_{m}(i)}) - \mu_{\omega}(X_{i}) - \sum_{\ell=1}^{k-1} \frac{1}{\ell!} \sum_{t \in \Lambda_{\ell}} \partial^{t} \mu_{\omega}(X_{i}) U_{m,i}^{t} \right| \leq \max_{t \in \Lambda_{k}} \left\| \partial^{t} \mu_{\omega} \right\|_{\infty} \frac{1}{k!} \sum_{t \in \Lambda_{k}} \|U_{m,i}\|^{k}. \quad (S4.7)$$

In the same way,

$$\left|\widehat{\mu}_{\omega}(X_{j_{m}(i)}) - \widehat{\mu}_{\omega}(X_{i}) - \sum_{\ell=1}^{k-1} \frac{1}{\ell!} \sum_{t \in \Lambda_{\ell}} \partial^{t} \widehat{\mu}_{\omega}(X_{i}) U_{m,i}^{t} \right| \leq \max_{t \in \Lambda_{k}} \left\| \partial^{t} \widehat{\mu}_{\omega} \right\|_{\infty} \frac{1}{k!} \sum_{t \in \Lambda_{k}} \left\| U_{m,i} \right\|^{k}. \quad (S4.8)$$

We also have

$$\left| \sum_{\ell=1}^{k-1} \frac{1}{\ell!} \sum_{t \in \Lambda_{\ell}} \left(\partial^{t} \widehat{\mu}_{\omega}(X_{i}) - \partial^{t} \mu_{\omega}(X_{i}) \right) U_{m,i}^{t} \right| \leq \sum_{\ell=1}^{k-1} \max_{t \in \Lambda_{\ell}} \left\| \partial^{t} \widehat{\mu}_{\omega} - \partial^{t} \mu_{\omega} \right\|_{\infty} \frac{1}{\ell!} \sum_{t \in \Lambda_{\ell}} \|U_{m,i}\|^{\ell}. \quad (S4.9)$$

Notice that $||U_{M,i}|| = \max_{m \in [M]} ||U_{m,i}||$ for any $i \in [n]$, $\omega \in \{0, 1\}$. Then for any $\omega \in \{0, 1\}$, plugging (S4.7), (S4.8), (S4.9) into (S4.6), we obtain

$$egin{aligned} |B_M - \widehat{B}_M| &\lesssim \Big(\max_{\omega \in \{0,1\}} \max_{t \in \Lambda_k} \left\| \partial^t \mu_\omega
ight\|_\infty + \max_{\omega \in \{0,1\}} \max_{t \in \Lambda_k} \left\| \partial^t \widehat{\mu}_\omega
ight\|_\infty \Big) igg(rac{1}{n} \sum_{i=1}^n & \|U_{M,i}\|^k igg) \ + \sum_{\ell=1}^{k-1} igg(\max_{\omega \in \{0,1\}} \max_{t \in \Lambda_\ell} \left\| \partial^t \widehat{\mu}_\omega - \partial^t \mu_\omega
ight\|_\infty \Big) igg(rac{1}{n} \sum_{i=1}^n & \|U_{M,i}\|^\ell igg). \end{aligned}$$

From Lemma C.2, all moments of $(n/M)^{p/d} \|U_{M,i}\|^p$ are bounded. Then for any positive integer p, using Markov's inequality, we have

$$\frac{1}{n}\sum_{i=1}^{n}\left\Vert U_{M,i}\right\Vert ^{p}=O_{\mathbb{P}}\bigg(\bigg(\frac{M}{n}\bigg)^{p/d}\bigg).$$

By Assumption 4.4 and Assumption 4.5, we then obtain

$$\begin{split} B_M - \widehat{B}_M &= O_{\mathbf{P}}(1) O_{\mathbf{P}} \bigg(\bigg(\frac{M}{n} \bigg)^{k/d} \bigg) + \max_{\ell \in \llbracket k-1 \rrbracket} O_{\mathbf{P}} \bigg(n^{-\gamma_\ell} \bigg) O_{\mathbf{P}} \bigg(\bigg(\frac{M}{n} \bigg)^{\ell/d} \bigg) \\ &= O_{\mathbf{P}} \bigg(\bigg(\frac{M}{n} \bigg)^{k/d} \bigg) + \max_{\ell \in \llbracket k-1 \rrbracket} O_{\mathbf{P}} \bigg(n^{-\gamma_\ell} \bigg(\frac{M}{n} \bigg)^{\ell/d} \bigg). \end{split}$$

The proof is thus complete by noticing the definition of γ and $M < n^{\gamma}$. Q.E.D.

S5. PROOFS OF RESULTS IN SUPPLEMENT

S5.1. Proof of Lemma S3.1

PROOF OF LEMMA S3.1: The first inequality is directly from the definition of Lebesgue points. The second inequality follows by

$$\left| \frac{\nu(B_{z,\|z-x\|})}{\lambda(B_{z,\|z-x\|})} - f(x) \right| \leq \frac{1}{\lambda(B_{z,\|z-x\|})} \int_{B_{z,\|z-x\|}} |f(y) - f(x)| \, \mathrm{d}y$$

$$\leq \frac{1}{\lambda(B_{z,\|z-x\|})} \int_{B_{x,2\|z-x\|}} |f(y) - f(x)| \, \mathrm{d}y$$

$$= \frac{\lambda(B_{x,2\|z-x\|})}{\lambda(B_{z,\|z-x\|})} \frac{1}{\lambda(B_{x,2\|z-x\|})} \int_{B_{x,2\|z-x\|}} |f(y) - f(x)| \, \mathrm{d}y$$

$$= 2^{d} \frac{1}{\lambda(B_{x,2\|z-x\|})} \int_{B_{x,2\|z-x\|}} |f(y) - f(x)| \, \mathrm{d}y,$$

and then the definition of Lebesgue points.

O.E.D.

S5.2. Proof of Lemma S3.3

PROOF OF LEMMA S3.3: Fix any $(\nu_0, \nu_1) \in \mathcal{P}_{x,p}(f_L, f_U, L, d, \delta)$.

We first prove the first claim. First, consider $f_1(x) > 0$. For any $\epsilon > 0$, there exists $\delta' > 0$ such that for any $z \in \mathbb{R}^d$ satisfying $\|z - x\| \le 2\delta'$, we have $|f_0(z) - f_0(x)| \le \epsilon f_0(x)$ and $|f_1(z) - f_1(x)| \le \epsilon f_1(x)$ from the local Lipschitz assumption. We take w > 0 sufficiently small such that $w < (1 - \epsilon)f_0(x)\lambda(B_{0,\delta'})$. Then $W \le w$ implies $\|x - Z\| \le \delta'$. Then for w > 0 sufficiently small,

$$\mathbf{P}(W \leq w) = \mathbf{P}(W \leq w, \|x - Z\| \leq \delta') \leq \mathbf{P}\left(\frac{1 - \epsilon}{1 + \epsilon} \frac{f_0(x)}{f_1(x)} \nu_1(B_{x, \|x - Z\|}) \leq w\right) = \frac{1 + \epsilon}{1 - \epsilon} \frac{f_1(x)}{f_0(x)} w,$$

and

$$P(W \le w) = P(W \le w, \|x - Z\| \le \delta') \ge P\left(\frac{1 + \epsilon}{1 - \epsilon} \frac{f_0(x)}{f_1(x)} \nu_1(B_{x, \|x - Z\|}) \le w, \|x - Z\| \le \delta'\right)$$

$$= P\left(\frac{1 + \epsilon}{1 - \epsilon} \frac{f_0(x)}{f_1(x)} \nu_1(B_{x, \|x - Z\|}) \le w\right) = \frac{1 - \epsilon}{1 + \epsilon} \frac{f_1(x)}{f_0(x)} w.$$

Then we have

$$\frac{1-\epsilon}{1+\epsilon}\frac{f_1(x)}{f_0(x)} \leq \liminf_{w\to 0} w^{-1} \mathrm{P}(W \leq w) \leq \limsup_{w\to 0} w^{-1} \mathrm{P}(W \leq w) \leq \frac{1+\epsilon}{1-\epsilon}\frac{f_1(x)}{f_0(x)}.$$

Since ϵ is arbitrary, we obtain

$$f_W(0) = \lim_{w \to 0} w^{-1} P(W \le w) = \frac{f_1(x)}{f_0(x)} = r(x).$$

The case for $f_1(x) = 0$ can be established in the same way. This completes the proof of the first claim.

For the second claim, for any $0 < \epsilon < f_L$, there exists $\delta' > 0$ such that for any $z \in \mathbb{R}^d$ satisfying $||z - x|| \le 2\delta'$, we have $|f_0(z) - f_0(x)| \le \epsilon$ and $|f_1(z) - f_1(x)| \le \epsilon$ from the local Lipschitz assumption. We take N_0 sufficiently large such that $2\frac{M}{N_0} < (f_L - \epsilon)\lambda(B_{0,\delta'})$. Then for any $0 < w \le 2\frac{M}{N_0}$, we have $w < (f_L - \epsilon)\lambda(B_{0,\delta'})$. We take t > 0 such that $w + t < (f_L - \epsilon)\lambda(B_{0,\delta'})$. Then for any $(\nu_0, \nu_1) \in \mathcal{P}_{x,p}(f_L, f_U, L, d, \delta)$,

$$\begin{split} \mathbf{P}(w \leq W \leq w + t) &= \nu_1 \big(\big\{ z \in \mathbb{R}^d : \nu_0(B_{z, \|x - z\|}) \in [w, w + t] \big\} \big) \\ &\leq \frac{f_1(x) + \epsilon}{f_0(x) - \epsilon} \nu_0 \big(\big\{ z \in \mathbb{R}^d : \nu_0(B_{z, \|x - z\|}) \in [w, w + t] \big\} \big). \end{split}$$

Notice that f_0 is lower bounded by f_L . Then for N_0 sufficiently large,

$$\limsup_{t\to 0} t^{-1} \mathrm{P}(w \le W \le w + t) \le \frac{f_1(x) + \epsilon}{f_0(x) - \epsilon} (1 + \epsilon).$$

This then completes the proof.

Q.E.D.

S5.3. Proof of Lemma S3.4

PROOF OF LEMMA S3.4: Due to the i.i.d.-ness of Z and \tilde{Z} ,

$$\left(\frac{N_0}{M}\right)^2 \left[P\left(W \le V, \widetilde{W} \le \widetilde{V}, W \le 2\frac{M}{N_0}, \widetilde{W} \le 2\frac{M}{N_0}\right) - P\left(W \le V, W \le 2\frac{M}{N_0}\right) P\left(\widetilde{W} \le \widetilde{V}, \widetilde{W} \le 2\frac{M}{N_0}\right) \right] \\
= \left(\frac{N_0}{M}\right)^2 \int_0^{2\frac{M}{N_0}} \int_0^{2\frac{M}{N_0}} \left[P(V \ge w_1, \widetilde{V} \ge w_2) - P(V \ge w_1) P(\widetilde{V} \ge w_2) \right] \\
\times f_W(w_1) f_W(w_2) dw_1 dw_2$$

$$\leq 4 \left(\frac{f_{U}}{f_{L}}\right)^{2} \left(\frac{N_{0}}{M}\right)^{2} \int_{0}^{2\frac{M}{N_{0}}} \int_{0}^{2\frac{M}{N_{0}}} \left| P(V \geq w_{1}, \widetilde{V} \geq w_{2}) - P(V \geq w_{1}) P(\widetilde{V} \geq w_{2}) \right| dw_{1} dw_{2}$$

$$= 4 \left(\frac{f_{U}}{f_{L}}\right)^{2} \int_{-1}^{1} \int_{-1}^{1} \left| P\left(V \geq \frac{M}{N_{0}}(1+t_{1}), \widetilde{V} \geq \frac{M}{N_{0}}(1+t_{2})\right) - P\left(V \geq \frac{M}{N_{0}}(1+t_{1})\right) P\left(\widetilde{V} \geq \frac{M}{N_{0}}(1+t_{2})\right) \right| dt_{1} dt_{2},$$

where the last step is from taking $w_1 = \frac{M}{N_0}(1+t_1)$ and $w_2 = \frac{M}{N_0}(1+t_2)$. Let

$$S(t_1, t_2) = \left| P\left(V \ge \frac{M}{N_0} (1 + t_1), \widetilde{V} \ge \frac{M}{N_0} (1 + t_2) \right) - P\left(V \ge \frac{M}{N_0} (1 + t_1) \right) P\left(\widetilde{V} \ge \frac{M}{N_0} (1 + t_2) \right) \right|.$$

If $t_1 \ge t_2 \ge 0$, $S(t_1, t_2) \le P(V \ge \frac{M}{N_0}(1 + t_1)) = P(U_{(M)} \ge \frac{M}{N_0}(1 + t_1))$. If $t_2 \ge t_1 \ge 0$, $S(t_1, t_2) \le P(\widetilde{V} \ge \frac{M}{N_0}(1 + t_2)) = P(U_{(M)} \ge \frac{M}{N_0}(1 + t_2))$. Then for $t_1, t_2 \ge 0$,

$$S(t_1, t_2) \leq P\left(U_{(M)} \geq \frac{M}{N_0}(1 + t_1 \vee t_2)\right).$$

If $t_1 \le t_2 \le 0$ and $P(V \ge \frac{M}{N_0}(1+t_1), \widetilde{V} \ge \frac{M}{N_0}(1+t_2)) \ge P(V \ge \frac{M}{N_0}(1+t_1))P(\widetilde{V} \ge \frac{M}{N_0}(1+t_2))$,

$$S(t_{1}, t_{2}) \leq P\left(\widetilde{V} \geq \frac{M}{N_{0}}(1 + t_{2})\right) - P\left(V \geq \frac{M}{N_{0}}(1 + t_{1})\right)P\left(\widetilde{V} \geq \frac{M}{N_{0}}(1 + t_{2})\right)$$

$$= P\left(V \leq \frac{M}{N_{0}}(1 + t_{1})\right)P\left(\widetilde{V} \geq \frac{M}{N_{0}}(1 + t_{2})\right) \leq P\left(V \leq \frac{M}{N_{0}}(1 + t_{1})\right)$$

$$= P\left(U_{(M)} \leq \frac{M}{N_{0}}(1 + t_{1})\right).$$

If $t_1 \le t_2 \le 0$ and $P(V \ge \frac{M}{N_0}(1+t_1), \widetilde{V} \ge \frac{M}{N_0}(1+t_2)) \le P(V \ge \frac{M}{N_0}(1+t_1))P(\widetilde{V} \ge \frac{M}{N_0}(1+t_2))$,

$$S(t_{1}, t_{2}) \leq P\left(\widetilde{V} \geq \frac{M}{N_{0}}(1 + t_{2})\right) - P\left(V \geq \frac{M}{N_{0}}(1 + t_{1}), \widetilde{V} \geq \frac{M}{N_{0}}(1 + t_{2})\right)$$

$$= P\left(V \leq \frac{M}{N_{0}}(1 + t_{1}), \widetilde{V} \geq \frac{M}{N_{0}}(1 + t_{2})\right) \leq P\left(V \leq \frac{M}{N_{0}}(1 + t_{1})\right)$$

$$= P\left(U_{(M)} \leq \frac{M}{N_{0}}(1 + t_{1})\right).$$

If $t_2 \le t_1 \le 0$, we can establish in the same way that

$$S(t_1, t_2) \le P\left(U_{(M)} \le \frac{M}{N_0}(1 + t_2)\right).$$

Then for $t_1, t_2 \leq 0$,

$$S(t_1, t_2) \leq P\left(U_{(M)} \leq \frac{M}{N_0}(1 + t_1 \wedge t_2)\right).$$

For $t_1 \ge 0 \ge t_2$, if $t_1 + t_2 \ge 0$, $S(t_1, t_2) \le P(U_{(M)} \ge \frac{M}{N_0}(1 + t_1))$, and if $t_1 + t_2 \le 0$, $S(t_1, t_2) \le P(U_{(M)} \le \frac{M}{N_0}(1 + t_2))$. Then

$$\left(\frac{N_{0}}{M}\right)^{2} \left[P\left(W \leq V, \widetilde{W} \leq \widetilde{V}, W \leq 2\frac{M}{N_{0}}, \widetilde{W} \leq 2\frac{M}{N_{0}}\right) - P\left(W \leq V, W \leq 2\frac{M}{N_{0}}\right) P\left(\widetilde{W} \leq \widetilde{V}, \widetilde{W} \leq 2\frac{M}{N_{0}}\right) \right]
\leq 4 \left(\frac{f_{U}}{f_{L}}\right)^{2} \int_{-1}^{1} \int_{-1}^{1} S(t_{1}, t_{2}) dt_{1} dt_{2}
= 4 \left(\frac{f_{U}}{f_{L}}\right)^{2} \left[\int_{0}^{1} \int_{0}^{1} S(t_{1}, t_{2}) dt_{1} dt_{2} + \int_{-1}^{0} \int_{-1}^{0} S(t_{1}, t_{2}) dt_{1} dt_{2} + 2 \int_{0}^{1} \int_{-1}^{0} S(t_{1}, t_{2}) dt_{1} dt_{2} \right],$$
(S5.1)

where the last step is from the symmetry of $S(t_1, t_2)$.

For the first term in (S5.1), by the symmetry of $S(t_1, t_2)$ and the Chernoff bound,

$$\int_{0}^{1} \int_{0}^{1} S(t_{1}, t_{2}) dt_{1} dt_{2}$$

$$\leq \int_{0}^{\infty} \int_{0}^{\infty} S(t_{1}, t_{2}) dt_{1} dt_{2} = 2 \int_{0}^{\infty} \int_{0}^{\infty} S(t_{1}, t_{2}) \mathbb{1}(t_{1} \geq t_{2}) dt_{1} dt_{2}$$

$$\leq 2 \int_{0}^{\infty} \int_{0}^{\infty} P\left(U_{(M)} \geq \frac{M}{N_{0}} (1 + t_{1} \vee t_{2})\right) \mathbb{1}(t_{1} \geq t_{2}) dt_{1} dt_{2}$$

$$= 2 \int_{0}^{\infty} t P\left(U_{(M)} \geq \frac{M}{N_{0}} (1 + t)\right) dt \leq 2 \int_{0}^{\infty} t (1 + t)^{M} e^{-Mt} dt.$$

Notice that since $M \to \infty$, by Stirling's approximation,

$$\int_0^\infty t(1+t)^M e^{-Mt} \, \mathrm{d}t = \frac{1}{M} + \frac{e^M}{M} \int_1^\infty t^M e^{-Mt} \, \mathrm{d}t \le \frac{1}{M} (1+o(1)). \tag{S5.2}$$

We then obtain

$$\int_0^1 \int_0^1 S(t_1, t_2) \, \mathrm{d}t_1 \, \mathrm{d}t_2 \le \frac{2}{M} (1 + o(1)). \tag{S5.3}$$

For the second term in (S5.1),

$$\int_{-1}^{0} \int_{-1}^{0} S(t_{1}, t_{2}) dt_{1} dt_{2} = 2 \int_{-1}^{0} \int_{-1}^{0} S(t_{1}, t_{2}) \mathbb{1}(t_{1} \leq t_{2}) dt_{1} dt_{2}$$

$$\leq 2 \int_{-1}^{0} \int_{-1}^{0} P\left(U_{(M)} \leq \frac{M}{N_{0}}(1 + t_{1} \wedge t_{2})\right) \mathbb{1}(t_{1} \leq t_{2}) dt_{1} dt_{2}$$

$$= 2 \int_{0}^{1} t P\left(U_{(M)} \leq \frac{M}{N_{0}}(1 - t)\right) dt \leq 2 \int_{0}^{1} t (1 - t)^{M} e^{Mt} dt.$$

Notice that

$$\int_0^1 t(1-t)^M e^{Mt} \, \mathrm{d}t \le \frac{1}{M}. \tag{S5.4}$$

We then obtain

$$\int_{-1}^{0} \int_{-1}^{0} S(t_1, t_2) dt_1 dt_2 \le \frac{2}{M}.$$
 (S5.5)

For the third term in (S5.1),

$$\begin{split} & \int_{0}^{1} \int_{-1}^{0} S(t_{1}, t_{2}) dt_{1} dt_{2} \\ & = \int_{0}^{1} \int_{-t_{1}}^{0} P\left(U_{(M)} \ge \frac{M}{N_{0}} (1 + t_{1})\right) dt_{1} dt_{2} + \int_{0}^{1} \int_{-1}^{-t_{1}} P\left(U_{(M)} \le \frac{M}{N_{0}} (1 + t_{2})\right) dt_{1} dt_{2} \\ & = \int_{0}^{1} t P\left(U_{(M)} \ge \frac{M}{N_{0}} (1 + t)\right) dt + \int_{-1}^{0} (-t) P\left(U_{(M)} \le \frac{M}{N_{0}} (1 + t)\right) dt \\ & \le \int_{0}^{\infty} t P\left(U_{(M)} \ge \frac{M}{N_{0}} (1 + t)\right) dt + \int_{-1}^{0} (-t) P\left(U_{(M)} \le \frac{M}{N_{0}} (1 + t)\right) dt \\ & \le \frac{1}{M} (1 + o(1)) + \frac{1}{M} = \frac{2}{M} (1 + o(1)), \end{split}$$

where the last step is from (S5.2) and (S5.4).

We then obtain

$$\int_0^1 \int_{-1}^0 S(t_1, t_2) \, \mathrm{d}t_1 \, \mathrm{d}t_2 \le \frac{2}{M} (1 + o(1)). \tag{S5.6}$$

Plugging (S5.3), (S5.5), (S5.6) into (S5.1) yields

$$\begin{split} & \left(\frac{N_0}{M}\right)^2 \left[P\left(W \leq V, \widetilde{W} \leq \widetilde{V}, W \leq 2 \frac{M}{N_0}, \widetilde{W} \leq 2 \frac{M}{N_0}\right) \\ & - P\left(W \leq V, W \leq 2 \frac{M}{N_0}\right) P\left(\widetilde{W} \leq \widetilde{V}, \widetilde{W} \leq 2 \frac{M}{N_0}\right) \right] \leq 32 \left(\frac{f_U}{f_L}\right)^2 \frac{1}{M} \left(1 + o(1)\right), \end{split}$$

and thus completes the proof.

Q.E.D.

REFERENCES

ABADIE, ALBERTO, AND GUIDO W. IMBENS (2006): "Large Sample Properties of Matching Estimators for Average Treatment Effects," *Econometrica*, 74, 235–267. [22]

BILLINGSLEY, PATRICK (2008): Probability and Measure. John Wiley and Sons. [21,22]

BOGACHEV, VLADIMIR I., AND MARIA A. S. RUAS (2007): Measure Theory, Vol. 1. Springer. [1]

Brown, Russell A. (2015): "Building k-d Tree in O(knlogn) Time," Journal of Computer Graphics Techniques, 4 (1), 50–68. [1]

CORMEN, THOMAS H., CHARLES E. LEISERSON, RONALD L. RIVEST, AND CLIFFORD STEIN (2009): Introduction to Algorithms. MIT press. [1]

FRIEDMAN, JEROME H., JON L. BENTLEY, AND RAPHAEL A. FINKEL (1977): "An Algorithm for Finding Best Matches in Logarithmic Expected Time," *ACM Transactions on Mathematical Software*, 3 (3), 209–226. [1]

ROMANOVSKY, VLADIMIR (1923): "Note on the Moments of a Binomial $(p+q)^n$ About Its Mean," *Biometrika*, 15 (3/4), 410–412. [10]

STEIN, ELIAS M. (2016): Singular Integrals and Differentiability Properties of Functions. Princeton University Press. [11]

TSYBAKOV, ALEXANDRE B. (2009): Introduction to Nonparametric Estimation. Springer. [19]

VAN DER VAART, AAD W. (1998): Asymptotic Statistics. Cambridge University Press. [24]

ZHAO, PUNING, AND LIFENG LAI (2022): "Analysis of KNN Density Estimation," *IEEE Transactions on Information Theory*, 68 (12), 7971–7995. [21]

Co-editor Guido Imbens handled this manuscript.

Manuscript received 22 February, 2022; final version accepted 7 September, 2023; available online 7 September, 2023.