

# Inference for Conditional Average Treatment Effects using Distributional Nearest Neighbors

Jakob R. Juergens  
University of Wisconsin - Madison

Last edited: February 22, 2025

---

## Abstract

This paper presents a computationally simple method for estimating heterogeneous treatment effects based on a weighted nearest-neighbor-type estimator. As part of this analysis, for a class of generalized U-statistics, I improve on conditions currently available in the literature required for consistent variance estimation. Furthermore, I provide corresponding results for asymptotically valid pointwise inference in a nonparametric regression setup. I extend the ideas of this nonparametric regression estimator to the estimation of conditional average treatment effects using ideas from double / debiased machine learning including Neyman-orthogonal moments and cross-fitting. This effectively results in a weighted-nearest-nearest-based DDML variant that is new to the literature and resembles recent Kernel-based approaches under a data-adaptive kernel.

---

Supplementary Material and R-Package available at: [https://github.com/JakobJuergens/Unif\\_Inf\\_TDNN](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. However, the development of inferential theory for these estimators 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 are twofold. First, this paper provides extended consistency results for the variance estimators 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, I 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 be at the heart of future iterations of this paper.

After short sections on notation and literature review, 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). To give readers a more economic understanding of the ideas presented in this paper, I will also introduce a running example that I will refer to throughout this paper. This running example will be a simplified version of the well-known problem of estimating treatment effects in a job-training program for individuals of differing characteristics. Furthermore, this section introduces and contextualizes most of the assumptions that I refer to at later stages of the paper. Section 3 is used to define the DNN and TDNN estimators in the context of nonparametric regression and introduces a novel estimator for the CATE setup. In addition, the main results of the distributional approximations of the estimators are introduced. Section 4 further embeds the DNN estimators and their CATE derivatives in the context of U-statistics and introduces essential notation for the analysis of the estimators. Section 5 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. The future iterations of this paper will then tackle the problem of simultaneous inference using techniques developed by Ritzwoller and Syrgkanis (2024). 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 that mimics an economic analysis. 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  $\ell$  have been removed. To simplify the notation in the case that a single observation (say the  $i$ 'th observation) is removed, I use the notation  $\mathbf{D}_{n,-i}$ . Similarly, given such a data set  $\mathbf{D}_{[n]}$  and index vector  $\ell$ ,

denote by  $\mathbf{D}_\ell$  the data set consisting only of the observations in  $\mathbf{D}_{[n]}$  corresponding to the indices in  $\ell$ . In an abuse of notation, when considering two index vectors  $\ell$  and  $\iota$  that do not share any entries, I 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 weak convergence, 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 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 particular interest in the context of this paper are the so-called “Weighted Nearest-Neighbor” methods for nonparametric regression. While this is a well-studied type of estimator in and of itself, I draw particular connections to bagged-nearest-neighbor type estimators. This class of estimators is built on the framework of “potential closest neighbors” as introduced by Lin and Jeon (2006). Relevant papers studying their properties are, among others, Biau, Cérou, and Guyader (2010), Biau and Devroye (2010), and Steele (2009). These papers also 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.

U-statistics were 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 that studies 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 been previously established in the literature. While being a relatively novel development, there is a significant body of literature concerning infinite-order U-statistics, which 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 that topic. Thus, articles such as Wager, Hastie, and Efron (2014) and Wager and Athey (2018) are of special interest, especially since causal forests are considered the state-of-the-art technique for estimating CATEs.

In the context of estimation and inference regarding CATEs using DDML, Victor Chernozhukov, Chetverikov, et al. (2018) should be pointed out first. By combining cross-fitting 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.

Several extensions to this highly influential idea have been proposed, some of which explicitly aim at estimating CATEs. An example is Semenova and Victor Chernozhukov (2021), who develops estimation and inference procedures for the best linear predictor of a class of causal functions that contain the CATE. Following a different approach, Victor Chernozhukov, Whitney K. Newey, and Syrgkanis (2024) introduce the concept of conditional influence functions and develop a Kernel-based method that similarly to the paper at hand aims at the nonparametric estimation of causal parameters such as the CATE. In a similar vein, Chernozhukov, W K Newey, and Singh (2022) is highly relevant, as it provides a very general analysis of DDML as a meta-algorithm, covering the estimation and inference for CATE.

## 2 Setup

Throughout this paper, I will consider two distinct setups focusing on separate but intertwined problems. As a running example to give immediate economic meaning to the statistical problems, I consider the problem of evaluating a job training program given only observational data in the style of LaLonde (1986).

**Example 1** (Average Hourly Earnings and IT Training).

*Imagine obtaining a large data set containing hourly earnings  $Y$  of workers with age  $X_1$  and  $X_2$  years of education. A typical quantity of interest is the average hourly earnings given a specific combination of worker characteristics, for example the average earnings for a 20-year-old worker with 12 years of education (a recent high school graduate without further education).*

To simplify the presentation, the running example ignores many widely discussed economic problems in this environment and instead focuses on a highly stylized problem with a limited set of characteristics. The first statistical setup is a pure nonparametric regression setup that closely mirrors the structure of Demirkaya et al. (2024) and will be useful to illustrate the inner workings of the estimator of interest.

**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  $Y \in \mathcal{Y} \subset \mathbb{R}$  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  $\varphi$  on  $\mathcal{X}$ , and  $\mu(x)$  is the unknown mean regression function.  $\varepsilon$  is the unobservable model error on which I 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  $Z_i = (X_i, Y_i) \stackrel{iid}{\sim} P$ .*

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

**Example 1** (Average Hourly Earnings and IT Training - continued).

*Consider the introduction of a job training program that covers basic IT skills that potentially influence the average earnings of a worker with given characteristics. Think of providing said IT training to a 20-year-old high-school graduate without further education or a 56-year-old with a PhD. It seems unreasonable to expect the effect of basic IT training on hourly wages to be similar in these cases. Furthermore, it seems unlikely for these participants to choose to participate in the training with equal probability, adding an additional complicating factor.*

**Assumption 2** (Heterogeneous Treatment Effect DGP).

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

$$\begin{aligned} \mathbf{D}_n &= \{Z_i = (X_i, W_i, Y_i)\}_{i=1}^n \quad \text{from the model} \quad Y = \mathbb{1}(W = 0)\mu_0^0(X) + \mathbb{1}(W = 1)\mu_0^1(X) + \varepsilon, \\ W_i &\sim \text{Bern}(\pi_0(X_i)) \end{aligned} \tag{2.3}$$

where  $Y \in \mathcal{Y} \subset \mathbb{R}$  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 an associated probability measure  $\varphi$  on  $\mathcal{X}$  and  $\varepsilon$  is the unobservable model error on which I impose the following conditions.

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

Furthermore,  $\mu_0^0 : \mathcal{X} \rightarrow \mathbb{R}$  and  $\mu_0^1 : \mathcal{X} \rightarrow \mathbb{R}$  are the two unknown potential outcome functions and  $\pi_0 : \mathcal{X} \rightarrow [0, 1]$  is a function describing the probability of treatment uptake, effectively corresponding to the propensity score.

- We denote the vector of functional nuisance parameters by  $\eta_0 = (\mu_0^0, \mu_0^1, \pi_0)'$ .
- Let the distribution induced by this model be denoted by  $Q$  and thus  $Z_i = (X_i, W_i, Y_i) \stackrel{iid}{\sim} Q$  supported on  $\mathcal{Z} \subset \mathcal{X} \times \{0, 1\} \times \mathcal{Y}$ .

In this second setting, I will use the notation  $\mathbf{D}^{(0)}$  and  $\mathbf{D}^{(1)}$  to refer to the data subsets that contain 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 way.

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

**Assumption 3** (Technical Assumptions).

In both settings (Setup 1 and Setup 2) the following conditions hold as applicable to the settings, respectively.

- 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  $\infty$ .

$$\forall x \in \mathcal{X} : \quad 0 < \underline{f} \leq f(x) \leq \bar{f} < \infty \quad (2.5)$$

- $f(\cdot)$ ,  $\mu(\cdot)$ ,  $\mu_0^0(\cdot)$ ,  $\mu_0^1(\cdot)$ , and  $\pi_0(\cdot)$  are four times continuously differentiable with bounded second, third, and fourth-order partial derivatives. Specifically, in mathematical terms:

$$\forall w \in \{0, 1\} \quad \forall x \in \mathcal{X} \quad \forall (i, j, k, l) \in [d]^4 : \quad (2.6)$$

$$\begin{aligned} -\infty &< \underline{f}' \leq \partial_{i,j} f(x), \quad \partial_{i,j,k} f(x), \quad \partial_{i,j,k,l} f(x) \leq \bar{f}' < \infty \\ -\infty &< \underline{m}' \leq \partial_{i,j} \mu(x), \quad \partial_{i,j,k} \mu(x), \quad \partial_{i,j,k,l} \mu(x) \leq \bar{m}' < \infty \\ -\infty &< \underline{m}' \leq \partial_{i,j} \mu_0^w(x), \quad \partial_{i,j,k} \mu_0^w(x), \quad \partial_{i,j,k,l} \mu_0^w(x) \leq \bar{m}' < \infty \\ -\infty &< \underline{p}' \leq \partial_{i,j} \pi_0(x), \quad \partial_{i,j,k} \pi_0(x), \quad \partial_{i,j,k,l} \pi_0(x) \leq \bar{p}' < \infty \end{aligned}$$

- $\mu(\cdot), \mu_0^0(\cdot), \mu_0^1(\cdot) \in L^2(\mathcal{X})$  are square-integrable functions on  $\mathcal{X}$ .

There is considerable potential to relax these assumptions at the cost of requiring both less interpretable conditions and more technically sophisticated proofs. For example, the bounded derivatives condition can be relaxed to hold only in a neighborhood of  $x$  while requiring a weaker, more complex condition on the behavior of the derivatives beyond that neighborhood. Furthermore, it is necessary to point out again that this setup considers a highly stylized example, as applying these conditions to the IT training program quickly runs into problems.

**Example 1** (Average Hourly Earnings and IT Training - continued).

Consider two types of workers: those with 12 years of education, i.e., a high school diploma and those who drop out slightly before finishing high school and therefore have slightly less than 12 years of education. Given the continuous differentiability assumption on the regression function, I do not allow for a discontinuity in average earnings between high school graduates and dropouts. This seems unreasonable as employers would *ceteris paribus* prefer the worker with a high school diploma in most circumstances. Similarly, I would expect discontinuities in the distribution of educational attainment, for example right after the standard time necessary to achieve common educational milestones. For the sake of a simplified exposition, I will ignore problems such as this going forward.

Additionally, I require a rather standard assumption in localized regression approaches, namely that the variance changes continuously.

**Assumption 4** (Error Distribution Assumptions).

The error terms  $\varepsilon$  defined in Setup 1 and Setup 2, respectively, fulfill the following conditions.

- $\varepsilon$  has continuously varying variance. In other words,  $\sigma_\varepsilon^2 : \mathcal{X} \rightarrow \mathbb{R}_{>0}$  is a continuous function.
- $\sigma_\varepsilon^2 \in L^2(\mathcal{X})$  is a square-integrable function on  $\mathcal{X}$

As  $\mathcal{X}$  is a compact and bounded 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 seconds moments of  $Y$  in both scenarios. Furthermore, to ensure that there are a sufficient number of treated and untreated observations local to each point of interest asymptotically, I 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), I assume that there exists a constant  $\mathfrak{p} \in (0, 1/2)$  such that

$$\forall x \in \mathcal{X} : \quad 0 < \mathfrak{p} \leq \pi_0(x) \leq 1 - \mathfrak{p} < 1. \quad (2.7)$$

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

**Example 1** (Average Hourly Earnings and IT Training - continued).

In the IT training example, this condition requires that for each combination of characteristics, there are workers who choose to participate in the IT training and those who do not. More precisely, I assume that as I consider larger and larger data sets, there are both types of workers with characteristics that are very similar to those we are currently interested in. Intuitively this ensures that we have a suitable basis of workers to achieve meaningful comparisons.

**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 a 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.

$$\begin{aligned} Y_i &= \mathbb{1}([\mathfrak{W}_n(\mathbf{X})]_i = 0) \mu_0^0(\mathbf{X}_i) + \mathbb{1}([\mathfrak{W}_n(\mathbf{X})]_i = 1) \mu_0^1(\mathbf{X}_i) + \epsilon_i \\ &= \mathbb{1}([\mathfrak{W}'_n(\mathbf{X})]_i = 0) \mu_0^0(\mathbf{X}_i) + \mathbb{1}([\mathfrak{W}'_n(\mathbf{X})]_i = 1) \mu_0^1(\mathbf{X}_i) + \epsilon_i = Y'_i \end{aligned} \quad (2.8)$$

Technically, since 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 Distributional Nearest Neighbor Estimators

---

While less economically attractive, I 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, I will begin by adapting it to tackle the problem of estimating conditional average treatment effects.

#### 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  $\mu$  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, Cérou, and Guyader (2010). Given a sample as described in Assumption 1 and a fixed feature vector  $x$ , I 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 I 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  $\kappa$  following the notation of Ritzwoller and Syrgkanis (2024).

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

Here,  $\xi$  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  $\xi$  necessary again. Using  $\kappa$ , 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 trade-off 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 is valid. Then, for any fixed  $x \in \mathcal{X}$ , we have for some positive sequence  $\omega_n$  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 set-up, 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 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 and 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  $\mathfrak{c} \in (0, 1/2)$ . Then, for any fixed  $x \in \text{supp}(X) \subset \mathbb{R}^d$ , it holds that for some positive sequence  $\sigma_n$  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 in Heterogeneous Treatment Effect Estimation

Motivated by the nonparametric regression setup, I set out to apply the underlying idea in the context of heterogeneous treatment effects. Similarly to before, I 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. This approach closely follows the leading example of Ritzwoller and Syrgkanis (2024). 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  $x$ ,  $\theta_0(x) = \mathbb{E}[Y_i(W_i = 1) - Y_i(W_i = 0) | X_i = x]$ , I will employ a Neyman-orthogonal score function to curtail the influence of the nuisance parameters on our estimation.

$$\begin{aligned} M(x; \theta_0, \eta) &= \mathbb{E}[m(Z_i; \theta, \eta) | X_i = x] = 0 \quad \text{where} \\ m(Z_i; \theta, \eta) &= \mu^1(X_i) - \mu^0(X_i) + \beta(W_i, X_i)(Y_i - \mu^{W_i}(X_i)) - \theta(X_i) \end{aligned} \quad (3.16)$$

Here, I 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, I will furthermore use  $m(Z_i; \eta) = m(Z_i; \theta_0, \eta) + \theta_0(X_i)$ . 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.

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

However, due to the presence of infinite-dimensional nuisance parameters, it becomes attractive to proceed using this weighted empirical moment equation embedded in the DML2 estimator of Victor Chernozhukov, Chetverikov, 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, I will assume that  $m = n/K$ , i.e. the desired number of observations in each fold, is an integer going forward.

**Definition 1.** *DNN-DML2 CATE-Estimator*

To estimate the Conditional Average Treatment Effect at a point of interest  $x \in \mathcal{X}$ , 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. In contrast, for observation  $j$  (**NOT** the observation being assigned rank  $j$ ) denote the corresponding fold by  $k_j$ .

2. For each  $k \in [K]$ , use a first-stage estimator for the functional nuisance parameters on the data set  $\mathbf{D}_{I_k^C} \dots$

(a) to estimate the nuisance parameters  $\mu_0^0$  and  $\mu_0^1$ :

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

(b) if  $\pi_0$  is unknown, i.e. we are not in a randomized experiment setting, additionally estimate  $\pi_0$

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

For each fold  $k$ , denote the vector of estimates by  $\hat{\eta}_k = (\hat{\mu}_k^0, \hat{\mu}_k^1, \hat{\pi}_k)'$ . For each observation  $i$ , define the combined estimator as follows  $\hat{\eta}(Z_i) = \hat{\eta}_{k_{(i)}}(Z_i)$ .

3. Construct the estimator  $\hat{\theta}(x)$  as the solution to the following equation.

$$0 = \sum_{k=1}^K \sum_{i \in I_k} K(x, X_i) m(Z_i; \hat{\theta}(x), \hat{\eta}_k) = \sum_{i=1}^{n-s+1} \frac{\binom{n-i}{s-1}}{\binom{n}{s}} m(Z_{(i)}; \hat{\theta}(x), \hat{\eta}_{k_{(i)}}) \quad (3.21)$$

Observe that the weights  $K(x, X_i)$  chosen in the second step are chosen according to the whole sample - not according to the chosen folds. 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 because a lower subsampling scale averages over a larger number of observations and thus can contribute to better smoothing properties for the propensity score. Using the score function in the equation that defines the estimator, we can observe the following.

$$\begin{aligned} \hat{\theta}(x) &= \sum_{i=1}^{n-s+1} \frac{\binom{n-i}{s-1}}{\binom{n}{s}} \left[ \hat{\mu}_{k_{(i)}}^1(X_{(i)}) - \hat{\mu}_{k_{(i)}}^0(X_{(i)}) + \hat{\beta}_{k_{(i)}}(W_{(i)}, X_{(i)}) (Y_{(i)} - \hat{\mu}_{k_{(i)}}^{W_{(i)}}(X_{(i)})) \right] \\ &= \sum_{i=1}^{n-s+1} \frac{\binom{n-i}{s-1}}{\binom{n}{s}} m(Z_{(i)}, \hat{\eta}_{k_{(i)}}) \end{aligned} \quad (3.22)$$

Thus, given first-stage estimates of nuisance parameters, I 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 re-weighting of the terms corresponding to different observations.

As is typical in the literature on doubly-robust inference, we need to make assumptions on the convergence rate of

the first-stage estimators to allow for the desired type of inference. Specifically, I take the following slightly modified assumptions from Assumption 3.2 of Victor Chernozhukov, Chetverikov, et al. (2018).

**Assumption 8** (DDML-Rate Conditions).

*LOREM IPSUM*

1. Given a random subset  $I$  of  $[n]$  of size  $m = n/k$ , the nuisance parameter estimator  $\hat{\eta}(\mathbf{D}_{IC})$  belongs to the realization set  $\mathcal{T}_n$  with probability  $1 - \Delta_n$ , where  $\mathcal{T}_n$  contains  $\eta_0$  and is constrained by the following conditions. Denote the event that  $\hat{\eta}(\mathbf{D}_{IC}) \in \mathcal{T}_n$  by  $\mathcal{E}_n$ .
2. The following moment condition holds:

$$m_n := \sup_{\eta \in \mathcal{T}_n} (\mathbb{E}_Z [|m(Z; \theta_0, \eta)|^q])^{1/q} \leq c_1 \quad (3.23)$$

3. The following conditions on the statistical rates  $r_n$ ,  $r'_n$ , and  $\lambda'_n$  hold:

$$r'_n := \sup_{\eta \in \mathcal{T}_n} \left( \mathbb{E}_Z \left[ |m(Z; \theta_0, \eta) - m(Z; \theta_0, \eta_0)|^2 \right] \right)^{1/2} \leq \delta_n, \quad (3.24)$$

$$\lambda'_n := \sup_{r \in (0,1), \eta \in \mathcal{T}_n} \left| \partial_r^2 \mathbb{E}_Z [m(Z; \theta_0, \eta_0 + r(\eta - \eta_0))] \right| \leq \delta_n / \sqrt{n}. \quad (3.25)$$

4. The variance of the score is non-zero.

*LOREM IPSUM* (3.26)

In contrast to the paper that these assumptions are taken from, there is an absence of conditions on what is named  $\psi^a$  in the original paper. This is due to the fact that in the case of the CATE, this term is  $-1$  and thus independent of the nuisance parameters. Thus, specific constraints on the effects of the first-stage estimation problem are unnecessary. In a potential generalization of the method proposed in this paper to other doubly-robust inference problems, this aspect should be addressable with relatively minimal additional technical arguments along the lines in Victor Chernozhukov, Chetverikov, et al. (2018).

## 4 DNN Estimators as Generalized U-Statistics

As most 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).

$$\tilde{\mu}_s(x; \mathbf{D}_n) = \binom{n}{s}^{-1} \sum_{\ell \in L_{n,s}} h_s(x; D_\ell) \quad \text{and} \quad \hat{\mu}_{s_1, s_2}(x; \mathbf{D}_n) = \binom{n}{s}^{-1} \sum_{\ell \in L_{n, s_2}} h_{s_1, s_2}(x; D_\ell) \quad (4.1)$$

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 4.1** (Kernel of TDNN Estimator - Adapted from Lemma 8 of Demirkaya et al. (2024)).

*The kernel of the TDNN estimator takes the following form.*

$$\begin{aligned} h_{s_1, s_2}(x; D) &= w_1^* \left[ \binom{s_2}{s_1}^{-1} \sum_{\ell \in L_{s_2, s_1}} h_{s_1}(x; D_\ell) \right] + w_2^* h_{s_2}(x; D) \\ &= w_1^* \tilde{\mu}_{s_1}(x; D) + w_2^* h_{s_2}(x; D) \end{aligned} \quad (4.2)$$

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

$$\psi_s^c(x; \mathbf{z}_1, \dots, \mathbf{z}_c) = \mathbb{E}_D [h_s(x; D) \mid Z_1 = \mathbf{z}_1, \dots, Z_c = \mathbf{z}_c] \quad (4.3)$$

$$h_s^{(1)}(x; \mathbf{z}_1) = \psi_s^1(x; \mathbf{z}_1) - \mu(x) \quad (4.4)$$

$$h_s^{(c)}(x; \mathbf{z}_1, \dots, \mathbf{z}_c) = \psi_s^c(x; \mathbf{z}_1, \dots, \mathbf{z}_c) - \sum_{j=1}^{c-1} \left( \sum_{\ell \in L_{n,j}} h_s^{(j)}(x; \mathbf{z}_\ell) \right) - \mu(x) \quad \text{for } c = 2, \dots, s \quad (4.5)$$

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

$$\xi_s^c(x) = \text{Var}_{1:c}(\psi_s^c(x; Z_1, \dots, Z_c)) \quad (4.6)$$

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

$$\zeta_{s_1, s_2}^c(x) = \text{Var}_{1:c}(\psi_{s_1, s_2}^c(x; Z_1, \dots, Z_c)) \quad (4.7)$$

with an analogous definition of  $Z'$ . As a byproduct (or main purpose depending on the perspective) these terms can be used to derive the Hoeffding decomposition of the TDNN estimator.

$$H_s^c(x; \mathbf{D}_n) = \binom{n}{c}^{-1} \sum_{\ell \in L_{n,c}} h_s^{(c)}(x; D_\ell) \quad \text{and} \quad H_{s_1, s_2}^c(x; \mathbf{D}_n) = \binom{n}{c}^{-1} \sum_{\ell \in L_{n,c}} h_{s_1, s_2}^{(c)}(x; D_\ell) \quad (4.8)$$

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

$$\tilde{\mu}_s(x; \mathbf{D}_n) = \mu(x) + \sum_{j=1}^s \binom{s}{j} H_s^j(x; \mathbf{D}_n) \quad \text{and} \quad \hat{\mu}_{s_1, s_2}(x; \mathbf{D}_n) = \mu(x) + \sum_{j=1}^{s_2} \binom{s_2}{j} H_{s_1, s_2}^j(x; \mathbf{D}_n) \quad (4.9)$$

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.

$$\forall c = 1, 2, \dots, j-1 : \mathbb{E}_D \left[ h_{s_1, s_2}^{(j)}(x; D) \mid Z_1 = \mathbf{z}_1, \dots, Z_c = \mathbf{z}_c \right] = 0 \quad \text{and} \quad \mathbb{E}_D \left[ h_{s_1, s_2}^{(j)}(x; D) \right] = 0 \quad (4.10)$$

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

$$\text{Var}_D(\hat{\mu}_{s_1, s_2}(x; D)) = \sum_{j=1}^{s_2} \binom{s_2}{j}^2 \text{Var}_D(H_{s_1, s_2}^j(x; D)) \quad (4.11)$$

$$\text{Var}_D(H_{s_1, s_2}^j(x; D)) = \binom{n}{j}^{-1} \text{Var}_D(h_{s_1, s_2}^{(j)}(x; D)) =: \binom{n}{j}^{-1} V_{s_1, s_2}^j(x) \quad (4.12)$$

Third, the following equivalent expression for the kernel variance.

$$\zeta_{s_1, s_2}^{s_2}(x) = \text{Var}_D(h_{s_1, s_2}(x; D)) = \sum_{j=1}^{s_2} \binom{s_2}{j} V_{s_1, s_2}^j(x) \quad (4.13)$$

#### 4.1 CATE-Estimators as Generalized U-Statistics

Given estimates of the functional nuisance parameters, the proposed CATE estimators can be analyzed as generalized U-statistics in the same way as in the nonparametric regression context. As most of the theoretical results on these estimators will similarly rely on Hoeffding projection arguments, I will introduce analogous notation in this more general scenario. First, observe that the DNN-DML2 CATE estimator can be rewritten as follows to explicitly show its construction as a generalized U-statistic.

$$\begin{aligned} \hat{\theta}(x; \mathbf{D}) &= \sum_{i=1}^{n-s+1} \frac{\binom{n-1}{s-1}}{\binom{n}{s}} \left[ \hat{\mu}_{k(i)}^1(X_{(i)}) - \hat{\mu}_{k(i)}^0(X_{(i)}) + \hat{\beta}_{k(i)}(W_{(i)}, X_{(i)}) \left( Y_{(i)} - \hat{\mu}_{k(i)}^{W_{(i)}}(X_{(i)}) \right) \right] \\ &= \binom{n}{s}^{-1} \sum_{\ell \in L_{n,s}} \sum_{i=1}^n \frac{\mathbb{1}(\text{rk}(x; Z_i, D_\ell) = 1)}{s!} \left[ \hat{\mu}_{k_i}^1(X_i) - \hat{\mu}_{k_i}^0(X_i) + \hat{\beta}_{k_i}(W_i, X_i) \left( Y_i - \hat{\mu}_{k_i}^{W_i}(X_i) \right) \right] \\ &= \binom{n}{s}^{-1} \sum_{\ell \in L_{n,s}} \underbrace{\sum_{i=1}^n \frac{\mathbb{1}(\kappa(x; Z_i, D_\ell) = 1)}{s!} m(Z_i, \hat{\eta}_{k_i})}_{\chi_s(x; \mathbf{D}_\ell, \hat{\eta})} \\ &= \binom{n}{s}^{-1} \sum_{\ell \in L_{n,s}} \chi_s(x; \mathbf{D}_\ell, \hat{\eta}) \end{aligned} \quad (4.14)$$

In the same fashion as in the previous case, I can now define the Hoeffding decomposition for the estimator. Note that in this step, the functional nuisance parameter estimates  $\hat{\mu}$  and  $\hat{\pi}$  are considered as depending on the data that



we form the expectation over and are thus itself random and part of the expectation.

$$\vartheta_s^c(x; \mathbf{z}_1, \dots, \mathbf{z}_c, \hat{\eta}) = \mathbb{E}_D [\chi_s(x; \mathbf{D}, \hat{\eta}) | Z_1 = \mathbf{z}_1, \dots, Z_c = \mathbf{z}_c] \quad (4.15)$$

$$\chi_s^{(1)}(x; \mathbf{z}_1, \hat{\eta}) = \vartheta_s^1(x; \mathbf{z}_1, \hat{\eta}) - \mathbb{E}_D [\chi_s(x; \mathbf{D}, \hat{\eta})] \quad (4.16)$$

As before, I define the higher-order projection terms, i.e. for  $c = 2, \dots, s$ , in the following way.

$$\chi_s^{(c)}(x; \mathbf{z}_1, \dots, \mathbf{z}_c, \hat{\eta}) = \vartheta_s^c(x; \mathbf{z}_1, \dots, \mathbf{z}_c, \hat{\eta}) - \sum_{j=1}^{c-1} \left( \sum_{\ell \in L_{n,j}} \chi_s^{(j)}(x; \mathbf{z}_\ell, \hat{\eta}) \right) - \mathbb{E}_D [\chi_s(x; \mathbf{D}, \hat{\eta})] \quad (4.17)$$

In anticipation of arguments involving empirical process theory, I define as follows.

$$\chi_{s,0}(x; \mathbf{D}_{[s]}) = \chi_s(x; \mathbf{D}_{[s]}, \eta_0) \quad (4.18)$$

$$\vartheta_{s,0}^c(x; \mathbf{z}_1, \dots, \mathbf{z}_c) = \vartheta_s^c(x; \mathbf{z}_1, \dots, \mathbf{z}_c, \eta_0) \quad (4.19)$$

$$\chi_{s,0}^{(1)}(x; \mathbf{z}_1) = \chi_s^{(1)}(x; \mathbf{z}_1, \eta_0) = \vartheta_{s,0}^1(x; \mathbf{z}_1) - \mathbb{E}_D [\chi_{s,0}(x; \mathbf{D})] \quad (4.20)$$

$$\begin{aligned} \chi_{s,0}^{(c)}(x; \mathbf{z}_1, \dots, \mathbf{z}_c) &= \chi_s^{(c)}(x; \mathbf{z}_1, \dots, \mathbf{z}_c, \eta_0) \\ &= \vartheta_{s,0}^c(x; \mathbf{z}_1, \dots, \mathbf{z}_c) - \sum_{j=1}^{c-1} \left( \sum_{\ell \in L_{n,j}} \chi_{s,0}^{(j)}(x; \mathbf{z}_\ell) \right) - \mathbb{E}_D [\chi_{s,0}^{(c)}(x; Z_1, \dots, Z_c)] \end{aligned} \quad (4.21)$$

These definitions allow us to use the following decomposition, where now I acknowledge the randomness in the first-stage estimation.

$$\chi_s^{(c)}(x; \mathbf{D}_\ell, \hat{\eta}) = \chi_{s,0}^{(c)}(x; \mathbf{D}_\ell) + \underbrace{\chi_s^{(c)}(x; \mathbf{D}_\ell, \hat{\eta}) - \chi_{s,0}^{(c)}(x; \mathbf{D}_\ell)}_{R_c(x; \mathbf{D}_\ell)} \quad (4.22)$$

Considering this decomposition will be of great usefulness when analyzing the asymptotic properties of the CATE estimators. Let  $\hat{\mu}_k$  and  $\hat{\pi}_k$  denote the nuisance parameter estimates calculated on the complement of fold  $k$ . Using the full expression for the CATE estimator, the following is obtained.

$$\begin{aligned} \hat{\theta}(x; \mathbf{D}) &= \underbrace{\mathbb{E}_D [\hat{\theta}(x; \mathbf{D})]}_{\text{Centering-Term}} + \underbrace{\frac{s}{n} \sum_{i=1}^n \chi_{s,0}^{(1)}(x; Z_i)}_{\text{Oracle-Hájek-Projection}} + \underbrace{\frac{s}{k} \sum_{l=1}^k \frac{1}{m} \sum_{i \in \mathcal{I}_k} \left( \underbrace{\chi_s^{(1)}(x; Z_i, \hat{\eta}_k) - \chi_{s,0}^{(1)}(x; Z_i)}_{R_{1,k}(x; Z_i)} \right)}_{\text{Oracle-Hájek-Projection Error}} \\ &\quad + \underbrace{\sum_{j=2}^s \binom{s}{j} \binom{n}{j}^{-1} \sum_{\ell \in L_{n,j}} \chi_{s,0}^{(j)}(x; \mathbf{D}_\ell)}_{\text{Oracle-Hájek-Residual}} + \underbrace{\sum_{j=2}^s \binom{s}{j} \binom{n}{j}^{-1} \sum_{\ell \in L_{n,j}} R_j(x; \mathbf{D}_\ell)}_{\text{Higher-Order Error Terms}} \end{aligned} \quad (4.23)$$

I call these terms Oracle-Hájek projection and Oracle-Hájek residual because of the additional randomness due to the first-stage estimation of the functional nuisance parameters that is reflected in the  $R_j$  terms. An integral part of showing the asymptotic properties of the estimator is to show that a number of terms in this representation are asymptotically negligible as long as the first-stage estimator converges fast enough.

## 5 Pointwise Inference for DNN Estimators

To perform inference in the regression setup, Demirkaya et al. (2024) introduce variance estimators based on Jackknife and Bootstrap. However, as they point out, their consistency results rely on a likely suboptimal rate condition for the subsampling scale. Although Theorem 3.2 allows  $s_2$  to be of the order  $o(n)$ , the variance estimation results rely on the considerably stronger condition  $s_2 = o(n^{1/3})$ . 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, I 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.

### 5.1 Jackknife Variance Estimators for Nonparametric Regression

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

$$\omega^2(x) = \text{Var}_D(\hat{\mu}_{s_1, s_2}(x; \mathbf{D}_n)) \quad (5.1)$$

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.

$$\hat{\omega}_{JK}^2(x; \mathbf{D}_n) = \frac{n-1}{n} \sum_{i=1}^n (\hat{\mu}_{s_1, s_2}(x; \mathbf{D}_{n, -i}) - \hat{\mu}_{s_1, s_2}(x; \mathbf{D}_n))^2 \quad (5.2)$$

**Theorem 5.1** (Closed Form Expression for the Jackknife-Variance Estimator).

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

$$\text{LOREMIPSUM} \quad (5.3)$$

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

$$\text{LOREMIPSUM} \quad (5.4)$$

As a generalization to the Jackknife, we can also consider the delete-d Jackknife that builds on the same working principle but averages over all possible d-subset removals. This leads to the following representation.

$$\hat{\omega}_{JKD}^2(x; d, \mathbf{D}_n) = \frac{n-d}{d} \binom{n}{d}^{-1} \sum_{\ell \in L_{n,d}} (\hat{\mu}_{s_1, s_2}(x; \mathbf{D}_{n, -\ell}) - \hat{\mu}_{s_1, s_2}(x; \mathbf{D}_n))^2 \quad (5.5)$$

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, because of the unwieldiness of the closed form, I refrain from deriving it.

In this section, I 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.

$$\hat{\omega}_{PI}^2(x; \mathbf{D}_n) = \frac{s_2^2}{n^2} \sum_{i=1}^n \left[ \left( \binom{n-1}{s-1}^{-1} \sum_{\ell \in L_{s_2-1}([n] \setminus \{i\})} h_{s_1, s_2}(x; D_{\ell \cup \{i\}}) \right) - \hat{\mu}_{s_1, s_2}(x; \mathbf{D}_n) \right]^2 \quad (5.6)$$

## 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, I obtain the following result.

**Theorem 5.2** (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

$$\frac{\hat{\omega}_{PI}^2(x; \mathbf{D}_n)}{\omega^2(x; \mathbf{D}_n)} \xrightarrow{p} 1. \quad (5.7)$$

In an analogous fashion to Theorems 5 and 6 of Demirkaya et al. (2024), I further obtain the following consistency results for the variance estimators presented. 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), I then obtain the following results on the Jackknife and Bootstrap variance estimators, respectively. As part of the proof of these results, I obtain general results on the consistency of Jackknife and Bootstrap variance estimators for infinite-order U-statistics beyond the TDNN estimator.

**Theorem 5.3** (Jackknife Variance Estimator Consistency).

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

$$\frac{\hat{\omega}_{JK}^2(x; \mathbf{D}_n)}{\omega^2(x; \mathbf{D}_n)} \xrightarrow{p} 1. \quad (5.8)$$

**Theorem 5.4** (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

$$\frac{\hat{\omega}_{JKD}^2(x; d, \mathbf{D}_n)}{\omega^2(x; \mathbf{D}_n)} \xrightarrow{p} 1. \quad (5.9)$$

## 5.2 Asymptotic Normality of the (T)DNN-DML2 CATE Estimator

---

As a first intermediate result in the analysis of the asymptotic properties of the CATE estimator, we can find the following concerning an oracle variant of the proposed estimator defined as the solution to the following problem.

$$0 = \sum_{i=1}^n K(x, X_i) m \left( Z_i; \hat{\theta}_0(x), \eta_0 \right) = \sum_{i=1}^{n-s+1} \frac{\binom{n-i}{s-1}}{\binom{n}{s}} m \left( Z_{(i)}; \hat{\theta}_0(x), \eta_0 \right) \quad (5.10)$$

This is a rather basic result that can be seen as a corollary to the nonparametric regression case.

**Theorem 5.5 (PRELIMINARY! THIS IS NOT YET FORMALLY SHOWN).**

*Assume that we observe data as described in Assumption 2 and that Assumption 3 is valid. Then, for any fixed  $x \in \mathcal{X}$ , we have for some positive sequence  $\omega_n$  of order  $\sqrt{s/n}$*

$$\frac{\hat{\theta}_0(x; \mathbf{D}_n) - \theta_0(x)}{\omega_n} \rightsquigarrow \mathcal{N}(0, 1) \quad (5.11)$$

*as  $n, s \rightarrow \infty$  with  $s = o(n)$ .*

LOREM IPSUM

## 5.3 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 (2024) 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 (2024). 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, I need to introduce some additional notation, where, for simplicity, I assume that  $m$ , i.e. the number of observations in each  $I_k$ , is even.

**Definition 2** (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.

$$|H_k| = \frac{|I_k|}{2} = m/2 \quad \text{and} \quad H_k \subset I_k \quad (5.12)$$

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.

$$R_{n,s}^*(x; \mathbf{D}_{[n]}, \mathcal{I}) = \bar{\theta}_{\mathcal{H}}(x) - \hat{\theta}(x) \quad (5.13)$$

Here,  $\bar{\theta}_{\mathcal{H}}(x)$  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, I 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.

$$\begin{aligned} 0 &= \sum_{k=1}^K \sum_{i \in H_k} K(x, X_i | \mathcal{H}) m(Z_i; \overline{CATE}_{\mathcal{H}}(x), \hat{\mu}_{k,s}, \hat{\pi}_{k,s}) \\ &= \sum_{i=1}^{n/2-s+1} \left[ \frac{\binom{n/2-i}{s-1}}{\binom{n/2}{s}} m(Z_{(i | \mathcal{H})}; \bar{\theta}_{\mathcal{H}}(x), \hat{\mu}_{k(i | \mathcal{H}),s}, \hat{\pi}_{k(i | \mathcal{H}),s}) \right] \end{aligned} \quad (5.14)$$

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

$$\begin{aligned} \bar{\theta}_{\mathcal{H}}(x) &= \sum_{i=1}^{n/2-s+1} \frac{\binom{n/2-i}{s-1}}{\binom{n/2}{s}} \left[ \hat{\mu}_{k(i | \mathcal{H}),s}^1(X_{(i | \mathcal{H})}) - \hat{\mu}_{k(i | \mathcal{H}),s}^0(X_{(i | \mathcal{H})}) \right. \\ &\quad \left. + \hat{\beta}_{k(i | \mathcal{H}),s}(W_{(i | \mathcal{H})}, X_{(i | \mathcal{H})})(Y_{(i | \mathcal{H})} - \mu_{W_{(i | \mathcal{H})}}(X_{(i | \mathcal{H})})) \right] \end{aligned} \quad (5.15)$$

Recognizing the similarity to  $\hat{\theta}(x)$ , I can further simplify in the following way.

$$R_{n,s}^*(x; \mathbf{D}_{[n]}, \mathcal{I}) = \textcolor{red}{LOREM IPSUM} \quad (5.16)$$

**Theorem 5.6** (Consistent Variance Estimation for the (T)DNN-DML2 CATE Estimator).

*LOREM IPSUM*

## 5.4 Pointwise Inference with the TDNN Estimator

**Theorem 5.7** (Pointwise Inference in Nonparametric Regression).

*LOREM IPSUM*

**Theorem 5.8** (Pointwise Inference in Heterogeneous Treatment Effect Estimation).

*LOREM IPSUM*

## 6 Simulations

Having developed theoretical results concerning simultaneous inference methods for the TDNN estimator, I 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, I consider a collection of setups. First, I focus on illustrating the bias correction 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 available by combining multiple subsampling scales.

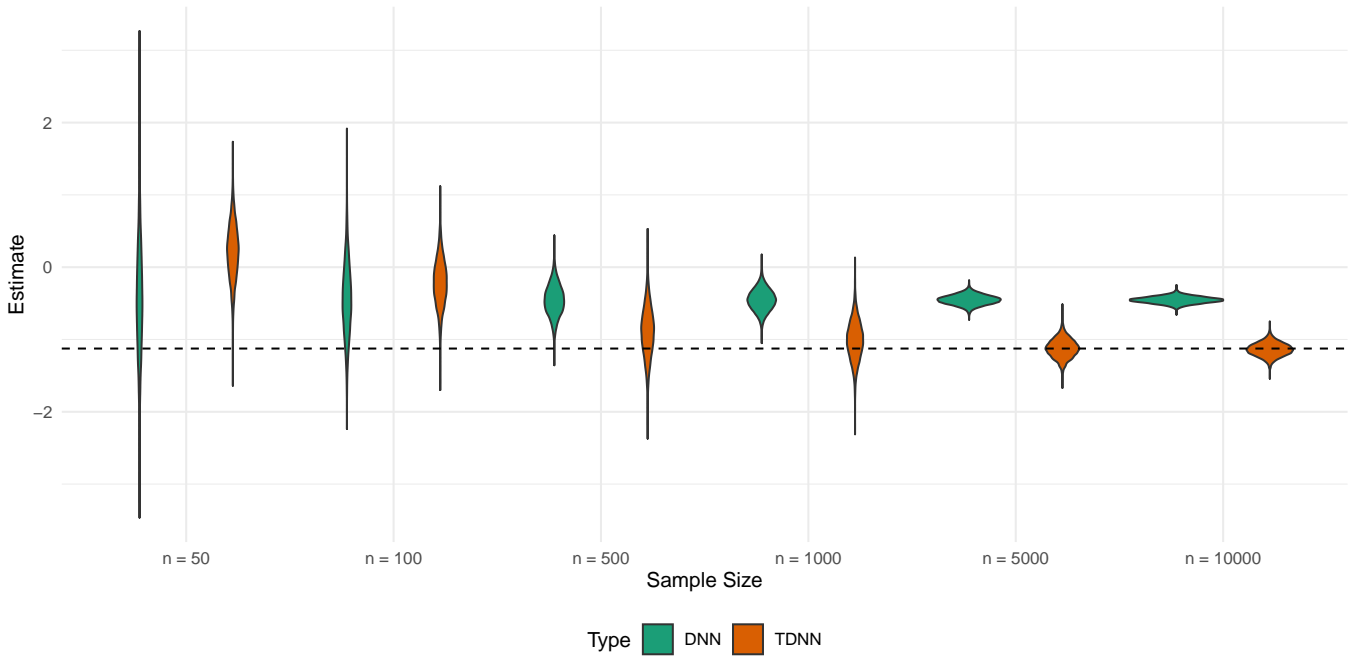


Figure 1: Comparison of the DNN ( $s = 20$ ) and TDNN ( $s_1 = 20, s_2 = 50$ ) Estimators for different sample sizes. The dashed line indicates the value of the unknown regression function at the point of interest. Simulation Setup replicates Setting 1 from Demirkaya et al. (2024) for 10000 Monte Carlo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

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

As a second, potentially more illustrative example, I consider the estimation of a function of two arguments. Specifically, I 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 shown in Figure 2.

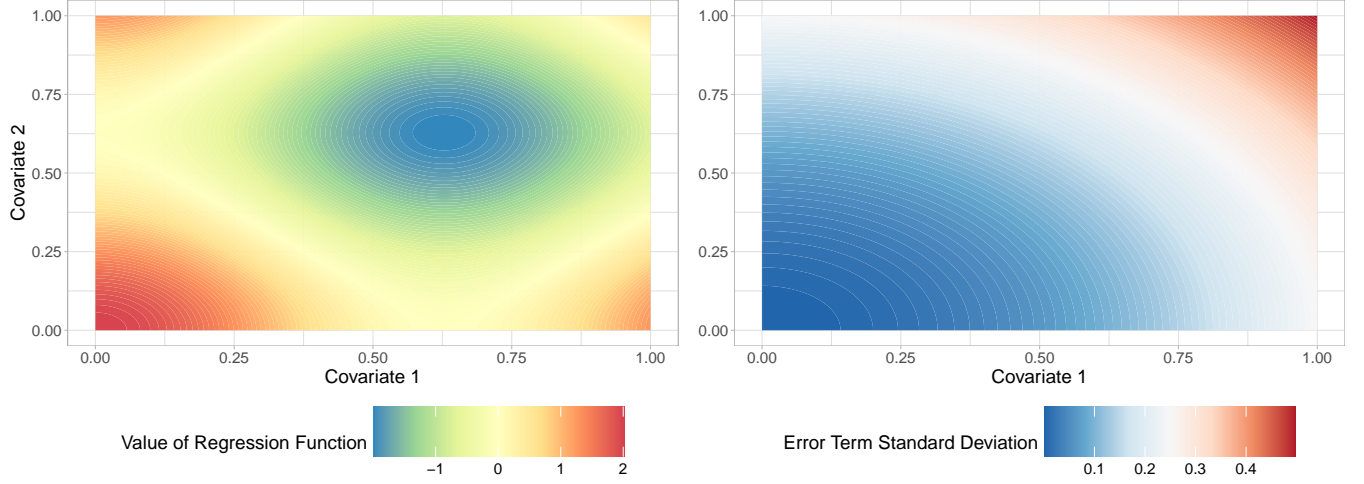


Figure 2: Value of the Regression Function (left) and Variance of the Error Term (right)

To analyze the behavior of the estimator in this setting, I run a number of Monte Carlo simulations each consisting of 10000 simulation runs. Although 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  $[0, 1]^2$ , and I find the following concerning the bias and variance of the estimators given different kernel orders  $s_1$  and  $s_2$ .

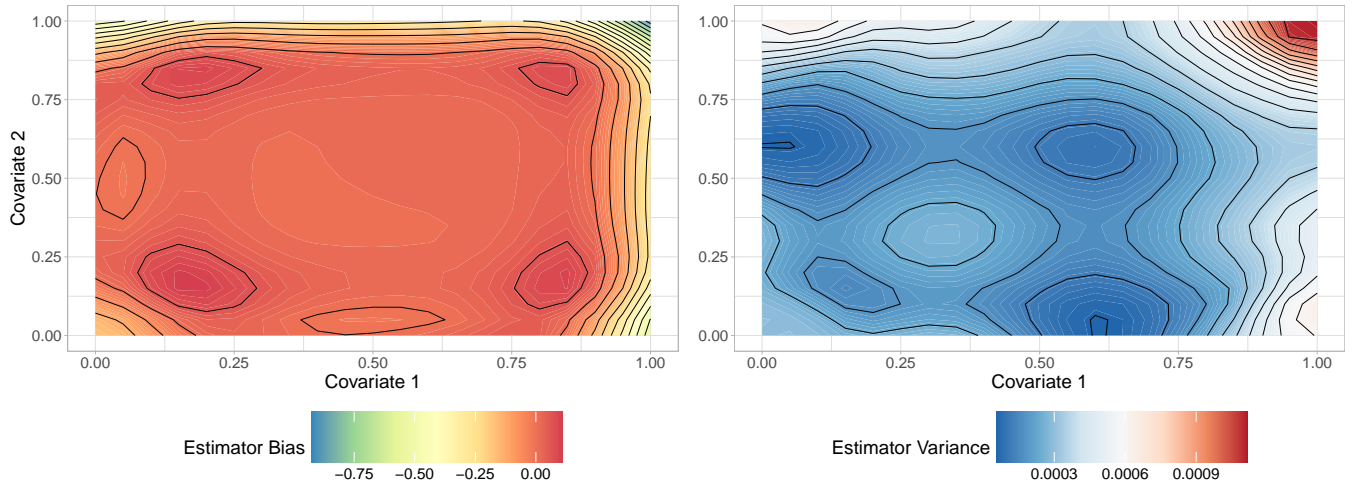


Figure 3: Approximate Bias (left) and Variance (right) of the TDNN Estimator with  $s_1 = 10$  and  $s_2 = 25$



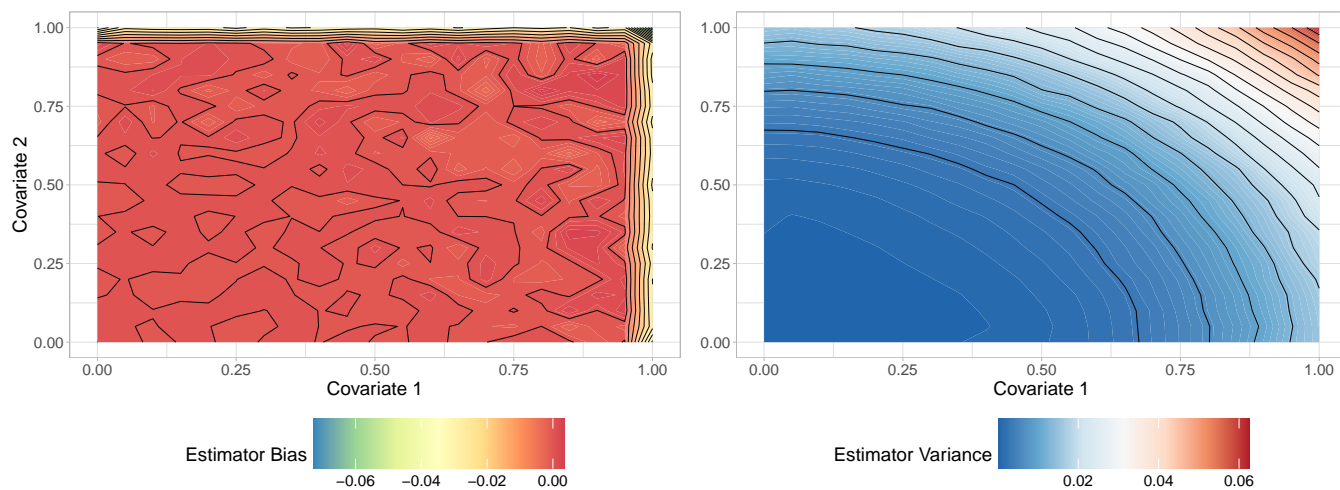


Figure 4: Approximate Bias (left) and Variance (right) of the TDNN Estimator with  $s_1 = 1000$  and  $s_2 = 2500$

## 6.2 CATE-Estimation

---

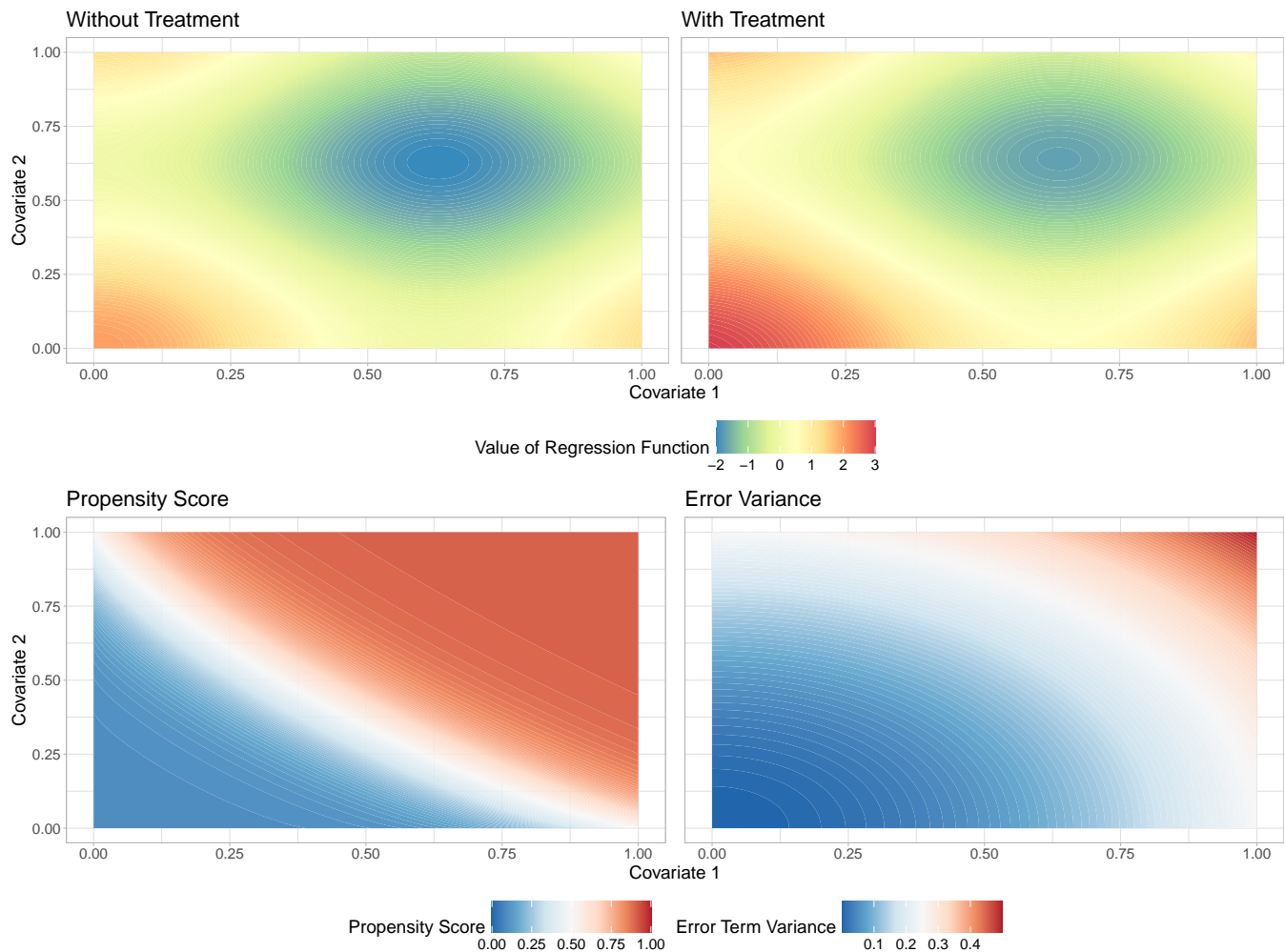


Figure 5: Value of the Regression Functions  $\mu_0$  (upper) and  $\mu_1$  (lower). Error term structure remains unchanged.

LOREM IPSUM

# 7 Application

---

LOREM IPSUM

## 8 Conclusion and Outlook

---

LOREM IPSUM

### 8.1 Outlook

---

The promising results of this paper lay out a few avenues for continued research. A first idea to potentially improve on the performance of the implemented weighting scheme is to use weights that correspond to a column-subsampled version of the (T)DNN weighting scheme. In the context of RF this is a commonly employed technique that helps to de-correlate the individual trees that are used to construct the RF. This decorrelation has been shown to be an integral part of the desirable performance that RF delivers in practice via simulation studies. In the context of this paper, a similar effect can be expected to take place while having the potential to still be expressible in a relatively simple, principled closed-form weighting scheme.

LOREM IPSUM

## References

- Arcones, Miguel A. and Evarist Gine (June 1992). “On the Bootstrap of  $U$ -Statistics”. In: *The Annals of Statistics* 20.2. DOI: 10.1214/aos/1176348650.
- Arvesen, James N. (Dec. 1969). “Jackknifing  $U$ -Statistics”. en. In: *The Annals of Mathematical Statistics* 40.6, pp. 2076–2100. DOI: 10.1214/aoms/1177697287.
- Biau, Gérard, Frédéric Cérou, and Arnaud Guyader (2010). “On the Rate of Convergence of the Bagged Nearest Neighbor Estimate”. en. In: *Journal of Machine Learning Research* 11.22, pp. 687–712.
- 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. Cham: Springer International Publishing. DOI: 10.1007/978-3-319-25388-6.
- Breiman, Leo (Oct. 2001). “Random Forests”. en. 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. 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, Denis Chetverikov, 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.
- Chernozhukov, Victor, Whitney K. Newey, and Vasilis Syrgkanis (Dec. 2024). *Conditional Influence Functions*. DOI: 10.48550/arXiv.2412.18080.
- 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.
- LaLonde, Robert J. (1986). “Evaluating the Econometric Evaluations of Training Programs with Experimental Data”. In: *The American Economic Review* 76.4, pp. 604–620.
- Lee, A J. (Mar. 2019). *U-Statistics*. en. 0th ed. Routledge. DOI: 10.1201/9780203734520.
- Lin, Yi and Yongho Jeon (June 2006). “Random Forests and Adaptive Nearest Neighbors”. en. 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. DOI: 10.1214/21-EJS1958.
- Peng, Wei, Lucas Mentch, and Leonard Stefanski (2021). *Bias, Consistency, and Alternative Perspectives of the Infinitesimal Jackknife*. DOI: 10.48550/ARXIV.2106.05918.
- Ritzwoller, David M. and Vasilis Syrgkanis (Sept. 2024). *Simultaneous Inference for Local Structural Parameters with Random Forests*. DOI: 10.48550/arXiv.2405.07860.
- Semenova, Vira and Victor Chernozhukov (May 2021). “Debiased machine learning of conditional average treatment effects and other causal functions”. 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”. en. 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”. en. In: *Journal of the American Statistical Association* 113.523, 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: *Journal of Machine Learning Research* 15.48, pp. 1625–1651.

## A Useful Results

### A.1 Properties of the $\kappa$ Function

**Lemma A.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,  $\varphi$  denotes the probability measure on  $\mathbb{R}^d$  induced by the random vector  $X$ .

**Lemma A.2** (Demirkaya et al. (2024) - Lemma 13).

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

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

As an additional tool, we will make use of the following analogous results concerning products of two kernel functions with nonzero expectation. These results will then be used to construct an analogon of Lemma A.2 for the corresponding cases. This will serve very similar purposes in the analysis of (conditional) covariance terms as the previous results serve for (conditional) expectations.

**Lemma A.3.**

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 three statements hold.

$$\forall i \in [c] : \mathbb{E}_{D, D'} [\kappa(x; Z_i, D) \kappa(x; Z_i, D')] = (2s - c)^{-1} = \omega(e^{-s}) \quad (\text{A.2})$$

$$\begin{aligned} \forall i \in [c] \forall j \in \{c+1, \dots, s\} : \mathbb{E}_{D, D'} [\kappa(x; Z_i, D) \kappa(x; Z'_j, D')] &= \frac{1}{(2s - c)(2s - c - 1)} \sum_{i=0}^{s-c-1} \binom{s-c-1}{i} \binom{2s-c-2}{i}^{-1} \\ &= \omega(e^{-s}) \end{aligned} \quad (\text{A.3})$$

$$\begin{aligned} \forall i, j \in \{c+1, \dots, s\} : \mathbb{E}_{D, D'} [\kappa(x; Z_i, D) \kappa(x; Z'_j, D')] &= \frac{2}{(2s - c)(2s - c - 1)} \sum_{i=0}^{s-c-1} \binom{s-c-1}{i} \binom{2s-c-1}{s-1+i}^{-1} \\ &= \omega(e^{-s}) \end{aligned} \quad (\text{A.4})$$

*Proof of Lemma A.3.*

Without loss of generality, we will consider the cases of  $i = 1$  and  $j = c + 1$  for the first two equations.

$$\begin{aligned}\mathbb{E}_{D,D'} [\kappa(x; Z_1, D) \kappa(x; Z_1, D')] &= \mathbb{E}_{D,D'} \left[ \kappa(x; Z_1, D_{1:c}) \kappa(x; Z_1, D_{(c+1):s}) \kappa(x; Z_1, D'_{(c+1):s}) \right] \\ &= \mathbb{E} [\kappa(x; Z_1, D_{[2s-c]})] = (2s - c)^{-1}\end{aligned}\tag{A.5}$$

Considering the second case, we find the following.

$$\begin{aligned}\mathbb{E}_{D,D'} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D')] &= \frac{1}{(2s - c)!} \sum_{i=0}^{s-c-1} \binom{s-c-1}{i} i! ((s-1) + (s-c-1-i))! \\ &= \frac{1}{(2s - c)!} \sum_{i=0}^{s-c-1} \binom{s-c-1}{i} i! (2s - c - 2 - i)! = \frac{(2s - c - 2)!}{(2s - c)!} \sum_{i=0}^{s-c-1} \binom{s-c-1}{i} \binom{2s - c - 2}{i}^{-1}\end{aligned}\tag{A.6}$$

While unintuitive at first, the terms in this expression have intuitive meaning when we consider this as a combinatoric problem. Consider lining up the observations in order of their distance to the point of interest and counting the cases for which the expression in the expectation is equal to one. First, there are  $(2s - c)!$  possible orderings of the observations with probability one, leading to the denominator. Next, notice that only those orderings where  $\|X'_{c+1} - x\| \leq \|X_1 - x\|$  and  $\|X_1 - x\| \leq \|X_i - x\|$  for any  $i = 2, \dots, c$  can possibly lead to a non-zero realization of the kernel term. Furthermore, out of the  $(s - c - 1)$  observations in  $D'_{(c+2):s}$ , it is possible for  $i = 0, \dots, s - c - 1$  observations to lie at a distance to the point of interest that is smaller than  $\|X_1 - x\|$  but larger than  $\|X'_{c+1} - x\|$  in any permutation. The sum adjusts for those possible configurations. Next, we can make the following observation concerning the expression we just derived.

$$\begin{aligned}\mathbb{E}_{D,D'} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D')] &\geq \frac{1}{(2s - c)(2s - c - 1)} \sum_{i=0}^{s-c-1} \frac{(s - c - i)^i}{i!} \frac{i!}{(2s - c - 2)^i} \\ &= \frac{1}{(2s - c)(2s - c - 1)} \left( 1 + \sum_{i=0}^{s-c-1} \left( \frac{s - c - i}{2s - c - 2} \right)^i \right) \\ &\geq \frac{1}{(2s - c)^2}\end{aligned}\tag{A.7}$$

We can now observe the following using the small Omega Bachmann-Landau notation.

$$\lim_{s \rightarrow \infty} \frac{(2s - c)^{-2}}{e^{-s}} = \infty \implies \mathbb{E}_{D,D'} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D')] = \omega(e^{-x})\tag{A.8}$$



Considering the third case, without loss of generality, we consider the case of  $i = j = c + 1$ . We find the following.

$$\begin{aligned}
& \mathbb{E}_{D,D'} [\kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D')] \\
&= \mathbb{E}_{D,D'} \left[ \kappa(x; Z_{c+1}, D_{1:c}) \kappa(x; Z'_{c+1}, D_{1:c}) \kappa(x; Z_{c+1}, D_{(c+1):s}) \kappa(x; Z'_{c+1}, D'_{(c+1):s}) \right] \\
&= \frac{2}{(2s-c)!} \sum_{i=0}^{s-c-1} \binom{s-c-1}{i} (s-1+i)! (s-c-1-i)! \\
&= \frac{2(2s-c-2)!}{(2s-c)!} \sum_{i=0}^{s-c-1} \binom{s-c-1}{i} \binom{2s-c-2}{s-1+i}^{-1} \\
&= \frac{2}{(2s-c)(2s-c-1)} \sum_{i=0}^{s-c-1} \binom{s-c-1}{i} \binom{2s-c-2}{s-1+i}^{-1}
\end{aligned} \tag{A.9}$$

The third case follows from a similar combinatorial logic as the second. We consider without loss of generality the case that  $\|X'_{c+1} - x\| \leq \|X_{c+1} - x\|$  and adjust for this fact by multiplying the whole expression by two. Notice now that any number  $i = 0, \dots, s-c-1$  of observations in  $D'_{(c+2):s}$  can be farther away from  $x$  than  $X_{c+1}$  or at a distance that is between  $\|X'_{c+1} - x\|$  and  $\|X_{c+1} - x\|$ . The summation adjusts for all possible permutations that fulfill this criterion. Furthermore, we can make the following observation.

$$\begin{aligned}
& \mathbb{E}_{D,D'} [\kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D')] \geq \frac{1}{(2s-c)^2} \sum_{i=0}^{s-c-1} \binom{s-c-1}{i} \binom{2s-c-2}{s-1+i}^{-1} \\
&= \frac{1}{(2s-c)^2} \sum_{i=0}^{s-c-1} \left( \frac{(s-c-1)!}{i!(s-c-1-i)!} \cdot \frac{(s-1+i)!(s-c-1-i)!}{(2s-c-2)!} \right) \\
&= \frac{1}{(2s-c)^2} \cdot \frac{(s-c-1)!}{(2s-c-2)!} \sum_{i=0}^{s-c-1} \frac{(s-1+i)!}{i!} \\
&= \frac{1}{(2s-c)^2} \cdot \frac{(s-c-1)!}{(2s-c-2)!} \left( \frac{(2s-c-2)!}{(s-c-1)!} + \sum_{i=0}^{s-c-2} \frac{(s-1+i)!}{i!} \right) \\
&\geq \frac{1}{(2s-c)^2}
\end{aligned} \tag{A.10}$$

We can now observe the following using the small Omega Bachmann-Landau notation.

$$\lim_{s \rightarrow \infty} \frac{(2s-c)^{-2}}{e^{-s}} = \infty \quad \implies \quad \mathbb{E}_{D,D'} [\kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D')] = \omega(e^{-x}) \tag{A.11}$$

■

**Lemma A.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 statements hold.

$$\forall i \in [c] : \quad \mathbb{E} [\kappa(x; Z_i, D) \kappa(x; Z_i, D') \mid X_i] = \{1 - \varphi(B(x, \|X_i - x\|))\}^{2s-c-1} \quad (\text{A.12})$$

$$\begin{aligned} \forall i \in [c] \forall j \in \{c+1, \dots, s\} : \quad & \mathbb{E} [\kappa(x; Z_i, D) \kappa(x; Z'_j, D') \mid X_i, X'_j] \\ &= \mathbb{1}(\|X'_j - x\| \leq \|X_i - x\|) \cdot \{1 - \varphi(B(x, \|X_i - x\|))\}^{s-1} \cdot \{1 - \varphi(B(x, \|X'_j - x\|))\}^{s-c-1} \end{aligned} \quad (\text{A.13})$$

$$\begin{aligned} \forall i, j \in \{c+1, \dots, s\} : \quad & \mathbb{E} [\kappa(x; Z_i, D) \kappa(x; Z'_j, D') \mid X_i, X'_j] \\ &= \{1 - \varphi(B(x, \min(\|X_i - x\|, \|X'_j - x\|)))\}^{s-c-1} \cdot \{1 - \varphi(B(x, \max(\|X_i - x\|, \|X'_j - x\|)))\}^{s-1} \end{aligned} \quad (\text{A.14})$$

*Proof of Lemma A.4.*

Without loss of generality, we will consider the cases of  $i = 1$  for the first equation.

$$\begin{aligned} \mathbb{E} [\kappa(x; Z_1, D) \kappa(x; Z_1, D') \mid X_1] &= \mathbb{E} [\kappa(x; Z_1, D_{[c]}) \kappa(x; Z_1, D_{(c+1):s}) \kappa(x; Z_1, D'_{(c+1):s}) \mid X_1] \\ &= \mathbb{E} [\kappa(x; Z_1, D_{[c]}) \mid X_1] \cdot \mathbb{E} [\kappa(x; Z_1, D_{(c+1):s}) \mid X_1] \cdot \mathbb{E} [\kappa(x; Z_1, D'_{(c+1):s}) \mid X_1] \\ &= \{1 - \varphi(B(x, \|X_1 - x\|))\}^{c-1} \cdot \{1 - \varphi(B(x, \|X_1 - x\|))\}^{s-c} \cdot \{1 - \varphi(B(x, \|X_1 - x\|))\}^{s-c} \\ &= \{1 - \varphi(B(x, \|X_1 - x\|))\}^{2s-c-1} \end{aligned} \quad (\text{A.15})$$

Without loss of generality, we will consider the cases of  $i = 1$  and  $j = c+1$  for the second equation.

$$\begin{aligned} & \mathbb{E} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D') \mid X_1, X'_{c+1}] \\ &= \mathbb{E} [\mathbb{E} [\kappa(x; Z_1, D_{1:c}) \kappa(x; Z_1, D_{(c+1):s}) \kappa(x; Z'_{c+1}, D_{1:c}) \kappa(x; Z'_{c+1}, D'_{(c+1):s}) \mid X_1, \dots, X_c, X'_{c+1}] \mid X_1, X'_{c+1}] \\ &= \mathbb{E} [\mathbb{E} [\kappa(x; Z_1, D_{(c+1):s}) \cdot \kappa(x; Z'_{c+1}, D'_{(c+1):s}) \mid X_1, \dots, X_c, X'_{c+1}] \cdot \kappa(x; Z_1, D_{1:c}) \cdot \kappa(x; Z'_{c+1}, D_{1:c}) \mid X_1, X'_{c+1}] \\ &= \mathbb{E} [\mathbb{E} [\kappa(x; Z_1, D_{(c+1):s}) \cdot \kappa(x; Z'_{c+1}, D'_{(c+1):s}) \mid X_1, X'_{c+1}] \cdot \kappa(x; Z_1, D_{1:c}) \cdot \kappa(x; Z'_{c+1}, D_{1:c}) \mid X_1, X'_{c+1}] \\ &= \mathbb{E} [\kappa(x; Z_1, D_{(c+1):s}) \cdot \kappa(x; Z'_{c+1}, D'_{(c+1):s}) \mid X_1, X'_{c+1}] \cdot \mathbb{E} [\kappa(x; Z_1, D_{1:c}) \cdot \kappa(x; Z'_{c+1}, D_{1:c}) \mid X_1, X'_{c+1}] \\ &= \mathbb{E} [\kappa(x; Z_1, D_{(c+1):s}) \mid X_1] \cdot \mathbb{E} [\kappa(x; Z'_{c+1}, D'_{(c+1):s}) \mid X'_{c+1}] \cdot \mathbb{1}(\|X'_{c+1} - x\| \leq \|X_1 - x\|) \cdot \mathbb{E} [\kappa(x; Z_1, D_{1:c}) \mid X_1] \\ &= \mathbb{1}(\|X'_{c+1} - x\| \leq \|X_1 - x\|) \cdot \mathbb{E} [\kappa(x; Z_1, D) \mid X_1] \cdot \mathbb{E} [\kappa(x; Z'_{c+1}, D'_{(c+1):s}) \mid X'_{c+1}] \\ &= \mathbb{1}(\|X'_{c+1} - x\| \leq \|X_1 - x\|) \cdot \{1 - \varphi(B(x, \|X_1 - x\|))\}^{s-1} \cdot \{1 - \varphi(B(x, \|X'_{c+1} - x\|))\}^{s-c-1} \end{aligned} \quad (\text{A.16})$$

For the third case, without loss of generality, we consider the case of  $i = j = c + 1$ .

$$\begin{aligned}
& \mathbb{E} [\kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D') \mid X_{c+1}, X'_{c+1}] \\
&= \mathbb{E} [\mathbb{E} [\kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D') \mid X_1, \dots, X_c, X_{c+1}, X'_{c+1}] \mid X_{c+1}, X'_{c+1}] \\
&= \mathbb{E} [\mathbb{E} [\kappa(x; Z_{c+1}, D_{1:(c+1)}) \kappa(x; Z'_{c+1}, D'_{1:(c+1)}) \\
&\quad \kappa(x; Z_{c+1}, D_{(c+1):s}) \kappa(x; Z'_{c+1}, D'_{(c+1):s}) \mid X_1, \dots, X_c, X_{c+1}, X'_{c+1}] \mid X_{c+1}, X'_{c+1}] \\
&= \mathbb{E} [\mathbb{E} [\kappa(x; Z_{c+1}, D_{1:(c+1)}) \kappa(x; Z'_{c+1}, D'_{1:(c+1)}) \mid X_1, \dots, X_c, X_{c+1}, X'_{c+1}] \\
&\quad \kappa(x; Z_{c+1}, D_{(c+1):s}) \kappa(x; Z'_{c+1}, D'_{(c+1):s}) \mid X_{c+1}, X'_{c+1}] \\
&= \mathbb{E} [\kappa(x; Z_{c+1}, D_{1:(c+1)}) \kappa(x; Z'_{c+1}, D'_{1:(c+1)}) \mid X_{c+1}, X'_{c+1}] \\
&\quad \cdot \mathbb{E} [\kappa(x; Z_{c+1}, D_{(c+1):s}) \mid X_{c+1}] \cdot \mathbb{E} [\kappa(x; Z'_{c+1}, D'_{(c+1):s}) \mid X'_{c+1}]
\end{aligned} \tag{A.17}$$

Without loss of generality, consider the case that  $\|X_{c+1} - x\| \leq \|X'_{c+1} - x\|$ .

$$\mathbb{E} [\kappa(x; Z_{c+1}, D_{1:(c+1)}) \kappa(x; Z'_{c+1}, D'_{1:(c+1)}) \mid X_{c+1}, X'_{c+1}] = \mathbb{E} [\kappa(x; Z'_{c+1}, D'_{1:(c+1)}) \mid X'_{c+1}] \tag{A.18}$$

Furthermore, observe the following.

$$\mathbb{E} [\kappa(x; Z'_{c+1}, D'_{1:(c+1)}) \mid X'_{c+1}] \cdot \mathbb{E} [\kappa(x; Z'_{c+1}, D'_{(c+1):s}) \mid X'_{c+1}] = \mathbb{E} [\kappa(x; Z'_{c+1}, D') \mid X'_{c+1}] \tag{A.19}$$

Thus, we can find the following.

$$\begin{aligned}
& \mathbb{E} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D') \mid X_1, X'_{c+1}] \\
&= \mathbb{1} (\|X'_{c+1} - x\| \leq \|X_{c+1} - x\|) \cdot \{1 - \varphi(B(x, \|X_{c+1} - x\|))\}^{s-1} \cdot \{1 - \varphi(B(x, \|X'_{c+1} - x\|))\}^{s-c-2} \\
&\quad + \mathbb{1} (\|X'_{c+1} - x\| > \|X_{c+1} - x\|) \cdot \{1 - \varphi(B(x, \|X_{c+1} - x\|))\}^{s-c-1} \cdot \{1 - \varphi(B(x, \|X'_{c+1} - x\|))\}^{s-1} \\
&= \{1 - \varphi(B(x, \min(\|X_{c+1} - x\|, \|X'_{c+1} - x\|)))\}^{s-c-1} \cdot \{1 - \varphi(B(x, \max(\|X_{c+1} - x\|, \|X'_{c+1} - x\|)))\}^{s-1}
\end{aligned} \tag{A.20}$$

■

**Lemma A.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  $D$ .

For any  $L^2(\mathcal{X})$  function  $f$  that is continuous at  $x$ , it holds that

$$\lim_{s \rightarrow \infty} \underbrace{\mathbb{E}_1 [f^2(X_1)(2s - c) \mathbb{E}_{2:s} [\kappa(x; Z_1, D) \kappa(x; Z_1, D')]]}_{(A)} = f^2(x) \quad (\text{A.21})$$

$$\lim_{s \rightarrow \infty} \underbrace{\mathbb{E}_{1, (c+1)'} \left[ f(X_1) f(X'_{c+1}) \cdot \frac{\mathbb{E}_{D, D'} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D') \mid Z_1, Z'_{c+1}]}{\mathbb{E}_{D, D'} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D')]} \right]}_{(B)} = f^2(x) \quad (\text{A.22})$$

$$\lim_{s \rightarrow \infty} \underbrace{\mathbb{E}_{c+1} \left[ f(X_{c+1}) f(X'_{c+1}) \cdot \frac{\mathbb{E}_{D, D'} [\kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D') \mid Z_{c+1}, Z'_{c+1}]}{\mathbb{E}_{D, D'} [\kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D')]} \right]}_{(C)} = f^2(x) \quad (\text{A.23})$$

*Proof of Lemma A.5.*

We will largely argue along the same lines as the original proof in Demirkaya et al. (2024). Thus, consider first the following inequalities.

$$\begin{aligned} |(A) - f^2(x)| &= |\mathbb{E}_1 [f^2(X_1)(2s - c) \mathbb{E}_{2:s} [\kappa(x; Z_1, D) \kappa(x; Z_1, D')]] - f^2(x)| \\ &\leq \mathbb{E}_1 [|f^2(X_1) - f^2(x)| (2s - c) \mathbb{E}_{2:s} [\kappa(x; Z_1, D) \kappa(x; Z_1, D')]] \end{aligned} \quad (\text{A.24})$$

$$\begin{aligned} |(B) - f^2(x)| &= \left| \mathbb{E}_{1, (c+1)'} \left[ f(X_1) f(X'_{c+1}) \cdot \frac{\mathbb{E}_{D, D'} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D') \mid Z_1, Z'_{c+1}]}{\mathbb{E}_{D, D'} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D')]} \right] - f^2(x) \right| \\ &\leq \mathbb{E}_{1, (c+1)'} \left[ |f(X_1) f(X'_{c+1}) - f^2(x)| \cdot \frac{\mathbb{E}_{D, D'} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D') \mid Z_1, Z'_{c+1}]}{\mathbb{E}_{D, D'} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D')]} \right] \end{aligned} \quad (\text{A.25})$$

$$\begin{aligned} |(C) - f^2(x)| &= \left| \mathbb{E}_{c+1} \left[ f(X_{c+1}) f(X'_{c+1}) \cdot \frac{\mathbb{E}_{D, D'} [\kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D') \mid Z_{c+1}, Z'_{c+1}]}{\mathbb{E}_{D, D'} [\kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D')]} \right] - f^2(x) \right| \\ &\leq \mathbb{E}_{c+1} \left[ |f(X_{c+1}) f(X'_{c+1}) - f^2(x)| \cdot \frac{\mathbb{E}_{D, D'} [\kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D') \mid Z_{c+1}, Z'_{c+1}]}{\mathbb{E}_{D, D'} [\kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D')]} \right] \end{aligned} \quad (\text{A.26})$$

Now, fix an arbitrary  $\epsilon > 0$ . By continuity of  $f$  at  $x$ , there exists a  $\delta > 0$ , such that the following holds.

$$\forall X, X' \in B(x, \delta) : \quad |f(X) \cdot f(X') - f^2(x)| < \epsilon \quad (\text{A.27})$$

We can consider decompositions of these terms in analogy to Demirkaya et al. (2024), i.e. by considering cases with

observations lying within this sphere or outside of it, and observe the following.

$$\begin{aligned}
& \mathbb{E}_1 \left[ |f^2(X_1) - f^2(x)| (2s - c) \mathbb{E}_{2:s} [\kappa(x; Z_1, D) \kappa(x; Z_1, D') \mathbb{1}(X_1 \in B(x, \delta))] \right] \\
& \leq \epsilon \cdot \mathbb{E}_1 [(2s - c) \mathbb{E}_{2:s} [\kappa(x; Z_1, D) \kappa(x; Z_1, D') \mathbb{1}(X_1 \in B(x, \delta))] ] \\
& \leq \epsilon \cdot \mathbb{E}_1 [(2s - c) \mathbb{E}_{2:s} [\kappa(x; Z_1, D) \kappa(x; Z_1, D')]] = \epsilon
\end{aligned} \tag{A.28}$$

$$\begin{aligned}
& \mathbb{E}_{1, (c+1)'} \left[ |f(X_1)f(X'_{c+1}) - f^2(x)| \cdot \frac{\mathbb{E}_{D, D'} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D') \mid Z_1, Z'_{c+1}]}{\mathbb{E}_{D, D'} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D')]} \mathbb{1}(X_1, X'_{c+1} \in B(x, \delta)) \right] \\
& \leq \epsilon \cdot \mathbb{E}_{1, (c+1)'} \left[ \frac{\mathbb{E}_{D, D'} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D') \mid Z_1, Z'_{c+1}]}{\mathbb{E}_{D, D'} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D')]} \mathbb{1}(X_1, X'_{c+1} \in B(x, \delta)) \right] \\
& \leq \epsilon \cdot \mathbb{E}_{1, (c+1)'} \left[ \frac{\mathbb{E}_{D, D'} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D') \mid Z_1, Z'_{c+1}]}{\mathbb{E}_{D, D'} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D')]} \right] = \epsilon
\end{aligned} \tag{A.29}$$

$$\begin{aligned}
& \mathbb{E}_{c+1} \left[ |f(X_{c+1})f(X'_{c+1}) - f^2(x)| \cdot \frac{\mathbb{E}_{D, D'} [\kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D') \mid Z_{c+1}, Z'_{c+1}]}{\mathbb{E}_{D, D'} [\kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D')]} \mathbb{1}(X_{c+1}, X'_{c+1} \in B(x, \delta)) \right] \\
& \leq \epsilon \cdot \mathbb{E}_{c+1} \left[ \frac{\mathbb{E}_{D, D'} [\kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D') \mid Z_{c+1}, Z'_{c+1}]}{\mathbb{E}_{D, D'} [\kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D')]} \mathbb{1}(X_{c+1}, X'_{c+1} \in B(x, \delta)) \right] \\
& \leq \epsilon \cdot \mathbb{E}_{c+1} \left[ \frac{\mathbb{E}_{D, D'} [\kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D') \mid Z_{c+1}, Z'_{c+1}]}{\mathbb{E}_{D, D'} [\kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D')]} \right] = \epsilon
\end{aligned} \tag{A.30}$$

Considering next the parts of the expectation that are not covered by the previous cases, we can find the following. As in the original proof, we use the fact that if  $X$  or  $X'$  do not lie within  $B(x, \delta)$ , then the following holds

$$B(x, \delta) \subseteq B(x, \max(\|X - x\|, \|X' - x\|)). \tag{A.31}$$

This allows us to find the following.

$$\begin{aligned}
& \mathbb{E}_1 \left[ |f^2(X_1) - f^2(x)| (2s - c) \mathbb{E}_{2:s} [\kappa(x; Z_1, D) \kappa(x; Z_1, D') (1 - \mathbb{1}(X_1 \in B(x, \delta)))] \right] \\
& \leq \mathbb{E}_1 \left[ |f^2(X_1) - f^2(x)| (2s - c) \{1 - \varphi(B(x, \delta))\}^{2s-c-1} (1 - \mathbb{1}(X_1 \in B(x, \delta))) \right] \\
& \leq (2s - c) 1 \cdot \{1 - \varphi(B(x, \delta))\}^{2s-c-1} \cdot \mathbb{E}_1 [|f^2(X_1) - f^2(x)|]
\end{aligned} \tag{A.32}$$

In the second case, first recall the form of the conditional expectation from Lemma A.4.

$$\begin{aligned}
& \mathbb{E} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D') \mid X_1, X'_{c+1}] \\
& = \mathbb{1}(\|X'_{c+1} - x\| \leq \|X_1 - x\|) \cdot \{1 - \varphi(B(x, \|X_1 - x\|))\}^{s-1} \cdot \{1 - \varphi(B(x, \|X'_{c+1} - x\|))\}^{s-c-1}
\end{aligned} \tag{A.33}$$

The indicator variable in this expression is only non-zero if  $\max(\|X_1 - x\|, \|X'_{c+1} - x\|) = \|X_1 - x\|$ . Thus, in light of the conditioning, we can observe the following.

$$B(x, \delta) \subseteq B(x, \|X_1 - x\|) \tag{A.34}$$

Thus, we can make the following observation.

$$\begin{aligned}
& \mathbb{E}_{1,(c+1)'} \left[ \left| f(X_1)f(X'_{c+1}) - f^2(x) \right| \cdot \frac{\mathbb{E}_{D,D'} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D') \mid Z_1, Z'_{c+1}]}{\mathbb{E}_{D,D'} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D')]} (1 - \mathbb{1}(X_1, X'_{c+1} \in B(x, \delta))) \right] \\
& \stackrel{\text{Lem A.3}}{\leq} (2s - c)^2 \cdot \mathbb{E}_{1,(c+1)'} \left[ \left| f(X_1)f(X'_{c+1}) - f^2(x) \right| \right. \\
& \quad \cdot \mathbb{E}_{D,D'} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D') \mid Z_1, Z'_{c+1}] (1 - \mathbb{1}(X_1, X'_{c+1} \in B(x, \delta))) \left. \right] \\
& \stackrel{\text{Lem A.4}}{=} (2s - c)^2 \cdot \mathbb{E}_{1,(c+1)'} \left[ \left| f(X_1)f(X'_{c+1}) - f^2(x) \right| \cdot \mathbb{1}(\|X'_{c+1} - x\| \leq \|X_1 - x\|) \right. \\
& \quad \cdot \{1 - \varphi(B(x, \|X_1 - x\|))\}^{s-1} \cdot \{1 - \varphi(B(x, \|X'_{c+1} - x\|))\}^{s-c-1} \cdot (1 - \mathbb{1}(X_1, X'_{c+1} \in B(x, \delta))) \left. \right] \\
& \leq (2s - c)^2 \cdot \{1 - \varphi(B(x, \delta))\}^{s-1} \cdot \mathbb{E}_{1,(c+1)'} \left[ \left| f(X_1)f(X'_{c+1}) - f^2(x) \right| \cdot \mathbb{1}(\delta < \|X'_{c+1} - x\| \leq \|X_1 - x\|) \right] \\
& \leq (2s - c)^2 \cdot \{1 - \varphi(B(x, \delta))\}^{s-1} \cdot \mathbb{E}_{1,(c+1)'} \left[ \left| f(X_1)f(X'_{c+1}) - f^2(x) \right| \right]
\end{aligned} \tag{A.35}$$

Similarly, considering the third case, we observe the following.

$$\begin{aligned}
& \mathbb{E}_{c+1} \left[ \left| f(X_{c+1})f(X'_{c+1}) - f^2(x) \right| \cdot \frac{\mathbb{E}_{D,D'} [\kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D') \mid Z_{c+1}, Z'_{c+1}]}{\mathbb{E}_{D,D'} [\kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D')]} (1 - \mathbb{1}(X_{c+1}, X'_{c+1} \in B(x, \delta))) \right] \\
& \stackrel{\text{Lem A.3}}{\leq} (2s - c)^2 \cdot \mathbb{E}_{c+1} \left[ \left| f(X_{c+1})f(X'_{c+1}) - f^2(x) \right| \right. \\
& \quad \cdot \mathbb{E}_{D,D'} [\kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D') \mid Z_{c+1}, Z'_{c+1}] (1 - \mathbb{1}(X_{c+1}, X'_{c+1} \in B(x, \delta))) \left. \right] \\
& \stackrel{\text{Lem A.4}}{=} (2s - c)^2 \cdot \mathbb{E}_{c+1} \left[ \left| f(X_{c+1})f(X'_{c+1}) - f^2(x) \right| \cdot \{1 - \varphi(B(x, \min(\|X_{c+1} - x\|, \|X'_{c+1} - x\|)))\}^{s-c-1} \right. \\
& \quad \cdot \{1 - \varphi(B(x, \max(\|X_{c+1} - x\|, \|X'_{c+1} - x\|)))\}^{s-1} \cdot (1 - \mathbb{1}(X_{c+1}, X'_{c+1} \in B(x, \delta))) \left. \right] \\
& \leq (2s - c)^2 \cdot \mathbb{E}_{c+1} \left[ \left| f(X_{c+1})f(X'_{c+1}) - f^2(x) \right| \cdot \{1 - \varphi(B(x, \min(\|X_{c+1} - x\|, \|X'_{c+1} - x\|)))\}^{s-c-1} \right. \\
& \quad \cdot \{1 - \varphi(B(x, \delta))\}^{s-1} \cdot (1 - \mathbb{1}(X_{c+1}, X'_{c+1} \in B(x, \delta))) \left. \right] \\
& \leq (2s - c)^2 \cdot \{1 - \varphi(B(x, \delta))\}^{s-1} \cdot \mathbb{E}_{c+1} \left[ \left| f(X_{c+1})f(X'_{c+1}) - f^2(x) \right| \right]
\end{aligned} \tag{A.36}$$

Concerning the resulting terms in these three expressions, we can then make the following observations.

$$\mathbb{E} \left[ \left| f^2(X) - f^2(x) \right| \right] \leq \mathbb{E} \left[ f^2(X) \right] + f^2(x) = \|f\|_{L_2}^2 + f^2(x) \tag{A.37}$$

$$\mathbb{E} \left[ \left| f(X)f(X') - f^2(x) \right| \right] \leq \mathbb{E} \left[ \left| f(X)f(X') \right| \right] + f^2(x) \leq \mathbb{E} \left[ \left| f(X) \right| \cdot \left| f(X') \right| \right] + f^2(x) = \|f\|_{L_1}^2 + f^2(x) \tag{A.38}$$

As  $f$  is an  $L^2(\mathcal{X})$  function on a bounded domain, observe that  $\|f\|_{L^1}$  is finite. Thus, we can find the following.

$$\mathbb{E}_1 \left[ \left| f^2(X_1) - f^2(x) \right| (2s - c) \mathbb{E}_{2,s} [\kappa(x; Z_1, D) \kappa(x; Z_1, D') (1 - \mathbb{1}(X_1 \in B(x, \delta)))] \right] \longrightarrow 0 \quad \text{as } s \rightarrow \infty \tag{A.39}$$

$$\begin{aligned}
& \mathbb{E}_{1,(c+1)'} \left[ \left| f(X_1)f(X'_{c+1}) - f^2(x) \right| \cdot \frac{\mathbb{E}_{D,D'} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D') \mid Z_1, Z'_{c+1}]}{\mathbb{E}_{D,D'} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D')]} (1 - \mathbb{1}(X_1, X'_{c+1} \in B(x, \delta))) \right] \\
& \longrightarrow 0 \quad \text{as } s \rightarrow \infty
\end{aligned} \tag{A.40}$$

$$\mathbb{E}_{c+1} \left[ \left| f(X_{c+1})f(X'_{c+1}) - f^2(x) \right| \cdot \frac{\mathbb{E}_{D,D'} [\kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D') \mid Z_{c+1}, Z'_{c+1}]}{\mathbb{E}_{D,D'} [\kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D')]} (1 - \mathbb{1}(X_{c+1}, X'_{c+1} \in B(x, \delta))) \right] \rightarrow 0 \quad \text{as } s \rightarrow \infty \quad (\text{A.41})$$

Combining these findings, for large enough  $s$  we can bound the terms  $\|(A) - f^2(x)\|$ ,  $\|(B) - f^2(x)\|$ , and  $\|(C) - f^2(x)\|$  by  $2\epsilon$ , respectively. As  $\epsilon$  was arbitrary this concludes the proof.  $\blacksquare$

**Lemma A.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 hold.

$$\forall i \in [c] \forall j \in \{c+1, \dots, s\} : \quad \mathbb{E}_{D,D'} [\kappa(x; Z_i, D) \kappa(x; Z'_j, D')] \leq \frac{s}{(2s-c)(2s-c-1)(c+1)} \quad (\text{A.42})$$

$$\forall i, j \in \{c+1, \dots, s\} : \quad \mathbb{E}_{D,D'} [\kappa(x; Z_i, D) \kappa(x; Z'_j, D')] \leq \frac{2(s-c)}{(2s-c)^2(2s-c-1)} \quad (\text{A.43})$$

*Proof of Lemma A.6.*

Recall the results of Lemma A.3 and make the following observations.

$$\begin{aligned} \mathbb{E}_{D,D'} [\kappa(x; Z_i, D) \kappa(x; Z'_j, D')] &= \frac{1}{(2s-c)(2s-c-1)} \sum_{i=0}^{s-c-1} \binom{s-c-1}{i} \binom{2s-c-2}{i}^{-1} \\ &\leq \frac{1}{(2s-c)(2s-c-1)} \sum_{i=0}^{s-c-1} \frac{(s-c-1)^i}{i!} \cdot \frac{i!}{(2s-c-1-i)^i} = \frac{1}{(2s-c)(2s-c-1)} \sum_{i=0}^{s-c-1} \left( \frac{s-c-1}{2s-c-1-i} \right)^i \\ &\leq \frac{1}{(2s-c)(2s-c-1)} \sum_{i=0}^{s-c-1} \left( \frac{s-c-1}{2s-c-1-i} \right)^i \leq \frac{1}{(2s-c)(2s-c-1)} \sum_{i=0}^{s-c-1} \left( \frac{s-c-1}{s} \right)^i \\ &\leq \frac{1}{(2s-c)(2s-c-1)} \sum_{i=0}^{\infty} \left( \frac{s-c-1}{s} \right)^i = \frac{s}{(2s-c)(2s-c-1)(c+1)} \end{aligned} \quad (\text{A.44})$$

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

$$\begin{aligned}
\mathbb{E}_{D,D'} [\kappa(x; Z_i, D) \kappa(x; Z'_j, D')] &= \frac{2}{(2s-c)(2s-c-1)} \sum_{i=0}^{s-c-1} \binom{s-c-1}{i} \binom{2s-c-1}{s-1+i}^{-1} \\
&= \frac{2}{(2s-c)(2s-c-1)} \sum_{i=0}^{s-c-1} \binom{s-c-1}{i} \binom{2s-c-1}{s-c-i}^{-1} \\
&= \frac{2}{(2s-c)(2s-c-1)} \sum_{i=0}^{s-c-1} \frac{(s-c-1)!}{(s-c-1-i)!i!} \frac{(s-1+i)!(s-c-i)!}{(2s-c-1)!} \\
&= \frac{2}{(2s-c)(2s-c-1)} \cdot \frac{(s-c-1)!(s-1)!}{(2s-c-1)!} \cdot \sum_{i=0}^{s-c-1} (s-c-i) \binom{s-1+i}{i} \\
&\leq \frac{2}{(2s-c)(2s-c-1)} \cdot \frac{(s-c)!(s-1)!}{(2s-c-1)!} \cdot \sum_{i=0}^{s-c-1} \binom{s-1+i}{i} \\
&= \frac{2}{(2s-c)(2s-c-1)} \cdot \binom{2s-c}{s-c}^{-1} \cdot \binom{2s-c-1}{s-c-1} = \frac{2(s-c)}{(2s-c)^2(2s-c-1)}
\end{aligned} \tag{A.45}$$

■

---

**Lemma A.7** (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

$$\sum [X_i + Y_i]^2 \xrightarrow{p} 1 \quad \text{and} \quad \mathbb{E} \left[ \sum (X_i + Y_i)^2 \right] \rightarrow 1. \tag{A.46}$$

---

**Lemma A.8** (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  $D_{\ell, -i} = \{Z_l \mid l \in \ell \setminus \{i\}\}$ .



## A.2 Consequences of DDML Rate-Assumptions

---

Recall first the conditions imposed in Assumption 8.

$$r'_n := \sup_{\eta \in \mathcal{T}_n} \left( \mathbb{E}_Z \left[ |m(Z; \theta_0, \eta) - m(Z; \theta_0, \eta_0)|^2 \right] \right)^{1/2} \leq \delta_n \quad (\text{A.47})$$

$$\lambda'_n := \sup_{r \in (0,1), \eta \in \mathcal{T}_n} \left| \partial_r^2 \mathbb{E}_Z [m(Z; \theta_0, \eta_0 + r(\eta - \eta_0))] \right| \leq \delta_n / \sqrt{n} \quad (\text{A.48})$$

We are interested in how these conditions translate into statements on the Oracle-Hoeffding decomposition errors. Thus, recall the form of the error terms as introduced in Equation 4.22.

$$R_c(x; \mathbf{D}_\ell) = \chi_s^{(c)}(x; \mathbf{D}_\ell, \hat{\eta}) - \chi_{s,0}^{(c)}(x; \mathbf{D}_\ell) \quad (\text{A.49})$$

More specifically, we are interested in the part absent the centering term as these mostly cancel out across the error terms of different order.

$$\tilde{R}_c(x; \mathbf{D}_\ell) = R_c(x; \mathbf{D}_\ell) - (-1)^c \cdot (\mathbb{E}_D [\chi_s(x; \mathbf{D}_{[s]}, \hat{\eta}) - \chi_{s,0}(x; \mathbf{D}_{[s]})]) \quad (\text{A.50})$$

---

**Lemma A.9.**

*LOREM IPSUM*

*Proof of Lemma A.9.*

Consider the following argument.

$$\begin{aligned} & \sup_{\eta \in \mathcal{T}_n} \left( \mathbb{E}_Z \left[ \left| \chi_s^{(1)}(x; Z, \eta) - \chi_{s,0}^{(1)}(x; Z) + (\mathbb{E}_D [\chi_s(x; \mathbf{D}_{[s]}, \hat{\eta}) - \chi_{s,0}(x; \mathbf{D}_{[s]})]) \right|^2 \right] \right)^{1/2} \\ &= \sup_{\eta \in \mathcal{T}_n} \left( \mathbb{E}_Z \left[ \left| \vartheta_s^1(x; Z, \eta) - \vartheta_{s,0}^1(x; Z) \right|^2 \right] \right)^{1/2} \\ &= \sup_{\eta \in \mathcal{T}_n} \left( \mathbb{E}_Z \left[ \left| \frac{m(Z, \theta_0, \eta) - m(Z, \theta_0, \eta_0)}{s!} \cdot \mathbb{E}_D [\kappa(x; Z_1, \mathbf{D}_{[s]}) \mid Z_1 = Z] \right. \right. \right. \\ & \quad \left. \left. + \frac{s-1}{s!} \cdot \mathbb{E}_D [\kappa(x; Z_2, \mathbf{D}_{[s]}) \cdot (m(Z_2, \theta_0, \eta) - m(Z_2, \theta_0, \eta_0)) \mid Z_1 = Z] \right|^2 \right] \right)^{1/2} \\ & \stackrel{\Delta\text{-ineq.}}{\leq} \sup_{\eta \in \mathcal{T}_n} \left\{ \frac{1}{s!} \left( \mathbb{E}_Z \left[ \left| (m(Z, \theta_0, \eta) - m(Z, \theta_0, \eta_0)) \cdot \mathbb{E}_D [\kappa(x; Z_1, \mathbf{D}_{[s]}) \mid Z_1 = Z] \right|^2 \right] \right)^{1/2} \right. \\ & \quad \left. + \frac{s-1}{s!} \left( \mathbb{E}_Z \left[ \left| \mathbb{E}_D [\kappa(x; Z_2, \mathbf{D}_{[s]}) \cdot (m(Z_2, \theta_0, \eta) - m(Z_2, \theta_0, \eta_0)) \mid Z_1 = Z] \right|^2 \right] \right)^{1/2} \right\} \end{aligned} \quad (\text{A.51})$$

Continuing from this step, we can further observe the following.

$$\begin{aligned}
& \sup_{\eta \in \mathcal{T}_n} \left( \mathbb{E}_Z \left[ \left| \chi_s^{(1)}(x; Z, \eta) - \chi_{s,0}^{(1)}(x; Z) + (\mathbb{E}_D [\chi_s(x; \mathbf{D}_{[s]}, \hat{\eta}) - \chi_{s,0}(x; \mathbf{D}_{[s]})] \right|^2 \right] \right)^{1/2} \\
& \leq \frac{1}{s!} \sup_{\eta \in \mathcal{T}_n} \left( \mathbb{E}_Z \left[ (m(Z, \theta_0, \eta) - m(Z, \theta_0, \eta_0))^2 \cdot \mathbb{E}_D [\kappa(x; Z_1, \mathbf{D}_{[s]} | Z_1 = Z)]^2 \right] \right)^{1/2} \\
& \quad + \frac{s-1}{s!} \sup_{\eta \in \mathcal{T}_n} \left( \mathbb{E}_Z \left[ \mathbb{E}_D [\kappa(x; Z_2, \mathbf{D}_{[s]}) \cdot (m(Z_2, \theta_0, \eta) - m(Z_2, \theta_0, \eta_0)) | Z_1 = Z]^2 \right] \right)^{1/2} \\
& \leq \frac{1}{s!} \sup_{\eta \in \mathcal{T}_n} \left( \mathbb{E}_Z \left[ (m(Z, \theta_0, \eta) - m(Z, \theta_0, \eta_0))^2 \right] \right)^{1/2} + \frac{s-1}{s!} \sup_{\eta \in \mathcal{T}_n} \left( \mathbb{E}_Z \left[ \mathbb{E}_D [m(Z_2, \theta_0, \eta) - m(Z_2, \theta_0, \eta_0) | Z_1 = Z]^2 \right] \right)^{1/2} \\
& \leq \frac{\delta_n}{(s-1)!}
\end{aligned} \tag{A.52}$$

LOREM IPSUM

■

**Lemma A.10.**

LOREM IPSUM

*Proof of Lemma A.10.*

$$\begin{aligned}
& \sup_{r \in (0,1), \eta \in \mathcal{T}_n} \left| \partial_r^2 \mathbb{E}_Z [\vartheta_s^1(x; Z, \eta_0 + r(\eta - \eta_0))] \right| = \sup_{r \in (0,1), \eta \in \mathcal{T}_n} \left| \partial_r^2 \mathbb{E}_Z [\mathbb{E}_D [\chi_s(x; \mathbf{D}_{[s]}, \eta_0 + r(\eta - \eta_0)) | Z_1 = Z]] \right| \\
& = \sup_{r \in (0,1), \eta \in \mathcal{T}_n} \left| \partial_r^2 \mathbb{E}_D [\chi_s(x; \mathbf{D}_{[s]}, \eta_0 + r(\eta - \eta_0))] \right| \\
& = \frac{1}{s!} \cdot \sup_{r \in (0,1), \eta \in \mathcal{T}_n} \left| \partial_r^2 \mathbb{E}_D [m(Z, \theta_0, \eta_0 + r(\eta - \eta_0)) \cdot s\kappa(x; Z, \mathbf{D}_{[s]})] \right| \\
& = \frac{1}{s!} \cdot \sup_{r \in (0,1), \eta \in \mathcal{T}_n} \left| \partial_r^2 \mathbb{E}_1 [m(Z_1, \theta_0, \eta_0 + r(\eta - \eta_0)) \cdot s\mathbb{E}_{2:s} [\kappa(x; Z_1, \mathbf{D}_{[s]})]] \right| \\
& = \text{LOREMIPSUM}
\end{aligned} \tag{A.53}$$

effectively, we need a decaying second derivative in a shrinking neighborhood of  $x$ . This might be easier to do with a slightly weaker but more explicit rate condition that I would need to show for the CATE moment. I might be able to do something with DCT here! LOREM IPSUM

■

## B Proofs for Results in Section 5

### B.1 Closed Form Representations

*Proof of Theorem 5.1.*

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

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

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  $s$ .

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

Even more simple, we can use the approximate weights to find the following representation. For this purpose recall that  $\alpha_s = s/n$  and define  $\tilde{\alpha}_s = s/(n-1)$ . Thus,  $\tilde{\alpha}_s = \frac{n}{n-1} \alpha_s$ .

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

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

LOREM IPSUM

■

## B.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 B.1** (NPR - DNN Kernel Expectation).

Let  $x$  denote a point of interest. Then

$$\mathbb{E}_D [h_s(x; D)] = \mathbb{E}_1 [Y_1 s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] \longrightarrow \mu(x) \quad \text{as } s \rightarrow \infty \quad (\text{B.4})$$


---

*Proof of Lemma B.1.* This result follows immediately from Lemma A.2 and the following observation.

$$\begin{aligned} \mathbb{E}_1 [Y_1 s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] &= \mathbb{E}_1 [(\mu(X_1) + \varepsilon_1) s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] \\ &= \mathbb{E}_1 [(\mu(X_1) + \mathbb{E}[\varepsilon_1 | X_1]) s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] \\ &= \mathbb{E}_1 [\mu(X_1) s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] \xrightarrow{(\text{Lem A.2})} \mu(x) \quad \text{as } s \rightarrow \infty \end{aligned} \quad (\text{B.5})$$

■

---

**Lemma B.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

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


---

*Proof of Lemma B.2.*

$$\begin{aligned} \psi_s^1(x; z_1) &= \mathbb{E}_D [h_s(x; D) | Z_1 = z_1] = \mathbb{E}_D \left[ \sum_{i=1}^s \kappa(x; Z_i, D) Y_i | Z_1 = z_1 \right] \\ &= \mathbb{E}_D \left[ (\mu(x_1) + \varepsilon_1) \kappa(x; Z_1, D) + \sum_{i=2}^s \kappa(x; Z_i, D) \mu(X_i) | Z_1 = z_1 \right] \\ &= \varepsilon_1 \mathbb{E}_D [\kappa(x; Z_1, D) | X_1 = x_1] + \mathbb{E}_D \left[ \sum_{i=2}^s \kappa(x; Z_i, D) \mu(X_i) | X_1 = x_1 \right] \end{aligned} \quad (\text{B.7})$$

■

### B.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 B.3** (CATE - DNN Kernel Expectation).

Let  $x$  denote a point of interest. Then

$$\begin{aligned}\mathbb{E}_D [h_s(x; D)] &= \mathbb{E}_1 [m(Z_1, \eta_0) s\mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] \\ &\longrightarrow \theta_0(x) \quad \text{as } s \rightarrow \infty\end{aligned}\tag{B.8}$$

*Proof of Lemma B.3.* This result follows immediately from Lemma A.2 and the following observation.

$$\begin{aligned}\mathbb{E}_1 [m(Z_1; \eta_0) s\mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] &= \mathbb{E}_1 [(\mu_0^1(X_1) - \mu_0^0(X_1) + \beta(W_1, X_1) \varepsilon_1) s\mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] \\ &= \mathbb{E}_1 [(\mu_0^1(X_1) - \mu_0^0(X_1) + \beta(W_1, X_1) \mathbb{E}[\varepsilon_1 | X_1]) s\mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] \\ &= \mathbb{E}_1 [(\mu_0^1(X_1) - \mu_0^0(X_1)) s\mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] \\ &\stackrel{(\text{Lem A.2})}{\longrightarrow} \mu_0^1(x) - \mu_0^0(x) = \theta_0(x) \quad \text{as } s \rightarrow \infty\end{aligned}\tag{B.9}$$

■

**Lemma B.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

$$\psi_s^1(x; z_1) = \beta(W_1, X_1) \varepsilon_1 \cdot \mathbb{E}[\kappa(x; Z_1, D) | X_1 = x_1] + \mathbb{E}_D \left[ \sum_{i=2}^s \kappa(x; Z_i, D) (\mu_0^1(X_i) - \mu_0^0(X_i)) \mid Z_1 = z_1 \right]\tag{B.10}$$

*Proof of Lemma B.4.*

$$\begin{aligned}\psi_s^1(x; z_1) &= \mathbb{E}_D [\chi_{s,0}(x; D) | Z_1 = z_1] \\ &= \mathbb{E}_D \left[ \sum_{i=1}^s \kappa(x; Z_i, D) m(Z_i, \eta_0) \mid Z_1 = z_1 \right] \\ &= (\mu_0^1(X_1) - \mu_0^0(X_1) + \beta(W_1, X_1) \varepsilon_1) \mathbb{E}[\kappa(x; Z_1, D) | X_1 = x_1] \\ &\quad + \mathbb{E}_D \left[ \sum_{i=2}^s \kappa(x; Z_i, D) (\mu_0^1(X_i) - \mu_0^0(X_i)) \mid Z_1 = z_1 \right] \\ &= \beta(W_1, X_1) \varepsilon_1 \cdot \mathbb{E}[\kappa(x; Z_1, D) | X_1 = x_1] + \mathbb{E}_D \left[ \sum_{i=2}^s \kappa(x; Z_i, D) (\mu_0^1(X_i) - \mu_0^0(X_i)) \mid Z_1 = z_1 \right]\end{aligned}\tag{B.11}$$

■

## B.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 B.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

$$\Omega_s(x) = \mathbb{E} [h_s^2(x; Z_1, \dots, Z_s)] . \quad (\text{B.12})$$

Then,

$$\Omega_s(x) = \mathbb{E}_1 \left[ (\mu(X_1) + \varepsilon_1)^2 s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)] \right] \lesssim \mu^2(x) + \bar{\sigma}_\varepsilon^2 + o(1) \quad \text{as } s \rightarrow \infty. \quad (\text{B.13})$$


---

*Proof of Lemma B.5.*

This result follows immediately from Lemma A.2 and the following observation.

$$\begin{aligned} \Omega_s(x) &= \mathbb{E} [h_s^2(x; Z_1, \dots, Z_s)] = \mathbb{E}_D \left[ \left( \sum_{i=1}^s \kappa(x; Z_i, D) Y_i \right)^2 \right] = \mathbb{E}_D \left[ \sum_{i=1}^s \sum_{j=1}^s (\kappa(x; Z_i, D) \kappa(x; Z_j, D) Y_i Y_j) \right] \\ &= \mathbb{E}_D [s \kappa(x; Z_1, D) Y_1^2] = \mathbb{E}_1 [Y_1^2 s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] = \mathbb{E}_1 [(\mu(X_1) + \varepsilon_1)^2 s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] \\ &= \mathbb{E}_1 [(\mu^2(X_1) + 2\mu(X_1)\varepsilon_1 + \varepsilon_1^2) s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] \\ &= \mathbb{E}_1 [(\mu^2(X_1) + 2\mu(X_1)\mathbb{E}[\varepsilon_1 | X_1] + \mathbb{E}[\varepsilon_1^2 | X_1]) s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] \\ &= \mathbb{E}_1 [(\mu^2(X_1) + \sigma_\varepsilon^2(X_1)) s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] \xrightarrow{(\text{Lem A.2})} \mu^2(x) + \sigma_\varepsilon^2(x) \quad \text{as } s \rightarrow \infty \end{aligned} \quad (\text{B.14})$$

Furthermore, we have the following inequality.

$$\mu^2(x) + \sigma_\varepsilon^2(x) \leq \mu^2(x) + \bar{\sigma}_\varepsilon^2 \quad (\text{B.15})$$

Thus, we obtain the desired result. ■

---

**Lemma B.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

$$\Omega_s^c(x) = \mathbb{E} [h_s(x; Z_1, \dots, Z_c, Z_{c+1}, \dots, Z_s) \cdot h_s(x; Z_1, \dots, Z_c, Z'_{c+1}, \dots, Z'_s)] . \quad (\text{B.16})$$

Then,

$$\Omega_s^c(x) \lesssim \mu^2(x) + \bar{\sigma}_\varepsilon^2 + o(1) \quad \text{as } s \rightarrow \infty. \quad (\text{B.17})$$


---

*Proof of Lemma B.6.*

$$\begin{aligned}
\Omega_s^c(x) &= \mathbb{E} [h_s(x; Z_1, \dots, Z_c, Z_{c+1}, \dots, Z_s) \cdot h_s(x; Z_1, \dots, Z_c, Z'_{c+1}, \dots, Z'_s)] \\
&= \mathbb{E}_{D, D'} \left[ \left( \sum_{i=1}^s \kappa(x; Z_i, D) Y_i \right) \left( \sum_{j=1}^c \kappa(x; Z_j, D') Y_j + \sum_{j=c+1}^s \kappa(x; Z'_j, D') Y'_j \right) \right] \\
&= \mathbb{E}_{D, D'} \left[ \sum_{i=1}^c \sum_{j=1}^c \kappa(x; Z_i, D) \kappa(x; Z_j, D') Y_i Y_j \right] + \mathbb{E}_{D, D'} \left[ \sum_{i=1}^c \sum_{j=c+1}^s \kappa(x; Z_i, D) \kappa(x; Z'_j, D') Y_i Y'_j \right] \\
&\quad + \mathbb{E}_{D, D'} \left[ \sum_{i=c+1}^s \sum_{j=1}^c \kappa(x; Z_i, D) \kappa(x; Z_j, D') Y_i Y_j \right] + \mathbb{E}_{D, D'} \left[ \sum_{i=c+1}^s \sum_{j=c+1}^s \kappa(x; Z_i, D) \kappa(x; Z'_j, D') Y_i Y'_j \right] \\
&= \underbrace{\mathbb{E}_{D, D'} [c \kappa(x; Z_1, D) \kappa(x; Z_1, D') Y_1^2]}_{(A)} + 2 \cdot \underbrace{\mathbb{E}_{D, D'} [c(s-c) \kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D') Y_1 Y'_{c+1}]}_{(B)} \\
&\quad + \underbrace{\mathbb{E}_{D, D'} [(s-c)^2 \kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D') Y_{c+1} Y'_{c+1}]}_{(C)}
\end{aligned} \tag{B.18}$$

Starting from this decomposition, we will analyze the terms one by one using Lemma A.5.

$$\begin{aligned}
(A) &= \mathbb{E}_{D, D'} [c \kappa(x; Z_1, D) \kappa(x; Z_1, D') Y_1^2] = \frac{c}{2s-c} \cdot \mathbb{E}_1 [(\mu(X_1) + \varepsilon_1)^2 \cdot (2s-c) \cdot \mathbb{E}_{2:s} [\kappa(x; Z_1, D) \kappa(x; Z_1, D')]] \\
&= \frac{c}{2s-c} \cdot \mathbb{E}_1 [(\mu^2(X_1) + \sigma^2(X_1)) \cdot (2s-c) \cdot \mathbb{E}_{2:s} [\kappa(x; Z_1, D) \kappa(x; Z_1, D')]] \\
&\stackrel{(\text{Lem A.5})}{\lesssim} \frac{c}{2s-c} (\mu^2(x) + \sigma_\varepsilon(x)) + o(1)
\end{aligned} \tag{B.19}$$

Similarly, we can find the following.

$$\begin{aligned}
(B) &= \mathbb{E}_{D, D'} [c(s-c) \kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D') Y_1 Y'_{c+1}] \\
&\stackrel{\text{Lem A.6}}{\leq} \frac{c(s-c)s}{(2s-c)(2s-c-1)(c+1)} \\
&\quad \cdot \mathbb{E}_{1, (c+1)'} \left[ (\mu(X_1) + \varepsilon_1) \cdot (\mu(X'_{c+1}) + \varepsilon'_{c+1}) \cdot \frac{\mathbb{E}_{D, D'} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D') \mid Z_1, Z'_{c+1}]}{\mathbb{E}_{D, D'} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D')]} \right] \\
&\leq \frac{(s-c)s}{(2s-c)(2s-c-1)} \cdot \mathbb{E}_{1, (c+1)'} \left[ \mu(X_1) \cdot \mu(X'_{c+1}) \cdot \frac{\mathbb{E}_{D, D'} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D') \mid Z_1, Z'_{c+1}]}{\mathbb{E}_{D, D'} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D')]} \right] \\
&\stackrel{\text{Lem A.5}}{\lesssim} \frac{(s-c)s}{(2s-c)(2s-c-1)} \cdot \mu^2(x) + o(1)
\end{aligned} \tag{B.20}$$

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

$$\begin{aligned}
(C) &= \mathbb{E}_{D, D'} [(s-c)^2 \kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D') Y_{c+1} Y'_{c+1}] \\
&\stackrel{\text{Lem A.6}}{\leq} \frac{2(s-c)^3}{(2s-c)^2(2s-c-1)} \cdot \mathbb{E}_{c+1} \left[ \mu(X_{c+1}) \cdot \mu(X'_{c+1}) \cdot \frac{\mathbb{E}_{D, D'} [\kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D') \mid Z_{c+1}, Z'_{c+1}]}{\mathbb{E}_{D, D'} [\kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D')]} \right] \\
&\stackrel{\text{Lem A.5}}{\lesssim} \frac{2(s-c)}{(2s-c-1)} \cdot \mu^2(x) + o(1)
\end{aligned} \tag{B.21}$$

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

**Lemma B.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

$$\Upsilon_{s_1, s_2}(x) = \mathbb{E}[h_{s_1}(x; Z_1, \dots, Z_{s_1}) \cdot h_{s_2}(x; Z_1, \dots, Z_{s_1}, \dots, Z_{s_2})]. \quad (\text{B.22})$$

Then,

$$\Upsilon_{s_1, s_2}(x) \lesssim \mu^2(x) + \bar{\sigma}_\varepsilon^2 + o(1) \quad \text{as } s_1, s_2 \rightarrow \infty \quad \text{with } 0 < \mathfrak{c} \leq s_1/s_2 \leq 1 - \mathfrak{c} < 1. \quad (\text{B.23})$$

*Proof of Lemma B.7.*

$$\begin{aligned} \Upsilon_{s_1, s_2}(x) &= \mathbb{E}[h_{s_1}(x; Z_1, \dots, Z_{s_1}) \cdot h_{s_2}(x; Z_1, \dots, Z_{s_1}, \dots, Z_{s_2})] \\ &= \mathbb{E}_D \left[ \left( \sum_{i=1}^{s_1} \kappa(x; Z_i, D_{[s_1]}) Y_i \right) \left( \sum_{j=1}^{s_1} \kappa(x; Z_j, D) Y_j + \sum_{j=s_1+1}^{s_2} \kappa(x; Z_j, D) Y_j \right) \right] \\ &= \mathbb{E}_D \left[ \sum_{i=1}^{s_1} \kappa(x; Z_i, D) Y_i^2 \right] + \mathbb{E}_D \left[ \sum_{i=1}^{s_1} \sum_{j=s_1+1}^{s_2} \kappa(x; Z_i, D_{[s_1]}) \kappa(x; Z_j, D) Y_i Y_j \right] \\ &= \mathbb{E}_D [Y_1^2 s_1 \kappa(x; Z_1, D)] + \mathbb{E}_D [Y_1 Y_{s_2} s_1 (s_2 - s_1) \kappa(x; Z_1, D_{[s_1]}) \kappa(x; Z_{s_2}, D)] \\ &= \mathbb{E}_D [(\mu^2(X_1) + \sigma_\varepsilon^2(X_1)) s_1 \kappa(x; Z_1, D)] + \mathbb{E}_D [\mu(X_1) \mu(X_{s_2}) s_1 (s_2 - s_1) \kappa(x; Z_1, D_{[s_1]}) \kappa(x; Z_{s_2}, D)] \\ &= \frac{s_1}{s_2} \mathbb{E}_D [(\mu^2(X_1) + \sigma_\varepsilon^2(X_1)) s_1 \kappa(x; Z_1, D)] + \frac{s_2 - s_1}{s_2} \mathbb{E}_D [\mu(X_1) \mu(X_{s_2}) s_1 s_2 \kappa(x; Z_1, D_{[s_1]}) \kappa(x; Z_{s_2}, D)] \\ &\leq \frac{s_1}{s_2} \mathbb{E}_D [(\mu^2(X_1) + \sigma_\varepsilon^2(X_1)) s_2 \kappa(x; Z_1, D)] \\ &\quad + \frac{s_2 - s_1}{s_2} \mathbb{E}_D [|\mu(X_1)| s_1 \kappa(x; Z_1, D_{[s_1]})] \mathbb{E}_D [|\mu(X_{s_2})| s_2 \kappa(x; Z_{s_2}, D)] \\ &\lesssim \mu^2(x) + \sigma_\varepsilon^2(x) + o(1) \leq \mu^2(x) + \bar{\sigma}_\varepsilon^2 + o(1). \end{aligned} \quad (\text{B.24})$$

■



**Lemma B.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

$$\Upsilon_{s_1, s_2}^c(x) = \mathbb{E} \left[ h_{s_1}(x; Z_1, \dots, Z_c, Z'_{c+1}, \dots, Z'_{s_1}) \cdot h_{s_2}(x; Z_1, \dots, Z_{s_2}) \right]. \quad (\text{B.25})$$

Then,

$$\Upsilon_{s_1, s_2}^c(x) \lesssim \frac{cs_2 - c^2 + s_1 s_2}{s_1 s_2} \mu^2(x) + (c/s_1) \bar{\sigma}_\varepsilon^2 + o(1) \quad (\text{B.26})$$

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

and thus

$$\Upsilon_{s_1, s_2}^c(x) \lesssim \mu^2(x) + o(1) \quad \text{as } s_1, s_2 \rightarrow \infty \quad \text{with } 0 < \mathfrak{c} \leq s_1/s_2 \leq 1 - \mathfrak{c} < 1. \quad (\text{B.27})$$

*Proof of Lemma B.8.*

$$\begin{aligned} \Upsilon_{s_1, s_2}^c(x) &= \mathbb{E} \left[ h_{s_1}(x; Z_1, \dots, Z_c, Z'_{c+1}, \dots, Z'_{s_1}) \cdot h_{s_2}(x; Z_1, \dots, Z_{s_2}) \right] \\ &= \mathbb{E}_{D, D'} \left[ \left( \sum_{i=1}^c \kappa(x; Z_i, D') Y_i + \sum_{i=c+1}^{s_1} \kappa(x; Z'_i, D') Y'_i \right) \left( \sum_{j=1}^c \kappa(x; Z_j, D) Y_j + \sum_{j=c+1}^{s_2} \kappa(x; Z_j, D) Y_j \right) \right] \\ &= \underbrace{\mathbb{E}_{D, D'} \left[ \sum_{i=1}^c \sum_{j=1}^c \kappa(x; Z_i, D') \kappa(x; Z_j, D) Y_i Y_j \right]}_{(A)} + \underbrace{\mathbb{E}_{D, D'} \left[ \left( \sum_{i=1}^c \kappa(x; Z_i, D') Y_i \right) \left( \sum_{j=c+1}^{s_2} \kappa(x; Z_j, D) Y_j \right) \right]}_{(B)} \\ &\quad + \underbrace{\mathbb{E}_{D, D'} \left[ \sum_{i=c+1}^{s_1} \sum_{j=1}^c \kappa(x; Z'_i, D') \kappa(x; Z_j, D) Y'_i Y_j \right]}_{(C)} + \underbrace{\mathbb{E}_{D, D'} \left[ \left( \sum_{i=c+1}^{s_1} \kappa(x; Z'_i, D') Y'_i \right) \left( \sum_{j=c+1}^{s_2} \kappa(x; Z_j, D) Y_j \right) \right]}_{(D)} \end{aligned} \quad (\text{B.28})$$

Again, we have four terms to analyze individually.

$$\begin{aligned} (A) &= \mathbb{E}_{D, D'} \left[ \sum_{i=1}^c \sum_{j=1}^c \kappa(x; Z_i, D') \kappa(x; Z_j, D) Y_i Y_j \right] \\ &= \mathbb{E}_{D, D'} \left[ \sum_{i=1}^c Y_i^2 \kappa(x; Z_i, D') \kappa(x; Z_i, D) \right] \\ &= \mathbb{E}_{D, D'} \left[ Y_1^2 c \kappa(x; Z_1, D') \kappa(x; Z_1, D) \right] = \mathbb{E}_{D, D'} \left[ (\mu^2(X_1) + \sigma_\varepsilon^2(X_1)) c \kappa(x; Z_1, D_{[c]}) \kappa(x; Z_1, D'_{c+1:s_1}) \right] \\ &= \mathbb{E}_D \left[ (\mu^2(X_1) + \sigma_\varepsilon^2(X_1)) c \kappa(x; Z_1, D) \right] = \frac{c}{s_1} \mathbb{E}_D \left[ (\mu^2(X_1) + \sigma_\varepsilon^2(X_1)) s_1 \kappa(x; Z_1, D) \right] \\ &\lesssim (c/s_1) (\mu^2(x) + \sigma_\varepsilon^2(x)) + o(1) \leq (c/s_1) (\mu^2(x) + \bar{\sigma}_\varepsilon^2) + o(1) \end{aligned} \quad (\text{B.29})$$

Considering the second term, we find the following.

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

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

$$\begin{aligned}
(C) &= \mathbb{E}_{D,D'} \left[ \sum_{i=c+1}^{s_1} \sum_{j=1}^c \kappa(x; Z'_i, D') \kappa(x; Z_j, D) Y'_i Y_j \right] = \mathbb{E}_{D,D'} [Y'_{s_1} Y_1 (s_1 - c) c \kappa(x; Z'_{s_1}, D') \kappa(x; Z_1, D)] \\
&= \frac{(s_1 - c)c}{s_1 s_2} \mathbb{E}_{D,D'} [\mu(X'_{s_1}) \mu(X_1) s_1 s_2 \kappa(x; Z'_{s_1}, D') \kappa(x; Z_1, D)] \\
&\leq \frac{(s_1 - c)c}{s_1 s_2} \mathbb{E}_D [|\mu(X'_{s_1})| s_1 \kappa(x; Z'_{s_1}, D')] \mathbb{E}_D [|\mu(X_1)| s_2 \kappa(x; Z_1, D)] \\
&\lesssim \frac{(s_1 - c)c}{s_1 s_2} \mu^2(x) + o(1)
\end{aligned} \tag{B.31}$$

Lastly, concerning the fourth term, observe the following.

$$\begin{aligned}
(D) &= \mathbb{E}_{D,D'} \left[ \left( \sum_{i=c+1}^{s_1} \kappa(x; Z'_i, D') Y'_i \right) \left( \sum_{j=c+1}^{s_2} \kappa(x; Z_j, D) Y_j \right) \right] = \mathbb{E}_{D,D'} \left[ \sum_{i=c+1}^{s_1} \sum_{j=c+1}^{s_2} \kappa(x; Z'_i, D') \kappa(x; Z_j, D) Y'_i Y_j \right] \\
&= \mathbb{E}_{D,D'} [\mu(X'_{s_1}) \mu(X_{s_2}) (s_1 - c)(s_2 - c) \kappa(x; Z'_{s_1}, D') \kappa(x; Z_{s_2}, D)] \\
&= \frac{(s_1 - c)(s_2 - c)}{s_1 s_2} \mathbb{E}_{D,D'} [\mu(X'_{s_1}) \mu(X_{s_2}) s_1 s_2 \kappa(x; Z'_{s_1}, D') \kappa(x; Z_{s_2}, D)] \\
&\leq \frac{(s_1 - c)(s_2 - c)}{s_1 s_2} \mathbb{E}_{D'} [|\mu(X'_{s_1})| s_1 \kappa(x; Z'_{s_1}, D')] \mathbb{E}_D [|\mu(X_{s_2})| s_2 \kappa(x; Z_{s_2}, D)] \\
&\lesssim \frac{(s_1 - c)(s_2 - c)}{s_1 s_2} \mu^2(x) + o(1)
\end{aligned} \tag{B.32}$$

■

**Lemma B.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*

$$\zeta_{s_1, s_2}^{s_2}(x) \lesssim \mu^2(x) + \bar{\sigma}_\varepsilon + o(1) \quad \text{as } s_1, s_2 \rightarrow \infty \quad \text{with } 0 < \mathfrak{c} \leq s_1/s_2 \leq 1 - \mathfrak{c} < 1. \quad (\text{B.33})$$

*Proof of Lemma B.9.* Consider first the following decomposition.

$$\begin{aligned} \zeta_{s_1, s_2}^{s_2}(x) &= \text{Var}(h_{s_1, s_2}(x; Z_1, \dots, Z_{s_2})) = \text{Var}_D(h_{s_1, s_2}(x; D)) \\ &\leq \mathbb{E}_D[h_{s_1, s_2}^2(x; D)] = \mathbb{E}_D[(w_1^* \tilde{\mu}_{s_1}(x; D) + w_2^* h_{s_2}(x; D))^2] \\ &= (w_1^*)^2 \mathbb{E}_D[\tilde{\mu}_{s_1}^2(x; D)] + 2w_1^* w_2^* \mathbb{E}_D[\tilde{\mu}_{s_1}(x; D) h_{s_2}(x; D)] + (w_2^*)^2 \Omega_{s_2} \end{aligned} \quad (\text{B.34})$$

Then, observe the following.

$$\begin{aligned} \mathbb{E}_D[\tilde{\mu}_{s_1}^2(x; D)] &= \mathbb{E}_D \left[ \left( \binom{s_2}{s_1}^{-1} \sum_{\ell \in L_{s_2, s_1}} h_{s_1}(x; D_\ell) \right)^2 \right] = \binom{s_2}{s_1}^{-2} \mathbb{E}_D \left[ \sum_{\ell, \ell' \in L_{s_2, s_1}} h_{s_1}(x; D_\ell) h_{s_1}(x; D_{\ell'}) \right] \\ &= \binom{s_2}{s_1}^{-2} \sum_{c=0}^{s_1} \binom{s_2}{s_1} \binom{s_1}{c} \binom{s_2 - s_1}{s_1 - c} \Omega_{s_1}^c = \binom{s_2}{s_1}^{-1} \sum_{c=0}^{s_1} \binom{s_1}{c} \binom{s_2 - s_1}{s_1 - c} \Omega_{s_1}^c \\ &\lesssim \Omega_{s_1} \lesssim \mu(x)^2 + \sigma_\varepsilon^2 + o(1) \quad \text{as } s \rightarrow \infty \end{aligned} \quad (\text{B.35})$$

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

$$\Omega_{s_2} \lesssim \mu(x)^2 + \sigma_\varepsilon^2 + o(1) \quad \text{as } s \rightarrow \infty \quad (\text{B.36})$$

Lastly, consider the following.

$$\begin{aligned} \mathbb{E}_D[\tilde{\mu}_{s_1}(x; D) h_{s_2}(x; D)] &= \mathbb{E}_D \left[ \binom{s_2}{s_1}^{-1} \sum_{\ell \in L_{s_2, s_1}} h_{s_1}(x; D_\ell) h_{s_2}(x; D) \right] \\ &= \mathbb{E}_D[h_{s_1}(x; D_{[s_1]}) h_{s_2}(x; D)] = \Upsilon_{s_1, s_2}(x) \end{aligned} \quad (\text{B.37})$$

Thus, we find the following.

$$\begin{aligned} \zeta_{s_2, s_2}(x) &\lesssim (w_1^*)^2 \Omega_{s_1} + 2w_1^* w_2^* \Upsilon_{s_1, s_2}(x) + (w_1^*)^2 \Omega_{s_2} \\ &\lesssim (w_1^* + w_2^*)^2 (\mu^2(x) + \sigma_\varepsilon) + o(1) = \mu^2(x) + \sigma_\varepsilon + o(1). \end{aligned} \quad (\text{B.38})$$

■

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

$$0 < \mathfrak{c} \leq s_1/s_2 \leq 1 - \mathfrak{c} < 1 \quad \text{and} \quad s_2 = o(n), \quad (\text{B.39})$$

*it holds that*

$$\zeta_{s_1, s_2}^1(x) \sim s_2^{-1}. \quad (\text{B.40})$$

## B.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 B.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

$$\Omega_s(x) = \mathbb{E} [\chi_{s,0}^2(x; Z_1, \dots, Z_s)] . \quad (\text{B.41})$$

Then,

$$\Omega_s(x) \lesssim (\mu_0^1(x) - \mu_0^0(x))^2 + \frac{\bar{\sigma}_\varepsilon^2}{\mathbf{p}(1-\mathbf{p})} + o(1) \quad (\text{B.42})$$

*Proof of Lemma B.11.*

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

$$\begin{aligned} \Omega_s(x) &= \mathbb{E} [\chi_{s,0}^2(x; Z_1, \dots, Z_s)] = \mathbb{E}_D \left[ \left( \sum_{i=1}^s \kappa(x; Z_i, D) m(Z_i; \eta_0) \right)^2 \right] \\ &= \mathbb{E}_D \left[ \sum_{i=1}^s \sum_{j=1}^s \kappa(x; Z_i, D) \kappa(x; Z_j, D) m(Z_i; \eta_0) m(Z_j; \eta_0) \right] = \mathbb{E}_D [s \kappa(x; Z_1, D) m^2(Z_1; \eta_0)] \\ &= \mathbb{E}_1 [m^2(Z_1; \eta_0) s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)]] \\ &= \mathbb{E}_1 \left[ (\mu_0^1(X_1) - \mu_0^0(X_1) + \beta(W_1, X_1) \varepsilon_1)^2 s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)] \right] \\ &= \mathbb{E}_1 \left[ (\mu_0^1(X_1) - \mu_0^0(X_1))^2 s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)] \right] + \mathbb{E}_1 \left[ (\beta(W_1, X_1) \varepsilon_1)^2 s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)] \right] \\ &= \mathbb{E}_1 \left[ (\mu_0^1(X_1) - \mu_0^0(X_1))^2 s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)] \right] + \mathbb{E}_1 \left[ \left( \frac{W_1}{\pi_0(X_1)} - \frac{1-W_1}{1-\pi_0(X_1)} \right)^2 \varepsilon_1^2 s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)] \right] \\ &= \underbrace{\mathbb{E}_1 \left[ (\mu_0^1(X_1) - \mu_0^0(X_1))^2 s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)] \right]}_{\xrightarrow{\text{Lem A.2}} (\mu_0^1(x) - \mu_0^0(x))^2 \text{ as } s \rightarrow \infty} + \underbrace{\mathbb{E}_1 \left[ \mathbb{E} \left[ \left( \frac{W_1}{\pi_0(X_1)} - \frac{1-W_1}{1-\pi_0(X_1)} \right)^2 \varepsilon_1^2 \middle| X_1 \right] s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)] \right]}_{(B)} \end{aligned} \quad (\text{B.43})$$

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

$$\begin{aligned} (B) &= \mathbb{E}_1 \left[ \mathbb{E} \left[ \left( \frac{W_1}{\pi_0(X_1)} - \frac{1-W_1}{1-\pi_0(X_1)} \right)^2 \varepsilon_1^2 \middle| X_1 \right] s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)] \right] \\ &= \mathbb{E}_1 \left[ \frac{\sigma_\varepsilon^2(X_1) \cdot s \mathbb{E}_{2:s} [\kappa(x; Z_1, D)]}{\pi_0^2(X_1) (1-\pi_0(X_1))^2} \cdot \mathbb{E} \left[ (W_1 (1-\pi_0(X_1)) - (1-W_1) \pi_0(X_1))^2 \middle| X_1 \right] \right] \end{aligned} \quad (\text{B.44})$$

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.

$$\begin{aligned}
(B) &= \mathbb{E}_1 \left[ \frac{\sigma_\varepsilon^2(X_1) \cdot s\mathbb{E}_{2:s}[\kappa(x; Z_1, D)]}{\pi_0^2(X_1)(1 - \pi_0(X_1))^2} \cdot \mathbb{E} \left[ W_1(1 - \pi_0(X_1))^2 + (1 - W_1)\pi_0^2(X_1) \mid X_1 \right] \right] \\
&= \mathbb{E}_1 \left[ \frac{\sigma_\varepsilon^2(X_1) \cdot s\mathbb{E}_{2:s}[\kappa(x; Z_1, D)]}{\pi_0^2(X_1)(1 - \pi_0(X_1))^2} \cdot \pi_0(X_1)(1 - \pi_0(X_1)) \cdot (1 - \pi_0(X_1) + \pi_0(X_1)) \right] \\
&= \mathbb{E}_1 \left[ \frac{\sigma_\varepsilon^2(X_1)}{\pi_0(X_1)(1 - \pi_0(X_1))} \cdot s\mathbb{E}_{2:s}[\kappa(x; Z_1, D)] \right] \xrightarrow{(\text{Lem A.2})} \frac{\sigma_\varepsilon^2(x)}{\pi_0(x)(1 - \pi_0(x))} \quad \text{as } s \rightarrow \infty
\end{aligned} \tag{B.45}$$

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

$$\mathbb{E}_1 \left[ m^2(Z_i; \eta_0) s\mathbb{E}_{2:s}[\kappa(x; Z_1, D)] \right] \xrightarrow{(\text{Lem A.2})} (\mu_0^1(x) - \mu_0^0(x))^2 + \frac{\sigma_\varepsilon^2(x)}{\pi_0(x)(1 - \pi_0(x))} \quad \text{as } s \rightarrow \infty \tag{B.46}$$

This gives us the desired result.

$$\Omega_s(x) \lesssim (\mu_0^1(x) - \mu_0^0(x))^2 + \frac{\bar{\sigma}_\varepsilon^2}{\mathfrak{p}(1 - \mathfrak{p})} + o(1) \tag{B.47}$$

■

---

**Lemma B.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

$$\Omega_s^c(x) = \mathbb{E} \left[ \chi_{s,0}(x; Z_1, \dots, Z_c, Z_{c+1}, \dots, Z_s) \cdot \chi_{s,0}(x; Z_1, \dots, Z_c, Z'_{c+1}, \dots, Z'_s) \right]. \tag{B.48}$$

Then,

$$\Omega_s^c(x) \lesssim C \left[ (\mu_0^1(x) - \mu_0^0(x))^2 + \frac{\bar{\sigma}_\varepsilon^2}{\mathfrak{p}(1 - \mathfrak{p})} \right] + o(1). \tag{B.49}$$


---

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

$$\begin{aligned}
\Omega_s^c(x) &= \mathbb{E} [\chi_{s,0}(x; Z_1, \dots, Z_c, Z_{c+1}, \dots, Z_s) \cdot \chi_{s,0}(x; Z_1, \dots, Z_c, Z'_{c+1}, \dots, Z'_s)] \\
&= \mathbb{E}_{D,D'} \left[ \left( \sum_{i=1}^s \kappa(x; Z_i, D) m(Z_i; \eta_0) \right) \left( \sum_{j=1}^c \kappa(x; Z_j, D') m(Z_j; \eta_0) + \sum_{j=c+1}^s \kappa(x; Z'_j, D') m(Z'_j; \eta_0) \right) \right] \\
&= \underbrace{\mathbb{E}_{D,D'} \left[ \left( \sum_{i=1}^c \kappa(x; Z_i, D) m(Z_i; \eta_0) \right) \left( \sum_{j=1}^c \kappa(x; Z_j, D') m(Z_j; \eta_0) \right) \right]}_{(A)} \\
&\quad + \underbrace{2 \mathbb{E}_{D,D'} \left[ \left( \sum_{i=1}^c \kappa(x; Z_i, D) m(Z_i; \eta_0) \right) \left( \sum_{j=c+1}^s \kappa(x; Z'_j, D') m(Z'_j; \eta_0) \right) \right]}_{(B)} \\
&\quad + \underbrace{\mathbb{E}_{D,D'} \left[ \left( \sum_{i=c+1}^s \kappa(x; Z_i, D) m(Z_i; \eta_0) \right) \left( \sum_{j=c+1}^s \kappa(x; Z'_j, D') m(Z'_j; \eta_0) \right) \right]}_{(C)}
\end{aligned} \tag{B.50}$$

Considering these terms one by one, we can make the following observations. Here, we rely on the same argument structure as in the proof of Lemma B.6 and observations from the proof of Lemma B.11.

$$\begin{aligned}
(A) &= \mathbb{E}_{D,D'} \left[ \left( \sum_{i=1}^c \kappa(x; Z_i, D) m(Z_i; \eta_0) \right) \left( \sum_{j=1}^c \kappa(x; Z_j, D') m(Z_j; \eta_0) \right) \right] \\
&= \mathbb{E}_{D,D'} \left[ \sum_{i=1}^c \sum_{j=1}^c \kappa(x; Z_i, D) \kappa(x; Z_j, D') m(Z_i; \eta_0) m(Z_j; \eta_0) \right] \\
&= \frac{c}{2s-c} \cdot \mathbb{E}_1 [m^2(Z_1; \eta_0) \cdot (2s-c) \cdot \mathbb{E}_{2:s} [\kappa(x; Z_1, D) \kappa(x; Z_1, D')]] \\
&\stackrel{\text{Lem A.5}}{\lesssim} \frac{c}{2s-c} \cdot \left( (\mu_0^1(x) - \mu_0^0(x))^2 + \frac{\sigma_\varepsilon^2(x)}{\pi_0(x)(1-\pi_0(x))} \right) + o(1) \\
&\leq \frac{c}{2s-c} \cdot \left( (\mu_0^1(x) - \mu_0^0(x))^2 + \frac{\bar{\sigma}_\varepsilon^2}{\mathfrak{p}(1-\mathfrak{p})} \right) + o(1)
\end{aligned} \tag{B.51}$$

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

$$\begin{aligned}
(B) &= \mathbb{E}_{D,D'} \left[ \left( \sum_{i=1}^c \kappa(x; Z_i, D) m(Z_i; \eta_0) \right) \left( \sum_{j=c+1}^s \kappa(x; Z'_j, D') m(Z'_j; \eta_0) \right) \right] \\
&= \mathbb{E}_{D,D'} \left[ \sum_{i=1}^c \sum_{j=c+1}^s \kappa(x; Z_i, D) \kappa(x; Z'_j, D') m(Z_i; \eta_0) m(Z'_j; \eta_0) \right] \\
&= \mathbb{E}_{D,D'} [c(s-c) \kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D') m(Z_1; \eta_0) m(Z'_{c+1}; \eta_0)] \\
&\stackrel{\text{Lem A.6}}{\leq} \frac{c(s-c)s}{(2s-c)(2s-c-1)(c+1)} \cdot \mathbb{E}_{1,(c+1)'} \left[ m(Z_1; \eta_0) m(Z'_{c+1}; \eta_0) \cdot \frac{\mathbb{E}_{D,D'} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D') \mid Z_1, Z'_{c+1}]}{\mathbb{E}_{D,D'} [\kappa(x; Z_1, D) \kappa(x; Z'_{c+1}, D')]} \right] \\
&\stackrel{\text{Lem A.5}}{\lesssim} \frac{c(s-c)s}{(2s-c)(2s-c-1)(c+1)} \cdot (\mu_0^1(x) - \mu_0^0(x))^2 + o(1)
\end{aligned} \tag{B.52}$$

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

$$\begin{aligned}
(C) &= \mathbb{E}_{D,D'} \left[ \left( \sum_{i=c+1}^s \kappa(x; Z_i, D) m(Z_i; \eta_0) \right) \left( \sum_{j=c+1}^s \kappa(x; Z'_j, D') m(Z'_j; \eta_0) \right) \right] \\
&= \mathbb{E}_{D,D'} [(s-c)^2 \kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D') m(Z_{c+1}; \eta_0) m(Z'_{c+1}; \eta_0)] \\
&\stackrel{\text{Lem A.6}}{\leq} \frac{2(s-c)^3}{(2s-c)^2(2s-c-1)} \cdot \mathbb{E}_{c+1} \left[ m(Z_{c+1}; \eta_0) m(Z'_{c+1}; \eta_0) \cdot \frac{\mathbb{E}_{D,D'} [\kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D') \mid Z_{c+1}, Z'_{c+1}]}{\mathbb{E}_{D,D'} [\kappa(x; Z_{c+1}, D) \kappa(x; Z'_{c+1}, D')]} \right] \\
&\stackrel{\text{Lem A.5}}{\lesssim} \frac{2(s-c)}{(2s-c-1)} \cdot (\mu_0^1(x) - \mu_0^0(x))^2 + o(1)
\end{aligned} \tag{B.53}$$

Thus, we find the desired result.

$$\begin{aligned}
\Omega_s^c(x) &= (A) + 2 \cdot (B) + (C) \\
&\lesssim \left( \frac{c}{2s-c} + \frac{4(s-c)}{2s-c-1} + \frac{2(s-c)}{2s-c-1} \right) \cdot (\mu_0^1(x) - \mu_0^0(x))^2 + \frac{c}{2s-c} \cdot \frac{\bar{\sigma}_\varepsilon^2}{\mathfrak{p}(1-\mathfrak{p})} + o(1) \\
&\leq \frac{6s-5c}{2s-c-1} \cdot (\mu_0^1(x) - \mu_0^0(x))^2 + \frac{c}{2s-c} \cdot \frac{\bar{\sigma}_\varepsilon^2}{\mathfrak{p}(1-\mathfrak{p})} + o(1)
\end{aligned} \tag{B.54}$$

■



## B.6 Asymptotic Normality Results

---

## B.7 NPR-Estimators - Asymptotic Normality

---

## B.8 CATE-Estimators - Asymptotic Normality

---

Recall the decomposition of the DNN-DML2 estimator introduced in equation 4.23.

$$\begin{aligned}
\hat{\theta}(x; \mathbf{D}) = & \underbrace{\mathbb{E}_D [\hat{\theta}(x; \mathbf{D})]}_{\text{Centering-Term}} + \underbrace{\frac{s}{n} \sum_{i=1}^n \chi_{s,0}^{(1)}(x; Z_i)}_{\text{Oracle-Hájek-Projection}} + \underbrace{\frac{s}{k} \sum_{l=1}^k \frac{1}{m} \sum_{i \in \mathcal{I}_k} \left( \underbrace{\chi_s^{(1)}(x; Z_i, \hat{\eta}_k) - \chi_{s,0}^{(1)}(x; Z_i)}_{R_{1,k}(x; Z_i)} \right)}_{\text{Oracle-Hájek-Projection Error}} \\
& + \underbrace{\sum_{j=2}^s \binom{s}{j} \binom{n}{j}^{-1} \sum_{\ell \in L_{n,j}} \chi_{s,0}^{(j)}(x; \mathbf{D}_\ell)}_{\text{Oracle-Hájek-Residual}} + \underbrace{\sum_{j=2}^s \binom{s}{j} \binom{n}{j}^{-1} \sum_{\ell \in L_{n,j}} R_j(x; \mathbf{D}_\ell)}_{\text{Higher-Order Error Terms}}
\end{aligned} \tag{B.55}$$

where we have the following definition from Equation 4.22.

$$\chi_s^{(c)}(x; \mathbf{D}_\ell, \hat{\eta}) = \chi_{s,0}^{(c)}(x; \mathbf{D}_\ell) + \underbrace{\chi_s^{(c)}(x; \mathbf{D}_\ell, \hat{\eta}) - \chi_{s,0}^{(c)}(x; \mathbf{D}_\ell)}_{R_c(x; \mathbf{D}_\ell)} \tag{B.56}$$

Recall that we are ultimately even more interested in the approximation errors of the following form.

$$\tilde{R}_c(x; \mathbf{D}_\ell) = R_c(x; \mathbf{D}_\ell) - (-1)^c \cdot (\mathbb{E}_D [\chi_s(x; \mathbf{D}_{[s]}, \hat{\eta}) - \chi_{s,0}(x; \mathbf{D}_{[s]})]) \tag{B.57}$$

---

**Lemma B.13** (Behavior of Oracle-Hájek-Projection Error).

*Proof of Lemma B.13.*

Consider first the average Oracle-error within a given fold  $k$  and observe the following.

$$\bar{R}_{1,k}(x) = \frac{1}{m} \sum_{l \in I_k} \tilde{R}_{1,k}(x, Z_l) = \frac{1}{m} \sum_{l \in I_k} (\vartheta_s^1(x; Z_l, \hat{\eta}_k) - \vartheta_{s,0}^1(x; Z_l)) \tag{B.58}$$

Define the following empirical process notation, where  $f$  is any  $Q$ -integrable function on  $\mathcal{Z}$ .

$$\mathbb{G}_{m,k}[f(Z)] = \sqrt{\frac{1}{m}} \sum_{i \in I_k} (f(Z_i) - \mathbb{E}_Z[f(Z)]) \tag{B.59}$$

Here, in analogy to step 3 in the proof of Theorem 3.1 in Victor Chernozhukov, Chetverikov, et al. (2018), we can now observe the following which follows from the triangle inequality.

$$|\bar{R}_{1,k}(x)| \leq \sqrt{\frac{1}{m}} \cdot (\mathcal{I}_{3,k}^{(1)} + \mathcal{I}_{4,k}^{(1)}) \tag{B.60}$$

where

$$\mathcal{I}_{3,k}^{(1)} := |\mathbb{G}_{m,k}[\vartheta_s^1(x; Z, \hat{\eta}_k)] - \mathbb{G}_{m,k}[\vartheta_{s,0}^1(x; Z)]| \tag{B.61}$$

$$\mathcal{I}_{4,k}^{(1)} := \sqrt{m} \cdot \left| \mathbb{E}_Z \left[ \vartheta_s^1(x; Z, \hat{\eta}_k) \mid \mathbf{D}_{I_k^C} \right] - \mathbb{E}_Z \left[ \vartheta_{s,0}^1(x; Z) \right] \right| \quad (\text{B.62})$$

Notice that conditional on  $\mathbf{D}_{I_k^C}$  the first-stage estimate  $\hat{\eta}_k$  is non-stochastic. This allows us to make the following observation given the event  $\mathcal{E}_n$  obtains.

$$\textcolor{red}{LOREMIPSUM} \quad (\text{B.63})$$

**LOREM IPSUM** ■

---

**Lemma B.14** (Behavior of Higher-Order Error Terms).

*Proof of Lemma B.14.*

■

---

## B.9 Variance Estimator Consistency Theorems

---

**Lemma B.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  $\zeta_s^s$  and  $\zeta_s^1$  be defined in analogy to Section 3. Assume that the following condition holds.

$$\frac{s}{n} \left( \frac{\zeta_s^s}{s\zeta_s^1} - 1 \right) \rightarrow 0 \quad (\text{B.64})$$

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

$$\frac{n}{s^2} \frac{\text{Var}(U_s(\mathbf{D}_{[n]}))}{\zeta_s^1} \rightarrow 1. \quad (\text{B.65})$$


---

*Proof.*

$$\begin{aligned} 1 &\leq \frac{n}{s^2} \frac{\text{Var}(U_s(\mathbf{D}_{[n]}))}{\zeta_s^1} = \left( \frac{s^2}{n} \zeta_s^1 \right)^{-1} \sum_{j=1}^s \binom{s}{j}^2 \binom{n}{j}^{-1} V_s^j \\ &\leq 1 + \left( \frac{s^2}{n} \zeta_s^1 \right)^{-1} \frac{s^2}{n^2} \sum_{j=2}^s \binom{s}{j} V_s^j \\ &\leq 1 + \frac{s}{n} \left( \frac{\zeta_s^s}{s\zeta_s^1} - 1 \right) \rightarrow 1. \end{aligned} \quad (\text{B.66})$$

■

---

**Lemma B.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 B.15.

---

*Proof.* Recall the results from Lemmas B.9 and B.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.

$$\frac{s_2}{n} \left( \frac{\zeta_{s_1, s_2}^{s_2}(x)}{s_2 \zeta_{s_1, s_2}^1(x)} - 1 \right) \sim \frac{s_2}{n} (\mu^2(x) + \sigma_\varepsilon + o(1) - 1) \sim \frac{s_2}{n} \rightarrow 0 \quad (\text{B.67})$$

■

---

*Proof of Theorem 5.2.*

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

*Proof of Theorem 5.3.*

Recall the definition of the Jackknife Variance estimator.

$$\hat{\omega}_{JK}^2(x; \mathbf{D}_n) = \frac{n-1}{n} \sum_{i=1}^n (\hat{\mu}_{s_1, s_2}(x; \mathbf{D}_{n, -i}) - \hat{\mu}_{s_1, s_2}(x; \mathbf{D}_n))^2 \quad (\text{B.68})$$

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

$$\begin{aligned} \hat{\omega}_{JK}^2(x; \mathbf{D}_n) &= \frac{n-1}{n} \sum_{i=1}^n \left( \sum_{j=1}^{s_2} \binom{s_2}{j} H_{s_1, s_2}^j(\mathbf{D}_{n, -i}) - \sum_{j=1}^{s_2} \binom{s_2}{j} H_{s_1, s_2}^j(\mathbf{D}_n) \right)^2 \\ &= \frac{n-1}{n} \sum_{j=1}^n \left( \sum_{j=1}^{s_2} \binom{s_2}{j} (H_{s_1, s_2}^j - H_{s_1, s_2}^j(\mathbf{D}_{n, -i})) \right)^2 \\ &= \frac{n-1}{n} \sum_{j=1}^n \left( \sum_{j=1}^{s_2} \binom{s_2}{j} \left( \binom{n}{j}^{-1} \sum_{\iota \in L_{n, j}} h_{s_1, s_2}^{(j)}(\mathbf{D}_\iota) - \binom{n-1}{j}^{-1} \sum_{\ell \in L_j([n] \setminus \{i\})} h_{s_1, s_2}^{(j)}(\mathbf{D}_\ell) \right) \right)^2 \\ &= \frac{n-1}{n} \sum_{j=1}^n \left[ \frac{s_2}{n} h_{s_1, s_2}^{(1)}(Z_i) + \sum_{j \neq i} \left( \frac{s_2}{n} - \frac{s_2}{n-1} \right) h_{s_1, s_2}^{(1)}(Z_j) \right. \\ &\quad \left. + \sum_{j=2}^{s_2} \binom{s_2}{j} \left( \binom{n}{j}^{-1} \sum_{\iota \in L_{n, j}} h_{s_1, s_2}^{(j)}(\mathbf{D}_\iota) - \binom{n-1}{j}^{-1} \sum_{\ell \in L_j([n] \setminus \{i\})} h_{s_1, s_2}^{(j)}(\mathbf{D}_\ell) \right) \right]^2 \\ &= \frac{n-1}{n} \frac{s_2^2}{n^2} \sum_{j=1}^n \left[ h_{s_1, s_2}^{(1)}(Z_i) - \frac{1}{n-1} \sum_{j \neq i} h_{s_1, s_2}^{(1)}(Z_j) \right. \\ &\quad \left. + \frac{n}{s} \sum_{j=2}^{s_2} \binom{s_2}{j} \left( \binom{n}{j}^{-1} \sum_{\iota \in L_{j-1}([n] \setminus \{i\})} h_{s_1, s_2}^{(j)}(\mathbf{D}_{\iota \cup \{i\}}) + \left[ \binom{n}{j}^{-1} - \binom{n-1}{j}^{-1} \right] \sum_{\ell \in L_j([n] \setminus \{i\})} h_{s_1, s_2}^{(j)}(\mathbf{D}_\ell) \right) \right] \\ &=: \frac{n-1}{n} \frac{s_2^2}{n^2} \sum_{j=1}^n [h_{s_1, s_2}^{(1)}(Z_i) + T_i]^2 \end{aligned} \quad (\text{B.69})$$

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.

$$\mathbb{E} \left[ \left( h_{s_1, s_2}^{(1)}(Z_i) \right)^2 \right] = V_{s_1, s_2}^1 = \zeta_{s_1, s_2}^1 \quad (\text{B.70})$$

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{B.71}$$

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

$$\zeta_{s_1, s_2}^{s_2}(x) \lesssim \mu^2(x) + \sigma_\varepsilon + o(1) \quad \text{and} \quad \zeta_{s_1, s_2}^1(x) \sim s_2^{-1} \tag{B.72}$$

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.

$$\begin{aligned}
\frac{\mathbb{E}[T_i^2]}{V_{s_1, s_2}^1} &\leq \frac{\left(\frac{1}{n-1} + \frac{2es_2n}{(n-es_2)^2}\right) \zeta_{s_1, s_2}^1 + 2e\left(\frac{1}{n} + \frac{n}{(n-es_2)^2}\right) (\zeta_{s_1, s_2}^{s_2} - s_2 \zeta_{s_1, s_2}^1)}{\zeta_{s_1, s_2}^1} \\
&= \frac{1}{n-1} + \frac{2es_2n}{(n-es_2)^2} + 2e\left(\frac{1}{n} + \frac{n}{(n-es_2)^2}\right) \left(\frac{\zeta_{s_1, s_2}^{s_2} - s_2 \zeta_{s_1, s_2}^1}{\zeta_{s_1, s_2}^1}\right) \rightarrow 0
\end{aligned} \tag{B.73}$$

Therefore, we can conclude that  $h_s^{(1)}(Z_i)$  dominates  $T_i^2$  in the expression of interest. Using Lemma A.7, we can thus conclude the following.

$$\begin{aligned}
\frac{\frac{n}{s_2^2} \hat{\omega}_{JK}^2(x; \mathbf{D}_n)}{V_{s_1, s_2}^1(x)} &\rightarrow_p \frac{n-1}{n} \frac{1}{n} \sum_{i=1}^n \frac{\left(h_{s_1, s_2}^{(1)}(x; Z_i)\right)^2}{V_{s_1, s_2}^1(x)} \\
&\rightarrow_p \frac{n-1}{n} \frac{\mathbb{E}\left[\left(h_{s_1, s_2}^{(1)}(x; Z_i)\right)^2\right]}{V_{s_1, s_2}^1(x)} \rightarrow 1
\end{aligned} \tag{B.74}$$

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

*Proof of Theorem 5.4.*

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

$$\hat{\omega}_{JKD}^2(x; d, \mathbf{D}_n) = \frac{n-d}{d} \binom{n}{d}^{-1} \sum_{\ell \in L_{n,d}} (\hat{\mu}_{s_1, s_2}(x; \mathbf{D}_{n,-\ell}) - \hat{\mu}_{s_1, s_2}(x; \mathbf{D}_n))^2 \quad (\text{B.75})$$

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

$$\begin{aligned} \hat{\omega}_{JKD}^2(x; d, \mathbf{D}_n) &= \frac{n-d}{d} \binom{n}{d}^{-1} \sum_{\ell \in L_{n,d}} \left( \sum_{j=1}^{s_2} \binom{s_2}{j} (H_{P_t}^j - H_{P_t}^j(\mathbf{D}_{n,-\ell})) \right)^2 \\ &= \frac{n-d}{d} \binom{n}{d}^{-1} \sum_{\ell \in L_{n,d}} \left( \sum_{j=1}^{s_2} \binom{s_2}{j} \left( \binom{n}{j}^{-1} \sum_{\iota \in L_{n,j}} h_{s_1, s_2}^{(j)}(\mathbf{D}_\iota) - \binom{n-d}{j}^{-1} \sum_{\iota \in L_j([n] \setminus \ell)} h_{s_1, s_2}^{(j)}(\mathbf{D}_\iota) \right) \right)^2 \\ &= \frac{n-d}{d} \binom{n}{d}^{-1} \sum_{\ell \in L_{n,d}} \left[ \frac{s_2}{n} \sum_{i \in \ell} h_{s_1, s_2}^{(1)}(Z_i) + \sum_{i \in [n] \setminus \ell} \left( \frac{s_2}{n} - \frac{s_2}{n-d} \right) h_{s_1, s_2}^{(1)}(Z_i) \right. \\ &\quad \left. + \sum_{j=2}^{s_2} \binom{s_2}{j} \left( \binom{n}{j}^{-1} \sum_{\iota \in L_{n,j}} h_{s_1, s_2}^{(j)}(\mathbf{D}_\iota) - \binom{n-d}{j}^{-1} \sum_{\iota \in L_j([n] \setminus \ell)} h_{s_1, s_2}^{(j)}(\mathbf{D}_\iota) \right) \right]^2 \\ &= \frac{n-d}{d} \binom{n}{d}^{-1} \left( \frac{s_2}{n} \right)^2 \sum_{\ell \in L_{n,d}} \left[ \sum_{i \in \ell} h_{s_1, s_2}^{(1)}(Z_i) - \frac{d}{n-d} \sum_{i \in [n] \setminus \ell} h_{s_1, s_2}^{(1)}(Z_i) \right. \\ &\quad \left. + \frac{n}{s_2} \sum_{j=2}^{s_2} \binom{s_2}{j} \left( \binom{n}{j}^{-1} \sum_{\iota \in L_{n,j}} h_{s_1, s_2}^{(j)}(\mathbf{D}_\iota) - \binom{n-d}{j}^{-1} \sum_{\iota \in L_j([n] \setminus \ell)} h_{s_1, s_2}^{(j)}(\mathbf{D}_\iota) \right) \right]^2 \\ &=: (n-d) \binom{n}{d}^{-1} \left( \frac{s_2}{n} \right)^2 \sum_{\ell \in L_{n,d}} \left[ \frac{1}{\sqrt{d}} \sum_{i \in \ell} h_{s_1, s_2}^{(1)}(Z_i) + T_\ell \right]^2 \end{aligned} \quad (\text{B.76})$$

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  $T_\ell$  in the sense of Lemma A.7. Luckily, since Lemma A.7 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 5.3. Thus, consider the following for an arbitrary fixed index-subset  $\ell$  with cardinality  $d$ .

$$\begin{aligned} \mathbb{E} \left[ \left( \frac{1}{\sqrt{d}} \sum_{i \in \ell} h_{s_1, s_2}^{(1)}(Z_i) \right)^2 \right] &= \frac{1}{d} \mathbb{E} \left[ \sum_{i \in \ell} \sum_{j \in \ell} h_{s_1, s_2}^{(1)}(Z_i) h_{s_1, s_2}^{(1)}(Z_j) \right] = \frac{1}{d} \sum_{i \in \ell} \sum_{j \in \ell} \mathbb{E} \left[ h_{s_1, s_2}^{(1)}(Z_i) h_{s_1, s_2}^{(1)}(Z_j) \right] \\ &= \frac{|\ell|}{d} \cdot \mathbb{E} \left[ \left( h_{s_1, s_2}^{(1)}(Z_1) \right)^2 \right] = \zeta_{P_t, 1} \end{aligned} \quad (\text{B.77})$$

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.

$$\begin{aligned}
T_\ell &= \frac{\sqrt{d}}{n-d} \sum_{i \in [n] \setminus \ell} h_{s_1, s_2}^{(1)}(Z_i) \\
&\quad + \frac{n}{s_2 \sqrt{d}} \left\{ \sum_{j=2}^d \binom{s_2}{j} \left( \binom{n}{j}^{-1} \left( \sum_{a=1}^j \sum_{\substack{\kappa \in L_a(\ell) \\ \varrho \in L_{j-a}([n] \setminus \ell)}} h_{s_1, s_2}^{(j)}(D_{\kappa \cup \varrho}) \right) + \left( \binom{n}{j}^{-1} - \binom{n-d}{j}^{-1} \right) \sum_{\iota \in L_j([n] \setminus \ell)} h_{s_1, s_2}^{(j)}(\mathbf{D}_\iota) \right) \right. \\
&\quad \left. + \sum_{j=d+1}^{s_2} \binom{s_2}{j} \left( \binom{n}{j}^{-1} \left( \sum_{a=1}^d \sum_{\substack{\kappa \in L_a(\ell) \\ \varrho \in L_{j-a}([n] \setminus \ell)}} h_{s_1, s_2}^{(j)}(D_{\kappa \cup \varrho}) \right) + \left( \binom{n}{j}^{-1} - \binom{n-d}{j}^{-1} \right) \sum_{\iota \in L_j([n] \setminus \ell)} h_{s_1, s_2}^{(j)}(\mathbf{D}_\iota) \right) \right\} \quad (\text{B.78})
\end{aligned}$$

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

$$\begin{aligned}
T_\ell &= \frac{\sqrt{d}}{n-d} \sum_{i \in [n] \setminus \ell} h_{s_1, s_2}^{(1)}(Z_i) \\
&\quad + \frac{n}{s_2 \sqrt{d}} \sum_{j=2}^{s_2} \binom{s_2}{j} \left( \binom{n}{j}^{-1} \left( \sum_{a=1}^j \sum_{\substack{\kappa \in L_a(\ell) \\ \varrho \in L_{j-a}([n] \setminus \ell)}} h_{s_1, s_2}^{(j)}(D_{\kappa \cup \varrho}) \right) + \left( \binom{n}{j}^{-1} - \binom{n-d}{j}^{-1} \right) \sum_{\iota \in L_j([n] \setminus \ell)} h_{s_1, s_2}^{(j)}(\mathbf{D}_\iota) \right) \quad (\text{B.79})
\end{aligned}$$

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

$$\begin{aligned}
\mathbb{E} \left[ (T_\ell)^2 \right] &= \frac{d}{n-d} V_{s_1, s_2}^1 \\
&\quad + \frac{n^2}{s_2^2 d} \sum_{j=2}^d \binom{s_2}{j}^2 \left( \binom{n}{j}^{-2} \sum_{a=1}^j \left[ \binom{d}{a} \binom{n-d}{j-a} \right] + \left[ \binom{n}{j}^{-1} - \binom{n-d}{j}^{-1} \right]^2 \binom{n-d}{j} \right) V_{s_1, s_2}^j \\
&\quad + \frac{n^2}{s_2^2 d} \sum_{j=d+1}^{s_2} \binom{s_2}{j}^2 \left( \binom{n}{j}^{-2} \sum_{a=1}^d \left[ \binom{d}{a} \binom{n-d}{j-a} \right] + \left[ \binom{n}{j}^{-1} - \binom{n-d}{j}^{-1} \right]^2 \binom{n-d}{j} \right) V_{s_1, s_2}^j \\
&\stackrel{(\star)}{=} \frac{d}{n-d} V_{s_1, s_2}^1 \\
&\quad + \frac{n^2}{s_2^2 d} \sum_{j=2}^d \binom{s_2}{j}^2 \binom{n}{j}^{-2} \left( \binom{n}{j} - \binom{n-d}{j} + \left[ 1 - \binom{n}{j} \binom{n-d}{j}^{-1} \right]^2 \binom{n-d}{j} \right) V_{s_1, s_2}^j \\
&\quad + \frac{n}{s_2 d} \sum_{j=d+1}^{s_2} \frac{\binom{s_2-1}{j-1} \binom{n-d}{j}}{\binom{n-1}{j-1} \binom{n}{j}} \left( \sum_{a=1}^d \frac{\binom{d}{a} \binom{n-d}{j-a}}{\binom{n-d}{j}} + \left[ 1 - \binom{n}{j} \binom{n-d}{j}^{-1} \right]^2 \right) \binom{s_2}{j} V_{s_1, s_2}^j \quad (\text{B.80})
\end{aligned}$$

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{B.81}$$

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.

$$\begin{aligned}
\mathbb{E} \left[ (T_\ell)^2 \right] &\lesssim \frac{d}{n-d} V_{s_1, s_2}^1 + \frac{n}{s_2 d} \sum_{j=2}^{s_2} \left( \frac{e(s_2-1)}{n-1} \right)^{j-1} \frac{j(n-2)}{(n-d-s_2+1)(n-d-2s_t)} \left[ \binom{s_2}{j} V_{s_1, s_2}^j \right] \\
&= \frac{d}{n-d} V_{s_1, s_2}^1 + \frac{2e \cdot n(n-2)}{(n-1)(n-d-s_2+1)(n-d-2s_t)d} \sum_{j=2}^{s_2} \frac{j}{s_2} \left( \frac{e(s_2-1)}{n-1} \right)^{j-1} \left[ \binom{s_2}{j} V_{s_1, s_2}^j \right] \\
&\lesssim \frac{d}{n-d} V_{s_1, s_2}^1 + \frac{4e}{(n-d-s_2)s_2 d} \sum_{j=2}^{s_2} j \left( \frac{e(s_2-1)}{n-1} \right)^{j-1} \left[ \binom{s_2}{j} V_{s_1, s_2}^j \right] \\
&\leq \frac{d}{n-d} V_{s_1, s_2}^1 + \frac{4e}{(n-d-s_2)s_2 d} \sum_{j=2}^{s_2} \left[ \binom{s_2}{j} V_{s_1, s_2}^j \right] + \frac{4e}{(n-d-s_2)s_2 d} \sum_{j=2}^{s_2} (j-1) \left( \frac{e(s_2-1)}{n-1} \right)^{j-1} \left[ \binom{s_2}{j} V_{s_1, s_2}^j \right] \\
&\leq \frac{d}{n-d} V_{s_1, s_2}^1 + \frac{4e}{(n-d-s_2)s_2 d} \sum_{j=2}^{s_2} \left[ \binom{s_2}{j} V_{s_1, s_2}^j \right] + \frac{4e \cdot \zeta_{s_1, s_2}^{s_2}}{(n-d-s_2)s_2 d} \sum_{j=1}^{\infty} j \left( \frac{e(s_2-1)}{n-1} \right)^j \\
&= \frac{d}{n-d} \zeta_{s_1, s_2}^1 + \frac{4e}{(n-d-s_2)s_2 d} (\zeta_{s_1, s_2}^{s_2} - s_2 \zeta_{s_1, s_2}^1) + \frac{4e \cdot \zeta_{s_1, s_2}^{s_2}}{(n-d-s_2)s_2 d} \cdot \frac{e(s_2-1)(n-1)}{(n-1-e(s_2-1))^2} \\
&= \left( \frac{d}{n-d} + \frac{e(s_2-1)(n-1)}{(n-1-e(s_2-1))^2} \right) \zeta_{s_1, s_2}^1 + \frac{4e}{(n-d-s_2)s_2 d} \left( 1 + \frac{e(s_2-1)(n-1)}{(n-1-e(s_2-1))^2} \right) (\zeta_{s_1, s_2}^{s_2} - s_2 \zeta_{s_1, s_2}^1)
\end{aligned} \tag{B.82}$$

We continue as in the default Jackknife case.

$$\begin{aligned}
\frac{\mathbb{E} [T_\ell^2]}{V_{s_1, s_2}^1} &\leq \frac{d}{n-d} + \frac{e(s_2-1)(n-1)}{(n-1-e(s_2-1))^2} + \frac{4e}{(n-d-s_2)s_2 d} \left( 1 + \frac{e(s_2-1)(n-1)}{(n-1-e(s_2-1))^2} \right) \frac{\zeta_{s_1, s_2}^{s_2} - s_2 \zeta_{s_1, s_2}^1}{\zeta_{s_1, s_2}^1} \\
&\rightarrow 0.
\end{aligned} \tag{B.83}$$

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. ■

## B.10 Pointwise Inference Results

*Proof of Theorem 5.7.*

■

*Proof of Theorem 5.8.*

■