

Inference for Conditional Average Treatment Effects using Distributional Nearest Neighbors

Jakob R. Juergens
University of Wisconsin - Madison

Last edited: December 11, 2024

Abstract

This paper presents a computationally simple method of estimating heterogeneous treatment effects based on the Two-Scale Distributional Nearest Neighbor (TDNN) estimator of Demirkaya et al. (2024). As part of this analysis, I improve on conditions required for consistent variance estimation presented in the original paper and provide results for asymptotically valid pointwise inference in a nonparametric regression setup and extend the analysis to the estimation of conditional average treatment effects. Building on the framework of Ritzwoller and Syrgkanis (2024b), I develop uniformly valid confidence bands for the TDNN estimator. I then show how to apply these to perform uniformly valid inference in both the nonparametric regression setup and the heterogeneous treatment effect setup. A main contribution is the development of a computationally simple method that leverages the theoretical results of the aforementioned papers.

Supplementary Material and R Package available at: https://github.com/JakobJuergens/Unif_Inf_TDNN

1 Introduction

Nearest-Neighbor type estimators and their derivatives are a popular class of estimators that is frequently used in fields such as Computer Science or Economics. The development of inferential theory for these estimators, however, is not yet up to par with their widespread adoption in practice. One such estimator with a particularly close connection to random forests (RF) is the “Two-Scale Distributional Nearest Neighbor Estimator” (TDNN) of Demirkaya et al. (2024). In the aforementioned paper, the authors develop a novel debiasing method that promises great improvements on the finite sample properties of the estimator and show its asymptotic normality. The main contributions of this paper in the current state are twofold. First, this paper provides extended consistency results for variance estimator for the DNN and TDNN estimators. These results show consistency for three Jackknife-based variance estimators in a broad class of asymptotically gaussian generalized U-statistics that extends considerably beyond the DNN-based regression approaches. Second, we propose a novel estimator for the CATE based on the ideas inherent to the DNN estimator and methods from DDML. Simulations show promising performance of the estimator and its simple structure when compared to competing estimators motivates further research into its use for pointwise and simultaneous inference. These extensions will lie at the heart of future iterations of this paper.

After short notation and literature review sections, the remainder of this paper is organized as follows. Section 2 introduces the two setups covered in this paper: first, a relatively simple nonparametric regression setup and second, a setup that mimics the problem of estimating conditional average treatment effects (CATE). Furthermore, this section introduces and contextualizes most of the assumptions that we refer to at later stages of the paper. Section 3 is used to define the DNN and TDNN estimators in the nonparametric regression context and introduces a novel estimator for the CATE setup. Additionally, main results on distributional approximations of the estimators are introduced. Section 4 introduces consistency results for variance estimators to allow for pointwise inference using the DNN estimator and its derivatives. While purely asymptotic in nature, these results improve on currently available results for generalized U-statistics and apply to a broader context than the one presented in this paper. Future iterations of this paper will then go on to tackle the problem of simultaneous inference using techniques developed by Ritzwoller and Syrgkanis (2024a). These novel developments have the potential to significantly extend the applicability of the estimator to scenarios where the treatment effect for a large number of subgroups is of importance. Section 6 contains multiple simulation experiments that show the performance of the methods presented in this paper in a setting mimicking an economic analysis. At this stage, these simulations are limited to the nonparametric regression case, but simulations concerning the estimation of CATE will follow in short time. Lastly, Section 8 concludes.

1.1 Notation

Let $[n] = \{1, \dots, n\}$. Given a finite index set $\mathcal{I} \subset \mathbb{N}$, I introduce the following notational conventions.

$$L_s(\mathcal{I}) = \{(l_1, \dots, l_s) \in \mathcal{I}^s \mid \forall i \neq j: l_i \neq l_j\} \quad \text{and} \quad L_{n,s} = L_s([n]) \quad (1.1)$$

For a data set $\mathbf{D}_{[n]} = (Z_1, \dots, Z_n)$ and a vector $\ell \in L_{n,s}$, denote by $\mathbf{D}_{[n],-\ell}$ the data set where the observations corresponding to indices in ℓ have been removed. To simplify the notation in the case that a single observation (say the i 'th observation) is removed, we use the notation $\mathbf{D}_{n,-i}$. Similarly, given such a data set $\mathbf{D}_{[n]}$ and index vector ℓ , denote by \mathbf{D}_ℓ the data set only consisting of the observations in $\mathbf{D}_{[n]}$ corresponding to indices in ℓ . In an abuse of notation, when considering two index vectors ℓ and ι that do not share any entries, we denote by $\ell \cup \iota$ the concatenation of the two vectors, e.g. if $\ell = (8, 2, 5)$ and $\iota = (1, 6)$, then $\ell \cup \iota = (8, 2, 5, 1, 6)$.

In the following, \rightsquigarrow denotes convergence in distribution, while \rightarrow_p denotes convergence in probability and $\rightarrow_{a.s.}$ denotes almost sure convergence. We will use the symbol \lesssim to denote an inequality that holds for sufficiently large sample sizes n or kernel orders s . As we will consider settings where these diverge together, the specific reference parameter will be clear from the context.

1.2 Related Literature

The related literature can be broadly categorized into three main strands: Nearest-Neighbor type estimators in nonparametric regression, variance estimation for (generalized) U-statistic type estimators including Random Forest, and estimation and inference for CATEs using Double/Debiased Machine Learning (DDML) Methods. A great introduction to the Nearest-Neighbor method is given in Biau and Devroye (2015), illustrating the potential of the method for classification and regression tasks. Of specific interest in the context of this paper are so-called “Weighted Nearest-Neighbor” methods for non-parametric regression. While this is a well-studied type of estimator in and of itself, we draw particular connections to bagged-nearest neighbor type estimators. This class of estimators is built on the framework of “Potential Nearest Neighbors” as introduced by Lin and Jeon (2006). Relevant papers studying their properties are, among others, Biau and Guyader (2010), Biau and Devroye (2010), and Steele (2009). These papers additionally point out the close connections to RF and illustrate why studying the bagged nearest neighbor method could potentially guide our analysis of RF. Recently, Demirkaya et al. (2024) developed a clever debiasing procedure for the bagged or as they coin it distributional nearest neighbor estimator by combining multiple subsampling scales. The resulting TDNN estimator lies at the heart of this paper and the results presented here should be seen in the context of the already established distributional approximations established in the paper.

The concept of U-statistics was introduced by Wassily Hoeffding in Hoeffding (1948) and have been a well-established tool in mathematical statistics for a long time. Thus there is a significant body of literature studying their properties, including outstanding introductions such as Lee (2019). Concerning variance estimation for U-statistics, two highly related papers are Arvesen (1969), exploring the theory of the Jackknife when applied to U-statistics, and Arcones and Gine (1992) which fulfills a similar role for the bootstrap. Building on the concept of U-statistics, Peng, Coleman, and Mentch (2022) introduced the notion of generalized U-statistics, unifying randomized, incomplete, and infinite-order U-statistics that have previously been established in the literature. While being a relatively novel development, there is a significant body of literature concerning infinite-order U-statistics, that share their structure with the TDNN estimator. As the purpose of variance estimation in the problem at hand is ultimately to employ distributional approximations, papers such as Chen and Kato (2019) and Song, Chen, and Kato (2019) are similarly of high relevance for potential applications. Due to the close connection to the random forest method introduced by Breiman (2001), there is also a relevant overlap with the literature on said topic. Thus, papers such as Wager, Hastie, and Efron (2014) and Wager and Athey (2018) are of special interest, especially as causal forest are considered the state of the art technique for estimating CATEs.

In the context of estimation and inference concerning CATEs using DDML, Victor Chernozhukov et al. (2018) should be pointed out first. By combining crossfitting with the use of Neyman-orthogonal moments, the authors built the foundation for many modern methods for estimation in the presence of high-dimensional nuisance parameters. A number of extensions have been proposed to this highly influential idea, some of which are explicitly aimed at estimating CATEs. One example being Semenova and Victor Chernozhukov (2021), who develop estimation and inference procedures for the best linear predictor of a class of causal functions containing the CATE. In a similar vein, Chernozhukov, Newey, and Singh (2022) is highly relevant as it provides a very general analysis of DDML as a meta

algorithm, covering the estimation of and inference for CATE.

2 Setup

Throughout this paper, we will consider two distinct setups. The first is a pure nonparametric regression setup closely mirroring the structure of Demirkaya et al. (2024). This setup will be very useful to illustrate the inner workings of the estimator of interest and serve as a leading example for the theoretical results.

Assumption 1 (Nonparametric Regression DGP).

The observed data consists of an i.i.d. sample taking the following form.

$$\mathbf{D}_n = \{Z_i = (X_i, Y_i)\}_{i=1}^n \quad \text{from the model} \quad Y = \mu(X) + \varepsilon, \quad (2.1)$$

where $Y \in \mathbb{R}$ is the response, $X \in \mathcal{X} \subset \mathbb{R}^k$ is a feature vector of fixed dimension k distributed according to a density function f with associated probability measure φ on \mathcal{X} , and $\mu(x)$ is the unknown mean regression function. ε is the unobservable model error on which we impose the following conditions.

$$\mathbb{E}[\varepsilon | X] = 0, \quad \text{Var}(\varepsilon | X = x) = \sigma_\varepsilon^2(x) \quad (2.2)$$

Let the distribution induced by this model be denoted by P and thus $Z_i = (X_i, Y_i) \stackrel{iid}{\sim} P$.

In contrast to this rather statistical setup, we will consider a setting with more immediate econometric relevance: estimation of and inference on heterogeneous treatment effects in the potential outcomes framework. This serves as a more immediately applicable version of the theoretical setup presented in Ritzwoller and Syrgkanis (2024a) and brings their results closer to practitioners in the field of economics.

Assumption 2 (Heterogeneous Treatment Effect DGP).

The observed data consists of an i.i.d. sample taking the following form.

$$\begin{aligned} \mathbf{D}_n &= \{Z_i = (X_i, W_i, Y_i)\}_{i=1}^n \quad \text{from the model} \quad Y = \mathbb{1}(W = 0)\mu_0(X) + \mathbb{1}(W = 1)\mu_1(X) + \varepsilon, \\ W_i &\sim \text{Bern}(\pi(X_i)) \end{aligned} \quad (2.3)$$

where $Y \in \mathbb{R}$ is the response and $W \in \{0, 1\}$ is an observed treatment indicator. $X \in \mathcal{X} \subset \mathbb{R}^k$ is a vector of covariates of fixed dimension k distributed according to a density function f with associated probability measure φ on \mathcal{X} and ε is the unobservable model error on which we impose the following conditions.

$$\varepsilon \perp\!\!\!\perp W | X, \quad \mathbb{E}[\varepsilon | X] = 0, \quad \text{Var}(\varepsilon | X = x) = \sigma_\varepsilon^2(x) \quad (2.4)$$

Furthermore, $\mu_0 : \mathcal{X} \rightarrow \mathbb{R}$ and $\mu_1 : \mathcal{X} \rightarrow \mathbb{R}$ are the two unknown potential outcome functions and $\pi : \mathcal{X} \rightarrow [0, 1]$ is a function describing the probability of treatment uptake, effectively corresponding to the propensity score. Let the distribution induced by this model be denoted by Q and thus $Z_i = (X_i, W_i, Y_i) \stackrel{iid}{\sim} Q$.

In this second setting, we will use the notation $\mathbf{D}^{(0)}$ and $\mathbf{D}^{(1)}$ to refer to the data subsets containing only observations with $W = 0$ and $W = 1$, respectively. Clearly, this model can be interpreted in the context of the potential outcomes framework in the usual manner.

Remark 1 (Potential Applications).

*From an microeconomic perspective, these two setups cover a wide array of applications. While nonparametric regression is itself often advantageous to answer economic questions, the real strengths show when considering the second setup. **LOREM IPSUM***

Throughout this paper, we will additionally rely on a number of assumptions that are more technical in nature.

Assumption 3 (Technical Assumptions).

In both settings (Assumption 1 and Assumption 2) the following conditions hold:

- *The feature space $\mathcal{X} = \text{supp}(X)$ is a bounded, compact subset of \mathbb{R}^k*
- *The density $f(\cdot)$ is bounded away from 0 and ∞*
- *$f(\cdot)$ and $\mu(\cdot)$ are four times continuously differentiable with bounded second, third, and fourth-order partial derivatives in a neighborhood of x*

In the Heterogeneous Treatment Effect setting (Assumption 2), the following additional condition holds:

- *$\mu_0(\cdot)$ and $\mu_1(\cdot)$ are four times continuously differentiable with bounded second, third, and fourth-order partial derivatives in a neighborhood of x*

There is potential to relax these assumptions at the cost of requiring both less interpretable conditions and more technically sophisticated proofs. Additionally, we require a rather standard assumption in localized regression approaches, namely that the variance changes continuously.

Assumption 4 (Error Distribution Assumptions).

The error terms ε defined in Setup 1 and Setup 2, respectively, have continuously varying variance. In other terms, $\sigma_\varepsilon^2 : \mathcal{X} \rightarrow \mathbb{R}_{>0}$ is a continuous function.

As \mathcal{X} is a bounded and compact set, this implies that there exists a $\bar{\sigma}_\varepsilon^2 > 0$ such that for any $x \in \mathcal{X}$ we have $\sigma_\varepsilon^2(x) \leq \bar{\sigma}_\varepsilon^2$. Readers of Demirkaya et al. (2024) will recognize that this setup, in contrast to the original paper, allows for heteroskedasticity of the error terms. This comes at basically no cost as the original proofs can be used nearly unchanged to prove the corresponding theorems on distributional approximations. Additionally, due to the assumptions on the regression functions, this ensures the existence of second moments of Y in both scenarios. Furthermore, to assure that there is a sufficient number of treated and untreated observations local to each point of interest asymptotically, we require the following condition on the treatment assignment and uptake mechanism.

Assumption 5 (Non-Trivial Treatment Overlap).

In the Heterogeneous Treatment Effect Setup (Assumption 2), we assume that there exist a constant $\mathbf{p} \in (0, 1/2)$ such that

$$\forall x \in \mathcal{X} : \quad 0 < \mathbf{p} \leq \pi(x) \leq 1 - \mathbf{p} < 1. \quad (2.5)$$

This assumption seems rather strong when considering a full universe of potential treatment recipients. In reality we can constrain this overlap assumption to neighborhoods of points of interests x . As long as there is sufficient overlap in those neighborhoods the ideas of our identification strategy continue to hold locally.

Assumption 6 (Stable Unit Treatment Value Assumption (SUTVA)).

For any n , let $\mathfrak{W}_n : \mathcal{X}^n \rightarrow \{0,1\}^n$ and $\mathfrak{W}'_n : \mathcal{X}^n \rightarrow \{0,1\}^n$ be two functions characterizing treatment assignment among a group of n potential observations. Fixing collection of potential observations corresponding to a collection of feature vectors $\mathbf{X} \in \mathcal{X}^n$ for the potential observations and $i \in [n]$, we impose that given $[\mathfrak{W}_n(\mathbf{X})]_i = [\mathfrak{W}'_n(\mathbf{X})]_i$, the following holds.

$$\begin{aligned} Y_i &= \mathbb{1}([\mathfrak{W}_n(\mathbf{X})]_i = 0) \mu_0(\mathbf{X}_i) + \mathbb{1}([\mathfrak{W}_n(\mathbf{X})]_i = 1) \mu_1(\mathbf{X}_i) + \epsilon_i \\ &= \mathbb{1}([\mathfrak{W}'_n(\mathbf{X})]_i = 0) \mu_0(\mathbf{X}_i) + \mathbb{1}([\mathfrak{W}'_n(\mathbf{X})]_i = 1) \mu_1(\mathbf{X}_i) + \epsilon_i = Y'_i \end{aligned} \tag{2.6}$$

Technically, as we are assuming i.i.d. observations in the characterization of the CATE setup, this is already implied. However, due to the importance of the SUTVA assumption in the treatment estimation literature, it seems appropriate to explicitly point out, that it is implicitly assumed that the assumption holds.

3 Two-Scale Distributional Nearest Neighbor Estimator

While less economically enticing, we will introduce the TDNN estimator using the simple nonparametric regression setup first. We will do this by first considering the simpler (one-scale) distributional nearest neighbor estimator, which naturally extends to its two-scale variant as shown in Demirkaya et al. (2024). Then, having established the method, we will commence by adapting it to tackle the problem of estimating heterogeneous treatment effects. As we will embed both estimation problems in the context of subsampled conditional moment regression to then build uniform inference procedures based on Ritzwoller and Syrgkanis (2024a), the approach might at first seem unnatural. However, due to the constructions that follow in Section 5, this approach will be well worth the slightly cumbersome initial presentation.

3.1 DNN and TDNN in Nonparametric Regression

We can rephrase the nonparametric regression problem in terms of estimating specific conditional moments. In the case at hand, this means that our problem can be phrased in the following way.

$$M(x; \mu) = \mathbb{E}[m(Z_i; \mu) | X_i = x] = 0 \quad \text{where} \quad m(Z_i; \mu) = Y_i - \mu(X_i). \quad (3.1)$$

Due to the absence of nuisance parameters, conditions such as local Neyman-orthogonality vacuously hold. We point this out to highlight a contrast that we will encounter when studying the treatment effect setting. In the simpler non-parametric regression setting, we can approach the problem by solving the corresponding empirical conditional moment equation.

$$M_n(x; \mu, \mathbf{D}_n) = \sum_{i=1}^n K(x, X_i) m(Z_i; \mu) = 0 \quad (3.2)$$

In this equation, $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is a data-dependent Kernel function measuring the “distance” between the point of interest and an observation. Notationally, this makes the local and data-dependent approach of this procedure explicit. One estimator that fulfills the purpose of estimating μ nonparametrically is the Distributional Nearest Neighbor (DNN) estimator. With a name coined by Demirkaya et al. (2024), the DNN estimator is based on important work by Steele (2009) and Biau and Guyader (2010). Given a sample as described in Assumption 1 and a fixed feature vector x , we first order the sample based on the distance to the point of interest.

$$\|X_{(1)} - x\|_2 \leq \|X_{(2)} - x\|_2 \leq \dots \leq \|X_{(n)} - x\|_2 \quad (3.3)$$

Here draws are broken according to the natural indices of the observations in a deterministic way to simplify the derivations going forward. While the distance induced by the euclidean norm is a useful tool for developing an intuition for the method, the idea is not inherently connected to it. In fact, any distance induced by a norm that captures the geometry of the feature space in a suitable way can be used to construct an analogous weighting scheme. The generated ordering implies an associated ordering on the response variables and we denote by $Y_{(i)}$ the response corresponding to $X_{(i)}$. Let $\text{rk}(x; X_i, D)$ denote the *rank* that is assigned to observation i in a sample D relative to a point of interest x , setting $\text{rk}(x; X_i, D) = \infty$ if $Z_i \notin D$. Similarly, let $Y_{(1)}(x; D)$ indicate the response value of the closest neighbor in set D . This enables us to define a data-driven kernel function κ following the notation of Ritzwoller and Syrgkanis (2024a).

$$\kappa(x; Z_i, D, \xi) = \mathbb{1}(\text{rk}(x; X_i, D) = 1) \quad (3.4)$$

Here, ξ is an additional source of randomness in the construction of the base learner that comes into play when analyzing, for example, random forests as proposed by Breiman (2001) using the CART-algorithm described in Breiman et al. (2017). As the DNN estimator does not incorporate such additional randomness, the term is omitted in further considerations. In future research, additional randomness such as, for example, column subsampling could be considered, in turn making the addition of ξ necessary again. Using κ , it is straightforward to find an expression for the distance function K in Equation 3.2 corresponding to the DNN estimator.

$$K(x, X_i) = \binom{n}{s}^{-1} \sum_{\ell \in L_{n,s}} \mathbb{1}(i \in \ell) \frac{\kappa(x; Z_i, D_\ell)}{s!} = \binom{n}{s}^{-1} \sum_{\ell \in L_{n,s}} \frac{\mathbb{1}(\text{rk}(x; Z_i, D_\ell) = 1)}{s!} \quad (3.5)$$

Inserting into Equation 3.2, this gives us the following empirical conditional moment equation.

$$M_n(x; \mu, \mathbf{D}_n) = \sum_{i=1}^n \left(\binom{n}{s}^{-1} \sum_{\ell \in L_{n,s}} \frac{\mathbb{1}(\text{rk}(x; Z_i, D_\ell) = 1)}{s!} \right) (Y_i - \mu(X_i)) = 0 \quad (3.6)$$

Solving this empirical conditional moment equation then yields the DNN estimator $\tilde{\mu}_s(x)$ with subsampling scale s . Defining the kernel function, $h_s(x; D_\ell) := (s!)^{-1} Y_{(1)}(x; D_\ell)$, it is given by the following U-statistic.

$$\tilde{\mu}_s(x; \mathbf{D}_n) = \binom{n}{s}^{-1} \sum_{\ell \in L_{n,s}} h_s(x; D_\ell) \quad (3.7)$$

Steele (2009) shows that the DNN estimator has a simple closed form representation based on the original ordered sample.

$$\tilde{\mu}_s(x; \mathbf{D}_n) = \binom{n}{s}^{-1} \sum_{i=1}^{n-s+1} \binom{n-i}{s-1} Y_{(i)} \quad (3.8)$$

This representation will allow me to derive computationally simple representations for the practical use of the procedures presented in this paper. This is in contrast to most U-statistic based methods that inherently rely on evaluating the kernel on individual subsets, incurring a potentially prohibitive computational cost. Furthermore, this representation motivates an asymptotic approximation of the weights assigned to each observation that starkly reduces the potentially computationally intensive computation of large binomial coefficients. For this purpose let $\alpha_s = s/n$ leading to the following approximation of the DNN estimator using asymptotic weights.

$$\tilde{\mu}_s(x; \mathbf{D}_n) \approx \sum_{i=1}^{n-s+1} \alpha_s (1 - \alpha_s)^{i-1} Y_{(i)} \quad (3.9)$$

It is worthwhile to point out that the role of s in the implicit bias-variance tradeoff of the DNN estimator runs counter to the role of k in the usual k-NN regression. Where a larger k is usually associated with a lower variance at the cost of a higher bias, a larger s does the opposite. This is due to the fact that a higher s reduces the number of observations that can occur as the closest observation in any given s -subset. As a special example that illustrates the relationship, consider the DNN estimator choosing $s = n$ recovering the simple 1-NN regression estimator. As part of their paper, Demirkaya et al. (2024) develop an explicit expression for the first-order bias term of the DNN estimator and the following distributional approximation result.

Theorem 3.1 (Demirkaya et al. (2024) - Theorem 2).

Assume that we observe data as described in Assumption 1 and that Assumption 3 holds. Then, for any fixed $x \in \mathcal{X}$, we have that for some positive sequence ω_n of order $\sqrt{s/n}$

$$\frac{\tilde{\mu}_s(x; \mathbf{D}_n) - \mu(x) - B(s) - R(s)}{\omega_n} \rightsquigarrow \mathcal{N}(0, 1) \quad (3.10)$$

as $n, s \rightarrow \infty$ with $s = o(n)$. Here, $B(s)$ and $R(s)$ are defined as the following bias terms.

$$B(s) = \Gamma(2/k + 1) \frac{f(x) \text{tr}(\mu''(x)) + 2\mu'(x)^T f'(x)}{2dV_d^{2/k} f(x)^{1+2/k}} s^{-2/k} \quad \text{and} \quad R(s) = \begin{cases} O(s^{-3}), & k = 1 \\ O(s^{-4/k}), & k \geq 2 \end{cases} \quad (3.11)$$

where...

- $V_d = \frac{k^{k/2}}{\Gamma(1+k/2)}$
- $\Gamma(\cdot)$ is the gamma function
- $\text{tr}(\cdot)$ stands for the trace of a matrix
- $f'(\cdot)$ and $\mu'(\cdot)$ denote the first-order gradients of $f(\cdot)$ and $\mu(\cdot)$, respectively
- $f''(\cdot)$ and $\mu''(\cdot)$ represent the $d \times d$ Hessian matrices of $f(\cdot)$ and $\mu(\cdot)$, respectively

Starting from this setup, Demirkaya et al. (2024) develop a novel bias-correction method for the DNN estimator that leads to appealing finite-sample properties of the resulting Two-Scale Distributional Nearest Neighbor (TDNN) estimator. Their method is based on the explicit formula for the first-order bias term of the DNN estimator, which in turn allows them to eliminate it through a clever combination of two DNN estimators. Choosing two subsampling scales $1 \leq s_1 < s_2 \leq n$ and two corresponding weights

$$w_1^*(s_1, s_2) = \frac{1}{1 - (s_1/s_2)^{-2/k}} \quad \text{and} \quad w_2^*(s_1, s_2) = 1 - w_1^*(s_1, s_2) \quad (3.12)$$

they define the corresponding TDNN estimator as follows.

$$\hat{\mu}_{s_1, s_2}(x; \mathbf{D}_n) = w_1^*(s_1, s_2) \tilde{\mu}_{s_1}(x; \mathbf{D}_n) + w_2^*(s_1, s_2) \tilde{\mu}_{s_2}(x; \mathbf{D}_n) \quad (3.13)$$

This leads to the elimination of the first-order bias term shown in Theorem 3.1 leading to desirable finite-sample properties. Furthermore, the authors show that this construction improves the quality of the normal approximation.

Assumption 7 (Bounded Ratio of Kernel-Orders).

There is a constant $\mathfrak{c} \in (0, 1/2)$ such that the ratio of the kernel orders is bounded in the following way.

$$\forall n : \quad 0 < \mathfrak{c} \leq s_1/s_2 \leq 1 - \mathfrak{c} < 1. \quad (3.14)$$

We make this assumption to avoid edge cases, where asymptotically the TDNN estimator converges to one of the DNN estimators that make it up. As this edge case is irrelevant in practice as it would be simpler to employ the corresponding DNN estimator in the first place, this is not a practically substantial restriction.

Theorem 3.2 (Demirkaya et al. (2024) - Theorem 3).

Assume that we observe data as described in Assumption 1 and that Assumption 3 holds. Furthermore, let $s_1, s_2 \rightarrow \infty$ with $s_1 = o(n)$ and $s_2 = o(n)$ be such that Assumption 7 holds for some $c \in (0, 1/2)$. Then, for any fixed $x \in \text{supp}(X) \subset \mathbb{R}^d$, it holds that for some positive sequence σ_n of order $(s_2/n)^{1/2}$,

$$\sigma_n^{-1} (\hat{\mu}_{s_1, s_2}(x; \mathbf{D}_n) - \mu(x) - \Lambda) \rightsquigarrow \mathcal{N}(0, 1) \quad (3.15)$$

as $n \rightarrow \infty$, where

$$\Lambda = \begin{cases} O(s_1^{-4/d} + s_2^{-4/d}) & \text{for } d \geq 2 \\ O(s_1^{-3} + s_2^{-3}) & \text{for } d = 1 \end{cases}.$$

3.2 DNN and TDNN in Heterogeneous Treatment Effect Estimation

Motivated by the nonparametric regression setup, we set out to apply the underlying idea in the context of heterogeneous treatment effects. Similar to before, we start by specifying a moment corresponding to our object of interest taking into account the additional factors that come into play. Due to the presence of a high-dimensional nuisance parameter in the form of the function q , it is natural to apply the concepts of DDML (DML). This approach closely follows the leading example of Ritzwoller and Syrgkanis (2024a). The main goal at this stage is to construct a highly practical method based on their ideas that leverages the computational simplicity of the distributional nearest neighbor framework.

While considering the problem of point-estimation of a conditional average treatment effect given a feature vector x , $\text{CATE}(x) = \mathbb{E}[Y_i(W_i = 1) - Y_i(W_i = 0) | X_i = x]$, we will employ a Neyman-orthogonal score function to curtail the influence of the nuisance parameters on our estimation.

$$\begin{aligned} M(x; \text{CATE}, \mu, p) &= \mathbb{E}[m(Z_i; \text{CATE}, \mu, \pi) | X_i = x] = 0 \quad \text{where} \\ m(Z_i; \text{CATE}, \mu, \pi) &= \mu_1(X_i) - \mu_0(X_i) + \beta(W_i, X_i)(Y_i - \mu_{W_i}(X_i)) - \text{CATE}(X_i) \end{aligned} \quad (3.16)$$

Here, we make use of the following notation, that is common in the potential outcomes framework, and the well-known Horvitz-Thompson weight.

$$\text{for } w = 1, 2: \quad \mu_w(x) = \mathbb{E}[Y_i | W_i = w, X_i = x] \quad \text{and} \quad \beta(w, x) = \frac{w}{\pi(x)} - \frac{1-w}{1-\pi(x)} \quad (3.17)$$

As a shorthand notation, we will furthermore use $m(Z_i; \mu, \pi) = m(Z_i; \text{CATE}, \mu, \pi) + \text{CATE}(X_i)$. This notation will mainly be used to shorten the presentation of proofs in the appendix. Proceeding in an analogous fashion to the nonparametric regression setup leads us to the following empirical moment equation, where $\hat{\mu}$ and $\hat{\pi}$ are first-stage estimators and K is the data-driven kernel function defined in Equation 3.5.

$$M_n(x; \hat{\mu}, \hat{\pi}) = \sum_{i=1}^n K(x, X_i) m(Z_i; \hat{\mu}, \hat{\pi}) = 0 \quad (3.18)$$

However, due to the presence of infinite-dimensional nuisance parameters, it becomes attractive to proceed by using this weighted empirical moment equation embedded into the DML2 estimator of Victor Chernozhukov et al. (2018).

Applying these ideas to the context of estimating the CATE has been previously explored, for example by Semenova and Victor Chernozhukov (2021). For the sake of simplicity, we will assume that $m = n/K$, i.e. the desired number of observations in each fold, is an integer going forward.

Definition 1. *(T)DNN-DML2 CATE-Estimator*

To estimate the Conditional Average Treatment Effect at a point of interest $x \in \mathcal{X}$, we proceed as follows.

1. Take a K -fold random partition $\mathcal{I} = (I_k)_{k=1}^K$ of the observation indices $[n]$ such that the size of each fold I_k is $m = n/K$. For each $k \in [K]$, define $I_k^C = [n] \setminus I_k$. Furthermore, for the observation being assigned rank $i \in [n]$, denote by $k(i)$ the fold that the observation appears in.

2. For each $k \in [K]$, use the DNN estimator on the data set $\mathbf{D}_{I_k^C} \dots$

(a) to estimate the nuisance parameters μ_0 and μ_1 :

$$\hat{\mu}_{k,s}^w(x) = \hat{\mu}_{w,s}(x; \mathbf{D}_{I_k^C}^{(w)}) \quad \text{for } w = 0, 1 \quad (3.19)$$

(b) if π is unknown, i.e. we are not in a randomized experiment setting, additionally estimate π

$$\hat{\pi}_{k,s}(x) = \hat{\mu}_s(x; \mathbf{D}_{I_k^C}) \quad \text{where the predicted variable is } W \quad (3.20)$$

3. Construct the estimator $\widehat{\text{CATE}}(x)$ as the solution to the following equation.

$$\begin{aligned} 0 &= \sum_{k=1}^K \sum_{i \in I_k} K(x, X_i) m \left(Z_i; \widehat{\text{CATE}}(x), \hat{\mu}_{k,s}, \hat{\pi}_{k,s} \right) \\ &= \sum_{i=1}^{n-s+1} \left[\frac{\binom{n-i}{s-1}}{\binom{n}{s}} \sum_{k=1}^K \mathbb{1}(i \in I_k) m \left(Z_{(i)}; \widehat{\text{CATE}}(x), \hat{\mu}_{k,s}, \hat{\pi}_{k,s} \right) \right] \\ &= \sum_{i=1}^{n-s+1} \left[\frac{\binom{n-i}{s-1}}{\binom{n}{s}} m \left(Z_{(i)}; \widehat{\text{CATE}}(x), \hat{\mu}_{k(i),s}, \hat{\pi}_{k(i),s} \right) \right] \end{aligned} \quad (3.21)$$

This description shows the case of the DNN estimator. Observe, that the weights $K(x, X_i)$ chosen in the second step are chosen according to the full sample - not according to the chosen folds. The corresponding TDNN-based estimator is defined analogously, employing the TDNN estimator in the first-stage estimation procedure and using the corresponding weights of the TDNN-estimator in the second stage. It should be pointed out that the use of the TDNN estimator in the estimation of the propensity score can have the potentially adverse property of generating estimates outside the unit interval. This is due to the presence of negative weights for specific combinations of subsampling scales. Thus, restricting the procedure to rely on the DNN estimator for the estimation of propensity scores in the first stage might be desirable. Specifically, using a lower choice of subsampling scale for this estimation step can help avoid extreme values in the Neyman orthogonal score function due to estimated propensity scores close to zero or one. This is due to the fact that a lower subsampling scale averages over a larger number of observations and can thus contribute to better smoothing properties for the propensity score.

Plugging in for the score function in the equation that defines the estimator, we can observe the following.

$$\widehat{\text{CATE}}(x) = \sum_{i=1}^{n-s+1} \frac{\binom{n-i}{s-1}}{\binom{n}{s}} \left[\widehat{\mu}_{k(i),s}^1(X_{(i)}) - \widehat{\mu}_{k(i),s}^0(X_{(i)}) + \widehat{\beta}_{k(i),s}(W_{(i)}, X_{(i)}) \left(Y_{(i)} - \widehat{\mu}_{k(i),s}^{W_{(i)}}(X_{(i)}) \right) \right] \quad (3.22)$$

Thus, given first-stage estimates of the nuisance parameters, we have a closed form representation of the CATE-estimator for a given partition of $[n]$. Furthermore, given these first-stage estimates, the evaluation of the CATE-estimator at a different point of interest is merely a reweighting of the terms corresponding to different observations. Considering the first stage estimates, we can recognize that the estimation of μ^0 and μ^1 is effectively a nonparametric regression problem as previously described where we used the reduced data sets $\mathbf{D}^{(0)}$ and $\mathbf{D}^{(1)}$, respectively. In contrast, the estimation of π that is necessary in nearly all contexts but randomized experiments can be described further due to the binary outcome. For that purpose, let $Z_{(i|k)}$ denote the i 'th closest observation in fold k akin to the construction shown in Equation 3.3 relative to x but with respect to the data in fold k .

$$\widehat{\pi}_{k,s}(x) = \sum_{i=1}^{n-m-s+1} \frac{\binom{n-m-i}{s-1}}{\binom{n-m}{s}} W_{(i|k)} \quad (3.23)$$

What these equations show is that the main computational cost associated with these methods comes from having to construct multiple orderings of the sample of interest. The essential strength of this approach: It is not necessary to solve any complex optimization problems to obtain the estimator. Furthermore, due to the prevalence of constructing orderings of data with respect to the euclidean norm, this is a well-studied problem with efficient algorithms available of the shelf.

As an extension, we can consider a leave-one-out estimation analog for the functional nuisance parameters, where these are estimated at each observation based on all other available observations. This approach eliminates the randomness inherent to the crossvalidation procedure while preserving the advantages obtained through the usage of DML ideas. This leads to the following estimator.

Definition 2. *(T)DNN-LOO-DML CATE-Estimator*

To estimate the Conditional Average Treatment Effect at a point of interest $x \in \mathcal{X}$, we proceed as follows.

1. For each observation i , use the (T)DNN estimator on the data set $\mathbf{D}_{n,-i} \dots$

(a) to estimate the nuisance parameters μ_0 and μ_1 :

$$\tilde{\mu}_s^w(x) = \widehat{\mu}_{w,s}(x; \mathbf{D}_{n,-i}^{(w)}) \quad \text{for } w = 0, 1 \quad (3.24)$$

(b) if π is unknown, i.e. we are not in a randomized experiment setting, additionally estimate π

$$\tilde{\pi}_s(x) = \widehat{\mu}_s(x; \mathbf{D}_{n,-i}) \quad \text{where the predicted variable is } W \quad (3.25)$$

2. Construct the estimator $\widehat{\text{CATE}}(x)$ as

$$\widehat{\text{CATE}}(x) = \sum_{i=1}^{n-s+1} \frac{\binom{n-i}{s-1}}{\binom{n}{s}} \left[\tilde{\mu}_s^1(X_{(i)}) - \tilde{\mu}_s^0(X_{(i)}) + \tilde{\beta}_s(W_{(i)}, X_{(i)}) \left(Y_{(i)} - \tilde{\mu}_s^{W_{(i)}}(X_{(i)}) \right) \right] \quad (3.26)$$

4 Pointwise Inference for the TDNN Estimator

To perform inference in the regression setup, Demirkaya et al. (2024) introduce variance estimators based on the Jackknife and Bootstrap. However, as they point out, their consistency results rely on a likely suboptimal rate condition for the subsampling scale. While Theorem 3.2 allows s_2 to be of the order $o(n)$, the variance estimators rely on the considerably stronger condition that $s_2 = o(n^{1/3})$. Establishing consistent variance estimation under weaker assumptions on the subsampling rates could broaden the scope of the TDNN estimator for inferential purposes considerably. Furthermore, it can contribute to a better balance between variance and bias as the choice of the kernel orders is crucial when considering the finite sample properties of the estimator. In this paper, we will focus specifically on variance estimators based on the Jackknife and show consistency results under $s = o(n)$. This is motivated by the closed form representation of the estimators in question leading to computationally simple formulas for the exact Jackknife variance estimators.

4.1 Jackknife Variance Estimators for Nonparametric Regression

Define the following variance we need to estimate to perform pointwise inference at a point of interest x .

$$\omega^2(x) = \text{Var}_D(\hat{\mu}_{s_1, s_2}(x; \mathbf{D}_n)) \quad (4.1)$$

We denote by $\mathbf{D}_{n, -i}$ the data set \mathbf{D}_n after removing the i 'th observation. Then, the proposed Jackknife variance estimator takes the following form.

$$\hat{\omega}_{JK}^2(x; \mathbf{D}_n) = \frac{n-1}{n} \sum_{i=1}^n (\hat{\mu}_{s_1, s_2}(x; \mathbf{D}_{n, -i}) - \hat{\mu}_{s_1, s_2}(x; \mathbf{D}_n))^2 \quad (4.2)$$

Theorem 4.1 (Closed Form Expression for the Jackknife-Variance Estimator).

The Jackknife variance estimator for the DNN estimator has the following convenient closed-form representations.

$$\text{LOREMIPSUM} \quad (4.3)$$

Similarly, the Jackknife variance estimator for the TDNN estimator admits the following representation.

$$\text{LOREMIPSUM} \quad (4.4)$$

As a generalization to the Jackknife, we can also consider the delete-d Jackknife that builds on the same working principle. Instead of removing one observation at a time, we remove d observations and average over all possible d -subsets removals. This leads to the following representation.

$$\hat{\omega}_{JKD}^2(x; d, \mathbf{D}_n) = \frac{n-d}{d} \binom{n}{d}^{-1} \sum_{\ell \in L_{n,d}} (\hat{\mu}_{s_1, s_2}(x; \mathbf{D}_{n, -\ell}) - \hat{\mu}_{s_1, s_2}(x; \mathbf{D}_n))^2 \quad (4.5)$$

Similar to the Jackknife, it is possible to derive a closed form representation for the delete-d Jackknife. The derivation would proceed along the exact same lines as in the Jackknife case. However, due to the unwieldiness of the closed form, we refrain from deriving it.

In this section, we will loosen that restrictive condition to make use of the attractive performance of U-statistics with large subsampling rates in the context of inference. The PIJK variance estimator applied to the TDNN estimator is as follows.

$$\hat{\omega}_{PI}^2(x; \mathbf{D}_n) = \frac{s_2^2}{n^2} \sum_{i=1}^n \left[\left(\binom{n-1}{s-1}^{-1} \sum_{\ell \in L_{s_2-1}([n] \setminus \{i\})} h_{s_1, s_2}(x; D_{\ell \cup \{i\}}) \right) - \hat{\mu}_{s_1, s_2}(x; \mathbf{D}_n) \right]^2 \quad (4.6)$$

LOREM IPSUM

Analyzing the kernel of the TDNN estimators, it can be shown that the conditions of Theorem 6 of Peng, Mentch, and Stefanski (2021) apply under the regime $s_2 = o(n)$. Thus, we obtain the following result.

Theorem 4.2 (Pseudo-Infinitesimal Jackknife Variance Estimator Consistency).

Let $0 < \mathfrak{c} \leq s_1/s_2 \leq 1 - \mathfrak{c} < 1$ and $s_2 = o(n)$, then

$$\frac{\hat{\omega}_{PI}^2(x; \mathbf{D}_n)}{\omega^2(x; \mathbf{D}_n)} \xrightarrow{p} 1. \quad (4.7)$$

In an analogous fashion to Theorems 5 and 6 from Demirkaya et al. (2024), we furthermore obtain the following consistency results for the presented variance estimators. As they point out, proving these results goes beyond the techniques presented in Arvesen (1969), instead relying on results for infinite-order U-statistics. Following the ideas from Peng, Mentch, and Stefanski (2021), we then obtain the following results on the Jackknife and Bootstrap variance estimators respectively. As part of the proof of these results, we obtain general results on the consistency of Jackknife and Bootstrap variance estimators for infinite-order U-statistics beyond the TDNN estimator.

Theorem 4.3 (Jackknife Variance Estimator Consistency).

Let $0 < \mathfrak{c} \leq s_1/s_2 \leq 1 - \mathfrak{c} < 1$ and $s_2 = o(n)$, then

$$\frac{\hat{\omega}_{JK}^2(x; \mathbf{D}_n)}{\omega^2(x; \mathbf{D}_n)} \xrightarrow{p} 1. \quad (4.8)$$

Theorem 4.4 (delete-d Jackknife Variance Estimator Consistency).

Let $0 < \mathfrak{c} \leq s_1/s_2 \leq 1 - \mathfrak{c} < 1$, $s_2 = o(n)$, and $d = o(n)$, then

$$\frac{\hat{\omega}_{JKD}^2(x; d, \mathbf{D}_n)}{\omega^2(x; \mathbf{D}_n)} \xrightarrow{p} 1. \quad (4.9)$$

4.2 Variance Estimation for the (T)DNN-DML2 CATE Estimator

Ideas:

- Ignoring the occurrence of left-out observation in nuisance parameter estimation and do basic Jackknife - does this lead to bias?
- Leave Fold-Out Bootstrap with slowly diverging number of folds ($k \rightarrow \infty$, $m = o(n)$) - Effectively a variant of

delete-d bootstrap

- Leave out two folds in the estimator's first step. Then use each previously left out fold for Jackknife construction to eliminate contamination from nuisance parameters
- Modify approach presented in Ritzwoller and Syrgkanis (2024b) Appendix F.4 - modified half-sample k-fold cross-split bootstrap root

A fitting variance estimator given the context of this paper in the literature can be obtained by modifying a construction presented in Ritzwoller and Syrgkanis (2024b). Specifically, the procedure is based on a variation of the approach presented in Appendix F.4 of the aforementioned paper and makes use of a carefully constructed bootstrap-root. Thus, we need to introduce some additional notation, where, for simplicity, we assume that m , i.e. the number of observations in each I_k , is even.

Definition 3 (Crossfitting Half-Sample).

Given a K -fold partition $\mathcal{I} = (I_k)_{k=1}^K$ of $[n]$, a corresponding half sample of \mathcal{I} is a collection of subsets $\mathcal{H} = (H_k)_{k=1}^K$ such that for all $k \in [K]$, the following holds.

$$|H_k| = \frac{|I_k|}{2} = m/2 \quad \text{and} \quad H_k \subset I_k \quad (4.10)$$

The set of all such half-samples of \mathcal{I} is denoted by $\mathfrak{H}(\mathcal{I})$.

This bootstrap root will take the following structure.

$$R_{n,s}^*(x; \mathbf{D}_{[n]}, \mathcal{I}) = \overline{\text{CATE}}_{\mathcal{H}}(x) - \widehat{\text{CATE}}(x) \quad (4.11)$$

Here, $\overline{\text{CATE}}_{\mathcal{H}}(x)$ is the solution to the following equation, where \mathcal{I} is a fixed partition of $[n]$ and \mathcal{H} is a fixed half-sample corresponding to \mathcal{I} . In analogy to the previously established notation, we let $K(x, X_i | \mathcal{H})$ denote the kernel as previously established but with respect to the chosen half-sample and $Z_{(i | \mathcal{H})}$ denote the i 'th closest observation to the point of interest \mathbf{x} that is contained in \mathcal{H} . Furthermore, $k(i | \mathcal{H})$ denotes the fold $k \in [K]$ that the i 'th closest observation in \mathcal{H} is contained in.

$$\begin{aligned} 0 &= \sum_{k=1}^K \sum_{i \in H_k} K(x, X_i | \mathcal{H}) m(Z_i; \overline{\text{CATE}}_{\mathcal{H}}(x), \hat{\mu}_{k,s}, \hat{\pi}_{k,s}) \\ &= \sum_{i=1}^{n/2-s+1} \left[\frac{\binom{n/2-i}{s-1}}{\binom{n/2}{s}} m(Z_{(i | \mathcal{H})}; \overline{\text{CATE}}_{\mathcal{H}}(x), \hat{\mu}_{k(i | \mathcal{H}),s}, \hat{\pi}_{k(i | \mathcal{H}),s}) \right] \end{aligned} \quad (4.12)$$

Plugging in for the moment under consideration once more, we find the following.

$$\begin{aligned} \overline{\text{CATE}}_{\mathcal{H}}(x) &= \sum_{i=1}^{n/2-s+1} \frac{\binom{n/2-i}{s-1}}{\binom{n/2}{s}} \left[\hat{\mu}_{k(i | \mathcal{H}),s}^1(X_{(i | \mathcal{H})}) - \hat{\mu}_{k(i | \mathcal{H}),s}^0(X_{(i | \mathcal{H})}) \right. \\ &\quad \left. + \hat{\beta}_{k(i | \mathcal{H}),s}(W_{(i | \mathcal{H})}, X_{(i | \mathcal{H})})(Y_{(i | \mathcal{H})} - \mu_{W_{(i | \mathcal{H})}}(X_{(i | \mathcal{H})})) \right] \end{aligned} \quad (4.13)$$

Recognizing the similarity to $\widehat{\text{CATE}}(x)$, we can further simplify in the following way.

$$R_{n,s}^*(x; \mathbf{D}_{[n]}, \mathcal{I}) = \text{LOREMIPSUM} \quad (4.14)$$

Theorem 4.5 (Consistent Variance Estimation for the (T)DNN-DML2 CATE Estimator).

LOREM IPSUM

4.3 Pointwise Inference with the TDNN Estimator

Theorem 4.6 (Pointwise Inference in Nonparametric Regression).

LOREM IPSUM

Theorem 4.7 (Pointwise Inference in Heterogeneous Treatment Effect Estimation).

LOREM IPSUM

5 Uniform Inference for the TDNN Estimator

Noteworthy properties of κ are its permutational symmetry in D_ℓ and that κ does not consider the response variable when assigning weights to the observations under consideration. The latter immediately implies a property that has been called “Honesty” by Wager and Athey (2018).

Definition 4 (Symmetry and Honesty - Adapted from Ritzwoller and Syrgkanis (2024a)).

1. The kernel $\kappa(\cdot, \cdot, D_\ell)$ is Honest in the sense that

$$\kappa(x, X_i, D_\ell) \perp\!\!\!\perp m(Z_i; \mu) \mid X_i, D_{\ell, -i},$$

where $\perp\!\!\!\perp$ denotes conditional independence.

2. The kernel $\kappa(\cdot, \cdot, D_\ell)$ is positive and satisfies the restriction $\sum_{i \in s} \kappa(\cdot, X_i, D_\ell) = 1$ almost surely. Moreover, the kernel $\kappa(\cdot, X_i, D_\ell)$ is invariant to permutations of the data D_ℓ .

Absent from Demirkaya et al. (2024) is a way to construct uniformly valid confidence bands around the TDNN estimator. Luckily, as a byproduct of considering the methods from Ritzwoller and Syrgkanis (2024a), procedures for uniform inference can be developed relatively easily.

To consider this problem in detail we first introduce additional notation. Instead of a single point of interest, previously denoted by x , we will consider a vector of p points of interest denoted by $x^{(p)} \in (\text{supp}(X))^p$. Consequently, the j -th entry of $x^{(p)}$ will be denoted by $x_j^{(p)}$. In an abuse of notation, let functions (such as μ or the DNN/TDNN estimators) evaluated at $x^{(p)}$ denote the vector of corresponding function values evaluated at the point, respectively. It should be pointed out that, due to the local definition of the kernel in the estimators, this does not translate to the evaluation of the same function at different points in the most immediate sense. To summarize the kind of object we want to construct, we define a uniform confidence region for the TDNN estimator in the following way following closely the notation of Ritzwoller and Syrgkanis (2024a).

Definition 5 (Uniform Confidence Regions).

A confidence region for the TDNN (or DNN) estimators that is uniformly valid at the rate $r_{n,d}$ is a family of random intervals

$$\widehat{\mathcal{C}}(x^{(p)}) := \left\{ \widehat{\mathcal{C}}(x_j^{(p)}) = [c_L(x_j^{(p)}), c_U(x_j^{(p)})] : j \in [p] \right\} \quad (5.1)$$

based on the observed data, such that

$$\sup_{P \in \mathbf{P}} \left| P \left(\mu(x^{(d)}) \in \widehat{\mathcal{C}}(x^{(d)}) \right) \right| \leq r_{n,d} \quad (5.2)$$

for some sequence $r_{n,d}$, where \mathbf{P} is some statistical family containing P .

5.1 Low-Level

In our pursuit of constructing uniform confidence regions for the TDNN estimator, we return to the results from Ritzwoller and Syrgkanis (2024a) in their high-dimensional form.

Theorem 5.1 (Ritzwoller and Syrgkanis (2024a) - Theorem 4.1).

For any sequence of kernel orders $b = b_n$, where

$$\frac{1}{n} \frac{\nu_j^2}{\sigma_{b,j}^2} \rightarrow 0 \quad \text{as } n \rightarrow \infty, \quad (5.3)$$

we have that

$$\sqrt{\frac{n}{\sigma_{b,j}^2 b^2}} \binom{n}{b}^{-1} \sum_{\mathbf{s} \in \mathbf{S}_{n,b}} u(x_j^{(p)}; D_{\mathbf{s}}) \rightsquigarrow \mathcal{N}(0, 1), \quad \text{as } n \rightarrow \infty. \quad (5.4)$$

Theorem 5.2 (Adapted from Ritzwoller and Syrgkanis (2024a) - Theorem 4.2).

Define the terms

$$\bar{\psi}_{s_2}^2 = \max_{j \in [p]} \{\nu_j^2 - s_2 \sigma_{s_2,j}^2\} \quad \text{and} \quad \underline{\sigma}_{s_2}^2 = \min_{j \in [p]} \sigma_{s_2,j}^2. \quad (5.5)$$

If the kernel function $h_{s_1,s_2}(x; D_\ell)$ satisfies the bound

$$\|h_{s_1,s_2}(x; D_\ell)\|_{\psi_1} \leq \phi \quad (5.6)$$

for each j in $[d]$, then

$$\sqrt{\frac{n}{s_2^2 \underline{\sigma}_{s_2}^2}} \left\| \hat{\mu}_{s_1,s_2}(x^{(p)}; \mathbf{D}_n) - \mu(x^{(p)}) - \frac{s_2}{n} \sum_{i=1}^n h_{s_1,s_2}^{(1)}(x^{(p)}; \mathbf{z}_i) \right\|_\infty = \sqrt{\frac{n}{s_2^2 \underline{\sigma}_{s_2}^2}} \left\| \text{HR}_{s_1,s_2}(x^{(p)}; \mathbf{D}_n) \right\|_\infty \lesssim \xi_{n,s_2}, \quad (5.7)$$

where

$$\xi_{n,s_2} = \left(\frac{C s_2 \log(pn)}{n} \right)^{s_2/2} \left(\left(\frac{n \bar{\psi}_{s_2}^2}{s_2^2 \underline{\sigma}_{s_2}^2} \right)^{1/2} + \left(\frac{\phi^2 s_2 \log^4(pn)}{\underline{\sigma}_{s_2}^2} \right)^{1/2} \right), \quad (5.8)$$

with probability greater than $1 - C/n$.

5.2 High-Level

Recent advances in the field of uniform inference for infinite-order U-statistics, specifically Ritzwoller and Syrgkanis (2024a), and careful analysis of the Hoeffding projections of different orders will be the cornerstones in developing uniform inference methods. The authors' approach to constructing uniform confidence regions is based on the half-sample bootstrap root.

Definition 6 (Half-Sample Bootstrap Root Approximation - Ritzwoller and Syrgkanis (2024a)).

The Half-Sample Bootstrap Root Approximation of the sampling distribution of the root

$$R\left(x^{(p)}; \mathbf{D}_n\right) := \hat{\mu}\left(x^{(p)}; \mathbf{D}_n\right) - \mu(x^{(p)}) \quad (5.9)$$

is given by the conditional distribution of the half-sample bootstrap root

$$R^*\left(x^{(p)}; \mathbf{D}_n\right) := \hat{\mu}\left(x^{(p)}; D_l\right) - \hat{\mu}\left(x^{(p)}; \mathbf{D}_n\right) \quad (5.10)$$

where l denotes a random element from $L_{n,n/2}$.

Next, to standardize the relevant quantities, we introduce a corresponding studentized process.

$$\hat{\lambda}_j^2\left(x^{(p)}; \mathbf{D}_n\right) = \text{Var}\left(\sqrt{n}R^*\left(x_j^{(p)}; \mathbf{D}_n\right) \mid \mathbf{D}_n\right) \quad \text{and} \quad \hat{\Lambda}_n\left(x^{(p)}; \mathbf{D}_n\right) = \text{diag}\left(\left\{\hat{\lambda}_j^2\left(x^{(p)}; \mathbf{D}_n\right)\right\}_{j=1}^p\right) \quad (5.11)$$

$$\hat{S}^*\left(x^{(p)}; \mathbf{D}_n\right) := \sqrt{n} \left\| \left(\hat{\Lambda}_n\left(x^{(p)}; \mathbf{D}_n\right)\right)^{-1/2} R^*\left(x^{(p)}; \mathbf{D}_n\right) \right\|_2 \quad (5.12)$$

Let $\text{cv}(\alpha; \mathbf{D}_n)$ denote the $1 - \alpha$ quantile of the distribution of $\hat{S}^*\left(x^{(p)}; \mathbf{D}_n\right)$. As the authors point out specifically, and as indicated by the more explicit notation chosen in this presentation, this is a quantile of the conditional distribution given the data \mathbf{D}_n . Given this construction, the uniform confidence region developed in Ritzwoller and Syrgkanis (2024a) adapted to the TDNN estimator takes the following form.

Theorem 5.3 (Uniform Confidence Region - Ritzwoller and Syrgkanis (2024a)).

Define the intervals

$$\hat{\mathcal{C}}\left(x_j^{(p)}; \mathbf{D}_n\right) := \hat{\mu}\left(x_j^{(p)}; \mathbf{D}_n\right) \pm n^{-1/2} \hat{\lambda}_j\left(x_j^{(p)}; \mathbf{D}_n\right) \text{cv}(\alpha; \mathbf{D}_n) \quad (5.13)$$

The α -level uniform confidence region for $\mu\left(x^{(p)}\right)$ is given by $\hat{\mathcal{C}}\left(x^{(p)}\right)$.

To justify the use of this uniform confidence region, it remains to be shown if and how the other conditions for the inner workings of this procedure apply to the TDNN estimator. This is substantially simplified due to the absence of a nuisance parameter. Thus, consider the following conditions from Ritzwoller and Syrgkanis, 2024b that are simplified to fit the problem at hand.

Definition 7 (Shrinkage and Incrementality - Adapted from Ritzwoller and Syrgkanis (2024a)).

We say that the kernel $\kappa(\cdot, \cdot, D_\ell)$ has a uniform shrinkage rate ϵ_b if

$$\sup_{P \in \mathbf{P}} \sup_{j \in [p]} \mathbb{E} \left[\max \left\{ \left\| X_i - x_j^{(p)} \right\|_2 : \kappa\left(x_j^{(p)}, X_i, D_\ell\right) > 0 \right\} \right] \leq \epsilon_b. \quad (5.14)$$

We say that a kernel $\kappa(\cdot, \cdot, D_\ell)$ is uniformly incremental if

$$\inf_{P \in \mathbf{P}} \sup_{j \in [p]} \text{Var} \left(\mathbb{E} \left[\sum_{i \in \ell} \kappa\left(x_j^{(p)}, X_i, D_\ell\right) m(Z_i; \mu) \mid l \in \ell, Z_l = Z \right] \right) \gtrsim b^{-1} \quad (5.15)$$

where Z is an independent random variable with distribution P .

Translating these properties to suit the TDNN regression problem, we obtain the following conditions that need to be verified. First, to verify uniform shrinkage at a rate ϵ_b , the following remains to be shown.

$$\sup_{P \in \mathbf{P}} \sup_{j \in [p]} \mathbb{E} \left[\max \left\{ \|X_i - x_j^{(p)}\|_2 : \text{rk}(x_j^{(p)}; X_i, D_\ell) = 1 \right\} \right] \leq \epsilon_b \quad (5.16)$$

Second, for uniform incrementality, we need to show the following.

$$\begin{aligned} & \inf_{P \in \mathbf{P}} \sup_{j \in [p]} \text{Var} \left(\mathbb{E} \left[\sum_{i \in \ell} \mathbb{1}(\text{rk}(x_j^{(p)}; X_i, D_\ell) = 1) (Y_i - \mu(X_i)) \mid l \in \ell, Z_l = Z \right] \right) \\ &= \inf_{P \in \mathbf{P}} \sup_{j \in [p]} \text{Var} \left(\sum_{i \in \ell} \mathbb{E} \left[\mathbb{1}(\text{rk}(x_j^{(p)}; X_i, D_\ell) = 1) \varepsilon_i \mid l \in \ell, Z_l = Z \right] \right) \\ &= \inf_{P \in \mathbf{P}} \sup_{j \in [p]} \text{Var} \left(\sum_{i=1}^s \mathbb{E} \left[\mathbb{1}(\text{rk}(x_j^{(p)}; X_i, D_{1:s}) = 1) \varepsilon_i \mid l \in [s], Z_l = Z \right] \right) \\ &= \inf_{P \in \mathbf{P}} \sup_{j \in [p]} s^2 \cdot \text{Var} \left(\mathbb{E} \left[\mathbb{1}(\text{rk}(x_j^{(p)}; X_1, D_{1:s}) = 1) \varepsilon_1 \mid l \in [s], Z_l = Z \right] \right) \gtrsim b^{-1} \end{aligned} \quad (5.17)$$

To verify these assumptions, recent theory developed in Peng, Coleman, and Mentch (2022) is of great help. Specifically, the following Proposition and its proof are helpful in showing the desired uniform incrementality property.

LOREM IPSUM

Assumption 8 (Boundedness - Adapted from Ritzwoller and Syrgkanis (2024a)).

The absolute value of the function $m(\cdot; \mu)$ is bounded by the constant $(\theta + 1)\phi$ almost surely.

$$|m(Z_i; \mu)| = |Y_i - \mu(X_i)| = |\varepsilon_i| \leq (\theta + 1)\phi \quad a.s. \quad (5.18)$$

To follow the notational conventions, we will further define the two functions $m^{(1)}(Z_i; \mu) = -\mu(X_i)$ and $m^{(2)}(Z_i) = Y_i$. As the authors point out, the boundedness condition can easily be replaced by a condition on the subexponential norm. This, being more in line with the assumptions of Demirkaya et al. (2024), is a desirable substitution. Thus, we will instead consider the following assumption and fill in parts of the proofs that hinge on boundedness for ease of exposition in the original paper.

Assumption 9 (Sub-Exponential Norm Bound).

LOREM IPSUM

Assumption 10 (Moment Smoothness - Adapted from Ritzwoller and Syrgkanis (2024a)).

Define the moments

$$M^{(1)}(x; \mu) = \mathbb{E} \left[m^{(1)}(Z_i; \mu) \mid X_i = x \right] \quad \text{and} \quad M^{(2)}(x) = \mathbb{E} \left[m^{(2)}(Z_i) \mid X_i = x \right], \quad (5.19)$$

associated with the functions $m^{(1)}(\cdot; \mu)$ and $m^{(2)}(\cdot)$. Plugging in yields the following functions.

$$M^{(1)}(x; \mu) = -\mu(x) \quad \text{and} \quad M^{(2)}(x) = \mu(x). \quad (5.20)$$

Both moments are uniformly Lipschitz in their first component, in the sense that

$$\forall x, x' \in \text{supp}(X) : \sup_{P \in \mathbf{P}} |\mu(x) - \mu(x')| \lesssim \|x - x'\|_2. \quad (5.21)$$

and $M^{(1)}$ is bounded below in the following sense

$$\inf_{P \in \mathbf{P}} \inf_{j \in [p]} \left| M^{(1)} \left(x_j^{(p)} \right) \right| = \inf_{P \in \mathbf{P}} \inf_{j \in [p]} \left| \mu \left(x_j^{(p)} \right) \right| \geq c \quad (5.22)$$

for some positive constant c .

The Lipschitz continuity part of this assumption translates directly into a Lipschitz continuity assumption on the unknown nonparametric regression function. The boundedness assumption is **LOREM IPSUM**

5.3 Uniform Inference with the TDNN Estimator

Theorem 5.4 (Uniform Inference in Nonparametric Regression).

LOREM IPSUM

Theorem 5.5 (Uniform Inference in Heterogeneous Treatment Effect Estimation).

LOREM IPSUM

6 Simulations

Having developed theoretical results concerning uniform inference methods for the TDNN estimator, we will proceed by testing their properties in several simulation studies.

6.1 Nonparametric Regression

To investigate the practicality of the nonparametric regression estimators presented in this paper, we consider a collection of setups. First, we focus on illustrating the bias correcting properties of the TDNN estimator by replicating some of the findings of Demirkaya et al. (2024). One such promising example is shown in Figure 1 highlighting the potential improvements obtainable by combining multiple subsampling scales.

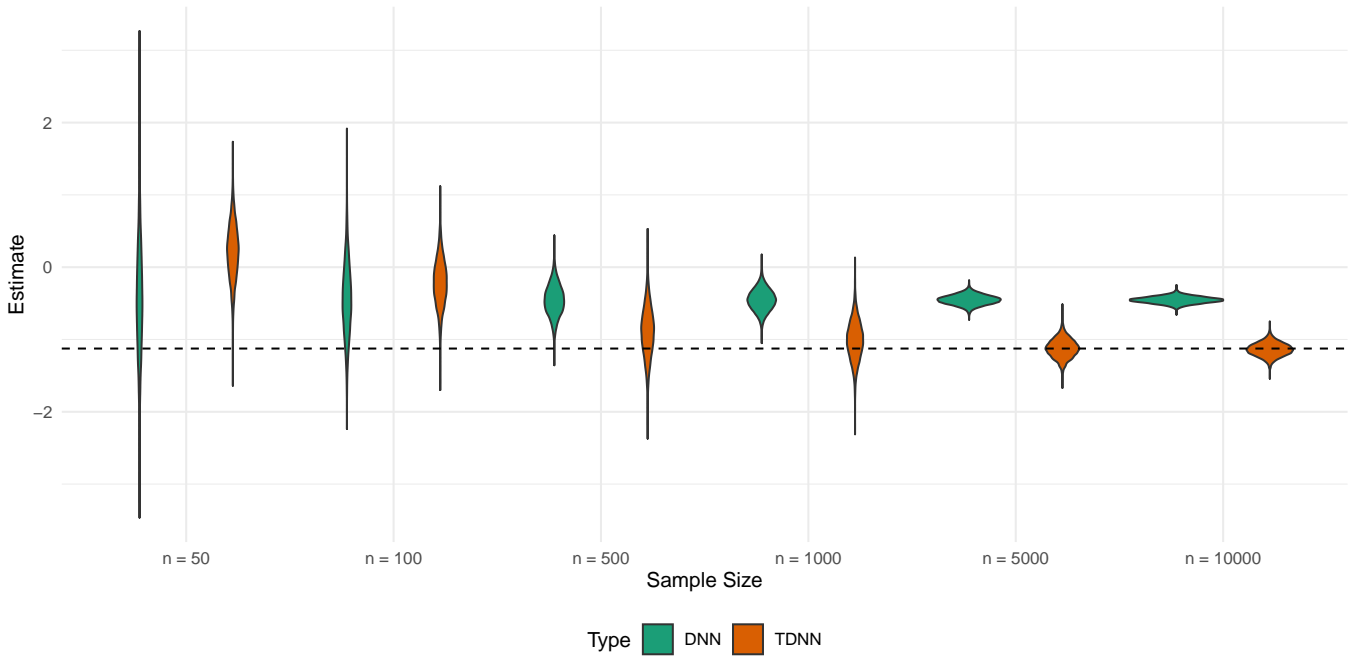


Figure 1: Comparison of the DNN ($s = 20$) and TDNN ($s_1 = 20, s_2 = 50$) Estimators for different sample sizes. The dashed line indicates the value of the unknown regression function at the point of interest. Simulation Setup replicates Setting 1 from Demirkaya et al. (2024) for 10000 Monte Carlo Replications.

As can be seen in Figure 1, a suitable choice of subsampling scales can effectively reduce the bias of the TDNN estimator when compared with the DNN estimator. This reinforces the idea that the TDNN estimator can be a useful tool in practice that has potential to improve on well-established nearest neighbor methods.

As a second, potentially more illustrative example, we consider the estimation of a function of two arguments. Specifically, we consider the function $\mu(x) = 5 \cdot (\cos(x_1) + \cos(x_2))$ on $[0, 1]^2$ with heteroskedastic error terms whose variance is determined by $\sigma_\varepsilon^2(x) = \frac{1}{16} (x_1^2 + x_2^2)^2$. The resulting surface is depicted in Figure 2.

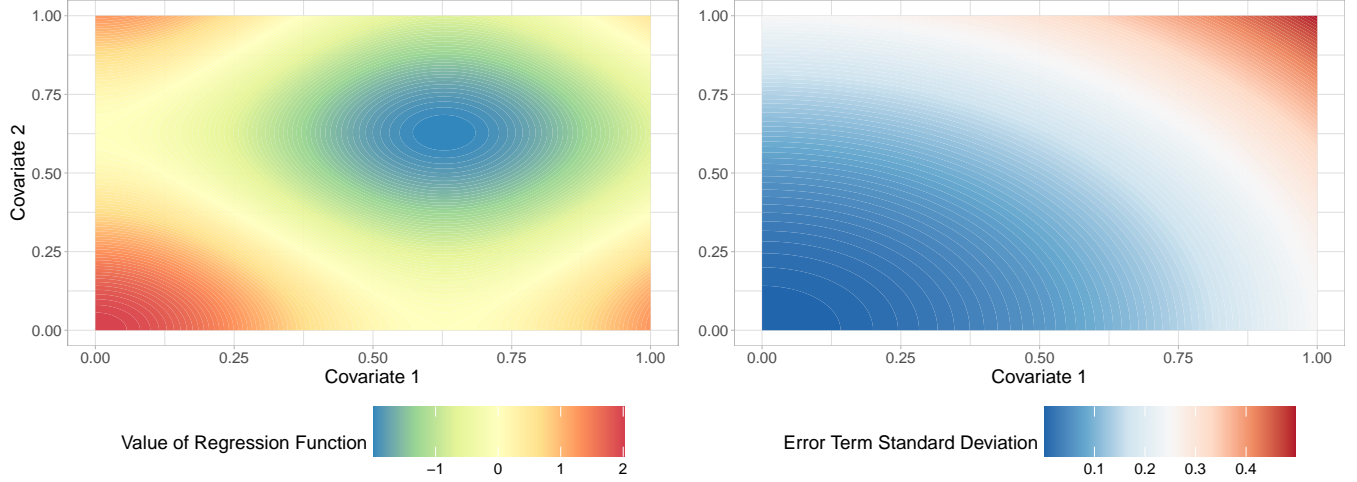


Figure 2: Value of the Regression Function (left) and Variance of the Error Term (right)

To analyze the behavior of the estimator in this setting, we run a number of Monte-Carlo simulations each consisting of 10000 simulation runs. While the theoretical analysis was of purely asymptotic nature, these simulation results can provide a modicum of guidance when it comes to choices such as the kernel orders employed in the estimation procedure. Each run consists of 10000 observations that are uniformly distributed on $[0, 1]^2$, we find the following concerning the estimators bias and variance given different kernel orders s_1 and s_2 . As these diagrams show, the estimator with the given subsampling scales can suffer from bias, specifically in regions that are close to local optima. This is to be expected mechanically as in a local optimum only values of the regression function that lie above or below the optimum, respectively, can occur.

6.2 CATE-Estimation

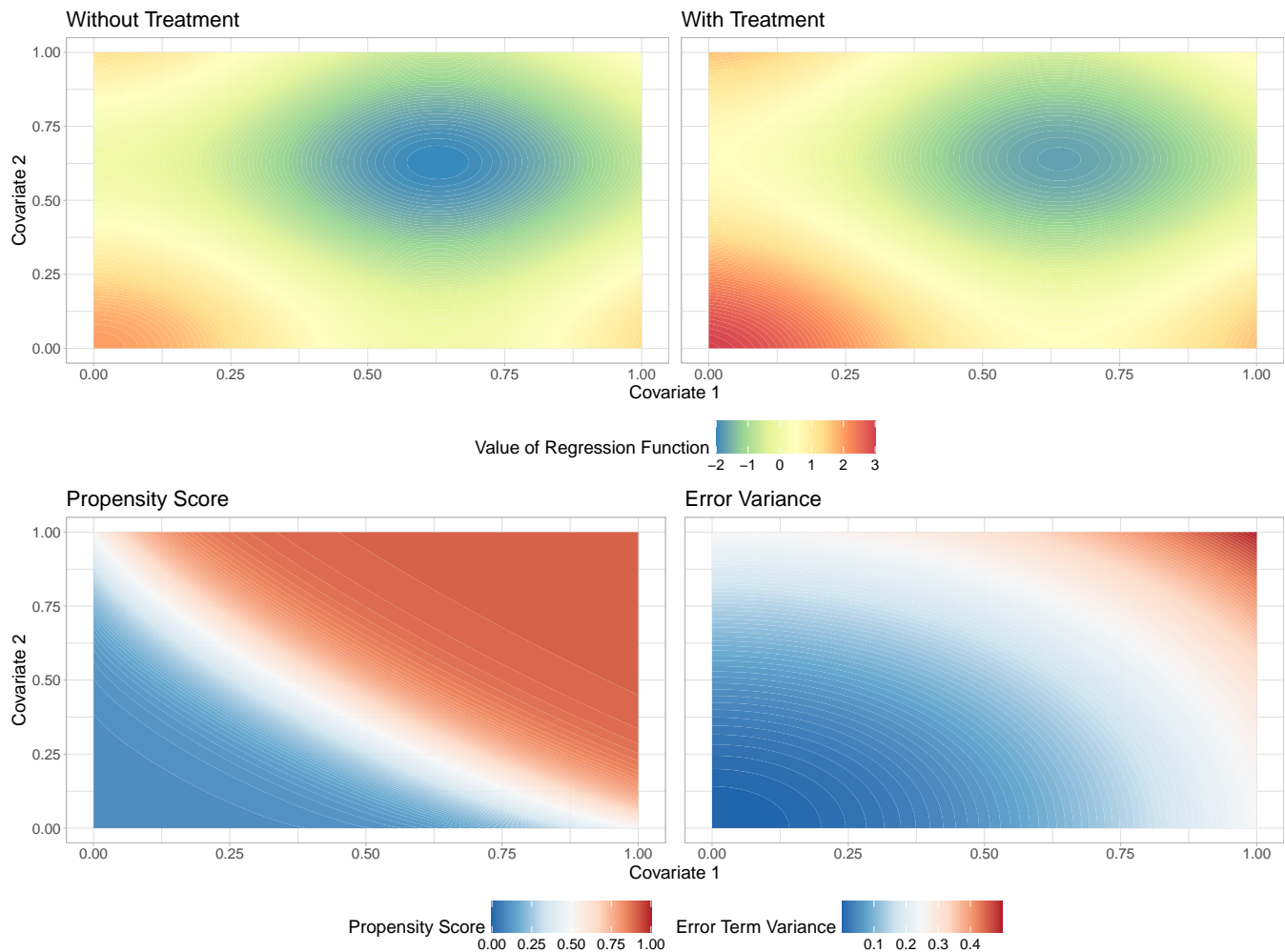


Figure 3: Value of the Regression Functions μ_0 (upper) and μ_1 (lower). Error term structure remains unchanged.

LOREM IPSUM

7 Application

LOREM IPSUM

8 Conclusion

LOREM IPSUM

References

- Arcones, Miguel A. and Evarist Gine (June 1992). “On the Bootstrap of U and V Statistics”. In: *The Annals of Statistics* 20.2, pp. 655–674. DOI: 10.1214/aos/1176348650.
- Arvesen, James N. (Dec. 1969). “Jackknifing U-Statistics”. In: *The Annals of Mathematical Statistics* 40.6, pp. 2076–2100. DOI: 10.1214/aoms/1177697287.
- Biau, Gérard and Luc Devroye (Nov. 2010). “On the layered nearest neighbour estimate, the bagged nearest neighbour estimate and the random forest method in regression and classification”. In: *Journal of Multivariate Analysis* 101.10, pp. 2499–2518. DOI: 10.1016/j.jmva.2010.06.019.
- (2015). *Lectures on the Nearest Neighbor Method*. Springer Series in the Data Sciences. Springer International Publishing. DOI: 10.1007/978-3-319-25388-6.
- Biau, Gérard and Arnaud Guyader (Mar. 2010). “On the Rate of Convergence of the Bagged Nearest Neighbor Estimate”. In: *The Journal of Machine Learning Research* 11, pp. 687–712.
- Breiman, Leo (Oct. 2001). “Random Forests”. In: *Machine Learning* 45.1, pp. 5–32. DOI: 10.1023/A:1010933404324.
- Breiman, Leo et al. (Oct. 2017). *Classification and Regression Trees*. New York: Chapman and Hall/CRC. ISBN: 978-1-315-13947-0. DOI: 10.1201/9781315139470.
- Chen, Xiaohui and Kengo Kato (Dec. 2019). “Randomized incomplete $\$U\$$ -statistics in high dimensions”. In: *The Annals of Statistics* 47.6, pp. 3127–3156. DOI: 10.1214/18-AOS1773.
- Chernozhukov, V, W K Newey, and R Singh (June 2022). “A simple and general debiased machine learning theorem with finite-sample guarantees”. In: *Biometrika* 110.1, pp. 257–264. DOI: 10.1093/biomet/asac033.
- Chernozhukov, Victor et al. (Feb. 2018). “Double/debiased machine learning for treatment and structural parameters”. In: *The Econometrics Journal* 21.1, pp. C1–C68. DOI: 10.1111/ectj.12097.
- Demirkaya, Emre et al. (Jan. 2024). “Optimal Nonparametric Inference with Two-Scale Distributional Nearest Neighbors”. In: *Journal of the American Statistical Association* 119.545, pp. 297–307. DOI: 10.1080/01621459.2022.2115375.
- Hoeffding, Wassily (Sept. 1948). “A Class of Statistics with Asymptotically Normal Distribution”. In: *The Annals of Mathematical Statistics* 19.3, pp. 293–325. DOI: 10.1214/aoms/1177730196.
- Lee, A. J. (Mar. 2019). *U-Statistics: Theory and Practice*. New York: Routledge. ISBN: 978-0-203-73452-0. DOI: 10.1201/9780203734520.
- Lin, Yi and Yongho Jeon (June 2006). “Random Forests and Adaptive Nearest Neighbors”. In: *Journal of the American Statistical Association* 101.474, pp. 578–590. DOI: 10.1198/016214505000001230.
- Peng, Wei, Tim Coleman, and Lucas Mentch (Jan. 2022). “Rates of convergence for random forests via generalized U-statistics”. In: *Electronic Journal of Statistics* 16.1, pp. 232–292. DOI: 10.1214/21-EJS1958.
- Peng, Wei, Lucas Mentch, and Leonard Stefanski (June 2021). *Bias, Consistency, and Alternative Perspectives of the Infinitesimal Jackknife*. arXiv:2106.05918 [math, stat]. DOI: 10.48550/arXiv.2106.05918.
- Ritzwoller, David M. and Vasilis Syrgkanis (Sept. 2024a). *Simultaneous Inference for Local Structural Parameters with Random Forests*. DOI: 10.48550/arXiv.2405.07860.
- (May 2024b). *Uniform Inference for Subsampled Moment Regression*. DOI: 10.48550/arXiv.2405.07860.
- Semenova, Vira and Victor Chernozhukov (June 2021). “Debiased machine learning of conditional average treatment effects and other causal functions”. en. In: *The Econometrics Journal* 24.2, pp. 264–289. DOI: 10.1093/ectj/utaa027.
- Song, Yanglei, Xiaohui Chen, and Kengo Kato (Jan. 2019). “Approximating high-dimensional infinite-order $\$U\$$ -statistics: Statistical and computational guarantees”. In: *Electronic Journal of Statistics* 13.2, pp. 4794–4848. DOI: 10.1214/19-EJS1643.

- Steele, Brian M. (Mar. 2009). “Exact bootstrap k-nearest neighbor learners”. In: *Machine Learning* 74.3, pp. 235–255. DOI: 10.1007/s10994-008-5096-0.
- Wager, Stefan and Susan Athey (July 2018). “Estimation and Inference of Heterogeneous Treatment Effects using Random Forests”. In: *Journal of the American Statistical Association* 113.523. Publisher: Taylor & Francis, pp. 1228–1242. DOI: 10.1080/01621459.2017.1319839.
- Wager, Stefan, Trevor Hastie, and Bradley Efron (2014). “Confidence intervals for random forests: The jackknife and the infinitesimal jackknife”. In: *The journal of machine learning research* 15.1, pp. 1625–1651.

A The (T)DNN Estimator as a Generalized U-Statistic

As the majority of the theoretical results in Demirkaya et al. (2024) rely on representations as a U-statistic, it is helpful to introduce additional concepts and notation at this stage. Recalling Equation 3.7, the DNN and TDNN estimators can be expressed in the following U-statistic form and are thus a type of generalized complete U-statistic as introduced by Peng, Coleman, and Mentch (2022).

$$\tilde{\mu}_s(x; \mathbf{D}_n) = \binom{n}{s}^{-1} \sum_{\ell \in L_{n,s}} h_s(x; D_\ell) \quad \text{and} \quad \hat{\mu}_{s_1, s_2}(x; \mathbf{D}_n) = \binom{n}{s}^{-1} \sum_{\ell \in L_{n, s_2}} h_{s_1, s_2}(x; D_\ell) \quad (\text{A.1})$$

It is worth pointing out that in contrast to the DNN estimator, the kernel for the TDNN estimator is of order $s_2 > s_1$. The authors derive an explicit formula for the kernel that shows the connection between the DNN and TDNN estimators. This connection will prove useful going forward.

Lemma A.1 (Kernel of TDNN Estimator - Adapted from Lemma 8 of Demirkaya et al. (2024)).
The kernel of the TDNN estimator takes the following form.

$$\begin{aligned} h_{s_1, s_2}(x; D) &= w_1^* \left[\binom{s_2}{s_1}^{-1} \sum_{\ell \in L_{s_2, s_1}} h_{s_1}(x; D_\ell) \right] + w_2^* h_{s_2}(x; D) \\ &= w_1^* \tilde{\mu}_{s_1}(x; D) + w_2^* h_{s_2}(x; D) \end{aligned} \quad (\text{A.2})$$

Borrowing the notational conventions from Lee (2019), additionally, introduce the following notation.

$$\psi_s^c(x; \mathbf{z}_1, \dots, \mathbf{z}_c) = \mathbb{E}_D [h_s(x; D) \mid Z_1 = \mathbf{z}_1, \dots, Z_c = \mathbf{z}_c] \quad (\text{A.3})$$

$$h_s^{(1)}(x; \mathbf{z}_1) = \psi_s^1(x; \mathbf{z}_1) - \mu(x) \quad (\text{A.4})$$

$$h_s^{(c)}(x; \mathbf{z}_1, \dots, \mathbf{z}_c) = \psi_s^c(x; \mathbf{z}_1, \dots, \mathbf{z}_c) - \sum_{j=1}^{c-1} \left(\sum_{\ell \in L_{n,j}} h_s^{(j)}(x; \mathbf{z}_\ell) \right) - \mu(x) \quad \text{for } c = 2, \dots, s \quad (\text{A.5})$$

In contrast to the notational inspiration, the subsampling size s is made explicit. Since we are dealing with an infinite-order U-statistic, s will be diverging with n . Completely analogous, define the corresponding objects for the TDNN estimator. For the DNN estimator and any $1 \leq c \leq s$, define

$$\xi_s^c(x) = \text{Var}_{1:c}(\psi_s^c(x; Z_1, \dots, Z_c)) \quad (\text{A.6})$$

where Z'_{c+1}, \dots, Z'_n are i.i.d. from P and independent of Z_1, \dots, Z_n and thus $\xi_s^s(x) = \text{Var}(h_s(x; Z_1, \dots, Z_s))$. Similarly, for the TDNN estimator and any $1 \leq c \leq s_2$, let

$$\zeta_{s_1, s_2}^c(x) = \text{Var}_{1:c}(\psi_{s_1, s_2}^c(x; Z_1, \dots, Z_c)) \quad (\text{A.7})$$

with an analogous definition of Z' .

A.1 Hoeffding-Decomposition

As a byproduct (or main purpose depending on the perspective) these terms can be used to derive the Hoeffding decomposition of the TDNN estimator.

$$H_s^c(x; \mathbf{D}_n) = \binom{n}{c}^{-1} \sum_{\ell \in L_{n,c}} h_s^{(c)}(x; D_\ell) \quad \text{and} \quad H_{s_1, s_2}^c(x; \mathbf{D}_n) = \binom{n}{c}^{-1} \sum_{\ell \in L_{n,c}} h_{s_1, s_2}^{(c)}(x; D_\ell) \quad (\text{A.8})$$

These projection terms can then be used to construct the following Hoeffding decompositions.

$$\tilde{\mu}_s(x; \mathbf{D}_n) = \mu(x) + \sum_{j=1}^s \binom{s}{j} H_s^j(x; \mathbf{D}_n) \quad \text{and} \quad \hat{\mu}_{s_1, s_2}(x; \mathbf{D}_n) = \mu(x) + \sum_{j=1}^{s_2} \binom{s_2}{j} H_{s_1, s_2}^j(x; \mathbf{D}_n) \quad (\text{A.9})$$

Standard results for U-statistics (see for example Lee (2019)) now give us a number of useful results. First, an immediate result on the expectations of the Hoeffding-projection kernels.

$$\forall c = 1, 2, \dots, j-1: \quad \mathbb{E}_D \left[h_{s_1, s_2}^{(j)}(x; D) \mid Z_1 = \mathbf{z}_1, \dots, Z_c = \mathbf{z}_c \right] = 0 \quad \text{and} \quad \mathbb{E}_D \left[h_{s_1, s_2}^{(j)}(x; D) \right] = 0 \quad (\text{A.10})$$

Second, we obtain a useful variance decomposition in terms of the Hoeffding-projection variances.

$$\text{Var}_D(\hat{\mu}_{s_1, s_2}(x; D)) = \sum_{j=1}^{s_2} \binom{s_2}{j}^2 \text{Var}_D(H_{s_1, s_2}^j(x; D)) \quad (\text{A.11})$$

$$\text{Var}_D(H_{s_1, s_2}^j(x; D)) = \binom{n}{j}^{-1} \text{Var}_D(h_{s_1, s_2}^{(j)}(x; D)) =: \binom{n}{j}^{-1} V_{s_1, s_2}^j(x) \quad (\text{A.12})$$

Third, the following equivalent expression for the kernel variance.

$$\zeta_{s_1, s_2}^{s_2}(x) = \text{Var}_D(h_{s_1, s_2}(x; D)) = \sum_{j=1}^{s_2} \binom{s_2}{j} V_{s_1, s_2}^j(x) \quad (\text{A.13})$$

B Useful Results

Lemma B.1 (Demirkaya et al. (2024) - Lemma 12).

Let $D = \{Z_1, \dots, Z_s\}$ an i.i.d. sample drawn from P . The indicator functions $\kappa(x; Z_i, D)$ satisfy the following properties.

1. For any $i \neq j$, we have $\kappa(x; Z_i, D) \kappa(x; Z_j, D) = 0$ with probability one;
2. $\sum_{i=1}^s \kappa(x; Z_i, D) = 1$;
3. $\forall i \in [s] : \mathbb{E}_{1:s} [\kappa(x; Z_i, D)] = s^{-1}$
4. $\mathbb{E}_{2:s} [\kappa(x; Z_1, D)] = \{1 - \varphi(B(x, \|X_1 - x\|))\}^{s-1}$

Here $\mathbb{E}_{i:s}$ denotes the expectation with respect to $\{Z_i, Z_{i+1}, \dots, Z_s\}$. Furthermore, φ denotes the probability measure on \mathbb{R}^d induced by the random vector X .

Lemma B.2 (Demirkaya et al. (2024) - Lemma 13).

For any L^1 function f that is continuous at x , it holds that

$$\lim_{s \rightarrow \infty} \mathbb{E}_1 [f(X_1) s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] = f(x). \quad (\text{B.1})$$

Lemma B.3.

As a consequence of Lemma B.2, we find the following limit results in the nonparametric regression setup.

$$\lim_{s \rightarrow \infty} \mathbb{E}_1 [Y_1 s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] = \mu(x) \quad (\text{B.2})$$

$$\lim_{s \rightarrow \infty} \mathbb{E}_1 [Y_1^2 s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] = \mu^2(x) + \sigma_\varepsilon^2(x) \leq \mu^2(x) + \bar{\sigma}_\varepsilon^2 \quad (\text{B.3})$$

Similarly, in the CATE estimation setup, we can make the following observations.

$$\lim_{s \rightarrow \infty} \mathbb{E}_1 [m(Z_1; \mu, \pi) s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] = \mu_1(x) - \mu_0(x) \quad (\text{B.4})$$

$$\begin{aligned} \lim_{s \rightarrow \infty} \mathbb{E}_1 [m^2(Z_1; \mu, \pi) s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] &= (\mu_1(x) - \mu_0(x))^2 + \frac{\sigma_\varepsilon^2(x)}{\pi(x)(1 - \pi(x))} \\ &\leq (\mu_1(x) - \mu_0(x))^2 + \frac{\bar{\sigma}_\varepsilon^2}{\mathfrak{p}(1 - \mathfrak{p})} \end{aligned} \quad (\text{B.5})$$

Proof of Lemma B.3. Starting with the first limit, we find the following.

$$\begin{aligned} \mathbb{E}_1 [Y_1 s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] &= \mathbb{E}_1 [(\mu(X_1) + \varepsilon_1) s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] \\ &= \mathbb{E}_1 [(\mu(X_1) + \mathbb{E}[\varepsilon_1 | X_1]) s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] \\ &= \mathbb{E}_1 [\mu(X_1) s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] \xrightarrow{(\text{Lem B.2})} \mu(x) \quad \text{as } s \rightarrow \infty \end{aligned} \quad (\text{B.6})$$

Similarly, when considering the second limit, we can make the following observation.

$$\begin{aligned}
\mathbb{E}_1 [Y_1^2 s\mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] &= \mathbb{E}_1 [(\mu(X_1) + \varepsilon_1)^2 s\mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] \\
&= \mathbb{E}_1 [(\mu^2(X_1) + 2\mu(X_1)\varepsilon_1 + \varepsilon_1^2) s\mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] \\
&= \mathbb{E}_1 [(\mu^2(X_1) + 2\mu(X_1)\mathbb{E}[\varepsilon_1 | X_1] + \mathbb{E}[\varepsilon_1^2 | X_1]) s\mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] \\
&= \mathbb{E}_1 [(\mu^2(X_1) + \sigma_\varepsilon^2(X_1)) s\mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] \\
&\xrightarrow{(\text{Lem B.2})} \mu^2(x) + \sigma_\varepsilon^2(x) \quad \text{as } s \rightarrow \infty
\end{aligned} \tag{B.7}$$

In the CATE estimation setting, we can proceed analogously.

$$\begin{aligned}
\mathbb{E}_1 [m(Z_1; \mu, \pi) s\mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] &= \mathbb{E}_1 [(\mu_1(X_1) - \mu_0(X_1) + \beta(W_1, X_1)\varepsilon_i) s\mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] \\
&= \mathbb{E}_1 [(\mu_1(X_1) - \mu_0(X_1) + \beta(W_1, X_1)\mathbb{E}[\varepsilon_i | X_1]) s\mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] \\
&= \mathbb{E}_1 [(\mu_1(X_1) - \mu_0(X_1)) s\mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] \\
&\xrightarrow{(\text{Lem B.2})} \mu_1(x) - \mu_0(x) \quad \text{as } s \rightarrow \infty
\end{aligned} \tag{B.8}$$

Similarly, we can find the following.

$$\begin{aligned}
\mathbb{E}_1 [m^2(Z_i; \mu, \pi) s\mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] &= \mathbb{E}_1 [(\mu_1(X_1) - \mu_0(X_1) + \beta(W_1, X_1)\varepsilon_i)^2 s\mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] \\
&= \mathbb{E}_1 [(\mu_1(X_i) - \mu_0(X_1))^2 s\mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] + \underbrace{\mathbb{E}_1 [(\mu_1(X_i) - \mu_0(X_1))\beta(W_1, X_1)\mathbb{E}[\varepsilon_i | X_1] s\mathbb{E}_{2:s} [\kappa(x; Z_1, D)]]}_{=0} \\
&\quad + \mathbb{E}_1 [(\beta(W_1, X_1)\varepsilon_i)^2 s\mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] \\
&= \mathbb{E}_1 [(\mu_1(X_1) - \mu_0(X_1))^2 s\mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] + \mathbb{E}_1 \left[\left(\frac{W_1}{\pi(X_1)} - \frac{1-W_1}{1-\pi(X_1)} \right)^2 \varepsilon_1^2 s\mathbb{E}_{2:s} [\kappa(x; Z_1, D)] \right] \\
&= \underbrace{\mathbb{E}_1 [(\mu_1(X_1) - \mu_0(X_1))^2 s\mathbb{E}_{2:s} [\kappa(x; Z_1, D)]]}_{\rightarrow (\mu_1(x) - \mu_0(x))^2 \text{ as } s \rightarrow \infty} + \underbrace{\mathbb{E}_1 \left[\mathbb{E} \left[\left(\frac{W_1}{\pi(X_1)} - \frac{1-W_1}{1-\pi(X_1)} \right)^2 \varepsilon_1^2 \middle| X_1 \right] s\mathbb{E}_{2:s} [\kappa(x; Z_1, D)] \right]}_{(B)}
\end{aligned} \tag{B.9}$$

Continuing with the second term, marked by (B), we find the following.

$$\begin{aligned}
(B) &= \mathbb{E}_1 \left[\mathbb{E} \left[\left(\frac{W_1}{\pi(X_1)} - \frac{1-W_1}{1-\pi(X_1)} \right)^2 \varepsilon_1^2 \middle| X_1 \right] s\mathbb{E}_{2:s} [\kappa(x; Z_1, D)] \right] \\
&= \mathbb{E}_1 \left[\frac{\sigma_\varepsilon^2(X_1) \cdot s\mathbb{E}_{2:s} [\kappa(x; Z_1, D)]}{\pi^2(X_1)(1-\pi(X_1))^2} \cdot \mathbb{E} \left[(W_1(1-\pi(X_1)) - (1-W_1)\pi(X_1))^2 \middle| X_1 \right] \right]
\end{aligned} \tag{B.10}$$

Observe that $W_1(1-W_1) = 0$, $W_1^2 = W_1$, and $(1-W_1)^2 = 1-W_1$, which allows us to use the following simplification.

$$\begin{aligned}
(B) &= \mathbb{E}_1 \left[\frac{\sigma_\varepsilon^2(X_1) \cdot s\mathbb{E}_{2:s} [\kappa(x; Z_1, D)]}{\pi^2(X_1)(1-\pi(X_1))^2} \cdot \mathbb{E} \left[W_1(1-\pi(X_1))^2 + (1-W_1)\pi^2(X_1) \middle| X_1 \right] \right] \\
&= \mathbb{E}_1 \left[\frac{\sigma_\varepsilon^2(X_1) \cdot s\mathbb{E}_{2:s} [\kappa(x; Z_1, D)]}{\pi^2(X_1)(1-\pi(X_1))^2} \cdot \left(\pi(X_1)(1-\pi(X_1))^2 + (1-\pi(X_1))\pi^2(X_1) \right) \right] \\
&= \mathbb{E}_1 \left[\frac{\sigma_\varepsilon^2(X_1) \cdot s\mathbb{E}_{2:s} [\kappa(x; Z_1, D)]}{\pi(X_1)(1-\pi(X_1))} \right] \xrightarrow{(\text{Lem B.2})} \frac{\sigma_\varepsilon^2(x)}{\pi(x)(1-\pi(x))} \quad \text{as } s \rightarrow \infty
\end{aligned} \tag{B.11}$$

Recombining the terms of interest, we find the desired limit bound.

$$\mathbb{E}_1 \left[m^2(Z_i; \mu, \pi) s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)] \right] \xrightarrow{(\text{Lem B.2})} (\mu_1(x) - \mu_0(x))^2 + \frac{\sigma_\varepsilon^2(x)}{\pi(x)(1 - \pi(x))} \quad \text{as } s \rightarrow \infty \quad (\text{B.12})$$

■

Lemma B.4.

Fix sample size n , subsampling scale s , and c such that $0 < c \leq s \leq n$. Let $D = \{Z_1, Z_2, \dots, Z_c, Z_{c+1}, \dots, Z_s\}$ be an i.i.d. data set drawn from P as described in Setup 1. Let $D' = \{Z_1, Z_2, \dots, Z_c, Z'_{c+1}, \dots, Z'_s\}$ be a second data set that shares the first c observations with D . The remaining $s - c$ observations of D' , i.e. $\{Z'_{c+1}, \dots, Z'_s\}$, are i.i.d. draws from P that are independent of D .

Then, the following inequalities holds for sufficiently large s

$$\mathbb{E}_{D, D'} [Y_1 Y'_{c+1} c(s - c) \kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D')] \lesssim \frac{c(s - c)}{s^2} \mu^2(x) + o(1) \quad (\text{B.13})$$

Proof of Lemma B.4.

Consider first the following argument.

$$\begin{aligned} & \mathbb{E}_{D, D'} [Y_1 Y'_{c+1} c(s - c) \kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D')] \\ &= \mathbb{E}_{1, c+1} [\mu(X_1) \mu(X'_{c+1}) c(s - c) \mathbb{E} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D') \mid X_1, X'_{c+1}]] \\ &\leq \frac{c(s - c)}{s^2} \cdot \mathbb{E}_{1, c+1} \left[|\mu(X_1)| |\mu(X'_{c+1})| s^2 \underbrace{\mathbb{E} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D') \mid X_1, X'_{c+1}]}_{(A)} \right] \end{aligned} \quad (\text{B.14})$$

Analyzing term (A) individually, we find the following.

$$\begin{aligned} (A) &= \mathbb{E} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D') \mid X_1, X'_{c+1}] \\ &= \mathbb{E} \left[\left(\prod_{i=2}^s \mathbb{1}(\|X_1 - x\| < \|X_i - x\|) \right) \left(\prod_{i=1}^c \mathbb{1}(\|X_i - x\| > \|X'_{c+1} - x\|) \right) \left(\prod_{i=c+2}^s \mathbb{1}(\|X'_i - x\| > \|X'_{c+1} - x\|) \right) \mid X_1, X'_{c+1} \right] \\ &= \mathbb{1}(\|X_1 - x\| > \|X'_{c+1} - x\|) \cdot \mathbb{E} \left[\prod_{i=2}^c \mathbb{1}(\|X_i - x\| > \max\{\|X_1 - x\|, \|X'_{c+1} - x\|\}) \mid X_1, X'_{c+1} \right] \\ &\quad \cdot \mathbb{E} \left[\prod_{i=c+1}^s \mathbb{1}(\|X_i - x\| > \|X_1 - x\|) \mid X_1 \right] \cdot \mathbb{E} \left[\prod_{i=c+2}^s \mathbb{1}(\|X'_i - x\| > \|X'_{c+1} - x\|) \mid X'_{c+1} \right] \\ &\leq \mathbb{E} [\kappa(x; Z_1, D) \mid X_1] \cdot \mathbb{E} [\kappa(x; Z'_{c+1}, D') \mid X'_{c+1}] \end{aligned} \quad (\text{B.15})$$

Plugging back into the expression of interest, we find the desired result.

$$\begin{aligned}
& \mathbb{E}_{D,D'} [Y_1 Y'_{c+1} c(s-c) \kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D')] \\
& \leq \frac{c(s-c)}{s^2} \cdot \mathbb{E}_{1,c+1} [|\mu(X_1)| |\mu(X'_{c+1})| s^2 \mathbb{E} [\kappa(x; Z_1, D) | X_1] \cdot \mathbb{E} [\kappa(x; Z'_{c+1}, D') | X'_{c+1}]] \\
& = \frac{c(s-c)}{s^2} \cdot \mathbb{E}_1 [|\mu(X_1)| s \mathbb{E} [\kappa(x; Z_1, D) | X_1]] \mathbb{E}_{c+1} [|\mu(X_{c+1})| s \mathbb{E} [\kappa(x; Z'_{c+1}, D') | X'_{c+1}]] \\
& = \frac{c(s-c)}{s^2} \cdot (\mathbb{E}_1 [|\mu(X_1)| s \mathbb{E} [\kappa(x; Z_1, D) | X_1]])^2 \stackrel{(\text{Lem B.3})}{\lesssim} \frac{c(s-c)}{s^2} \mu^2(x) + o(1) \quad \text{as } s \rightarrow \infty.
\end{aligned} \tag{B.16}$$

■

Lemma B.5.

Fix sample size n , subsampling scale s , and c such that $0 < c \leq s \leq n$. Let $D = \{Z_1, Z_2, \dots, Z_c, Z_{c+1}, \dots, Z_s\}$ be an i.i.d. data set drawn from P as described in Setup 1. Let $D' = \{Z_1, Z_2, \dots, Z_c, Z'_{c+1}, \dots, Z'_s\}$ be a second data set that shares the first c observations with D . The remaining $s - c$ observations of D' , i.e. $\{Z'_{c+1}, \dots, Z'_s\}$, are i.i.d. draws from P that are independent of D .

Then, the following inequalities holds for sufficiently large s

$$\mathbb{E}_{D,D'} [Y_1^2 c \kappa(x; Z_1, D) \kappa(x; Z'_1, D')] \lesssim \frac{c}{2s-c} (\mu^2(x) + \sigma_\varepsilon^2(x)) + o(1) \leq \frac{c}{2s-c} (\mu^2(x) + \bar{\sigma}_\varepsilon^2) + o(1) \tag{B.17}$$

Proof of Lemma B.5.

We can make the following observation.

$$\begin{aligned}
\mathbb{E}_{D,D'} [Y_1^2 c \kappa(x; Z_1, D) \kappa(x; Z'_1, D')] &= \mathbb{E}_1 \left[\mathbb{E} \left[(\mu(X_1) + \varepsilon_1)^2 \mid X_1 \right] c^2 \mathbb{E} [\kappa(x; Z_1, D) \kappa(x; Z'_1, D') \mid X_1] \right] \\
&= \frac{c}{2s-c} \cdot \mathbb{E}_1 [(\mu^2(X_1) + \sigma_\varepsilon^2(X_1)) (2s-c) \mathbb{E} [\kappa(x; Z_1, D) \kappa(x; Z'_1, D') \mid X_1]] \\
&\stackrel{(\text{Lem B.3})}{\lesssim} \frac{c}{2s-c} (\mu^2(x) + \sigma_\varepsilon^2(x)) + o(1)
\end{aligned} \tag{B.18}$$

■

Lemma B.6.

Fix sample size n , subsampling scale s , and c such that $0 < c \leq s \leq n$. Let $D = \{Z_1, Z_2, \dots, Z_c, Z_{c+1}, \dots, Z_s\}$ be an i.i.d. data set drawn from P as described in Setup 1. Let $D' = \{Z_1, Z_2, \dots, Z_c, Z'_{c+1}, \dots, Z'_s\}$ be a second data set that shares the first c observations with D . The remaining $s - c$ observations of D' , i.e. $\{Z'_{c+1}, \dots, Z'_s\}$, are i.i.d. draws from P that are independent of D .

Then, the following inequalities holds for sufficiently large s

$$\mathbb{E}_{D,D'} [Y_{c+1} Y'_{c+1} (s-c)^2 \kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D')] \lesssim \textcolor{red}{LOREMIPSUM} \tag{B.19}$$

Proof of Lemma B.6.

We can make a similar argument as before.

$$\begin{aligned}
& \mathbb{E}_{D,D'} [Y_{c+1} Y'_{c+1} (s-c)^2 \kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D')] \\
&= \mathbb{E}_{D,D'} [\mathbb{E}[(\mu(X_{c+1}) + \varepsilon_{c+1}) \cdot (\mu(X'_{c+1}) + \varepsilon'_{c+1}) \mid X_{c+1}, X'_{c+1}] (s-c)^2 \kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D')] \\
&= \mathbb{E}_{D,D'} [\mu(X_{c+1}) \mu(X'_{c+1}) (s-c)^2 \mathbb{E}[\kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D') \mid X_{c+1}, X'_{c+1}]] \\
&= \textcolor{red}{LOREMIPSUM}
\end{aligned} \tag{B.20}$$

■

Lemma B.7.

Fix sample size n , subsampling scale s , and c such that $0 < c \leq s \leq n$. Let $D = \{Z_1, Z_2, \dots, Z_c, Z_{c+1}, \dots, Z_s\}$ be an i.i.d. data set drawn from Q as described in Setup 2. Let $D' = \{Z_1, Z_2, \dots, Z_c, Z'_{c+1}, \dots, Z'_s\}$ be a second data set that shares the first c observations with D . The remaining $s - c$ observations of D' , i.e. $\{Z'_{c+1}, \dots, Z'_s\}$, are i.i.d. draws from Q that are independent of D .

Then, the following three inequalities hold.

Similarly, consider the CATE estimation setting (Setup 2), i.e. replacing observations drawn from P by observations drawn from Q . Then, analogous inequalities hold, where we replace...

- Y_i by $m(Z_i, \mu, \pi)$
- $\sigma_\varepsilon^2(x)$ by $\frac{\sigma_\varepsilon^2(x)}{\pi(x)(1-\pi(x))}$
- $\mu^2(x)$ by $(\mu_1(x) - \mu_0(x))^2$
- $\bar{\sigma}_\varepsilon^2$ by $\frac{\bar{\sigma}_\varepsilon^2}{p(1-p)}$

Proof of Lemma B.7.

The inequalities follow analogous to the proofs of Lemma B.4, Lemma B.5, and Lemma B.6.

■

Lemma B.8 (Peng, Mentch, and Stefanski (2021) - Lemma 1).

Suppose that $\sum X_i^2 \xrightarrow{p} 1$, $\sum \mathbb{E}[X_i^2] \rightarrow 1$, and $\sum_{i=1}^n \mathbb{E}[Y_i^2] \rightarrow 0$, then

$$\sum [X_i + Y_i]^2 \xrightarrow{p} 1 \quad \text{and} \quad \mathbb{E} \left[\sum (X_i + Y_i)^2 \right] \rightarrow 1. \tag{B.21}$$

Lemma B.9 (Honesty of the DNN/TDNN Estimators).

The DNN and TDNN estimator kernels $\kappa(\cdot, \cdot, D_\ell)$ are Honest in the sense of Wager and Athey (2018).

$$\kappa(x, X_i, D_\ell) \perp\!\!\!\perp Y_i \mid X_i, D_{\ell, -i},$$

where $\perp\!\!\!\perp$ denotes conditional independence and $D_{\ell, -i} = \{Z_l \mid l \in \ell \setminus \{i\}\}$.

C Proofs for Results in Section 3

D Proofs for Results in Section 4

D.1 Closed Form Representations

Proof of Theorem 4.1.

Recall the closed form representation of the DNN estimator as presented in Equation 3.8 and its asymptotic approximation in Equation 3.9.

$$\tilde{\mu}_s(x; \mathbf{D}_n) = \binom{n}{s}^{-1} \sum_{i=1}^{n-s+1} \binom{n-i}{s-1} Y_{(i)} \approx \sum_{i=1}^{n-s+1} \alpha_s (1 - \alpha_s)^{i-1} Y_{(i)} \quad (\text{D.1})$$

Plugging into the Jackknife variance estimator for the DNN estimator now gives us the following where we assume that n is sufficiently large for $n - s + 1$ to be larger than s .

$$\begin{aligned} \hat{\omega}_{\text{JK}}^2 &= \frac{n-1}{n} \sum_{i=1}^n (\tilde{\mu}_s(x; \mathbf{D}_{n,-i}) - \tilde{\mu}_s(x; \mathbf{D}_n))^2 \\ &= \end{aligned} \quad (\text{D.2})$$

Even more simple, we can use the approximate weights to find the following representation.

$$\begin{aligned} \hat{\omega}_{\text{JK}}^2 &= \frac{n-1}{n} \sum_{i=1}^n (\tilde{\mu}_s(x; \mathbf{D}_{n,-i}) - \tilde{\mu}_s(x; \mathbf{D}_n))^2 \\ &\approx \frac{n-1}{n} \alpha_s^2 \sum_{i=1}^{n-s+1} \left(\frac{n}{n-1} \left[\sum_{j=1}^{i-1} \left(1 - \frac{n}{n-1} \alpha_s \right)^{j-1} Y_{(j)} + \sum_{j=i+1}^{n-s+2} \left(1 - \frac{n}{n-1} \alpha_s \right)^{j-2} Y_{(j)} \right] - \sum_{j=1}^{n-s+1} (1 - \alpha_s)^{j-1} Y_{(j)} \right)^2 \\ &= \frac{s^2}{n(n-1)} \sum_{i=1}^{n-s+1} \left(\sum_{j=1}^{i-1} \left[\left(\frac{n-1-s}{n-1} \right)^{j-1} - \left(\frac{n-s}{n} \right)^{j-1} \right] Y_{(j)} + \sum_{j=i+1}^{n-s+2} \left[\left(\frac{n-1-s}{n-1} \right)^{j-2} - \left(\frac{n-s}{n} \right)^{j-1} \right] Y_{(j)} \right. \\ &\quad \left. - \frac{n-1}{n} (1 - \alpha_s)^{i-1} Y_{(i)} \right)^2 \end{aligned} \quad (\text{D.3})$$

The closed form of the Jackknife variance estimator for the TDNN estimator follows from the same approach.

LOREM IPSUM ■

D.2 NPR - Kernel (Conditional) Expectations

As part of deriving consistency results for the variance estimators under consideration, we need to do a careful analysis of the Kernel of the DNN and TDNN estimators. In this section of the appendix we will thus derive the expectations of the kernel and its corresponding Hájek projection. First, we start with the nonparametric regression setup.

Lemma D.1 (NPR - DNN Kernel Expectation).

Let x denote a point of interest. Then

$$\mathbb{E}_D [h_s(x; D)] = \mathbb{E}_1 \left[\mu(X_1) s (1 - \psi(B(x, \|X_1 - x\|)))^{s-1} \right] \longrightarrow \mu(x) \quad \text{as } s \rightarrow \infty \quad (\text{D.4})$$

Proof of Lemma D.1. This result follows immediately from Lemma B.3. ■

Lemma D.2 (NPR - DNN Hajék Kernel Expectation).

Let $z_1 = (x_1, y_1)$ denote a specific realization of Z and x denote a point of interest. Then

$$\psi_s^1(x; z_1) = \varepsilon_1 \mathbb{E}_D \left[\kappa(x; Z_1, D) \mid X_1 = x_1 \right] + \mathbb{E}_D \left[\sum_{i=1}^s \kappa(x; Z_i, D) \mu(X_i) \mid X_1 = x_1 \right] \quad (\text{D.5})$$

Proof of Lemma D.2.

$$\begin{aligned} \psi_s^1(x; z_1) &= \mathbb{E}_D [h_s(x; D) \mid Z_1 = z_1] = \mathbb{E}_D \left[\sum_{i=1}^s \kappa(x; Z_i, D) Y_i \mid Z_1 = z_1 \right] \\ &= \mathbb{E}_D \left[(\mu(x_1) + \varepsilon_1) \kappa(x; Z_1, D) + \sum_{i=2}^s \kappa(x; Z_i, D) \mu(X_i) \mid Z_1 = z_1 \right] \\ &= \varepsilon_1 \mathbb{E}_D \left[\kappa(x; Z_1, D) \mid X_1 = x_1 \right] + \mathbb{E}_D \left[\sum_{i=1}^s \kappa(x; Z_i, D) \mu(X_i) \mid X_1 = x_1 \right] \end{aligned} \quad (\text{D.6})$$

■

D.3 CATE - Kernel (Conditional) Expectations

Next, we address the CATE estimation setup, where we first consider the scenario where the nuisance parameters are assumed to be known a priori. In a second step, we will show that asymptotically, the estimation of nuisance parameters as described in Definition 1, does not alter the asymptotic analysis of the estimator. For clarity, we point out that in contexts relating to the estimation of the conditional average treatment effect, the kernel or score function h_s could hypothetically signify the first or second stage kernel. As the first stage is effectively covered by the nonparametric regression setup, we will take h_s in these contexts to mean the kernel weighted Neyman-orthogonal score associated with the CATE.

Lemma D.3 (CATE - DNN Kernel Expectation).

Let x denote a point of interest. Then

$$\begin{aligned} \mathbb{E}_D [h_s(x; D)] &= \mathbb{E}_1 \left[(\mu_1(X_1) - \mu_0(X_1) + \beta(W_1, X_1)(Y_1 - \mu_{W_1}(X_1))) s (1 - \psi(B(x, \|X_1 - x\|)))^{s-1} \right] \\ &\longrightarrow \text{CATE}(x) \quad \text{as } s \rightarrow \infty \end{aligned} \quad (\text{D.7})$$

Proof of Lemma D.3. This result follows immediately from Lemma B.3. ■

Lemma D.4 (CATE - DNN Hajék Kernel Expectation).

Let $z_1 = (x_1, W_1, y_1)$ denote a specific realization of Z and x denote a point of interest. Then

$$\psi_s^1(x; z_1) = \beta(W_1, X_1) \varepsilon_1 \cdot \mathbb{E}[\kappa(x; Z_1, D) \mid X_1 = x_1] + \mathbb{E}_D \left[\sum_{i=1}^s \kappa(x; Z_i, D) (\mu_1(X_i) - \mu_0(X_i)) \mid Z_1 = z_1 \right] \quad (\text{D.8})$$

Proof of Lemma D.4.

$$\begin{aligned} \psi_s^1(x; z_1) &= \mathbb{E}_D [h_s(x; D) \mid Z_1 = z_1] \\ &= \mathbb{E}_D \left[\sum_{i=1}^s \kappa(x; Z_i, D) (\mu_1(X_i) - \mu_0(X_i) + \beta(W_i, X_i)(Y_i - \mu_{W_i}(X_i))) \mid Z_1 = z_1 \right] \\ &= (\mu_1(X_1) - \mu_0(X_1) + \beta(W_1, X_1) \varepsilon_1) \mathbb{E}[\kappa(x; Z_1, D) \mid X_1 = x_1] \\ &\quad + \mathbb{E}_D \left[\sum_{i=2}^s \kappa(x; Z_i, D) (\mu_1(X_i) - \mu_0(X_i)) \mid Z_1 = z_1 \right] \\ &= \beta(W_1, X_1) \varepsilon_1 \cdot \mathbb{E}[\kappa(x; Z_1, D) \mid X_1 = x_1] + \mathbb{E}_D \left[\sum_{i=1}^s \kappa(x; Z_i, D) (\mu_1(X_i) - \mu_0(X_i)) \mid Z_1 = z_1 \right] \end{aligned} \quad (\text{D.9})$$

■

D.4 NPR - Kernel Variances & Covariances

Similar to the previous section of proofs, we will continue by analyzing the variances and covariances of the kernels under consideration. These results will play an important role in the derivation of consistency properties for the variance estimators. Similar to the previous part, we will first consider the nonparametric regression setup and then proceed to the conditional average treatment effect setup.

Lemma D.5 (Adapted from Demirkaya et al. (2024)).

Let $D = \{Z_1, \dots, Z_s\}$ be a vector of i.i.d. random variables drawn from P . Furthermore, let

$$\Omega_s(x) = \mathbb{E} [h_s^2(x; Z_1, \dots, Z_s)] . \quad (\text{D.10})$$

Then,

$$\Omega_s(x) = \mathbb{E}_1 \left[(\mu(X_1) + \varepsilon_1)^2 s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)] \right] \lesssim \mu^2(x) + \bar{\sigma}_\varepsilon^2 + o(1) \quad \text{as } s \rightarrow \infty. \quad (\text{D.11})$$

Proof of Lemma D.5.

$$\begin{aligned} \Omega_s(x) &= \mathbb{E} [h_s^2(x; Z_1, \dots, Z_s)] = \mathbb{E}_D \left[\left(\sum_{i=1}^s \kappa(x; Z_i, D) Y_i \right)^2 \right] = \mathbb{E}_D \left[\sum_{i=1}^s \sum_{j=1}^s (\kappa(x; Z_i, D) \kappa(x; Z_j, D) Y_i Y_j) \right] \\ &= \mathbb{E}_D [s \kappa(x; Z_1, D) Y_1^2] = \mathbb{E}_1 [Y_1^2 s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] \xrightarrow{(\text{Lem B.3})} \mu^2(x) + \sigma_\varepsilon^2(x) \quad \text{as } s \rightarrow \infty \end{aligned} \quad (\text{D.12})$$

Thus, we obtain the desired result. ■

Lemma D.6.

Let $D = \{Z_1, \dots, Z_s\}$ be a vector of i.i.d. random variables drawn from P . Let $D' = \{Z_1, \dots, Z_c, Z'_{c+1}, \dots, Z'_s\}$ where Z'_{c+1}, \dots, Z'_s are i.i.d. draws from P that are independent of D . Furthermore, let

$$\Omega_s^c(x) = \mathbb{E} [h_s(x; Z_1, \dots, Z_c, Z_{c+1}, \dots, Z_s) \cdot h_s(x; Z_1, \dots, Z_c, Z'_{c+1}, \dots, Z'_s)] . \quad (\text{D.13})$$

Then,

$$\Omega_s^c(x) \lesssim \frac{s^2 + cs - c^2}{s^2} \mu^2(x) + (c/s) \bar{\sigma}_\varepsilon^2 + o(1) \quad \text{for } s \text{ sufficiently large} \quad (\text{D.14})$$

and thus

$$\Omega_s^c(x) \lesssim \mu^2(x) + \bar{\sigma}_\varepsilon^2 + o(1) \quad \text{as } s \rightarrow \infty. \quad (\text{D.15})$$

Proof of Lemma D.6.

$$\begin{aligned}
\Omega_s^c(x) &= \mathbb{E} [h_s(x; Z_1, \dots, Z_c, Z_{c+1}, \dots, Z_s) \cdot h_s(x; Z_1, \dots, Z_c, Z'_{c+1}, \dots, Z'_s)] \\
&= \mathbb{E}_{D, D'} \left[\left(\sum_{i=1}^s \kappa(x; Z_i, D) Y_i \right) \left(\sum_{j=1}^c \kappa(x; Z_j, D') Y_j + \sum_{j=c+1}^s \kappa(x; Z'_j, D') Y'_j \right) \right] \\
&= \mathbb{E}_{D, D'} \left[\sum_{i=1}^c \sum_{j=1}^c \kappa(x; Z_i, D) \kappa(x; Z_j, D') Y_i Y_j \right] + \mathbb{E}_{D, D'} \left[\sum_{i=1}^c \sum_{j=c+1}^s \kappa(x; Z_i, D) \kappa(x; Z'_j, D') Y_i Y'_j \right] \\
&\quad + \mathbb{E}_{D, D'} \left[\sum_{i=c+1}^s \sum_{j=1}^c \kappa(x; Z_i, D) \kappa(x; Z_j, D') Y_i Y_j \right] + \mathbb{E}_{D, D'} \left[\sum_{i=c+1}^s \sum_{j=c+1}^s \kappa(x; Z_i, D) \kappa(x; Z'_j, D') Y_i Y'_j \right] \\
&= \underbrace{\mathbb{E}_{D, D'} [c \kappa(x; Z_1, D) \kappa(x; Z_1, D') Y_1^2]}_{(A)} + \underbrace{\mathbb{E}_{D, D'} [c(s-c) \kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D') Y_1 Y'_{c+1}]}_{(B)} \\
&\quad + \underbrace{\mathbb{E}_{D, D'} [c(s-c) \kappa(x; Z_{c+1}, D) \kappa(x; Z_1, D') Y_{c+1} Y_1]}_{(C)} \\
&\quad + \underbrace{\mathbb{E}_{D, D'} [(s-c)^2 \kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D') Y_{c+1} Y'_{c+1}]}_{(D)}
\end{aligned} \tag{D.16}$$

Starting from this decomposition, we will analyze the terms one by one. First, by Lemma B.3, we find the following.

$$\begin{aligned}
(A) &= \mathbb{E}_{D, D'} [c \kappa(x; Z_1, D) \kappa(x; Z_1, D') Y_1^2] = (c/s) \mathbb{E}_1 [Y_1^2 s \mathbb{E}_{2:s} [\kappa(x; Z_1, D) \kappa(x; Z_1, D')]] \\
&\leq (c/s) \mathbb{E}_1 [Y_1^2 s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] \stackrel{(\text{Lem B.3})}{\lesssim} (c/s) (\mu^2(x) + \sigma_\varepsilon(x)) + o(1)
\end{aligned} \tag{D.17}$$

Similarly, we can find that:

$$(B) = \mathbb{E}_{D, D'} [c(s-c) \kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D') Y_1 Y'_{c+1}] \stackrel{\text{Lem B.4}}{\lesssim} \frac{c(s-c)}{s^2} \mu^2(x) + o(1) \tag{D.18}$$

Following analogous steps, we find the same result for the third term.

$$(C) = \mathbb{E}_{D, D'} [c(s-c) \kappa(x; Z_{c+1}, D) \kappa(x; Z_1, D') Y_{c+1} Y_1] \stackrel{\text{Lem B.4}}{\lesssim} \frac{c(s-c)}{s^2} \mu^2(x) + o(1) \tag{D.19}$$

The fourth term can be asymptotically bounded in the following way.

$$\begin{aligned}
(D) &= \mathbb{E}_{D, D'} [(s-c)^2 \kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D') Y_{c+1} Y'_{c+1}] \\
&= \mathbb{E}_{D, D'} [\mu(X_{c+1}) \mu(X'_{c+1}) (s-c)^2 \kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D')] \\
&\leq \mathbb{E}_D [|\mu(X_{c+1})| (s-c) \kappa(x; Z_{c+1}, D)] \mathbb{E}_{D'} [|\mu(X'_{c+1})| (s-c) \kappa(x; Z'_{c+1}, D')] \\
&= \frac{(s-c)^2}{s^2} \mathbb{E}_D [|\mu(X_{c+1})| s \kappa(x; Z_{c+1}, D)] \mathbb{E}_{D'} [|\mu(X'_{c+1})| s \kappa(x; Z'_{c+1}, D')] \\
&= \frac{(s-c)^2}{s^2} (\mathbb{E}_D [|\mu(X_{c+1})| s \kappa(x; Z_{c+1}, D)])^2 \lesssim \frac{(s-c)^2}{s^2} \mu^2(x) + o(1)
\end{aligned} \tag{D.20}$$

The result of Lemma D.6 follows immediately by summing up the asymptotic bounds for the individual terms. \blacksquare

Lemma D.7.

Let $D = \{Z_1, \dots, Z_{s_2}\}$ be a vector of i.i.d. random variables drawn from P for $s_2 > s_1$. Furthermore, let

$$\Upsilon_{s_1, s_2}(x) = \mathbb{E}[h_{s_1}(x; Z_1, \dots, Z_{s_1}) \cdot h_{s_2}(x; Z_1, \dots, Z_{s_1}, \dots, Z_{s_2})]. \quad (\text{D.21})$$

Then,

$$\Upsilon_{s_1, s_2}(x) \lesssim \mu^2(x) + \bar{\sigma}_\varepsilon^2 + o(1) \quad \text{as } s_1, s_2 \rightarrow \infty \quad \text{with } 0 < \mathfrak{c} \leq s_1/s_2 \leq 1 - \mathfrak{c} < 1. \quad (\text{D.22})$$

Proof of Lemma D.7.

$$\begin{aligned} \Upsilon_{s_1, s_2}(x) &= \mathbb{E}[h_{s_1}(x; Z_1, \dots, Z_{s_1}) \cdot h_{s_2}(x; Z_1, \dots, Z_{s_1}, \dots, Z_{s_2})] \\ &= \mathbb{E}_D \left[\left(\sum_{i=1}^{s_1} \kappa(x; Z_i, D_{[s_1]}) Y_i \right) \left(\sum_{j=1}^{s_1} \kappa(x; Z_j, D) Y_j + \sum_{j=s_1+1}^{s_2} \kappa(x; Z_j, D) Y_j \right) \right] \\ &= \mathbb{E}_D \left[\sum_{i=1}^{s_1} \kappa(x; Z_i, D) Y_i^2 \right] + \mathbb{E}_D \left[\sum_{i=1}^{s_1} \sum_{j=s_1+1}^{s_2} \kappa(x; Z_i, D_{[s_1]}) \kappa(x; Z_j, D) Y_i Y_j \right] \\ &= \mathbb{E}_D [Y_1^2 s_1 \kappa(x; Z_1, D)] + \mathbb{E}_D [Y_1 Y_{s_2} s_1 (s_2 - s_1) \kappa(x; Z_1, D_{[s_1]}) \kappa(x; Z_{s_2}, D)] \\ &= \mathbb{E}_D [(\mu^2(X_1) + \sigma_\varepsilon^2(X_1)) s_1 \kappa(x; Z_1, D)] + \mathbb{E}_D [\mu(X_1) \mu(X_{s_2}) s_1 (s_2 - s_1) \kappa(x; Z_1, D_{[s_1]}) \kappa(x; Z_{s_2}, D)] \\ &= \frac{s_1}{s_2} \mathbb{E}_D [(\mu^2(X_1) + \sigma_\varepsilon^2(X_1)) s_1 \kappa(x; Z_1, D)] + \frac{s_2 - s_1}{s_2} \mathbb{E}_D [\mu(X_1) \mu(X_{s_2}) s_1 s_2 \kappa(x; Z_1, D_{[s_1]}) \kappa(x; Z_{s_2}, D)] \\ &\leq \frac{s_1}{s_2} \mathbb{E}_D [(\mu^2(X_1) + \sigma_\varepsilon^2(X_1)) s_2 \kappa(x; Z_1, D)] \\ &\quad + \frac{s_2 - s_1}{s_2} \mathbb{E}_D [|\mu(X_1)| s_1 \kappa(x; Z_1, D_{[s_1]})] \mathbb{E}_D [|\mu(X_{s_2})| s_2 \kappa(x; Z_{s_2}, D)] \\ &\lesssim \mu^2(x) + \sigma_\varepsilon^2(x) + o(1) \leq \mu^2(x) + \bar{\sigma}_\varepsilon^2 + o(1). \end{aligned} \quad (\text{D.23})$$

■

Lemma D.8.

Let $D = \{Z_1, \dots, Z_{s_2}\}$ be a vector of i.i.d. random variables drawn from P for $s_2 > s_1$. Let $D' = \{Z_1, \dots, Z_c, Z'_{c+1}, \dots, Z'_{s_1}\}$ where $Z'_{c+1}, \dots, Z'_{s_1}$ are i.i.d. draws from P that are independent of D . Furthermore, let

$$\Upsilon_{s_1, s_2}^c(x) = \mathbb{E} [h_{s_1}(x; Z_1, \dots, Z_c, Z'_{c+1}, \dots, Z'_{s_1}) \cdot h_{s_2}(x; Z_1, \dots, Z_{s_2})]. \quad (\text{D.24})$$

Then,

$$\Upsilon_{s_1, s_2}^c(x) \lesssim \frac{cs_2 - c^2 + s_1 s_2}{s_1 s_2} \mu^2(x) + (c/s_1) \bar{\sigma}_\varepsilon^2 + o(1) \quad (\text{D.25})$$

for s_1, s_2 sufficiently large with $0 < \mathfrak{c} \leq s_1/s_2 \leq 1 - \mathfrak{c} < 1$

and thus

$$\Upsilon_{s_1, s_2}^c(x) \lesssim \mu^2(x) + o(1) \quad \text{as } s_1, s_2 \rightarrow \infty \quad \text{with } 0 < \mathfrak{c} \leq s_1/s_2 \leq 1 - \mathfrak{c} < 1. \quad (\text{D.26})$$

Proof of Lemma D.8.

$$\begin{aligned} \Upsilon_{s_1, s_2}^c(x) &= \mathbb{E} [h_{s_1}(x; Z_1, \dots, Z_c, Z'_{c+1}, \dots, Z'_{s_1}) \cdot h_{s_2}(x; Z_1, \dots, Z_{s_2})] \\ &= \mathbb{E}_{D, D'} \left[\left(\sum_{i=1}^c \kappa(x; Z_i, D') Y_i + \sum_{i=c+1}^{s_1} \kappa(x; Z'_i, D') Y'_i \right) \left(\sum_{j=1}^c \kappa(x; Z_j, D) Y_j + \sum_{j=c+1}^{s_2} \kappa(x; Z_j, D) Y_j \right) \right] \\ &= \underbrace{\mathbb{E}_{D, D'} \left[\sum_{i=1}^c \sum_{j=1}^c \kappa(x; Z_i, D') \kappa(x; Z_j, D) Y_i Y_j \right]}_{(A)} + \underbrace{\mathbb{E}_{D, D'} \left[\left(\sum_{i=1}^c \kappa(x; Z_i, D') Y_i \right) \left(\sum_{j=c+1}^{s_2} \kappa(x; Z_j, D) Y_j \right) \right]}_{(B)} \\ &\quad + \underbrace{\mathbb{E}_{D, D'} \left[\sum_{i=c+1}^{s_1} \sum_{j=1}^c \kappa(x; Z'_i, D') \kappa(x; Z_j, D) Y'_i Y_j \right]}_{(C)} + \underbrace{\mathbb{E}_{D, D'} \left[\left(\sum_{i=c+1}^{s_1} \kappa(x; Z'_i, D') Y'_i \right) \left(\sum_{j=c+1}^{s_2} \kappa(x; Z_j, D) Y_j \right) \right]}_{(D)} \end{aligned} \quad (\text{D.27})$$

Again, we have four terms to analyze individually.

$$\begin{aligned} (A) &= \mathbb{E}_{D, D'} \left[\sum_{i=1}^c \sum_{j=1}^c \kappa(x; Z_i, D') \kappa(x; Z_j, D) Y_i Y_j \right] \\ &= \mathbb{E}_{D, D'} \left[\sum_{i=1}^c Y_i^2 \kappa(x; Z_i, D') \kappa(x; Z_i, D) \right] \\ &= \mathbb{E}_{D, D'} [Y_1^2 c \kappa(x; Z_1, D') \kappa(x; Z_1, D)] = \mathbb{E}_{D, D'} [(\mu^2(X_1) + \sigma_\varepsilon^2(X_1)) c \kappa(x; Z_1, D_{[c]}) \kappa(x; Z_1, D'_{c+1:s_1})] \\ &= \mathbb{E}_D [(\mu^2(X_1) + \sigma_\varepsilon^2(X_1)) c \kappa(x; Z_1, D)] = \frac{c}{s_1} \mathbb{E}_D [(\mu^2(X_1) + \sigma_\varepsilon^2(X_1)) s_1 \kappa(x; Z_1, D)] \\ &\lesssim (c/s_1)(\mu^2(x) + \sigma_\varepsilon^2(x)) + o(1) \leq (c/s_1)(\mu^2(x) + \bar{\sigma}_\varepsilon^2) + o(1) \end{aligned} \quad (\text{D.28})$$

Considering the second term, we find the following.

$$\begin{aligned}
(B) &= \mathbb{E}_{D,D'} \left[\left(\sum_{i=1}^c \kappa(x; Z_i, D') Y_i \right) \left(\sum_{j=c+1}^{s_2} \kappa(x; Z_j, D) Y_j \right) \right] = \mathbb{E}_{D,D'} \left[\sum_{i=1}^c \sum_{j=c+1}^{s_2} Y_i Y_j \kappa(x; Z_i, D') \kappa(x; Z_j, D) \right] \\
&= \mathbb{E}_{D,D'} [c(s_2 - c) Y_1 Y_{s_1} \kappa(x; Z_1, D') \kappa(x; Z_{s_2}, D)] = \frac{c(s_2 - c)}{s_1 s_2} \mathbb{E}_{D,D'} [Y_1 Y_{s_2} s_1 s_2 \kappa(x; Z_1, D') \kappa(x; Z_{s_2}, D)] \\
&\leq \frac{c(s_2 - c)}{s_1 s_2} \mathbb{E}_{D'} [|\mu(X_1)| s_1 \kappa(x; Z_1, D')] \mathbb{E}_D [|\mu(X_{s_2})| s_2 \kappa(x; Z_{s_2}, D)] \\
&\lesssim \frac{c(s_2 - c)}{s_1 s_2} \mu^2(x) + o(1)
\end{aligned} \tag{D.29}$$

Similarly, by simplifying the third term, we find the following.

$$\begin{aligned}
(C) &= \mathbb{E}_{D,D'} \left[\sum_{i=c+1}^{s_1} \sum_{j=1}^c \kappa(x; Z'_i, D') \kappa(x; Z_j, D) Y'_i Y_j \right] = \mathbb{E}_{D,D'} [Y'_{s_1} Y_1 (s_1 - c) c \kappa(x; Z'_{s_1}, D') \kappa(x; Z_1, D)] \\
&= \frac{(s_1 - c)c}{s_1 s_2} \mathbb{E}_{D,D'} [\mu(X'_{s_1}) \mu(X_1) s_1 s_2 \kappa(x; Z'_{s_1}, D') \kappa(x; Z_1, D)] \\
&\leq \frac{(s_1 - c)c}{s_1 s_2} \mathbb{E}_D [|\mu(X'_{s_1})| s_1 \kappa(x; Z'_{s_1}, D')] \mathbb{E}_D [|\mu(X_1)| s_2 \kappa(x; Z_1, D)] \\
&\lesssim \frac{(s_1 - c)c}{s_1 s_2} \mu^2(x) + o(1)
\end{aligned} \tag{D.30}$$

Lastly, concerning the fourth term, observe the following.

$$\begin{aligned}
(D) &= \mathbb{E}_{D,D'} \left[\left(\sum_{i=c+1}^{s_1} \kappa(x; Z'_i, D') Y'_i \right) \left(\sum_{j=c+1}^{s_2} \kappa(x; Z_j, D) Y_j \right) \right] = \mathbb{E}_{D,D'} \left[\sum_{i=c+1}^{s_1} \sum_{j=c+1}^{s_2} \kappa(x; Z'_i, D') \kappa(x; Z_j, D) Y'_i Y_j \right] \\
&= \mathbb{E}_{D,D'} [\mu(X'_{s_1}) \mu(X_{s_2}) (s_1 - c)(s_2 - c) \kappa(x; Z'_{s_1}, D') \kappa(x; Z_{s_2}, D)] \\
&= \frac{(s_1 - c)(s_2 - c)}{s_1 s_2} \mathbb{E}_{D,D'} [\mu(X'_{s_1}) \mu(X_{s_2}) s_1 s_2 \kappa(x; Z'_{s_1}, D') \kappa(x; Z_{s_2}, D)] \\
&\leq \frac{(s_1 - c)(s_2 - c)}{s_1 s_2} \mathbb{E}_{D'} [|\mu(X'_{s_1})| s_1 \kappa(x; Z'_{s_1}, D')] \mathbb{E}_D [|\mu(X_{s_2})| s_2 \kappa(x; Z_{s_2}, D)] \\
&\lesssim \frac{(s_1 - c)(s_2 - c)}{s_1 s_2} \mu^2(x) + o(1)
\end{aligned} \tag{D.31}$$

■

Lemma D.9 (Kernel Variance of the TDNN Kernel). *For the kernel of the TDNN estimator with subsampling scales s_1 and s_2 , it holds that*

$$\zeta_{s_1, s_2}^{s_2}(x) \lesssim \mu^2(x) + \bar{\sigma}_\varepsilon + o(1) \quad \text{as } s_1, s_2 \rightarrow \infty \quad \text{with } 0 < \mathfrak{c} \leq s_1/s_2 \leq 1 - \mathfrak{c} < 1. \quad (\text{D.32})$$

Proof of Lemma D.9. Consider first the following decomposition.

$$\begin{aligned} \zeta_{s_1, s_2}^{s_2}(x) &= \text{Var}(h_{s_1, s_2}(x; Z_1, \dots, Z_{s_2})) = \text{Var}_D(h_{s_1, s_2}(x; D)) \\ &\leq \mathbb{E}_D[h_{s_1, s_2}^2(x; D)] = \mathbb{E}_D[(w_1^* \tilde{\mu}_{s_1}(x; D) + w_2^* h_{s_2}(x; D))^2] \\ &= (w_1^*)^2 \mathbb{E}_D[\tilde{\mu}_{s_1}^2(x; D)] + 2w_1^* w_2^* \mathbb{E}_D[\tilde{\mu}_{s_1}(x; D) h_{s_2}(x; D)] + (w_2^*)^2 \Omega_{s_2} \end{aligned} \quad (\text{D.33})$$

Then, observe the following.

$$\begin{aligned} \mathbb{E}_D[\tilde{\mu}_{s_1}^2(x; D)] &= \mathbb{E}_D\left[\left(\binom{s_2}{s_1}^{-1} \sum_{\ell \in L_{s_2, s_1}} h_{s_1}(x; D_\ell)\right)^2\right] = \binom{s_2}{s_1}^{-2} \mathbb{E}_D\left[\sum_{\ell, \ell' \in L_{s_2, s_1}} h_{s_1}(x; D_\ell) h_{s_1}(x; D_{\ell'})\right] \\ &= \binom{s_2}{s_1}^{-2} \sum_{c=0}^{s_1} \binom{s_2}{s_1} \binom{s_1}{c} \binom{s_2 - s_1}{s_1 - c} \Omega_{s_1}^c = \binom{s_2}{s_1}^{-1} \sum_{c=0}^{s_1} \binom{s_1}{c} \binom{s_2 - s_1}{s_1 - c} \Omega_{s_1}^c \\ &\lesssim \Omega_{s_1} \lesssim \mu(x)^2 + \sigma_\varepsilon^2 + o(1) \quad \text{as } s \rightarrow \infty \end{aligned} \quad (\text{D.34})$$

Recall that by Lemma D.5, we have the following.

$$\Omega_{s_2} \lesssim \mu(x)^2 + \sigma_\varepsilon^2 + o(1) \quad \text{as } s \rightarrow \infty \quad (\text{D.35})$$

Lastly, consider the following.

$$\begin{aligned} \mathbb{E}_D[\tilde{\mu}_{s_1}(x; D) h_{s_2}(x; D)] &= \mathbb{E}_D\left[\binom{s_2}{s_1}^{-1} \sum_{\ell \in L_{s_2, s_1}} h_{s_1}(x; D_\ell) h_{s_2}(x; D)\right] \\ &= \mathbb{E}_D[h_{s_1}(x; D_{[s_1]}) h_{s_2}(x; D)] = \Upsilon_{s_1, s_2}(x) \end{aligned} \quad (\text{D.36})$$

Thus, we find the following.

$$\begin{aligned} \zeta_{s_2, s_2}(x) &\lesssim (w_1^*)^2 \Omega_{s_1} + 2w_1^* w_2^* \Upsilon_{s_1, s_2}(x) + (w_1^*)^2 \Omega_{s_2} \\ &\lesssim (w_1^* + w_2^*)^2 (\mu^2(x) + \sigma_\varepsilon) + o(1) = \mu^2(x) + \sigma_\varepsilon + o(1). \end{aligned} \quad (\text{D.37})$$

■

Lemma D.10 (Lemma 10 - Demirkaya et al. (2024)). *For the kernel of the TDNN estimator with subsampling scales s_1 and s_2 satisfying*

$$0 < \mathfrak{c} \leq s_1/s_2 \leq 1 - \mathfrak{c} < 1 \quad \text{and} \quad s_2 = o(n), \quad (\text{D.38})$$

it holds that

$$\zeta_{s_1, s_2}^1(x) \sim s_2^{-1}. \quad (\text{D.39})$$

D.5 CATE - Kernel Variances & Covariances

Next, we will continue by showing analogous properties in the CATE setting. Similar to before, we will start under the assumption that the functional nuisance parameters are known a priori, to then show that the estimation of said parameters does not impact the asymptotic behavior of the estimator.

Lemma D.11.

Let $D = \{Z_1, \dots, Z_s\}$ be a vector of i.i.d. random variables generated by the setup shown in Assumption 2. Furthermore, let

$$\Omega_s(x) = \mathbb{E} \left[h_s^2(x; Z_1, \dots, Z_s) \right]. \quad (\text{D.40})$$

Then,

$$\Omega_s(x) \lesssim (\mu_1(x) - \mu_0(x))^2 + \frac{\bar{\sigma}_\varepsilon^2}{\mathbf{p}(1-\mathbf{p})} + o(1) \quad (\text{D.41})$$

Proof of Lemma D.11.

First, notice that we can decompose the quantity of interest in the following way.

$$\begin{aligned} \Omega_s(x) &= \mathbb{E} \left[h_s^2(x; Z_1, \dots, Z_s) \right] = \mathbb{E}_D \left[\left(\sum_{i=1}^s \kappa(x; Z_i, D) m(Z_i; \mu, \pi) \right)^2 \right] \\ &= \mathbb{E}_D \left[\sum_{i=1}^s \sum_{j=1}^s \kappa(x; Z_i, D) \kappa(x; Z_j, D) m(Z_i; \mu, \pi) m(Z_j; \mu, \pi) \right] \\ &= \mathbb{E}_D \left[s \kappa(x; Z_1, D) m^2(Z_1; \mu, \pi) \right] = \mathbb{E}_1 \left[m^2(Z_1; \mu, \pi) s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)] \right] \\ &\stackrel{(\text{Lem B.3})}{\longrightarrow} (\mu_1(x) - \mu_0(x))^2 + \frac{\sigma_\varepsilon^2(x)}{\pi(x)(1-\pi(x))} \quad \text{as } s \rightarrow \infty \end{aligned} \quad (\text{D.42})$$

This gives us the desired result.

$$\Omega_s(x) \lesssim (\mu_1(x) - \mu_0(x))^2 + \frac{\bar{\sigma}_\varepsilon^2}{\mathbf{p}(1-\mathbf{p})} + o(1) \quad (\text{D.43})$$

■

Lemma D.12.

Let $D = \{Z_1, \dots, Z_s\}$ be a vector of i.i.d. random variables drawn from as described in Setup 2. Let $D' = \{Z_1, \dots, Z_c, Z'_{c+1}, \dots, Z'_s\}$ where Z'_{c+1}, \dots, Z'_s are i.i.d. draws from the model that are independent of D . Furthermore, let

$$\Omega_s^c(x) = \mathbb{E} \left[h_s(x; Z_1, \dots, Z_c, Z_{c+1}, \dots, Z_s) \cdot h_s(x; Z_1, \dots, Z_c, Z'_{c+1}, \dots, Z'_s) \right]. \quad (\text{D.44})$$

Then,

$$\Omega_s^c(x) \lesssim (\mu_1(x) - \mu_0(x))^2 + \frac{\bar{\sigma}_\varepsilon^2}{\mathbf{p}(1-\mathbf{p})} + o(1). \quad (\text{D.45})$$

Proof of Lemma D.12. First, we decompose the term of interest in a similar fashion to before.

$$\begin{aligned} \Omega_s^c(x) &= \mathbb{E} \left[h_s(x; Z_1, \dots, Z_c, Z_{c+1}, \dots, Z_s) \cdot h_s(x; Z_1, \dots, Z_c, Z'_{c+1}, \dots, Z'_s) \right] \\ &= \mathbb{E}_{D, D'} \left[\left(\sum_{i=1}^s \kappa(x; Z_i, D) m(Z_i; \mu, p) \right) \left(\sum_{j=1}^c \kappa(x; Z_j, D') m(Z_j; \mu, p) + \sum_{j=c+1}^s \kappa(x; Z'_j, D') m(Z'_j; \mu, p) \right) \right] \\ &= \mathbb{E}_{D, D'} \left[\underbrace{\left(\sum_{i=1}^c \kappa(x; Z_i, D) m(Z_i; \mu, p) \right) \left(\sum_{j=1}^c \kappa(x; Z_j, D') m(Z_j; \mu, p) \right)}_{(A)} \right] \\ &\quad + \mathbb{E}_{D, D'} \left[\underbrace{\left(\sum_{i=1}^c \kappa(x; Z_i, D) m(Z_i; \mu, p) \right) \left(\sum_{j=c+1}^s \kappa(x; Z'_j, D') m(Z'_j; \mu, p) \right)}_{(B)} \right] \\ &\quad + \mathbb{E}_{D, D'} \left[\underbrace{\left(\sum_{i=c+1}^s \kappa(x; Z_i, D) m(Z_i; \mu, p) \right) \left(\sum_{j=1}^c \kappa(x; Z_j, D') m(Z_j; \mu, p) \right)}_{(C)} \right] \\ &\quad + \mathbb{E}_{D, D'} \left[\underbrace{\left(\sum_{i=c+1}^s \kappa(x; Z_i, D) m(Z_i; \mu, p) \right) \left(\sum_{j=c+1}^s \kappa(x; Z'_j, D') m(Z'_j; \mu, p) \right)}_{(D)} \right] \end{aligned} \quad (\text{D.46})$$

Considering these terms one by one, we can make the following observations.

$$\begin{aligned} (A) &= \mathbb{E}_{D, D'} \left[\left(\sum_{i=1}^c \kappa(x; Z_i, D) m(Z_i; \mu, p) \right) \left(\sum_{j=1}^c \kappa(x; Z_j, D') m(Z_j; \mu, p) \right) \right] \\ &= \mathbb{E}_{D, D'} \left[\sum_{i=1}^c \sum_{j=1}^c \kappa(x; Z_i, D) \kappa(x; Z_j, D') m(Z_i; \mu, p) m(Z_j; \mu, p) \right] \\ &= \mathbb{E}_1 \left[m^2(Z_1; \mu, p) c \mathbb{E}_{2:s} [\kappa(x; Z_1, D) \kappa(x; Z_1, D')] \right] \leq (c/s) \cdot \mathbb{E}_1 \left[m^2(Z_1; \mu, p) s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)] \right] \\ &\stackrel{(\text{Lem B.3})}{\lesssim} (c/s) \left[(\mu_1(x) - \mu_0(x))^2 + \frac{\sigma_\varepsilon^2(x)}{\pi(x)(1-\pi(x))} \right] + o(1) \end{aligned} \quad (\text{D.47})$$

Similarly, for the second term, we can make the following observation.

$$\begin{aligned}
(B) &= \mathbb{E}_{D,D'} \left[\left(\sum_{i=1}^c \kappa(x; Z_i, D) m(Z_i; \mu, \pi) \right) \left(\sum_{j=c+1}^s \kappa(x; Z'_j, D') m(Z'_j; \mu, \pi) \right) \right] \\
&= \mathbb{E}_{D,D'} \left[\sum_{i=1}^c \sum_{j=c+1}^s \kappa(x; Z_i, D) \kappa(x; Z'_j, D') m(Z_i; \mu, \pi) m(Z'_j; \mu, \pi) \right] \\
&= \mathbb{E}_{D,D'} [c(s-c) \kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D') m(Z_1; \mu, \pi) m(Z'_{c+1}; \mu, \pi)] \\
&\leq \mathbb{E}_D [c \kappa(x; Z_1, D) |m(Z_1; \mu, \pi)|] \mathbb{E}_{D'} [(s-c) \kappa(x; Z'_{c+1}, D') |m(Z'_{c+1}; \mu, \pi)|] \\
&= \frac{c(s-c)}{s^2} \cdot (\mathbb{E}_1 [|m(Z_1; \mu, \pi)| s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)]])^2 \\
&\lesssim \frac{c(s-c)}{s^2} (\mu_1(x) - \mu_0(x))^2 + o(1)
\end{aligned} \tag{D.48}$$

Applying the same principles to the third term we find a similar result.

$$\begin{aligned}
(C) &= \mathbb{E}_{D,D'} \left[\left(\sum_{i=c+1}^s \kappa(x; Z_i, D) m(Z_i; \mu, p) \right) \left(\sum_{j=1}^c \kappa(x; Z_j, D') m(Z_j; \mu, p) \right) \right] \\
&\lesssim \frac{c(s-c)}{s^2} (\mu_1(x) - \mu_0(x))^2 + o(1) \\
(D) &= \mathbb{E}_{D,D'} \left[\left(\sum_{i=c+1}^s \kappa(x; Z_i, D) m(Z_i; \mu, p) \right) \left(\sum_{j=c+1}^s \kappa(x; Z'_j, D') m(Z'_j; \mu, p) \right) \right] \\
&= \mathbb{E}_{D,D'} [(s-c)^2 \kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D') m(Z_{c+1}; \mu, p) m(Z'_{c+1}; \mu, p)] \\
&\leq \frac{(s-c)^2}{s^2} \cdot \mathbb{E}_D [|m(Z_{c+1}; \mu, p)| s \kappa(x; Z_{c+1}, D)] \mathbb{E}_{D'} [|m(Z'_{c+1}; \mu, p)| s \kappa(x; Z'_{c+1}, D')] \\
&= \frac{(s-c)^2}{s^2} (\mathbb{E}_1 [|m(Z_1; \mu, p)| s \mathbb{E}_{2:s} [\kappa(x; Z_{c+1}, D)]])^2 \\
&\lesssim \frac{(s-c)^2}{s^2} (\mu_1(x) - \mu_0(x))^2 + o(1)
\end{aligned} \tag{D.50}$$

Thus, we find the desired result.

$$\begin{aligned}
\Omega_s^c(x) &= (A) + (B) + (C) + (D) \\
&\lesssim \frac{c}{s} \left[(\mu_1(x) - \mu_0(x))^2 + \frac{\bar{\sigma}_\varepsilon^2}{\mathbf{p}(1-\mathbf{p})} \right] + 2 \frac{c(s-c)}{s^2} (\mu_1(x) - \mu_0(x))^2 + \frac{(s-c)^2}{s^2} (\mu_1(x) - \mu_0(x))^2 + o(1) \\
&= \left[\frac{cs + 2c(s-c) + (s-c)^2}{s^2} \right] (\mu_1(x) - \mu_0(x))^2 + \frac{c}{s} \frac{\bar{\sigma}_\varepsilon^2}{\mathbf{p}(1-\mathbf{p})} + o(1) \\
&\lesssim (\mu_1(x) - \mu_0(x))^2 + \frac{\bar{\sigma}_\varepsilon^2}{\mathbf{p}(1-\mathbf{p})} + o(1)
\end{aligned} \tag{D.51}$$

■

Lemma D.13.

Let $D = \{Z_1, \dots, Z_{s_2}\}$ be a vector of i.i.d. random variables drawn from Q for $s_2 > s_1$. Furthermore, let

$$\Upsilon_{s_1, s_2}(x) = \mathbb{E}[h_{s_1}(x; Z_1, \dots, Z_{s_1}) \cdot h_{s_2}(x; Z_1, \dots, Z_{s_1}, \dots, Z_{s_2})]. \quad (\text{D.52})$$

Then,

$$\Upsilon_{s_1, s_2}(x) \lesssim 2(\mu_1(x) - \mu_0(x))^2 + \frac{\bar{\sigma}_\varepsilon^2}{\mathfrak{p}(1 - \mathfrak{p})} + o(1) \quad \text{as } s_1, s_2 \rightarrow \infty \quad \text{with } 0 < \mathfrak{c} \leq s_1/s_2 \leq 1 - \mathfrak{c} < 1. \quad (\text{D.53})$$

Proof of Lemma D.13. Consider first the following.

$$\begin{aligned} \Upsilon_{s_1, s_2}(x) &= \mathbb{E}[h_{s_1}(x; Z_1, \dots, Z_{s_1}) \cdot h_{s_2}(x; Z_1, \dots, Z_{s_1}, \dots, Z_{s_2})] \\ &= \mathbb{E}_D \left[\left(\sum_{i=1}^{s_1} \kappa(x; Z_i, D_{[s_1]}) m(Z_i; \mu, p) \right) \left(\sum_{j=1}^{s_1} \kappa(x; Z_j, D) m(Z_j; \mu, p) + \sum_{j=s_1+1}^{s_2} \kappa(x; Z_j, D) m(Z_j; \mu, p) \right) \right] \\ &= \mathbb{E}_D \left[\underbrace{\sum_{i=1}^{s_1} \sum_{j=1}^{s_1} \kappa(x; Z_i, D_{[s_1]}) \kappa(x; Z_j, D) m(Z_i; \mu, p) m(Z_j; \mu, p)}_{(A)} \right] \\ &\quad + \mathbb{E}_D \left[\underbrace{\sum_{i=1}^{s_1} \sum_{j=s_1+1}^{s_2} \kappa(x; Z_i, D_{[s_1]}) \kappa(x; Z_j, D) m(Z_i; \mu, p) m(Z_j; \mu, p)}_{(B)} \right] \end{aligned} \quad (\text{D.54})$$

Using this decomposition, we can make the following findings.

$$\begin{aligned} (A) &= \mathbb{E}_D \left[\sum_{i=1}^{s_1} \sum_{j=1}^{s_1} \kappa(x; Z_i, D_{[s_1]}) \kappa(x; Z_j, D) m(Z_i; \mu, p) m(Z_j; \mu, p) \right] \\ &= \mathbb{E}_D \left[\sum_{i=1}^{s_1} \kappa(x; Z_i, D_{[s_1]}) m^2(Z_i; \mu, p) \right] = \mathbb{E}_D [m^2(Z_1; \mu, p) s_1 \kappa(x; Z_1, D_{[s_1]})] \\ &= \mathbb{E}_1 [m^2(Z_1; \mu, p) s_1 \mathbb{E}_{2:s_2} [\kappa(x; Z_1, D_{[s_1]})]] = \mathbb{E}_1 [m^2(Z_1; \mu, p) s_1 \mathbb{E}_{2:s_1} [\kappa(x; Z_1, D_{[s_1]})]] \\ &\stackrel{(\text{Lem B.3})}{\lesssim} (\mu_1(x) - \mu_0(x))^2 + \frac{\bar{\sigma}_\varepsilon^2}{\mathfrak{p}(1 - \mathfrak{p})} + o(1) \end{aligned} \quad (\text{D.55})$$

$$\begin{aligned} (B) &= \mathbb{E}_D \left[\sum_{i=1}^{s_1} \sum_{j=s_1+1}^{s_2} \kappa(x; Z_i, D_{[s_1]}) \kappa(x; Z_j, D) m(Z_i; \mu, p) m(Z_j; \mu, p) \right] \\ &= \mathbb{E}_D [s_1(s_2 - s_1) \kappa(x; Z_1, D_{[s_1]}) \kappa(x; Z_{s_2}, D) m(Z_1; \mu, p) m(Z_{s_2}; \mu, p)] \\ &\leq \frac{(s_2 - s_1)}{s_2} \mathbb{E}_D [|m(Z_1; \mu, p)| s_1 \kappa(x; Z_1, D_{[s_1]})] \mathbb{E}_D [|m(Z_{s_2}; \mu, p)| s_2 \kappa(x; Z_{s_2}, D)] \\ &\lesssim \frac{(s_2 - s_1)}{s_2} (\mu_1(x) - \mu_0(x))^2 + o(1) \end{aligned} \quad (\text{D.56})$$

Thus, we obtain the desired result.

$$\Upsilon_{s_1, s_2}(x) = (A) + (B) \lesssim 2(\mu_1(x) - \mu_0(x))^2 + \frac{\bar{\sigma}_\varepsilon^2}{\mathfrak{p}(1-\mathfrak{p})} + o(1) \quad (\text{D.57})$$

■

Lemma D.14.

Let $D = \{Z_1, \dots, Z_{s_2}\}$ be a vector of i.i.d. random variables drawn from Q for $s_2 > s_1$. Let $D' = \{Z_1, \dots, Z_c, Z'_{c+1}, \dots, Z'_{s_1}\}$ where $Z'_{c+1}, \dots, Z'_{s_1}$ are i.i.d. draws from P that are independent of D . Furthermore, let

$$\Upsilon_{s_1, s_2}^c(x) = \mathbb{E} \left[h_{s_1}(x; Z_1, \dots, Z_c, Z'_{c+1}, \dots, Z'_{s_1}) \cdot h_{s_2}(x; Z_1, \dots, Z_{s_2}) \right]. \quad (\text{D.58})$$

Then,

$$\Upsilon_{s_1, s_2}^c(x) \lesssim 4(\mu_1(x) - \mu_0(x))^2 + \frac{\sigma_\varepsilon^2(x)}{\mathfrak{p}(1-\mathfrak{p})} + o(1) \quad (\text{D.59})$$

for s_1, s_2 sufficiently large with $0 < \mathfrak{c} \leq s_1/s_2 \leq 1 - \mathfrak{c} < 1$.

Proof of Lemma D.14.

$$\begin{aligned} \Upsilon_{s_1, s_2}^c(x) &= \mathbb{E} \left[h_{s_1}(x; Z_1, \dots, Z_c, Z'_{c+1}, \dots, Z'_{s_1}) \cdot h_{s_2}(x; Z_1, \dots, Z_{s_2}) \right] \\ &= \mathbb{E}_{D, D'} \left[\left(\sum_{i=1}^c \kappa(x; Z_i, D'_{[s_1]}) m(Z_i; \mu, p) + \sum_{i=c+1}^{s_1} \kappa(x; Z'_i, D'_{[s_1]}) m(Z'_i; \mu, p) \right) \left(\sum_{j=1}^{s_2} \kappa(x; Z_j, D) m(Z_j; \mu, p) \right) \right] \\ &= \mathbb{E}_{D, D'} \left[\underbrace{\sum_{i=1}^c \sum_{j=1}^c \kappa(x; Z_i, D'_{[s_1]}) \kappa(x; Z_j, D) m(Z_i; \mu, p) m(Z_j; \mu, p)}_{(A)} \right] \\ &\quad + \mathbb{E}_{D, D'} \left[\underbrace{\sum_{i=1}^c \sum_{j=c+1}^{s_2} \kappa(x; Z_i, D'_{[s_1]}) \kappa(x; Z_j, D) m(Z_i; \mu, p) m(Z_j; \mu, p)}_{(B)} \right] \\ &\quad + \mathbb{E}_{D, D'} \left[\underbrace{\sum_{i=c+1}^{s_1} \sum_{j=1}^c \kappa(x; Z'_i, D'_{[s_1]}) \kappa(x; Z_j, D) m(Z'_i; \mu, p) m(Z_j; \mu, p)}_{(C)} \right] \\ &\quad + \mathbb{E}_{D, D'} \left[\underbrace{\sum_{i=c+1}^{s_1} \sum_{j=c+1}^{s_2} \kappa(x; Z'_i, D'_{[s_1]}) \kappa(x; Z_j, D) m(Z'_i; \mu, p) m(Z_j; \mu, p)}_{(D)} \right] \end{aligned} \quad (\text{D.60})$$

Now, considering the terms individually, we find the following.

$$\begin{aligned}
(A) &= \mathbb{E}_{D,D'} \left[c\kappa(x; Z_1, D'_{[s_1]})\kappa(x; Z_1, D)m^2(Z_1; \mu, p) \right] = \frac{c}{s_2} \cdot \mathbb{E}_1 \left[m^2(Z_1; \mu, p) s_2 \mathbb{E}_{2:s_2} \left[\kappa(x; Z_1, D'_{[s_1]})\kappa(x; Z_1, D) \right] \right] \\
&\leq \frac{c}{s_2} \cdot \mathbb{E}_1 \left[m^2(Z_1; \mu, p) s_2 \mathbb{E}_{2:s_2} [\kappa(x; Z_1, D)] \right] \lesssim \frac{c}{s_2} \left((\mu_1(x) - \mu_0(x))^2 + \frac{\sigma_\varepsilon^2(x)}{\mathfrak{p}(1-\mathfrak{p})} \right) + o(1)
\end{aligned} \tag{D.61}$$

Similarly, we find the following.

$$\begin{aligned}
(B) &= \mathbb{E}_{D,D'} \left[c(s_2 - c)\kappa(x; Z_1, D'_{[s_1]})\kappa(x; Z_{c+1}, D)m(Z_1; \mu, p)m(Z_{c+1}; \mu, p) \right] \\
&= \frac{c(s_2 - c)}{s_1 s_2} \mathbb{E}_{D,D'} \left[m(Z_1; \mu, p)m(Z_{c+1}; \mu, p) s_1 s_2 \kappa(x; Z_1, D'_{[s_1]})\kappa(x; Z_{c+1}, D) \right] \\
&\leq \frac{c(s_2 - c)}{s_1 s_2} \mathbb{E}_{D,D'} \left[|m(Z_1; \mu, p)| s_1 \kappa(x; Z_1, D'_{[s_1]}) \right] \mathbb{E}_{D,D'} [|m(Z_{c+1}; \mu, p)| s_2 \kappa(x; Z_{c+1}, D)] \\
&\lesssim \frac{c(s_2 - c)}{s_1 s_2} (\mu_1(x) - \mu_0(x))^2 + o(1)
\end{aligned} \tag{D.62}$$

Applying the same argument to the third term, we find an analogous result.

$$\begin{aligned}
(C) &= \mathbb{E}_{D,D'} \left[(s_1 - c)c\kappa(x; Z'_{c+1}, D'_{[s_1]})\kappa(x; Z_1, D)m(Z'_{c+1}; \mu, p)m(Z_1; \mu, p) \right] \\
&\lesssim \frac{c(s_1 - c)}{s_1 s_2} (\mu_1(x) - \mu_0(x))^2 + o(1)
\end{aligned} \tag{D.63}$$

Finally, for the fourth term, we can make the following observation.

$$\begin{aligned}
(D) &= \mathbb{E}_{D,D'} \left[(s_1 - c)(s_2 - c)\kappa(x; Z'_{c+1}, D'_{[s_1]})\kappa(x; Z_{c+1}, D)m(Z'_{c+1}; \mu, p)m(Z_{c+1}; \mu, p) \right] \\
&= \frac{(s_1 - c)(s_2 - c)}{s_1 s_2} \mathbb{E}_{D,D'} \left[m(Z'_{c+1}; \mu, p)m(Z_{c+1}; \mu, p) s_1 s_2 \kappa(x; Z'_{c+1}, D'_{[s_1]})\kappa(x; Z_{c+1}, D) \right] \\
&\lesssim \frac{(s_1 - c)(s_2 - s_1)}{s_1 s_2} (\mu_1(x) - \mu_0(x))^2 + o(1)
\end{aligned} \tag{D.64}$$

By combining these asymptotic bounds, we find the desired result.

$$\Upsilon_{s_1, s_2}^c(x) = (A) + (B) + (C) + (D) \lesssim 4(\mu_1(x) - \mu_0(x))^2 + \frac{\sigma_\varepsilon^2(x)}{\mathfrak{p}(1-\mathfrak{p})} + o(1) \tag{D.65}$$

■

D.6 Variance Estimator Consistency Theorems

Lemma D.15 (Asymptotic Dominance of Hájek Projection).

Let $U_s(\mathbf{D}_{[n]})$ be a non-randomized complete generalized U -statistic with kernel h_s . Let the kernel variance terms ζ_s^s and ζ_s^1 be defined in analogy to Section 3. Assume that the following condition holds.

$$\frac{s}{n} \left(\frac{\zeta_s^s}{s\zeta_s^1} - 1 \right) \rightarrow 0 \quad (\text{D.66})$$

Then, asymptotically, the Hájek projection term dominates the variance of the U -statistic in the following sense.

$$\frac{n}{s^2} \frac{\text{Var}(U_s(\mathbf{D}_{[n]}))}{\zeta_s^1} \rightarrow 1. \quad (\text{D.67})$$

Proof.

$$\begin{aligned} 1 &\leq \frac{n}{s^2} \frac{\text{Var}(U_s(\mathbf{D}_{[n]}))}{\zeta_s^1} = \left(\frac{s^2}{n} \zeta_s^1 \right)^{-1} \sum_{j=1}^s \binom{s}{j}^2 \binom{n}{j}^{-1} V_s^j \\ &\leq 1 + \left(\frac{s^2}{n} \zeta_s^1 \right)^{-1} \frac{s^2}{n^2} \sum_{j=2}^s \binom{s}{j} V_s^j \\ &\leq 1 + \frac{s}{n} \left(\frac{\zeta_s^s}{s\zeta_s^1} - 1 \right) \rightarrow 1. \end{aligned} \quad (\text{D.68})$$

■

Lemma D.16 (Hájek Dominance for TDNN Estimator).

Let $0 < \mathfrak{c} \leq s_1/s_2 \leq 1 - \mathfrak{c} < 1$ and $s_2 = o(n)$, then under Assumptions ??, ?? and ??, then the TDNN estimator fulfills the asymptotic Hájek dominance condition shown in Lemma D.15.

Proof. Recall the results from Lemmas D.9 and D.10.

$$\zeta_{s_1, s_2}^{s_2}(x) \lesssim \mu^2(x) + \sigma_\varepsilon + o(1) \quad \text{and} \quad \zeta_{s_1, s_2}^1(x) \sim s_2^{-1}$$

Using these results, we can find the following.

$$\frac{s_2}{n} \left(\frac{\zeta_{s_1, s_2}^{s_2}(x)}{s_2 \zeta_{s_1, s_2}^1(x)} - 1 \right) \sim \frac{s_2}{n} (\mu^2(x) + \sigma_\varepsilon + o(1) - 1) \sim \frac{s_2}{n} \rightarrow 0 \quad (\text{D.69})$$

■

Proof of Theorem 4.2.

The desired result immediately follows from an application of Theorem 6 from Peng, Mentch, and Stefanski (2021). ■

Proof of Theorem 4.3.

Recall the definition of the Jackknife Variance estimator.

$$\hat{\omega}_{JK}^2(x; \mathbf{D}_n) = \frac{n-1}{n} \sum_{i=1}^n (\hat{\mu}_{s_1, s_2}(x; \mathbf{D}_{n, -i}) - \hat{\mu}_{s_1, s_2}(x; \mathbf{D}_n))^2 \quad (\text{D.70})$$

Using the Hoeffding-decomposition of the original U-statistic, we can reformulate this expression in the following way.

$$\begin{aligned} \hat{\omega}_{JK}^2(x; \mathbf{D}_n) &= \frac{n-1}{n} \sum_{i=1}^n \left(\sum_{j=1}^{s_2} \binom{s_2}{j} H_{s_1, s_2}^j(\mathbf{D}_{n, -i}) - \sum_{j=1}^{s_2} \binom{s_2}{j} H_{s_1, s_2}^j(\mathbf{D}_n) \right)^2 \\ &= \frac{n-1}{n} \sum_{j=1}^n \left(\sum_{j=1}^{s_2} \binom{s_2}{j} (H_{s_1, s_2}^j - H_{s_1, s_2}^j(\mathbf{D}_{n, -i})) \right)^2 \\ &= \frac{n-1}{n} \sum_{j=1}^n \left(\sum_{j=1}^{s_2} \binom{s_2}{j} \left(\binom{n}{j}^{-1} \sum_{\iota \in L_{n, j}} h_{s_1, s_2}^{(j)}(\mathbf{D}_\iota) - \binom{n-1}{j}^{-1} \sum_{\ell \in L_j([n] \setminus \{i\})} h_{s_1, s_2}^{(j)}(\mathbf{D}_\ell) \right) \right)^2 \\ &= \frac{n-1}{n} \sum_{j=1}^n \left[\frac{s_2}{n} h_{s_1, s_2}^{(1)}(Z_i) + \sum_{j \neq i} \left(\frac{s_2}{n} - \frac{s_2}{n-1} \right) h_{s_1, s_2}^{(1)}(Z_j) \right. \\ &\quad \left. + \sum_{j=2}^{s_2} \binom{s_2}{j} \left(\binom{n}{j}^{-1} \sum_{\iota \in L_{n, j}} h_{s_1, s_2}^{(j)}(\mathbf{D}_\iota) - \binom{n-1}{j}^{-1} \sum_{\ell \in L_j([n] \setminus \{i\})} h_{s_1, s_2}^{(j)}(\mathbf{D}_\ell) \right) \right]^2 \\ &= \frac{n-1}{n} \frac{s^2}{n^2} \sum_{j=1}^n \left[h_{s_1, s_2}^{(1)}(Z_i) - \frac{1}{n-1} \sum_{j \neq i} h_{s_1, s_2}^{(1)}(Z_j) \right. \\ &\quad \left. + \frac{n}{s} \sum_{j=2}^{s_2} \binom{s_2}{j} \left(\binom{n}{j}^{-1} \sum_{\iota \in L_{j-1}([n] \setminus \{i\})} h_{s_1, s_2}^{(j)}(\mathbf{D}_{\iota \cup \{i\}}) + \left[\binom{n}{j}^{-1} - \binom{n-1}{j}^{-1} \right] \sum_{\ell \in L_j([n] \setminus \{i\})} h_{s_1, s_2}^{(j)}(\mathbf{D}_\ell) \right) \right] \\ &=: \frac{n-1}{n} \frac{s^2}{n^2} \sum_{j=1}^n [h_{s_1, s_2}^{(1)}(Z_i) + T_i]^2 \end{aligned} \quad (\text{D.71})$$

Observe that due to the independence of the observations and the uncorrelatedness of Hoeffding projections of differing orders, $h_{s_1, s_2}^{(1)}(Z_i)$ and T_i are uncorrelated and both have mean zero. Now, continuing to follow the line of argument in Peng, Mentch, and Stefanski (2021), observe the following.

$$\mathbb{E} \left[\left(h_{s_1, s_2}^{(1)}(Z_i) \right)^2 \right] = V_{s_1, s_2}^1 = \zeta_{s_1, s_2}^1 \quad (\text{D.72})$$

Furthermore, as a consequence of the independence of the observations and the uncorrelatedness of Hoeffding projec-

tions of differing order, we find that

$$\begin{aligned}
\mathbb{E}[T_i^2] &= \frac{1}{n-1} V_{s_1, s_2}^1 + \frac{n^2}{s_2^2} \sum_{j=2}^{s_2} \binom{s_2}{j}^2 \left\{ \binom{n}{j}^{-2} \binom{n-1}{j-1} V_{s_1, s_2}^j + \left[\binom{n}{j}^{-1} - \binom{n-1}{j}^{-1} \right]^2 \binom{n-1}{j} V_{s_1, s_2}^j \right\} \\
&= \frac{1}{n-1} V_{s_1, s_2}^1 + \frac{n^2}{s_2^2} \sum_{j=2}^{s_2} \binom{s_2}{j}^2 \left\{ \binom{n}{j}^{-2} \frac{j}{n-j} \binom{n-1}{j} V_{s_1, s_2}^j + \binom{n}{j}^{-2} \left[1 - \binom{n}{j} \binom{n-1}{j}^{-1} \right]^2 \binom{n-1}{j} V_{s_1, s_2}^j \right\} \\
&= \frac{1}{n-1} V_{s_1, s_2}^1 + \frac{n^2}{s_2^2} \sum_{j=2}^{s_2} \binom{s_2}{j}^2 \binom{n}{j}^{-2} \left(\frac{j}{n-j} + \left(1 - \frac{n}{n-j} \right)^2 \right) \binom{n-1}{j} V_{s_1, s_2}^j \\
&= \frac{1}{n-1} V_{s_1, s_2}^1 + \frac{n^2}{s_2^2} \sum_{j=2}^{s_2} \binom{s_2}{j} \binom{n}{j}^{-2} \binom{n-1}{j} \cdot \left(\frac{j}{n-j} + \frac{j^2}{(n-j)^2} \right) \left[\binom{s_2}{j} V_{s_1, s_2}^j \right] \\
&= \frac{1}{n-1} V_{s_1, s_2}^1 + \frac{n^2}{s_2^2} \sum_{j=2}^{s_2} \binom{s_2}{j} \binom{n}{j}^{-1} \frac{n-j}{n} \cdot \frac{j}{n} \left(\frac{n}{n-j} + \frac{jn}{(n-j)^2} \right) \left[\binom{s_2}{j} V_{s_1, s_2}^j \right] \\
&= \frac{1}{n-1} V_{s_1, s_2}^1 + \sum_{j=2}^{s_2} \frac{j}{s_2} \binom{s_2-1}{j-1} \binom{n-1}{j-1}^{-1} \frac{n-j}{n} \left(\frac{n}{n-j} + \frac{j}{n} \right) \left[\binom{s_2}{j} V_{s_1, s_2}^j \right] \\
&\leq \frac{1}{n-1} V_{s_1, s_2}^1 + \sum_{j=2}^{s_2} \frac{j}{s_2} \left(e \frac{s_2-1}{n-1} \right)^{j-1} \frac{n-j}{n} \left(\frac{n}{n-j} + \frac{j}{n} \right) \left[\binom{s_2}{j} V_{s_1, s_2}^j \right] \\
&\lesssim \frac{1}{n-1} V_{s_1, s_2}^1 + 2 \sum_{j=2}^{s_2} \frac{j}{s_2} \left(e \frac{s_2-1}{n-1} \right)^{j-1} \left[\binom{s_2}{j} V_{s_1, s_2}^j \right] \\
&\leq \frac{1}{n-1} V_{s_1, s_2}^1 + 2e \sum_{j=2}^{s_2} \frac{1}{s_2} \frac{s_2-1}{n-1} \left[\binom{s_2}{j} V_{s_1, s_2}^j \right] + 2 \sum_{j=2}^{s_2} \frac{j-1}{s_2} \left(e \frac{s_2-1}{n-1} \right)^{j-1} \left[\binom{s_2}{j} V_{s_1, s_2}^j \right] \\
&\leq \frac{1}{n-1} V_{s_1, s_2}^1 + \frac{2e}{n-1} \sum_{j=2}^{s_2} \frac{s_2-1}{s_2} \left[\binom{s_2}{j} V_{s_1, s_2}^j \right] + 2 \sum_{j=2}^{s_2} \frac{j-1}{s_2} \left(e \frac{s_2-1}{n-1} \right)^{j-1} \zeta_{s_1, s_2}^{s_2} \\
&= \frac{1}{n-1} V_{s_1, s_2}^1 + \frac{2e}{n} \sum_{j=2}^{s_2} \frac{n(s_2-1)}{(n-1)s_2} \left[\binom{s_2}{j} V_{s_1, s_2}^j \right] + 2\zeta_{s_1, s_2}^{s_2} \sum_{j=2}^{s_2} \frac{j-1}{s_2} \left(e \frac{s_2-1}{n-1} \right)^{j-1} \\
&\leq \frac{1}{n-1} V_{s_1, s_2}^1 + \frac{2e}{n} \sum_{j=2}^{s_2} \binom{s_2}{j} V_{s_1, s_2}^j + \frac{2\zeta_{s_1, s_2}^{s_2}}{s_2} \sum_{j=1}^{\infty} j \left(e \frac{s_2-1}{n-1} \right)^j \\
&\leq \frac{1}{n-1} V_{s_1, s_2}^1 + \frac{2e}{n} \sum_{j=2}^{s_2} \binom{s_2}{j} V_{s_1, s_2}^j + \frac{2\zeta_{s_1, s_2}^{s_2}}{s_2} \sum_{j=1}^{\infty} j \left(e \frac{s_2}{n} \right)^j \\
&= \frac{1}{n-1} \zeta_{s_1, s_2}^1 + \frac{2e}{n} (\zeta_{s_1, s_2}^{s_2} - s_2 \zeta_{s_1, s_2}^1) + \frac{2en}{(n-es_2)^2} \zeta_{s_1, s_2}^{s_2} \\
&= \left(\frac{1}{n-1} + \frac{2es_2n}{(n-es_2)^2} \right) \zeta_{s_1, s_2}^1 + 2e \left(\frac{1}{n} + \frac{n}{(n-es_2)^2} \right) (\zeta_{s_1, s_2}^{s_2} - s_2 \zeta_{s_1, s_2}^1)
\end{aligned} \tag{D.73}$$

Recall the results of Lemmas D.9 and D.10.

$$\zeta_{s_1, s_2}^{s_2}(x) \lesssim \mu^2(x) + \sigma_\varepsilon + o(1) \quad \text{and} \quad \zeta_{s_1, s_2}^1(x) \sim s_2^{-1} \tag{D.74}$$

This immediately implies that $\frac{s_2}{n} \left(\frac{\zeta_{s_1, s_2}^{s_2}}{s_2 \zeta_{s_1, s_2}^1} - 1 \right) \rightarrow 0$. Using this result and the previous asymptotic upper bound, we

can find the following.

$$\begin{aligned}
\frac{\mathbb{E}[T_i^2]}{V_{s_1, s_2}^1} &\leq \frac{\left(\frac{1}{n-1} + \frac{2es_2n}{(n-es_2)^2}\right) \zeta_{s_1, s_2}^1 + 2e\left(\frac{1}{n} + \frac{n}{(n-es_2)^2}\right) (\zeta_{s_1, s_2}^{s_2} - s_2 \zeta_{s_1, s_2}^1)}{\zeta_{s_1, s_2}^1} \\
&= \frac{1}{n-1} + \frac{2es_2n}{(n-es_2)^2} + 2e\left(\frac{1}{n} + \frac{n}{(n-es_2)^2}\right) \left(\frac{\zeta_{s_1, s_2}^{s_2} - s_2 \zeta_{s_1, s_2}^1}{\zeta_{s_1, s_2}^1}\right) \rightarrow 0
\end{aligned} \tag{D.75}$$

Therefore, we can conclude that $h_s^{(1)}(Z_i)$ dominates T_i^2 in the expression of interest. Using Lemma B.8, we can thus conclude the following.

$$\begin{aligned}
\frac{\frac{n}{s_2^2} \hat{\omega}_{JK}^2(x; \mathbf{D}_n)}{V_{s_1, s_2}^1(x)} &\rightarrow_p \frac{n-1}{n} \frac{1}{n} \sum_{i=1}^n \frac{\left(h_{s_1, s_2}^{(1)}(x; Z_i)\right)^2}{V_{s_1, s_2}^1(x)} \\
&\rightarrow_p \frac{n-1}{n} \frac{\mathbb{E}\left[\left(h_{s_1, s_2}^{(1)}(x; Z_i)\right)^2\right]}{V_{s_1, s_2}^1(x)} \rightarrow 1
\end{aligned} \tag{D.76}$$

The desired rate-consistency then immediately follows from an application of Lemma D.15. ■

Proof of Theorem 4.4.

Consider first the case absent additional randomization in the form of ω and recall the definition of the delete-d Jackknife Variance estimator.

$$\hat{\omega}_{JKD}^2(x; d, \mathbf{D}_n) = \frac{n-d}{d} \binom{n}{d}^{-1} \sum_{\ell \in L_{n,d}} (\hat{\mu}_{s_1, s_2}(x; \mathbf{D}_{n,-\ell}) - \hat{\mu}_{s_1, s_2}(x; \mathbf{D}_n))^2 \quad (\text{D.77})$$

Now, as in the proof for the conventional Jackknife variance estimator, we make use of the Hoeffding-decomposition in the following way.

$$\begin{aligned} \hat{\omega}_{JKD}^2(x; d, \mathbf{D}_n) &= \frac{n-d}{d} \binom{n}{d}^{-1} \sum_{\ell \in L_{n,d}} \left(\sum_{j=1}^{s_2} \binom{s_2}{j} \left(H_{P_t}^j - H_{P_t}^j(\mathbf{D}_{n,-\ell}) \right) \right)^2 \\ &= \frac{n-d}{d} \binom{n}{d}^{-1} \sum_{\ell \in L_{n,d}} \left(\sum_{j=1}^{s_2} \binom{s_2}{j} \left(\binom{n}{j}^{-1} \sum_{\iota \in L_{n,j}} h_{s_1, s_2}^{(j)}(\mathbf{D}_\iota) - \binom{n-d}{j}^{-1} \sum_{\iota \in L_j([n] \setminus \ell)} h_{s_1, s_2}^{(j)}(\mathbf{D}_\iota) \right) \right)^2 \\ &= \frac{n-d}{d} \binom{n}{d}^{-1} \sum_{\ell \in L_{n,d}} \left[\frac{s_2}{n} \sum_{i \in \ell} h_{s_1, s_2}^{(1)}(Z_i) + \sum_{i \in [n] \setminus \ell} \left(\frac{s_2}{n} - \frac{s_2}{n-d} \right) h_{s_1, s_2}^{(1)}(Z_i) \right. \\ &\quad \left. + \sum_{j=2}^{s_2} \binom{s_2}{j} \left(\binom{n}{j}^{-1} \sum_{\iota \in L_{n,j}} h_{s_1, s_2}^{(j)}(\mathbf{D}_\iota) - \binom{n-d}{j}^{-1} \sum_{\iota \in L_j([n] \setminus \ell)} h_{s_1, s_2}^{(j)}(\mathbf{D}_\iota) \right) \right]^2 \\ &= \frac{n-d}{d} \binom{n}{d}^{-1} \left(\frac{s_2}{n} \right)^2 \sum_{\ell \in L_{n,d}} \left[\sum_{i \in \ell} h_{s_1, s_2}^{(1)}(Z_i) - \frac{d}{n-d} \sum_{i \in [n] \setminus \ell} h_{s_1, s_2}^{(1)}(Z_i) \right. \\ &\quad \left. + \frac{n}{s_2} \sum_{j=2}^{s_2} \binom{s_2}{j} \left(\binom{n}{j}^{-1} \sum_{\iota \in L_{n,j}} h_{s_1, s_2}^{(j)}(\mathbf{D}_\iota) - \binom{n-d}{j}^{-1} \sum_{\iota \in L_j([n] \setminus \ell)} h_{s_1, s_2}^{(j)}(\mathbf{D}_\iota) \right) \right]^2 \\ &=: (n-d) \binom{n}{d}^{-1} \left(\frac{s_2}{n} \right)^2 \sum_{\ell \in L_{n,d}} \left[\frac{1}{\sqrt{d}} \sum_{i \in \ell} h_{s_1, s_2}^{(1)}(Z_i) + T_\ell \right]^2 \end{aligned} \quad (\text{D.78})$$

We want to proceed in an analogous way to the proof of the pure Jackknife result. Thus, we want to show that $\sum_{i \in \ell} h_{s_1, s_2}^{(1)}(Z_i)$ dominates T_ℓ in the sense of Lemma B.8. Luckily, since Lemma B.8 does not depend on any particular independence assumptions of summands etc. this is a relatively straightforward adaptation of the strategy shown in the proof of Theorem 4.3. Thus, consider the following for an arbitrary fixed index-subset ℓ with cardinality d .

$$\begin{aligned} \mathbb{E} \left[\left(\frac{1}{\sqrt{d}} \sum_{i \in \ell} h_{s_1, s_2}^{(1)}(Z_i) \right)^2 \right] &= \frac{1}{d} \mathbb{E} \left[\sum_{i \in \ell} \sum_{j \in \ell} h_{s_1, s_2}^{(1)}(Z_i) h_{s_1, s_2}^{(1)}(Z_j) \right] = \frac{1}{d} \sum_{i \in \ell} \sum_{j \in \ell} \mathbb{E} \left[h_{s_1, s_2}^{(1)}(Z_i) h_{s_1, s_2}^{(1)}(Z_j) \right] \\ &= \frac{|\ell|}{d} \cdot \mathbb{E} \left[\left(h_{s_1, s_2}^{(1)}(Z_1) \right)^2 \right] = \zeta_{P_t, 1} \end{aligned} \quad (\text{D.79})$$

For the error term we introduce a case distinction. Case one corresponds to parameter choices where $s_2 \geq d$ and thus takes the following form.

$$\begin{aligned}
T_\ell &= \frac{\sqrt{d}}{n-d} \sum_{i \in [n] \setminus \ell} h_{s_1, s_2}^{(1)}(Z_i) \\
&+ \frac{n}{s_2 \sqrt{d}} \left\{ \sum_{j=2}^d \binom{s_2}{j} \left(\binom{n}{j}^{-1} \left(\sum_{a=1}^j \sum_{\substack{\kappa \in L_a(\ell) \\ \varrho \in L_{j-a}([n] \setminus \ell)}} h_{s_1, s_2}^{(j)}(D_{\kappa \cup \varrho}) \right) + \left(\binom{n}{j}^{-1} - \binom{n-d}{j}^{-1} \right) \sum_{\iota \in L_j([n] \setminus \ell)} h_{s_1, s_2}^{(j)}(\mathbf{D}_\iota) \right) \right. \\
&\left. + \sum_{j=d+1}^{s_2} \binom{s_2}{j} \left(\binom{n}{j}^{-1} \left(\sum_{a=1}^d \sum_{\substack{\kappa \in L_a(\ell) \\ \varrho \in L_{j-a}([n] \setminus \ell)}} h_{s_1, s_2}^{(j)}(D_{\kappa \cup \varrho}) \right) + \left(\binom{n}{j}^{-1} - \binom{n-d}{j}^{-1} \right) \sum_{\iota \in L_j([n] \setminus \ell)} h_{s_1, s_2}^{(j)}(\mathbf{D}_\iota) \right) \right\} \quad (\text{D.80})
\end{aligned}$$

Case two covers setups of the form $s_2 < d$ and thus takes the following form.

$$\begin{aligned}
T_\ell &= \frac{\sqrt{d}}{n-d} \sum_{i \in [n] \setminus \ell} h_{s_1, s_2}^{(1)}(Z_i) \\
&+ \frac{n}{s_2 \sqrt{d}} \sum_{j=2}^{s_2} \binom{s_2}{j} \left(\binom{n}{j}^{-1} \left(\sum_{a=1}^j \sum_{\substack{\kappa \in L_a(\ell) \\ \varrho \in L_{j-a}([n] \setminus \ell)}} h_{s_1, s_2}^{(j)}(D_{\kappa \cup \varrho}) \right) + \left(\binom{n}{j}^{-1} - \binom{n-d}{j}^{-1} \right) \sum_{\iota \in L_j([n] \setminus \ell)} h_{s_1, s_2}^{(j)}(\mathbf{D}_\iota) \right) \quad (\text{D.81})
\end{aligned}$$

Having separated these two cases, we continue by investigating the expectation of their respective squares. Beginning with case one, we find the following.

$$\begin{aligned}
\mathbb{E} \left[(T_\ell)^2 \right] &= \frac{d}{n-d} V_{s_1, s_2}^1 \\
&+ \frac{n^2}{s_2^2 d} \sum_{j=2}^d \binom{s_2}{j}^2 \left(\binom{n}{j}^{-2} \sum_{a=1}^j \left[\binom{d}{a} \binom{n-d}{j-a} \right] + \left[\binom{n}{j}^{-1} - \binom{n-d}{j}^{-1} \right]^2 \binom{n-d}{j} \right) V_{s_1, s_2}^j \\
&+ \frac{n^2}{s_2^2 d} \sum_{j=d+1}^{s_2} \binom{s_2}{j}^2 \left(\binom{n}{j}^{-2} \sum_{a=1}^d \left[\binom{d}{a} \binom{n-d}{j-a} \right] + \left[\binom{n}{j}^{-1} - \binom{n-d}{j}^{-1} \right]^2 \binom{n-d}{j} \right) V_{s_1, s_2}^j \\
&\stackrel{(\star)}{=} \frac{d}{n-d} V_{s_1, s_2}^1 \\
&+ \frac{n^2}{s_2^2 d} \sum_{j=2}^d \binom{s_2}{j}^2 \binom{n}{j}^{-2} \left(\binom{n}{j} - \binom{n-d}{j} + \left[1 - \binom{n}{j} \binom{n-d}{j}^{-1} \right]^2 \binom{n-d}{j} \right) V_{s_1, s_2}^j \\
&+ \frac{n}{s_2 d} \sum_{j=d+1}^{s_2} \frac{\binom{s_2-1}{j-1} \binom{n-d}{j}}{\binom{n-1}{j-1} \binom{n}{j}} \left(\sum_{a=1}^d \frac{\binom{d}{a} \binom{n-d}{j-a}}{\binom{n-d}{j}} + \left[1 - \binom{n}{j} \binom{n-d}{j}^{-1} \right]^2 \right) \binom{s_2}{j} V_{s_1, s_2}^j \quad (\text{D.82})
\end{aligned}$$

The equality marked by (\star) holds by the Chu-Vandermonde identity - specifically with respect to the equivalent expression for the sum in the second term.

Continuing the analysis, we find the following.

$$\begin{aligned}
\mathbb{E} \left[(T_\ell)^2 \right] &= \frac{d}{n-d} V_{s_1, s_2}^1 \\
&+ \frac{n}{s_2 d} \sum_{j=2}^d \frac{\binom{s_2-1}{j-1} \binom{n-d}{j}}{\binom{n-1}{j-1} \binom{n}{j}} \left(\binom{n}{j} \binom{n-d}{j}^{-1} - 1 + \left[1 - \binom{n}{j} \binom{n-d}{j}^{-1} \right]^2 \right) \left[\binom{s_2}{j} V_{s_1, s_2}^j \right] \\
&+ \frac{n}{s_2 d} \sum_{j=d+1}^{s_2} \frac{\binom{s_2-1}{j-1} \binom{n-d}{j}}{\binom{n-1}{j-1} \binom{n}{j}} \left(\frac{\binom{n}{j}}{\binom{n-d}{j}} \sum_{a=1}^d \frac{\binom{d}{a} \binom{n-d}{j-a}}{\binom{n}{j}} + \left[1 - \binom{n}{j} \binom{n-d}{j}^{-1} \right]^2 \right) \left[\binom{s_2}{j} V_{s_1, s_2}^j \right] \\
&= \frac{d}{n-d} V_{s_1, s_2}^1 \\
&+ \frac{n}{s_2 d} \sum_{j=2}^d \frac{\binom{s_2-1}{j-1} \binom{n-d}{j}}{\binom{n-1}{j-1} \binom{n}{j}} \left(\binom{n}{j}^2 \binom{n-d}{j}^{-2} - \binom{n}{j} \binom{n-d}{j}^{-1} \right) \left[\binom{s_2}{j} V_{s_1, s_2}^j \right] \\
&+ \frac{n}{s_2 d} \sum_{j=d+1}^{s_2} \frac{\binom{s_2-1}{j-1} \binom{n-d}{j}}{\binom{n-1}{j-1} \binom{n}{j}} \left(\frac{\binom{n}{j}}{\binom{n-d}{j} \binom{n}{d}} \sum_{a=1}^d \binom{j}{a} \binom{n-j}{d-a} + \left[1 - \binom{n}{j} \binom{n-d}{j}^{-1} \right]^2 \right) \left[\binom{s_2}{j} V_{s_1, s_2}^j \right] \\
&\stackrel{(\star\star)}{=} \frac{d}{n-d} V_{s_1, s_2}^1 \\
&+ \frac{n}{s_2 d} \sum_{j=2}^d \frac{\binom{s_2-1}{j-1}}{\binom{n-1}{j-1}} \left(\binom{n}{j} \binom{n-d}{j}^{-1} - 1 \right) \left[\binom{s_2}{j} V_{s_1, s_2}^j \right] \\
&+ \frac{n}{s_2 d} \sum_{j=d+1}^{s_2} \frac{\binom{s_2-1}{j-1} \binom{n-d}{j}}{\binom{n-1}{j-1} \binom{n}{j}} \left(\frac{\binom{n}{j}}{\binom{n-d}{j}} \left[1 - \binom{n-j}{d} \binom{n}{d}^{-1} \right] + \left[1 - \binom{n}{j} \binom{n-d}{j}^{-1} \right]^2 \right) \left[\binom{s_2}{j} V_{s_1, s_2}^j \right] \\
&= \frac{d}{n-d} V_{s_1, s_2}^1 + \frac{n}{s_2} \sum_{j=2}^d \frac{\binom{s_2-1}{j-1}}{\binom{n-1}{j-1}} \left(\binom{n}{j} \binom{n-d}{j}^{-1} - 1 \right) \left[\binom{s_2}{j} V_{s_1, s_2}^j \right] \\
&+ \frac{n}{s_2 d} \sum_{j=d+1}^{s_2} \frac{\binom{s_2-1}{j-1} \binom{n-d}{j}}{\binom{n-1}{j-1} \binom{n}{j}} \left(\frac{\binom{n}{j}}{\binom{n-d}{j}} - 1 + \left[1 - \binom{n}{j} \binom{n-d}{j}^{-1} \right]^2 \right) \left[\binom{s_2}{j} V_{s_1, s_2}^j \right] \\
&= \frac{d}{n-d} V_{s_1, s_2}^1 + \frac{n}{s_2 d} \sum_{j=2}^{s_2} \frac{\binom{s_2-1}{j-1}}{\binom{n-1}{j-1}} \left(\binom{n}{j} \binom{n-d}{j}^{-1} - 1 \right) \left[\binom{s_2}{j} V_{s_1, s_2}^j \right] \\
&= \frac{d}{n-d} V_{s_1, s_2}^1 + \frac{n}{s_2 d} \sum_{j=2}^{s_2} \frac{\binom{s_2-1}{j-1}}{\binom{n-1}{j-1}} \left(\prod_{i=0}^{d-1} \left(1 + \frac{j}{n-i-j} \right) - 1 \right) \left[\binom{s_2}{j} V_{s_1, s_2}^j \right] \\
&\stackrel{(\star\star\star)}{\leq} \frac{d}{n-d} V_{s_1, s_2}^1 + \frac{n}{s_2 d} \sum_{j=2}^{s_2} \frac{\binom{s_2-1}{j-1}}{\binom{n-1}{j-1}} \frac{\sum_{i=0}^{d-1} \frac{j}{n-i-j}}{1 - \sum_{i=0}^{d-1} \frac{j}{n-i-j}} \left[\binom{s_2}{j} V_{s_1, s_2}^j \right] \\
&\leq \frac{d}{n-d} V_{s_1, s_2}^1 + \frac{n}{s_2 d} \sum_{j=2}^{s_2} \frac{\binom{s_2-1}{j-1}}{\binom{n-1}{j-1}} \frac{j(n-j)}{(n-d-j+1)(n-d-2j)} \left[\binom{s_2}{j} V_{s_1, s_2}^j \right] \\
&\leq \frac{d}{n-d} V_{s_1, s_2}^1 + \frac{n}{s_2 d} \sum_{j=2}^{s_2} \left(\frac{e(s_2-1)}{n-1} \right)^{j-1} \frac{j(n-2)}{(n-d-s_2+1)(n-d-2s_t)} \left[\binom{s_2}{j} V_{s_1, s_2}^j \right]
\end{aligned} \tag{D.83}$$

The equality marked by $(\star\star)$ holds by the Chu-Vandermonde identity applied to the third summand, whereas the inequality marked by the equality marked by $(\star\star\star)$ follows from a Weierstrass-Product type inequality. Furthermore, this derivation shows that we do not really need to distinguish between the two described cases for the error term.

Proceeding this way allows us to continue our analysis similar to the proof for the simple leave-one-out Jackknife.

$$\begin{aligned}
\mathbb{E} \left[(T_\ell)^2 \right] &\lesssim \frac{d}{n-d} V_{s_1, s_2}^1 + \frac{n}{s_2 d} \sum_{j=2}^{s_2} \left(\frac{e(s_2-1)}{n-1} \right)^{j-1} \frac{j(n-2)}{(n-d-s_2+1)(n-d-2s_t)} \left[\binom{s_2}{j} V_{s_1, s_2}^j \right] \\
&= \frac{d}{n-d} V_{s_1, s_2}^1 + \frac{2e \cdot n(n-2)}{(n-1)(n-d-s_2+1)(n-d-2s_t)d} \sum_{j=2}^{s_2} \frac{j}{s_2} \left(\frac{e(s_2-1)}{n-1} \right)^{j-1} \left[\binom{s_2}{j} V_{s_1, s_2}^j \right] \\
&\lesssim \frac{d}{n-d} V_{s_1, s_2}^1 + \frac{4e}{(n-d-s_2)s_2 d} \sum_{j=2}^{s_2} j \left(\frac{e(s_2-1)}{n-1} \right)^{j-1} \left[\binom{s_2}{j} V_{s_1, s_2}^j \right] \\
&\leq \frac{d}{n-d} V_{s_1, s_2}^1 + \frac{4e}{(n-d-s_2)s_2 d} \sum_{j=2}^{s_2} \left[\binom{s_2}{j} V_{s_1, s_2}^j \right] + \frac{4e}{(n-d-s_2)s_2 d} \sum_{j=2}^{s_2} (j-1) \left(\frac{e(s_2-1)}{n-1} \right)^{j-1} \left[\binom{s_2}{j} V_{s_1, s_2}^j \right] \\
&\leq \frac{d}{n-d} V_{s_1, s_2}^1 + \frac{4e}{(n-d-s_2)s_2 d} \sum_{j=2}^{s_2} \left[\binom{s_2}{j} V_{s_1, s_2}^j \right] + \frac{4e \cdot \zeta_{s_1, s_2}^{s_2}}{(n-d-s_2)s_2 d} \sum_{j=1}^{\infty} j \left(\frac{e(s_2-1)}{n-1} \right)^j \\
&= \frac{d}{n-d} \zeta_{s_1, s_2}^1 + \frac{4e}{(n-d-s_2)s_2 d} (\zeta_{s_1, s_2}^{s_2} - s_2 \zeta_{s_1, s_2}^1) + \frac{4e \cdot \zeta_{s_1, s_2}^{s_2}}{(n-d-s_2)s_2 d} \cdot \frac{e(s_2-1)(n-1)}{(n-1-e(s_2-1))^2} \\
&= \left(\frac{d}{n-d} + \frac{e(s_2-1)(n-1)}{(n-1-e(s_2-1))^2} \right) \zeta_{s_1, s_2}^1 + \frac{4e}{(n-d-s_2)s_2 d} \left(1 + \frac{e(s_2-1)(n-1)}{(n-1-e(s_2-1))^2} \right) (\zeta_{s_1, s_2}^{s_2} - s_2 \zeta_{s_1, s_2}^1)
\end{aligned} \tag{D.84}$$

We continue as in the default Jackknife case.

$$\begin{aligned}
\frac{\mathbb{E} [T_\ell^2]}{V_{s_1, s_2}^1} &\leq \frac{d}{n-d} + \frac{e(s_2-1)(n-1)}{(n-1-e(s_2-1))^2} + \frac{4e}{(n-d-s_2)s_2 d} \left(1 + \frac{e(s_2-1)(n-1)}{(n-1-e(s_2-1))^2} \right) \frac{\zeta_{s_1, s_2}^{s_2} - s_2 \zeta_{s_1, s_2}^1}{\zeta_{s_1, s_2}^1} \\
&\rightarrow 0.
\end{aligned} \tag{D.85}$$

Now, following the exact same logic as in the proof for the consistency of the Jackknife variance estimator, we obtain consistency of the delete-d Jackknife variance estimator. ■

D.7 Pointwise Inference Results

Proof of Theorem 4.6.

■

Proof of Theorem 4.7.

■

E Proofs for Results in Section 5
