# Balancing Flexibility and Interpretability: A Conditional Linear Model Estimation via Random Forest

Ricardo Masini [1*]        Marcelo Medeiros[2]

February 20, 2025

**Abstract**

Traditional parametric econometric models often rely on rigid functional forms, while nonparametric techniques, despite their flexibility, frequently lack interpretability. This paper proposes a parsimonious alternative by modeling the outcome $Y$ as a linear function of a vector of variables of interest $\boldsymbol{X}$, conditional on additional covariates $\boldsymbol{Z}$. Specifically, the conditional expectation is expressed as $\mathbb{E}[Y|\boldsymbol{X},\boldsymbol{Z}] = \boldsymbol{X}^{\mathsf{T}}\boldsymbol{\beta}(\boldsymbol{Z})$, where $\boldsymbol{\beta}(\cdot)$ is an unknown Lipschitz-continuous function. We introduce an adaptation of the Random Forest (RF) algorithm to estimate this model, balancing the flexibility of machine learning methods with the interpretability of traditional linear models. This approach addresses a key challenge in applied econometrics by accommodating heterogeneity in the relationship between covariates and outcomes. Furthermore, the heterogeneous partial effects of $\boldsymbol{X}$ on $Y$ are represented by $\boldsymbol{\beta}(\cdot)$ and can be directly estimated using our proposed method. Our framework effectively unifies established parametric and nonparametric models, including varying-coefficient, switching regression, and additive models. We provide theoretical guarantees, such as pointwise and $\mathcal{L}^p$-norm rates of convergence for the estimator, and establish a pointwise central limit theorem through subsampling, aiding inference on the function $\boldsymbol{\beta}(\cdot)$. We present Monte Carlo simulation results to assess the finite-sample performance of the method.

**Keywords**: random forests; heterogeneous partial effects; machine learning.

---

[1]Department of Statistics, University of California, Davis.

[2]Department of Economics, University of Illinois, Urbana-Champaign.

[*]Corresponding author: rmasini@ucdavis.edu

# Contents

# 1 Introduction

In economics and related social sciences, epidemiology and medicine, psychology, and many other areas, estimating the partial effects (causal or not) of one or more factors on a target variable is extremely important. In general, partial effects estimation is conducted either on parametric models with strong functional-form assumptions or in overly simplified semi-parametric alternatives. These models provide interpretive clarity and computational efficiency but often impose restrictive assumptions that may inadequately capture the complexity of real-world data. In recent years, machine learning (ML) methods have broadened the statistician/econometrician's toolkit by offering more flexible approaches to modeling relationships within potentially high-dimensional datasets without imposing stringent parametric assumptions; see, for example, Athey and Imbens (2019) and Masini et al. (2023) for recent review papers. Despite these advances, the challenge of balancing flexibility with interpretability remains a significant issue for applied researchers.

This paper presents a locally linear model that addresses the challenge of integrating the flexibility of machine learning (ML) methods with the interpretability of linear models. The proposed model aligns with several well-established parametric and semi-parametric specifications in the literature, including switching regression (Dagenais, 1969; Goldfeld and Quandt, 1972), varying-coefficient models, and additive models (Hastie and Tibishirani, 1993; Chen and Tsay, 1993). It also incorporates recent advancements in machine learning, such as Random Forests (RF) introduced by Breiman (2001), Generalized Random Forests (GRF) presented by Athey et al. (2019), and local linear forests (LLF) discussed in Friedberg et al. (2021). Our central contribution is a robust framework that enables the nonparametric estimation of heterogeneous partial effects of one or more variables of interest on an outcome. Our approach is intuitive and computationally simple, achieved through an adaptation of the RF method.

## 1.1 Motivation

Let $Y$ be a response random variable and define $h(\boldsymbol{x}, \boldsymbol{z}) := \mathbb{E}(Y | \boldsymbol{X} = \boldsymbol{x}, \boldsymbol{Z} = \boldsymbol{z})$, where $h(\cdot)$ is an unknown Lipschitz continuous function of two vectors of random covariates $\boldsymbol{X}$ and $\boldsymbol{Z}$. Suppose the goal is to estimate $h(\boldsymbol{x}, \boldsymbol{z})$ for arbitrary data points $\boldsymbol{x}$ and $\boldsymbol{z}$, as well as the partial effects $\partial h(\boldsymbol{x}, \boldsymbol{z}) / \partial \boldsymbol{x}$ assuming $h$ is differentiable with respect to its first argument. Modern machine learning methods provide a range of nonparametric algorithms to estimate $h(\boldsymbol{x}, \boldsymbol{z})$, enabling greater flexibility in modeling intricate relationships. However, these models often lack interpretability, and the

computation of partial effects can be both intensive and non-trivial, particularly in the context of deep learning models.[1] Moreover, conducting inference on partial effects within the framework of general machine learning models remains an open problem, necessitating further research and development to enhance methodological rigor and applicability in statistical analysis.

This paper presents an alternative approach rooted in contemporary machine learning literature, specifically aimed at enhancing interpretability. We consider covariates $\boldsymbol{X}$ and $\boldsymbol{Z}$, which are not necessarily mutually independent, with a focus on the estimation $\partial h(\boldsymbol{x}, \boldsymbol{z})/\partial \boldsymbol{x}$. We propose a model in which the conditional expectation of $Y$, given $\boldsymbol{X}$ and $\boldsymbol{Z}$, is represented as follows:

$$h(\boldsymbol{x}, \boldsymbol{z}) := \mathbb{E}(Y|\boldsymbol{X} = \boldsymbol{x}, \boldsymbol{Z} = \boldsymbol{z}) = \boldsymbol{\beta}(\boldsymbol{z})^{\mathsf{T}}\boldsymbol{x}.$$

This is a varying-coefficient model where the coefficients of the linear relationship between $\boldsymbol{X}$ and $Y$ vary as a function of $\boldsymbol{Z}$. Importantly, the partial effect of $\boldsymbol{X}$ on $Y$, $\partial h(\boldsymbol{x}, \boldsymbol{z})/\partial \boldsymbol{x}$, is given directly by $\boldsymbol{\beta}(\boldsymbol{z})$, and thus exhibits heterogeneity across different values of $\boldsymbol{Z}$. Estimating the map $\boldsymbol{z} \mapsto \boldsymbol{\beta}(\boldsymbol{z})$ is equivalent to recovering the heterogeneous partial effects, providing a natural and intuitive model interpretation. Clearly, this is only the case when $\boldsymbol{Z}$ and $\boldsymbol{X}$ do not share the same covariates. We consider the estimation of $\boldsymbol{\beta}(\boldsymbol{z})$ by modification of the Random Forest (RF) method, where the trees have a linear model on each of their leaves. If $\boldsymbol{X}$ is a fixed scalar, the model equals the Random Forest (RF) method, where the trees have an intercept model on each of their leaves.

This model presents two primary advantages. First, it facilitates flexible, data-driven estimation of complex relationships, all while ensuring the interpretability of partial effects, which is often a pivotal focus in applied research. Second, by integrating the flexibility of the Random Forest (RF) algorithm within a locally linear framework, we effectively capture heterogeneity in partial effects in a computationally feasible manner. This characteristic renders the approach particularly well-suited for large-scale empirical applications.

In many empirical applications, the dimensionality of $\boldsymbol{X}$ is expected to be small relative to the sample size (it is not uncommon to have an univariate $\boldsymbol{X}$), simplifying the application of our model. Partial effects of covariates can be estimated directly using modern machine learning and nonparametric techniques. We choose the RF algorithm for its robustness and flexibility, mainly because it requires fewer hyperparameter choices, mitigating the risks of overfitting and cherry-picking. Additionally, RF's efficient estimation algorithms make it computationally attractive.

---

[1]Deep learning methods usually require the definition of too many hyperparameters, the network architecture, and also require large amounts of data for effective training.

While alternative approaches like deep learning offer comparable flexibility, they are more sensitive to hyperparameter tuning and initial conditions, require larger datasets, and are computationally intensive. By contrast, RF offers a more stable and efficient solution.

## 1.2 Main contributions and comparison to the literature

The concept of introducing nonlinearity through varying coefficients in linear models is well-established. It originated in threshold regression models by Dagenais (1969), Quandt (1972), Goldfeld and Quandt (1972), Kiefer (1978), Tong and Lim (1980), and later work by Tsay (1989). These studies focus on sharp parameter changes based on a univariate threshold variable, capturing structural shifts in the data. Smooth transitions (shifts) were introduced by Chan and Tong (1986) and Teräsvirta (1994), but these models relied on a single transition variable. Extensions to multiple transition variables inspired by the neural network (NN) literature were explored by Medeiros and Veiga (2000), Medeiros and Veiga (2005), and Suarez-Fariñas et al. (2004), though these approaches were parametric and limited to low-dimensional settings. Nonparametric alternatives have been proposed by Hastie and Tibishirani (1993), Fan and Zhang (1999), and Cai et al. (2000). However, estimating these models becomes challenging with an increase in the number of covariates. In contrast, our semi-parametric approach based on random forests accommodates a large set of covariates, provided the dimensionality remains smaller than the sample size.

Estimating regression trees with linear models in the terminal nodes is not new and is nested in the more general Generalized Random Forest model by Athey et al. (2019). Friedberg et al. (2021) formalized the locally linear random forests and recommended a local Ridge regression to estimate the parameters. Nevertheless, this paper differs and complements the ones cited above in a few ways. First, we estimate the local linear models with the usual ordinary least-squares method, making our algorithm simpler than the ones in Athey et al. (2019) and Friedberg et al. (2021). Second, as our primary focus is on the estimation of the partial effects $\boldsymbol{\beta}(\boldsymbol{z})$ and not $\mathbb{E}(Y|\boldsymbol{X} = \boldsymbol{x}, \boldsymbol{Z} = \boldsymbol{z})$, we complement the previous papers by deriving a new consistency and asymptotic normality results for $\widehat{\boldsymbol{\beta}}(\boldsymbol{z})$. We also worked out the rates of convergence and derived the asymptotic covariance matrix for the estimator of $\widehat{\boldsymbol{\beta}}(\boldsymbol{Z})$. Third, we proposed a consistent estimator for the covariance matrix of $\widehat{\boldsymbol{\beta}}(\boldsymbol{Z})$. Fourth, unlike most papers in the literature, we show that our results are also valid when $\boldsymbol{Z}$ contains discrete random variables. Note that having discrete-valued elements of $\boldsymbol{Z}$ is not the same as introducing interaction effects between $\boldsymbol{X}$ and dummy variables constructed from the classes in $\boldsymbol{Z}$. Our approach allows the partial effect heterogeneity to be determined by unknown interactions

4

among the elements of $\boldsymbol{Z}$. Finally, based on our new convergence results, we derived two tests to conduct inference on the partial effects and test whether the partial effects are homogeneous.

## 1.3 Outline of the Paper

Following this introduction, the paper is structured as follows. In Section 2, we define the model, outline assumptions related to the data-generating mechanism, and present special cases of our proposal. In Section 3, we present the theoretical results of this paper. More specifically, Section 3.1 details the main convergence results, while Sections 3.2.1 and 3.2.2 provide descriptions of the specification tests proposed in the paper. The case of discrete random variables is examined in Section 3.3. Monte Carlo simulations are illustrated in Section 4, and empirical samples are presented in Section 5. Finally, Section 6 wraps up the paper. All technical derivations are provided in the Supplemental Material.

## 1.4 Notation

Vectors are denoted by bold lowercase $\boldsymbol{x}$ and matrices by bold uppercase $\boldsymbol{M}$. $\boldsymbol{x}^\mathsf{T}$ denotes the transpose of the vector $\boldsymbol{x}$. Similarly $\boldsymbol{M}^\mathsf{T}$ denotes the transpose for matrix $\boldsymbol{M}$. We write $\|\boldsymbol{x}\|_p$ for $p \in [1, \infty]$ to denote the $\ell^p$-norm if $\boldsymbol{x}$ is a (possibly random) vector or the induced operator $\ell^p$–$\ell^p$-norm if $\boldsymbol{M}$ is a matrix. For a matrix $\boldsymbol{M}$, we write $\|\boldsymbol{M}\|_{\max}$ for the maximum absolute entry and $\|\boldsymbol{M}\|_\mathrm{F}$ for the Frobenius norm. We denote positive semi-definiteness by $\boldsymbol{M} \succeq 0$ and write $\boldsymbol{I}_d$ for the $d \times d$ identity matrix.

For scalar sequences $x_n$ and $y_n$, we write $x_n \lesssim y_n$ if there exists a positive constant $C$ such that $|x_n| \leq C|y_n|$ for sufficiently large $n$. We write $x_n \asymp y_n$ to indicate both $x_n \lesssim y_n$ and $y_n \lesssim x_n$. Similarly, for random variables $X_n$ and $Y_n$, we write $X_n \lesssim_\mathbb{P} Y_n$ if for every $\varepsilon > 0$ there exists a positive constant $C$ such that $\mathbb{P}(|X_n| \geq C|Y_n|) \leq \varepsilon$, and write $X_n \xrightarrow{\mathbb{P}} X$ and $X_n \xrightarrow{d} X$ for limits in probability and in distribution, respectively. For real numbers $a$ and $b$ we use $a \wedge b = \max\{a, b\}$ and $a \vee b = \min\{a, b\}$.

## 2 Setup

We consider the following model.

**Assumption 1** (Locally Linear Model)
*Let $Y$ be an integrable random variable and $\boldsymbol{X}, \boldsymbol{Z}$ be random vectors taking value on $\mathbb{R}^{d_Z}$ and $\mathbb{R}^{d_X}$*

respectively that do not share common variables. We assume that for some Lipschitz function $\boldsymbol{\beta} : \mathcal{Z} \subseteq \mathbb{R}^{d_Z} \to \mathbb{R}^{d_X}$

$$h(\boldsymbol{x}, \boldsymbol{z}) := \mathbb{E}[Y|\boldsymbol{X} = \boldsymbol{x}, \boldsymbol{Z} = \boldsymbol{z}] = \boldsymbol{x}^{\mathsf{T}} \boldsymbol{\beta}(\boldsymbol{z}). \tag{1}$$

The local-linear model (1) has the advantage of having the marginal treatment effect built-in since $h(\boldsymbol{x}, \boldsymbol{z})/\partial \boldsymbol{x} = \boldsymbol{\beta}(\boldsymbol{z})$. At the same time, the model is flexible enough in terms of the covariates $\boldsymbol{Z}$ to accommodate complex heterogeneous partial effects, which include higher-order interactions between $\boldsymbol{X}$ and $\boldsymbol{Z}$. Model (1) nests several notable cases of interest.

Since $\mathbb{E}[Y|\boldsymbol{X}, \boldsymbol{Z}]$ is the minimizer of $\mathbb{E}[(Y - f(\boldsymbol{X}, \boldsymbol{Z}))^2]$ over the class of functions $f$ such that $\mathbb{E}[f(\boldsymbol{X}, \boldsymbol{Z})^2] < \infty$, we can explicitly characterize the function $\boldsymbol{z} \mapsto \beta(\boldsymbol{z})$ such that $\boldsymbol{\beta}(\boldsymbol{z}) = \boldsymbol{\Omega}(\boldsymbol{z})^{-1} \boldsymbol{\gamma}(\boldsymbol{z})$, where $\boldsymbol{\Omega}(\boldsymbol{z}) := \mathbb{E}[\boldsymbol{X}\boldsymbol{X}^{\mathsf{T}}|\boldsymbol{Z} = \boldsymbol{z}]$ and $\boldsymbol{\gamma}(\boldsymbol{z}) := \mathbb{E}[\boldsymbol{X}Y|\boldsymbol{Z} = \boldsymbol{z}]$, and provided that $\boldsymbol{\Omega}(\boldsymbol{z})$ is almost sure positive definite.

**Example 1** (Heterogeneous treatment effects)

*Consider a collection of potential outcomes $\{Y^*(t) : t \in \mathcal{T}\}$ where $Y(t)$ is a random variable and $\mathcal{T} \subseteq \mathbb{R}$. In a binary treatment case $\mathcal{T} = \{0, 1\}$. Here we allow continuous treatment such as $\mathcal{T} = [a, b]$ for $a < b$ or $\mathcal{T} = \mathbb{R}$. In the continuous treatment literature, the function $t \mapsto Y^*(t)$ is called the dose-response function, and $t \mapsto \mathbb{E}[Y^*(t)]$ is the average dose function. For each unit in the sample, we observe the realization of the treatment status $T$, a vector of covariates $\boldsymbol{Z}$ taking values on $\mathbb{R}^{d_Z}$, and the potential outcome corresponding to the treatment received $Y := Y^*(t)$.*

*The interest relies on estimating the (potentially) heterogeneous (marginal) treatment effect.*

$$\tau(\boldsymbol{z}) := \begin{cases} \frac{d\mathbb{E}[Y^*(t)|\boldsymbol{Z}=\boldsymbol{z}]}{dt} & \text{if } \mathcal{T} \text{ is an interval,} \\ \mathbb{E}[Y^*(1) - Y^*(0)|\boldsymbol{Z} = \boldsymbol{z}] & \text{if } \mathcal{T} = \{0, 1\}. \end{cases}$$

*By definition we have that $\mathbb{E}[Y^*(t)|T = t] = \mathbb{E}[Y|T = t]$ for $t \in \mathcal{T}$. However, in general, we have $\mathbb{E}[Y^*(t)] \neq \mathbb{E}[Y|T = t]$ due to confounding effects. If conditional on the covariates $\boldsymbol{Z}$ the average dose-response function is mean independent of the treatment, we can consider the following model*

$$Y = \beta(\boldsymbol{Z})T + U, \tag{2}$$

*where $\mathbb{E}(U|\boldsymbol{Z}, T) = 0$. We can also include additional controls in the above model. For example, let $\boldsymbol{W}$ be a set of control variables and write*

$$Y = \beta(\boldsymbol{Z})T + \boldsymbol{\delta}^{\mathsf{T}}\boldsymbol{W} + U, \tag{3}$$

6

where $U$ is an error term such that $\mathbb{E}(U|\boldsymbol{Z}, T, \boldsymbol{W}) = 0$. Equations (2) and (3) are models where the treatment effects are not constant and may vary with subject characteristics defined by the vector $\boldsymbol{Z}$. Examples of such specifications can be found at Hansen (2000), Athey and Imbens (2016), or Chernozhukov et al. (2018), among many others.

**Example 2** (Smooth Transition Regression)

*Suppose that the coefficients of a regression model change smoothly and monotonically between zero and one according to a scalar variable $Z_i$ such that*

$$Y = \boldsymbol{\beta}_0^\mathsf{T} \boldsymbol{X} + g(\boldsymbol{Z})\boldsymbol{\beta}_1^\mathsf{T} \boldsymbol{X} + U, \tag{4}$$

*where $\lim_{z \to \infty} g(\boldsymbol{z}) = 1$ and $\lim_{z \to -\infty} g(\boldsymbol{z}) = 0$ and $\mathbb{E}(U|\boldsymbol{X}, Z) = 0$. In this case, $\boldsymbol{\beta}(\boldsymbol{z}) = \boldsymbol{\beta}_0 + g(\boldsymbol{z})\boldsymbol{\beta}_1$. This is a general case of a two-regime smooth transition regression model. An interesting example of such specification is the nonlinear Phillips curve model discussed in Areosa et al. (2011).*

**Example 3** (Grouped patterns of heterogeneity)

*Consider a model where the intercept varies according to a set of variables $\boldsymbol{Z}$. In this case,*

$$Y = \boldsymbol{\beta}(\boldsymbol{Z}) + \boldsymbol{\delta}^\mathsf{T} \boldsymbol{X} + U, \tag{5}$$

*where $\mathbb{E}(U|\boldsymbol{X}, Z) = 0$. This type of specification is commonly employed to model clustering behavior and has gained traction in the empirical economic growth literature. See, for example, Bonhomme and Manresa (2015). This specification is also related to the binscatter methods discussed in Cattaneo et al. (2024).*

Given a random sample $\{(Y_i, \boldsymbol{X}_i, \boldsymbol{Z}_i) : 1 \leq i \leq n\}$ of $(Y, \boldsymbol{X}, \boldsymbol{Z})$ we propose to estimate (1) by a modification of the Random Forest procedure (refer to Algorithm 1). We start by defining for any nonempty set of indices $\mathcal{I} \subseteq [n]$, the residual sum of squares ($\mathsf{RSS}$) of a least-square estimation (we will impose conditions such that the estimator is well-defined). Hence, write

$$\mathsf{RSS}(\mathcal{I}) := \sum_{i \in \mathcal{I}} \left[ Y_i - \boldsymbol{X}_i^\mathsf{T} \widehat{\boldsymbol{\beta}}(\mathcal{I})^2 \right], \qquad \text{and} \qquad \widehat{\boldsymbol{\beta}}(\mathcal{I}) := \left[ \sum_{i \in \mathcal{I}} \boldsymbol{X}_i \boldsymbol{X}_i^\mathsf{T} \right]^{-1} \sum_{i \in \mathcal{I}} \boldsymbol{X}_i Y_i.$$

Also, for any element of $\boldsymbol{Z}$ indexed by $j \in [d_Z]$, and $\delta \in [0, 1]$ define

$$\Delta(\mathcal{I}, j, \delta) := \mathsf{RSS}(\{i \in \mathcal{I} : Z_{ij} \leq \delta\}) + \mathsf{RSS}(\{i \in \mathcal{I} : Z_{ij} > \delta\}). \tag{6}$$

The optimum splitting point $\delta^*$ of the index set $\mathcal{I}$ along direction $j$ is given by

$$\delta^* := \delta^*(\mathcal{I}, j) \in \underset{\delta \in [0,1]}{\arg\min} \, \Delta(R, j, \delta). \tag{7}$$

Note that we might take $\delta^* \in \{Z_{i,j} : i \in \mathcal{I}\}$. Finally, define the left and right child nodes of $\mathcal{I}$ by

$$\mathcal{I}^- := \{i \in \mathcal{I} : Z_{ij} \le \delta^*\}, \qquad \text{and} \qquad \mathcal{I}^+ := \{i \in \mathcal{I} : Z_{ij} > \delta^*\}.$$

Starting for an initial index set $\mathcal{B} \subseteq [n]$ and repeating the steps above, we end up with a partition of $[0,1]^{d_Z}$ into disjoint rectangles with axis-aligned sides (leaves). Furthermore, since the initial index set (root node) and the slipt direction $j$ are randomly chosen independent of the data, we denote by $\omega$ this independent source of randomness in the algorithm.

For $\boldsymbol{z} \in [0,1]^{d_Z}$, let $R(\boldsymbol{z}, \omega)$ denotes the unique (random) leaf containing $\boldsymbol{z}$, i.e., $R(\boldsymbol{z}, \omega)$ is a random rectangle that depends on the observation indexed by $\mathcal{B}$ and an external source of randomness. For an nonempty $\mathcal{A} \subseteq [n]$ such that $\mathcal{A} \cap \mathcal{B} = \emptyset$, define $\mathcal{A}(\boldsymbol{z}, \omega) = \{i \in \mathcal{A} : \boldsymbol{Z}_i \in R(\boldsymbol{z}, \omega)\}$, $\widehat{m}(\boldsymbol{x}, \boldsymbol{z}, \omega) := \boldsymbol{x}^\mathsf{T} \widehat{\boldsymbol{\beta}}(\boldsymbol{z}, \omega)$, and $\widehat{\boldsymbol{\beta}}(\boldsymbol{z}, \omega) := \widehat{\boldsymbol{\beta}}(\mathcal{A}(\boldsymbol{z}, \omega))$.

More explicitly, write

$$\widehat{\boldsymbol{\beta}}(\boldsymbol{z}, \omega) = \widehat{\boldsymbol{\Omega}}(\boldsymbol{z}, \omega)^{-1} \widehat{\boldsymbol{\gamma}}(\boldsymbol{z}, \omega), \tag{8}$$

where $\widehat{\boldsymbol{\Omega}}(\boldsymbol{z}, \omega) := \frac{1}{|\mathcal{A}(\boldsymbol{z}, \omega)|} \sum_{i \in \mathcal{A}(\boldsymbol{z}, \omega)} \boldsymbol{X}_i \boldsymbol{X}_i^\mathsf{T}$ and $\widehat{\boldsymbol{\gamma}}(\boldsymbol{z}, \omega) := \frac{1}{|\mathcal{A}(\boldsymbol{z}, \omega)|} \sum_{i \in \mathcal{A}(\boldsymbol{z}, \omega)} \boldsymbol{X}_i Y_i$. Finally, for $B \ge 1$, we propose to estimate $m(\boldsymbol{x}, \boldsymbol{z})$ using $\overline{m}(\boldsymbol{x}, \boldsymbol{z}) := \frac{1}{B} \sum_{b=1}^{B} \widehat{m}(\boldsymbol{X}, \boldsymbol{z}, \omega_b) = \boldsymbol{X}^\mathsf{T} \overline{\boldsymbol{\beta}}(\boldsymbol{z})$ with

$$\overline{\boldsymbol{\beta}}(\boldsymbol{z}) := \frac{1}{B} \sum_{b=1}^{B} \widehat{\boldsymbol{\beta}}(\boldsymbol{z}, \omega_b) \tag{9}$$

where $\{\omega_b : b \in [B]\}$ is an independent and identically sequence independent of $\{(Y_i, \boldsymbol{X}_i, \boldsymbol{Z}_i) : i \in [n]\}$. We also investigate the properties of the following estimator:

$$\check{\boldsymbol{\beta}}(\boldsymbol{z}) := \overline{\boldsymbol{\Omega}}(\boldsymbol{z})^{-1} \overline{\boldsymbol{\gamma}}(\boldsymbol{z}),$$

where $\overline{\boldsymbol{\Omega}}(\boldsymbol{z}) := \frac{1}{B} \sum_{b=1}^{B} \widehat{\boldsymbol{\Omega}}(\boldsymbol{z}, \omega_b)$ and $\overline{\boldsymbol{\gamma}}(\boldsymbol{z}) := \frac{1}{B} \sum_{b=1}^{B} \widehat{\boldsymbol{\gamma}}(\boldsymbol{z}, \omega_b)$.

---
**Algorithm 1** Local-Linear Double-sample tree
---
**Require:** $n \geq 2$ training samples of the form $(y_i, \boldsymbol{x}_i, \boldsymbol{z}_i) \in (\mathbb{R} \times \mathbb{R}^{d_X} [0,1]^{d_Z})$, a subsample size $s < n$,

   minimum leaf size $k \leq s$, number of trees $B \in \mathbb{N}$, minimum fraction of observation in each child

   node $\alpha \in (0, 0.5)$ and the probability of selecting a variable to split is at least $\pi/d_Z$ for $\pi \in (0, 1]$.

   **for** $b = 1, \ldots, B$ **do**

        1. Draw a random subsample of size $s$ for $\{1, \ldots, n\}$ without replacement

        2. Partition the subsample into $\mathcal{A} \cup \mathcal{B}$ such that $|\mathcal{A}| = \lfloor s/2 \rfloor$ and $|\mathcal{B}| = \lceil s/2 \rceil$

        3. Grow a tree using only data from $\mathcal{B}$ via standard CART algorithm

        **while** the number of observations in every leaf is not between $k$ and $2k - 1$ **do**

            Sample a variable $Z_j$ from $\{Z_1, \ldots, Z_{d_Z}\}$

            Decide the cutoff in $Z_j$ based on criteria (7) imposing the $\alpha$-regularity condition above

        **end while**

        4. Estimate the leaves responses using only data from $\mathcal{A}$ as in (8)

   **end for**

   Compute the random forest estimator given in (9)
---

# 3  Main Results

In this section, we state the main results of the paper and their underlying assumptions.

## 3.1  Estimation

**Assumption 2**

*We consider the following conditions:*

(a) *(Sampling) The training sample $\{(Y_i, \boldsymbol{X}_i, \boldsymbol{Z}_i) : i \in [n]\}$ is $n$ independent copies of $(Y, \boldsymbol{X}, \boldsymbol{Z})$ taking values on $\mathbb{R} \times \mathbb{R}^{d_X} \times [0,1]^{d_Z}$ and $\boldsymbol{X}_i$ has finite moments up to order $4$.*

(b) *(Density) $\boldsymbol{Z}$ admits a density $f$ that is bounded away from zero and infinity*

(c) *(Smoothness) $\boldsymbol{z} \mapsto \boldsymbol{\beta}(\boldsymbol{z})$, $\boldsymbol{z} \mapsto \boldsymbol{\Omega}(\boldsymbol{z})$, $\boldsymbol{z} \mapsto \mathbb{E}[\boldsymbol{X}(Y - \mathbb{E}[Y|\boldsymbol{X}, \boldsymbol{Z}])|\boldsymbol{Z} = \boldsymbol{z}]$ and $\boldsymbol{z} \mapsto \mathbb{V}[\boldsymbol{X}(Y - \mathbb{E}[Y|\boldsymbol{X}, \boldsymbol{Z}]|\boldsymbol{Z} = \boldsymbol{z}]$ are Lipschitz-continuous, and $\boldsymbol{z} \mapsto \mathbb{V}[X_j X_k|\boldsymbol{Z} = \boldsymbol{z}]$ and $\boldsymbol{z} \mapsto \mathbb{V}[X_j(Y - \mathbb{E}[Y|\boldsymbol{X}, \boldsymbol{Z}]|\boldsymbol{Z} = \boldsymbol{z}]$ are bounded for $j, k \in [d_Z]$;*

(d) *(Conditional Moment Lower bounds) The matrices $\mathbb{E}[\boldsymbol{X}\boldsymbol{X}^{\mathsf{T}}|\boldsymbol{Z} = \boldsymbol{z}]$ and $\mathbb{V}[\boldsymbol{X}(Y - \mathbb{E}[Y|\boldsymbol{X}, \boldsymbol{Z}])|\boldsymbol{Z} = \boldsymbol{z}]$ are non-singular for each $\boldsymbol{z} \in [0,1]^d$;*

9

*(e) The trees are generated by Algorithm (1) with with $\alpha \in (0, 0.5)$ and $\pi \in (0, 1]$.*

In Assumption 2(e), we require the minimum of observations per leaf to increase with the same size to ensure that each tree is consistent, as opposed to the consistency of the random forest. Specifically, we need the variance of $\widehat{\boldsymbol{\Omega}}(\boldsymbol{z})$ to vanish on each tree as the sample size increases. At the same time, $k$ must grow slower than the subsampling rate $s$ so that each cell is split enough times to make its diameter shrink toward zero, and the tree bias vanishes as the sample size increases.

Recall that, for the tree constructed with randomness $\omega$, $R(\boldsymbol{z}, \omega)$ denotes the leaf containing $\boldsymbol{z}$. Define the map $\widetilde{\boldsymbol{\beta}} : [0, 1]^d \to \mathbb{R}^{d_X}$ by

$$\widetilde{\boldsymbol{\beta}}(\boldsymbol{z}) := \nabla_x \mathbb{E}[Y | \boldsymbol{X} = \boldsymbol{x}, \boldsymbol{Z} \in R(\boldsymbol{z}, \omega)] = \mathbb{E}[\boldsymbol{\beta}(\boldsymbol{Z}) | \boldsymbol{Z} \in R(\boldsymbol{z}, \omega)]; \qquad \boldsymbol{z} \in [0, 1]^d.$$

Note that in general we expect $\widetilde{\boldsymbol{\beta}}(\cdot) \neq \boldsymbol{\beta}(\cdot)$. However, due to the sample split we have that $\widehat{\boldsymbol{\beta}}(\cdot, \omega)$ is an unbiased estimator for $\widetilde{\boldsymbol{\beta}}(\cdot)$ even when $\widehat{\boldsymbol{\beta}}(\cdot, \omega)$ is biased for $\boldsymbol{\beta}(\cdot)$ (and $\widetilde{\boldsymbol{\beta}}(\cdot)$).

**Theorem 1** (Unbiasedness)

*Under Assumptions 1 and 2, $\mathbb{E}[\widehat{\boldsymbol{\beta}}(\boldsymbol{z})] = \widetilde{\boldsymbol{\beta}}(\boldsymbol{z})$ for all $\boldsymbol{z} \in [0, 1]^d$.*

**Theorem 2** (Rate of Convergence)

*Under Assumptions 1 and 2, if $k \gtrsim s^\epsilon$ for some $\epsilon \in (0, 1)$ then for any $\delta \in (0, 1)$*

$$|\widehat{\boldsymbol{\beta}}(\boldsymbol{z}, \omega) - \boldsymbol{\beta}(\boldsymbol{z})| \lesssim_{\mathbb{P}} k^{-\frac{1}{2}} \left(\frac{s}{k}\right)^{\epsilon/2} + \left(\frac{k}{s}\right)^{(1-\delta)\frac{K(\alpha)\pi}{d_Z}}; \qquad \boldsymbol{z} \in [0, 1]^d,$$

*where $K(\alpha) := \frac{\log((1-\alpha)^{-1})}{\log(1/\alpha)} \in (0, 1)$. If further then , for $2 \le p \le q$,*

$$\left[\int_{[0,1]^d} |\widehat{\boldsymbol{\beta}}(\boldsymbol{z}, \omega) - \boldsymbol{\beta}(\boldsymbol{z})|^p f(\boldsymbol{z}) dz\right]^{1/p} \lesssim_{\mathbb{P}} k^{-\frac{1}{2}} \left(\frac{s}{k}\right)^{\epsilon/2} + \left(\frac{k}{s}\right)^{(1-\delta^*)\frac{K(\alpha)\pi}{d_Z}},$$

*where $\delta^* \in (0, 1)$ is defined in Lemma 1.*

*Remark.* The second term in the right-hand size is due to the tree's bias. Given the Lipschitz condition, asymptotic unbiasedness is obtained by shrinking the leaf diameter towards zero, which is a consequence of subsequent splits in the node, which, in turn, depend on $k/s$. The second one is due to the variance within each leaf and, not surprisingly, improves for larger $k$. The bias-variance trade-off could be optimized with select $k \asymp s^\eta$ for $\eta > 0$ as a function of $d_Z$, $\alpha$, and $\pi$ such that those two terms are of the same order.

Although a tree construct using Algorithm 1 is indeed honest and symmetric, there is no guarantee to be a $k$-PNN predictor. To see that, fix a test point $\boldsymbol{z} \in [0, 1]^d$ and let $R(\boldsymbol{z}, \omega)$ denote the unique

leaf containing $z$. Recall that $R(\boldsymbol{z}, \omega)$ is independent of $\mathcal{A}$-sample due to honesty. Then the number of observations in the $R(\boldsymbol{z}, \omega)$ wich we denoted by $|\mathcal{A}(\omega, z)|$ follows a Binomial distribution conditional on $R(\boldsymbol{z}, \omega)$ with $s$ trial an probability of success $p(\boldsymbol{z}, \omega) := \mathbb{P}(\boldsymbol{Z} \in R(\boldsymbol{z}, \omega) | R(\boldsymbol{z}, \omega))$.

Therefore, the expected number of $\mathcal{A}$-sample observations in $R(\boldsymbol{z}, \omega)$ conditional on $R(\boldsymbol{z}, \omega)$ is given by $sp(\boldsymbol{z}, \omega)$ while the number of $\mathcal{B}$-sample in $R(\boldsymbol{z}, \omega)$ is between $k$ and $2k-1$ by construction. The question becomes how $|\mathcal{A}(\boldsymbol{z}, \omega)|$ relates to $k$. As it is shown in Lemma 2, $|\mathcal{A}(\omega, z)|$ can be upper and lower bound in probability as

$$\mathbb{P}\left[s\left(\frac{s}{k}\right)^{-1/K(\alpha)} \lesssim |\mathcal{A}(\boldsymbol{z}, \omega)| \lesssim s\left(\frac{s}{k}\right)^{-K(\alpha)}\right] \gtrsim 1 - \frac{1}{s(s/k)^{1/K(\alpha)}}.$$

Set $k \asymp s^\eta$ for $\eta \in [0, 1)$. For $\eta \in (1 - K(\alpha), 1)$ we have that $|A(\boldsymbol{z}, \omega)|$ diverges as $s \to \infty$ with high probability. Precisely.

$$\mathbb{P}\left[s^{1-\frac{1-\eta}{K(\alpha)}} \lesssim |\mathcal{A}(\boldsymbol{z}, \omega)| \lesssim s^{1-(1-\eta)K(\alpha)}\right] \to 1$$

For $\eta \in [0, 1 - K(\alpha)]$ and, in particular, $\eta = 0$ (fixed $k$), the above bound is vacuous. It is not clear whether, for any fixed $k > 0$, it is possible to claim that $|\mathcal{A}(\boldsymbol{z}, \omega)| \gtrsim k$ with high probability.

**Theorem 3** (Random Forest Asymptotic Normality)
*Under Assumptions 1 and 2, if $k \asymp s^\eta$ for $\eta \in (0, 1)$ and $s \asymp n^\omega$. where $\omega \in \left(\left[(1 - \eta)\left(\frac{\pi K(\alpha)}{d_Z} + \frac{1}{K(\alpha)}\right)\right]^{-1}, 1\right)$, then as $n \to \infty$*

$$\boldsymbol{\Sigma}(\boldsymbol{z})^{-1/2}\left[\overline{\boldsymbol{\beta}}(\boldsymbol{z}) - \boldsymbol{\beta}(\boldsymbol{z})\right] = \frac{s}{n}\sum_{i=1}^{n} h_i(\boldsymbol{z}) \xrightarrow{d} \mathsf{N}(0, \boldsymbol{I}_{d_X}), \qquad \boldsymbol{z} \in [0, 1]^{d_Z},$$

*where $\boldsymbol{\Sigma}(\boldsymbol{z}) := \boldsymbol{\Omega}^{-1}(\boldsymbol{z})\boldsymbol{\Lambda}(\boldsymbol{z})\boldsymbol{\Omega}^{-1}$ with*

$$\boldsymbol{\Lambda}(\boldsymbol{z}) := \frac{s^2}{n}\mathbb{E}\left[\epsilon^2 \theta(\boldsymbol{z}, \boldsymbol{Z})^2 \boldsymbol{X}\boldsymbol{X}^\mathsf{T}\right], \tag{10}$$

*and $\theta(\boldsymbol{z}, \boldsymbol{z}') = \mathbb{E}_\omega[S(\boldsymbol{z}, \boldsymbol{z}', \omega)]$ with $S(\boldsymbol{z}, \boldsymbol{z}', \omega) = |\{j \in \mathcal{A} : Z_j \in R(\boldsymbol{z}, \omega)\}|^{-1}$ if $\boldsymbol{z}' \in R(\boldsymbol{z}, \omega)$ and zero otherwise, for $\boldsymbol{z}, \boldsymbol{z}' \in [0, 1]^{d_Z}$. In particular*

$$\|\overline{\boldsymbol{\beta}}(\boldsymbol{z}) - \boldsymbol{\beta}(\boldsymbol{z})\| \lesssim_{\mathbb{P}} n^{-\frac{1}{2}\left[\beta(1-\eta)\left(\frac{\pi K(\alpha)}{d_Z} + \frac{1}{K(\alpha)}\right) - 1\right]}(\log n)^{d_Z/2} \to 0; \qquad \boldsymbol{z} \in [0, 1]^{d_Z}.$$

*Remark.* Converge of each tree is necessary for the random forest convergence. Hence, the requirement to set $k \asymp s^\eta$ according to Theorem 3. It is interesting to compare the rate of convergence of a (single) tree with $s = n$ with the rate of convergence of the forest with the restrictions to ensure asymptotic normality.

We propose to estimate the covariance matrix $\boldsymbol{\Sigma}(\boldsymbol{z})$ appearing in Theorem 3 using

$$\widehat{\boldsymbol{\Sigma}}(\boldsymbol{z}) := [\widehat{\boldsymbol{\Omega}}(\boldsymbol{z})]^{-1}\widehat{\boldsymbol{\Lambda}}(\boldsymbol{z})[\widehat{\boldsymbol{\Omega}}^{-1}(\boldsymbol{z})]^{-1}, \tag{11}$$

where $\widehat{\boldsymbol{\Omega}}(\boldsymbol{z})$ is given by (8). As for an estimator for $\widehat{\boldsymbol{\Lambda}}(\boldsymbol{z})$, define the residual function of the RF by $\boldsymbol{z} \mapsto \widehat{\epsilon}_i(\boldsymbol{z}) := Y_i - \boldsymbol{X}_i^\mathsf{T}\overline{\boldsymbol{\beta}}(\boldsymbol{z})$ for $i \in [n]$ and $\boldsymbol{z} \in [0,1]^{d_Z}$. Note any pair $\boldsymbol{z}, \boldsymbol{z}' \in [0,1]^{d_Z}$ we observe $B$ independent realizations of $S(\boldsymbol{z}, \boldsymbol{z}', \omega_b)$. So $\theta(\boldsymbol{z}, \boldsymbol{z}')$ can be unbiased estimated by $\widehat{\theta}(\boldsymbol{z}, \boldsymbol{z}') := \frac{1}{B}\sum_{b=1}^B S(\boldsymbol{z}, \boldsymbol{z}', \omega_b)$. So we propose to estimate $\widehat{\boldsymbol{\Lambda}}(\boldsymbol{z})$ using

$$\widehat{\boldsymbol{\Lambda}}(\boldsymbol{z}) := \frac{s^2}{n}\sum_{i=1}^n \widehat{\epsilon}_i(\boldsymbol{z})^2\widehat{\theta}(\boldsymbol{z}, \boldsymbol{Z}_i)^2 \boldsymbol{X}_i\boldsymbol{X}_i^\mathsf{T} = \frac{s^2}{n}\sum_{i=1}^n \widehat{\epsilon}_i(\boldsymbol{z})^2 \left[\frac{1}{B}\sum_{b=1}^B S(\boldsymbol{z}, \boldsymbol{Z}_i, \omega_b)\right]^2 \boldsymbol{X}_i\boldsymbol{X}_i^\mathsf{T}. \tag{12}$$

Even when an observation $\boldsymbol{Z}_i$ is not used to grow a tree $\omega_b$, we can compute $S(\boldsymbol{z}, \boldsymbol{Z}_i, \omega)$ by checking whether $\boldsymbol{Z}_i$ lies in $R(\boldsymbol{z}, \omega)$ and if so count how many observations are in $R(\boldsymbol{z}, \omega)$. Since we have $\mathbb{E}[\epsilon\theta(\boldsymbol{z}, \boldsymbol{Z})X] = 0$, we might consider the centered version of the estimator below given by

$$\widehat{\boldsymbol{\Lambda}}_c(\boldsymbol{z}) := \widehat{\boldsymbol{\Lambda}}(\boldsymbol{z}) - s^2. \tag{13}$$

## 3.2 Testing for Heterogeneity

### 3.2.1 Generalized Likelihood Ratio Test

Suppose we wish to test that the conditional expectation does not depend on $\boldsymbol{Z}$ (homogeneous partial effect). Specifically, we are interested in the parametric null.

$$\mathcal{H}_0 : \boldsymbol{\beta}(\boldsymbol{z}) = \boldsymbol{\beta}_0 \qquad \text{for some } \boldsymbol{\beta}_0 \in \mathbb{R}^{d_X} \text{ and all } \boldsymbol{z} \in [0,1]^d$$

against the non-parametric alternative hypothesis $\mathcal{H}_1 : \boldsymbol{\beta}(\boldsymbol{z})$ is Lipschitz. Following Fan et al. (2001); we propose to test $\mathcal{H}_0$ using the Generalized LRT, which is given as (after taking logs)

$$\Lambda(\mathcal{H}_0) := \frac{n}{2}\log\frac{\mathsf{RSS}_0}{\mathsf{RSS}},$$

where $\mathsf{RSS}_0 := \sum_{i=1}^n (Y_i - \boldsymbol{X}_i^\mathsf{T}\widetilde{\boldsymbol{\beta}})^2$, $\mathsf{RSS} := \sum_{i=1}^n \left[Y_i - \boldsymbol{X}_i^\mathsf{T}\overline{\boldsymbol{\beta}}(\boldsymbol{Z}_i)\right]^2$ and $\widetilde{\boldsymbol{\beta}}$ is the OLS estimator.

**Theorem 4** (Generalized LRT)

*Under the same condition of Theorem 3 and $\mathcal{H}_0$, if further $\mathbb{V}[Y|\boldsymbol{X}, \boldsymbol{Z}] = \sigma^2$ for some positive constant $\sigma^2$, then $\Lambda \overset{a}{\sim} \chi_\mu^2$ where $\chi_\mu^2$ is a chi-squared distribution with $\mu$ degrees of freedom, in the sense that,*

$$\nu^{-1}\left[\Lambda(\mathcal{H}_0) - \mu + o(\mu)\right] \overset{d}{\longrightarrow} \mathsf{N}(0,1)$$

*where*

$$\mu := d_X \left[ 2s\mathbb{E}[\theta(\mathbf{Z}, \mathbf{Z})] + \frac{s^2}{n}\mathbb{E}[\theta(\mathbf{Z}, \mathbf{Z})^2] + s^2\mathbb{E}[\theta(\mathbf{Z}, \mathbf{Z}')^2] \right]$$

$$\nu^2 := 4d_X \mathbb{E}\left[ \left( s\theta(\mathbf{Z}, \mathbf{Z}') + s^2\mathbb{E}[\theta(\mathbf{Z}, \mathbf{Z}')\theta(\mathbf{Z}, \mathbf{Z}'')] \right)^2 \right].$$

*and $\theta(\mathbf{z}, \mathbf{z}')$ is defined in Theorem 3; and $\mathbf{Z}'$ and $\mathbf{Z}''$ are independent copies of $\mathbf{Z}$.*

### 3.2.2 Langrage Multiplier Test

The idea is to explore the fact that for a Lagrange Multiplier (LM) type test, we are only required to estimate the model under the null and obtain a consistent estimator under both the null and alternative. When the null completely characterizes $\boldsymbol{\beta}(\mathbf{z})$ or when $\boldsymbol{\beta}(\mathbf{z})$ is known up to finite unknowns (parametric), we propose a somewhat canonical test.

Under $\mathcal{H}_0 : \boldsymbol{\beta}(\mathbf{z})$ is not a function of $z$, we have that $\mathbb{E}[\mathbf{Z}(Y - \mathbf{X}^\mathsf{T}\widetilde{\boldsymbol{\beta}}) = 0$ for some unknown $\widetilde{\boldsymbol{\beta}} \in \mathbb{R}^{d_X}$. Let $\widehat{\epsilon}_i^{OLS} := Y_i - \mathbf{X}_i^\mathsf{T}\widehat{\boldsymbol{\beta}}_{OLS}$ for $i \in [b]$ where $\widehat{\boldsymbol{\beta}}_{OLS}$ is the OLS estimator of $Y$ regressed onto $X$

So, we can use the sample moment condition below as a basis for constructing our test statistic.

$$M = \sum_{i=1}^{n} \widehat{\epsilon}_i^{OLS} \mathbf{Z}_i,$$

because under $\mathcal{H}_0$

$$M/\sqrt{n} = L\frac{1}{\sqrt{n}}\sum_{i=1}^{n} \epsilon_i W_i + o_\mathbb{P}(1); \qquad W_i := (\mathbf{Z}_i^\mathsf{T}, \mathbf{X}_i^\mathsf{T})^\mathsf{T}$$

where $L := (I_d : -\mathbb{E}[Z\mathbf{X}^T](\mathbb{E}[\mathbf{X}\mathbf{X}^T])^{-1})$ is an $(d_Z(d_Z + d_X))$ matrix. Since $\frac{1}{\sqrt{n}}\sum_{i=1}^{n}\epsilon_i W_i \xrightarrow{d} \mathsf{N}(0, \mathbb{E}[\epsilon^2 WW^\mathsf{T}])$, we have that $M/\sqrt{n} \xrightarrow{d} \mathsf{N}(0, L\mathbb{E}[\epsilon^2 WW^\mathsf{T}]L^\mathsf{T})$ therefore

$$M^\mathsf{T}(nV)^{-1}M \xrightarrow{d} \chi^2_{d_Z}; \qquad V := L\mathbb{E}[\epsilon^2 WW^\mathsf{T}]L^\mathsf{T}$$

Let $\widehat{\epsilon}_i^{RF} := Y_i - \mathbf{X}_i^\mathsf{T}\overline{\boldsymbol{\beta}}$ be the residuals of the RF. Under the alternative (and the null) $\overline{\boldsymbol{\beta}}(\mathbf{z}) \xrightarrow{\mathbb{P}} \boldsymbol{\beta}(\mathbf{z})$ then $\epsilon_i^{RF} \xrightarrow{\mathbb{P}} \epsilon_i$ for $i \in [n]$. Define the plug-in estimators

$$\widehat{L} := \left[ \boldsymbol{I}_d : -\sum_{i=1}^{n} Z_i \mathbf{X}_i^\mathsf{T} \left( \sum_{i=1}^{n} \mathbf{X}_i \mathbf{X}_i^\mathsf{T} \right)^{-1} \right], \qquad \widehat{V} := \widehat{L}\frac{1}{n}\sum_{i=1}^{n} \left[ (\widehat{\epsilon}_i^{RF})^2 W_i W_i^\mathsf{T} \right] \widehat{L}^\mathsf{T}$$

Hence, the test statistics become

$$T := M^\mathsf{T}(n\widehat{V})^{-1}M \xrightarrow{d} \chi^2_{d_Z}$$

as $n \to \infty$ under $\mathcal{H}_0$.

Note that the same process works to test $\mathcal{H}_0 : \boldsymbol{\beta}(\mathbf{z}) = \boldsymbol{\beta}_0(\mathbf{z})$ for some known function $\mathbf{z} \mapsto \boldsymbol{\beta}_0(\mathbf{z})$.

## 3.3 Extension to discrete controls

Partition the control variables as $\boldsymbol{Z} = (\boldsymbol{Z}', \boldsymbol{Z}'')$ where $\boldsymbol{Z}'' = (Z_1'', \dots Z_{d_{\boldsymbol{Z}''}}'')$ and $\boldsymbol{Z}_j''$ are discrete random variables with $m_j \geq 2$ categories for $j \in [d_{\boldsymbol{Z}''}]$. Without loss of generality we may assume that $Z_j''$ is supported on $\mathcal{S}_j = \{0, 1/(m_j - 1), 1/(m_j - 2), \dots, 1\}$ and thus $\boldsymbol{Z}''$ has support on $\mathcal{S} = \bigtimes_{j=1}^{d_{\boldsymbol{Z}''}} \mathcal{S}_j$.

Let $R = \bigtimes_{j=1}^{d_Z} [a_j, b_j]$ denote a rectangle in $[0,1]^{d_Z}$. For convenience, we can define the length of a rectangle along a discrete variable $\boldsymbol{Z}_j''$ as the fraction of the categories in the rectangle. Specifically $0 \leq \operatorname{diam}(R)_j := (|\mathcal{S}_j \cap [a_j, b_j]| - 1)/(m_j - 1) \leq 1$.

When a discrete variable $Z_j$ is chosen to split we replace (6) by the condition

$$\delta^* \in \underset{\delta \in \mathcal{S}_j}{\arg\min} \left[ \mathsf{RSS}(\{i \in \mathcal{I} : Z_{ij} = \delta)\}) + \mathsf{RSS}(\{i \in \mathcal{I} : Z_{ij} \neq \delta)\}) \right], \tag{14}$$

and identify the left and right child nodes by the index sets

$$\mathcal{I}^- := \{i \in \mathcal{I} : Z_{ij} = \delta^*\}, \quad \mathcal{I}^+ := \{i \in \mathcal{I} : Z_{ij} \neq \delta^*\}. \tag{15}$$

Fix a test point $\boldsymbol{z} = (\boldsymbol{z}', \boldsymbol{z}'') \in [0,1]^{d_{\boldsymbol{Z}'}} \times \mathcal{S}$ and let $M_j(\boldsymbol{z})$ denote the number of splits along the $j$-th discrete variable to form the (unique) leaf containing $z$. Since $s/k \to \infty$ we have that some large $s$

$$\mathbb{P}(M_j(\boldsymbol{x}) < m_j - 1) \leq \mathbb{P}\left( M_j(\boldsymbol{z}) \leq \frac{(1-\delta)\pi}{d_Z} \frac{\log(s/(2k-1))}{\log(1/(1-\alpha))} \right) \leq \left( \frac{s}{2k-1} \right)^{-\frac{\delta^2 \pi^2}{2 d_Z^2 \log(1/(1-\alpha))}}.$$

Let $\mathcal{E}_D(\boldsymbol{x}) = \bigcap_{j \in [d_Z'']} \{M_j(\boldsymbol{x}) = m_j - 1\}$, then by the union bound we have

$$\mathbb{P}(\mathcal{E}_D(\boldsymbol{x})) \gtrsim 1 - d_Z'' \left( \frac{s}{k} \right)^{-\frac{\delta^2 \pi^2}{2 d_Z^2 \log(1/(1-\alpha))}}$$

i.e., all discrete variables have a single class in the leaf $R(\boldsymbol{z}, \omega)$ with probability at least $1 - \left( \frac{s}{k} \right)^{\frac{-\delta^2 \pi^2}{2 d_Z^2 \log(1/(1-\alpha))}}$. Hence, the bias on the leaf $R(\boldsymbol{z}, \omega)$ can be upper bounded on $\mathcal{E}_D(\boldsymbol{x})$ as

$$\max_{(u', u'') \in R((z', z''), \omega)} \|\boldsymbol{\beta}(u', u'') - \boldsymbol{\beta}(z', z'')\| \leq \|\boldsymbol{\beta}(u', z'') - \boldsymbol{\beta}(z', z'')\| \leq C \operatorname{diam}(R(\boldsymbol{z}, \omega))$$

**Assumption 3**

*We consider the following conditions:*

    *(b) (Conditional Density) We can partition $\boldsymbol{Z} = (\boldsymbol{Z}', \boldsymbol{Z}'')$ such that $\boldsymbol{Z}'$ has a conditional density with respect to $\boldsymbol{Z}''$, which we denote by $f_{\boldsymbol{Z}'|\boldsymbol{Z}''}$, and $1/C \leq f_{\boldsymbol{Z}'|\boldsymbol{Z}''} \leq C$ for some constant $C > 0$. Finallly, $\mathbb{P}(\boldsymbol{Z}'' = z'') > 0$ for all $z'' \in \mathcal{S}$.*

*(c) (Smoothness) Let $z = (z', z'') \in [0,1]^{d_{Z_C}} \mathcal{S}$. We assume that $z' \mapsto \boldsymbol{\beta}(z', z'')$ $z' \mapsto \boldsymbol{\Omega}(z', z'')$,*

$z' \mapsto \mathbb{E}[\boldsymbol{X}(Y - \mathbb{E}[Y|\boldsymbol{X}, \boldsymbol{Z}])) | Z_C = \boldsymbol{z}', Z_D = \boldsymbol{z}'']$ *and* $\boldsymbol{z} \mapsto \mathbb{V}[\boldsymbol{X}(Y - \mathbb{E}[Y|\boldsymbol{X}, \boldsymbol{Z}] | Z_C = \boldsymbol{z}', Z_D =$

$\boldsymbol{z}'']$ *are Lipschitz-continuous for every* $z'' \in \mathcal{S}$, *and* $\boldsymbol{z} \mapsto \mathbb{V}[X_j X_k | \boldsymbol{Z} = \boldsymbol{z}]$ *and* $\boldsymbol{z} \mapsto \mathbb{V}[X_j(Y -$

$\mathbb{E}[Y|\boldsymbol{X}, \boldsymbol{Z}]|\boldsymbol{Z} = \boldsymbol{z}]$ *are bounded for* $j, k \in [d_Z]$;

**Theorem 5** (Discrete Control Variables)

*Under Assumptions 1 and 2 with 2(b) replaced by 3(b) and 2(c) replaced by 3(c), Theorems 1-3 hold.*

# 4 Monte Carlo Simulation

We conducted a simulation study to test the validity of our results in finite sample. We consider several specifications with different sample sizes. Specifically, we set the number of observations to be $n \in \{250, 500, 1000\}$. The dimension $d_Z$ of the variables determining model heterogeneity is $d_Z = \{1, 2, 3, 5\}$. Each simulated model is estimated by the Random Forest as per the Algorithm 1 with $B = 3000$ subsampling replications. We consider a fraction of $s = 0.8$ observations in each subsample, $k = s^{1/6}$ and $\alpha = 0.005$. The number of Monte Carlo replications is 500.

## 4.1 Model Estimation

The simulated models are determined by the general equation

$$Y_i = \beta_0(\boldsymbol{Z}_i) + \beta_1(\boldsymbol{Z}_i)X_i + U_i, \tag{16}$$

where $U_i$ is an independent and normally distributed random variable with zero mean and standard deviation equal to 0.5. The scalar continuous treatment variable is $X_i \sim \mathsf{Uniform}(0, 1)$. $\boldsymbol{Z}_i$ is a random vector of $d_Z$ mutually independent random variables uniformly distributed between 0 and 1.

We set $\beta_0(Z_i) = 0$ for all $i = 1, \ldots, n$, and consider the following specifications for the slope parameter $\beta_1(\boldsymbol{Z}_i)$:

- Model I: $p = 1$

$$\beta_1(\boldsymbol{Z}_i) = 5 \left[ 1 + \frac{1}{1 + e^{-100(Z_i - 0.3)}} - \frac{1}{1 + e^{-100(Z_i - 0.7)}} \right].$$

Figure 1(a) illustrates the evolution of $\beta_1(Z_i)$ as a function of $Z_i$. As can be observed from the plot, the slope coefficient smoothly changes between zero and one according to the source of heterogeneity $Z_i$.

- Model II: $p = 2$

$$\beta_1(\boldsymbol{Z}_i) = 2f_{i,1}f_{i,2} - 2f_{i,1}(1 - f_{i,2}),$$

where

$$f_{i,1} = 1 + \frac{1}{1 + e^{-100(Z_{1,i}-0.3)}} - \frac{1}{1 + e^{-100(Z_{1,i}-0.7)}} \quad \text{and}$$
$$f_{i,2} = 1 + \frac{1}{1 + e^{-100(Z_{2,i}-0.3)}} - \frac{1}{1 + e^{-100(Z_{2,i}-0.7)}}$$

Figure 1(b) illustrates the evolution of $\beta_1(\boldsymbol{Z}_i)$ as a function of $\boldsymbol{Z}_i$. As in the case of the first simulated model, the slope coefficient smoothly changes between different regimes.

- Model III: $p = 3$

$$\beta_1(\boldsymbol{Z}_i) = 2f_{i,1}f_{i,2}f_{i,3} - f_{i,1}f_{i,2}(1 - f_{i,3}) - 1.5f_{i,1}(1 - f_{i,2}),$$

where $f_{i,j} = \frac{1}{1+e^{-100(Z_{j,i}-0.5)}}$, for $j = 1, 2, 3$ and $i = 1, \ldots, n$.

The resulting model is a product of logistic functions. Figure 1(c) illustrates the heterogeneity in the slope coefficient as a function of $Z_1$ and $Z_2$ holding $Z_3$ fixed.

- Model IV: $p = 5$

$$\beta_1(\boldsymbol{Z}_i) = f_{i,1}f_{i,2}f_{i,3} - f_{i,1}f_{i,2}(1 - f_{i,3}) - 1.5 * f_{i,1}(1 - f_{i,2}) + 1.5(1 - f_{i,1})f_{i,4}$$
$$- 0.8(1 - f_{i,1})(1 - f_{i,4})f_{i,5} + 0.7(1 - f_{i,1})(1 - f_{i,4})(1 - f_{i,5})$$

where $f_{i,j} = \frac{1}{1+e^{-100(Z_{j,i}-0.5)}}$, for $j = 1, \ldots, 5$ and $i = 1, \ldots, n$.

As in the previous models, Model IV also implies a heterogeneity pattern where the slope coefficient smoothly changes among six limiting regimes.

The results are reported in Figure 2 and Tables 1–5. Figure 2 reports results for Model I where $p = 1$. It illustrates the median slope estimation across the Monte Carlo simulations, as well as the 95% confidence bands. It is evident the estimation is precise in this case.

Table 1 presents the averages derived from the Monte Carlo simulations for various descriptive statistics pertinent to goodness-of-fit assessments. The evaluated statistics include the mean, standard deviation, kurtosis, and skewness. For Model I, Panel (I.a) examines the estimated residuals from the model fit, Panel (I.b) addresses the estimation error associated with the varying

intercept of the model, and Panel (I.c) pertains to the estimation error for the varying slope coefficient. The subsequent panels replicate these findings for Models II, III, and IV.

Several facts emerge from the table. First, the model approximation is satisfactory even in samples as small as 250 observations. Notably, the average of the residual standard deviation is close to 0.5, which is the true value. Furthermore, the average estimated kurtosis is approximately three, and the average estimated skewness is close to zero, indicating that the residuals are approximately normally distributed. Additionally, the results suggest that the estimation of the varying intercept is more precise than the slope estimation, as the average standard deviation of the errors is significantly smaller for the former than for the latter. Finally, as expected, the performance of the method slightly deteriorates as the dimension of $\boldsymbol{Z}$ increases.

Tables 2 and 3 present results for several test values pertinent to the vector $\boldsymbol{Z}$. We analyze 11 points where $\boldsymbol{Z}$ constitutes a 51 that spans from a vector of zeros to a vector of ones, with 0.1 increments. For simplicity and without loss of generality, we assume that all elements of the test points are equal. Table 2 reports the average bias and the mean squared error (MSE) for the varying intercept coefficient evaluated at each test point for different sample sizes. It is evident that the intercept estimation is precise and improves with increasing sample size. Table 3 reports the average bias and the mean squared error (MSE) for the varying slope coefficient evaluated at each test point for different sample sizes. The results in the table corroborate our previous conclusion that as the sample size increases, the method's performance improves. Finally, there is evidence that the function approximation is better for points closer to the center of the distribution of $\boldsymbol{Z}$.

## 4.2 Inference

We present coverage results in Tables 4 and 5. Specifically, we provide the average 90% and 95% coverage across Monte Carlo simulations for both the intercept and the slope parameters. Table 4 details the coverage for the intercept, while 5 presents the corresponding results for the slope parameter. It is evident from the tables that the coverage probabilities for the intercept are close to the expected values, even with small sample sizes and an increased number of variables. This finding supports our previous simulation results, indicating minimal bias in intercept estimation. Conversely, the coverage for the slope parameter is significantly underestimated, particularly for values that lie farther from the center of the covariate distribution. This observation aligns with the larger biases reported in Table 3. It is important to mention that such narrow coverage probabilities have also been reported in the simulations in Friedberg et al. (2021).

## 4.3 Linarity Testing

Finally, to evaluate the finite sample performance of the Lagrange Multiplier homogeneity (linearity) test described in Section 3.2.2. We simulated linear (homogeneous) models with $p \in \{1, 2, 3, 5\}$ covariates. The results are reported in Figure 3, which illustrates the size discrepancy relative to the nominal size. As observed, the size distortions remain negligible.

## 5 Empirical Illustrations

In this section, we illustrate our proposed methodology utilizing two distinct datasets. The first is a synthetic dataset generated by a widely used model specifically designed to evaluate nonparametric methods. The second is a real dataset concerning the economic convergence of Brazilian municipalities.

### 5.1 Simulated data

The first dataset consists of observations generated according to the following model:

$$Y_i = \beta(\boldsymbol{Z}_i)X_i + U_i, \quad i = 1, \ldots, n \tag{17}$$

$$\beta(\boldsymbol{Z}_i) = 10\sin(\pi Z_{1,i}Z_{2i}) + 20(Z_{3,i} - 0.5)^2 + 10Z_{4,i} + 5Z_{5,i}, \tag{18}$$

where $\boldsymbol{Z}_i$ is a vector of mutually independent uniform random variables taking values on $[0, 1]^5$, and $U_i$ is a zero-mean normally distributed random variable with unit variance. The model for $\beta(\boldsymbol{Z}_i)$ has been widely employed in the literature to evaluate semi-parametric models. In this context, heterogeneity is jointly influenced by interactions, quadratic forms, and a robust linear signal. Figure 4 illustrates the three sources of heterogeneity. See, for example, Friedberg et al. (2021) for a similar data-generating process.

To evaluate the performance of the model presented in this paper, we consider various sample sizes (2000, 4000, 8000, 16000, 32000). Figure 5 illustrates the empirical distribution of $\boldsymbol{\beta}(\boldsymbol{Z}_i)$, $i = 1, \ldots, 32000$. We estimate the locally linear random forest model using the same hyperparameter settings as those utilized in the Monte Carlo simulation.

Figure 6 shows results concerning the estimation parameters. Panels (a) and (b) illustrate the scatter plot of the fitted $\boldsymbol{\beta}(\boldsymbol{Z})$ against the true values for $n = 2000$ and $n = 36000$, respectively. A linear regression line is also included. Figure 7 reports the evolution of the bias and the MSE as a function of the sample size. As we can observe from both figures, the estimation improves as the

sample size increases, which aligns with established statistical theory. However, as anticipated by our theoretical results, the convergence to the true values is notably slow.

## 5.2 Economic growth and convergence among Brazilian municipalities

We illustrate our methodology by testing heterogeneity in the convergence among Brazilian municipalities between 1970 and 2000. Our starting point is the simplified convergence equation presented in Barro and Sala-i-Martin (1992):

$$\log\left(\frac{Y_{i,t}}{Y_{i,t-1}}\right) = a_i + \gamma_i \log\left(Y_{i,t-1}\right) + \phi_i\left(t-1\right) + U_{i,t}, \tag{19}$$

where $Y_{i,t}$ is the per capita income of region $i$ in period $t$, $a_i$ is associated with the steady-state level of per capita income and the rate of technological progress, $\phi_i$ is a parameter related to the time trend determined by the technological progress, and $U_{i,t}$ is the random term. Convergence corresponds to the parameter $\gamma_i$.

From a conceptual point of view, two alternative assumptions determine the most important distinction of convergence concepts. First, we can assume that $a_i = a$, $\gamma_i = \gamma$, and $\phi_i = \phi$, i.e., that the basic parameters of preference and technology are the same for all economies represented in the sample. This is when $\gamma < 0$ represents *unconditional convergence* - a situation where poorer regions tend unconditionally to grow more quickly than richer ones. Alternatively, we can state a weaker assumption, allowing for possible differences in the steady state across the economies considered and heterogeneity in the convergence rate. In terms of equation (19), $a_i$, $\gamma_i$ and $\phi_i$ are allowed to vary across different regions.

We estimate (19) in a cross-section setup, where there is no identifiable time trend, and we are not able to distinguish between $\phi_i$ and $a_i$. Thus, we estimate the following equation:

$$\log\left(\frac{Y_{i,2000}}{Y_{i,1970}}\right) = \alpha_i + \gamma_i \log\left(Y_{i,1970}\right) + U_{i,2000}. \tag{20}$$

The data originate from the Brazilian Demographic Censuses conducted in 1970 and 2000. The geographical units were adjusted to account for the reorganization of Brazilian territory throughout this time frame. In 1970, Brazil was composed of 3,951 municipalities. By 2000, the count had increased to 5,507 municipalities. Consequently, all data collected in 2000 were aggregated to align with the municipal structure as it existed in 1970. Our dependent variable is defined as the average growth in per capita income from 1970 to 2000 for each municipality. The independent variable utilized in this analysis is the logarithm of the per capita income level recorded in 1970.

Figure 8 illustrates the differences across municipalities in terms of growth rate. Figure 8 shows that the variations in growth do not coincide with the administrative state frontiers. There are substantial variations within many of the Brazilian states.

We employ the semi-parametric approach outlined in the previous section to assess conditional convergence. Variations in preferences and technological parameters are endogenously incorporated into the analysis through geographical proximities. The underlying assumption suggests that cities situated in close proximity experience similar steady states. Within our modeling framework, we estimate (20), recognizing that $\alpha_i$ represents a semi-parametric function of the latitude and longitude coordinates. The results are displayed in Figures 9 and 10, which report the heterogeneity patterns in the intercept and the convergence parameter. The estimated geographic heterogeneity varies remarkably, with notable clusters of high-income municipalities in the central and southern parts of the country. Significant geographic differences across municipalities do not coincide with the geographic structure of the Brazilian states.

# 6    Conclusions

This paper presents a semi-parametric framework for estimating heterogeneous partial effects, combining the flexibility of machine learning techniques with the interpretability of traditional parametric models. By utilizing a Random Forest-based methodology, we provide a robust and adaptable approach to capturing complex heterogeneity in the relationship between explanatory variables and outcomes.

Our theoretical contributions establish key consistency and asymptotic normality results, ensuring the reliability of our estimator. Importantly, our framework accommodates both continuous and discrete covariates while maintaining desirable statistical properties. The Monte Carlo simulations demonstrate the method's accuracy, even in moderate sample sizes, and highlight the precision of our approach in recovering varying intercepts and slopes. Additionally, our empirical analysis of Brazilian municipal economic convergence underscores the practical relevance of our method, revealing substantial geographic heterogeneity in growth dynamics.

Future research can extend this methodology to settings with dependent data and high-dimensional covariates, broadening its applicability to dynamic panel models and network data. The current approach can accommodate high-dimensional controls ($\boldsymbol{Z}$) under additional conditions on the underlying data-generating process. Specifically, it would require that (i) only a small

number of controls are relevant to the model (sparsity assumption) and (ii) the relevant controls are independent of the irrelevant controls and the outcome variable. The latter is a strong assumption in most applications, so we did not pursue this route here.

Regarding the dependent data across observations, the subsampling and sample split steps must be adjusted to preserve both the data dependence structure and the independence of the two samples. A block-subsampling technique is a natural candidate; however, implementing this would significantly alter the proof techniques for the subsequent steps and might obscure the main idea of the paper. Furthermore, refinements in inference procedures for heterogeneous effects, including the construction of a uniformly valid confidence interval, would enhance the model's applicability in applied research. Additionally, considering Lemma 1(d), a significant modification to the algorithm would be necessary to achieve uniform convergence of the proposed estimators.

Overall, this paper offers a flexible, interpretable, and computationally efficient tool for studying heterogeneity in the partial effect of a variable of interest, bridging the gap between parametric and nonparametric estimation.

## Table 1: Godness-of-Fit

This table reports the average across the Monte Carlo simulations for several descriptive statistics related to goodness-of-fit. The statistics are the mean, the standard deviation, the kurtosis, and skewness. Panel (I.a) considers the estimated residuals from the model fit. Panel (I.b) considers the estimation error for the varying intercept of the model. Panel (I.c) is related to the estimation error for the varying slope coefficient. The remaining panels show the same results for models II, III, and IV.

|  | Panel (I.a): $p=1$ error term Sample Size | | | Panel (II.a): $p=2$ error term Sample Size | | | Panel (III.a): $p=3$ error term Sample Size | | | Panel (IV.a): $p=5$ error term Sample Size | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | 250 | 500 | 1000 | 250 | 500 | 1000 | 250 | 500 | 1000 | 250 | 500 | 1000 |
| mean | -0.0019 | -0.0002 | 0.0001 | -0.0087 | -0.0069 | -0.0047 | 0.0024 | 0.0014 | 0.0009 | 0.0004 | 0.0002 | 0.0001 |
| standard deviation | 0.4506 | 0.4458 | 0.4568 | 0.6207 | 0.5058 | 0.4807 | 0.4750 | 0.4533 | 0.4571 | 0.5638 | 0.5235 | 0.5154 |
| kurtosis | 3.0154 | 3.0182 | 3.0346 | 3.9624 | 3.3529 | 3.0946 | 3.0793 | 3.1119 | 3.0315 | 2.8920 | 2.9941 | 2.9761 |
| skewness | 0.0128 | 0.0015 | 0.0027 | 0.5457 | 0.2077 | 0.0683 | 0.0710 | 0.0337 | 0.0182 | 0.0102 | 0.0095 | 0.0126 |

|  | Panel (I.b): $p=1$ intercept Sample Size | | | Panel (II.b): $p=2$ intercept Sample Size | | | Panel (III.b): $p=3$ intercept Sample Size | | | Panel (IV.b): $p=5$ intercept Sample Size | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | 250 | 500 | 1000 | 250 | 500 | 1000 | 250 | 500 | 1000 | 250 | 500 | 1000 |
| mean | -0.0077 | -0.0004 | -0.0019 | 0.0044 | -0.0037 | 0.0028 | -0.0049 | 0.0010 | 0.0001 | 0.0023 | -0.0024 | -0.0044 |
| standard deviation | 0.2396 | 0.2383 | 0.1967 | 0.2483 | 0.2179 | 0.1642 | 0.1609 | 0.1544 | 0.1256 | 0.1540 | 0.1460 | 0.1089 |
| kurtosis | 3.0068 | 3.2621 | 3.0542 | 3.3724 | 3.4301 | 3.0503 | 2.9366 | 3.1829 | 3.0051 | 3.0813 | 3.3344 | 3.0628 |
| skewness | -0.0370 | 0.0671 | -0.0195 | 0.0905 | -0.0007 | 0.0072 | -0.0084 | 0.0119 | 0.0058 | -0.0173 | -0.0068 | -0.0001 |

|  | Panel (I.c): $p=1$ slope Sample Size | | | Panel (II.c): $p=2$ slope Sample Size | | | Panel (III.c): $p=3$ slope Sample Size | | | Panel (IV.c): $p=5$ slope Sample Size | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | 250 | 500 | 1000 | 250 | 500 | 1000 | 250 | 500 | 1000 | 250 | 500 | 1000 |
| mean | 0.0137 | -0.0013 | 0.0021 | -0.0237 | -0.0061 | -0.0101 | 0.0168 | 0.0003 | 0.0034 | -0.0020 | 0.0000 | 0.0126 |
| standard deviation | 0.4567 | 0.4239 | 0.3499 | 0.9791 | 0.6639 | 0.4672 | 0.4841 | 0.3962 | 0.3019 | 0.7661 | 0.6474 | 0.5455 |
| kurtosis | 3.7132 | 3.3192 | 3.0741 | 4.0835 | 4.1379 | 3.8797 | 3.2178 | 3.3337 | 3.3105 | 1.8627 | 2.0687 | 1.9109 |
| skewness | - -0.0132 | -0.0308 | -0.0121 | 0.8756 | 0.5659 | 0.3234 | 0.2993 | 0.2660 | 0.2157 | 0.0338 | 0.0724 | 0.0331 |

## Table 2: Test Points (Goodness-of-Fit – Intercept)

The table shows the average bias and mean squared error for the intercept estimation at different test points (MSE). The first column shows the points considered. For the case where $p > 1$, all the values of $\boldsymbol{Z}$ are equal to the number indicated in the first column. Panel (I) considers the case of Model I ($p = 1$). Panel (II) considers the case of Model II ($p = 2$). Panel (III) considers the case of Model III ($p = 3$). Panel (IV) considers the case of Model IV ($p = 5$).

| | Panel I: $p = 1$ | | | | | | Panel II: $p = 2$ | | | | | |
| | (a): Bias | | | (b): MSE | | | (a): Bias | | | (b): MSE | | |
| | Sample Size | | | Sample Size | | | Sample Size | | | Sample Size | | |
| Point | 250 | 500 | 1000 | 250 | 500 | 1000 | 250 | 500 | 1000 | 250 | 500 | 1000 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.0561 | 0.0062 | -0.0111 | 0.0927 | 0.0702 | 0.0563 | 0.0028 | 0.0016 | 0.0051 | 0.0766 | 0.0644 | 0.0321 |
| 0.1 | 0.0024 | 0.0102 | -0.0312 | 0.0563 | 0.0510 | 0.0424 | 0.0037 | -0.0014 | 0.0049 | 0.0531 | 0.0462 | 0.0245 |
| 0.2 | 0.0081 | 0.0139 | 0.0044 | 0.0533 | 0.0566 | 0.0407 | -0.0080 | -0.0053 | 0.0269 | 0.0646 | 0.0416 | 0.0291 |
| 0.3 | -0.0196 | -0.0109 | 0.0289 | 0.1139 | 0.0811 | 0.0441 | 0.0161 | -0.0263 | 0.0114 | 0.1096 | 0.0701 | 0.0320 |
| 0.4 | -0.0004 | -0.0257 | 0.0134 | 0.0574 | 0.0503 | 0.0348 | 0.0035 | 0.0017 | 0.0075 | 0.1097 | 0.0543 | 0.0243 |
| 0.5 | 0.0001 | -0.0113 | -0.0008 | 0.0487 | 0.0590 | 0.0373 | 0.0145 | -0.0040 | 0.0276 | 0.0722 | 0.0404 | 0.0219 |
| 0.6 | -0.0090 | 0.0064 | 0.0146 | 0.0645 | 0.0529 | 0.0425 | -0.0115 | -0.0137 | -0.0057 | 0.1027 | 0.0626 | 0.0250 |
| 0.7 | -0.0143 | -0.0076 | -0.0235 | 0.1413 | 0.0598 | 0.0447 | -0.0252 | 0.0074 | 0.0050 | 0.0917 | 0.0696 | 0.0254 |
| 0.8 | -0.0086 | -0.0264 | -0.0281 | 0.0644 | 0.0623 | 0.0399 | 0.0177 | -0.0182 | 0.0001 | 0.0578 | 0.0334 | 0.0255 |
| 0.9 | -0.0313 | 0.0170 | 0.0023 | 0.0650 | 0.0493 | 0.0505 | 0.0205 | -0.0040 | 0.0173 | 0.0539 | 0.0380 | 0.0232 |
| 1 | -0.0393 | 0.0144 | -0.0029 | 0.0761 | 0.0779 | 0.0569 | 0.0426 | -0.0113 | 0.0159 | 0.0600 | 0.0567 | 0.0311 |

| | Panel III: $p = 3$ | | | | | | Panel IV: $p = 5$ | | | | | |
| | (a): Bias | | | (b): MSE | | | (a): Bias | | | (b): MSE | | |
| | Sample Size | | | Sample Size | | | Sample Size | | | Sample Size | | |
| Point | 250 | 500 | 1000 | 250 | 500 | 1000 | 250 | 500 | 1000 | 250 | 500 | 1000 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | -0.0195 | -0.0147 | -0.0041 | 0.0282 | 0.0269 | 0.0160 | -0.0064 | 0.0031 | -0.0153 | 0.0344 | 0.0238 | 0.0123 |
| 0.1 | -0.0159 | -0.0130 | 0.0034 | 0.0249 | 0.0216 | 0.0157 | -0.0096 | 0.0060 | -0.0123 | 0.0323 | 0.0245 | 0.0128 |
| 0.2 | -0.0035 | -0.0044 | -0.0055 | 0.0222 | 0.0223 | 0.0127 | -0.0041 | 0.0064 | -0.0038 | 0.0264 | 0.0201 | 0.0095 |
| 0.3 | 0.0066 | 0.0028 | -0.0015 | 0.0209 | 0.0251 | 0.0138 | -0.0060 | 0.0046 | 0.0024 | 0.0198 | 0.0164 | 0.0087 |
| 0.4 | - 0.0110 | 0.0109 | 0.0026 | 0.0189 | 0.0155 | 0.0119 | -0.0018 | 0.0073 | -0.0055 | 0.0185 | 0.0123 | 0.0069 |
| 0.5 | 0.0172 | 0.0062 | -0.0023 | 0.0194 | 0.0123 | 0.0083 | 0.0146 | -0.0033 | -0.0092 | 0.0211 | 0.0095 | 0.0048 |
| 0.6 | -0.0011 | -0.0174 | 0.0009 | 0.0254 | 0.0180 | 0.0140 | 0.0215 | 0.0076 | 0.0033 | 0.0189 | 0.0141 | 0.0087 |
| 0.7 | -0.0233 | -0.0113 | -0.0103 | 0.0265 | 0.0171 | 0.0156 | 0.0164 | 0.0051 | 0.0152 | 0.0192 | 0.0136 | 0.0099 |
| 0.8 | -0.0341 | -0.0205 | -0.0124 | 0.0316 | 0.0221 | 0.0163 | 0.0035 | -0.0195 | 0.0062 | 0.0237 | 0.1428 | 0.0104 |
| 0.9 | -0.0216 | -0.0147 | -0.0046 | 0.0409 | 0.0254 | 0.0204 | 0.0115 | -0.0038 | 0.0027 | 0.0255 | 0.0206 | 0.0112 |
| 1 | -0.0178 | -0.0125 | -0.0067 | 0.0380 | 0.0342 | 0.0199 | 0.0072 | -0.0036 | 0.0003 | 0.0268 | 0.0199 | 0.0110 |

## Table 3: Test Points (Goodness-of-Fit – Slope)

The table shows the average bias and mean squared error for the intercept estimation at different test points (MSE). The first column shows the points considered. For the case where $p > 1$, all the values of $\boldsymbol{Z}$ are equal to the number indicated in the first column. Panel (I) considers the case of Model I ($p = 1$). Panel (II) considers the case of Model II ($p = 2$). Panel (III) considers the case of Model III ($p = 3$). Panel (IV) considers the case of Model IV ($p = 5$).

| | Panel I: $p = 1$ | | | | | | Panel II: $p = 2$ | | | | | |
| | (a): Bias | | | (b): MSE | | | (a): Bias | | | (b): MSE | | |
| | Sample Size | | | Sample Size | | | Sample Size | | | Sample Size | | |
| Point | 250 | 500 | 1000 | 250 | 500 | 1000 | 250 | 500 | 1000 | 250 | 500 | 1000 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | -0.0693 | 0.0344 | 0.0421 | 0.2735 | 0.2201 | 0.1663 | -0.6764 | -0.4972 | -0.4115 | 0.2912 | 0.2247 | 0.1126 |
| 0.1 | -0.0095 | -0.0085 | 0.0464 | 0.1532 | 0.1674 | 0.1096 | -0.5533 | -0.2788 | -0.1789 | 0.1990 | 0.1399 | 0.0800 |
| 0.2 | -0.0310 | -0.0413 | 0.0198 | 0.1829 | 0.1606 | 0.1260 | -0.7367 | -0.3826 | -0.2830 | 0.2374 | 0.1470 | 0.0787 |
| 0.3 | -0.4974 | -0.2378 | -0.1641 | 0.5734 | 0.3932 | 0.2077 | -0.8960 | -0.9171 | -1.0405 | 0.9729 | 0.7729 | 0.5702 |
| 0.4 | 0.0131 | 0.0427 | -0.0127 | 0.1655 | 0.1768 | 0.0957 | 1.8125 | 0.9784 | 0.5650 | 0.5971 | 0.2518 | 0.0853 |
| 0.5 | 0.0092 | 0.0211 | 0.0164 | 0.1530 | 0.1779 | 0.1128 | 1.3114 | 0.6609 | 0.3321 | 0.3396 | 0.1632 | 0.0669 |
| 0.6 | 0.0331 | -0.0032 | -0.0280 | 0.1818 | 0.1641 | 0.1041 | 1.8445 | 0.9141 | 0.5676 | 0.6678 | 0.2635 | 0.0929 |
| 0.7 | 0.2560 | 0.0292 | -0.0460 | 0.8037 | 0.3561 | 0.2422 | 0.7413 | 0.3897 | -0.0711 | 0.7615 | 0.8367 | 0.5368 |
| 0.8 | -0.0146 | 0.0259 | 0.0563 | 0.1765 | 0.1954 | 0.1028 | -0.8750 | -0.4580 | -0.2882 | 0.2634 | 0.1475 | 0.0786 |
| 0.9 | 0.0513 | -0.0293 | -0.0013 | 0.1968 | 0.1386 | 0.1509 | -0.6706 | -0.3357 | -0.2046 | 0.2004 | 0.1341 | 0.0766 |
| 1 | 0.0660 | -0.0088 | 0.0298 | 0.2547 | 0.2396 | 0.1601 | -0.7933 | -0.5480 | -0.4556 | 0.2332 | 0.1989 | 0.1218 |

| | Panel III: $p = 3$ | | | | | | Panel IV: $p = 5$ | | | | | |
| | (a): Bias | | | (b): MSE | | | (a): Bias | | | (b): MSE | | |
| | Sample Size | | | Sample Size | | | Sample Size | | | Sample Size | | |
| Point | 250 | 500 | 1000 | 250 | 500 | 1000 | 250 | 500 | 1000 | 250 | 500 | 1000 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.2818 | 0.1940 | 0.1568 | 0.0847 | 0.0828 | 0.0457 | 0.5893 | 0.4986 | 0.5042 | 0.1206 | 0.0853 | 0.0447 |
| 0.1 | 0.2228 | 0.1240 | 0.0670 | 0.0819 | 0.0647 | 0.0477 | 0.5565 | 0.4389 | 0.4202 | 0.1170 | 0.0895 | 0.0423 |
| 0.2 | 0.1771 | 0.0910 | 0.0604 | 0.0727 | 0.0639 | 0.0361 | 0.5083 | 0.3906 | 0.3521 | 0.1004 | 0.0705 | 0.0331 |
| 0.3 | 0.1643 | 0.1040 | 0.0662 | 0.0647 | 0.0769 | 0.0428 | 0.5059 | 0.4038 | 0.3476 | 0.0744 | 0.0585 | 0.0321 |
| 0.4 | - 0.2247 | 0.1486 | 0.1013 | 0.0607 | 0.0468 | 0.0331 | 0.5655 | 0.4792 | 0.4418 | 0.0767 | 0.0535 | 0.0317 |
| 0.5 | -0.0243 | -0.0329 | 0.0215 | 0.0703 | 0.0607 | 0.0491 | 0.0192 | 0.0441 | 0.0392 | 0.0905 | 0.0815 | 0.0586 |
| 0.6 | 1.0165 | 0.8154 | 0.5452 | 0.1167 | 0.0768 | 0.0478 | 0.9135 | 0.8141 | 0.6863 | 0.0959 | 0.0683 | 0.0428 |
| 0.7 | 0.7940 | 0.5273 | 0.3219 | 0.1142 | 0.0590 | 0.0499 | 0.8217 | 0.6647 | 0.5278 | 0.0884 | 0.0575 | 0.0414 |
| 0.8 | 0.7404 | 0.4528 | 0.2789 | 0.1047 | 0.0748 | 0.0483 | 0.8040 | 0.6659 | 0.5314 | 0.1012 | 0.2341 | 0.0425 |
| 0.9 | 0.7893 | 0.5353 | 0.3340 | 0.1477 | 0.0838 | 0.0621 | 0.8085 | 0.6992 | 0.6035 | 0.1139 | 0.0817 | 0.0469 |
| 1 | 0.8816 | 0.7217 | 0.5808 | 0.1392 | 0.1114 | 0.0648 | 0.8349 | 0.7665 | 0.7086 | 0.1165 | 0.0792 | 0.0465 |

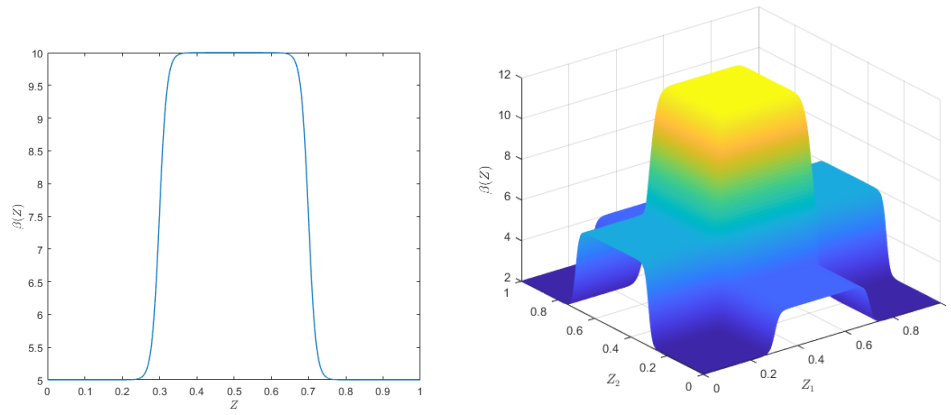## Table 4: Test Points (Coverage – Intercept)

The table shows the average coverage for the intercept estimation at different test points (MSE). The first column shows the points considered. For the case where $p > 1$, all the values of $\boldsymbol{Z}$ are equal to the number indicated in the first column. Panel (I) considers the case of Model I ($p = 1$). Panel (II) considers the case of Model II ($p = 2$). Panel (III) considers the case of Model III ($p = 3$). Panel (IV) considers the case of Model IV ($p = 5$).

|  | Panel I: $p = 1$ | | | | | | Panel II: $p = 2$ | | | | | |
|  | (a): 90% | | | (b): 95% | | | (a): 90% | | | (b): 95% | | |
|  | Sample Size | | | Sample Size | | | Sample Size | | | Sample Size | | |
| Point | 250 | 500 | 1000 | 250 | 500 | 1000 | 250 | 500 | 1000 | 250 | 500 | 1000 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.9000 | 0.8900 | 0.8750 | 0.9400 | 0.9650 | 0.9300 | 0.9050 | 0.9100 | 0.9000 | 0.9350 | 0.9500 | 0.9650 |
| 0.1 | 0.9000 | 0.9150 | 0.8800 | 0.9500 | 0.9500 | 0.9300 | 0.8750 | 0.9200 | 0.8900 | 0.9450 | 0.9550 | 0.9250 |
| 0.2 | 0.9150 | 0.8750 | 0.9050 | 0.9650 | 0.9500 | 0.9500 | 0.9100 | 0.9050 | 0.9000 | 0.9550 | 0.9450 | 0.9500 |
| 0.3 | 0.8950 | 0.9000 | 0.9000 | 0.9450 | 0.9350 | 0.9550 | 0.8950 | 0.8900 | 0.8950 | 0.9450 | 0.9450 | 0.9600 |
| 0.4 | 0.9250 | 0.8850 | 0.9000 | 0.9700 | 0.9800 | 0.9500 | 0.9350 | 0.8950 | 0.9300 | 0.9450 | 0.9550 | 0.9650 |
| 0.5 | 0.9150 | 0.9000 | 0.8950 | 0.9600 | 0.9600 | 0.9400 | 0.9050 | 0.9050 | 0.8850 | 0.9500 | 0.9500 | 0.9450 |
| 0.6 | 0.8850 | 0.8900 | 0.8850 | 0.9650 | 0.9500 | 0.9700 | 0.9000 | 0.9100 | 0.8850 | 0.9550 | 0.9500 | 0.9450 |
| 0.7 | 0.9350 | 0.9100 | 0.9200 | 0.9600 | 0.9350 | 0.9650 | 0.8950 | 0.8700 | 0.9100 | 0.9350 | 0.9400 | 0.9550 |
| 0.8 | 0.9050 | 0.9100 | 0.8950 | 0.9350 | 0.9450 | 0.9450 | 0.8850 | 0.8850 | 0.9100 | 0.9500 | 0.9500 | 0.9550 |
| 0.9 | 0.8950 | 0.8900 | 0.8850 | 0.9350 | 0.9500 | 0.9500 | 0.8900 | 0.9050 | 0.8950 | 0.9500 | 0.9600 | 0.9600 |
| 1 | 0.8950 | 0.9200 | 0.9000 | 0.9400 | 0.9550 | 0.9550 | 0.9000 | 0.9100 | 0.8900 | 0.9350 | 0.9550 | 0.9550 |

|  | Panel III: $p = 3$ | | | | | | Panel IV: $p = 5$ | | | | | |
|  | (a): 90% | | | (b): 95% | | | (a): 90% | | | (b): 95% | | |
|  | Sample Size | | | Sample Size | | | Sample Size | | | Sample Size | | |
| Point | 250 | 500 | 1000 | 250 | 500 | 1000 | 250 | 500 | 1000 | 250 | 500 | 1000 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.9100 | 0.8800 | 0.9000 | 0.9300 | 0.9250 | 0.9500 | 0.9000 | 0.8750 | 0.9150 | 0.9300 | 0.9400 | 0.9400 |
| 0.1 | 0.9150 | 0.8950 | 0.8900 | 0.9400 | 0.9450 | 0.9450 | 0.8900 | 0.9000 | 0.9100 | 0.9250 | 0.9500 | 0.9500 |
| 0.2 | 0.9050 | 0.8950 | 0.9000 | 0.9600 | 0.9450 | 0.9600 | 0.8850 | 0.8950 | 0.9050 | 0.9450 | 0.9650 | 0.9550 |
| 0.3 | 0.8950 | 0.9200 | 0.8950 | 0.9500 | 0.9500 | 0.9300 | 0.8950 | 0.8850 | 0.9050 | 0.9350 | 0.9550 | 0.9600 |
| 0.4 | 0.9050 | 0.9050 | 0.9100 | 0.9300 | 0.9550 | 0.9450 | 0.8850 | 0.8900 | 0.9150 | 0.9550 | 0.9450 | 0.9450 |
| 0.5 | 0.9150 | 0.8900 | 0.9050 | 0.9650 | 0.9550 | 0.9600 | 0.9300 | 0.9100 | 0.8800 | 0.9700 | 0.9400 | 0.9500 |
| 0.6 | 0.9100 | 0.9100 | 0.8900 | 0.9550 | 0.9650 | 0.9700 | 0.8800 | 0.9050 | 0.9050 | 0.9600 | 0.9450 | 0.9500 |
| 0.7 | 0.9100 | 0.9050 | 0.9050 | 0.9600 | 0.9600 | 0.9450 | 0.9100 | 0.9150 | 0.9000 | 0.9450 | 0.9600 | 0.9450 |
| 0.8 | 0.8850 | 0.9100 | 0.9050 | 0.9450 | 0.9550 | 0.9550 | 0.9000 | 0.9950 | 0.8900 | 0.9350 | 0.9950 | 0.9400 |
| 0.9 | 0.9350 | 0.9150 | 0.9200 | 0.9750 | 0.9350 | 0.9600 | 0.9050 | 0.9250 | 0.9000 | 0.9400 | 0.9550 | 0.9600 |
| 1 | 0.8750 | 0.8950 | 0.9050 | 0.9500 | 0.9350 | 0.9650 | 0.8850 | 0.9050 | 0.9050 | 0.9300 | 0.9650 | 0.9600 |

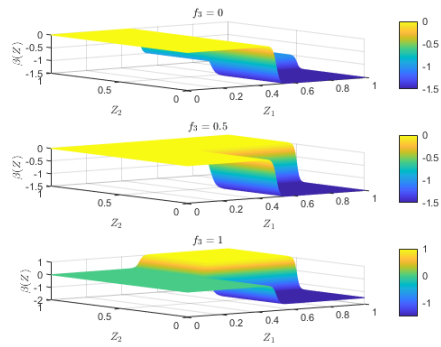## Table 5: Test Points (Coverage – Slope)

The table shows the average coverage for the slope estimation at different test points (MSE). The first column shows the points considered. For the case where $p > 1$, all the values of $\boldsymbol{Z}$ are equal to the number indicated in the first column. Panel (I) considers the case of Model I ($p = 1$). Panel (II) considers the case of Model II ($p = 2$). Panel (III) considers the case of Model III ($p = 3$). Panel (IV) considers the case of Model IV ($p = 5$).

|  | Panel I: $p = 1$ | | | | | | Panel II: $p = 2$ | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | (a): 90% | | | (b): 95% | | | (a): 90% | | | (b): 95% | | |
|  | Sample Size | | | Sample Size | | | Sample Size | | | Sample Size | | |
| Point | 250 | 500 | 1000 | 250 | 500 | 1000 | 250 | 500 | 1000 | 250 | 500 | 1000 |
| 0 | 0.9100 | 0.9150 | 0.8650 | 0.9500 | 0.9450 | 0.9450 | 0.6350 | 0.7250 | 0.6600 | 0.7200 | 0.8150 | 0.7800 |
| 0.1 | 0.9000 | 0.9000 | 0.8950 | 0.9450 | 0.9400 | 0.9250 | 0.6350 | 0.8200 | 0.8500 | 0.7350 | 0.8750 | 0.9100 |
| 0.2 | 0.9250 | 0.8950 | 0.9050 | 0.9550 | 0.9450 | 0.9600 | 0.5650 | 0.7500 | 0.7700 | 0.6800 | 0.8350 | 0.8450 |
| 0.3 | 0.8450 | 0.8700 | 0.8800 | 0.9150 | 0.9250 | 0.9450 | 0.7700 | 0.7050 | 0.6150 | 0.8350 | 0.8050 | 0.7400 |
| 0.4 | 0.8900 | 0.9100 | 0.9000 | 0.9500 | 0.9550 | 0.9650 | 0.2350 | 0.3900 | 0.4100 | 0.3450 | 0.5550 | 0.5300 |
| 0.5 | 0.9200 | 0.8900 | 0.8850 | 0.9650 | 0.9400 | 0.9400 | 0.3200 | 0.5300 | 0.6800 | 0.4100 | 0.6700 | 0.7700 |
| 0.6 | 0.8900 | 0.8750 | 0.9150 | 0.9550 | 0.9500 | 0.9800 | 0.2850 | 0.4550 | 0.3900 | 0.4050 | 0.5750 | 0.5150 |
| 0.7 | 0.8650 | 0.8900 | 0.8800 | 0.9400 | 0.9650 | 0.9600 | 0.7450 | 0.8550 | 0.8950 | 0.8450 | 0.9250 | 0.9500 |
| 0.8 | 0.8950 | 0.9000 | 0.8900 | 0.9300 | 0.9500 | 0.9650 | 0.5150 | 0.6950 | 0.7100 | 0.5950 | 0.7750 | 0.7950 |
| 0.9 | 0.8950 | 0.9050 | 0.9100 | 0.9400 | 0.9500 | 0.9500 | 0.5950 | 0.7350 | 0.8300 | 0.7200 | 0.8450 | 0.9100 |
| 1 | 0.8950 | 0.9050 | 0.8850 | 0.9550 | 0.9550 | 0.9450 | 0.5400 | 0.6650 | 0.6500 | 0.6450 | 0.7550 | 0.7600 |

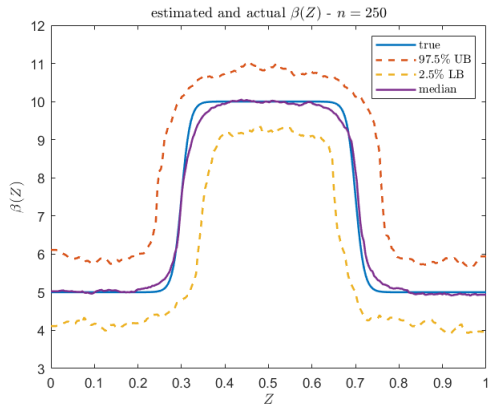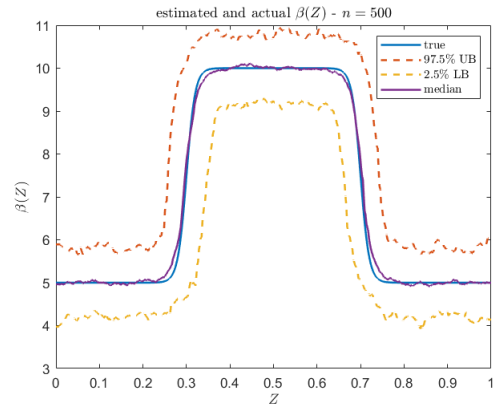|  | Panel III: $p = 3$ | | | | | | Panel IV: $p = 5$ | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | (a): 90% | | | (b): 95% | | | (a): 90% | | | (b): 95% | | |
|  | Sample Size | | | Sample Size | | | Sample Size | | | Sample Size | | |
| Point | 250 | 500 | 1000 | 250 | 500 | 1000 | 250 | 500 | 1000 | 250 | 500 | 1000 |
| 0 | 0.7800 | 0.7950 | 0.8100 | 0.8500 | 0.9200 | 0.8750 | 0.4450 | 0.4400 | 0.2400 | 0.5900 | 0.5750 | 0.3400 |
| 0.1 | 0.8050 | 0.8550 | 0.8900 | 0.8850 | 0.9250 | 0.9400 | 0.4850 | 0.5700 | 0.3700 | 0.6300 | 0.6850 | 0.4750 |
| 0.2 | 0.8250 | 0.8850 | 0.8800 | 0.9150 | 0.9250 | 0.9550 | 0.5200 | 0.5700 | 0.3600 | 0.6200 | 0.7000 | 0.5500 |
| 0.3 | 0.8500 | 0.9200 | 0.8800 | 0.9200 | 0.9500 | 0.9400 | 0.4000 | 0.4800 | 0.3450 | 0.5500 | 0.6100 | 0.4700 |
| 0.4 | 0.7700 | 0.8250 | 0.8500 | 0.8650 | 0.8900 | 0.9000 | 0.3850 | 0.3500 | 0.1750 | 0.5050 | 0.4700 | 0.2500 |
| 0.5 | 0.8900 | 0.8900 | 0.8850 | 0.9400 | 0.9450 | 0.9400 | 0.9000 | 0.8900 | 0.9000 | 0.9550 | 0.9450 | 0.9600 |
| 0.6 | 0.1000 | 0.0900 | 0.1500 | 0.1650 | 0.1650 | 0.2650 | 0.0850 | 0.0550 | 0.0550 | 0.1350 | 0.1450 | 0.1050 |
| 0.7 | 0.2450 | 0.2900 | 0.5900 | 0.3400 | 0.4250 | 0.6950 | 0.1250 | 0.1500 | 0.1750 | 0.2050 | 0.2150 | 0.2750 |
| 0.8 | 0.2350 | 0.4950 | 0.6100 | 0.3600 | 0.5950 | 0.7450 | 0.2000 | 0.7600 | 0.2150 | 0.2900 | 0.8900 | 0.3000 |
| 0.9 | 0.3200 | 0.4150 | 0.6250 | 0.4550 | 0.5350 | 0.7400 | 0.2400 | 0.2150 | 0.1150 | 0.3550 | 0.2950 | 0.2000 |
| 1 | 0.2050 | 0.3000 | 0.2500 | 0.3000 | 0.3950 | 0.3900 | 0.2250 | 0.1150 | 0.0600 | 0.3000 | 0.2100 | 0.0900 |

(a) Model I

(b) Model II

(c) Model III

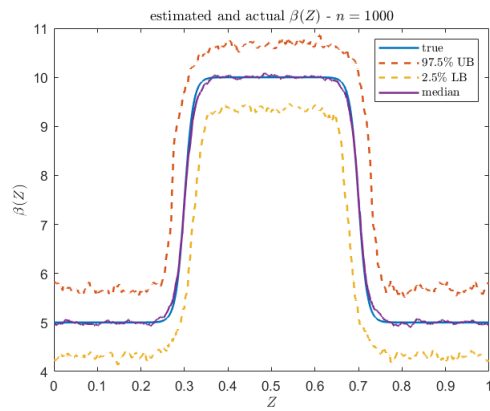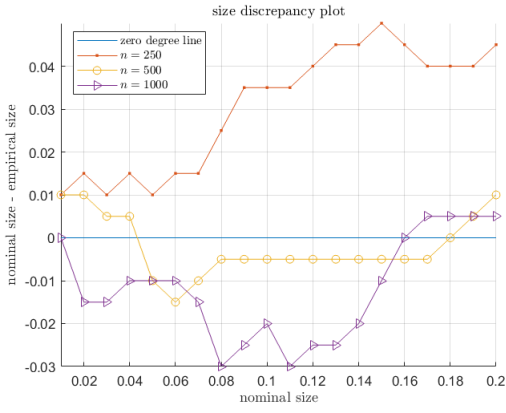Figure 1: Heterogeneity pattern in $\beta_1(\boldsymbol{Z}_i)$ for the simulated models.
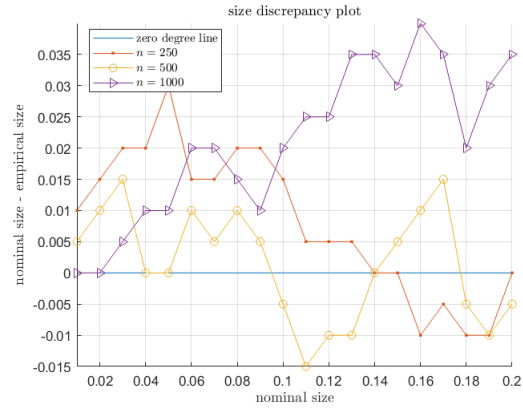
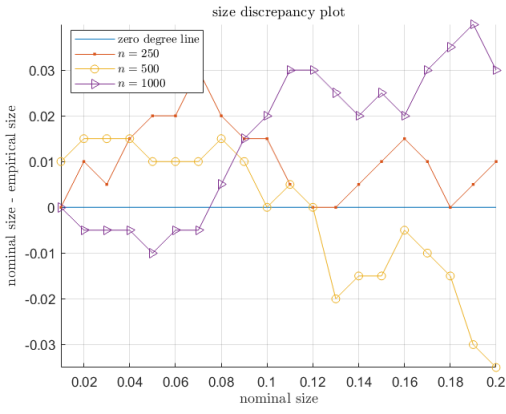(a) $n = 250$

(b) $n = 500$

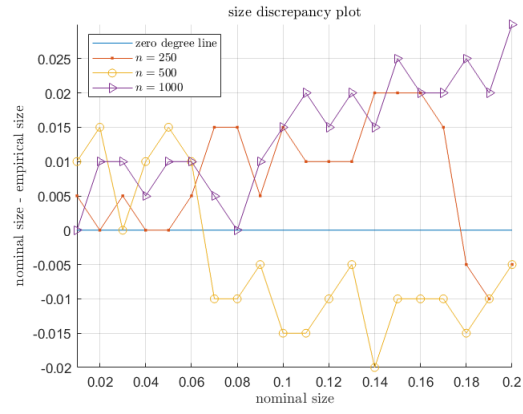(c) $n = 1000$

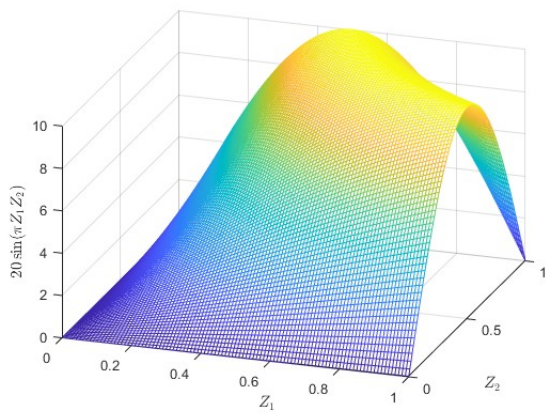Figure 2: Actual and fitted $\beta(Z)$

(a) $p = 1$
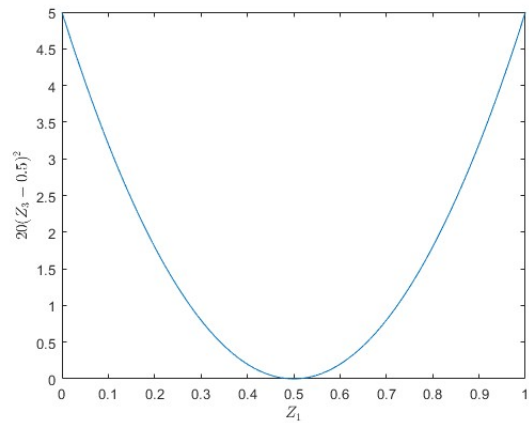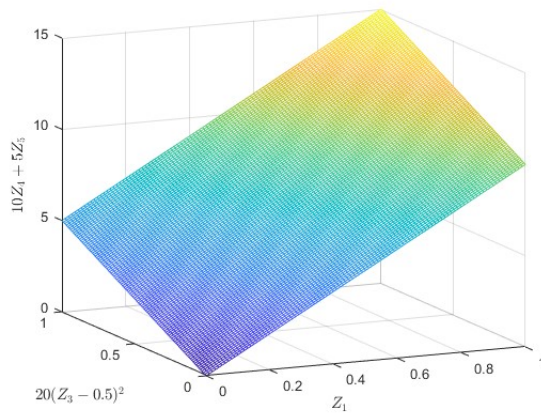
(b) $p = 2$

(c) $p = 3$

(d) $p = 5$

Figure 3: Size discrepancy plots

29

(a) Interaction
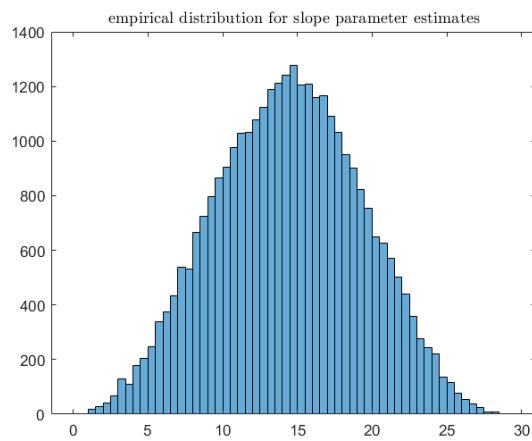
(b) Quadratic



(c) Linear

Figure 4: Function componets
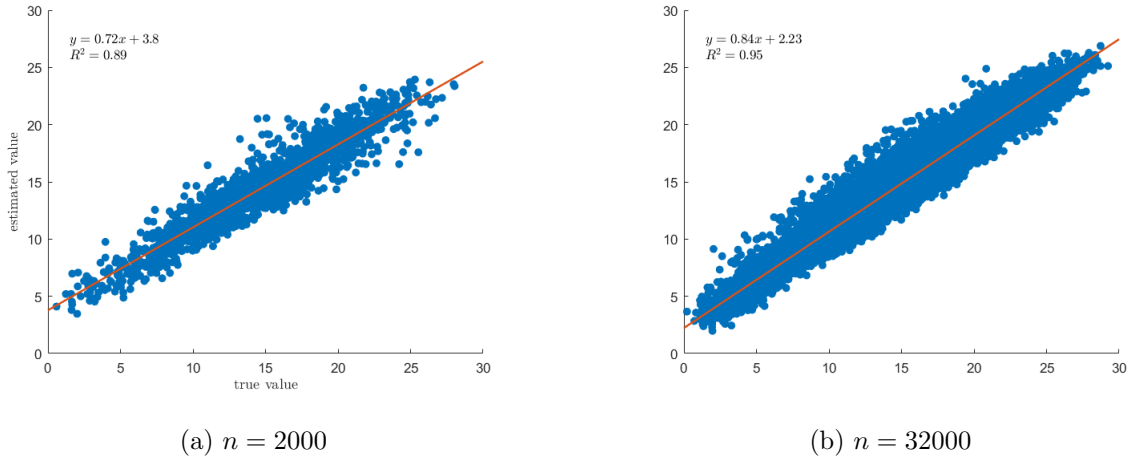


Figure 5: $\gamma$ heterogeneity

(a) $n = 2000$                    (b) $n = 32000$

Figure 6: Scatter plot of the true slope parameter versus the estimated one



Figure 7: Bias and MSE for the estimation of the slope parameter

31

Figure 8: Growth heterogeneity

Figure 9: $\alpha$ heterogeneity

Figure 10: $\gamma$ heterogeneity

# References

Areosa, W., McAleer, M., and Medeiros, M. (2011). Moment-based estimation of smooth transition regression models with endogenous variables. *Journal of Econometrics*, 165:100–111.

Athey, S. and Imbens, G. (2016). Recursive partitioning for heterogeneous causal effects. *Proceedings of the National Academy of Sciences*, 113:7353–7360.

Athey, S. and Imbens, G. (2019). Machine learning methods that economists should know about. *Annual Review of Economics*, 11:685–725.

Athey, S., Tibishirani, J., and Wager, S. (2019). Generalized random forests. *Annals of Statistics*, 47:1148–1178.

Barro, R. and Sala-i-Martin, X. (1992). Convergence. *Journal of Political Economy*, 100:223–251.

Bonhomme, S. and Manresa, E. (2015). Grouped patterns of heterogeneity in panel data. *Econometrica*, 83:1147–1184.

Breiman, L. (2001). Random forests. *Machine Learning*, 45:5–32.

Cai, Z., Fan, J., and Yao, Q. (2000). Functional-coefficient regression models for nonlinear time series. *Journal of the American Statistical Association*, 95:941–956.

Cattaneo, M., Crump, R., Farrell, M., and Feng, Y. (2024). On binscatter. *American Economic Review*, 114:1488–1514.

Chan, K. S. and Tong, H. (1986). On estimating thresholds in autoregressive models. *Journal of Time Series Analysis*, 7:179–190.

Chen, R. and Tsay, R. S. (1993). Functional coefficient autoregressive models. *Journal of the American Statistical Association*, 88:298–308.

Chernozhukov, V., Fernández-Val, I., and Luo, Y. (2018). The sorted effects method: Discovering heterogeneous effects beyond their averages. *Econometrica*, 86:1911–1938.

Cribari-Neto, F., Garcia, N. L., and Vasconcellos, K. L. P. (2000). A note on inverse moments of binomial variates. *Brazilian Review of Econometrics*, 20(2).

Dagenais, M. G. (1969). A threshold regression model. *Econometrica*, 37:193–203.

Fan, J., Zhang, C., and Zhang, J. (2001). Generalized likelihood ratio statistics and wilks phenomenon. *Annals of Statistics*, 29:153–193.

Fan, J. and Zhang, W. (1999). Statistical estimation in varying coefficient models. *Annals of Statistics*, 27:1491–1518.

Feller, W. (1957). *An Introduction to Probability Theory and Its Applications*. Number v. 1-2 in An Introduction to Probability Theory and Its Applications. Wiley.

Friedberg, R., Tibshirani, J., Athey, S., and Wager, S. (2021). Local linear forests. *Journal of Computational and Graphical Statistics*, 30:503–517.

Goldfeld, S. M. and Quandt, R. (1972). *Nonlinear Methods in Econometrics*. North Holland, Amsterdam.

Hansen, B. (2000). Sample splitting and threshold estimation. *Econometrica*, 575:575–603.

Hastie, T. and Tibishirani, R. (1993). Varying-coefficient models. *Journal of the Royal Statistical Society. Series B (Methodological)*, 55:757–796.

Kiefer, N. M. (1978). Discrete parameter variation: Efficient estimation of a switching regression model. *Econometrica*, 46:427–434.

Masini, R., Mendes, E., and Medeiros, M. (2023). Machine learning advances for time series forecasting. *Journal of Economic Surveys*, 37:76–111.

Medeiros, M. C. and Veiga, A. (2000). A hybrid linear-neural model for time series forecasting. *IEEE Transactions on Neural Networks*, 11:1402–1412.

Medeiros, M. C. and Veiga, A. (2005). A flexible coefficient smooth transition time series model. *IEEE Transactions on Neural Networks*, 16:97–113.

Peng, W., Coleman, T., and Mentch, L. (2022). Rates of convergence for random forests via generalized U-statistics. *Electronic Journal of Statistics*, 16(1):232 – 292.

Quandt, R. (1972). A new approach to estimating switching regression. *Journal of the American Statistical Association*, 67:306–310.

Suarez-Fariñas, Pedreira, C., and Medeiros, M. C. (2004). Local-global neural networks: A new approach for nonlinear time series modelling. *Journal of the American Statistical Association*, 99:1092–1107.

Teräsvirta, T. (1994). Specification, estimation, and evaluation of smooth transition autoregressive models. *Journal of the American Statistical Association*, 89:208–218.

Tong, H. and Lim, K. (1980). Threshold autoregression, limit cycles and cyclical data (with discussion). *Journal of the Royal Statistical Society, Series B*, 42:245–292.

Tsay, R. (1989). Testing and modeling threshold autoregressive processes. *Journal of the American Statistical Association*, 84:431–452.

Vaart, A. W. v. d. (1998). *Asymptotic Statistics*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press.

Wager, S. and Athey, S. (2018). Estimation and inference of heterogeneous treatment effects using random forests. *Journal of the American Statistical Association*, 113:1228–1242.

Wager, S. and Walther, G. (2016). Adaptive concentration of regression trees, with application to random forests.

# A    Proofs

## A.1    Proof of Theorem 1

Decompose

$$\widehat{\boldsymbol{\beta}}(\boldsymbol{z}) - \widetilde{\boldsymbol{\beta}}(\boldsymbol{z}) = \widehat{\boldsymbol{\Omega}}(\boldsymbol{z})^{-1}\widehat{\boldsymbol{\gamma}}(\boldsymbol{z}) - \widetilde{\boldsymbol{\beta}}(\boldsymbol{z}) = \widehat{\boldsymbol{\Omega}}(\boldsymbol{z})^{-1}\left[\sum_{i\in R(\boldsymbol{z})} \boldsymbol{X}_i\boldsymbol{X}_i^\top\big(\boldsymbol{\beta}(\boldsymbol{Z}_i) - \widetilde{\boldsymbol{\beta}}(\boldsymbol{z})\big) + \sum_{i\in R(\boldsymbol{z})} \boldsymbol{X}_i\epsilon_i\right].$$

Denote by $\mathcal{G}$ the sigma-algebra generated by $X_1, \ldots, X_n$ and $1\{Z_1 \in R(\boldsymbol{z},\omega)\}, \ldots, 1\{Z_n \in R(\boldsymbol{z},\omega)\}$. Then for $i \in [n]$

$$\mathbb{E}\left[\frac{1}{|\mathcal{A}(\boldsymbol{z})|}\sum_{i\in R(\boldsymbol{z})} \boldsymbol{X}_i\epsilon_i | \mathcal{G}\right] = \frac{1}{|A(\boldsymbol{z})|}\sum_{i\in R(\boldsymbol{z})} \boldsymbol{X}_i\mathbb{E}[\epsilon_i | \boldsymbol{X}_i, 1\{\boldsymbol{Z}_i \in R(\boldsymbol{z},\omega)\}] = 0,$$

where the last equality holds because $R(\boldsymbol{z},\omega)$ is independent of $\epsilon_i$ due to the sample split. Furthermore, the second term is also zero since $\mathbb{E}[\boldsymbol{\beta}(\boldsymbol{Z}_i) - \widetilde{\boldsymbol{\beta}}(\boldsymbol{z})|\mathcal{G}] = \mathbb{E}[\boldsymbol{\beta}(\boldsymbol{Z}) - \widetilde{\boldsymbol{\beta}}(\boldsymbol{z})|\boldsymbol{Z} \in R(\boldsymbol{z},\omega)] = 0$ by the definition of $\widetilde{\boldsymbol{\beta}}$. Hence $\mathbb{E}[\widehat{\boldsymbol{\beta}}(\boldsymbol{z})|\mathcal{G}] = \widetilde{\boldsymbol{\beta}}(\boldsymbol{z})$ and, therefore $\mathbb{E}[\widehat{\boldsymbol{\beta}}(\boldsymbol{z})] = \widetilde{\boldsymbol{\beta}}(\boldsymbol{z})$ for all $\boldsymbol{z} \in [0,1]^d$.

## A.2    Proof of Theorem 2

Recall that based on a subsample $\mathcal{S} \subseteq [n]$, we have the partition $\mathcal{A} \cup \mathcal{B} = \mathcal{S}$ and we define $\mathcal{A}(\boldsymbol{z},\omega) \subseteq \mathcal{A}$ as the set of indices $i \in \mathcal{A}$ such that $\boldsymbol{Z}_i \in R(\boldsymbol{z},\omega)$. By Assumption 2(b) we have $Y_i = \boldsymbol{X}_i^\top\boldsymbol{\beta}(\boldsymbol{Z}_i) + \epsilon_i$ where $\mathbb{E}[\epsilon_i|\boldsymbol{X}_i, \boldsymbol{Z}_i] = 0$, and write

$$\widehat{\boldsymbol{\beta}}(\boldsymbol{z},w) - \boldsymbol{\beta}(\boldsymbol{z}) = \left[\sum_{i\in\mathcal{A}(\boldsymbol{z},\omega)} \boldsymbol{X}_i\boldsymbol{X}_i^\top\right]^{-1}\left\{\sum_{i\in\mathcal{A}(\boldsymbol{z},\omega)} \boldsymbol{X}_i\boldsymbol{X}_i^\top\boldsymbol{\beta}(\boldsymbol{Z}_i) + \sum_{i\in\mathcal{A}(\boldsymbol{z},\omega)} \boldsymbol{X}_i\epsilon_i\right\} - \boldsymbol{\beta}(\boldsymbol{z})$$

$$= \widehat{\boldsymbol{\Omega}}(\boldsymbol{z})^{-1}\left\{\frac{1}{|\mathcal{A}(\boldsymbol{z},\omega)|}\sum_{i\in\mathcal{A}(\boldsymbol{z},\omega)} \boldsymbol{X}_i\boldsymbol{X}_i^\top\big[\boldsymbol{\beta}(\boldsymbol{Z}_i) - \boldsymbol{\beta}(\boldsymbol{z})\big] + \frac{1}{|\mathcal{A}(\boldsymbol{z},\omega)|}\sum_{i\in\mathcal{A}(\boldsymbol{z},\omega)} \boldsymbol{X}_i\epsilon_i\right\}.$$

Let $\Delta(\boldsymbol{z},\omega) := \widehat{\boldsymbol{\Omega}}(\boldsymbol{z},\omega) - \boldsymbol{\Omega}(\boldsymbol{z})$, then

$$\widehat{\boldsymbol{\Omega}}(\boldsymbol{z},\omega)^{-1} = [\boldsymbol{\Omega}(\boldsymbol{z}) + \Delta(\boldsymbol{z})]^{-1} = \boldsymbol{\Omega}(\boldsymbol{z})^{-1} - \boldsymbol{\Omega}(\boldsymbol{z})^{-1}\Delta(\boldsymbol{z},\omega)\widehat{\boldsymbol{\Omega}}(\boldsymbol{z},\omega),$$

Therefore, we have the following decomposition

$$\widehat{\boldsymbol{\beta}}(\boldsymbol{z},w) - \boldsymbol{\beta}(\boldsymbol{z}) = \boldsymbol{\Omega}(\boldsymbol{z})^{-1}\Big\{I_{d_X} - \Delta(\boldsymbol{z},\omega)\big[\Delta(\boldsymbol{z},\omega) + \boldsymbol{\Omega}(\boldsymbol{z})\big]\Big\} \tag{S.21}$$

$$\left\{\frac{1}{|\mathcal{A}(\boldsymbol{z},\omega)|}\sum_{i\in\mathcal{A}(\boldsymbol{z},\omega)} \boldsymbol{X}_i\boldsymbol{X}_i^\top\big[\boldsymbol{\beta}(\boldsymbol{Z}_i) - \boldsymbol{\beta}(\boldsymbol{z})\big] + \frac{1}{|\mathcal{A}(\boldsymbol{z},\omega)|}\sum_{i\in\mathcal{A}(\boldsymbol{z},\omega)} \boldsymbol{X}_i\epsilon_i\right\}.$$

Apply Lemma 5 with $W_i = X_{i,j} X_{i,k}$ for $j, k \in [d_X]$ to conclude that

$$\Delta(\boldsymbol{z}, \omega) \lesssim_{\mathbb{P}} k^{-\frac{1}{2}} \left(\frac{s}{k}\right)^{\epsilon/2} + \left(\frac{k}{s}\right)^{(1-\delta)\frac{K(\alpha)\pi}{d_Z}}. \tag{S.22}$$

Similarly, Lemma 5 with $W_i = X_{i,j} \epsilon_i$ for $j \in [d_X]$ yields

$$\frac{1}{|\mathcal{A}(\boldsymbol{z}, \omega)|} \sum_{i \in \mathcal{A}(\boldsymbol{z}, \omega)} \boldsymbol{X}_i \epsilon_i \lesssim_{\mathbb{P}} k^{-\frac{1}{2}} \left(\frac{s}{k}\right)^{\epsilon/2} + \left(\frac{k}{s}\right)^{(1-\delta)\frac{K(\alpha)\pi}{d_Z}}.$$

Also, for $j \in [d_X]$, we have

$$\frac{1}{|\mathcal{A}(\boldsymbol{z}, \omega)|} \sum_{i \in \mathcal{A}(\boldsymbol{z}, \omega)} \left[\boldsymbol{X}_i \boldsymbol{X}_i^{\mathsf{T}} \big[\boldsymbol{\beta}(\boldsymbol{Z}_i) - \boldsymbol{\beta}(\boldsymbol{z})\big]\right]_j = \sum_{\ell=1}^{d_X} \frac{1}{|\mathcal{A}(\boldsymbol{z}, \omega)|} \sum_{i \in \mathcal{A}(\boldsymbol{z}, \omega)} X_{i,j} X_{i\ell} \big[\beta_\ell(\boldsymbol{Z}_i) - \beta_\ell(\boldsymbol{z})\big].$$

By Cauchy-Schwartz inequality, (B) and Markov's inequality we have

$$\frac{1}{|\mathcal{A}(\boldsymbol{z}, \omega)|} \sum_{i \in \mathcal{A}(\boldsymbol{z}, \omega)} X_{i,j} X_{i\ell} \big[\beta_\ell(\boldsymbol{Z}_i) - \beta_\ell(\boldsymbol{z})\big] \leq \left(\frac{1}{|\mathcal{A}(\boldsymbol{z}, \omega)|} \sum_{i \in \mathcal{A}(\boldsymbol{z}, \omega)} (X_{i,j} X_{i\ell})^2 \frac{1}{|\mathcal{A}(\boldsymbol{z}, \omega)|} \sum_{i \in \mathcal{A}(\boldsymbol{z}, \omega)} \big[\beta_\ell(\boldsymbol{Z}_i) - \beta_\ell(\boldsymbol{z})\big]^2\right)^{1/2}$$

$$\leq \left(\frac{1}{|\mathcal{A}(\boldsymbol{z}, \omega)|} \sum_{i \in \mathcal{A}(\boldsymbol{z}, \omega)} (X_{i,j} X_{i\ell})^2\right)^{1/2} \max_{u \in R(\boldsymbol{z}, \omega)} \|\boldsymbol{\beta}(u) - \boldsymbol{\beta}(\boldsymbol{z})\|$$

$$\lesssim_{\mathbb{P}} \big[\mathbb{E}(X_j^2 X_k^2 | \boldsymbol{Z} \in R(\boldsymbol{z}, \omega)\big]^{1/2} \max_{u \in R(\boldsymbol{z}, \omega)} \|\boldsymbol{\beta}(u) - \boldsymbol{\beta}(\boldsymbol{z})\|$$

$$\lesssim \mathrm{diam}(R(\boldsymbol{z}, \omega)).$$

Finally, the union bound followed by Lemma 1(b) yields

$$\frac{1}{|\mathcal{A}(\boldsymbol{z}, \omega)|} \sum_{i \in \mathcal{A}(\boldsymbol{z}, \omega)} \boldsymbol{X}_i \boldsymbol{X}_i^{\mathsf{T}} \big[\boldsymbol{\beta}(\boldsymbol{Z}_i) - \boldsymbol{\beta}(\boldsymbol{z})\big] \lesssim_{\mathbb{P}} \mathrm{diam}(R(\boldsymbol{z}, \omega))) \lesssim_{\mathbb{P}} \left(\frac{k}{s}\right)^{(1-\delta)\frac{K(\alpha)\pi}{d_Z}}.$$

### A.3   Proof of Theorem 3

Taking the average of all trees using decomposition (S.21), we are left with

$$\overline{\boldsymbol{\beta}}(\boldsymbol{z}, w) - \boldsymbol{\beta}(\boldsymbol{z}) = \boldsymbol{\Omega}(\boldsymbol{z})^{-1} \frac{1}{B} \sum_{b=1}^{B} \left\{ \left[\boldsymbol{I}_{d_X} - \Delta(\boldsymbol{z}, \omega)\big(\Delta(\boldsymbol{z}, \omega) + \boldsymbol{\Omega}(\boldsymbol{z})\big)\right] \right. \tag{S.23}$$

$$\left. \left[\frac{1}{|\mathcal{A}(\boldsymbol{z}, \omega)|} \sum_{i \in \mathcal{A}(\boldsymbol{z}, \omega)} \boldsymbol{X}_i \boldsymbol{X}_i^{\mathsf{T}} \big[\boldsymbol{\beta}(\boldsymbol{Z}_i) - \boldsymbol{\beta}(\boldsymbol{z})\big] + \frac{1}{|\mathcal{A}(\boldsymbol{z}, \omega)|} \sum_{i \in \mathcal{A}(\boldsymbol{z}, \omega)} \boldsymbol{X}_i \epsilon_i\right] \right\}.$$

From the expression above, we note that the term in the form $\frac{1}{B} \sum_{b=1}^{B} \Delta(\boldsymbol{z}, \omega) \Delta(\boldsymbol{z}, \omega)$ will not vanish in probability unless each tree is consistent. In other words, the tree consistency is necessary for the random forest consistency in our setup.

From (S.22) we have

$$\overline{\beta}(z, w) - \beta(z) = \Omega(z)^{-1} \frac{1}{B} \sum_{b=1}^{B} \left\{ \left[ I_{d_X} + o_{\mathbb{P}}(1) \right] \right.$$

$$\left. \left[ \frac{1}{|\mathcal{A}(z,\omega)|} \sum_{i \in \mathcal{A}(z,\omega)} X_i X_i^\mathsf{T} [\beta(Z_i) - \beta(z)] + \frac{1}{|\mathcal{A}(z,\omega)|} \sum_{i \in \mathcal{A}(z,\omega)} X_i \epsilon_i \right] \right\}.$$

Consider the term $\mathcal{T}_1 := \frac{1}{B} \sum_{b=1}^{B} \frac{1}{|\mathcal{A}(z,\omega)|} \sum_{i \in \mathcal{A}(z,\omega)} X_i \epsilon_i$. By Lemma 5, we have that

$$\mathbb{E}[\mathcal{T}_1] = \sum_{i=1}^{s} \mathbb{E}[S_i X_i \epsilon_i] = \mathbb{E}[X\epsilon | Z \in R(z, \omega)] = \mathbb{E}\big[ X \mathbb{E}[\epsilon | X, Z \in R(z, \omega)] \big] = 0.$$

Below, we show that, as $n \to \infty$,

(i) $\Lambda(z)^{-1/2} \frac{1}{B} \sum_{b=1}^{B} \frac{1}{|\mathcal{A}(z,\omega_b)|} \sum_{i \in \mathcal{A}(z,\omega_b)} X_i \epsilon_i \xrightarrow{d} \mathsf{N}(0, I)$ for some covariance matrix $\Lambda(z)$;

(ii) $\Sigma(z)^{-1/2} \Omega(z)^{-1} \frac{1}{|\mathcal{A}(z,\omega)|} \sum_{i \in \mathcal{A}(z,\omega)} X_i X_i^\mathsf{T} [\beta(Z_i) - \beta(z)] \xrightarrow{\mathbb{P}} 0$ where $\Sigma(x) := \Omega^{-1}(z) \Lambda(z) \Omega^{-1}$,

the result then follows by Slutsky's lemma.

**Proof of** $(i)$

Recall that, for a sequence of random vectors $(Z_n)$, $Z_n \xrightarrow{d} \mathsf{N}(0, I)$ is equivalent to $a^\mathsf{T} Z_n \xrightarrow{d} \mathsf{N}(0, 1)$ for all $\|a\| = 1$. Set $W_i = a^\mathsf{T} X_i \epsilon_i$ in Lemma 6 to conclude for all $\|a\| = 1$, provided that $k \asymp s^\eta$ for $\eta \in (1 - K(\alpha, \pi), 1) \subseteq (0, 1)$, $s \to \infty$ and $s(\log n)^{d_Z} = o(n)$, we have

$$\overline{T}_B(z) := \frac{1}{\sqrt{a^\mathsf{T} \Lambda(z) a}} \frac{1}{B} \sum_{b=1}^{B} T_W(z, \omega_b) \xrightarrow{d} \mathsf{N}(0, 1),$$

where let $\mathcal{H}_1(v) := \mathbb{E}[T_W(z, \omega; Z_1, \ldots, Z_s) | Z_1 = v]$

$$\Lambda(z) := \frac{s^2}{n^2} \sum_{i=1}^{n} \mathbb{V}[\mathcal{H}_1(Z_i)] = \frac{s^2}{n} \mathbb{V}[X_1 \epsilon_1 \mathbb{E}[S_1(z, \omega) | Z_1]],$$

and $\frac{s^{\frac{1-\eta}{K(\alpha)}}}{n(\log s)^{d_Z}} \lesssim \min_{a:\|a\|=1} \lambda^2(z) a^\mathsf{T} \Lambda(z) a.$

**Proof of** $(ii)$

For $j \in [d_X]$ we wirte

$$\frac{1}{|\mathcal{A}(z,\omega)|} \sum_{i \in \mathcal{A}(z,\omega)} [X_i X_i^\mathsf{T} [\beta(Z_i) - \beta(z)]]_j = \sum_{\ell=1}^{d_X} \frac{1}{|\mathcal{A}(z,\omega)|} \sum_{i \in \mathcal{A}(z,\omega)} X_{i,j} X_{i\ell} [\beta_\ell(Z_i) - \beta_\ell(z)]$$

40

By Cauchy-Schwartz inequality and Lemma 5

$$\frac{1}{|\mathcal{A}(\boldsymbol{z},\omega)|}\sum_{i\in\mathcal{A}(\boldsymbol{z},\omega)}X_{i,j}X_{i\ell}\big[\beta_\ell(\boldsymbol{Z}_i)-\beta_\ell(\boldsymbol{z})\big]\le\left(\frac{1}{|\mathcal{A}(\boldsymbol{z},\omega)|}\sum_{i\in\mathcal{A}(\boldsymbol{z},\omega)}(X_{i,j}X_{i\ell})^2\frac{1}{|\mathcal{A}(\boldsymbol{z},\omega)|}\sum_{i\in\mathcal{A}(\boldsymbol{z},\omega)}\big[\beta_\ell(\boldsymbol{Z}_i)-\beta_\ell(\boldsymbol{z})\big]^2\right)^{1/2}$$

$$\le\left(\frac{1}{|\mathcal{A}(\boldsymbol{z},\omega)|}\sum_{i\in\mathcal{A}(\boldsymbol{z},\omega)}(X_{i,j}X_{i\ell})^2\right)^{1/2}\max_{u\in R(\boldsymbol{z},\omega)}\|\boldsymbol{\beta}(u)-\boldsymbol{\beta}(\boldsymbol{z})\|$$

$$\le\big[\mathbb{E}(X_j^2X_k^2|\boldsymbol{Z}=\boldsymbol{z})\big]^{1/2}\max_{u\in R(\boldsymbol{z},\omega)}\|\boldsymbol{\beta}(u)-\boldsymbol{\beta}(\boldsymbol{z})\|+o_\mathbb{P}(1).$$

Thus

$$\left\|\frac{1}{|\mathcal{A}(\boldsymbol{z},\omega)|}\sum_{i\in\mathcal{A}(\boldsymbol{z},\omega)}\boldsymbol{X}_i\boldsymbol{X}_i^{\mathsf{T}}\big[\boldsymbol{\beta}(\boldsymbol{Z}_i)-\boldsymbol{\beta}(\boldsymbol{z})\big]\right\|\lesssim_\mathbb{P}\mathrm{diam}(R(\boldsymbol{z},\omega)).$$

By the Lipschitz condition on $\boldsymbol{z}\mapsto\boldsymbol{\beta}(\boldsymbol{z})$ (Assumption 2(b)), we have

$$\|\Sigma(\boldsymbol{z})^{-1/2}\boldsymbol{\Omega}(\boldsymbol{z})^{-1}\frac{1}{|\mathcal{A}(\boldsymbol{z},\omega)|}\sum_{i\in\mathcal{A}(\boldsymbol{z},\omega)}\boldsymbol{X}_i\boldsymbol{X}_i^{\mathsf{T}}\big[\boldsymbol{\beta}(\boldsymbol{Z}_i)-\boldsymbol{\beta}(\boldsymbol{z})\big]\|\le C_{p_X}\frac{\|\boldsymbol{\Omega}(\boldsymbol{z})^{-1}\|^2}{\|\Lambda(\boldsymbol{z})\|^{1/2}}\max_{u\in R(\boldsymbol{z},\omega)}\|\boldsymbol{\beta}(u)-\boldsymbol{\beta}(\boldsymbol{z})\|$$

$$\le C_{p_X}C_{\boldsymbol{\beta}}\frac{\|\boldsymbol{\Omega}(\boldsymbol{z})^{-1}\|^2}{\|\Lambda(\boldsymbol{z})\|^{1/2}}\mathrm{diam}(R(\boldsymbol{z},\omega)),$$

where $C_{p_X}$ is a constant only depending on $p_X$ and $C_{\boldsymbol{\beta}}$ is the Lipschitz constant.

Since we assume that $\|\boldsymbol{\Omega}(\boldsymbol{z})^{-1}\|_2\lesssim 1$ by Assumption 2(d), it suffices for $(ii)$ that $\mathrm{diam}(R(\boldsymbol{z},\omega))=o_\mathbb{P}(\|\Lambda(\boldsymbol{z})\|^{1/2})$. For that, we have from Lemma 2 and the proof of part $(i)$,

$$\frac{\mathrm{diam}(R(\boldsymbol{z},\omega))}{\|\Lambda(\boldsymbol{z})\|^{1/2}}\lesssim_\mathbb{P}s^{-\frac{(1-\eta)\pi K(\alpha)}{2d_Z}}\left(\frac{s^{\frac{1-\eta}{K(\alpha)}}}{n(\log s)^{d_Z}}\right)^{-1/2}\lesssim n^{-\frac{1}{2}\left[\boldsymbol{\beta}(1-\eta)\left(\frac{\pi K(\alpha)}{d_Z}+\frac{1}{K(\alpha)}\right)-1\right]}(\log n)^{d_Z/2}.$$

By assumption $\boldsymbol{\beta}(1-\eta)\big(\frac{\pi K(\alpha)}{d_Z}+\frac{1}{K(\alpha)}\big)>1$ hence the right-hand side converges to 0 is probability which proves $(ii)$.

## A.4 Proof of Theorem 4

Recall the random forest decomposition.

$$\overline{\boldsymbol{\beta}}(\boldsymbol{z},w)-\boldsymbol{\beta}(\boldsymbol{z})=\boldsymbol{\Omega}(\boldsymbol{z})^{-1}\frac{1}{B}\sum_{b=1}^B\left\{\left[I_{d_X}-\Delta(\boldsymbol{z},\omega)\big(\Delta(\boldsymbol{z},\omega)+\boldsymbol{\Omega}(\boldsymbol{z})\big)\right]\right.$$

$$\left.\left[\frac{1}{|\mathcal{A}(\boldsymbol{z},\omega)|}\sum_{i\in\mathcal{A}(\boldsymbol{z},\omega)}\boldsymbol{X}_i\boldsymbol{X}_i^{\mathsf{T}}\big[\boldsymbol{\beta}(\boldsymbol{Z}_i)-\boldsymbol{\beta}(\boldsymbol{z})\big]+\frac{1}{|\mathcal{A}(\boldsymbol{z},\omega)|}\sum_{i\in\mathcal{A}(\boldsymbol{z},\omega)}\boldsymbol{X}_i\epsilon_i\right]\right\}.$$

Under $\mathcal{H}_0$, the bias term (the first term in the square brackets) vanishes, and we are left with

$$\overline{\boldsymbol{\beta}}(\boldsymbol{z}, w) - \beta_0 = \boldsymbol{\Omega}(\boldsymbol{z})^{-1} \frac{1}{B} \sum_{b=1}^{B} \left\{ \left[ I_{d_X} - \Delta(\boldsymbol{z}, \omega)\big(\Delta(\boldsymbol{z}, \omega) + \boldsymbol{\Omega}(\boldsymbol{z})\big) \right] \frac{1}{|\mathcal{A}(\boldsymbol{z}, \omega)|} \sum_{i \in \mathcal{A}(\boldsymbol{z}, \omega)} \boldsymbol{X}_i \epsilon_i \right\}$$

$$= \boldsymbol{\Omega}(\boldsymbol{z})^{-1} \left[ \frac{1}{B} \sum_{b=1}^{B} \frac{1}{|\mathcal{A}(\boldsymbol{z}, \omega)|} \sum_{i \in \mathcal{A}(\boldsymbol{z}, \omega)} \boldsymbol{X}_i \epsilon_i + \zeta_1(\boldsymbol{z}) \right]$$

$$= \boldsymbol{\Omega}(\boldsymbol{z})^{-1} \left[ \frac{s}{n} \sum_{i=1}^{n} \mathcal{H}_1(\boldsymbol{z}, W_i) + \zeta_1(\boldsymbol{z}) + \zeta_2(\boldsymbol{z}) \right],$$

where $\zeta(\boldsymbol{z}) = \zeta_1(\boldsymbol{z}) + \zeta_2(\boldsymbol{z})$ with

$$\zeta_1(\boldsymbol{z}) := \frac{1}{B} \sum_{b=1}^{B} \left\{ -\Delta(\boldsymbol{z}, \omega)\big(\Delta(\boldsymbol{z}, \omega) + \boldsymbol{\Omega}(\boldsymbol{z})\big) \frac{1}{|\mathcal{A}(\boldsymbol{z}, \omega)|} \sum_{i \in \mathcal{A}(\boldsymbol{z}, \omega)} \boldsymbol{X}_i \epsilon_i \right\}$$

$$\zeta_2(\boldsymbol{z}) := \frac{1}{B} \sum_{b=1}^{B} \frac{1}{|\mathcal{A}(\boldsymbol{z}, \omega_b)|} \sum_{i \in \mathcal{A}(\boldsymbol{z}, \omega_b)} \boldsymbol{X}_i \epsilon_i - \frac{s}{n} \sum_{i=1}^{n} \mathcal{H}_1(\boldsymbol{z}, W_i)$$

$$= \sum_{i=1}^{n} \left[ \omega(\boldsymbol{z}, \boldsymbol{Z}_i) - s\theta(\boldsymbol{z}, \boldsymbol{Z}_i)/n \right] \epsilon_i \boldsymbol{X}_i$$

where $\omega(\boldsymbol{z}, \boldsymbol{Z}_i) := \frac{1}{B} \sum_{b=1}^{B} \frac{1\{\boldsymbol{Z}_i \in \mathcal{A}(\boldsymbol{z}, \omega_b)\}}{|\mathcal{A}(\boldsymbol{z}, \omega_b)|}$. Note that $\mathbb{E}[\zeta_2(Z_j)|Z_1, \ldots, Z_n, \omega] = 0$ for $j \in [n]$ because $\omega(Z_j, \boldsymbol{Z}_i) - s\theta(Z_j, \boldsymbol{Z}_i)/n$ is $Z_1, \ldots, Z_n, \omega$ measurable and $\mathbb{E}[\epsilon_i|\boldsymbol{Z}_i, \boldsymbol{X}_i] = 0$. Also

$$\mathbb{V}[\zeta_2(\boldsymbol{Z}_i|\boldsymbol{Z}) = \sum_{i=1}^{n} \mathbb{V}[\epsilon_i \boldsymbol{X}_i]$$

$$\sum_{j=1}^{n} \zeta_2(Z_j) = \sum_{i=1}^{n} \epsilon_i \boldsymbol{X}_i \sum_{j=1}^{n} \left[ \omega(Z_j, \boldsymbol{Z}_i) - s\theta(Z_j, \boldsymbol{Z}_i)/n \right] = \sum_{i=1}^{n} \epsilon_i \boldsymbol{X}_i q_i(\boldsymbol{Z}^{(n)})$$

and

$$\mathbb{V}[\sum_{j=1}^{n} \zeta_2(Z_j)|\boldsymbol{Z}^{(n)}] = \sigma^2 \sum_{j=1}^{n} \mathbb{E}[\boldsymbol{\Omega}(Z_j)] q_j^2$$

$$\mathbb{E}[\zeta_2(Z_\ell)\zeta_2(Z_k)^{\mathsf{T}}] = \sum_{i=1}^{n} \sum_{j=1}^{n} \mathbb{E}\left\{ \epsilon_i \epsilon_j \left[ \omega(Z_\ell, \boldsymbol{Z}_i) - s\theta(Z_\ell, \boldsymbol{Z}_i)/n \right] \left[ \omega(Z_k, \boldsymbol{Z}_i) - s\theta(Z_k, \boldsymbol{Z}_i)/n \right] \boldsymbol{X}_i \boldsymbol{X}_j^{\mathsf{T}} \right\}$$

$$= \sigma^2 \sum_{i=1}^{n} \left[ \omega(Z_\ell, \boldsymbol{Z}_i) - s\theta(Z_\ell, \boldsymbol{Z}_i)/n \right] \left[ \omega(Z_k, \boldsymbol{Z}_i) - s\theta(Z_k, \boldsymbol{Z}_i)/n \right] \boldsymbol{X}_i \boldsymbol{X}_j^{\mathsf{T}}$$

$$\mathbb{V}[\sum_{i=1}^{n} \zeta_2(\boldsymbol{Z}_i)] \sum_{i=1}^{n} \sum_{j=1}^{n} \left[ \omega(\boldsymbol{z}, \boldsymbol{Z}_i) - s\theta(\boldsymbol{z}, \boldsymbol{Z}_i)/n \right] \epsilon_i \boldsymbol{X}_i$$

42

write

Recall that

$$\mathbb{V}[\zeta_2(\boldsymbol{z})] = \mathbb{V}[T] + \mathbb{V}[T_0]$$

Then

$$\Delta := \mathsf{RSS}_0 - \mathsf{RSS} = 2\underbrace{\sum_{i=1}^{n} \epsilon_i \boldsymbol{X}_i^{\mathsf{T}}(\overline{\boldsymbol{\beta}}(\boldsymbol{Z}_i) - \widetilde{\boldsymbol{\beta}})}_{=:\Delta_1} + \underbrace{\sum_{i=1}^{n}(\overline{\boldsymbol{\beta}}(\boldsymbol{Z}_i) - \widetilde{\boldsymbol{\beta}})^{\mathsf{T}}\boldsymbol{X}_i\boldsymbol{X}_i^{\mathsf{T}}(\overline{\boldsymbol{\beta}}(\boldsymbol{Z}_i) - \widetilde{\boldsymbol{\beta}})}_{=:\Delta_2}$$

For $\boldsymbol{Z}_i, Z_j$ and $\omega$ define

$$S(\boldsymbol{Z}_i, Z_j, \omega) := \begin{cases} |\{k \in \mathcal{A} : Z_k \in R(Z_j, \omega)\}|^{-1} & \text{if } \boldsymbol{Z}_i \in R(Z_j, \omega) \\ 0 & \text{otherwise.} \end{cases}$$

and

$$\theta(\boldsymbol{Z}_i, Z_j) := \mathbb{E}[S(\boldsymbol{Z}_i, Z_j, \omega)|\boldsymbol{Z}_i, Z_j].$$

Note that $\theta$ is symmetric in its two arguments because $S$ is symmetric with respect to its first two arguments. Also, for $i = j$ we have that $S(\boldsymbol{Z}_i, \boldsymbol{Z}_i, \omega) = (1 + B)^{-1}$ where $B \sim \mathrm{Binom}(s - 1, p_i)$ conditional on $\omega$ and $\boldsymbol{Z}_i$ where $p_i \asymp R(\boldsymbol{Z}_i, \omega)$. Then from (S.33) on a event with probability approach 1

$$\mathbb{E}[S(\boldsymbol{Z}_i, \boldsymbol{Z}_i, \omega)|\omega, \boldsymbol{Z}_i] \asymp \frac{1}{(s-1)p_i} \asymp \frac{1}{s(k/s)^{1+o(1)}}$$

therefore

$$\theta(\boldsymbol{Z}_i, \boldsymbol{Z}_i) \asymp_{\mathbb{P}} \frac{1}{s(k/s)^{1+o(1)}}.$$

For $i \neq j$ we have $S(\boldsymbol{Z}_i, Z_j, \omega) \leq S(\boldsymbol{Z}_i, \boldsymbol{Z}_i, \omega)$ then $\theta(\boldsymbol{Z}_i, Z_j) \lesssim_{\mathbb{P}} \frac{1}{s(k/s)^{1+o(1)}}$.

Also, $\mathcal{H}_1(\boldsymbol{Z}_i, Z_j) = \mathbb{E}[\boldsymbol{X}_i\epsilon_i S(\boldsymbol{Z}_i, Z_j, \omega)|\boldsymbol{Z}_i, \boldsymbol{X}_i, Y_i] = \epsilon_i \boldsymbol{X}_i \theta(\boldsymbol{Z}_i, Z_j)$ thus

$$\overline{\boldsymbol{\beta}}(\boldsymbol{Z}_i) - \widetilde{\boldsymbol{\beta}} = \overline{\boldsymbol{\beta}}(\boldsymbol{Z}_i) - \boldsymbol{\beta}(\boldsymbol{Z}_i) + \boldsymbol{\beta}(\boldsymbol{Z}_i) - \widetilde{\boldsymbol{\beta}}$$

$$= \boldsymbol{\Omega}(\boldsymbol{Z}_i)^{-1}\frac{s}{n}\sum_{j=1}^{n}\mathcal{H}_1(\boldsymbol{Z}_i, Z_j) + \zeta(\boldsymbol{Z}_i) - \left(\sum_{j=1}^{n}X_jX_j^{\mathsf{T}}\right)^{-1}\sum_{j=1}^{n}\epsilon_j X_j - \left(\sum_{j=1}^{n}X_jX_j^{\mathsf{T}}\right)^{-1}\sum_{j=1}^{n}(\boldsymbol{\beta}(Z_j) - \boldsymbol{\beta}(\boldsymbol{Z}_i))$$

$$= \boldsymbol{\Omega}(\boldsymbol{z})^{-1}\frac{s}{n}\sum_{j=1}^{n}\epsilon_j X_j\theta(\boldsymbol{z}, Z_j) - \left(\mathbb{E}[\boldsymbol{X}\boldsymbol{X}^{\mathsf{T}}]\right)^{-1}\frac{1}{n}\sum_{j=1}^{n}\epsilon_j X_j + \zeta(\boldsymbol{z}) + O_{\mathbb{P}}(n^{-1})$$

Under $\mathcal{H}_0$, we have $\boldsymbol{\Omega} := \mathbb{E}[\boldsymbol{X}\boldsymbol{X}^{\mathsf{T}}] = \boldsymbol{\Omega}(\boldsymbol{z})$ for $\boldsymbol{z} \in [0, 1]^d$, then

$$\overline{\boldsymbol{\beta}}(\boldsymbol{Z}_i) - \widetilde{\boldsymbol{\beta}} = \boldsymbol{\Omega}^{-1}\frac{1}{n}\sum_{j=1}^{n}\epsilon_j X_j\big[s\theta(\boldsymbol{Z}_i, Z_j) - 1\big] + \zeta(\boldsymbol{Z}_i) + O_{\mathbb{P}}(n^{-1}); \quad i \in [d] \qquad \text{(S.24)}$$

43

Now for $\Delta_1$ we have

$$\Delta_1 := 2\sum_{i=1}^{n} \epsilon_i \boldsymbol{X}_i^{\mathsf{T}}(\overline{\boldsymbol{\beta}}(\boldsymbol{Z}_i) - \widetilde{\boldsymbol{\beta}})$$

$$= \frac{2}{n}\sum_{i=1}^{n}\sum_{j=1}^{n} \epsilon_i\epsilon_j\big[s\theta(\boldsymbol{Z}_i, Z_j) - 1\big]\boldsymbol{X}_i^{\mathsf{T}}\boldsymbol{\Omega}^{-1}X_j + 2\sum_{i=1}^{n}\epsilon_i\boldsymbol{X}_i^{\mathsf{T}}\zeta(\boldsymbol{Z}_i) + O_{\mathbb{P}}(n^{-1})\sum_{i=1}^{n}\epsilon_i\boldsymbol{X}_i$$

$$= \frac{2}{n}\sum_{i=1}^{n}\epsilon_i^2\big[s\theta(\boldsymbol{Z}_i, \boldsymbol{Z}_i) - 1\big]\boldsymbol{X}_i^{\mathsf{T}}\boldsymbol{\Omega}^{-1}\boldsymbol{X}_i + \frac{2}{n}\sum_{i\neq j}^{n}\epsilon_i\epsilon_j\big[s\theta(\boldsymbol{Z}_i, Z_j) - 1\big]\boldsymbol{X}_i^{\mathsf{T}}\boldsymbol{\Omega}^{-1}X_j$$

$$\quad + 2\sum_{i=1}^{n}\epsilon_i\boldsymbol{X}_i^{\mathsf{T}}\zeta(\boldsymbol{Z}_i) + O_{\mathbb{P}}(n^{-1/2})$$

$$=: \Delta_{11} + \Delta_{12} + \mathcal{E}_1$$

.

Note that $\mathbb{E}[\Delta_{11}] = 2\sigma^2 d_X\big(s\mathbb{E}[\theta(Z_1, Z_1)] - 1\big)$ and $\mathbb{V}[\Delta_{11}] = O(\mathbb{V}[\theta(\boldsymbol{Z}, \boldsymbol{Z})]s^2/n) = O(\mathbb{E}[\theta^2(\boldsymbol{Z}, \boldsymbol{Z})]s^2/n) \asymp [(k/s)^{2(1+o(1))}n]^{-1}$. Hence

$$\Delta_{11} = 2\sigma^2 d_X\big(s\mathbb{E}[\theta(Z_1, Z_1)] - 1\big) + O_{\mathbb{P}}\left(\frac{(s/k)^{1+o(1)}}{\sqrt{n}}\right),$$

and therefore,

$$\Delta_1 = 2\sigma^2 d_X\big(s\mathbb{E}[\theta(Z_1, Z_1)] - 1\big) + \frac{2}{n}\sum_{i\neq j}^{n}\epsilon_i\epsilon_j\big[s\theta(\boldsymbol{Z}_i, Z_j) - 1\big]\boldsymbol{X}_i^{\mathsf{T}}\boldsymbol{\Omega}^{-1}X_j$$

$$\quad + O_{\mathbb{P}}\left(\frac{(s/k)^{1+o(1)}}{\sqrt{n}}\right) + \mathcal{E}_1.$$

.

For $\Delta_2$, let $T_{ij} := s\theta(\boldsymbol{Z}_i, Z_j) - 1$ for $i, j \in [d]$, use (S.24) and collect terms we are left with

$$\Delta_2 := \sum_{i=1}^{n}(\overline{\boldsymbol{\beta}}(\boldsymbol{Z}_i) - \widetilde{\boldsymbol{\beta}})^{\mathsf{T}}\boldsymbol{X}_i\boldsymbol{X}_i^{\mathsf{T}}(\overline{\boldsymbol{\beta}}(\boldsymbol{Z}_i) - \widetilde{\boldsymbol{\beta}})$$

$$= \frac{1}{n^2}\sum_{i=1}^{n}\sum_{j=1}^{n}\sum_{k=1}^{n}\epsilon_j\epsilon_k T_{ij}T_{ik}X_j^{\mathsf{T}}\boldsymbol{\Omega}^{-1}\boldsymbol{X}_i\boldsymbol{X}_i^{\mathsf{T}}\boldsymbol{\Omega}^{-1}X_k + \mathcal{E}_2.$$

Decompose the first term in the last expression as $\Delta_{21} + \Delta_{22}$ where

$$\Delta_{21} := \frac{1}{n^2}\sum_{i=1}^{n}\sum_{j=1}^{n}\epsilon_j^2 T_{ij}^2 X_j^{\mathsf{T}}\boldsymbol{\Omega}^{-1}\boldsymbol{X}_i\boldsymbol{X}_i^{\mathsf{T}}\boldsymbol{\Omega}^{-1}X_j$$

$$\Delta_{22} := \frac{1}{n^2}\sum_{i=1}^{n}\sum_{j\neq k}^{n}\epsilon_j\epsilon_k T_{ij}T_{ik}X_j^{\mathsf{T}}\boldsymbol{\Omega}^{-1}\boldsymbol{X}_i\boldsymbol{X}_i^{\mathsf{T}}\boldsymbol{\Omega}^{-1}X_k.$$

Also
$$\mathbb{E}[\Delta_{21}] = \frac{\sigma^2}{n^2} \sum_{i=1}^{n} \sum_{j=1}^{n} \mathbb{E}[T_{ij}^2 \mathsf{tr}(\mathbf{\Omega}(Z_j)\mathbf{\Omega}(\mathbf{Z}_i)^{-1})].$$

Under $\mathcal{H}_0$, we have $\mathbf{\Omega}(\mathbf{Z}_i) = \mathbf{\Omega}$, hence

$$\mathbb{E}[\Delta_{21}] = \frac{\sigma^2 d_X}{n^2} \sum_{i,j} \mathbb{E}[T_{ij}^2] = \frac{\sigma^2 d_X}{n^2}(n\mathbb{E}[T_{11}^2] + n(n-1)\mathbb{E}[T_{12}^2]).$$

Decompose further $\Delta_{22} = \Delta_{221} + \Delta_{222}$ where

$$\Delta_{221} := \frac{1}{n^2} \sum_{i \notin \{j,k\}} \sum_{j \neq k} \epsilon_j \epsilon_k T_{ij} T_{ik} X_j^\mathsf{T} \mathbf{\Omega}(\mathbf{Z}_i)^{-1} \mathbf{X}_i \mathbf{X}_i^\mathsf{T} \mathbf{\Omega}(\mathbf{Z}_i)^{-1} X_k$$

$$\Delta_{222} := \frac{2}{n^2} \sum_{j \neq k} \epsilon_j \epsilon_k T_{jj} T_{jk} X_j^\mathsf{T} \mathbf{\Omega}(Z_j)^{-1} \mathbf{X}_i \mathbf{X}_i^\mathsf{T} \mathbf{\Omega}(Z_j)^{-1} X_k$$

Define

$$H_{jk} := \mathbb{E}[T_{ij} T_{ik} \mathbf{\Omega}(\mathbf{Z}_i)^{-1} \mathbf{X}_i \mathbf{X}_i^\mathsf{T} \mathbf{\Omega}(\mathbf{Z}_i)^{-1} | Z_j, Z_k]$$

$$= \mathbb{E}\Big[T_{ij} T_{ik} \mathbf{\Omega}(\mathbf{Z}_i)^{-1} \mathbb{E}\big[\mathbf{X}_i \mathbf{X}_i^\mathsf{T} | \mathbf{Z}_i, Z_j, Z_k\big] \mathbf{\Omega}(\mathbf{Z}_i)^{-1} | Z_j, Z_k\Big]$$

$$= \mathbb{E}\Big[T_{ij} T_{ik} \mathbf{\Omega}(\mathbf{Z}_i)^{-1} | Z_j, Z_k\Big]$$

where $\mathbf{Z}'$ is an independent copy of $Z_1, \ldots Z_n$. Then

$$\Delta_{221} = \frac{2(n-2)}{n^2} \sum_{1 \leq j < k \leq n} \epsilon_j \epsilon_k X_j^\mathsf{T} H_{jk} X_k + O_\mathbb{P}()$$

Furthermore, under the $\mathcal{H}_0$,

$$H_{jk} = \mathbf{\Omega}^{-1} \mathbb{E}[T_{ij} T_{ik} | Z_j, Z_k] = \mathbf{\Omega}^{-1} \tau(Z_j, Z_k),$$

where $\tau(\mathbf{z}, \mathbf{z}') := \mathbb{E}[(s\theta(\mathbf{Z}, z) - 1)(s\theta(\mathbf{Z}, z') - 1)]$.

Putting everything together, under $\mathcal{H}_0$,

$$\Delta = 2\sigma^2 d_X \mathbb{E}[T_{11}] + \frac{\sigma^2 d_X}{n^2} \sum_{i,j} \mathbb{E}[T_{ij}^2]$$

$$+ \frac{2}{n} \sum_{i \neq j}^{n} \epsilon_i \epsilon_j T_{ij} \mathbf{X}_i^\mathsf{T} \mathbf{\Omega}^{-1} X_j + \frac{2(n-2)}{n^2} \sum_{1 \leq j < k \leq n} \epsilon_j \epsilon_k X_j^\mathsf{T} H_{jk} X_k + O_\mathbb{P}()$$

$$= 2\sigma^2 d_X \mathbb{E}[T_{11}] + \frac{\sigma^2 d_X}{n^2} \sum_{i,j} \mathbb{E}[T_{ij}^2]$$

$$+ \frac{2}{n} \sum_{i \neq j}^{n} \epsilon_i \epsilon_j \big[T_{ij} + \tfrac{n-2}{n} \tau(\mathbf{Z}_i, Z_j)\big] \mathbf{X}_i^\mathsf{T} \mathbf{\Omega}^{-1} X_j$$

$$+ o_\mathbb{P}(1).$$

45

Let $\eta(\boldsymbol{z}, \boldsymbol{z}') := s\theta(\boldsymbol{z}, \boldsymbol{z}') - 1 + \frac{n-2}{n}\tau(\boldsymbol{z}, \boldsymbol{z}')$ for $\boldsymbol{z}, \boldsymbol{z}' \in [0,1]^{d_Z}$ and

$$
\begin{aligned}
\mathbb{E}\big[\epsilon_1^2 \epsilon_2^2 \eta(Z_1, Z_2)^2 (X_1^\mathsf{T} \boldsymbol{\Omega}^{-1} X_2)^2\big] &= \sigma^4 \mathbb{E}\big[\eta(Z_1, Z_2)^2 (X_1^\mathsf{T} \boldsymbol{\Omega}^{-1} X_2)^2\big] \\
&= \sigma^4 \mathbb{E}\Big[\eta(Z_1, Z_2)^2 X_1^\mathsf{T} \boldsymbol{\Omega}^{-1} \mathbb{E}\big[X_2 X_2^\mathsf{T} | Z_1, Z_2, X_1\big] \boldsymbol{\Omega}^{-1} X_1\Big] \\
&= \sigma^4 \mathbb{E}\Big[\eta(Z_1, Z_2)^2 X_1^\mathsf{T} \boldsymbol{\Omega}^{-1} X_1\Big] \\
&= \sigma^4 d_X \mathbb{E}\Big[\eta(Z_1, Z_2)^2\Big].
\end{aligned}
$$

Then $\nu^{-1}(\Delta - \mu) \xrightarrow{d} \mathsf{N}(0,1)$ where

$$
\mu := 2\sigma^2 d_X (s\mathbb{E}[\theta(\boldsymbol{Z}, \boldsymbol{Z})] - 1) + \frac{\sigma^2 d_X}{n^2}(n\mathbb{E}[T_{11}^2] + n(n-1)/2\mathbb{E}[T_{12}^2])
$$

and

$$
\nu^2 := \frac{4n(n-1)}{n^2}\sigma^4 d_X \mathbb{E}\Big[\eta(Z_1, Z_2)^2\Big].
$$

Finally, note that $\tau(\boldsymbol{Z}, \boldsymbol{Z}') \lesssim (s/k)^{2(1+o(1))}$.

# B   Auxiliary Lemmas and Proofs

**Lemma 1** (Upper bounds on the leaf diameter)

*For $\boldsymbol{z} \in [0,1]^d$, let $R(\boldsymbol{z}, \omega)$ denote the unique leaf of a tree grown by Algorithm 1 that contains $z$ and $M_j(\boldsymbol{z})$ is the total number of splits along $Z_j$ to form $R(\boldsymbol{z}, \omega)$ for $j \in [d_Z]$. Then*

*(a) For $\delta \in (0,1)$*

$$
\mathbb{P}\left(M_j(\boldsymbol{z}) \leq \frac{(1-\delta)\pi}{d_Z}\frac{\log(s/(2k-1))}{\log(1/(1-\alpha))}\right) \leq \left(\frac{s}{2k-1}\right)^{-\frac{\delta^2 \pi^2}{2d_Z^2 \log(1/(1-\alpha))}}.
$$

*(b) For $\delta \in (0,1)$,*

$$
\mathbb{P}\left(\operatorname{diam}(R(\boldsymbol{z}, \omega)) \gtrsim \left(\frac{k}{s}\right)^{(1-\delta)\frac{K(\alpha)\pi}{d_Z}}\right) \lesssim \left(\frac{k}{s}\right)^{\frac{\delta^2 \pi^2}{2(1+o(1))^2 d_Z^2 \log(1/\alpha)}} + \frac{1}{s}, \qquad \text{(S.25)}
$$

*where $K(\alpha) := \frac{\log((1-\alpha)^{-1})}{\log(1/\alpha)}$ for $\alpha \in (0, 0.5)$ and $\pi \in (0,1]$;*

*(c) For $t > 0$ and $1 \leq p \leq q < \infty$*

$$
\mathbb{P}\left(\left(\mathbb{E}_{\boldsymbol{Z}'}\Big[\operatorname{diam}(R(\boldsymbol{Z}', \omega))^p\Big]\right)^{1/p} \gtrsim t \left(\tfrac{k}{s}\right)^{(1-\delta^*)\frac{K(\alpha)\pi}{d_Z}}\right) \lesssim \frac{1}{t^q} + \frac{1}{s}, \qquad \text{(S.26)}
$$

*for some $\delta^* \in (0,1)$ defined by (S.29), where $\mathbb{E}_{\boldsymbol{Z}'}(\cdot)$ denotes the expectation with respect $\boldsymbol{Z}'$ which equals in distribution $\boldsymbol{Z}$ but in independent of everything else;*

(d) $\mathbb{P}\left(\sup_{bsz\in[0,1]^d} \operatorname{diam}(R(\boldsymbol{z},\omega)) \geq 1\right) \to 1$ *for* $d_Z \geq 2$ *provided that* $k = o(s)$ *as* $s \to \infty$.

*Remark.* In the proof of part $(c)$ in Lemma 1, we show that with probability approaching 1, there exists a leaf $R_\ell$ that is never split along the variable $Z_j$, which in turn means that the diameter of this leaf is at least 1 with probability approaching 1. Hence, a uniform convergence for a tree grown by Algorithm 1 is not possible with high probability as long as it depends on the diameter to shrink to zero as $s \to \infty$.

**Proof** (of Lemma 1)

Let $M(\boldsymbol{z}) := \sum_{j\in[d_Z]} M_j(\boldsymbol{z})$ denote the total number of (random) splits to form the leaf $R(\boldsymbol{z},\omega)$. Since the variable to split is chosen at random, we have

$$M_j(\boldsymbol{z})|M(\boldsymbol{z}) \sim \operatorname{Binomial}(M(\boldsymbol{z}), q_j); \qquad j \in [d_Z],\ bsz \in [0,1]^{d_Z},$$

where $q_j$ is the probability of a split in $Z_j$ which is lower bounded by $\pi/d_Z > 0$ by assumption.

Let $N^{\mathcal{B}}(\boldsymbol{z})$ be the number of observations of the sample $\mathcal{B}$ on the $R(\boldsymbol{z},\omega)$. Then, by the minimum leaf size condition on the tree construction, we have

$$k \leq N^{\mathcal{B}}(\boldsymbol{z}) \leq 2k - 1; \qquad bsz \in [0,1]^{d_Z}.$$

Also, $N^{\mathcal{B}}(\boldsymbol{z}) = s \prod_{\ell=1}^{M(\boldsymbol{z})} \widetilde{\alpha}_\ell$ for some sequence $\{\widetilde{\alpha}_\ell\}$ such that $\alpha \leq \widetilde{\alpha}_\ell \leq (1-\alpha)$ for $\ell \in [M(\boldsymbol{z})]$ by the $\alpha$-regularity condition. Thus

$$k \leq s\alpha^{M(\boldsymbol{z})} \leq s(1-\alpha)^{M(\boldsymbol{z})} \leq 2k - 1; \qquad bsz \in [0,1]^{d_Z}.$$

Hence, from the last expressions, we obtain a lower bound on the number of splits of any leaf given by

$$\underline{M} := \frac{\log(s/(2k-1))}{\log(1/(1-\alpha))} \leq M(\boldsymbol{z}); \qquad ; \qquad bsz \in [0,1]^{d_Z}. \tag{S.27}$$

Recall the Chernoff's inequality. Let $S_m = W_1 + \cdots + W_m$ where $W_1,\ldots,W_m$ are independent and $W_j \in [0,1]$ for $j \in [m]$, then for every $t > 0$,

$$\mathbb{P}(S \leq \mathbb{E}[S] - t) \leq \exp(-t^2/2m). \tag{S.28}$$

Use (S.27) twice combined with (S.28) to obtain for $\delta \in (0, 1)$,

$$
\begin{aligned}
\mathbb{P}\left(M_{j,\ell} \leq \tfrac{(1-\delta)\pi}{d_Z}\underline{M}\right) &\leq \mathbb{E}\left[\mathbb{P}\left(M_{j,\ell} \leq \tfrac{(1-\delta)\pi}{d_Z}M_\ell \Big| M_\ell\right)\right] \\
&\leq \mathbb{E}\left[\exp\left(-\frac{\delta^2\pi^2 M_\ell}{2d_Z^2}\right)\right] \\
&\leq \exp\left(-\frac{\delta^2\pi^2\underline{M}}{2d_Z^2}\right) \\
&= \left(\frac{s}{2k-1}\right)^{-\frac{\delta^2\pi^2}{2d_Z^2\log(1/(1-\alpha))}},
\end{aligned}
$$

which concludes the proof of part $(a)$.

For $(b)$, we use (S.32) to write that on $\mathcal{E}_j$ and for $x > 0$,

$$
\mathbb{P}\left(\mathrm{diam}_j(R_\ell) \geq x | \mathcal{E}_j\right) \leq \mathbb{P}\left(M_{j,\ell} \leq \frac{-\log x}{(1+o(1))\log(1/(1-\alpha+o(1)))}\right).
$$

Set $\boldsymbol{X} = \left(\frac{2k-1}{s}\right)^{(1-\delta)\pi\frac{K(\alpha)}{d_Z}}$ and use the last inequality to write

$$
\begin{aligned}
\mathbb{P}\left(\mathrm{diam}_j(R_\ell) \geq \left(\tfrac{2k-1}{s}\right)^{(1-\delta)\pi\frac{K(\alpha)}{d_Z}} \Big| \mathcal{E}_j\right) &\leq \mathbb{P}\left(M_{j,\ell} \leq \frac{(1-\delta)\pi\underline{M}}{\mathbb{I}_1 d_Z} \Big| \mathcal{E}_j\right) \\
&= \mathbb{E}\left[\mathbb{P}\left(M_{j,\ell} \leq \frac{(1-\delta)\pi\underline{M}}{\mathbb{I}_1 d_Z} \Big| \mathcal{E}_j, M_\ell\right)\right] \\
&= \mathbb{E}\left[\mathbb{P}\left(M_{j,\ell} \leq \frac{(1-\delta)\pi\underline{M}}{\mathbb{I}_1 d_Z} \Big| M_\ell\right)\right] \\
&\leq \left(\frac{s}{2k-1}\right)^{-\frac{\delta^2\pi^2}{2\mathbb{I}_1(\alpha)^2 d_Z^2\log(1/\alpha)}},
\end{aligned}
$$

and from the union bound

$$
\mathbb{P}\left(\mathrm{diam}_j R((\boldsymbol{z},\omega)) \geq \left(\tfrac{s}{2k-1}\right)^{(1-\delta)\frac{K(\alpha,\pi)}{d_Z}}\right) \leq \left(\frac{s}{2k-1}\right)^{-\frac{\delta^2\pi^2}{2\mathbb{I}_1(\alpha)^2 d_Z^2\log(1/\alpha)}} + \frac{1}{s}.
$$

Therefore, $(b)$ follows from the union bound and the last expression because

$$
\begin{aligned}
\mathbb{P}\left(\mathrm{diam}(R(\boldsymbol{z},\omega)) \geq \sqrt{d_Z}\left(\tfrac{s}{2k-1}\right)^{(1-\delta)\frac{K(\alpha,\pi)}{d_Z}}\right) &\leq \mathbb{P}\left(\max_{j\in[d_Z]}\mathrm{diam}_j(R(\boldsymbol{z},\omega)) \geq \sqrt{d_Z}\left(\tfrac{s}{2k-1}\right)^{(1-\delta)\frac{K(\alpha,\pi)}{d_Z}}\right) \\
&\leq d_Z\max_{j\in[d_Z]}\mathbb{P}\left(\sup_{bsz\in[0,1]^d}\mathrm{diam}_j(R(\boldsymbol{z},\omega)) \geq \sqrt{d_Z}\left(\tfrac{s}{2k-1}\right)^{(1-\delta)\frac{K(\alpha,\pi)}{d_Z}}\right).
\end{aligned}
$$

For $(c)$, first note that, for $q \geq 1$ and $r > 0$,

$$
\begin{aligned}
\mathbb{E}\left[\mathrm{diam}(R(\boldsymbol{z},\omega))^q\right] &= \int_0^1 \mathbb{P}\left[\mathrm{diam}(R(\boldsymbol{z},\omega)) \geq t^{1/q}\right]dt \\
&= \int_0^{r^q} \mathbb{P}\left[\mathrm{diam}(R(\boldsymbol{z},\omega)) \geq t^{1/q}\right]dt + \int_{r^q}^1 \mathbb{P}\left[\mathrm{diam}(R(\boldsymbol{z},\omega)) \geq t^{1/q}\right]dt \\
&\leq r^q + \mathbb{P}\left[\mathrm{diam}(R(\boldsymbol{z},\omega)) \geq r\right].
\end{aligned}
$$

Set $r = \sqrt{d_Z} \left( \frac{2k-1}{s} \right)^{(1-\delta)\frac{K(\alpha)\pi}{d_Z}}$ and use (S.25) to write

$$\mathbb{E}\big[\mathrm{diam}(R(\boldsymbol{z},\omega))^q | \mathcal{E}\big] \lesssim \left( \frac{k}{s} \right)^{q(1-\delta)\frac{K(\alpha)\pi}{d_Z}} + \left( \frac{k}{s} \right)^{\frac{\delta^2 \pi^2}{2\mathbb{I}_1(\alpha)^2 d_Z^2 \log(1/\alpha)}}.$$

Note that $\delta \mapsto q(1-\delta)\frac{K(\alpha)\pi}{d_Z}$ is decreasing and vanishes as $\delta \uparrow 1$ and $\delta \mapsto \frac{\delta^2 \pi^2}{2\mathbb{I}_1(\alpha)^2 d_Z^2 \log(1/\alpha)}$ is non-negative increasing, so there is a $\delta^* \in (0,1)$ function of $q, \alpha, \pi, d_Z$ such that that those two functions agree, i.e.,

$$q(1-\delta^*)\frac{K(\alpha))}{d_Z} = \frac{\delta^{*2}\pi^2}{2\mathbb{I}_1(\alpha)^2 d_Z^2 \log(1/\alpha)}. \tag{S.29}$$

Therefore

$$\mathbb{E}\big[\mathrm{diam}(R(\boldsymbol{z},\omega))^q | \mathcal{E}\big] \lesssim \left( \frac{k}{s} \right)^{q(1-\delta^*)\frac{K(\alpha)\pi}{d_Z}}.$$

Set temporarily $D_p := \left( \mathbb{E}_{\boldsymbol{Z}'}\big[ \mathrm{diam}(\boldsymbol{Z}',\omega)^p \big] \right)^{1/p}$. For $1 \leq p \leq q < \infty$ and $x > 0$, we have, by Markov's inequality followed by Jensen's inequality and Fubini theorem,

$$\mathbb{P}(D_p \geq x | \mathcal{E}) = \mathbb{P}(D_p^q \geq \boldsymbol{X}^q | \mathcal{E}) \leq \frac{\mathbb{E}\big[(D_p^p)^{q/p}|\mathcal{E}\big]}{\boldsymbol{X}^q}$$

$$\leq \frac{\mathbb{E}\Big[\mathbb{E}_{\boldsymbol{Z}'}\big[\mathrm{diam}(R(\boldsymbol{Z}',\omega))^q\big]|\mathcal{E}\Big]}{\boldsymbol{X}^q} = \frac{\mathbb{E}_{\boldsymbol{Z}'}\Big[\mathbb{E}\big[\mathrm{diam}(R(\boldsymbol{Z}',\omega))^q|\mathcal{E}\big]\Big]}{\boldsymbol{X}^q}.$$

Set $x = t\sqrt{d_Z} \left( \frac{2k-1}{s} \right)^{(1-\delta^*)\frac{K(\alpha)\pi}{d_Z}}$ for $t > 0$ to conclude that for $1 \leq p \leq q < \infty$,

$$\mathbb{P}\left( D_p \gtrsim t \left( \tfrac{k}{s} \right)^{(1-\delta^*)\frac{K(\alpha)\pi}{d_Z}} \right) \leq \mathbb{P}(D_p \geq x | \mathcal{E}) + \mathbb{P}(\mathcal{E}^c) \lesssim \frac{1}{t^q} + \frac{d_Z}{s}.$$

which completes the proof of $(c)$.

Finally for $(d)$, let $M_j(\boldsymbol{z},\omega)$ denote the number of splits in $R(\boldsymbol{z},\omega)$ along the variable $j$. Now, if $\inf_{bsz\in[0,1]^d} M_j(\boldsymbol{z},\omega) = \min_{\ell\in[L]} M_{j,\ell} = 0$ for some $j \in [d_Z]$ then $\sup_{bsz\in[0,1]^d} \mathrm{diam}(R(\boldsymbol{z},\omega)) \geq 1$ thus

$$\mathbb{P}\left( \sup_{bsz\in[0,1]^d} \mathrm{diam}(R(\boldsymbol{z},\omega)) \geq 1 \right) \geq \mathbb{P}\left( \inf_{bsz\in[0,1]^d} M_{j,\ell} = 0 \right)$$

Also, note that $M_{j,\ell}$ and $M_{j,\ell'}$ are independent conditonally on $(M_\ell, M_{\ell'})$ for $\ell \neq \ell'$.

$$\mathbb{P}\left( \min_{\ell\in[L]} M_{j,\ell} = 0 \right) = 1 - \prod_{\ell=1}^{L} \big[1 - \mathbb{P}(M_{j,\ell} = 0)\big]$$

$$= 1 - \prod_{\ell=1}^{L} \big[1 - (1 - q_j)^{M_\ell}\big]$$

$$\geq 1 - \big[1 - (1 - \pi/d_Z)^{\overline{M}}\big]^L$$

$$= 1 - \left[1 - (s/k)^{\frac{\log(1-\pi/d_z)}{\log 1/\alpha}}\right]^L$$

49

Since $L \leq s/2k$ let $\beta := \frac{\log 1/\alpha}{\log(1/(1-\pi/d_z))}$ and note that since $1 - \pi/d_Z \geq 1 - 1/dz > \alpha$ for $d_Z \geq 2$ we have that $\beta > 1$. Set $v = (s/k)^{1/\beta} \to \infty$ then

$$\liminf_{s \to \infty} \mathbb{P}\left(\min_{\ell \in [L]} M_{j,\ell} = 0\right) \geq 1 - \lim_{s \to \infty} (1 - (s/k)^{-1/\beta})^{s/(2k)} = 1 - \lim_{v \to \infty} (1 - 1/v)^{v^\beta/2}$$

Since $1 - x \leq \exp(-x)$ we have

$$\liminf_{s \to \infty} \mathbb{P}\left(\min_{\ell \in [L]} M_{j,\ell} = 0\right) \geq 1 - \lim_{v \to \infty} \exp(-v^{(\beta-1)/2}) = 1,$$

which completes the proof of part $(d)$. $\qquad \square$

**Lemma 2** (Lower bounds on the number of observations in the leaves)

*Let $L$ denote the number of leaves in a tree constructed according to Algorithm 1 and $N :=$ $(N_1, \ldots, N_L)$ where $N_\ell$ is the number of observation of the sample $\mathcal{A}$ the $\ell$-th leave for $\ell \in [L]$. Then for $k \in \mathbb{N}$*

(a) $\liminf_{s \to \infty} \mathbb{P}(N_\ell = 0) \geq e^{-2\overline{f}(2k-1)}$ *for $\ell \in [L]$;*

(b) $\lim_{s \to \infty} \mathbb{P}(\min_{\ell \in [L]} N_\ell = 0) = 1$;

*Also,*

(c) *If $k \gtrsim s^\epsilon$ for some $\epsilon \in (0,1)$ then $k(k/s)^\epsilon \lesssim_\mathbb{P} N_\ell \lesssim_\mathbb{P} k(k/s)^{-\epsilon}$ for $\ell \in [L]$;*

(d) *If $k \gtrsim s^{1/2+\epsilon}$ for some $\epsilon \in (0, 1/2)$ then $\min_{\ell \in [L]} N_\ell \gtrsim_\mathbb{P} k(k/s)^\epsilon$.*

*Let $K(\alpha) := \frac{\log((1-\alpha)^{-1})}{\log(1/\alpha)}$ for $\alpha \in (0, 0.5)$ and $\pi \in (0, 1]$. For $k \in [s]$ and $\delta \in (0, 1)$ and*

$$\mathbb{P}\left(|\mathcal{A}(z, \omega)| \leq (1 - \delta)\underline{f}s\left(\frac{s}{k}\right)^{-\frac{\mathbb{I}_2(\alpha)}{K(\alpha)}}\right) \leq \frac{1}{\delta^2 \underline{f}s(s/k)^{-\frac{\mathbb{I}_2(\alpha)}{K(\alpha)}}}; \tag{S.30}$$

$$\mathbb{P}\left(|\mathcal{A}(z, \omega)| \geq (1 + \delta)s\overline{f}\left(\frac{s}{2k-1}\right)^{-\mathbb{I}_1(\alpha)K(\alpha)}\right) \leq \frac{1 + \delta}{\delta^2 s\underline{f}\left(\frac{s}{k}\right)^{-\frac{\mathbb{I}_2(\alpha)}{K(\alpha)}}}, \tag{S.31}$$

*where $\underline{f}, \overline{f}$ are the lower and upper bound on the density of $\mathbf{Z}$ respectively, $\mathbb{I}_1(\alpha) := (1+o(1))\frac{\log(1/(1-\alpha+o(1)))}{\log(1/(1-\alpha))} \to 1$, $\mathbb{I}_2(\alpha) := (1 + o(1))\frac{\log(1/(\alpha+o(1)))}{\log(1/\alpha)} \to 1$.*

*In particular, if $s \to \infty$ and $k \asymp s^\eta$ for $\eta \in (1 - K(\alpha, \pi), 1) \subseteq (0, 1)$ then the right-hand side of (S.30)-(S.31) vanishes and we conclude*

$$\mathrm{diam}(R(z, \omega)) \lesssim_\mathbb{P} s^{-\frac{(1-\eta)\pi K(\alpha)}{2d_Z}} \quad and \quad s^{1-\frac{1-\eta}{K(\alpha)}} \lesssim_\mathbb{P} |\mathcal{A}(z, \omega)| \lesssim_\mathbb{P} s^{1-(1-\eta)K(\alpha)}.$$

*Remark.* From parts (a) and (b) in Lemma 2, we conclude that every leaf from a tree grown by Algorithm 1 with a minimum leaf parameter $k \in \mathbb{N}$ is empty with probability at least $e^{-2(2k-1)}$ as $s \to \infty$. Therefore, to ensure that the leaves have a minimum number of observations with high probability, we need to have $k \to \infty$ as $s \to \infty$. In fact, by ignoring the $o(1)$ terms in (S.32), we have that $N_\ell \asymp k$ as $k \to \infty$ and $s \to \infty$.

**Proof** (of Lemma 2)

Let $L$ denote the number of leaves in a tree and $N^\mathcal{A} := (N_1^\mathcal{A}, \ldots, N_L^\mathcal{A})$ where $N_\ell^\mathcal{A}$ is the number of observation of the sample $\mathcal{A}$ on the $\ell$-th leaf for $\ell \in [L]$. Also let $p = (p_1, \ldots, p_L)$ where $p_\ell := \mathbb{P}(\boldsymbol{Z} \in R_\ell)$. Then

$$N^\mathcal{A} \sim \text{Multinomial}_L(s, p); \qquad N_\ell^\mathcal{A} \sim \text{Binomial}(s, p_\ell) \quad \ell \in [L].$$

Let $F_{n', p'}$ denote the cdf of a Binomial distribution with $n'$ trials and probability of success $p'$. Then

$$\mathbb{P}(N_\ell^\mathcal{A} \le x) = F_{s, p_\ell}(\boldsymbol{x}); \quad x \in \mathbb{R}, \ \ell \in [L].$$

Also, since the density of $\boldsymbol{Z}$ is bounded away from zero and infinity, we have that

$$\underline{f} \prod_{j=1}^{d_Z} \text{diam}_j(R_\ell) = \underline{f}\mu(R_\ell) \le p_\ell \le \overline{f}\mu(R_\ell) = \overline{f} \prod_{j=1}^{d_Z} \text{diam}_j(R_\ell); \qquad \ell \in [L].$$

Let $M_j(\boldsymbol{z}, \omega)$ denote the number of splits in $Z_j$ forming $R(\boldsymbol{z}, \omega)$ for $j \in [d_Z]$. By $\alpha$-regularity we have, from Lemma 12 and 13 in Wager and Walther (2016), that $\mathbb{P}(\mathcal{E}_j) \ge 1 - 1/s$ for $j \in [d_Z]$ where

$$\mathcal{E}_j := \left\{ (\alpha + o(1))^{(1+o(1))M_j(\boldsymbol{z}, \omega)} \le \text{diam}_j(R(\boldsymbol{z}, \omega)) \le (1 - \alpha + o(1))^{(1+o(1))M_j(\boldsymbol{z}, \omega)} : bsz \in [0, 1]^d \right\}.$$
(S.32)

Then, on $\mathcal{E} := \bigcap_j \mathcal{E}_j$,

$$\alpha^{(1+\zeta_1)M} \le \mu(R_\ell) \le (1 - \alpha)^{(1+\zeta_2)M}$$

where $M := M(\boldsymbol{z}, \omega) := \sum_{j=1}^{d_Z} M_j(\boldsymbol{z}, \omega)$ is the total number of splits forming $R(\boldsymbol{z}, \omega)$ and

$$\zeta_1 := (1 + o(1))\frac{\log(\alpha + o(1))}{\log(\alpha)} - 1 = o(1); \qquad \zeta_2 := (1 + o(1))\frac{\log(1 - \alpha + o(1))}{\log(1 - \alpha)} - 1 = o(1).$$

Also, by $\alpha$-regularity

$$k \le s\alpha^M \le s(1 - \alpha)^M \le 2k - 1.$$

Then, on $\mathcal{E}$ which occurs with probability at least $1 - d_Z/s$,

$$\underline{f}(k/s)^{1+\zeta_1} \le p_\ell \le \overline{f}((2k-1)/s)^{1+\zeta_2}; \qquad \ell \in [L].$$
(S.33)

Recall that the Binomial distribution is decreasing in $p$, in the sense that $F_{n,p} \leq F_{n,p'}$ pointwise for $p' \leq p$. Then, using the inequality above

$$F_{s, \frac{\overline{f}(2k-1)((2k-1)/s)^{o(1)}}{s}} \leq F_{s,p_\ell} \leq F_{s, \frac{\underline{f}k(k/s)^{o(1)}}{s}}.$$

Let $G_\Lambda$ denote the Poisson cdf with mean $\lambda$. Note that $0 \leq k/s \leq 1$, then if $\underline{\Lambda} := \liminf_{s\to\infty}(k/s)^{o(1)} > 0$ we have, by the pointwise limit $F_{n',p'} \to G_{\lambda'}$ as $n'p' \to \lambda'$,

$$\limsup_{s\to\infty} F_{s, \frac{\underline{f}k(k/s)^{o(1)}}{s}} \leq G_{\underline{\Lambda}\underline{f}k}.$$

Othewise if $\underline{\Lambda} = 0$ we get a trivial bound $\limsup_{s\to\infty} F_{s, \frac{\underline{f}k(k/s)^{o(1)}}{s}} \leq 1$. Similarly, we have $0 \leq (2k-1)/s \leq 2$ and if $\overline{\Lambda} := \limsup_{s\to\infty}((2k-1)/s)^{o(1)} > 0$ we have

$$\liminf_{s\to\infty} F_{s, \frac{\overline{f}(2k-1)((2k-1)/s)^{o(1)}}{s}} \geq G_{\overline{\Lambda}\overline{f}(2k-1)}.$$

Othewise when $\overline{\Lambda} = 0$ we get a *non* trivial bound $\liminf_{s\to\infty} F_{s, \frac{\overline{f}(2k-1)((2k-1)/s)^{o(1)}}{s}} \geq 1$. Therefore, defining pointwise $G_0(m) := \lim_{\lambda\downarrow 0} G_\lambda(m) = \mathbf{1}\{m \geq 0\}$, we have shown that for $k \in \mathbb{N}$ and $\ell \in [L]$:

$$G_{\overline{f\Lambda}(2k-1)}(\cdot) \leq \liminf_{s\to\infty} \mathbb{P}(N_\ell \leq \cdot) \leq \limsup_{s\to\infty} \mathbb{P}(N_\ell \leq \cdot) \leq G_{\underline{f\Lambda}k}(\cdot).$$

We then obtain the result $(a)$ by evaluating the last inequality at 0 and using the fact that $\overline{\Lambda} \leq 2$ to conclude

$$\liminf_{s\to\infty} \mathbb{P}(N_\ell = 0) \geq G_{\overline{f\Lambda}(2k-1)}(0) \geq G_{2\overline{f}(2k-1)}(0) = \frac{1}{e^{2\overline{f}(2k-1)}}. \tag{S.34}$$

For $(b)$, suppose that there are only two leaves $L = 2$ indexed by $\ell, \ell'$, then

$$\begin{aligned}
\mathbb{P}(N_\ell \geq k, N_{\ell'} \geq k) &= \mathbb{P}(N_\ell \geq k | N_{\ell'} \geq k)\mathbb{P}(N_{\ell'} \geq k) \\
&= \mathbb{P}(N_\ell \geq k | N_\ell \leq s - k)\mathbb{P}(N_{\ell'} \geq k) \\
&= \mathbb{P}(k \leq N_\ell \leq s - k)\mathbb{P}(N_\ell \leq s - k)\mathbb{P}(N_{\ell'} \geq k) \\
&\leq \mathbb{P}(N_\ell \geq k)\mathbb{P}(N_{\ell'} \geq k),
\end{aligned}$$

where we use the fact that $N_\ell + N_{\ell'} = s$. Applying induction, it is easy to verify that

$$\mathbb{P}(N_1 \geq k, \ldots, N_L \geq k) \leq \prod_{l=1}^{L} \mathbb{P}(N_\ell \geq k); \quad k \in \mathbb{N}_0. \tag{S.35}$$

Let $N_{\min} := \min_{\ell\in[L]} N_\ell$. Use (S.35) with $k = 1$, the fact that $1 - x \leq e^{-x}$ for $x \in \mathbb{R}$ and part $(a)$

to write for $k \in \mathbb{N}$

$$\mathbb{P}(N_{\min} = 0) = 1 - \mathbb{P}(N_1 \geq 1, \ldots, N_L \geq 1)$$

$$\geq 1 - \prod_{l=1}^{L} \mathbb{P}(N_\ell \geq 1)$$

$$= 1 - \prod_{l=1}^{L} \left[ 1 - \mathbb{P}(N_\ell = 0) \right]$$

$$\geq 1 - \exp \left[ - \sum_{l=1}^{L} \mathbb{P}(N_\ell = 0) \right]$$

$$\geq 1 - \exp \left[ - L(e^{-2\overline{f}(2k-1)} + o(1)) \right] \to 1,$$

because $L \asymp s/k \to \infty$ for each $k \in \mathbb{N}$ as $s \to \infty$, which concludes the proof of part $(b)$. Note that if $k \geq \log s/4$ then $\mathbb{P}(N_{\min} = 0) \not\to 1$.

For $(c)$, we have by Markov's inequality $N_\ell \lesssim_{\mathbb{P}} \mathbb{E}[N_\ell] = sp_\ell$ hence $N_\ell \lesssim_{\mathbb{P}} sp_\ell$. For the other direction, we apply tail bounds for the binomial distribution. Refer to (Feller, 1957, pg 151). Let $S_n \sim \text{Binom}(m, p)$, then

$$\mathbb{P}(S \leq r) \leq \frac{(m-r)p}{(mp-r)^2}; \qquad r \leq mp \tag{S.36}$$

$$\mathbb{P}(S \geq r) \leq \frac{r(1-p)}{(mp-r)^2}; \qquad r \geq mp. \tag{S.37}$$

Then using (S.36), we write for some $C > 1$

$$\mathbb{P}(N_\ell \leq sp_\ell/C) \leq \frac{(s - sp_\ell/C)p_\ell}{(1 - 1/C)^2 (sp_\ell)^2} \leq \frac{sp_\ell}{(1 - 1/C)^2 (sp_\ell)^2} = \frac{1}{(1 - 1/C)^2 sp_\ell}.$$

Also, from (S.33) we have for $\epsilon > 0$

$$k \left( \frac{k}{s} \right)^\epsilon \lesssim k \left( \frac{k}{s} \right)^{\zeta_1} \lesssim \underline{f} s \left( \frac{k}{s} \right)^{1+\zeta_1} \leq sp_\ell \leq \overline{f} s \left( \frac{2k-1}{s} \right)^{1+\zeta_2} \lesssim k \left( \frac{2k-1}{s} \right)^{\zeta_2} \lesssim k \left( \frac{k}{s} \right)^{-\epsilon} \tag{S.38}$$

because $\zeta_1 = o(1)$ and $\zeta_2 = o(1)$. By assumption $k \gtrsim s^\epsilon$ then the left-hand side is at least of order $k^\epsilon \to \infty$ as $s \to \infty$ which ensures that $sp_\ell \to \infty$. Then $N_\ell \gtrsim_{\mathbb{P}} sp_\ell$ and $N \asymp_{\mathbb{P}} sp_\ell$. Therefore, from (S.33) we have for $\epsilon > 0$

$$k \left( \frac{k}{s} \right)^\epsilon \lesssim_{\mathbb{P}} N_\ell \lesssim_{\mathbb{P}} k \left( \frac{k}{s} \right)^{-\epsilon},$$

which concludes the proof of part $(c)$.

For $(d)$, by the union bound, followed by (S.36), we have for $m < sp_{\min}$ where $p_{\min} := \min_{\ell \in [L]} p_\ell$,

$$\mathbb{P}(N_{\min} \leq m) = \mathbb{P} \left( \bigcup_{\ell=1}^{L} \{ N_\ell \leq m \} \right) \leq \sum_{\ell=1}^{L} \mathbb{P}(X_\ell \leq m) \leq \sum_{\ell=1}^{L} \frac{(s-m)p_\ell}{(sp_\ell - m)^2} \leq \frac{(s-m)}{(sp_{\min} - m)^2}.$$

Set $m = sp_{\min}/C$ for a constant $C > 1$ and use (S.38) to write

$$\mathbb{P}(N_{\min} \le sp_{\min}/C) \le \frac{s - sp_{\min}/C}{(1-1/C)^2 sp_{\min}} \le \frac{s}{(1-1/C)^2 (sp_{\min})^2} \lesssim \frac{1}{(1-1/C)^2 K^{2\epsilon}},$$

where we use the assumption that $k \gtrsim s^{1/2+\epsilon}$. Then $N_{\min} \gtrsim_{\mathbb{P}} sp_{\min} \gtrsim_{\mathbb{P}} k(k/s)^\epsilon$ which demostrate part $(d)$ and concludes the proof of the lemma.

$\square$

**Lemma 3** (Tree Variance-type bound)

*Let $\{W_i : i \in [n]\}$ be $n$ independent copies of the random variable $W$ and define the tree $T_W(\boldsymbol{z}, \omega) := \frac{1}{|\mathcal{A}(\boldsymbol{z},\omega)|} \sum_{i \in \mathcal{A}(\boldsymbol{z},\omega)} W_i$ and the map $V(\boldsymbol{z}) := T_W(\boldsymbol{z}, \omega) - \mathbb{E}[W | \boldsymbol{Z} = \boldsymbol{z}]$ for $\boldsymbol{z} \in [0,1]^{d_Z}$. If $\boldsymbol{z} \mapsto \mathbb{E}[|W|^q | \boldsymbol{Z} = \boldsymbol{z}] \le M_q < \infty$ for some $q \ge 2$ and constant $M_q \ge 0$ and $k \gtrsim s^\epsilon$ for some $\epsilon \in (0,1)$ then for $t > 0$, as $s \to \infty$,*

$$\mathbb{P}\left(|V(\boldsymbol{z})| \gtrsim t k^{-\frac{1}{2}} (s/k)^{\epsilon/2}\right) \lesssim t^{-q} + s^{-1}, \qquad \boldsymbol{z} \in [0,1]^d.$$

*Also, let $\Delta_p := \left(\mathbb{E}_{\boldsymbol{Z}'}\left[V(\boldsymbol{Z}')^p\right]\right)^{1/p}$ for $p \in [2,\infty)$ where $\boldsymbol{Z}'$ and independent copy of $\boldsymbol{Z}$ and $\Delta_\infty := \sup_{bs\boldsymbol{z} \in [0,1]^{d_Z}} V(\boldsymbol{z})$, then*

$$\mathbb{P}\left(\Delta_p \gtrsim t k^{-\frac{1}{2}} (s/k)^{\epsilon/2}\right) \lesssim t^{-q} + s^{-1}; \qquad 2 \le p \le q < \infty,$$
$$\mathbb{P}\left(\Delta_\infty \gtrsim t k^{-\frac{1}{2}} (s/k)^{1/q+\epsilon/2}\right) \lesssim t^{-q} + k^{-1}.$$

**Proof** (Lemma 3)

It is convenient re-write $T_W(\boldsymbol{z}, \omega)$ for a subsample of size $|\mathcal{S}| =: s \le n$ as

$$T_W(\boldsymbol{z}, \omega) = \sum_{i=1}^s S_i W_i; \qquad S_i := S_i(\boldsymbol{z}, \omega) := \begin{cases} |\{j \in \mathcal{A} : Z_j \in R(\boldsymbol{z}, \omega)\}|^{-1} & \text{if } \boldsymbol{Z}_i \in R(\boldsymbol{z}, \omega) \\ 0 & \text{otherwise.} \end{cases}$$

First we note that $T_W(\boldsymbol{z}, \omega)$ is unbiased for $\mathbb{E}[W | \boldsymbol{Z} \in R(\boldsymbol{z}, \omega)]$. In fact, $T_W(\boldsymbol{z}, \omega)$ is even conditional (on $\mathbb{S} := (S_1, \ldots, S_d)$) unbiased since

$$\mathbb{E}[T_W(\boldsymbol{z}, \omega) | \mathbb{S}] = \left[\sum_{i=1}^s S_i \mathbb{E}[W_i | S_i]\right] = \mathbb{E}[W | \boldsymbol{Z} \in R(\boldsymbol{z}, \omega)] \sum_{i=1}^s S_i = \mathbb{E}[W | \boldsymbol{Z} \in R(\boldsymbol{z}, \omega)],$$

where we use the fact that, due to the sample split, $S_i \mathbb{E}[W_i | S_i] = |\mathcal{A}(\boldsymbol{z}, \omega)|^{-1} \mathbb{E}[W | \boldsymbol{Z} \in R(\boldsymbol{z}, \omega)]$ if $\boldsymbol{Z}_i \in R(\boldsymbol{z}, \omega)$ and 0 otherwise, and $\sum_{i=1}^s S_i = 1$. Then $\mathbb{E}[T_W(\boldsymbol{z}, \omega)] = \mathbb{E}_{\mathbb{S}}\left[\mathbb{E}[T_W(\boldsymbol{z}, \omega) | \mathbb{S}]\right] = \mathbb{E}[W | \boldsymbol{Z} \in R(\boldsymbol{z}, \omega)]$.

For convenience define $U_i := W_i - \mathbb{E}[W|\mathbf{Z} \in R(\mathbf{z}, \omega)]$ for $i \in [s]$ and note that the sequence $\{S_i U_i : i \in [s]\}$ is zero mean and independent conditional on $\mathbb{S}$. Then, by Marcinkiewicz–Zygmund inequality followed by Jensen's inequality, we have, for $q \geq 2$,

$$\mathbb{E}\left[\left|\sqrt{|\mathcal{A}(\mathbf{z}, \omega)|} \sum_{i=1}^{s} S_i U_i\right|^q \Bigg| \mathbb{S}\right] \leq C_q \mathbb{E}\left[\left|\sum_{i=1}^{s} |\mathcal{A}(\mathbf{z}, \omega)| S_i^2 U_i^2\right|^{q/2} \Bigg| \mathbb{S}\right]$$

$$= C_q \mathbb{E}\left[\left|\sum_{i=1}^{s} S_i U_i^2\right|^{q/2} \Bigg| \mathbb{S}\right]$$

$$\leq C_q \left(\sum_{i=1}^{s} S_i^{q/2} \mathbb{E}[|U_i|^q | S_i]\right)$$

$$\leq C_q \left(\sum_{i=1}^{s} S_i \mathbb{E}[|U_i|^q | S_i]\right) = C_q \mathbb{E}[|U|^q | \mathbf{Z} \in R(\mathbf{z}, \omega)],$$

where $C_q$ is constant only depending on $q$. We also use $|\mathcal{A}(\mathbf{z}, \omega)|^{k-1} S_i^k = S_i$ and $S_i^k \leq S_i$ for $k \geq 1$. Hence

$$\mathbb{E}\left[\left|\sum_{i=1}^{s} S_i U_i\right|^q \Bigg| \mathbb{S}\right] \lesssim \frac{1}{|\mathcal{A}(\mathbf{z}, \omega)|^{q/2} \vee 1} \mathbb{E}\left[\left|\sqrt{|\mathcal{A}(\mathbf{z}, \omega)|} \sum_{i=1}^{s} S_i U_i\right|^q \Bigg| \mathbb{S}\right]$$

$$\lesssim \frac{\mathbb{E}[|U|^q | \mathbf{Z} \in R(\mathbf{z}, \omega)]}{|\mathcal{A}(\mathbf{z}, \omega)|^{q/2} \vee 1}.$$

From Cribari-Neto et al. (2000) we have that

$$\mathbb{E}[(1 + B)^{-\alpha}] \lesssim (np)^{-\alpha}; \qquad B \sim \text{Binomial}(n, p), \quad \alpha \in \mathbb{R}.$$

Recall that $|\mathcal{A}(\mathbf{z}, \omega)| \sim \text{Binomial}(s, p(\mathbf{z}, \omega))$ conditional on $R(\mathbf{z}, \omega)$ and $p(\mathbf{z}, \omega)$ is bounded by below on $\mathcal{E}$ as per (S.32). Then conditional on $R(\mathbf{z}, \omega)$ and $\mathcal{E}$, for $\alpha \in R$, we have

$$\mathbb{E}\left[\frac{1}{|\mathcal{A}(\mathbf{z}, \omega)|^{q/2} \vee 1}\right] \leq 2\mathbb{E}\left[\frac{1}{(1 + |\mathcal{A}(\mathbf{z}, \omega)|)^{q/2}}\right] \lesssim \frac{1}{(sp(\mathbf{z}, \omega))^{q/2}} \leq \frac{1}{(s\underline{f}(k/s)^{1+\zeta_1})^{q/2}}$$

Therefore

$$\mathbb{E}\left[\left|\sum_{i=1}^{s} S_i U_i\right|^q \Bigg| \mathcal{E}\right] = \mathbb{E}\left\{\mathbb{E}\left[\left|\sum_{i=1}^{s} S_i U_i\right|^q \Bigg| \mathbb{S}, \mathcal{E}\right] \Bigg| \mathcal{E}\right\}$$

$$\lesssim \mathbb{E}\left\{\frac{\mathbb{E}[|U|^q | \mathbf{Z} \in R(\mathbf{z}, \omega)]}{|\mathcal{A}(\mathbf{z}, \omega)|^{q/2} \vee 1} \Bigg| \mathcal{E}\right\}$$

$$= \mathbb{E}\left\{\mathbb{E}[|U|^q | \mathbf{Z} \in R(\mathbf{z}, \omega)] \mathbb{E}\left[\frac{1}{|\mathcal{A}(\mathbf{z}, \omega)|^{q/2} \vee 1} \Bigg| \mathcal{E}, R(\mathbf{z}, \omega)\right] \Bigg| \mathcal{E}\right\}$$

$$\lesssim \frac{\mathbb{E}[|U|^q | \mathbf{Z} \in R(\mathbf{z}, \omega) | \mathcal{E}]}{(s\underline{f}(k/s)^{1+\zeta_1})^{q/2}}$$

$$\leq M_q \left(\frac{1}{(s\underline{f}(k/s)^{1+\zeta_1})^{q/2}}\right).$$

Finally, by Markov inequality and the last bound, we have

$$\mathbb{P}\left(V(\boldsymbol{z}) \gtrsim \frac{t}{\sqrt{s(k/s)^{1+\zeta_1}}}\right) \leq \mathbb{P}(V(\boldsymbol{z}) \geq x | \mathcal{E}) + \mathbb{P}(\mathcal{E}^c) \lesssim \frac{1}{t^q} + \frac{d_Z}{s},$$

which demonstrates the first result.

For the second result in the case $2 \leq q < \infty$, we have, for $p \leq q$ and $x > 0$, by Markov's inequality followed by Jensen's inequality and Fubini theorem.

$$\mathbb{P}(\Delta_p \geq x | \mathcal{E}) = \mathbb{P}(\Delta_p^q \geq X^q | \mathcal{E}) \leq \frac{\mathbb{E}\left[(\Delta_p^p)^{q/p} | \mathcal{E}\right]}{X^q} \leq \frac{\mathbb{E}\left[\mathbb{E}_{\boldsymbol{Z}'}\left[V(\boldsymbol{Z}')^q\right] | \mathcal{E}\right]}{X^q} = \frac{\mathbb{E}_{\boldsymbol{Z}'}\left[\mathbb{E}\left[V(\boldsymbol{Z}')^q | \mathcal{E}\right]\right]}{X^q}.$$

Set $x = t(s\underline{f}(k/s)^{1+\zeta_1})^{-1/2}$ for $t > 0$ to conclude that for $2 \leq p \leq q < \infty$,

$$\mathbb{P}\left(\Delta_p \gtrsim \frac{\delta}{\sqrt{s(k/s)^{1+\zeta_1}}}\right) \leq \mathbb{P}(\Delta_p \geq x | \mathcal{E}) + \mathbb{P}(\mathcal{E}^c) \lesssim \frac{1}{t^q} + \frac{d_Z}{s}.$$

For the case $p = \infty$, we have that the number of leaves is bounded by $s/k$ and then by the union bound for $x = t(s/k)^{1/q} k^{-\frac{1}{2}}(s/k)^{\epsilon/2}$

$$\mathbb{P}\left(\sup_{bsz \in [0,1]^d} |V(\boldsymbol{z})| \geq x\right) \leq \frac{s}{k} \sup_{bsz \in [0,1]^d} \mathbb{P}\left(V(\boldsymbol{z}) \geq x\right) \lesssim \frac{s}{k}\left(\frac{1}{(s/k)t^q} + \frac{d_Z}{s}\right) = \frac{1}{t^q} + \frac{d_Z}{k}$$

$\square$

**Lemma 4** (Tree Bias-type bound)

*Consider the same setup of Lemma 3 and define $B(\boldsymbol{z}) := \mathbb{E}[W | \boldsymbol{Z} \in R(\boldsymbol{z}, \omega)] - \mathbb{E}[W | \boldsymbol{Z} = \boldsymbol{z}]$ for $\boldsymbol{z} \in [0,1]^{d_Z}$. If $\boldsymbol{z} \mapsto \mathbb{E}[W | \boldsymbol{Z} = \boldsymbol{z}]$ is Lipschitz then for $\delta \in (0,1)$,*

$$\mathbb{P}\left(|B(\boldsymbol{z})| \gtrsim \left(\frac{k}{s}\right)^{(1-\delta)\frac{K(\alpha)\pi}{d_Z}}\right) \lesssim \left(\frac{k}{s}\right)^{\frac{\delta^2 \pi^2}{2(1+o(1))^2 d_Z^2 \log(1/\alpha)}} + \frac{1}{s}. \tag{S.39}$$

*Also, let $\Pi_p := \left(\mathbb{E}_{\boldsymbol{Z}'}\left[B(\boldsymbol{Z}')^p\right]\right)^{1/p}$ for $p \in [2, \infty)$ where $\boldsymbol{Z}'$ and independent copy of $\boldsymbol{Z}$ and then*

$$\mathbb{P}\left(\Pi_p \gtrsim t\left(\frac{k}{s}\right)^{(1-\delta^*)\frac{K(\alpha)\pi}{d_Z}}\right) \lesssim \frac{1}{t^q} + \frac{1}{s}; \qquad 1 \leq p \leq q < \infty. \tag{S.40}$$

**Proof** (Lemma 4)

From the Lipschitz condition and Lemma 2, we have as $s \to \infty$

$$|B(\boldsymbol{z})| \lesssim \sup_{\boldsymbol{z}, \boldsymbol{z}' \in R(\boldsymbol{z}, \omega)} \|\boldsymbol{z} - \boldsymbol{z}'\| =: \text{diam}(R(\boldsymbol{z}, \omega))$$

Use (S.25) to upper bound the leaf diameter for each $\boldsymbol{z} \in [0,1]^d$ and obtain (S.39). Similarly, use (S.26) to upper bound the leaf diameter for each $\boldsymbol{z} \in [0,1]^d$ and obtain (S.40). $\square$

**Lemma 5** (Tree Rate of Convergence)

*Let $\{W_i : i \in [n]\}$ be $n$ independent copies of the random variable $W$ and define the tree $T_W(\boldsymbol{z}, \omega) :=$*
$\frac{1}{|\mathcal{A}(\boldsymbol{z}, \omega)|} \sum_{i \in \mathcal{A}(\boldsymbol{z}, \omega)} W_i$. *If $\boldsymbol{z} \mapsto T_0(\boldsymbol{z}) := \mathbb{E}[W | \boldsymbol{Z} = \boldsymbol{z}]$ is Lipschitz, $\boldsymbol{z} \mapsto \mathbb{E}[|W|^q | \boldsymbol{Z} = \boldsymbol{z}] \leq M_q < \infty$*
*for some $q \geq 2$ and constant $M_q \geq 0$ and $k \gtrsim s^\epsilon$ for some $\epsilon \in (0, 1)$ then for $\delta \in (0, 1)$, as $s \to \infty$,*

$$|T_W(\boldsymbol{z}, \omega) - T_0(\boldsymbol{z})| \lesssim_{\mathbb{P}} k^{-\frac{1}{2}} \left(\frac{s}{k}\right)^{\epsilon/2} + \left(\frac{k}{s}\right)^{(1-\delta)\frac{K(\alpha)\pi}{d_Z}} ; \qquad \boldsymbol{z} \in [0, 1]^d.$$

*Also, for $2 \leq p \leq q$,*

$$\left[\int_{[0,1]^d} |T_W(\boldsymbol{z}, \omega) - T_0(\boldsymbol{z})|^p f(\boldsymbol{z}) d\boldsymbol{z}\right]^{1/p} \lesssim_{\mathbb{P}} k^{-\frac{1}{2}} \left(\frac{s}{k}\right)^{\epsilon/2} + \left(\frac{k}{s}\right)^{(1-\delta^*)\frac{K(\alpha)\pi}{d_Z}},$$

*where $\delta^* \in (0, 1)$ is defined in Lemma 4.*

**Proof** (Lemma 5)

By the triangle inequality combined with Lemma 3 and 4 we obtain the first result since

$$|T(\boldsymbol{z}, \omega) - \mathbb{E}[W | \boldsymbol{Z} = \boldsymbol{z}]| \leq |T(\boldsymbol{z}, \omega) - \mathbb{E}[W | \boldsymbol{Z} \in R(\boldsymbol{z}, \omega)]| + |\mathbb{E}[W | \boldsymbol{Z} \in R(\boldsymbol{z}, \omega) - \mathbb{E}[W | \boldsymbol{Z} = \boldsymbol{z}]|$$

$$\lesssim_{\mathbb{P}} k^{-\frac{1}{2}} \left(\frac{s}{k}\right)^{\epsilon/2} + \left(\frac{k}{s}\right)^{(1-\delta)\frac{K(\alpha)\pi}{d_Z}}.$$

Similarly, the second result follows from the triangle inequality, Lemma 3 and 4 with $p = q$. $\qquad \square$

**Lemma 6**

*Let $\{V_i := (W_i, \boldsymbol{Z}_i) : i \in [n]\}$ be $n$ independent copies of the random vector $V := (W, \boldsymbol{Z})$ where $W$*
*is a random variable . If $\boldsymbol{z} \mapsto \mathbb{E}[W | \boldsymbol{Z} = \boldsymbol{z}]$ and $\boldsymbol{z} \mapsto \mathbb{E}[W^2 | \boldsymbol{Z} = \boldsymbol{z}]$ are Lipschitz, $\mathbb{V}[W | \boldsymbol{Z} = \boldsymbol{z}] \geq 0$*
*and $k \asymp s^\eta$ for $\eta \in (1 - K(\alpha, \pi), 1) \subseteq (0, 1)$ then, if $s \to \infty$ and $s(\log n)^{d_Z} = o(n)$*

$$\frac{\overline{T}(\boldsymbol{z}) - \mathbb{E}[\overline{T}(\boldsymbol{z})]}{\lambda(\boldsymbol{z})} \xrightarrow{d} \mathsf{N}(0, 1) \quad and \quad \frac{\mathbb{V}[\overline{T}(\boldsymbol{z})]}{\lambda^2(\boldsymbol{z})} \to 1, \tag{S.41}$$

*where $\overline{T}(\boldsymbol{z}) := \frac{1}{B} \sum_{b=1}^B T_W(\boldsymbol{z}, \omega_b)$, $T_W(\boldsymbol{z}, \omega_b) := \frac{1}{|\mathcal{A}(\boldsymbol{z}, \omega_b)|} \sum_{i \in \mathcal{A}(\boldsymbol{z}, \omega_b)} W_i$, $\{\omega_b : b \in [B]\}$ is indepen-*
*dent of $\{V_i : i \in [n]\}$ and $\lambda^2(\boldsymbol{z})$ is the variance of the Háyek Projection of $T_W(\boldsymbol{z}, \omega_b)$, which can be*
*lower bounded as*

$$\frac{s^{\frac{1-\eta}{K(\alpha)}}}{n(\log s)^{d_Z}} \lesssim \lambda^2(\boldsymbol{z}). \tag{S.42}$$

**Proof** (Lemma 6)

We start by writing $\overline{T}(\boldsymbol{z})$ as a generalized U-statistics following Peng et al. (2022). For a fixed
$\boldsymbol{z} \in [0, 1]^d$ define the randomized symmetric (in its $s$ arguments) kernel $h$ by,

$$h(v_1, \ldots, v_s) = \sum_{i=1}^s \chi_i w_i; \qquad \chi_i(\boldsymbol{z}, \xi) := \begin{cases} |\{j \in \mathcal{A} : z_j \in R(\boldsymbol{z}, \omega)\}|^{-1} & \text{if } z_i \in R(\boldsymbol{z}, \xi) \\ 0 & \text{otherwise.} \end{cases}$$

where $\{\xi : i \in [(n, s)]\}$ are independent copies of $\xi$, which incorporates the randomness of generating the tree with subsample $\mathcal{B}$ and $\omega$.

Then generalized U-statistic of order $s \in [n]$ is defined as

$$\overline{T} = \binom{n}{s}^{-1} \sum_{(n,s)} h(Z_{i_1}, \ldots, Z_{i_s}, \omega),$$

which have the H-decomposition expressed as

$$\overline{T} - \mathbb{E}[\overline{T}] = \sum_{j=1}^{s} U_j; \qquad U_j := \binom{s}{j}\binom{n}{j}^{-1} \sum_{(n,j)} h^{(j)}(\boldsymbol{X}, V_{i_1}, \ldots, V_{i_s}, \omega); \quad j \in [s],$$

where $U_j$'s are zero-mean and pairwise uncorrelated, and

$$h_i(v_1, \ldots, v_i) := \mathbb{E}[h(v_1, \ldots, v_i, V_{i+1}, \ldots, V_s, \xi)] - \mathbb{E}[h]; \qquad i \in [s]$$

$$h^{(i)}(v_1, \ldots, v_i) := h_i(v_1, \ldots, v_i) - \sum_{j=1}^{i-1} \sum_{(s,j)} h^{(j)}(v_{i_1}, \ldots, v_{i_j}); \qquad i \in [s-1]$$

$$h^{(s)}(v_1, \ldots, v_i) := h(v_1, \ldots, v_s, \xi) - \sum_{j=1}^{s-1} \sum_{(s,j)} h^{(j)}(v_{i_1}, \ldots, v_{i_j}).$$

In particular, $h^{(1)}(v) = \mathcal{H}_1(v) = \mathbb{E}[h(V_1, \cdots V_s, \xi | V_1 = v)] - \mathbb{E}[h]$ and

$$U_1 = \frac{s}{n} \sum_{i=1}^{n} \mathcal{H}_1(V_i).$$

Note that the sequence $\{(s/n)\mathcal{H}_1(V_i) : i \in [n]\}$ is centered and $i.i.d$ (for fixed $z$ and $n$ and $s$) hence then Linderberg condition is satisfied. Let $\sigma_1^2 := \mathbb{V}[\mathcal{H}_1(Z_1)]$ then

$$J_* := \frac{\sqrt{n}U_1}{s\sigma_1} \xrightarrow{d} \mathsf{N}(0, 1) \quad \text{in distribution as} \quad n \to \infty.$$

Consider the following decomposition.

$$J := \frac{\overline{T} - \mathbb{E}[\overline{T}]}{\sqrt{\mathbb{V}[\overline{T}]}} = J_* + (J - J_*).$$

If the second term vanishes in probability, we have that the term on the left-hand side weakly converges to a standard Gaussian. The second term is zero-mean, and its variance equals $2 - 2\mathbb{C}[J, J_*]$. Then, by Theorem 11.2 in Vaart (1998), it is sufficient for convergence in the second mean that

$$\rho := \frac{\mathbb{V}[\overline{T}]}{\mathbb{V}[U_1]} \to 1 \quad \text{as} \quad n \to \infty.$$

58

Since $\mathbb{V}[\overline{T}] \geq \mathbb{V}[U_1(\boldsymbol{x})]$, it suffices to state condition under which $\limsup_{s\to\infty} \rho \leq 1$. Moreover Peng et al. (2022) show that

$$\rho \leq 1 + \frac{s}{n}\frac{\mathbb{V}[h]}{s\sigma_1}.$$

If we apply Lemma 7 with $m \asymp s^{1-\frac{1-\eta}{K(\alpha)}}$ then by Lemma 2 we have that $\mathbb{P}(\mathcal{G}) \to 1$ and $m/s \to 0$. Therefore there is a constant $C_{f,d_Z}$ depending only $f$, the density of $\boldsymbol{Z}$ (which is assumed to bounded away from 0 and $\infty$ by Assumption 2(a)), and $d_Z$ such that with probability approaching 1

$$\limsup_{s\to\infty} \frac{C_{f_Z,d_Z}}{(\log s)^{d_Z}}\frac{\mathbb{V}[h]}{s\sigma_1} \lesssim 1.$$

Hence, if $s \to \infty$ and $s(\log n)^{d_Z} = o(n)$ as $n \to \infty$

$$J \xrightarrow{d} \mathsf{N}(0,1) \quad \text{and} \quad \sqrt{\rho}J \xrightarrow{d} \mathsf{N}(0,1) \quad \text{in distribution,}$$

and the proof of (S.41) is complete.

Also, from Lemma 7 we have that

$$\sigma_1^2 = \mathbb{V}[\mathbb{E}[T_W(x;z)|V_1] \gtrsim \mathbb{V}\big[\mathbb{E}[S_1|V_1]\big]\mathbb{V}[W|\boldsymbol{Z}=\boldsymbol{z}] \gtrsim \frac{C_{f,d_Z}\mathbb{V}[W|\boldsymbol{Z}=\boldsymbol{z}]}{sm(\log s)^d} \gtrsim \frac{1}{s(\log s)^d}$$

then

$$\lambda^2 := \frac{s^2}{n}\sigma_1^2 = \left(\frac{s}{n}\right)^2 n\sigma_1^2 \gtrsim \frac{s}{n}\frac{1}{s^{1-\frac{1-\eta}{K(\alpha)}}(\log s)^d} = \frac{s^{\frac{1-\eta}{K(\alpha)}}}{n(\log s)^{d_Z}}.$$

Moreover, there exist a constant $C$ (depending on $C_{f_Z,d_Z}$ and $\underline{\sigma}^2 := \min_{j\in d_X} \mathbb{V}[W|\boldsymbol{Z}] > 0$ by Assumption 2(d)) such that

$$\sigma_1^2(\boldsymbol{z}) = \mathbb{V}[\mathbb{E}[T_W(x;\boldsymbol{Z})|V_1] \gtrsim \mathbb{V}\big[\mathbb{E}[S_1|V_1]\big]\mathbb{V}[W|\boldsymbol{Z}=\boldsymbol{z}] \gtrsim \frac{C_{f,d_Z}\mathbb{V}[W|\boldsymbol{Z}=\boldsymbol{z}]}{sk(\log s)^d} \gtrsim \frac{1}{sk(\log s)^d}$$

Therefore, there are constants $C_1, C_2$ such that

$$C_1 s/(nk(\log s)^d) \leq \sigma_1(\boldsymbol{x}) \leq C_2(s/n)$$

Therefore we might take $\Lambda(\boldsymbol{z}) = \mathbb{V}[\mathbb{E}[T(\boldsymbol{z},\omega)|(Y_1,X_1,Z_1) = (\boldsymbol{X},y,x)]]$ to conclude that

$$C_1'\frac{s^{1-\eta}}{n(\log s)^d} \leq \|\Lambda(\boldsymbol{z})\| \leq C_2'(s/n) \to 0$$

hence $\|\Lambda(\boldsymbol{z})\| \xrightarrow{\mathbb{P}} 0$ as $n \to \infty$ and the proof of $(i)$ is complete.

Now, recall

$$\mathcal{H}_1(w, z; z') := \mathbb{E}[h(V_1, \ldots, V_s, \xi | V_1 = v)] - \mathbb{E}[h]$$

$$= \mathbb{E}\left[\sum_{i=1}^{s} W_i S_i(z', \zeta) \Big| W_1 = w, Z_1 = \boldsymbol{z}\right] - \mathbb{E}[W | \boldsymbol{Z} \in R(z', \omega)]$$

$$= w \mathbb{E}[S_1(z', \omega) | Z_1 = \boldsymbol{z}] - \mathbb{E}[W | \boldsymbol{Z} \in R(z', \omega)]$$

$$= w \mathbb{E}[S(z', \omega) | \boldsymbol{Z} = \boldsymbol{z}] - \mathbb{E}[W | \boldsymbol{Z} \in R(z', \omega)].$$

Note that when $z', z'' \in [0,1]^d$ belong to the same leaf $R_\ell$ then $(w, z) \mapsto \mathcal{H}_1(w, z; z')$ equals $(w, z) \mapsto \mathcal{H}_1(w, z; z'')$ so we define $\kappa_\ell(w, z) := w \mathbb{E}[S^{(\ell)} | \boldsymbol{Z} = \boldsymbol{z}] - \mathbb{E}[W | \boldsymbol{Z} \in R_\ell]$ where $S^{(\ell)}$ is defined as above with $R(z', \omega)$ replaced by $R_\ell$ for $l \in [L]$. Finally, we define

$$S = \frac{1}{\sqrt{n}} \sum_{i=1}^{n} Y_i; \qquad Y_i := (Y_{i,1}, \ldots, Y_{i,L})^\mathsf{T}; \qquad Y_{i,\ell} := \kappa_\ell(W_i, \boldsymbol{Z}_i); \quad \ell \in [L].$$

So $S$ is a sum of centered iid random vectors with covariance structure $\Gamma := [\sigma_{\ell,\ell'}]_{\ell,\ell' \in [L]}$

$$\gamma_{\ell,\ell'} := \mathbb{E}[Y_{1,\ell}, Y_{1,\ell'}] = \mathbb{E}[\kappa_\ell(W, \boldsymbol{Z}) \kappa_{\ell'}(W, \boldsymbol{Z})] = \mathbb{C}[W \mathbb{E}[S^{(\ell)} | \boldsymbol{Z}], W \mathbb{E}[S^{(\ell')} | \boldsymbol{Z}]].$$

We need an upper bound for

$$\mathbb{E}[W^2 \mathbb{E}[S^{(\ell)} | \boldsymbol{Z}] \mathbb{E}[S^{(\ell')} | \boldsymbol{Z}]]$$

Define $S_G := \sum_{i=1}^{n} G_i$ where $G_i \sim N(0, \Gamma)$ are iid and $S_G^*$ such that $S_G^* | \text{data} \sim N(0, \widehat{\Gamma})$

$$\widetilde{\Gamma} = \frac{1}{n} \sum_{i=1}^{n} (Y_i - \overline{Y})(Y_i - \overline{Y})^\mathsf{T}; \qquad \overline{Y} := \frac{1}{n} \sum_{i=1}^{n} Y_i.$$

By the triangle inequality

$$\sup_{t \in \mathbb{R}} |\mathbb{P}(\|S\| \le t) - \mathbb{P}(\|S_G^*\| \le t | D)| \le \sup_{t \in \mathbb{R}} |\mathbb{P}(\|S\| \le t) - \mathbb{P}(\|S_G\| \le t)|$$

$$+ \sup_{t \in \mathbb{R}} |\mathbb{P}(\|S_G\| \le t) - \mathbb{P}(\|S_G^*\| \le t | D)|$$

The first term can be bounded by

$$\frac{\mu_3 (\log L)^2}{\sqrt{n}}, \qquad \mu_3 := \mathbb{E}[Y^\mathsf{T} \Gamma^{-1} Y \|Y\|_2] \le \lambda(\Gamma_{\min}) \mathbb{E}[\|Y\|^3]$$

$\square$

**Lemma 7**

Let $V_i = (Y_i, \boldsymbol{X}_i, \boldsymbol{Z}_i)$ for $i \in [n]$ and define the event $\mathcal{G} := \mathcal{G}(m, z) = \{m \le |A(\omega, z)|\}$. Then for $m \in [s]$ and $bsz \in [0,1]^d$ there is for constant $C_{f,d_Z}$ depending only on $f$ and $d_Z$ such that

$$\mathbb{V}\big[\mathbb{E}[S_1|V_1]\big] \gtrsim \frac{C_{f,d_Z}}{ms(\log s)^{d_Z}}\mathbb{P}(\mathcal{G}) - 1/s^2.$$

Furthermore, if $\mathbb{V}[W|\boldsymbol{Z} = \boldsymbol{z}] \le M < \infty$ and $P(\mathcal{G}) \to 1$ and $m/s \to 0$ as $s \to \infty$ then

$$\frac{\mathbb{V}[T_W(\boldsymbol{z}, \omega)]}{s\mathbb{V}\big[\mathbb{E}[Z_1|V_1]\big]} \lesssim_{\mathbb{P}} \frac{M(\log s)^{d_Z}}{C_{f,d_Z}}.$$

If further $\boldsymbol{z} \mapsto \mathbb{E}[W|\boldsymbol{Z} = \boldsymbol{z}]$ and $\boldsymbol{z} \mapsto \mathbb{E}[W^2|\boldsymbol{Z} = \boldsymbol{z}]$ are Lipschitz, and $\mathbb{V}[W|\boldsymbol{Z} = \boldsymbol{z}] \ge 0$ then

$$\frac{\mathbb{V}[T_W(\boldsymbol{z}, \omega)]}{\mathbb{V}\big[\mathbb{E}[T_W(\boldsymbol{z}, \omega)|V_1]\big]} \lesssim_{\mathbb{P}} \frac{M(\log s)^{d_Z}}{C_{f,d_Z}\mathbb{V}[W|\boldsymbol{Z} = \boldsymbol{z}]}.$$

**Proof** (Lemma 7)

Recall that $\boldsymbol{Z}_i$ is a $m$-potential nearest neighbor (PNN) of $z$ if there is an axis-aligned rectangle containing *only* $z$ and a subset of $Z_1, \ldots Z_s$ containing $\boldsymbol{Z}_i$ with size between $m$ and $2m - 1$ and nothing else. Define

$$P_i := P_i(m, z) := \mathbf{1}\{\boldsymbol{Z}_i \text{ is a } m\text{-PNN of } z\}; \quad bsz \in [0,1]^d.$$

Note that $P_i = 0$ implies $S_i = 0$ conditional on $\mathcal{G}$ and then

$$\mathbb{E}[S_1|Z_1, \mathcal{G}] \le \frac{1}{m}\mathbb{E}[P_1|Z_1].$$

Expression (33) in the proof of Lemma 3.2 in Wager and Athey (2018) gives us.

$$\mathbb{P}\left(\mathbb{E}[P_1|V_1] \ge \frac{1}{s^2}\right) \lesssim m\frac{2^{d+1}(\log s)^d}{(d-1)!s},$$

which implies that

$$\mathbb{E}\big[(\mathbb{E}[S_1|Z_1, \mathcal{G}])^2\big] \gtrsim \frac{C_{f,d_Z}}{ms(\log s)^{d_Z}}.$$

Therefore,

$$
\begin{aligned}
\mathbb{V}\big[\mathbb{E}[S_1|V_1]\big] &= \mathbb{E}\big[(\mathbb{E}[S_1|V_1])^2\big] - \big(\mathbb{E}[S_1]\big)^2 \\
&= \mathbb{E}\big[(\mathbb{E}[S_1|V_1])^2\big] - 1/s^2 \\
&= \mathbb{E}\Big[\big(\mathbb{E}[S_1|Z_1, \mathcal{G}]\mathbb{P}(\mathcal{G}) + \mathbb{E}[S_1|Z_1, \mathcal{G}^c]\mathbb{P}(\mathcal{G}^c)\big)^2\Big] - 1/s^2 \\
&\gtrsim \frac{C_{f,d_Z}}{ms(\log s)^{d_Z}}\mathbb{P}(\mathcal{G}) - 1/s^2.
\end{aligned}
$$

Now on $\mathcal{G}$, we have

$$m\mathbb{V}[T_W(\boldsymbol{z}, \omega)] \leq |\mathcal{A}(\omega, z)|\mathbb{V}[T_W(\boldsymbol{z}, \omega)] = \mathbb{E}\big[|\mathcal{A}(\omega, z)|\mathbb{V}[T_W(\boldsymbol{z}, \omega)|\mathbb{S}]\big] = \mathbb{E}\big[\mathbb{V}[W|\boldsymbol{Z} = \boldsymbol{z}]\big] \leq M.$$

Therefore, on $\mathcal{G}$ such that $\mathbb{P}(\mathcal{G}) \to 1$ and $m/s \to 0$ we have with probability approaching 1

$$\frac{\mathbb{V}[T_W(\boldsymbol{z}, \omega)]}{s\mathbb{V}\big[\mathbb{E}[S_1|V_1]\big]} \lesssim \frac{M}{\frac{C_{f,d_Z}}{(\log s)^{d_Z}}\mathbb{P}(\mathcal{G}) - m/s} \to \frac{M(\log s)^{d_Z}}{C_{f,d_Z}}.$$

$\square$