

A General Framework: Inference on High-Dimensional Nonparametric Models

Abstract

We propose a novel nonparametric model which naturally generalizes the sparse additive model and sparse varying coefficient model. Given covariate of interest X_1 and covariate of nuisance \mathbf{X}_2 , our model assumes the mean function to be a sum of the product of a known function $\eta(\mathbf{X}_2)$ and the function of interest $f(X_1)$, and the nuisance function which can be estimated in the form of basis expansion. We propose to infer the function of interest $f(\cdot)$ using a two step plug-in approach, given an initial estimator of nuisance. We study the issue of bias propagation, which describes the dependence structure of bias when estimating the function of interest $f(\cdot)$ on the bias when estimating nuisance. We quantify such issue by establishing the rate of convergence estimating f under the supremum norm, in terms of estimation/approximation rate of nuisance. We also propose a confidence band for f via Gaussian multiplier bootstrap, with its coverage probability converging to the nominal probability at the polynomial rate. Thorough numerical results on synthetic data are provided to demonstrate the efficacy of the theory.

1 Introduction

Nonparametric regression investigates the relationship between a target variable Y and many input variables $\mathbf{X} = (X_1, \dots, X_d)^\top$ without imposing strong assumptions. Consider a model

$$Y = f(\mathbf{X}) + \varepsilon, \quad (1.1)$$

where $\mathbf{X} \in \mathbb{R}^d$ is a d -dimensional random vector in \mathbb{R}^d , ε is random error satisfying $\mathbb{E}[\varepsilon|\mathbf{X}] = 0$, and Y is a target variable. When d is small, fitting a fully nonparametric model (1.1) is feasible (Wasserman, 2006). However, the interpretation of these models is challenging. When d is large, consistently fitting $f(\cdot)$ must rely on additional structural assumptions due to curse of dimensionality.

Among simple and powerful nonparametric models for dimensionality reduction, most important two of them are additive models (Friedman and Stuetzle, 1981; Stone, 1985; Hastie and Tibshirani, 1990) and varying coefficient models (Hastie and Tibshirani, 1993; Fan and Zhang, 1999).

Specifically, in additive models we assume that

$$Y = \sum_{j=1}^d f_j(X_j) + \varepsilon, \quad \text{and} \quad \mathbb{E}_{X_j}[f_j(X_j)] = 0, \quad (1.2)$$

where $f_j(\cdot)$ s are smooth univariate functions. Under the assumption that only s components are nonzero ($s \ll d$), significant progress has been made to understand additive models in high dimensions

(Lin et al., 2006; Ravikumar et al., 2009; Meier et al., 2009; Huang et al., 2010; Koltchinskii et al., 2010; Kato, 2012).

Varying coefficient model assumes the following conditional linear structure

$$Y = \sum_{j=1}^d X_j \beta_j(Z) + \varepsilon, \quad (1.3)$$

for given covariates $(Z, X_1, \dots, X_d)^T$ and response variable Y . It is particularly appealing in longitudinal studies where it allows one to examine the extent to which covariates affect response over time. Under the assumption that only s components are nonzero ($s \ll d$), significant progress has been made to understand varying coefficient models in high dimensions (Lian, 2012; Fan et al., 2014) as well.

This paper proposes a novel model family which naturally generalizes both additive models and varying coefficient models. We assume

$$Y = f(X_1)\eta(\mathbf{X}_{\setminus 1}) + g(\mathbf{X}) + \varepsilon, \quad (1.4)$$

where $f(\cdot)$ is the function of interest, $\eta(\cdot)$ is a known function, and $g(\cdot)$ is nuisance. Once the function of interest is selected, we will observe such structure in additive model (1.2) and varying coefficient model (1.3). Under our model, we aim to estimate the marginal influence of the first input variable X_1 on the target Y . In particular, we are interested in estimating $f(\cdot)$ in a high dimensional setting and provide a confidence band for it.

The main results of our paper is twofold. First, we propose a unified framework for studying the estimation problem under both sparse additive model and sparse varying coefficient model. We propose a kernel estimator of f via a plug-in procedure, and establish a uniform rate of convergence of such estimator under our model. Intuitively, if the nuisance g in our model can be initially estimated with rate r_n , the function of interest f can then be estimated with a certain rate depending on r_n . We give a explicit form of this rate of convergence depending on r_n , and conditions that r_n , the rate of convergence of nuisance, needs to satisfy to yield a fast rate for estimating f . Therefore, our result serves as an interface to determine the estimation quality of function of interest f , given knowledge on the initial estimation quality of nuisance g .

Second, we propose one way to construct confidence bands for the function of interest f , and therefore provide a unified framework for studying inference problem under both models. To establish the validity of the proposed confidence bands, we develop three technical ingredients: (1) the analysis of suprema of a high dimensional empirical process that arises from kernel estimator, (2) a debiasing method for the proposed estimator, and (3) the approximation analysis for the Gaussian multiplier bootstrap procedure. The supremum norm for our estimator is derived by applying results on the suprema of empirical processes (Koltchinskii, 2011; Van Der Vaart and Wellner, 1996; Bousquet, 2002). The de-biasing procedure for the kernel estimator extends the approach used in Ning and Liu (2014) to a inferring a nonparametric component of interest under a high dimensional semiparametric model. To prove the validity of the confidence band constructed by the Gaussian multiplier bootstrap, we generalize the method proposed in Chernozhukov et al. (2014a) and Chernozhukov et al. (2014b) to our high dimensional model.

1.1 Related Literature

Our work contributes to two different areas. For both areas, we made new methodological and technical contributions.

First, we contribute to the literature on high dimensional nonparametric estimation, which has recently seen a lot of activity. [Lafferty and Wasserman \(2008\)](#), [Bertin et al. \(2008\)](#), [Comminges et al. \(2012\)](#), [Yang et al. \(2015\)](#) and [Linero \(2018\)](#) study variable selection in a high dimensional nonparametric regression setting without assuming structural assumptions on $f(\cdot)$ beyond that it depends only on a subset of variables. A large number of papers have studied the sparse additive model in (1.2) ([Sardy and Tseng, 2004](#); [Lin et al., 2006](#); [Avalos et al., 2007](#); [Ravikumar et al., 2009](#); [Meier et al., 2009](#); [Huang et al., 2010](#); [Koltchinskii et al., 2010](#); [Raskutti et al., 2012](#); [Kato, 2012](#); [Petersen et al., 2014](#); [Rosasco et al., 2013](#); [Lou et al., 2014](#); [Wahl, 2014](#); [Gregory et al., 2016](#); [Yuan and Zhou, 2016](#); [Guo and Zhang, 2019a](#); [Yao and Zhang, 2020](#)). In addition, [Yao and Zhang \(2020\)](#) proposes a multi-resolution group Lasso method without the knowledge of sparsity or smoothness, [Yuan and Zhou \(2016\)](#) study minimax optimal rates of convergence for estimation. [Xu et al. \(2014\)](#) study a high dimensional convex nonparametric regression. [Dalalyan et al. \(2014\)](#) study the compound model, which includes the additive model as a special case. [Meier et al. \(2009\)](#), [Huang et al. \(2010\)](#), [Koltchinskii et al. \(2010\)](#), [Raskutti et al. \(2012\)](#), and [Kato \(2012\)](#) develop estimation schemes mainly based on the basis approximation and sparsity-smoothness regularization.

Second, we contribute to a growing literature on high dimensional inference. Initial work on high dimensional statistics has focused on estimation and prediction (see, for example, [Bühlmann and Van De Geer, 2011](#), for a recent overview) and much less work has been done on quantifying uncertainty, for example, hypothesis testing and confidence intervals. Recently, the focus has started to shift towards the latter problems. Initial work on construction of p-values in high dimensional models relied on correct inclusion of the relevant variables ([Wasserman and Roeder, 2009](#); [Meinshausen et al., 2009](#)). [Meinshausen and Bühlmann \(2010\)](#) and [Shah and Samworth \(2013\)](#) study stability selection procedure, which provides the family-wise error rate for any selection procedure. Hypothesis testing and confidence intervals for low dimensional parameters in high dimensional linear and generalized linear models are studied in [Guo and Zhang \(2019b\)](#), [Cai and Guo \(2020\)](#), [Cai and Guo \(2017\)](#), [Belloni et al. \(2018\)](#), [Belloni et al. \(2017\)](#), [Belloni et al. \(2014a\)](#), [Belloni et al. \(2013b\)](#), [Van de Geer et al. \(2014\)](#), [Javanmard and Montanari \(2014\)](#), [Javanmard and Montanari \(2013\)](#), and [Farrell \(2013\)](#). These methods construct honest, uniformly valid confidence intervals and hypothesis test based on the ℓ_1 penalized estimator in the first stage. Similar results are obtained in the context of ℓ_1 penalized least absolute deviation and quantile regression ([Belloni et al., 2014b, 2013a](#)). [Kozbur \(2013\)](#) extends the approach developed in [Belloni et al. \(2014a\)](#) to a nonparametric regression setting, where a pointwise confidence interval is obtained based on the penalized series estimator. [Meinshausen \(2014\)](#) studies construction of one-sided confidence intervals for groups of variables under weak assumptions on the design matrix. [Lockhart et al. \(2014\)](#) studies significance of the input variables that enter the model along the lasso path. [Lee et al. \(2013\)](#) and [Taylor et al. \(2014\)](#) perform post-selection inference conditional on the selected model. [Chatterjee et al. \(2013\)](#), [Liu et al. \(2013\)](#), [Chernozhukov et al. \(2013\)](#) and [Lopes \(2014\)](#) study properties of the bootstrap in high-dimensions. Our work is different to the existing literature as it enables statisticians to make global inference under a nonparametric high dimensional regression setting for the first time.

1.2 Outline of the Paper

This paper is organized as follows. In the next section, we formally describe our nonparametric model, point out the issue of bias propagation when conducting inference within our model, and give two concrete examples of functions in the model family. In Section 3.2, we introduce the debiasing procedure as a remedy of bias propagation. In Section 4, we provide the theoretical results on the

statistical rate of convergence of the estimator and proposes a method for constructing confidence bands via Gaussian multiplier bootstrap. In Section 5, we apply our theoretical results to sparse additive model and sparse varying coefficient model. The numerical experiments for synthetic and real data are collected in Section 6.

1.3 Notation

Let $[n]$ denote the set $\{1, \dots, n\}$ and let $\mathbb{1}\{\cdot\}$ denote the indicator function. For a vector $\mathbf{a} \in \mathbb{R}^d$, we let $\text{supp}(\mathbf{a}) = \{j \mid a_j \neq 0\}$ be the support set (with an analogous definition for matrices $\mathbf{A} \in \mathbb{R}^{n_1 \times n_2}$), $\|\mathbf{a}\|_q$, for $q \in [1, \infty)$, the ℓ_q -norm defined as $\|\mathbf{a}\|_q = (\sum_{i \in [n]} |a_i|^q)^{1/q}$ with the usual extensions for $q \in \{0, \infty\}$, that is, $\|\mathbf{a}\|_0 = |\text{supp}(\mathbf{a})|$ and $\|\mathbf{a}\|_\infty = \max_{i \in [n]} |a_i|$. If the vector $\mathbf{a} \in \mathbb{R}^d$ is decomposed into groups such that $\mathbf{a} = (\mathbf{a}_{\mathcal{G}_1}, \dots, \mathbf{a}_{\mathcal{G}_g})^T$, where $\mathcal{G}_1, \dots, \mathcal{G}_g \subset [d]$ are disjoint sets, we denote $\|\mathbf{a}\|_{p,q}^q = \sum_{k=1}^g \|\mathbf{a}_{\mathcal{G}_k}\|_p^q$ and $\|\mathbf{a}\|_{p,\infty} = \max_{k \in [g]} \|\mathbf{a}_{\mathcal{G}_k}\|_p$ for any $p, q \in [1, \infty)$. We also denote the set $\{1, \dots, j-1, j+1, \dots, d\}$ as $\setminus j$ and the vector $\mathbf{a}_{\setminus j} = (a_1, \dots, a_{j-1}, a_{j+1}, \dots, a_d)^T$. For the function $f \in L^2(\mathbb{R})$, we define the L^2 norm $\|f\|_2 = [\int f^2(x) dx]^{1/2}$ and the supremum norm $\|f\|_\infty = \sup_{x \in \mathbb{R}} |f(x)|$. For a matrix $\mathbf{A} \in \mathbb{R}^{n_1 \times n_2}$, we use the notation $\text{vec}(\mathbf{A})$ to denote the vector in $\mathbb{R}^{n_1 n_2}$ formed by stacking the columns of \mathbf{A} . We denote the Frobenius norm of \mathbf{A} by $\|\mathbf{A}\|_F^2 = \sum_{i \in [n_1], j \in [n_2]} A_{ij}^2$ and denote the operator norm as $\|\mathbf{A}\|_2 = \sup_{\|\mathbf{v}\|_2=1} \|\mathbf{A}\mathbf{v}\|_2$. For two sequences of numbers $\{\alpha_n\}_{n=1}^\infty$ and $\{\beta_n\}_{n=1}^\infty$, we use $\alpha_n = O(\beta_n)$ to denote that $\alpha_n \leq C\beta_n$ for some finite positive constant C , and for all n large enough. If $\alpha_n = O(\beta_n)$ and $\beta_n = O(\alpha_n)$, we use the notation $\alpha_n \asymp \beta_n$. The notation $\alpha_n = o(\beta_n)$ is used to denote that $\alpha_n \beta_n^{-1} \xrightarrow{n \rightarrow \infty} 0$. Throughout the paper, we let c, C be two generic absolute constants, whose values may change from line to line.

2 A Family of General Nonparametric Models

We consider the problem of conducting global inference in a general nonparametric model. Specifically, we want to construct a confidence band for a nonparametric function of interest. The model we are considering here includes SpAM (**S**parse **A**dditive **M**odel) and sparse varying coefficient model as special cases.

First, we introduce the Hölder class of functions.

Definition 2.1 (Hölder function class). The γ -th Hölder class $\mathcal{H}(\gamma, L)$ on \mathcal{X} is the set of $\ell = \lfloor \gamma \rfloor$ times differentiable functions $f : \mathcal{X} \rightarrow \mathbb{R}$, where $\lfloor \gamma \rfloor$ represents the largest integer smaller than γ . The derivative $f^{(\ell)}$ satisfies

$$|f^{(\ell)}(x) - f^{(\ell)}(y)| \leq L|x - y|^{\gamma - \ell}, \quad \text{for any } x, y \in \mathcal{X}.$$

Let $\mathbf{X} = (X_1, \mathbf{X}_2^T)^T$ be a covariate vector in \mathcal{X}^{d+1} , where X_1 is one dimensional random variable and \mathbf{X}_2 is d dimensional random vector. We consider the following general nonparametric model

$$Y = f(X_1)\eta(\mathbf{X}_2) + g(X_1, \mathbf{X}_2) + \varepsilon, \quad (2.1)$$

where $f \in \mathcal{H}(2, L)$ is a function of a low dimensional covariate X_1 and is treated as the component of interest, η is a known function of possibly much higher dimensional covariates \mathbf{X}_2 , and g is a function of both covariates (X_1, \mathbf{X}_2) and is treated as a nuisance. We will see that both the sparse additive model and sparse varying coefficient model are special instances of this generic nonparametric model (2.1).

We first consider the sparse additive model.

Definition 2.2 (Sparse additive model). Let $\mathcal{S} \in [d]$ be of size $s = |\mathcal{S}| \ll d$. The sparse additive model can be written as

$$Y = \sum_{j \in \mathcal{S}} f_j(X_j) + \varepsilon, \quad (2.2)$$

with unknown $f_j \in \mathcal{H}(2, L)$, and ε being the error term that satisfies $\mathbb{E}[\varepsilon | \mathbf{X}] = 0$ almost surely.

Under the sparse additive model, we are interested in estimating and constructing a confidence band for f_1 . We note that the sparse additive model can be represented as in (2.1) as follows: X_1 is the variable of interest, $(X_2, \dots, X_d)^\top$ are represented by the nuisance variable “ \mathbf{X}_2 ,” f_1 in the sparse additive model is treated as the component of interest “ f ,” the constant 1 as the known function “ η ,” and $\sum_{j=2}^d f_j(X_j)$ as the nuisance “ g .”

Next, we consider the sparse varying coefficient model.

Definition 2.3 (Sparse varying coefficient model). Let $\mathcal{S} \in [d]$ be of size $s = |\mathcal{S}| \ll d$. The sparse varying coefficient model can be written as

$$Y = \sum_{j \in \mathcal{S}} X_j \beta_j(Z) + \varepsilon = \mathbf{X}^\top \boldsymbol{\beta}(Z) + \varepsilon, \quad (2.3)$$

with unknown $\beta_j \in \mathcal{H}(2, L)$, and ε being the error term that satisfies $\mathbb{E}[\varepsilon | \mathbf{X}] = 0$ almost surely. Here Z is a scalar variable that takes values in $[0, 1]$, $\mathbf{X} := (X_1, \dots, X_d)^\top$ is a d dimensional vector of covariates, and $\boldsymbol{\beta}(z) = (\beta_1(z), \dots, \beta_d(z))^\top \in \mathcal{R}^d$ is an unknown vector-valued function of coefficients, with only $s (s \ll d)$ of its components nonzero.

Under the sparse varying coefficient model, we are interested in estimating and constructing a confidence band for β_1 . We note that the sparse varying coefficient model can be represented as in (2.1) as follows: we view Z as the variable of interest “ X_1 ,” $(X_1, \dots, X_d)^\top$ are represented by the nuisance variable “ \mathbf{X}_2 ,” β_1 is treated as the component of interest “ f ,” X_1 as the known function “ η ,” and $\sum_{j=2}^d X_j \beta_j(Z)$ as the nuisance “ g .”

Our goal is to obtain a consistent estimator of f in the model (2.1), and to construct a confidence band for f based on the estimator. The first goal seems easier. Had we known exactly the nuisance function g , we should be able to construct an oracle estimator of f based on minimizing the following kernel loss function

$$\hat{f}^o(z) := \min_{\alpha} \frac{1}{2n} \sum_{i=1}^n K_h(X_{i1} - z) (Y_i - \alpha \cdot \eta(\mathbf{X}_{i2}) - g(X_{i1}, \mathbf{X}_{i2}))^2. \quad (2.4)$$

A practical estimator does not have access to the unknown nuisance g and has to estimate it. We assume that a consistent initial estimator of the nuisance, which we denote as \hat{g} , is obtainable. Later, we will provide a concrete estimator \hat{g} for sparse additive models and sparse varying coefficient models. With an initial estimator \hat{g} , we estimate f via the following plug-in procedure

$$\hat{f}(z) := \arg \min_{\alpha} \mathcal{L}_z(\alpha) = \frac{1}{2n} \sum_{i=1}^n K_h(X_{i1} - z) (Y_i - \alpha \cdot \eta(\mathbf{X}_{i2}) - \hat{g}(X_{i1}, \mathbf{X}_{i2}))^2. \quad (2.5)$$

We estimate the parameter of interest f by minimizing the loss function $\mathcal{L}_z(\alpha)$ in terms of α , where $\hat{g}(\mathbf{X}_i)$ is the initial estimator of the nuisance.

The quality of the estimator \hat{f} depends on the rate of convergence of the initial estimator \hat{g} . If the initial estimator \hat{g} is subject to bias introduced by regularization methods in high dimensional models, such bias will be propagated into \hat{f} during the plug-in estimation procedure (2.5), and therefore affecting the quality of estimator \hat{f} . If the bias propagated into \hat{f} is not asymptotically ignorable, it will later affect the covering probability of confidence band constructed using \hat{f} , preventing it from being an asymptotically honest one. We assume that the estimator of the nuisance $\hat{g}(X_1, \mathbf{X}_2)$ in the model (2.1) has the following form

$$\hat{g}(X_1, \mathbf{X}_2) = \hat{\gamma}^\top \psi(\mathbf{X}), \quad (2.6)$$

where $\psi(\mathbf{x})^\top = (\psi_1(\mathbf{x})^\top, \dots, \psi_d(\mathbf{x})^\top)$ is a series of known basis functions, and $\hat{\gamma}^\top = (\hat{\gamma}_1^\top, \dots, \hat{\gamma}_d^\top)$ is the corresponding coefficient vector. Here both ψ_j and $\hat{\gamma}_j$ are vectors of length m , that is, we use m basis functions when estimating each component of g . With the basis functions ψ_j ($j = 1, \dots, d$) fixed, we define the oracle vector γ^* as

$$\gamma^* = \arg \min_{\gamma} \|g(\cdot, \cdot) - \gamma^\top \psi(\cdot, \cdot)\|_{L_2}, \quad (2.7)$$

that is, we find the best approximation $\gamma^{*\top} \psi(\mathbf{x})$ of the nuisance g in the L_2 norm. The approximation rate in the supremum norm $\|g - \gamma^{*\top} \psi\|_\infty$ is naturally determined by the choice of basis functions ψ and the number of bases (i.e., m) we use to approximate.

We are now ready to introduce the debiasing procedure. The debiasing procedure serves as a remedy to the previous bias propagation issue, which we will discuss more in the following section.

3 Methodology

In Section 3.1, we present the plug-in estimator \hat{f} and analyze the bias of $\hat{f} - f$. In Section 3.2, we present the debiasing procedure. Specifically, we introduce the Immunization condition (Ning and Liu, 2014). In Section 3.3, we present the bootstrap procedures for constructing confidence bands.

3.1 Plug-in Estimator

We solve the following optimization to get plug-in estimator \hat{f} :

Input Dataset (\mathbf{X}, \mathbf{Y}) , an initial estimate $\hat{\gamma}$, basis function ψ , function η , parameter h , kernel density function K_h .

Output The plug-in estimator \hat{f} for function of interest f .

Step. Solve the follow quadratic optimization:

$$\hat{f}(z) \leftarrow \arg \min_{\alpha} \frac{1}{2n} \sum_{i=1}^n K_h(X_{i1} - z) (Y_i - \alpha \cdot \eta(\mathbf{X}_{i2}) - \hat{\gamma}^\top \psi(\mathbf{X}_i))^2. \quad (3.1)$$

In fact, we can have an explicit formula for \hat{f} . With assumption on structure of g and \hat{g} , we are now ready to study how the quality of estimation of f depends on the quality of initial estimator \hat{g} . We rewrite the loss function in (2.5) as a function of (α, γ) as

$$\mathcal{L}_z(\alpha, \gamma) := \frac{1}{2n} \sum_{i=1}^n K_h(X_{i1} - z) (Y_i - \alpha \cdot \eta(\mathbf{X}_{i2}) - \gamma^\top \psi(\mathbf{X}_i))^2. \quad (3.2)$$

Then the plug-in estimator $\hat{f}(z) := \arg \min_{\alpha} \mathcal{L}_z(\alpha, \hat{\gamma})$ can be viewed as the (unique) zero of $\nabla_{\alpha} \mathcal{L}_z(\alpha, \hat{\gamma})$. As the loss function is quadratic, we have the following representation of the error

$$\hat{f}(z) - f(z) = I(z)^{-1} (\mathbf{I}_1(z)^T (\gamma - \hat{\gamma}) + I_2(z) + I_3(z)),$$

where

$$\begin{aligned} I(z) &:= \mathbb{E}_n [K_h(X_1 - z) \eta(\mathbf{X}_2)^2], & \mathbf{I}_1(z) &:= \mathbb{E}_n [K_h(X_1 - z) \eta(\mathbf{X}_2) \psi(X)], \\ I_2(z) &:= \mathbb{E}_n [K_h(X_1 - z) \eta(\mathbf{X}_2) \varepsilon], & \text{and} & \quad I_3(z) := \mathbb{E}_n [K_h(X_1 - z) \eta(\mathbf{X}_2)^2 (f(X_1) - f(z))]. \end{aligned}$$

We notice that the initial estimation error $\hat{\gamma} - \gamma$ comes into estimation error $\hat{f}(z) - f(z)$ by a factor of $I(z)^{-1} \mathbf{I}_1(z)$. However, the estimator $\hat{f}(z)$ is not satisfactory for the purpose of constructing a confidence band. Let's take an analytic look at error $\hat{f}(z) - f(z)$. To construct a confidence band of f using \hat{f} , we want the error introduced by $\hat{\gamma} - \gamma$ in $\hat{f} - f$ to be asymptotically negligible up to the proper scaling, while here the factor $I(z)^{-1} \mathbf{I}_1(z)$ is obviously positive and not asymptotically negligible.

3.2 The Debiasing Procedure

In this section, we will discuss constructing confidence band based on the estimator in (3.2). In order to construct the confidence band, we will start with establish a score statistic of the formality $S(\alpha, \gamma) := \nabla_{\alpha} \mathcal{L}_z(\alpha, \gamma) - \mathbf{w}^T \nabla_{\gamma} \mathcal{L}_z(\alpha, \gamma)$, where $\mathbf{w} = \mathbf{w}(z)$ is a vector of the same length as γ and will be specified later. In particular, we will choose the vector \mathbf{w} so that the following immunization condition (Ning and Liu, 2014)

$$\mathbb{E}[\nabla_{\gamma} S(\alpha^*, \gamma^*)] = 0, \quad (3.3)$$

holds under the model in (2.1). We obtain our debiased estimator $\hat{\alpha} = \hat{f}(z)$ through the solution of $S(\alpha, \hat{\gamma}) = 0$. In general, this condition ensures that $S(\alpha, \gamma)$ is robust to the perturbation of high dimensional nuisance parameter γ . Specifically, we have

$$S(\alpha, \hat{\gamma}) \approx S(\alpha, \gamma)$$

the difference compared to the estimator obtained by minimizing (3.2), is that the bias introduced by estimation of the nuisance component will be negligible with respect to the bandwidth of the confidence band and will not affect the limiting distribution of $\hat{\alpha}$.

The choice of \mathbf{w} is determined by the condition (3.3). From the definition of $S(\alpha, \gamma)$, we have

$$\begin{aligned} \mathbb{E}[\nabla_{\gamma} S(\alpha, \gamma)] &= \mathbb{E} [\nabla_{\gamma}^2 \mathcal{L}_z(\alpha, \gamma)] - \mathbf{w}^T \mathbb{E} [\nabla_{\gamma}^2 \mathcal{L}_z(\alpha, \gamma)] \\ &= \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n K_h(X_{i1} - z) \eta(\mathbf{X}_{i2}) \psi(\mathbf{X}) \right] - \mathbf{w}^T \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n K_h(X_{i1} - z) \psi(\mathbf{X}) \psi(\mathbf{X})^T \right] \\ &= \mathbb{E} [K_h(X_1 - z) \eta(\mathbf{X}_2) \psi(\mathbf{X})] - \mathbf{w}^T \mathbb{E} [K_h(X_1 - z) \psi(\mathbf{X}) \psi(\mathbf{X})^T]. \end{aligned}$$

Therefore, the immunization condition (3.3) holds if \mathbf{w} is set as \mathbf{w}^* , where

$$\mathbf{w}^* := \mathbb{E} [K_h(X_1 - z) \psi(\mathbf{X}) \psi(\mathbf{X})^T]^{-1} \mathbb{E} [K_h(X_1 - z) \eta(\mathbf{X}_2) \psi(\mathbf{X})]. \quad (3.4)$$

We estimate \mathbf{w}^* as a solution $\tilde{\mathbf{w}} = \tilde{\mathbf{w}}_\lambda(z)$ of the following optimization problem

$$\tilde{\mathbf{w}} = \arg \min_{\mathbf{w}} \|\mathbf{w}\|_1 \quad \text{s.t.} \quad \left\| n^{-1} \sum_{i=1}^n K_h(X_{i1} - z) (\eta(\mathbf{X}_{i2}) - \mathbf{w}^T \boldsymbol{\psi}(\mathbf{X}_i)) \boldsymbol{\psi}(\mathbf{X}_i)^T \right\|_{2,\infty} \leq \lambda. \quad (3.5)$$

This is a convex optimization and we can solve it with some optimization software. As for λ , we just set λ of size $\mathcal{O}\left(\frac{m + \sqrt{m \log(dnh^{-1})}}{\sqrt{nh}}\right)$. Finally, we propose the following debiased estimator \tilde{f} of f ,

$$\tilde{f}(z) := \arg \min_{\alpha} \tilde{S}(\alpha, \hat{\gamma}) = \arg \min_{\alpha} \nabla_{\alpha} \mathcal{L}_z(\alpha, \hat{\gamma}) - \tilde{\mathbf{w}}^T \nabla_{\gamma} \mathcal{L}_z(\alpha, \hat{\gamma}), \quad (3.6)$$

where the estimated weight vector $\tilde{\mathbf{w}}$ is given in (3.5), the loss function \mathcal{L}_z is given in (3.2), and $\hat{\gamma}$ is assumed in (2.6) to be the coefficients of the basis expansion in the initial estimator \hat{g} . For this debiased estimator, the bias introduced by $\gamma - \hat{\gamma}$ can be controlled by $o(1/\sqrt{nh})$, which asymptotically negligible with respect to the bandwidth of our confidence band. Besides, we apply the "untangle and chord" procedure introduced in Ma et al. (2017). In detail, we split the whole dataset into 2 subsets randomly. We use the first subset to estimate $\tilde{\mathbf{w}}$ and the other one to estimate \tilde{f} . The benefit of this procedure is that $\tilde{\mathbf{w}}$ is independent to \tilde{f} , which makes $\tilde{\mathbf{w}}$ unable to effect the normality introduced \tilde{f} .

The debiasing procedure consists of the following two steps:

Input Datasets $(\mathbf{X}', \mathbf{Y}')$, (\mathbf{X}, \mathbf{Y}) , an initial estimate $\hat{\gamma}$, basis function $\boldsymbol{\psi}$, function η , parameter h , kernel density function K_h .

Output The debiased estimator \tilde{f} for function of interest f .

Step 1. Estimate a row of the inverse of Hessian by solving the following optimization program:

$$\tilde{\mathbf{w}} \leftarrow \arg \min_{\mathbf{w}} \|\mathbf{w}\|_1 \quad \text{s.t.} \quad \left\| n^{-1} \sum_{i=1}^n K_h(X'_{i1} - z) (\eta(\mathbf{X}'_{i2}) - \mathbf{w}^T \boldsymbol{\psi}(\mathbf{X}'_i)) \boldsymbol{\psi}(\mathbf{X}'_i)^T \right\|_{2,\infty} \leq \lambda. \quad (3.7)$$

Step 2. Obtain \tilde{f} with debiased score function:

$$\tilde{f}(z) \leftarrow \arg \min_{\alpha} \nabla_{\alpha} \mathcal{L}_z(\alpha, \hat{\gamma}) - \tilde{\mathbf{w}}^T \nabla_{\gamma} \mathcal{L}_z(\alpha, \hat{\gamma}). \quad (3.8)$$

3.3 Bootstrap Confidence Band

We construct the confidence band via Gaussian multiplier bootstrap with the following steps:

Input Datasets (\mathbf{X}, \mathbf{Y}) , an initial estimate $\hat{\gamma}$, basis function $\boldsymbol{\psi}$, function η , parameter h , kernel density function K_h , debiased estimator \tilde{f} , plug-in estimator \hat{f} , debiased parameter \tilde{w} , quantile parameter α .

Output The confidence band $\mathcal{C}_{n,\alpha}(z)$ for function of interest $f(z)$.

Step 1. Estimate 3 statistics with debiased parameter \tilde{w} :

$$\hat{q}^2(z) \leftarrow \frac{1}{n} \sum_{i=1}^n K_h^2(X_{i1} - z) (\eta(\mathbf{X}_{i2}) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{X}_i))^2, \quad (3.9)$$

$$\hat{p}(z) \leftarrow \frac{1}{n} \sum_{i=1}^n K_h(X_{i1} - z) (\eta(\mathbf{X}_{i2}) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{X}_i)) \eta(\mathbf{X}_{i2}), \quad (3.10)$$

$$\hat{\sigma}^2(z) \leftarrow \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{f}(X_{i1}) \eta(\mathbf{X}_{i2}) - \hat{\gamma} \boldsymbol{\psi}(\mathbf{X}_i))^2. \quad (3.11)$$

Step 2. Generate $\xi_1, \dots, \xi_{n_1} \sim N(0, I_{n \times n})$. Then for a series of z_1, \dots, z_{n_2} , we get $c_n(\alpha)$ as the $1 - \alpha$ quantile of

$$\left\{ \sup_j \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n \xi_{ji} \frac{\hat{\epsilon}_i K_h(X_{i1} - z_k) (\eta(\mathbf{X}_{i2}) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{X}_i))}{\hat{\sigma} \hat{q}(z_k)} \right| : k = 1, 2, \dots, n_2 \right\},$$

where $\hat{\epsilon}_i = Y_i - \hat{f}(X_{i1})$. Then our confidence band is:

$$\mathcal{C}_{n,\alpha}(z) := \tilde{f}(z) \pm \hat{c}_n(\alpha) n^{-1/2} \hat{\sigma} \hat{q}(z) / \hat{p}(z). \quad (3.12)$$

4 Theoretical Results

This section establishes the rate of convergence for the proposed estimator in Section 3.2. Based on the debiased estimator in (3.6), a confidence band for f is constructed via a Gaussian multiplier bootstrap.

4.1 Assumptions

We start with stating the required assumptions. Let $p_{\mathbf{X}}(\mathbf{x})$ denote the joint density of $\mathbf{X} = (X_1, \mathbf{X}_2)$ and let $p_{X_1}(x_1)$ denote the marginal density of X_1 . The b and B below refer to constants such that $0 < b \leq B < \infty$, while their values are subject to change in different assumptions.

- (A1) (Density function) The support \mathcal{X} is compact and the density function $p_{\mathbf{X}}(\mathbf{x})$ is continuous on \mathcal{X}^d . There exist fixed constants b, B such that $b \leq p_{X_1}(x_1) \leq B$ for all $x_1 \in \mathcal{X}$. The function $\eta(\cdot)$ depends only finite (say, s) components of \mathbf{x}_2 (say, (x_{21}, \dots, x_{2s})), and the joint density of $(X_1, X_{21}, \dots, X_{2s})$ is bounded away from 0: $p_{X_1, \mathbf{X}_{2,1:s}}(x_1, \mathbf{x}_{2,1:s}) \geq b$.
- (A2) (Kernel function) The kernel $K(u)$ is a continuous density function with bounded variation, satisfying

$$\int_{\mathcal{X}} K(u) du = 1 \quad \text{and} \quad \int_{\mathcal{X}} u K(u) du = 0.$$

- (A3) (Basis function) The basis functions $\boldsymbol{\psi}$ are uniformly bounded on \mathcal{X}

$$\max_{j \in [d], k \in [m]} \|\boldsymbol{\psi}_{jk}\|_{\infty} \leq B.$$

(A4) (Noise term) The error term ε satisfies $\mathbb{E}[\varepsilon | \mathbf{X}] = 0$, $\mathbb{E}[\varepsilon^2 | \mathbf{X}] = \sigma^2$, and is a subgaussian random variable such that $\mathbb{E}[\exp(\lambda\varepsilon)] \leq \exp(\lambda^2\sigma^2/2)$ for any λ .

(A5) (Nonparametric model) The nonparametric function $\mathbb{E}[Y | \mathbf{X}]$ has the form in (2.1), with the function $\eta(\cdot)$ satisfying

$$\|\eta\|_\infty \leq B \quad \text{and} \quad \|\eta\|_{L_2} \geq b.$$

(A6) (Rate of approximation and estimation) The nuisance component $g(\cdot)$ can be well approximated with the basis functions $\boldsymbol{\psi} = (\psi_1, \dots, \psi_d)$ as the rate r_a , that is

$$\|g - \boldsymbol{\gamma}^{*\top} \boldsymbol{\psi}\|_\infty \leq Cr_a. \quad (4.1)$$

The initial estimator $\hat{\boldsymbol{\gamma}}$ satisfies

$$\mathbb{P}(\|\hat{\boldsymbol{\gamma}} - \boldsymbol{\gamma}^*\|_{2,1} > Cr_e) \leq 1/n. \quad (4.2)$$

(A7) (Debiased term) To guarantee the existence of \mathbf{w}^* in (3.4), we assume the minimal eigenvalue of the Hessian matrix

$$\Sigma_z^* = \begin{bmatrix} \mathbb{E}K_h(X_1 - z)\boldsymbol{\psi}(\mathbf{X})\boldsymbol{\psi}(\mathbf{X})^\top & \mathbb{E}K_h(X_1 - z)\eta(\mathbf{X}_2)\boldsymbol{\psi}(\mathbf{X})^\top \\ \mathbb{E}K_h(X_1 - z)\eta(\mathbf{X}_2)\boldsymbol{\psi}(\mathbf{X}) & \mathbb{E}K_h(X_1 - z)\eta(\mathbf{X}_2)^2 \end{bmatrix}$$

is larger than a constant b for all z .

(A8) (Nonparametric sparsity) $\|\mathbf{w}^*\|_1 \leq L$.

Assumption (A1) is on the density function of covariates. Under sparse additive model $h \equiv 1$, hence h involves 0 component of \mathbf{x}_2 , and the assumption $p_{X_1, \mathbf{X}_{2,1:s}}(x_1, \mathbf{x}_{2,1:s}) \geq b$ automatically holds when $b = 1$ is chosen. Under sparse varying coefficient model $h(x_1) = x_1$, hence h involves only 1 variable of \mathbf{x}_2 , and the assumption $p_{X_1, \mathbf{X}_{2,1:s}}(x_1, \mathbf{x}_{2,1:s}) \geq b$ is equivalent to the joint density of (Z, X_1) is bounded away from 0, i.e., $p_{Z, X_1}(z, x_1) \leq b$. Under both cases, (A1) seems reasonable.

Assumption (A2) is standard in the literature on local linear regression (Fan, 1993). Assumption (A3) automatically holds if we use normalized B-spline basis as $\boldsymbol{\psi}$, and is hence reasonable. Assumption (A4) is standard for the noise variable in SpAM (Meier et al., 2009; Huang et al., 2010; Koltchinskii et al., 2010; Zhang et al., 2011; Raskutti et al., 2012; Kato, 2012).

Assumption (A5) is on the boundedness of function η . It is reasonable as we want the multiplicative factor in our model, function η , to be strong enough for f to be estimable.

Assumption (A6) is on the rate of basis approximation and initial estimation. For sparse additive model, we can choose $r_a = sm^{-2}$, and $r_e = n^{-1/2}m\sqrt{\log(md)} + m^{-3/2}$. For sparse varying coefficient model, we can choose $r_a = sm^{-2}$, and $r_e = n^{-1/2}m \log(md)$.

Assumption (A7) are conditions to guarantee the existence of \mathbf{w}^* . This is the same as the Assumption 4.2 in Ning and Liu (2014). Recall that $\mathbf{w}^* := \mathbb{E}[K_h(X_1 - z)\boldsymbol{\psi}(\mathbf{X})\boldsymbol{\psi}(\mathbf{X})^\top]^{-1} \mathbb{E}[K_h(X_1 - z)\eta(\mathbf{X}_2)\boldsymbol{\psi}(\mathbf{X})]$. In terms of Schur components, if Σ_z^* is nonsingular, $\mathbb{E}[K_h(X_1 - z)\boldsymbol{\psi}(\mathbf{X})\boldsymbol{\psi}(\mathbf{X})^\top]^{-1}$ is nonsingular too, which guarantees the existence of \mathbf{w}^* . Besides, we also have

$$\begin{aligned} \mathbb{E}K_h(X_1 - z)\eta(\mathbf{X}_2)^2 - \mathbf{w}^{*T} \mathbb{E}K_h(X_1 - z)\eta(\mathbf{X}_2)\boldsymbol{\psi}(\mathbf{X}) &= [-\mathbf{w}^{*T} \quad 1] \Sigma_z \begin{bmatrix} -\mathbf{w}^* \\ 1 \end{bmatrix} \\ &\geq b\sqrt{\|\mathbf{w}^*\|_2 + 1} \geq b. \end{aligned}$$

We will mainly use this inequality in our proofs instead of (A7).

Assumption (A8)

4.2 Estimation Consistency

We present the rate of convergence of the plug-in estimator.

Theorem 4.1 (Convergence rate of \hat{f} to f). Under Assumption 4.1, there exists a constant C such that if $(nh)^{-1} \log n = \mathcal{O}(1)$ and $h = o(1)$ $m \rightarrow \infty$ as $n \rightarrow \infty$ and we have that the plug-in estimator \hat{f} defined in (2.5) satisfies

$$\sup_{z \in \mathcal{X}} |\hat{f}(z) - f^*(z)| \leq C \left(r_a + \sqrt{n^{-1} h \log n} + h^2 + (nh)^{-1/2} \log n + \sqrt{mr_e} \right), \quad (4.3)$$

with probability $1 - c/n$ for a constant $c > 0$.

4.3 Asymptotically Honest Confidence Band via Gaussian Multiplier Bootstrap

Recently, Chernozhukov et al. (2014a) and Chernozhukov et al. (2014b) developed a novel framework to analyze the Gaussian multiplier bootstrap for approximating the distribution of empirical processes and apply it to kernel density estimation. We extend their approach and consider the Gaussian multiplier bootstrap to estimate the distribution of the supremum of a studentized process

$$Z_n(z) := \frac{\sqrt{n}(\tilde{f}(z) - \mathbb{E}_{f^*}[\tilde{f}(z)])}{\sigma(z)}, \quad (4.4)$$

where $\sigma(z)$ is a normalizing factor to be determined.

The idea here is straightforward. As long as the bias term $\mathbb{E}_{f^*}[\tilde{f}(z)] - f^*(z)$ is asymptotically negligible up to the normalizing factor $\sigma(z)$, we can construct a confidence band of the form $\tilde{f}(z) \pm c_n(\alpha) n^{-1/2} \sigma(z)$, where $c_n(\alpha)$ is the $(1 - \alpha)$ -quantile of supremum of the process $Z_n(z)$ (or its approximated version). Besides, Ma et al. (2017) also proposed a novel method for debiased estimator. Thus, we separate the dataset into subsets and use one to estimate $\tilde{\mathbf{w}}$ and another one to solve our debiased estimator. The benefit of this method is that $\tilde{\mathbf{w}}$ can be regarded as a constant in our proof and it won't damage our normality since it's estimated by another set. By setting $\sigma(z)$ as $\sigma q(z)/\hat{p}(z)$ and removing all bias terms, we obtain an empirical process

$$\mathbb{G}_n(z) := \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\varepsilon_i K_h(X_{i1} - z) (\eta(X_{i2}) - \tilde{\mathbf{w}}(z)^\top \boldsymbol{\psi}(\mathbf{X}_i))}{\sigma q(z)}, \quad (4.5)$$

where $q^2(z) := \mathbb{E}[K_h(X_1 - z)^2 (\eta(\mathbf{X}_2) - \tilde{\mathbf{w}}(z)^\top \boldsymbol{\psi}(\mathbf{X}))^2]$, and $\hat{p}(z) = n^{-1} \sum_{i=1}^n K_h(X_{i1} - z) \eta(\mathbf{X}_{i2}) (\eta(\mathbf{X}_{i2}) - \tilde{\mathbf{w}}(z)^\top \boldsymbol{\psi}(\mathbf{X}_i))$.

In order to estimate the distribution of supremum of \mathbb{G}_n , in particular its $(1 - \alpha)$ -quantile, we consider the following calculable Gaussian multiplier process

$$\hat{\mathbb{G}}_n(z) := \frac{1}{\sqrt{n}} \sum_{i=1}^n \xi_i \cdot \frac{\hat{\varepsilon}_i K_h(X_{i1} - z) (\eta(\mathbf{X}_{i2}) - \tilde{\mathbf{w}}(z)^\top \boldsymbol{\psi}(\mathbf{X}_i))}{\hat{\sigma} \hat{q}(z)}, \quad (4.6)$$

where $\{\xi_i\}_{i=1}^n$ are independent standard Gaussian random variables which are independent of $(\mathbf{X}, \boldsymbol{\varepsilon})$, and $\hat{\varepsilon}_i := Y_i - \hat{f}(X_{i1}) \eta(\mathbf{X}_{i2}) - \hat{g}(X_{i1}, \mathbf{X}_{i2})$ is the estimator of the noise term ε_i . The terms σ and $q(z)$ in denominator of $\mathbb{G}_n(z)$ are also replaced by their estimators $\hat{\sigma}^2 := n^{-1} \sum_{i=1}^n (Y_i - \hat{f}(X_{i1}) \eta(\mathbf{X}_{i2}) - \hat{g}(X_{i1}, \mathbf{X}_{i2}))^2$ and $\hat{q}(z)^2 := n^{-1} \sum_{i=1}^n K_h(X_{i1} - z)^2 (\eta(\mathbf{X}_{i2}) - \tilde{\mathbf{w}}(z)^\top \boldsymbol{\psi}(\mathbf{X}_i))^2$.

Finally, let $\widehat{c}_n(\alpha)$ be the conditional $(1 - \alpha)$ -quantile of $\sup_{z \in \mathcal{X}} |\widehat{\mathbb{G}}_n(z)|$, on fixed $\{(\mathbf{X}_i, Y_i)\}_{i=1}^n$. Based on the debiased estimator $\widetilde{f}(z)$, we construct the confidence band at level $(1 - \alpha)$ as $\mathcal{C}_{n,\alpha} = \{\mathcal{C}_{n,\alpha}(z) | z \in \mathcal{X}\}$ where

$$\mathcal{C}_{n,\alpha}(z) := \widetilde{f}(z) \pm \widehat{c}_n(\alpha) n^{-1/2} \widehat{\sigma} \widehat{q}(z) / \widehat{p}(z). \quad (4.7)$$

We can show that this confidence band is asymptotically honest with polynomial convergence rate.

Here we present several conditions about parameters r_e, r_a, m, h, d to guarantee our confidence band is honest.

Theorem 4.2 (Asymptotically honest confidence band via Gaussian multiplier bootstrap). Define the function class $\mathcal{F} = \{f \in \mathcal{F}(2, L) | \|f\|_\infty < M\}$. Under Assumption 4.1, assume $m \asymp n^{\delta_1}, h \asymp n^{-\delta_2}$, where $0 < \delta_1 < 1, \frac{1}{5} < \delta_2 < \frac{1}{4}$. If we also have $2\delta_1 + \delta_2 < 1, \log(dn) \cdot n^{\delta_1 - 1 + \delta_2} \rightarrow 0, r_a \cdot n^{(1 - \delta_2)/2} \rightarrow 0, r_e \cdot n^{\delta_1} \rightarrow 0, \log(dn) \cdot r_e \rightarrow 0$, then there exist constants $C, c > 0$ such that we have the coverage probability of $f^*(z)$

$$\inf_{f^*(z) \in \mathcal{F}} \mathbb{P}(f^*(z) \in \mathcal{C}_{n,\alpha}(z), \forall z \in \mathcal{X}) \geq 1 - \alpha - Cn^{-c}. \quad (4.8)$$

In particular, the confidence band $\mathcal{C}_{n,\alpha}$ is asymptotically honest, that is,

$$\liminf_{n \rightarrow \infty} \inf_{f^*(z) \in \mathcal{F}} \mathbb{P}(f^*(z) \in \mathcal{C}_{n,\alpha}(z), \forall z \in \mathcal{X}) \geq 1 - \alpha.$$

Remark 4.3. The confidence band constructed by the multiplier bootstrap neither requires ε to be symmetric nor relies on any conditions on the existence of limiting distribution for the supremum of the studentized empirical process in (4.4). We will give a proof of Theorem 4.2 in Section 7.

5 Convergence rates for Specific Model Classes

In this section, we give some theoretical results about convergence of \widehat{f} and confidence bands $\mathcal{C}_{n,\alpha}$ when applied to particular subclasses of the nonparametric model in (2.1). In particular, we provide debiasing estimators for sparse additive models and sparse varying coefficient models and corollaries for the convergences of their estimators and confidence bands.

5.1 Sparse Additive Models

We first present the debiased estimator for the sparse additive model (2.2). We approximate $f_j, j = 2, \dots, d$ by a linear combination of the normalized B-spline basis functions $\{\phi_1, \dots, \phi_m\}$. Given basis functions, we denote f_{jm} as the projection (in terms of L_2 norm) of f_j onto the space spanned by the basis. We have

$$f_{jm} := \arg \min_{f \in \mathcal{S}_m} \|f - f_j\|_2 = \sum_{k=1}^m \gamma_{jk}^* \psi_{jk},$$

where $\mathcal{S}_m = \text{Span}(\psi_{j1}, \dots, \psi_{jm})$ and $\psi_{jk} := \phi_k(x) - \mathbb{E}_n[\phi_k(X_j)]$ are centered basis functions. With this, we approximate the nuisance $g(\mathbf{X}) = \sum_{j=2}^d f_j(X_j)$ as

$$g(\mathbf{X}) \approx \sum_{j=2}^d f_{jm}(X_j) = \sum_{j=2}^d \sum_{k=1}^m \gamma_{jk}^* \psi_{jk}(X_j) = (\gamma_2^{*\top}, \dots, \gamma_d^{*\top})(\psi_2(X_2)^\top, \dots, \psi_d(X_d)^\top)^\top,$$

where $\boldsymbol{\gamma}_j^* = (\gamma_{j1}^*, \dots, \gamma_{jm}^*)^\top$ and $\boldsymbol{\psi}_j(x_j) = (\psi_{j1}(x_j), \dots, \psi_{jm}(x_j))^\top$. For functions in the Hölder class $\mathcal{H}(2, L)$, we have that $\|f_j - f_{jm}\|_\infty = \mathcal{O}(m^{-2})$ (see equation (20) in [Zhou et al. \(1998\)](#)). Furthermore, due to the sparsity assumption in our sparse additive model, we have that

$$\left\| \sum_{j=2}^d (f_j - f_{jm}) \right\|_\infty = \mathcal{O}(sm^{-2}). \quad (5.1)$$

We estimate $\boldsymbol{\gamma}^*$ using the group lasso estimator

$$\hat{\boldsymbol{\gamma}} := \arg \min_{\boldsymbol{\gamma}} \|\mathbf{Y} - \boldsymbol{\Psi}^\top \boldsymbol{\gamma}\|_2^2 + \lambda \sum_{j=1}^d \|\boldsymbol{\gamma}_j\|_2, \quad \text{where } \boldsymbol{\Psi}^{(n \times md)} = [\boldsymbol{\psi}(\mathbf{X}_1), \dots, \boldsymbol{\psi}(\mathbf{X}_n)]^\top, \quad (5.2)$$

which was proposed in [Huang et al. \(2010\)](#). When the number of nonzero components s is fixed and the tuning parameter $\lambda \asymp \sqrt{n \log(md)}$, we have that (see Theorem 1 in [Huang et al., 2010](#))

$$\|\hat{\boldsymbol{\gamma}} - \boldsymbol{\gamma}^*\|_2^2 = \mathcal{O}_{\mathbb{P}}\left(\frac{m^2 \log(md)}{n} + \frac{1}{m^3} + \frac{m^2 \lambda^2}{n^2}\right)$$

and

$$\|\hat{\boldsymbol{\gamma}} - \boldsymbol{\gamma}^*\|_{2,1} := \sum_{j=2}^d \|\hat{\boldsymbol{\gamma}}_j - \boldsymbol{\gamma}_j^*\|_2 = \mathcal{O}_{\mathbb{P}}\left(m \sqrt{\frac{\log(md)}{n}} + m^{-3/2}\right). \quad (5.3)$$

Applying Theorem 4.2 to the SpAM model, we have:

Corollary 5.1 (Confidence band under sparse additive model). For sparse additive model, we choose $r_a = sm^{-2}$, $r_e = n^{-1/2}m\sqrt{\log(md)} + m^{-3/2}$. Under Assumption 4.1, if $m \asymp n^{\delta_1}$, $3/16 < \delta_1 < 1/4$, $h \asymp n^{-\delta_2}$, $1/5 < \delta_2 < 1/4$ and as $n \rightarrow \infty$, $s \cdot n^{1/2-\delta_2/2-2\delta_1} \rightarrow 0$, $n^{-1+4\delta_1} \cdot \log(nd) \rightarrow 0$, there exists constants $C, c > 0$ such that we have the coverage probability of $f^*(z)$

$$\inf_{f^*(z) \in \mathcal{F}} \mathbb{P}(f^*(z) \in \mathcal{C}_{n,\alpha}(z), \forall z \in \mathcal{X}) \geq 1 - \alpha - Cn^{-c}. \quad (5.4)$$

In particular, the confidence band $\mathcal{C}_{n,\alpha}$ is asymptotically honest, that is,

$$\liminf_{n \rightarrow \infty} \inf_{f^*(z) \in \mathcal{F}} \mathbb{P}(f^*(z) \in \mathcal{C}_{n,\alpha}(z), \forall z \in \mathcal{X}) \geq 1 - \alpha.$$

Here we compare our estimator and confidence bands to those in [Lu et al. \(2019\)](#). There are three major differences between them: Firstly, their paper proposed a one-step plug-in estimator while our estimator is two-step. Secondly, their assumption for design matrix is specialized for B-splines, while we only assume our Hessian matrices' minimal eigenvalues are larger than a constant in Assumption 4.1 and our theoretic results can be applied to other first-step estimators for nuisance. Thirdly, about the biasing process, we apply the ‘‘untangle and chord’’ procedure introduced in [Ma et al. \(2017\)](#). In detail, we split the dataset into two subsets and use one subset to estimate $\tilde{\mathbf{w}}$ and the other one to estimate \hat{f} . While in [Lu et al. \(2019\)](#), they estimate $\tilde{\mathbf{w}}$ ($\hat{\theta}$ in their paper) and \hat{f} with the same dataset.

5.2 Sparse Varying Coefficient Model

In sparse varying coefficient model (2.3), again we can approximate β_j by B-spline basis. Given the normalized basis functions $\{\phi_1, \dots, \phi_m\}$, we denote β_{jm} as the projection (in terms of L_2 norm) of β_j onto the space spanned by the basis

$$\beta_{jm} := \arg \min_{\beta \in \mathcal{S}_m} \|\beta - \beta_j\|_2 = \sum_{k=1}^m \gamma_{jk}^* \phi_k, \quad \text{where } \mathcal{S}_m = \text{Span}(\phi_1, \dots, \phi_m).$$

Hence the nuisance $\sum_{j=2}^d \beta_j(Z) X_j$ can be approximated by

$$\begin{aligned} \sum_{j=2}^d \beta_j(Z) X_j &\asymp \sum_{j=2}^d \beta_{jm}(Z) X_j = \sum_{j=2}^d \sum_{k=1}^m \gamma_{jk}^* \phi_k(Z) X_j \\ &= (\gamma_2^{*\top}, \dots, \gamma_d^{*\top}) (\psi_2(Z, X_2)^\top, \dots, \psi_d(Z, X_d)^\top)^\top, \end{aligned}$$

where $\gamma_j^* = (\gamma_{j1}^*, \dots, \gamma_{jm}^*)^\top$ and $\psi_j(z, x_j) = (x_j \phi_1(z), \dots, x_j \phi_m(z))^\top$. For the same reason as above, we have the approximation rate (in $\|\cdot\|_\infty$) using B-spline basis on Hölder class $\mathcal{H}(2, L)$

$$\|\beta_j - \beta_{jm}\|_\infty = \mathcal{O}(m^{-2}), \quad \forall f_j \in \mathcal{H}(2, L),$$

and similarly as in (5.1)

$$\left\| \sum_{j=1}^d (x_j \beta_j(z) - x_j \beta_{jm}(z)) \right\|_\infty = \mathcal{O}(sm^{-2}). \quad (5.5)$$

Now let us consider the estimation rate of γ^* . We assume the number of nonzero components s is fixed. According to Theorem 3 in Lian (2012), the group lasso estimator $\hat{\gamma}$

$$\hat{\gamma} := \arg \min_{\gamma} \|Y - \Psi^\top \gamma\|_2^2 + \lambda \sum_{j=1}^d \|\gamma_j\|_1, \quad \text{where } \Psi^{(n \times md)} = [\psi(Z_1, \mathbf{X}_1), \dots, \psi(Z_n, \mathbf{X}_n)]^\top \quad (5.6)$$

with tuning parameter $\lambda/\sqrt{n \log d \vee n} \rightarrow \infty$ and $\lambda m^2/n \rightarrow 0$, satisfies

$$\sum_{i=1}^d \|\hat{\gamma}_i - \gamma_i^*\|_2^2 = \mathcal{O}_{\mathbb{P}}(m^2 \lambda^2 / n^2).$$

This implies we have the estimation rate in terms of $\|\cdot\|_{2,1}$ norm

$$\|\hat{\gamma} - \gamma^*\|_{2,1} := \sum_{i=2}^d \|\hat{\gamma}_i - \gamma_i^*\|_2 = \mathcal{O}_{\mathbb{P}}(m\lambda/n). \quad (5.7)$$

Applying Theorem 4.2 to the SVCMM model, we have:

Corollary 5.2 (Confidence band under sparse varying model). For sparse varying model, we choose $r_a = sm^{-2}$, $r_e = n^{-1/2} m \log(md)$. Under Assumption 4.1, if $m \asymp n^{\delta_1}$, $3/16 < \delta_1 < 1/4$,

$h \asymp n^{-\delta_2}$, $1/5 < \delta_2 < 1/4$ and as $n \rightarrow \infty$, $s \cdot n^{1/2-\delta_2/2-2\delta_1} \rightarrow 0$, $n^{-1/2+2\delta_1} \cdot \log(nd) \rightarrow 0$, there exist constants $C, c > 0$ such that we have the coverage probability of $f^*(z)$

$$\inf_{f^*(z) \in \mathcal{F}} \mathbb{P}(f^*(z) \in \mathcal{C}_{n,\alpha}(z), \forall z \in \mathcal{X}) \geq 1 - \alpha - Cn^{-c}. \quad (5.8)$$

In particular, the confidence band $\mathcal{C}_{n,\alpha}$ is asymptotically honest, that is,

$$\liminf_{n \rightarrow \infty} \inf_{f^*(z) \in \mathcal{F}} \mathbb{P}(f^*(z) \in \mathcal{C}_{n,\alpha}(z), \forall z \in \mathcal{X}) \geq 1 - \alpha.$$

6 Numerical Experiments

In this section, we study the finite sample properties of confidence bands for sparse additive model and sparse varying coefficient model.

Example 6.1. We generate data from the following sparse additive model

$$Y_i = f(\mathbf{X}_i) + \varepsilon_i = \sum_{j=1}^d f_j(X_{ij}) + \varepsilon_i, \quad i = 1, \dots, n,$$

where the additive components are designed as follows

$$\begin{aligned} f_1(t) &= 6(0.1 \sin(2\pi t) + 0.2 \cos(2\pi t) + 0.3 \sin(2\pi t)^2 + 0.4 \cos(2\pi t)^3 + 0.5 \sin(2\pi t)^3), \\ f_2(t) &= 5t, \quad f_3(t) = 3(2t - 1)^2, \quad f_4(t) = 4 \sin(2\pi t)/(2 - \sin(2\pi t)), \end{aligned}$$

and $f_5(t) = \dots = f_d(t) = 0$. The noise $\varepsilon_i \sim N(0, 1.5^2)$ for $i = 1, \dots, n$. The number of nonzero functions is thus $s = 4$. This generating model is the same as Example 1 of [Huang et al. \(2010\)](#). Let W_1, \dots, W_d and U follow iid Uniform[0, 1] and

$$X_j = \frac{W_j + tU}{1 + t} \quad \text{for } j = 1, \dots, d.$$

The covariates $X_{1j}, X_{2j}, \dots, X_{nj}$ are iid copies of X_j . We set the dimensions of covariates $d = 600$, $t = 0.3$, and consider three different sample sizes $n = 100, 300, 500$. In the initial group lasso estimator (2.6), we use the cubic B-splines with six evenly distributed knots and $m = 9$. The tuning parameter λ and bandwidth h are chosen by cross validation according to the BIC criterion defined as

$$\text{BIC}(\lambda, h) = \log(\text{RSS}_{(\lambda, h)}) + \text{df}_{(\lambda, h)} \cdot \frac{\log n}{n},$$

where RSS is the residual sum of squares for a given (λ, h) , and the degrees of freedom is defined as $\text{df} = \hat{s} \cdot m$ with \hat{s} being the number of nonzero estimated components selected for given (λ, h) . We aim at constructing confidence band for $f_1^*(t) = f_1(t) - \mathbb{E}[f_1(X_1)]$. To test the coverage probability of confidence bands for inactive covariates, we also construct confidence band for $f_5(t) \equiv 0$. The confidence band in (4.7) is constructed at the significance level 95% and the quantile estimator $\hat{c}_n(\alpha)$ is computed by bootstrap with 1,000 repetitions. To measure the coverage probability of the confidence band, we compute empirical probability that the confidence band covers the true function on the first 100 data points based on 500 repetitions. The numerical results are reported in Figure 1 and Table 1.

Example 6.2. We generate data from the following sparse varying coefficient model

$$Y_i = \mathbf{X}_i^T \boldsymbol{\beta}(Z_i) + \varepsilon_i = a_1 \cdot X_{i1} \beta_1(Z_i) + \sum_{j=1}^d X_{ij} \beta_j(Z_i) + \varepsilon_i, \quad i = 1, \dots, n,$$

where we set

$$a_1 \in \{0, 1\}, \quad \beta_1(z) = 3 + \sin(2\pi z), \quad \beta_2(z) = \exp(z + 1), \quad \beta_3(z) = 4(2z - 1)^2,$$

and $\beta_4(z) = \dots = \beta_d(z) = 0$. Here two values of $a_1 \in \{0, 1\}$ correspond to two scenarios that the true function is zero and nonzero. The noise $\varepsilon_i \sim N(0, 1.5^2)$ for $i = 1, \dots, n$. The number of nonzero functions is thus $s = 3$. This generating model is adapted from Example 1 of Lian (2012) by adding a_1 . The index variable Z is sampled uniformly on $[0, 1]$, and the covariates $X_{i1} = 1$ fixed and X_{i2}, \dots, X_{id} are independently and generated from $U(-1, 1)$, $i = 1, \dots, n$. We set the dimensions of covariates $d = 1000$ and consider three different sample sizes $n = 100, 300, 500$. In the initial group lasso estimator (2.6), we use the cubic B-splines with six evenly distributed knots and $m = 9$. Again we tune λ and h through cross validation via minimizing the BIC criterion. We aim to construct the confidence band for $a_1 f_1(t)$. The confidence band in (4.7) is constructed at the significance level 95% and the quantile estimator $\hat{c}_n(\alpha)$ is computed by bootstrap with 1,000 repetitions. To measure the coverage probability of the confidence band, we compute empirical probability that the confidence band covers the true function on the first 100 data points based on 500 repetitions. The numerical results are reported in Figure 2 and Table 1.

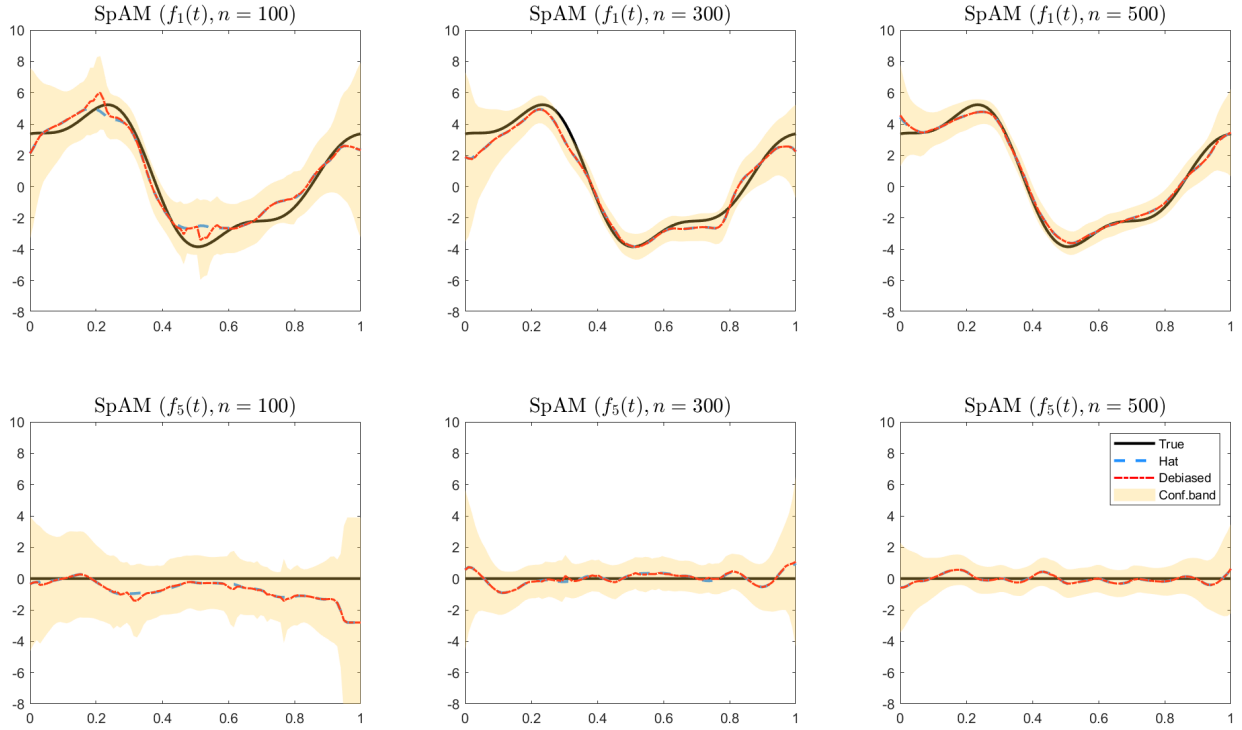


Figure 1: Debiased estimators for the $d = 600$ dimensional SpAM model $Y = \sum_{j=1}^4 f_j(X_j) + \varepsilon$, for $n = 100, 300, 500$ and the noise $\varepsilon \sim N(0, 1.5^2)$. The confidence bands at significant bands at significant level 95% cover $f_1(t)$ on the first row and $f_5(t) = 0$ on the second row.

| | | Scheme 1: f_1 | | | Scheme 2: f_5 | | |
|--------------|-------------------|---------------------|------|------|---------------------|------|------|
| Sample size: | | 100 | 300 | 500 | 100 | 300 | 500 |
| SpAM | Cover probability | 0.68 | 0.93 | 0.94 | 0.82 | 0.92 | 0.96 |
| | | Scheme 1: $a_1 = 1$ | | | Scheme 2: $a_1 = 0$ | | |
| Sample size: | | 100 | 300 | 500 | 100 | 300 | 500 |
| SpVC | Cover probability | 0.85 | 0.88 | 0.96 | 0.89 | 0.92 | 0.97 |

Table 1: Coverage probabilities for confidence bands at significant level 95% for SpAM model $Y = \sum_{j=1}^4 f_j(X_j) + \varepsilon$ and sparse varying coefficient model $Y = \sum_{j=1}^3 X_j \beta_j(Z) + \varepsilon$ with dimension $d = 1000$, sample size $n = 100, 300, 500$, and $\varepsilon \sim N(0, 1.5^2)$.

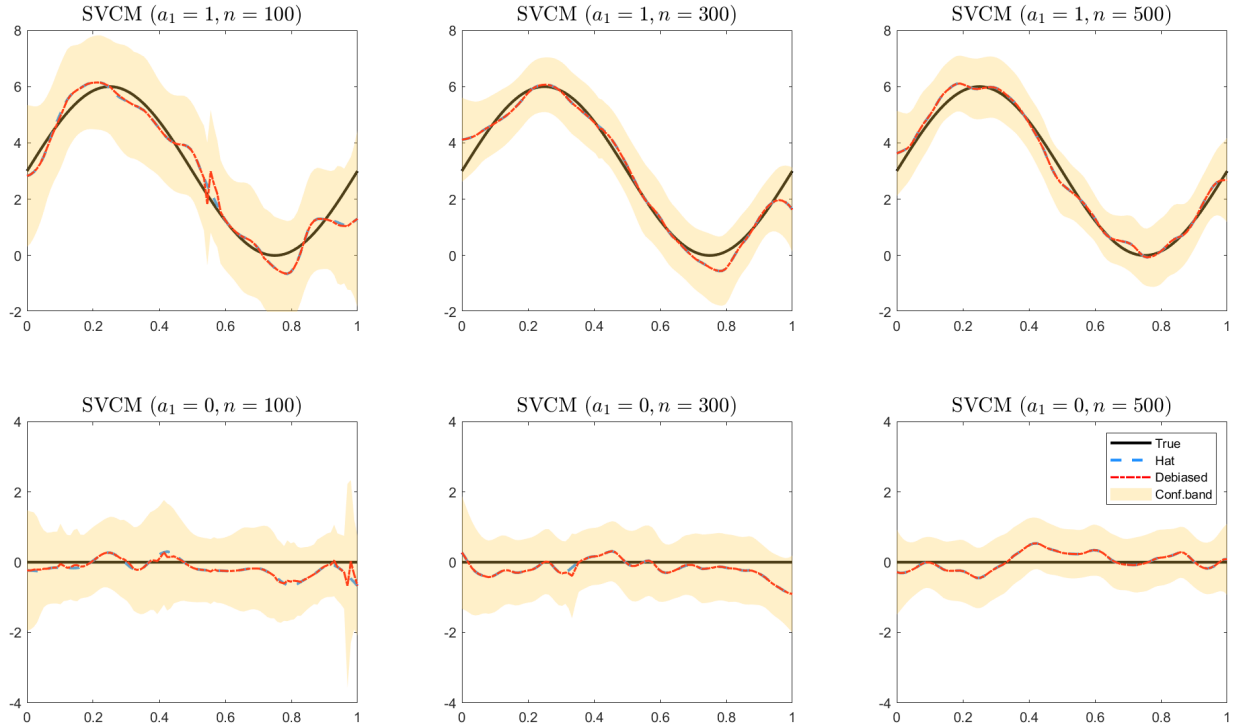


Figure 2: Debiased estimators for the $d = 600$ dimensional SVCVM model $Y = \sum_{j=1}^3 X_j \beta_j(Z) + \varepsilon$, for $n = 100, 300, 500$ and the noise $\varepsilon \sim N(0, 1.5^2)$. The confidence bands at significant level 95% cover $f_1^* = a_1 f_1$ for $a_1 \in \{0, 1\}$, respectively.

7 Proofs of Main Theoretical Results

We collect the proofs of main theorems in this section. The proofs of Theorem 4.1 and Theorem 4.2 are presented in Section 7.1 and Section 7.2 respectively.

7.1 Proof of the Statistical Rate of the Plug-in Estimator

This section outlines the proof of Theorem 4.1 on the statistical estimation rate of the plug-in estimator in (2.5).

The estimation error for the estimator \hat{f} comes from four sources: (1) noise ε , (2) approximation error by basis expansion, (3) estimation error introduced by kernel methods, and (4) error introduced by initial estimator $\hat{\gamma}$. Before presenting the main proof, we list several technical lemmas whose proofs are deferred to Appendix A in the supplementary material. The following lemma provides the rate for basis approximation. We denote $\mathbf{W}_z := \text{diag}(K_h(X_{11} - z), \dots, K_h(X_{n1} - z)) \in \mathbb{R}^{n \times n}$ as the kernel matrix, and $\boldsymbol{\eta} := (\eta(\mathbf{X}_{12}), \dots, \eta(\mathbf{X}_{n2}))^\top$ below.

Lemma 7.1. Let $\boldsymbol{\delta} = (\delta_1, \dots, \delta_n)^\top$, where $\delta_i = g(X_{i1}, \mathbf{X}_{i2}) - \gamma^{*\top} \boldsymbol{\psi}(\mathbf{X}_i)$ for $i = 1, \dots, n$. Under Assumptions 4.1 and $(nh)^{-1} \log n = \mathcal{O}(1)$, there exists a constant $C > 0$, such that the following inequality holds with probability at least $1 - 1/n$,

$$\sup_{z \in \mathcal{X}} \frac{1}{n} |\boldsymbol{\eta}^\top \mathbf{W}_z \boldsymbol{\delta}| \leq Cr_a, \quad (7.1)$$

where r_a is the approximation rate of basis $\boldsymbol{\psi}$: $\|g - \gamma^{*\top} \boldsymbol{\psi}\|_\infty \leq Cr_a$.

Our next lemma bound the estimation error introduced by kernel methods. We can see that the smoothing parameter (i.e., bandwidth) h play a role in the estimation.

Lemma 7.2. Let $\boldsymbol{\xi}_z = (\xi_1, \dots, \xi_n)^\top$, where $\xi_i(z) = (f(X_{i1}) - f(z))\eta(\mathbf{X}_{i2})$ for $i = 1, \dots, n$. Under Assumptions 4.1 and $(nh)^{-1} \log n = \mathcal{O}(1)$, there exists a constant $C > 0$, such that the following inequality holds with probability at least $1 - 1/n$,

$$\sup_{z \in \mathcal{X}} \frac{1}{n} |\boldsymbol{\eta}^\top \mathbf{W}_z \boldsymbol{\xi}_z| \leq Ch^2 + C\sqrt{n^{-1}h \log n}. \quad (7.2)$$

The following lemma quantifies the statistical rate from the noise ε .

Lemma 7.3. Let $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^\top$. Under Assumption 4.1 and $(nh)^{-1} \log n = \mathcal{O}(1)$, there exists a constant $C > 0$, such that the following inequality holds with probability at least $1 - 4/n$,

$$\sup_{z \in \mathcal{X}} \frac{1}{n} |\boldsymbol{\eta}^\top \mathbf{W}_z \boldsymbol{\varepsilon}| \leq C(nh)^{-1/2} \log n. \quad (7.3)$$

The last lemma quantifies the statistical rate introduced by initial estimation error of $\hat{\gamma}$.

Lemma 7.4. Let $\boldsymbol{\zeta} = (\zeta_1(z), \dots, \zeta_n(z))^\top$, where $\zeta_i = (\gamma^* - \hat{\gamma})^\top \boldsymbol{\psi}(\mathbf{X}_i)$. Under Assumption 4.1 and $(nh)^{-1} \log n = \mathcal{O}(1)$, there exists a constant $C > 0$, such that the following inequality holds with probability at least $1 - 1/n$,

$$\sup_{z \in \mathcal{X}} \frac{1}{n} |\boldsymbol{\eta}^\top \mathbf{W}_z \boldsymbol{\zeta}| \leq C\sqrt{mr_e}, \quad (7.4)$$

where r_e is the estimation rate of γ^* in ℓ_1 norm: $\mathbb{P}(\|\hat{\gamma} - \gamma^*\|_{2,1} > Cr_e) \leq 1/n$.

We are now ready to present the main proof of Theorem 4.1.

Proof of Theorem 4.1. Recall that $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^\top$, $\boldsymbol{\delta} = (\delta_1, \dots, \delta_n)^\top$, $\boldsymbol{\xi}_z = (\xi_1(z), \dots, \xi_n(z))^\top$, $\boldsymbol{\zeta} = (\zeta_1, \dots, \zeta_n)^\top$, where $\delta_i = g(X_{i1}, \mathbf{X}_{i2}) - \gamma^\top \boldsymbol{\psi}(\mathbf{X}_i)$, $\xi_i(z) = (f(X_{i1}) - f(z))\eta(\mathbf{X}_{i2})$, $\zeta_i = (\gamma - \hat{\gamma})^\top \boldsymbol{\psi}(\mathbf{X}_i)$, for $i = 1, \dots, n$. We denote $\boldsymbol{\eta}_z := \boldsymbol{\varepsilon} + \boldsymbol{\delta} + \boldsymbol{\xi}_z + \boldsymbol{\zeta}$.

$$\begin{aligned}\hat{f}(z) - f(z) &= \frac{\sum_{i=1}^n K_h(X_{i1} - z)\eta(\mathbf{X}_{i2})(Y_i - \hat{g}(X_{i1}, \mathbf{X}_{i2}))}{\sum_{i=1}^n K_h(X_{i1} - z)\eta(\mathbf{X}_{i2})^2} - f(z) \\ &= \frac{\sum_{i=1}^n K_h(X_{i1} - z)\eta(\mathbf{X}_{i2})(Y_i - f(z)\eta(\mathbf{X}_{i2}) - \hat{g}(X_{i1}, \mathbf{X}_{i2}))}{\sum_{i=1}^n K_h(X_{i1} - z)\eta(\mathbf{X}_{i2})^2}.\end{aligned}$$

Plugging in (2.1) and (2.6), we have

$$\begin{aligned}\hat{f}(z) - f(z) &= \frac{\sum_{i=1}^n K_h(X_{i1} - z)\eta(\mathbf{X}_{i2})\left((f(X_{i1}) - f(z))\eta(\mathbf{X}_{i2}) + g(X_{i1}, \mathbf{X}_{i2}) - \hat{g}(X_{i1}, \mathbf{X}_{i2}) + \varepsilon_i\right)}{\sum_{i=1}^n K_h(X_{i1} - z)\eta(\mathbf{X}_{i2})^2}, \\ &= \frac{\sum_{i=1}^n K_h(X_{i1} - z)\eta(\mathbf{X}_{i2})\left[(f(X_{i1}) - f(z))\eta(\mathbf{X}_{i2}) + g(X_{i1}, \mathbf{X}_{i2}) - \hat{\gamma}^\top \boldsymbol{\psi}(\mathbf{X}_i) + \varepsilon_i\right]}{\sum_{i=1}^n K_h(X_{i1} - z)\eta(\mathbf{X}_{i2})^2} \\ &= \frac{\boldsymbol{\eta}^\top \mathbf{W}_z (\boldsymbol{\xi}_z + \boldsymbol{\delta} + \boldsymbol{\zeta} + \boldsymbol{\varepsilon})}{\boldsymbol{\eta}^\top \mathbf{W}_z \boldsymbol{\eta}} = \frac{\boldsymbol{\eta}^\top \mathbf{W}_z \boldsymbol{\eta}_z}{\boldsymbol{\eta}^\top \mathbf{W}_z \boldsymbol{\eta}}.\end{aligned}$$

Using Lemma 7.1, Lemma 7.2, Lemma 7.3, and Lemma 7.4, there exist a constant C' such that with probability $1 - 8/n$,

$$\sup_{z \in \mathcal{X}} \frac{1}{n} |\boldsymbol{\eta}^\top \mathbf{W}_z \boldsymbol{\eta}_z| \leq C' r_n, \quad (7.5)$$

where $r_n := r_a + \sqrt{n^{-1}h \log n} + h^2 + (nh)^{-1/2} \log n + \sqrt{mr_e}$.

By McDiarmid's inequality, there exists a constant C , such that

$$\mathbb{P}\left(\left|n^{-1} \sum_{i=1}^n K_h(X_{i1} - z)\eta(\mathbf{X}_{i2})^2 - \mathbb{E}[K_h(X_1 - z)\eta(\mathbf{X}_2)^2]\right| \geq C\sqrt{\log n}/(\sqrt{nh})\right) \leq \frac{1}{n}.$$

Also by (A1) and (A5) in Assumption 4.1, we have

$$\begin{aligned}\mathbb{E}[\eta(\mathbf{X}_2)^2 | X_1 = x] &= \int \eta(\mathbf{X}_{2,1:s})^2 \frac{p_{X_1, \mathbf{X}_{2,1:s}}(x, \mathbf{x}_{2,1:s})}{p_{X_1}(x)} d\mathbf{x}_{2,1:s} \geq B^{-1}b \|h\|_{L_2}^2 \geq B^{-1}b^3. \\ \mathbb{E}[K_h(X_1 - z)\eta(\mathbf{X}_2)^2] &= \mathbb{E}[K_h(X_1 - z)\mathbb{E}[\eta(\mathbf{X}_2)^2 | X_1]] \geq B^{-1}b^3 \cdot \mathbb{E}[K_h(X_1 - z)],\end{aligned}$$

where $\mathbb{E}[K_h(X_1 - z)] = \int_{\mathcal{X}} K(t)p_{X_1}(z + th)dt \geq b/2$. Hence with probability at least $1 - 1/n$, we have

$$\frac{1}{n} \boldsymbol{\eta}^\top \mathbf{W}_z \boldsymbol{\eta} \geq B^{-1}b^4/4, \quad \forall z \in \mathcal{X}. \quad (7.6)$$

To sum up, on the event where both (7.5) and (7.6) hold (with probability at least $1 - 9/n$), we have

$$\sup_{z \in \mathcal{X}} \left| \hat{f}(z) - f^*(z) \right| = \sup_{z \in \mathcal{X}} \frac{n^{-1} |\boldsymbol{\eta}^\top \mathbf{W}_z \boldsymbol{\eta}_z|}{n^{-1} \boldsymbol{\eta}^\top \mathbf{W}_z \boldsymbol{\eta}} \leq 4BC'/b^4 \cdot r_n = Cr_n.$$

□

7.2 Proof of Covering Properties of the Bootstrap Confidence Bands

In this section, we prove Theorem 4.2 on the coverage probabilities for the Gaussian Multiplier bootstrap confidence bands $\mathcal{C}_{n,\alpha}$ in (4.7).

In our proof, we will apply the “untangle and chord” procedure introduced in Ma et al. (2017) to get our $\tilde{\mathbf{w}}$. Thus, we use $\{(\mathbf{X}'_i, Y'_i)\}, i = 1, \dots, n$ to denote the subdataset for \mathbf{w} and $\{(\mathbf{X}_i, Y_i)\}$ to denote the subdataset for our debiased estimator. Recall that we have two optimization problems for getting \tilde{f} :

$$\begin{aligned} \tilde{\mathbf{w}} &= \arg \min \|\mathbf{w}\|_1 \\ \text{s.t.} \quad & \|\nabla_{\alpha\gamma}^2 L_z - \mathbf{w}^T \nabla_{\gamma\gamma}^2 L_z\|_{2,\infty} \leq \lambda, \end{aligned} \quad (7.7)$$

where we use the subdataset $\{(\mathbf{X}'_i, Y'_i)\}$ and $\lambda_1 = C \frac{m + \sqrt{m \log(dnh^{-1})}}{\sqrt{nh}}$. Then,

$$\begin{aligned} \tilde{f} &= \arg \text{zero} S_z(\alpha, \hat{\gamma}) \\ \text{s.t.} \quad & S_z(\alpha, \hat{\gamma}) = \nabla_{\alpha} L_z(\alpha, \hat{\gamma}) - \tilde{\mathbf{w}} \nabla_{\gamma} L_z(\alpha, \hat{\gamma}), \end{aligned} \quad (7.8)$$

where we use the subdataset $\{(\mathbf{X}_i, Y_i)\}$. Now we list all the lemmas for the main proof. The first lemma guarantees the existence of solution to Optimization 7.7.

Lemma 7.5. Under Assumptions 4.1 and conditions in Theorem 4.2, with probability $1 - \frac{c}{n}$, there exists a solution $\tilde{\mathbf{w}}$ for the optimization problem 7.7, and $\|\tilde{\mathbf{w}}\|_1 \leq L$.

For convenience, we also use $\tilde{\alpha}, \alpha^*$ to denote $\tilde{f}(z), f^*(z)$. Then the following lemma gives the difference between $\tilde{\alpha}$ and α^* .

Lemma 7.6. According to the optimization problems 7.7 and 7.8, under the condition that $(nh)^{-1} \log n = o(1)$, we have

$$\tilde{\alpha} - \alpha^* = \frac{(\nabla_{\alpha\gamma}^2 L_z - \tilde{\mathbf{w}}^T \nabla_{\gamma\gamma}^2 L_z)(\gamma^* - \hat{\gamma}) + \tilde{\mathbf{w}} \nabla_r L_z(\alpha^*, \gamma^*) - \nabla_{\alpha} L_z}{\nabla_{\alpha\alpha}^2 L_z - \tilde{\mathbf{w}} \nabla_{\alpha\gamma}^2 L_z}.$$

Then we use I_1, I_2 to denote

$$\frac{(\nabla_{\alpha\gamma}^2 L_z - \tilde{\mathbf{w}}^T \nabla_{\gamma\gamma}^2 L_z)(\gamma^* - \hat{\gamma})}{\nabla_{\alpha\alpha}^2 L_z - \tilde{\mathbf{w}} \nabla_{\alpha\gamma}^2 L_z}, \quad \frac{\tilde{\mathbf{w}} \nabla_r L_z(\alpha^*, \gamma^*) - \nabla_{\alpha} L_z(\alpha^*, \gamma^*)}{\nabla_{\alpha\alpha}^2 L_z - \tilde{\mathbf{w}} \nabla_{\alpha\gamma}^2 L_z}.$$

Thus, $\tilde{\alpha} - \alpha^* = I_1 + I_2$. The following lemma upper bounds I_1 .

Lemma 7.7. Under Assumption 4.1, there exists a constant C such that

$$\mathbb{P} \left(I_1 \leq Cr_e \cdot \frac{m + \sqrt{m \log(dnh^{-1})}}{\sqrt{nh}} \right) \geq 1 - \frac{7}{n}.$$

We can rewrite I_2 as

$$\begin{aligned} I_2 &= \frac{\frac{1}{n} \sum_{i=1}^n K_h(X_{i1} - z)(\eta(\mathbf{X}_{i2}) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{X}_i))(Y_i - \alpha^* \eta(\mathbf{X}_{i2}) - \gamma^{*T} \boldsymbol{\psi}(\mathbf{X}_i))}{\frac{1}{n} \sum_{i=1}^n K_h(X_{i1} - z)(\eta(\mathbf{X}_{i2}) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{X}_i))\eta(\mathbf{X}_{i2})} \\ &= \frac{\frac{1}{n} \sum_{i=1}^n t(\mathbf{X}_i, z)[(f^*(X_{i1}) - f(z))\eta(\mathbf{X}_{i2}) + g(\mathbf{X}_i) - \gamma^{*T} \boldsymbol{\psi}(\mathbf{X}_i)]}{\hat{p}(z)} + \frac{\frac{1}{n} \sum_{i=1}^n t(\mathbf{X}_i, z)\varepsilon_i}{\hat{p}(z)}, \end{aligned}$$

where $t(\mathbf{X}_i, z) := K_h(X_{i1} - z)(\eta(\mathbf{X}_{i2}) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{X}_i))$ and $\hat{p}(z) := n^{-1} \sum_{i=1}^n K_h(X_{i1} - z) \eta(\mathbf{X}_{i2}) (\eta(\mathbf{X}_{i2}) - \tilde{\mathbf{w}}(z)^T \boldsymbol{\psi}(\mathbf{X}_i))$. Now we use I_{21}, I_{22} to denote,

$$\frac{\frac{1}{n} \sum_{i=1}^n t(\mathbf{X}_i, z) [(f^*(X_{i1}) - f(z)) \eta(\mathbf{X}_{i2}) + g(\mathbf{X}_i) - \gamma^{*T} \boldsymbol{\psi}(\mathbf{X}_i)]}{\hat{p}(z)}, \quad \frac{\frac{1}{n} \sum_{i=1}^n t(\mathbf{X}_i, z) \varepsilon_i}{\hat{p}(z)}.$$

The following lemma gives an upper bound of I_{21} .

Lemma 7.8. Under Assumption 4.1, there exist constants C and c , such that

$$\mathbb{P} \left(I_{21} \leq C(r_a + h^2 + \sqrt{n^{-1}h \log n}) \right) \geq 1 - \frac{c}{n}.$$

In fact $I_{21} + I_1$ is the bias term in our debiased estimator. We will proof it's asymptotically negligible with respect to our bandwidth of confidence band later. As for I_{22} , we will established a sequence of processes to construct confidence band for its numerator. Recall that

$$\begin{aligned} \mathbb{G}_n(z) &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\varepsilon_i K_h(X_{i1} - z) (\eta(\mathbf{X}_{i2}) - \tilde{\mathbf{w}}(z)^T \boldsymbol{\psi}(\mathbf{X}_i))}{\sigma q(z)} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\varepsilon_i t(\mathbf{X}_i, z)}{\sigma q(z)}, \\ \hat{\mathbb{G}}_n(z) &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \xi_i \cdot \frac{\hat{\varepsilon}_i K_h(X_{i1} - z) (\eta(\mathbf{X}_{i2}) - \tilde{\mathbf{w}}(z)^T \boldsymbol{\psi}(\mathbf{X}_i))}{\hat{\sigma} \hat{q}(z)} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\xi_i \hat{\varepsilon}_i t(\mathbf{X}_i, z)}{\hat{\sigma} \hat{q}(z)}, \end{aligned}$$

where $\hat{\sigma}^2 = n^{-1} \sum_{i=1}^n \hat{\varepsilon}_i^2$ is the estimator of $\sigma^2 = \text{Var}(\varepsilon_1)$, and $\hat{q}(z)^2 = n^{-1} \sum_{i=1}^n K_h(X_{i1} - z)^2 (\eta(\mathbf{X}_{i2}) - \tilde{\mathbf{w}}(z)^T \boldsymbol{\psi}(\mathbf{X}_i))^2$ is the estimator of $q^2(z) = \mathbb{E}[K_h(X_1 - z)^2 (\eta(\mathbf{X}_2) - \tilde{\mathbf{w}}(z)^T \boldsymbol{\psi}(\mathbf{X}))^2]$. The following lemma gives a bound for the difference between $\mathbb{G}_n(z)$'s and $\hat{\mathbb{G}}_n(z)$'s denominators. Let

$$r_n := r_a + \sqrt{n^{-1}h \log n} + h^2 + (nh)^{-1/2} \log n + r_e \sqrt{m}.$$

Lemma 7.9. Let the estimator for $\text{Var}(\varepsilon) = \sigma^2$ be $\hat{\sigma}^2 = n^{-1} \sum_{i=1}^n \hat{\varepsilon}_i^2$. Under Assumption 4.1, there exist constants C and c such that $\mathbb{P}(|\hat{\sigma}^2 - \sigma^2| \geq Cr_n) \leq 6/n$ and

$$\mathbb{P} \left(\sup_{z \in \mathcal{X}} \left| \frac{\hat{\sigma} \hat{q}(z)}{\sigma q(z)} - 1 \right| \leq C \left(r_n + \sqrt{\log n / (nh)} \right) \right) \geq 1 - \frac{c}{n}.$$

Then, to approximate $\mathbb{G}_n(z)$ by $\hat{\mathbb{G}}_n(z)$, we need two intermediate processes

$$\begin{aligned} \mathbb{G}_n^{(1)}(z) &:= \frac{1}{\sqrt{n}} \sum_{i=1}^n \left[\frac{\varepsilon_i I(|\varepsilon_i| \leq b_n)}{\sigma} \cdot \frac{K_h(X_{i1} - z) (\eta(\mathbf{X}_{i2}) - \tilde{\mathbf{w}}(z)^T \boldsymbol{\psi}(\mathbf{X}_i))}{q(z)} \right. \\ &\quad \left. - \mathbb{E} \left(\frac{\varepsilon_i I(|\varepsilon_i| \leq b_n)}{\sigma} \cdot \frac{K_h(X_{i1} - z) (\eta(\mathbf{X}_{i2}) - \tilde{\mathbf{w}}(z)^T \boldsymbol{\psi}(\mathbf{X}_i))}{q(z)} \right) \right], \end{aligned} \quad (7.9)$$

$$\mathbb{G}_n^{(2)}(z) := \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\xi_i \varepsilon_i I(|\varepsilon_i| \leq b_n) \cdot K_h(X_{i1} - z) (\eta(\mathbf{X}_{i2}) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{X}_i))}{\sigma q(z)}, \quad (7.10)$$

where $b_n := C' \sqrt{\log n}$ such that subgaussian random variables ε_i 's satisfy $\mathbb{P}(\max_{i \in [n]} \varepsilon_i \leq b_n) \geq 1 - n^{-1}$. The next two lemmas show how well $\mathbb{G}_n^{(1)}(z), \mathbb{G}_n^{(2)}(z)$ approximates $\mathbb{G}_n(z)$ and $\hat{\mathbb{G}}_n(z)$. Assume $\mathbf{X}_1^n = \{\mathbf{X}_1, \dots, \mathbf{X}_n\}$, $\varepsilon_1^n = \{\varepsilon_1, \dots, \varepsilon_n\}$ and $\mathbf{J}_n = (\mathbf{X}_1^n, \varepsilon_1^n)$. We denote

$$W_n := \sup_{z \in \mathcal{X}} \mathbb{G}_n(z), \quad W_n^{(1)} := \sup_{z \in \mathcal{X}} \mathbb{G}_n^{(1)}(z), \quad W_n^{(2)}(\mathbf{j}_n) := \sup_{z \in \mathcal{X}} \mathbb{G}_n^{(2)}(z) \quad \text{and} \quad \widehat{W}_n(\mathbf{j}_n) := \sup_{z \in \mathcal{X}} \hat{\mathbb{G}}_n(z),$$

for $\mathbf{j}_n \in \Omega := \mathbb{R}^n \times (\mathcal{X}^d)^n$. Then we have the following lemmas:

Lemma 7.10. Under conditions in Theorem 4.2, there exists constants c_1, c_2, c such that

$$\mathbb{P}(|W_n - W_n^{(1)}| > c_1 n^{-c}) \leq c_2 n^{-c}. \quad (7.11)$$

Lemma 7.11. Under conditions in Theorem 4.2, there exists constants c_1, c_2, c and a set $S_n \subset \Omega$ satisfying $\mathbb{P}(\mathbf{J}_n \in S_n) \geq 1 - \frac{3}{n}$. For any fixed $\mathbf{j}_n \in S_n$, we have

$$\mathbb{P}(|W_n^{(1)} - W_n^{(2)}(\mathbf{j}_n)| > c_1 n^{-c}) \leq c_2 n^{-c}. \quad (7.12)$$

Lemma 7.12. Under conditions in Theorem 4.2, there exists constants c_1, c_2, c and a set $S_n \subset \Omega$ satisfying $\mathbb{P}(\mathbf{J}_n \in S_n) \geq 1 - \frac{3}{n}$. For any fixed $\mathbf{j}_n \in S_n$, we have

$$\mathbb{P}(|W_n^{(2)}(\mathbf{j}_n) - \widehat{W}_n(\mathbf{j}_n)| > c_1 n^{-c}) \leq c_2 n^{-c}. \quad (7.13)$$

Now we can give the main proof concerning the convergence of our confidence band.

Proof of Theorem 4.2. Recall our confidence band is

$$\mathcal{C}_{n,\alpha}(z) := \tilde{f}(z) \pm \hat{c}_n(\alpha) n^{-1/2} \hat{\sigma} \hat{q}(z) / \hat{p}(z).$$

Firstly, we want to show the biased term in \tilde{f} is asymptotically negligible with respect to the bandwidth of $\mathcal{C}_{n,\alpha}(z)$. Applying Lemma B.8, Lemma B.9 and $\hat{p}(z) > b/4$, we could see the width of the band satisfying

$$2\hat{c}_n(\alpha) n^{-1/2} \hat{\sigma} \hat{q}(z) / \hat{p}(z) \geq C \frac{1}{\sqrt{nh}}. \quad (7.14)$$

Applying Lemma 7.7 and Lemma 7.8, the biased term in \tilde{f} is $I_1 + I_{21}$, and under conditions in Theorem 4.2, with probability $1 - \frac{c}{n}$,

$$I_1 + I_{21} \leq C \left(r_e \cdot \frac{m + \sqrt{m \log(dnh^{-1})}}{\sqrt{nh}} + r_a + h^2 + \sqrt{n^{-1}h \log n} \right) = o\left(\frac{1}{\sqrt{nh}}\right). \quad (7.15)$$

Combine (7.14) and (7.15), the bias in \tilde{f} is asymptotically negligible. Since $\tilde{f} - f^* = I_1 + I_{21} + I_{22}$, and

$$I_{22} = \frac{\frac{1}{n} \sum_{i=1}^n t(\mathbf{X}_i, z) \varepsilon_i}{\hat{p}(z)},$$

we will construct the confidence band for the numerator of I_{22} , $\frac{1}{n} \sum_{i=1}^n t(\mathbf{X}_i, z) \varepsilon_i$. After normalizing the variance of the numerator, we get

$$\mathbb{G}_n(z) := \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\varepsilon_i t(\mathbf{X}_i, z)}{\sigma q(z)}.$$

We hope to approximate $\mathbb{G}_n(z)$ with $\widehat{\mathbb{G}}_n(z)$,

$$\widehat{\mathbb{G}}_n(z) := \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\xi_i \widehat{\varepsilon}_i t(\mathbf{X}_i, z)}{\widehat{\sigma} \hat{q}(z)}.$$

We use the Corollary 3.1 in Chernozhukov et al. (2014a) to establish this approximation. Corollary 3.1 of Chernozhukov et al. (2014a) provides sufficient conditions for the confidence band to be asymptotically honest. Specifically, we verify the following high-level conditions (Let C_1 be some positive constant):

H1 (Gaussian approximation). There exists a sequence of random variables W_n^0 such that $W_n^0 \stackrel{d}{=} \sup_{z \in \mathcal{X}} \mathbb{G}_n^0(z)$, where \mathbb{G}_n^0 is a Gaussian process in $\ell^\infty(\mathcal{X})$. Furthermore, $\mathbb{E}[\sup_{z \in \mathcal{X}} \mathbb{G}_n^0(z)] \leq c_1 \sqrt{\log n}$ and

$$\sup_{f^* \in \mathcal{F}} \mathbb{P}(|W_n - W_n^0| > \epsilon_{1n}) \leq \delta_{1n},$$

for some ϵ_{1n} and δ_{1n} bounded from above by $c_1 n^{-c}$.

H2 (Anti-concentration). For any $\epsilon > 0$, the anti-concentration inequality

$$\sup_{x \in \mathbb{R}} \mathbb{P}\left(\left|\sup_{z \in \mathcal{X}} |\mathbb{G}_n(z)| - x\right| \leq \epsilon\right) \leq C_1 \epsilon \sqrt{\log n}$$

holds.

H3 (Estimation error of $\hat{c}_n(\alpha)$). Let $c_n(\alpha)$ be the $(1 - \alpha)$ -quantile of $W_n^0 = \sup_{z \in \mathcal{X}} G_n(z)$ and $\hat{c}_n(\alpha)$ be the $(1 - \alpha)$ -quantile of $\widehat{W}_n = \sup_{z \in \mathcal{X}} \widehat{G}_n(z)$. For some τ_n , ε_{2n} and δ_{2n} bounded from above by $c_1 n^{-c}$, we have

$$\sup_{f^* \in \mathcal{F}} \mathbb{P}\left(\hat{c}_n(\alpha) < c_n(\alpha + \tau_n) - \varepsilon_{2n}\right) \leq \delta_{2n},$$

and

$$\sup_{f^* \in \mathcal{F}} \mathbb{P}\left(\hat{c}_n(\alpha) > c_n(\alpha - \tau_n) + \varepsilon_{2n}\right) \leq \delta_{2n}.$$

H4 (Estimation error of $\hat{\sigma}_n(\cdot)$). For some ε_{3n} and δ_{3n} bounded from above by $c_1 n^{-c}$, we have

$$\sup_{f^* \in \mathcal{F}} \mathbb{P}\left(\sup_{z \in \mathcal{X}} \left|\frac{\widehat{\sigma} \widehat{q}(z)}{\sigma q(z)} - 1\right| > \varepsilon_{3n}\right) \leq \delta_{3n}.$$

First, we verify the condition **H1**. Under conditions in Theorem 4.2, Lemma 7.10 gives us the the rate of difference between $W_n^{(1)}$ and W_n

$$\mathbb{P}(|W_n - W_n^{(1)}| > c_1 n^{-c}) \leq c_2 n^{-c}. \quad (7.16)$$

Then we will construct a good Gaussian approximation to the intermediate random variable $W_n^{(1)}$. As for $\mathbb{G}_n^{(1)}$, the function class of interest is:

$$\mathcal{S}_h := \left\{ S_z(\varepsilon, \mathbf{X}) := \frac{\varepsilon I(|\varepsilon| \leq b_n)}{\sigma} \cdot \frac{K_h(X_1 - z)(\eta(\mathbf{X}_2) - \widetilde{\mathbf{w}}(z)^\top \boldsymbol{\psi}(\mathbf{X}))}{q(z)} \middle| z \in \mathcal{X} \right\},$$

where $b_n = C\sqrt{\log n}$. We are going to prove \mathcal{S}_h is a VC-type class (Definition C.4). As for the upper bound,

$$\|S_z(\varepsilon, \mathbf{X})\|_\infty \leq \frac{C\sqrt{\log n}(L+1)B}{\sqrt{h}\sigma b} = C\sqrt{\frac{\log n}{h}}.$$

Therefore, one parameter for VC-type class is $b_{VC} = C\sqrt{\frac{\log n}{h}}$. According to Lemma C.2, for any $\tau > 0$ and all finite measure \mathbb{Q} ,

$$N(\mathcal{S}_h, L_2(\mathbb{Q}), b_{VC}\tau) \leq \left(\frac{2\|K\|_{tv}A}{\tau}\right)^4,$$

where A is an absolute constant. Thus the parameters of \mathcal{S}_h are

$$\begin{cases} v_{VC} = 4 = \mathcal{O}(1), \\ a_{VC} = 2\|K\|_{tv}A = \mathcal{O}(1), \\ b_{VC} = \mathcal{O}(\sqrt{\log n \cdot h^{-1}}), \\ \sigma_{VC}^2 = \mathbb{E}S_z^2(\varepsilon, \mathbf{X}) = \mathcal{O}(1), \end{cases}$$

and

$$K_h := Av_{VC}(\log n \vee \log(a_{VC}b_{VC}/\sigma_{VC})) = \mathcal{O}(\log n).$$

Applying Lemma C.5, there exists a tight Gaussian random element $B = \{B(g) : g \in S_h\}$ such that

$$\mathbb{E}[B(g_1)B(g_2)] = \mathbb{E}[g_1(\mathbf{X}_1, \varepsilon_1)g_2(\mathbf{X}_1, \varepsilon_1)] - \mathbb{E}[g_1(\mathbf{X}_1, \varepsilon_1)]\mathbb{E}[g_2(\mathbf{X}_1, \varepsilon_1)], \quad (7.17)$$

for all $g_1, g_2 \in S_h$. Let $\mathbb{G}_n^0(z) = B(S_z)$ and $W^0 := \sup_z B(S_z)$. We have for any $\gamma \in (0, 1)$,

$$\mathbb{P}\left(|W_n^{(1)} - W^0| > C \cdot \left(\frac{\log^2 n}{\gamma^{1/2}hn^{1/2}} + \frac{\log^{5/4} n}{\gamma^{1/2}n^{1/4}h^{3/4}} + \frac{\log n}{\gamma^{1/3}n^{1/6}h^{2/3}}\right)\right) \leq A'\left(\gamma + \frac{\log n}{n}\right),$$

where A' is an absolute constant. Under conditions in Theorem 4.2, since $h \asymp n^{-\delta_2}$, $\delta_2 < 1/4$, there exists c such that $\frac{1}{n^{1/6}h^{2/3}} \leq n^{-c}$. Hence, with $\gamma = n^{-3c/2}$, there exist constants c_1, c_2 such that

$$\mathbb{P}\left(|W_n^{(1)} - W^0| > c_1n^{-c}\right) \leq c_2n^{-c}. \quad (7.18)$$

Then combining (7.16) and (7.18), we have

$$\mathbb{P}\left(|W_n - W^0| > c_1n^{-c}\right) \leq c_2n^{-c}.$$

Now we are going to verify $\mathbb{E}[\sup_{z \in \mathcal{X}} \mathbb{G}_n^0(z)] = \mathbb{E}[\sup_{z \in \mathcal{X}} B(S_z)] \leq c_1\sqrt{\log n}$. With (7.17), we know the covering number for $\mathbb{G}_n^0(z)$ is the same as $\mathbb{G}_n^{(1)}(z)$. Therefore, applying Dudley's inequality (Dudley, 2010), we have

$$\begin{aligned} \mathbb{E}[\sup_{z \in \mathcal{X}} \mathbb{G}_n^0(z)] &\leq C \int_0^\infty \sqrt{\log N(\mathcal{S}_h, L^2(\mathbb{Q}), \tau)} d\tau \\ &\leq C \int_0^C \sqrt{2 \log C \frac{\log n}{h} + 4 \log \frac{1}{\tau}} d\tau = C\sqrt{\log n}. \end{aligned}$$

The condition **H2** follows from **H1** and the anti-concentration inequality in Corollary 2.1 of Chernozhukov et al. (2014a).

Next, we verify the condition **H3**. Recall that $\widehat{c}_n(\alpha)$ is the $(1 - \alpha)$ -quantile of \widehat{W}_n . From the approximation result in (7.12), (7.13) and (7.18), there exists a subset S_n , such that $\mathbb{P}(\mathbf{J}_n \in S_n) \geq 1 - \frac{3}{n}$. For any $\mathbf{j}_n \in S_n$, we have

$$\mathbb{P}\left(|\widehat{W}_n(\mathbf{j}_n) - W_n^0| > c_1n^{-c} + \epsilon'_{1n}\right) \leq c_2n^{-c} + \delta'_{1n},$$

and hence with another two constants c'_1 and c'_2 , we have,

$$\mathbb{P}(|\widehat{W}_n(\mathbf{j}_n) - W_n^0| > c'_1 n^{-c}) \leq c'_2 n^{-c}.$$

Therefore, we can bound the probability

$$\begin{aligned} \mathbb{P}(W_n^0 \leq \widehat{c}_n(\alpha) + c'_1 n^{-c}) &\geq \mathbb{P}(\widehat{W}_n(\mathbf{j}_n) \leq \widehat{c}_n(\alpha)) - \mathbb{P}(|\widehat{W}_n(\mathbf{j}_n) - W_n^0| > c'_1 n^{-c}) \\ &\geq 1 - \alpha - c'_2 n^{-c}, \end{aligned} \quad (7.19)$$

which implies that the estimated quintile has the following lower bound

$$\widehat{c}_n(\alpha) \geq c_n(\alpha + c'_2 n^{-c}) - c'_1 n^{-c}. \quad (7.20)$$

Similarly, we also have $\widehat{c}_n(\alpha) \leq c_n(\alpha - c'_2 n^{-c}) + c'_1 n^{-c}$. By setting $\tau = c'_2 n^{-c}$, $\epsilon_{2n} = c'_1 n^{-c}$, and $\delta_{2n} = 0$, condition **H3** holds.

Finally, using Lemma 7.9, the condition **H4** is satisfied for $\epsilon_{3n} = C(r_n + \sqrt{\log n/(nh)}) = o(n^{-c})$ and $\delta_{3n} = 8/n$. Since the high-level condition **H1** – **H4** are verified, the result follows from Corollary 3.1 in Chernozhukov et al. (2014a).

□

References

- AVALOS, M., GRANDVALET, Y. and AMBROISE, C. (2007). Parsimonious additive models. *Computational statistics & data analysis* **51** 2851–2870.
- BELLONI, A., CHERNOZHUKOV, V., CHETVERIKOV, D. and WEI, Y. (2018). Uniformly valid post-regularization confidence regions for many functional parameters in z-estimation framework. *Ann. Statist.* **46** 3643–3675.
- BELLONI, A., CHERNOZHUKOV, V., FERNÁNDEZ-VAL, I. and HANSEN, C. (2017). Program evaluation and causal inference with high-dimensional data. *Econometrica* **85** 233–298.
- BELLONI, A., CHERNOZHUKOV, V. and HANSEN, C. (2014a). Inference on treatment effects after selection among high-dimensional controls. *The Review of Economic Studies* **81** 608–650.
- BELLONI, A., CHERNOZHUKOV, V. and KATO, K. (2013a). Robust inference in high-dimensional approximately sparse quantile regression models. *arXiv preprint arXiv:1312.7186*.
- BELLONI, A., CHERNOZHUKOV, V. and KATO, K. (2014b). Uniform post-selection inference for least absolute deviation regression and other z-estimation problems. *Biometrika* asu056.
- BELLONI, A., CHERNOZHUKOV, V. and WEI, Y. (2013b). Honest confidence regions for a regression parameter in logistic regression with a large number of controls. *arXiv preprint arXiv:1304.3969*.
- BERNSTEIN, S. (1964). On a modification of chebyshev’s inequality and on the error in laplace formula. *Collected Works, Izd-vo’Nauka’, Moscow (in Russian)* **4** 71–80.
- BERTIN, K., LECUÉ, G. ET AL. (2008). Selection of variables and dimension reduction in high-dimensional non-parametric regression. *Electronic Journal of Statistics* **2** 1224–1241.

- BORELL, C. (1975). The Brunn-Minkowski inequality in Gauss space. *Invent. Math.* **30** 207–216.
- BOUSQUET, O. (2002). A bennett concentration inequality and its application to suprema of empirical processes. *Comptes Rendus Mathematique* **334** 495–500.
- BÜHLMANN, P. and VAN DE GEER, S. (2011). *Statistics for high-dimensional data: methods, theory and applications*. Springer Science & Business Media.
- CAI, T. T. and GUO, Z. (2017). Confidence intervals for high-dimensional linear regression: minimax rates and adaptivity. *Ann. Statist.* **45** 615–646.
- CAI, T. T. and GUO, Z. (2020). Semisupervised inference for explained variance in high dimensional linear regression and its applications. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **82** 391–419.
- CHATTERJEE, A., LAHIRI, S. ET AL. (2013). Rates of convergence of the adaptive lasso estimators to the oracle distribution and higher order refinements by the bootstrap. *The Annals of Statistics* **41** 1232–1259.
- CHERNOZHUKOV, V., CHETVERIKOV, D., KATO, K. ET AL. (2013). Gaussian approximations and multiplier bootstrap for maxima of sums of high-dimensional random vectors. *The Annals of Statistics* **41** 2786–2819.
- CHERNOZHUKOV, V., CHETVERIKOV, D., KATO, K. ET AL. (2014a). Anti-concentration and honest, adaptive confidence bands. *The Annals of Statistics* **42** 1787–1818.
- CHERNOZHUKOV, V., CHETVERIKOV, D., KATO, K. ET AL. (2014b). Gaussian approximation of suprema of empirical processes. *The Annals of Statistics* **42** 1564–1597.
- COMMINGES, L., DALALYAN, A. S. ET AL. (2012). Tight conditions for consistency of variable selection in the context of high dimensionality. *The Annals of Statistics* **40** 2667–2696.
- DALALYAN, A., INGSTER, Y. and TSYBAKOV, A. B. (2014). Statistical inference in compound functional models. *Probability Theory and Related Fields* **158** 513–532.
- DUDLEY, R. M. (2010). The sizes of compact subsets of Hilbert space and continuity of Gaussian processes [mr0220340]. In *Selected works of R. M. Dudley*. Sel. Works Probab. Stat., Springer, New York, 125–165.
- FAN, J. (1993). Local linear regression smoothers and their minimax efficiencies. *The Annals of Statistics* 196–216.
- FAN, J., MA, Y. and DAI, W. (2014). Nonparametric independence screening in sparse ultra-high-dimensional varying coefficient models. *Journal of the American Statistical Association* **109** 1270–1284.
- FAN, J. and ZHANG, W. (1999). Statistical estimation in varying coefficient models. *Annals of Statistics* 1491–1518.
- FARRELL, M. H. (2013). Robust inference on average treatment effects with possibly more covariates than observations. *Available at SSRN 2324292* .

- FRIEDMAN, J. H. and STUETZLE, W. (1981). Projection pursuit regression. *Journal of the American statistical Association* **76** 817–823.
- GINÉ, E. and NICKL, R. (2009). An exponential inequality for the distribution function of the kernel density estimator, with applications to adaptive estimation. *Probability Theory and Related Fields* **143** 569–596.
- GREGORY, K., MAMMEN, E. and WAHL, M. (2016). Optimal estimation of sparse high-dimensional additive models. *Preprint* .
- GUO, Z. and ZHANG, C.-H. (2019a). Extreme nonlinear correlation for multiple random variables and stochastic processes with applications to additive models. *arXiv preprint arXiv:1904.12897* .
- GUO, Z. and ZHANG, C.-H. (2019b). Local inference in additive models with decorrelated local linear estimator. *arXiv preprint arXiv:1907.12732* .
- HASTIE, T. and TIBSHIRANI, R. (1993). Varying-coefficient models. *Journal of the Royal Statistical Society. Series B (Methodological)* 757–796.
- HASTIE, T. J. and TIBSHIRANI, R. J. (1990). *Generalized additive models*, vol. 43. CRC Press.
- HUANG, J., HOROWITZ, J. L. and WEI, F. (2010). Variable selection in nonparametric additive models. *Annals of statistics* **38** 2282.
- JAVANMARD, A. and MONTANARI, A. (2013). Nearly optimal sample size in hypothesis testing for high-dimensional regression. *arXiv preprint arXiv:1311.0274* .
- JAVANMARD, A. and MONTANARI, A. (2014). Confidence intervals and hypothesis testing for high-dimensional regression. *The Journal of Machine Learning Research* **15** 2869–2909.
- KATO, K. (2012). Two-step estimation of high dimensional additive models. *arXiv preprint arXiv:1207.5313* .
- KOLTCHINSKII, V. (2011). *Oracle Inequalities in Empirical Risk Minimization and Sparse Recovery Problems: Ecole d’Eté de Probabilités de Saint-Flour XXXVIII-2008*, vol. 2033. Springer Science & Business Media.
- KOLTCHINSKII, V., YUAN, M. ET AL. (2010). Sparsity in multiple kernel learning. *The Annals of Statistics* **38** 3660–3695.
- KOZBUR, D. (2013). Inference in additively separable models with a high dimensional component. Tech. rep., Working Paper. Booth School of Business, University of Chicago, Chicago, IL.
- LAFFERTY, J. and WASSERMAN, L. (2008). Rodeo: sparse, greedy nonparametric regression. *The Annals of Statistics* 28–63.
- LEE, J. D., SUN, D. L., SUN, Y. and TAYLOR, J. E. (2013). Exact post-selection inference with the lasso. *arXiv preprint arXiv:1311.6238* .
- LIAN, H. (2012). Variable selection for high-dimensional generalized varying-coefficient models. *Statistica Sinica* **22** 1563.

- LIN, Y., ZHANG, H. H. ET AL. (2006). Component selection and smoothing in multivariate nonparametric regression. *The Annals of Statistics* **34** 2272–2297.
- LINERO, A. R. (2018). Bayesian regression trees for high-dimensional prediction and variable selection. *J. Amer. Statist. Assoc.* **113** 626–636.
- LIU, H., YU, B. ET AL. (2013). Asymptotic properties of lasso+ mls and lasso+ ridge in sparse high-dimensional linear regression. *Electronic Journal of Statistics* **7** 3124–3169.
- LOCKHART, R., TAYLOR, J., TIBSHIRANI, R. J. and TIBSHIRANI, R. (2014). A significance test for the lasso. *Annals of statistics* **42** 413.
- LOPES, M. (2014). A residual bootstrap for high-dimensional regression with near low-rank designs. In *Advances in Neural Information Processing Systems*.
- LOU, Y., BIEN, J., CARUANA, R. and GEHRKE, J. (2014). Sparse partially linear additive models. *arXiv preprint arXiv:1407.4729* .
- LU, J., KOLAR, M. and LIU, H. (2019). Kernel meets sieve: Post-regularization confidence bands for sparse additive model. *Journal of the American Statistical Association* 1–40.
- MA, C., LU, J. and LIU, H. (2017). Inter-subject analysis: Inferring sparse interactions with dense intra-graphs. *arXiv preprint arXiv:1709.07036* .
- MEIER, L., VAN DE GEER, S., BÜHLMANN, P. ET AL. (2009). High-dimensional additive modeling. *The Annals of Statistics* **37** 3779–3821.
- MEINSHAUSEN, N. (2014). Group bound: confidence intervals for groups of variables in sparse high dimensional regression without assumptions on the design. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* .
- MEINSHAUSEN, N. and BÜHLMANN, P. (2010). Stability selection. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **72** 417–473.
- MEINSHAUSEN, N., MEIER, L. and BÜHLMANN, P. (2009). P-values for high-dimensional regression. *Journal of the American Statistical Association* **104**.
- NING, Y. and LIU, H. (2014). Sparc: Optimal estimation and asymptotic inference under semi-parametric sparsity. *arXiv preprint arXiv:1412.2295* .
- PETERSEN, A., WITTEN, D. and SIMON, N. (2014). Fused lasso additive model. *arXiv preprint arXiv:1409.5391* .
- RASKUTTI, G., WAINWRIGHT, M. J. and YU, B. (2012). Minimax-optimal rates for sparse additive models over kernel classes via convex programming. *The Journal of Machine Learning Research* **13** 389–427.
- RAVIKUMAR, P., LAFFERTY, J., LIU, H. and WASSERMAN, L. (2009). Sparse additive models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **71** 1009–1030.
- ROSASCO, L., VILLA, S., MOSCI, S., SANTORO, M. and VERRI, A. (2013). Nonparametric sparsity and regularization. *The Journal of Machine Learning Research* **14** 1665–1714.

- SARDY, S. and TSENG, P. (2004). Amlet, ramlet, and gamlet: automatic nonlinear fitting of additive models, robust and generalized, with wavelets. *Journal of Computational and Graphical Statistics* **13** 283–309.
- SHAH, R. D. and SAMWORTH, R. J. (2013). Variable selection with error control: another look at stability selection. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **75** 55–80.
- STONE, C. J. (1985). Additive regression and other nonparametric models. *The annals of Statistics* 689–705.
- TAYLOR, J., LOCKHART, R., TIBSHIRANI, R. J. and TIBSHIRANI, R. (2014). Exact post-selection inference for forward stepwise and least angle regression. *arXiv preprint arXiv:1401.3889* .
- VAART, A. v. D. and WELLNER, J. A. (1997). Weak convergence and empirical processes with applications to statistics. *Journal of the Royal Statistical Society-Series A Statistics in Society* **160** 596–608.
- VAN DE GEER, S., BÜHLMANN, P., RITOV, Y., DEZEURE, R. ET AL. (2014). On asymptotically optimal confidence regions and tests for high-dimensional models. *The Annals of Statistics* **42** 1166–1202.
- VAN DER VAART, A. W. and WELLNER, J. A. (1996). *Weak Convergence*. Springer.
- WAHL, M. (2014). Variable selection in high-dimensional additive models based on norms of projections. *arXiv preprint arXiv:1406.0052* .
- WASSERMAN, L. (2006). *All of nonparametric statistics*. Springer Science & Business Media.
- WASSERMAN, L. and ROEDER, K. (2009). High dimensional variable selection. *Annals of statistics* **37** 2178.
- XU, M., CHEN, M. and LAFFERTY, J. (2014). Faithful variable screening for high-dimensional convex regression. *arXiv preprint arXiv:1411.1805* .
- YANG, Y., TOKDAR, S. T. ET AL. (2015). Minimax-optimal nonparametric regression in high dimensions. *The Annals of Statistics* **43** 652–674.
- YAO, Y. and ZHANG, C.-H. (2020). Adaptive estimation in high-dimensional additive models with multi-resolution group lasso. *arXiv preprint arXiv:2011.06765* .
- YUAN, M. and ZHOU, D.-X. (2016). Minimax optimal rates of estimation in high dimensional additive models. *Ann. Statist.* **44** 2564–2593.
- ZHANG, H. H., CHENG, G. and LIU, Y. (2011). Linear or nonlinear? automatic structure discovery for partially linear models. *Journal of the American Statistical Association* **106**.
- ZHOU, S., SHEN, X., WOLFE, D. ET AL. (1998). Local asymptotics for regression splines and confidence regions. *The annals of statistics* **26** 1760–1782.

A Auxiliary Lemmas for Estimation Results

Throughout the appendix we use C to denote an absolute constant, whose value may change from line to line. The following technical result will be useful throughout the proofs.

Lemma A.1. With probability $1 - \frac{1}{n}$, there exist a constant C such that

$$\sup_z \frac{1}{n} \sum_{i=1}^n K_h(X_{i1} - z) \leq C.$$

Proof. Consider the process and the function class of interest:

$$U_n(z) = \frac{1}{n} \sum_{i=1}^n K_h(X_{i1} - z) - \mathbb{E}K_h(X_{i1} - z)$$

$$\mathcal{S}_h = \{S_z(x_1, \mathbf{x}_2) := K_h(x_1 - z) | z \in \mathcal{X}\}.$$

According to Lemma C.2,

$$N(\mathcal{S}_h, L^2(\mathbb{P}_n), \epsilon) \leq \left(\frac{2\|K\|_{TV}A}{h\epsilon} \right)^4,$$

where A is an absolute constant. Then we can upper bound the variance of the process as:

$$\sigma_p^2 := \sup_z \mathbb{E}K_h(X_{i1} - z)^2 = \mathcal{O}(h^{-1}).$$

Therefore, according to Lemma C.1,

$$\mathbb{E} \left[\sup_{z \in \mathcal{X}} U_n(z) \right] \leq Cn^{-1/2} \sigma_p \sqrt{\log \frac{2\|K\|_{TV}}{h\sigma_p}} \leq C \sqrt{\frac{\log h^{-1}}{nh}}. \quad (\text{A.1})$$

Since $\sup_z |K_h(x_1 - z)| \leq U := Ch^{-1}$, by Lemma C.3, we have

$$\mathbb{P} \left(\sup_z U_n(z) \geq \mathbb{E} \left[\sup_z U_n(z) \right] + t \sqrt{2 \left(\sigma_p^2 + 2U \mathbb{E} \left[\sup_z U_n(z) \right] \right)} + \frac{2Ut^2}{3} \right) \leq \exp(-nt^2). \quad (\text{A.2})$$

Combine (A.1) and (A.2), set $t = \sqrt{\frac{\log n}{n}}$, then, with probability at least $1 - \frac{1}{n}$,

$$\begin{aligned} \sup_z U_n(z) &\leq \mathbb{E} \left[\sup_z U_n(z) \right] + \sqrt{2\sigma_p^2 + 4U \mathbb{E} \left[\sup_z U_n(z) \right]} \cdot \sqrt{\frac{\log n}{n}} + \frac{2U}{3} \cdot \frac{\log n}{n} \\ &\leq C \sqrt{\frac{\log n}{nh}} \leq C, \end{aligned}$$

where we used the inequality $\sqrt{xy} \leq 2^{-1}(x + y)$ and the assumption that $(nh)^{-1} \log n = \mathcal{O}(1)$. \square

The estimation error for the plug-in estimator \hat{f} comes from four sources: (1) noise, (2) approximation error by finite B-spline bases, (3) bias introduced by the local constant estimator, and (4) estimation error from the initial estimator $\hat{\gamma}$.

The following lemma provides the rate of the B-spline approximation error. It shows how the number of B-spline basis functions m influences the final estimation rate.

Proof of Lemma 7.1. By Assumption 4.1, for any $i \in [n]$, $|\eta(\mathbf{X}_{i2})| \leq \|h\|_\infty$ and $|\delta_i| = |g(X_{i1}, \mathbf{X}_{i2}) - \gamma^{*\top} \psi(\mathbf{X}_i)| \leq \|g - \gamma^{*\top} \psi\|_\infty \leq r_a$ are uniformly bounded random variables. First,

$$\sup_{z \in \mathcal{X}} \frac{1}{n} |\mathbf{h}^\top \mathbf{W}_z \boldsymbol{\delta}| \leq \sup_{z \in \mathcal{X}} \frac{1}{n} \sum_{i=1}^n K_h(X_{i1} - z) \eta(\mathbf{X}_{i2}) |\delta_i|.$$

We denote the empirical process and the function class of interest as

$$U_{n,\boldsymbol{\delta}}(z) = \frac{1}{n} \sum_{i=1}^n \eta(\mathbf{X}_{i2}) K_h(X_{i1} - z) |\delta_i| - \mathbb{E}[\eta(\mathbf{X}_{12}) K_h(X_{11} - z) |\delta_1|],$$

$$\mathcal{S}_{h,\boldsymbol{\delta}} = \{s_z(x_1, \mathbf{x}_2) = \eta(\mathbf{x}_2) h^{-1} K(h^{-1}(x_1 - z)) |\delta(x_1, \mathbf{x}_2)| \mid z \in \mathcal{X}\},$$

where $\delta(\mathbf{x}) = g(\mathbf{x}) - \gamma^{*\top} \psi(\mathbf{x})$ satisfies $|\delta(\mathbf{x})| \leq r_a$ for any $\mathbf{x} = (x_1, \mathbf{x}_2)$. For simplicity, the notation $U_n(z)$ and \mathcal{S}_h is used instead of $U_{n,\boldsymbol{\delta}}(z)$ and $\mathcal{S}_{h,\boldsymbol{\delta}}$, when it causes no confusion. Let $\mathcal{F}_h = \{h^{-1} K(h^{-1}(\cdot - z)) \mid z \in \mathcal{X}\}$ and let $\|K\|_{\text{TV}}$ be the total variation of $K(\cdot)$. From Lemma C.2, for all probability measures Q on \mathcal{X} ,

$$N(\mathcal{F}_h, L^2(Q), \epsilon) \leq \left(\frac{2\|K\|_{\text{TV}} A}{h\epsilon} \right)^4,$$

where A is an absolute constant and $0 < \epsilon < 1$. Let $\tilde{\mathcal{F}}_h$ be an $\|h\|_\infty^{-1} r_a^{-1} \epsilon$ -cover of \mathcal{F}_h with respect to $\hat{\mathbb{P}}_{X_1} = n^{-1} \sum_{i=1}^n \delta_{X_{i1}}$. We then construct an ϵ -covering for \mathcal{S}_h with respect to $\mathbb{P}_n = n^{-1} \sum_{i=1}^n \delta_{X_{i1}, \mathbf{X}_{i2}}$ as

$$\tilde{\mathcal{S}}_h = \{\eta(\mathbf{x}_2) f_1(x_1) |\delta(x_1, \mathbf{x}_2)| \mid f_1 \in \tilde{\mathcal{F}}_h\}.$$

For a function $s_z = \eta(\mathbf{X}_2) h^{-1} K(h^{-1}(x_1 - z)) |\delta(x_1, \mathbf{x}_2)| \in \mathcal{S}_h$, let

$$s_{\tilde{z}} = \eta(\mathbf{x}_2) h^{-1} K(h^{-1}(x_1 - \tilde{z})) |\delta(x_1, \mathbf{x}_2)| \in \tilde{\mathcal{S}}_h$$

be the corresponding element in the cover, that is $h^{-1} K(h^{-1}(x_1 - \tilde{z})) \in \tilde{\mathcal{F}}_h$ is the corresponding element in the cover for $h^{-1} K(h^{-1}(x_1 - z)) \in \mathcal{F}_h$. Then

$$\begin{aligned} \|s_z - s_{\tilde{z}}\|_{L^2(\mathbb{P}_n)}^2 &= n^{-1} \sum_{i=1}^n [(K_h(X_{i1} - z) - K_h(X_{i1} - \tilde{z})) \eta(\mathbf{X}_{i2}) |\delta(X_{i1}, \mathbf{X}_{i2})|]^2 \\ &\leq n^{-1} \sum_{i=1}^n [K_h(X_{i1} - z) - K_h(X_{i1} - \tilde{z})]^2 \|\eta\|_\infty^2 r_a^2 \leq \epsilon^2. \end{aligned}$$

Therefore, the covering number can be bounded by

$$N(\mathcal{S}_h, L^2(\mathbb{P}_n), \epsilon) \leq N(\mathcal{F}_h, L^2(\hat{\mathbb{P}}_{X_1}), \|\eta\|_\infty^{-1} r_a^{-1} \epsilon) \leq \left(\frac{2\|K\|_{\text{TV}} A \|\eta\|_\infty r_a}{h\epsilon} \right)^4. \quad (\text{A.3})$$

The variance of the process

$$\sigma_{\text{P}}^2 := \sup_{z \in \mathcal{X}} \mathbb{E}[(\eta(\mathbf{X}_{12}) K_h(X_{11} - z) |\delta_1|)^2] \leq \|\eta\|_\infty^2 r_a^2 \mathbb{E}[K_h^2] = \mathcal{O}(r_a^2 h^{-1}).$$

Applying Lemma C.1, we have

$$\mathbb{E} \left[\sup_{z \in \mathcal{X}} U_n(z) \right] \leq C n^{-1/2} \sigma_P \sqrt{\log \left(\frac{2 \|K\|_{\text{TV}} A \|\eta\|_{\infty} r_a}{h \sigma_P} \right)} = \mathcal{O} \left(r_a \sqrt{\frac{\log(h^{-1})}{nh}} \right), \quad (\text{A.4})$$

where in the last step we use $\sigma_P = \mathcal{O}(r_a h^{-1/2})$. Next, we bound the supremum of the expectation

$$\begin{aligned} \mathbb{E} [\eta(\mathbf{X}_{12}) K_h(X_{11} - z) \delta_1] &= \int \eta(\mathbf{X}_{12}) K_h(X_{11} - z) |\delta_1| d\mathbb{P} \\ &\leq \|\eta\|_{\infty} r_a \int K_h(x - z) p_{X_1}(x) dx = \mathcal{O}(r_a), \end{aligned} \quad (\text{A.5})$$

where in the last step we used the assumption that $\|p_{X_1}\|_{\infty}$ is bounded from above. Using Lemma C.3, we have

$$\mathbb{P} \left(\sup_{z \in \mathcal{X}} U_n(z) \geq \mathbb{E} \left[\sup_{z \in \mathcal{X}} U_n(z) \right] + t \sqrt{2 \left(\sigma_P^2 + 2U \mathbb{E} \left[\sup_{z \in \mathcal{X}} U_n(z) \right] \right)} + \frac{2Ut^2}{3} \right) \leq \exp(-nt^2), \quad (\text{A.6})$$

where $\sup_{s_z \in \mathcal{S}_h} \|s_z\|_{\infty} \leq U := h^{-1} \|K\|_{\infty} \|\eta\|_{\infty} r_a = \mathcal{O}(r_a h^{-1})$.

Combining (A.4), (A.5) and (A.6), with $t = \sqrt{\log n/n}$, we have with probability at least $1 - 1/n$,

$$\begin{aligned} \sup_{z \in \mathcal{X}} \frac{1}{n} |\mathbf{h}^T \mathbf{W}_z \boldsymbol{\delta}| &\leq \sup_{z \in \mathcal{X}} \frac{1}{n} \sum_{i=1}^n \eta(\mathbf{X}_{i2}) K_h(X_{i1} - z) |\delta_i| \\ &\leq \sup_{z \in \mathcal{X}} U_n(z) + \sup_{z \in \mathcal{X}} \mathbb{E} [\eta(\mathbf{X}_{12}) K_h(X_{11} - z) \delta_1] = \mathcal{O}(r_a), \end{aligned}$$

where in the last step we used the assumption that $\log n/(nh) = \mathcal{O}(1)$. \square

Proof of Lemma 7.2. We rewrite $\sup_{z \in \mathcal{X}} n^{-1} |\mathbf{h}^T \mathbf{W}_z \boldsymbol{\xi}_z|$ as

$$\sup_{z \in \mathcal{X}} \frac{1}{n} \left| \sum_{i=1}^n K_h(X_{i1} - z) \eta(\mathbf{X}_{i2})^2 (f(X_{i1}) - f(z)) \right|.$$

Denote the empirical process and the function class of interest here as

$$\begin{aligned} U_{n,\boldsymbol{\xi}}(z) &= \frac{1}{n} \sum_{i=1}^n K_h(X_{i1} - z) \eta(\mathbf{X}_{i2})^2 (f(X_{i1}) - f(z)) - \mathbb{E} [K_h(X_{11} - z) \eta(\mathbf{X}_{12})^2 (f(X_{11}) - f(z))], \\ \mathcal{S}_{h,\boldsymbol{\xi}} &= \{s_z(x_1, x_2) = h^{-1} K(h^{-1}(x_1 - z)) \eta^2(\mathbf{x}_2) (f(x_1) - f(z)) \mid z \in \mathcal{X}\}. \end{aligned}$$

For simplicity, we will use the notation $U_n(z)$ and \mathcal{S}_h instead, when it leads to no confusion. We note that \mathcal{S}_h is a subset of $\mathcal{F}_{h,1} \times \mathcal{F}_{h,2}$, where $\mathcal{F}_{h,1} := \{h^{-1} K(h^{-1}(x_1 - z)) \eta^2(\mathbf{x}_2) \mid z \in \mathcal{X}\}$ and $\mathcal{F}_{h,2} := \{f(x_1) - f(z) \mid z \in \mathcal{X}\}$. A similar argument to that used to obtain (A.3) gives us

$$N(\mathcal{F}_{h,1}, L^2(\mathbb{P}_n), \epsilon) \leq \left(\frac{2 \|K\|_{\text{TV}} A \|\eta\|_{\infty}^2}{h \epsilon} \right)^4.$$

For $f \in \mathcal{H}(2, L)$ we have $|(f(x_1) - f(z_1)) - (f(x_1) - f(z_2))| \leq \|f'\|_\infty |z_1 - z_2|$, and

$$N(\mathcal{F}_{h,2}, L^2(\mathbb{P}_n), \epsilon) \leq N(\mathcal{X}, d(\cdot, \cdot), \|f'\|_\infty^{-1} \epsilon) \leq \frac{\|f'\|_\infty \text{Diam}(\mathcal{X})}{\epsilon}.$$

By Lemma C.7 and Lemma C.8, we obtain

$$N(\mathcal{S}_h, L^2(\mathbb{P}_n), \epsilon) \leq \|f'\|_\infty \text{Diam}(\mathcal{X}) \left(\frac{2\|K\|_{\text{TV}} A \|\eta\|_\infty^2}{h} \right)^4 \left(\frac{\|\mathcal{F}_{h,1}\|_\infty \vee \|\mathcal{F}_{h,2}\|_\infty}{\epsilon/2} \right)^5 = O(h^{-9} \epsilon^{-5}).$$

Furthermore, we bound the variance by applying the Taylor expansion as follows

$$\begin{aligned} \sigma_P^2 &:= \sup_{z \in \mathcal{X}} \mathbb{E} \left[\left(K_h(X_1 - z) \eta(\mathbf{X}_2)^2 (f(X_1) - f(z)) \right)^2 \right] \\ &\leq \sup_{z \in \mathcal{X}} \|\eta\|_\infty^2 h^{-2} \int K^2(h^{-1}(x_1 - z)) (f(x_1) - f(z))^2 p_{X_1}(x) dx \\ &= \sup_{z \in \mathcal{X}} \|\eta\|_\infty^2 h^{-1} \int K^2(u) (f(z + uh) - f(z))^2 p_{X_1}(z + uh) du \\ &= \sup_{z \in \mathcal{X}} \|\eta\|_\infty^2 h^{-1} \int K^2(u) (f'(z)uh + o(uh))^2 (p_{X_1}(z) + p'_{X_1}(z)uh + o(uh)) du \\ &= \sup_{z \in \mathcal{X}} \left[(f'(z))^2 p_{X_1}(z) \right] \cdot \int u^2 K^2(u) du \cdot \|\eta\|_\infty^2 h + o(h) = \mathcal{O}(h). \end{aligned}$$

Next, we study the uniform upper bound of functions in \mathcal{S}_h . Depending on whether $h^{-1}(x_1 - z)$ lies in the support, we have two cases

- if $h^{-1}(x_1 - z) \notin \mathcal{X}$, then $K_h(x_1 - z) \eta^2(\mathbf{x}_2) (f(x_1) - f(z)) = 0$;
- if $h^{-1}(x_1 - z) \in \mathcal{X}$, then by mean value theorem,

$$|K_h(x_1 - z) \eta^2(\mathbf{x}_2) (f(x_1) - f(z))| \leq h^{-1} \|K\|_\infty \|f'\|_\infty h \cdot \text{Diam}(\mathcal{X}) = \|K\|_\infty \|f'\|_\infty \text{Diam}(\mathcal{X}).$$

Therefore $\sup_{s_z \in \mathcal{S}_h} \|s_z\|_\infty \leq U := \|K\|_\infty \|f'\|_\infty \text{Diam}(\mathcal{X}) = \mathcal{O}(1)$. Combining the results above, by Lemma C.1 we have

$$\mathbb{E} \left[\sup_{z \in \mathcal{X}} U_n(z) \right] \leq C n^{-1/2} h^{1/2} \sqrt{\log \frac{h^{-9/5}}{h^{1/2}}} = \mathcal{O}(\sqrt{n^{-1} h \log h^{-1}}). \quad (\text{A.7})$$

Similar to σ_P^2 , we bound the expectation as follows

$$\begin{aligned} &\mathbb{E}[K_h(X_1 - z) \eta(\mathbf{X}_2)^2 (f(X_1) - f(z))] \\ &= h^{-1} \int K(h^{-1}(x - z)) (f(x) - f(z)) p_{X_1}(x) \cdot \mathbb{E}[\eta(\mathbf{X}_2)^2 | X_1 = x] dx \\ &= \int K(u) [f'(z)uh + f''(z)(uh)^2/2 + o(u^2 h^2)] [p_{X_1}(z) + p'_{X_1}(z)uh + o(uh)] \\ &\quad \cdot \left(\mathbb{E}[\eta(\mathbf{X}_2)^2 | X_1 = z] + uh \frac{d}{dz} \mathbb{E}[\eta(\mathbf{X}_2)^2 | X_1 = z] + o(uh) \right) du = \mathcal{O}(h^2), \end{aligned} \quad (\text{A.8})$$

where in the last step we used the fact that $K(\cdot)$ is an even function and $f \in \mathcal{H}(2, L)$ on a bounded set \mathcal{X} . Using Lemma C.3, we have

$$\mathbb{P} \left(\sup_{z \in \mathcal{X}} U_n(z) - \mathbb{E} \left[\sup_{z \in \mathcal{X}} U_n(z) \right] \geq t \sqrt{2 \left(\sigma_P^2 + 2U \mathbb{E} \left[\sup_{z \in \mathcal{X}} U_n(z) \right] \right)} + \frac{2Ut^2}{3} \right) \leq \exp(-nt^2), \quad (\text{A.9})$$

where $\sup_{z \in \mathcal{S}_h} \|s_z\|_\infty \leq U := \|K\|_\infty \|f'\|_\infty \text{Diam}(\mathcal{X}) = \mathcal{O}(1)$.

Combining (A.7), (A.8), and (A.9) with $t = \sqrt{\log n/n}$, we have with probability at least $1 - 1/n$,

$$\begin{aligned} \sup_{z \in \mathcal{X}} n^{-1} |\mathbf{h}^\top \mathbf{W}_z \boldsymbol{\xi}_z| &= \sup_{z \in \mathcal{X}} \frac{1}{n} \sum_{i=1}^n K_h(X_{i1} - z) \eta(\mathbf{X}_{i2})^2 (f(X_{i1}) - f(z)) \\ &= \sup_{z \in \mathcal{X}} U_n(z) + \sup_{z \in \mathcal{X}} \mathbb{E}[K_h(X_1 - z) \eta(\mathbf{X}_2)^2 (f(X_1) - f(z))] \\ &= \mathcal{O}(h^2 + \sqrt{n^{-1} h \log n}), \end{aligned}$$

where in the last step we used the assumption that $(nh)^{-1} \log n = \mathcal{O}(1)$. \square

Proof of Lemma 7.3. We denote the empirical process and function class of interest as

$$\begin{aligned} U_{n,\varepsilon}(z) &= \frac{1}{n} \sum_{i=1}^n \eta(\mathbf{X}_{i2}) K_h(X_{i1} - z) \varepsilon_i, \\ \mathcal{S}_{h,\varepsilon} &= \{g_z(x_1, \mathbf{x}_2) = h^{-1} K(h^{-1}(x_1 - z)) \eta(\mathbf{x}_2) | z \in \mathcal{X}\}. \end{aligned}$$

For simplicity we will use the notation $U_n(z)$ and \mathcal{S}_h instead, when it leads to no confusion. Since ε_i 's are subgaussian, we have $\mathbb{P}(\max_i |\varepsilon_i| \geq C\sqrt{\log n}) \leq 1/n$. In the following, we condition on the event $\mathcal{A} = \{\max_i |\varepsilon_i| < C\sqrt{\log n}\}$.

Using an argument similar to that in Lemma 7.1 and 7.2, Lemma C.2 gives us

$$N(\mathcal{S}_h, L^2(\mathbb{P}_n), \epsilon) \leq \left(\frac{2\|K\|_{\text{TV}} \|A\|_\infty \|h\|_\infty}{h\epsilon} \right)^4.$$

We upper bound the variance of the process as

$$\sigma_P^2 := \sup_{z \in \mathcal{X}} \mathbb{E} [K_h^2(X_1 - z) \eta^2(\mathbf{X}_2) \varepsilon_i^2 | \mathcal{A}] \leq \log n \|\eta\|_\infty^2 \sup_{z \in \mathcal{X}} \mathbb{E} [K_h^2(X_1 - z)] = \mathcal{O}(h^{-1} \log n).$$

By Lemma C.1, we have

$$\mathbb{E} \left[\sup_{z \in \mathcal{X}} U_n(z) | \mathcal{A} \right] \leq C \sigma_P n^{-1/2} \sqrt{\log(2\|K\|_{\text{TV}} \|A\|_\infty \|h\|_\infty / (h\sigma_P))} = \mathcal{O} \left(\sqrt{\frac{\log n}{nh}} \right). \quad (\text{A.10})$$

Using Lemma C.3, we have

$$\mathbb{P} \left(\sup_{z \in \mathcal{X}} U_n(z) \geq \mathbb{E} \left[\sup_{z \in \mathcal{X}} U_n(z) | \mathcal{A} \right] + t \sqrt{2 \left(\sigma_P^2 + 2U \mathbb{E} \left[\sup_{z \in \mathcal{X}} U_n(z) | \mathcal{A} \right] \right)} + \frac{2Ut^2}{3} \middle| \mathcal{A} \right) \leq \exp(-nt^2), \quad (\text{A.11})$$

where $\sup_{z \in \mathcal{X}} \|\eta(\mathbf{X}_{12})K_h(X_{11} - z)\varepsilon_1|_{\mathcal{A}}\|_{\infty} \leq U := C\sqrt{\log nh}^{-1}\|K\|_{\infty}\|\eta\|_{\infty} = \mathcal{O}(\sqrt{\log nh}^{-1})$. Combining (A.10) and (A.11), with $t = \sqrt{\log n/n}$, we have

$$\begin{aligned} & \mathbb{P}\left(\sup_{z \in \mathcal{X}} U_n(z) \geq C\sqrt{\log n/(nh)} + C(nh)^{-1/2} \log n\right) \\ & \leq \mathbb{P}\left(\sup_{z \in \mathcal{X}} U_n(z) \geq \mathbb{E}\left[\sup_{z \in \mathcal{X}} U_n(z)|_{\mathcal{A}}\right] + C(nh)^{-1/2} \log n\right) + P(\mathcal{A}^c) \leq 1/n + 1/n = 2/n, \end{aligned}$$

where we used the assumption that $\log n/(nh) = \mathcal{O}(1)$ and the constant C is large enough. Similarly, we have $\mathbb{P}(\sup_z (-U_n(z)) \geq C(nh)^{-1/2} \log n) \leq 2/n$. Therefore, with probability at least $1 - 4/n$,

$$\sup_{z \in \mathcal{X}} n^{-1}|\mathbf{h}^T \mathbf{W}_z \boldsymbol{\varepsilon}| = \sup_{z \in \mathcal{X}} |U_n(z)| \leq C(nh)^{-1/2} \log n,$$

which completes the proof. \square

Proof of Lemma 7.4. Note that

$$\begin{aligned} \frac{1}{n}|\mathbf{h}^T \mathbf{W}_z \boldsymbol{\zeta}| &= \left| \frac{1}{n} \sum_{i=1}^n K_h(X_{i1} - z) \eta(\mathbf{X}_{i2}) \boldsymbol{\psi}(X_{i1}, \mathbf{X}_{i2})^T (\boldsymbol{\gamma}^* - \hat{\boldsymbol{\gamma}}) \right| \\ &\leq \left\| \frac{1}{n} \sum_{i=1}^n K_h(X_{i1} - z) \eta(\mathbf{X}_{i2}) \boldsymbol{\psi}(X_{i1}, \mathbf{X}_{i2})^T \right\|_{2,\infty} \cdot \|\boldsymbol{\gamma}^* - \hat{\boldsymbol{\gamma}}\|_{2,1}. \end{aligned}$$

We can rewrite the first term in the last equation as

$$\begin{aligned} \max_{j \in [d]} \sup_{\mathbf{v} \in \mathbb{B}^m} n^{-1} \sum_{i=1}^n K_h(X_{i1} - z) \eta(\mathbf{X}_{i2}) \mathbf{v}^T \boldsymbol{\psi}_j(X_{i1}, \mathbf{X}_{i2}) &\leq n^{-1} \sum_{i=1}^n K_h(X_{i1} - z) B^2 \sqrt{m} \\ &\leq C\sqrt{m}, \end{aligned}$$

where in the last step we use Lemma A.1 and $\|\eta\|_{\infty}, \|\boldsymbol{\psi}_{jk}\|_{\infty} \leq B$. Thus, with probability $1 - c/n$,

$$\frac{1}{n}|\mathbf{h}^T \mathbf{W}_z \boldsymbol{\zeta}| \leq C\sqrt{m}r_e,$$

which completes the proof. \square

B Auxiliary Lemmas for Bootstrap Confidence Bands

We give proofs of technical lemmas stated in Section 4.3. We start by stating with two technical lemmas.

Lemma B.1. Let $P_1(\mathbf{w}) = \mathbb{E}K_h(X_{i1} - z)(\eta(\mathbf{X}_{i2}) - \mathbf{w}^T \boldsymbol{\psi}(\mathbf{X}_i))\boldsymbol{\psi}(\mathbf{X}_i)^T$. If $(nh)^{-1} \log n = \mathcal{O}(1)$, for any fixed \mathbf{w} , with $\|\mathbf{w}\|_1 \leq L$, there exist a constant C such that

$$\left\| \frac{1}{n} \sum_{i=1}^n K_h(X_{i1} - z)(\eta(\mathbf{X}_{i2}) - \mathbf{w}^T \boldsymbol{\psi}(\mathbf{X}_i))\boldsymbol{\psi}(\mathbf{X}_i)^T - P_1(\mathbf{w}) \right\|_{2,\infty} \leq C \cdot \frac{m + \sqrt{m \log(dnh^{-1})}}{\sqrt{nh}},$$

with probability at least $1 - 1/n$.

Proof. Note that

$$\begin{aligned} & \left\| \frac{1}{n} \sum_{i=1}^n K_h(X_{i1} - z) (\eta(\mathbf{X}_{i2}) - \mathbf{w}^T \boldsymbol{\psi}(\mathbf{X}_i)) \boldsymbol{\psi}(\mathbf{X}_i)^T - P_1(\mathbf{w}) \right\|_{2,\infty} \\ &= \max_{j \in [d]} \sup_{\mathbf{v} \in \mathbb{B}^m} n^{-1} \sum_{i=1}^n K_h(X_{i1} - z) (\eta(\mathbf{X}_{i2}) - \mathbf{w}^T \boldsymbol{\psi}(\mathbf{X}_i)) \mathbf{v}^T \boldsymbol{\psi}_j(\mathbf{X}_{i1}, \mathbf{X}_{i2}) - \mathbf{v}^T P_1(\mathbf{w}). \end{aligned}$$

Let $N_{\mathbf{v}}$ be a $1/2$ -covering of the unit ball $\mathbb{B}^m = \{\mathbf{v} \in \mathbb{R}^m \mid \|\mathbf{v}\|_2 \leq 1\}$ with cardinality $|N_{\mathbf{v}}| \leq 6^m$. For any $\mathbf{v} \in \mathbb{B}^m$, there exists $\pi(\mathbf{v}) \in N_{\mathbf{v}}$ such that $\|\mathbf{v} - \pi(\mathbf{v})\|_2 \leq 1/2$. Therefore,

$$\begin{aligned} & \max_{j \in [d]} \sup_{\mathbf{v} \in \mathbb{B}^m} \sum_{i=1}^n K_h(X_{i1} - z) (\eta(\mathbf{X}_{i2}) - \mathbf{w}^{*T} \boldsymbol{\psi}(\mathbf{X}_i)) \mathbf{v}^T \boldsymbol{\psi}_j(\mathbf{X}_i) - \mathbf{v}^T P_1(\mathbf{w}) \\ & \leq 2 \max_{j \in [d]} \sup_{\mathbf{v} \in N_{\mathbf{v}}} \sum_{i=1}^n K_h(X_{i1} - z) (\eta(\mathbf{X}_{i2}) - \mathbf{w}^{*T} \boldsymbol{\psi}(\mathbf{X}_i)) \mathbf{v}^T \boldsymbol{\psi}_j(\mathbf{X}_i) - \mathbf{v}^T P_1(\mathbf{w}). \end{aligned}$$

We consider the following process and function class of interest:

$$\begin{aligned} U_n(j, \mathbf{v}, z) &= \frac{1}{n} \sum_{i=1}^n K_h(X_{i1} - z) (\eta(\mathbf{X}_{i2}) - \mathbf{w}^T \boldsymbol{\psi}(\mathbf{X}_i)) \mathbf{v}^T \boldsymbol{\psi}_j(\mathbf{X}_i) - \mathbf{v}^T P(\mathbf{w}), \\ S_h &= \{S_z(x_1, x_2) = h^{-1} K(h^{-1}(x_1 - z)) (\eta(x_2) - \mathbf{w}^T \boldsymbol{\psi}(x)) \mathbf{v}^T \boldsymbol{\psi}_j(x) \mid z \in \mathcal{X}, j \in [d], \mathbf{v} \in N_{\mathbf{v}}\}. \end{aligned}$$

Note that S_h is a subset of $\mathcal{F}_{h,1} \times \mathcal{F}_{h,2}$, where

$$\mathcal{F}_{\eta,1} := \{h^{-1} K(h^{-1}(x_1 - z)) (\eta(x_2) - \mathbf{w}^T \boldsymbol{\psi}(x)) \mid z \in \mathcal{X}\}$$

and

$$\mathcal{F}_{h,2} := \{\mathbf{v}^T \boldsymbol{\psi}_j(x_1, x_2) \mid j \in [d], \mathbf{v} \in N_{\mathbf{v}}\}.$$

Since $(\eta(x_2) - \mathbf{w}^T \boldsymbol{\psi}(x)) \leq (L+1)B$, Lemma C.2 gives us

$$N(\mathcal{F}_{h,1}, L^2(\mathbb{P}_n), \epsilon) \leq \left(\frac{2\|K\|_{\text{TV}} A(L+1)B}{h\epsilon} \right)^4.$$

By Lemma C.8 and Lemma C.7, we then obtain

$$N(\mathcal{S}_h, L^2(\mathbb{P}_n), \epsilon) \leq 6^m d \cdot \left(\frac{2\|K\|_{\text{TV}} A(L+1)B}{h} \right)^4 \left(\frac{\sqrt{m} \vee h^{-1}}{\epsilon/2} \right)^4.$$

We upper bound the variance of the process as

$$\begin{aligned} \sigma_p^2 &:= \sup_{j, \mathbf{v}, z} \mathbb{E}[K_h^2(X_{i1} - z) (\eta(\mathbf{X}_{i2}) - \mathbf{w}^T \boldsymbol{\psi}(\mathbf{X}_i))^2 (\mathbf{v}^T \boldsymbol{\psi}_j(\mathbf{X}_i))^2] \\ &\leq B^4 (L+1)^2 m \sup_{z \in \mathcal{X}} \mathbb{E}[K_h^2(X_1 - z)] = \mathcal{O}(mh^{-1}). \end{aligned}$$

Then, by Lemma C.1, we have

$$\begin{aligned}\mathbb{E} \left[\sup_{j, \mathbf{v}, z} U_n(j, \mathbf{v}, z) \right] &\leq C \sigma_P n^{-1/2} \sqrt{\log(6^m d \cdot \|K\|_{\text{TV}} A(\sqrt{m} \vee h^{-1}) / (h \sigma_P))} \\ &= \mathcal{O} \left(\frac{m + \sqrt{m \log(dh^{-1})}}{\sqrt{nh}} \right). \quad (\text{B.1})\end{aligned}$$

By Lemma C.3, we further have

$$\begin{aligned}\mathbb{P} \left(\sup_{j, \mathbf{v}, z} U_n(j, \mathbf{v}, z) \geq \mathbb{E} \left[\sup_{j, \mathbf{v}, z} U_n(j, \mathbf{v}, z) \right] + t \sqrt{2 \left(\sigma_P^2 + 2U \mathbb{E} \left[\sup_{j, \mathbf{v}, z} U_n(j, \mathbf{v}, z) \right] \right)} + \frac{2Ut^2}{3} \right) \\ \leq \exp(-nt^2), \quad (\text{B.2})\end{aligned}$$

where

$$\begin{aligned}\sup_{j, \mathbf{v}, z} \|K_h(x_1 - z) (\eta(\mathbf{x}_2) - \mathbf{w}^T \boldsymbol{\psi}(\mathbf{x})) \mathbf{v}^T \boldsymbol{\psi}_j(x_1, \mathbf{x}_2)\|_\infty \\ \leq U := B^2(L+1) \sqrt{mh}^{-1} \|K\|_\infty = \mathcal{O}(\sqrt{mh}^{-1}).\end{aligned}$$

Combining (B.1) and (B.2), with $t = \sqrt{\log n/n}$, then, with probability at least $1 - 1/n$,

$$\begin{aligned}\sup_{j, \mathbf{v}, z} U_n(j, \mathbf{v}, z) &\leq \mathbb{E} \left[\sup_{j, \mathbf{v}, z} U_n(j, \mathbf{v}, z) \right] + \sqrt{2\sigma_P^2 + 4U \mathbb{E} \left[\sup_{j, \mathbf{v}, z} U_n(j, \mathbf{v}, z) \right]} \cdot \sqrt{\frac{\log n}{n}} + \frac{2U}{3} \cdot \frac{\log n}{n} \\ &= \mathcal{O} \left(\frac{m + \sqrt{m \log(dh^{-1})}}{\sqrt{nh}} + \sqrt{\frac{m \log n}{nh}} \right) = \mathcal{O} \left(\frac{m + \sqrt{m \log(dnh^{-1})}}{\sqrt{nh}} \right),\end{aligned}$$

where we used inequality $\sqrt{xy} \leq 2^{-1}(x+y)$ and the assumption that $(nh)^{-1} \log n = \mathcal{O}(1)$. \square

Applying Lemma B.1 to \mathbf{w}^* and $\tilde{\mathbf{w}}$, we have the following corollaries. We will proof that $\|\tilde{\mathbf{w}}\|_1 \leq L$ later. Note that $P_1(\mathbf{w}^*) = 0$.

Corollary B.2. With probability at least $1 - 1/n$, we have

$$\sup_{z \in \mathcal{X}} \left\| \frac{1}{n} \sum_{i=1}^n K_h(X_{i1} - z) (\eta(\mathbf{X}_{i2}) - \mathbf{w}^{*T} \boldsymbol{\psi}(\mathbf{X}_i)) \boldsymbol{\psi}(\mathbf{X}_i)^T \right\|_{2, \infty} \leq C \left(\frac{m + \sqrt{m \log(dnh^{-1})}}{\sqrt{nh}} \right).$$

Corollary B.3. With probability at least $1 - 2/n$, we have

$$\sup_{z \in \mathcal{X}} \left\| \frac{1}{n} \sum_{i=1}^n t(\mathbf{X}_i, z) \boldsymbol{\psi}(\mathbf{X}_i)^T - \frac{1}{n} \sum_{i=1}^n t(\mathbf{X}'_i, z) \boldsymbol{\psi}(\mathbf{X}'_i)^T \right\|_{2, \infty} \leq C \left(\frac{m + \sqrt{m \log(dnh^{-1})}}{\sqrt{nh}} \right),$$

where $t(\mathbf{X}_i, z) := K_h(X_{i1} - z)(\eta(\mathbf{X}_{i2}) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{X}_i))$ and $\{\mathbf{X}'_i\}$ is another dataset for estimating $\tilde{\mathbf{w}}$.

We also need the following technical result.

Lemma B.4. Let $P_2(\mathbf{w}) = \mathbb{E}K_h(X_{i1} - z)(\eta(\mathbf{X}_{i2}) - \mathbf{w}^\top \boldsymbol{\psi}(\mathbf{X}_i))\eta(X_{i2})$. If $(nh)^{-1} \log n = \mathcal{O}(1)$, for any fixed \mathbf{w} , with $\|\mathbf{w}\|_1 \leq L$, there exist a constant C such that

$$\sup_{z \in \mathcal{X}} \left| \frac{1}{n} K_h(X_{i1} - z)(\eta(\mathbf{X}_{i2}) - \mathbf{w}^\top \boldsymbol{\psi}(\mathbf{X}_i))\eta(X_{i2}) - P_2(\mathbf{w}) \right| \leq C \cdot \sqrt{\frac{\log n}{nh}},$$

with probability at least $1 - 1/n$.

Proof. The process and the function class of interest are:

$$U_n(z) = \frac{1}{n} \sum_{i=1}^n K_h(X_{i1} - z)(\eta(\mathbf{X}_{i2}) - \mathbf{w}^\top \boldsymbol{\psi}(\mathbf{X}_i))\eta(X_{i2}) - P_2(\mathbf{w}),$$

$$\mathcal{S}_h = \{S_z(x_1, \mathbf{x}_2) := K_h(x_1 - z)(\eta(\mathbf{x}_2) - \mathbf{w}^\top \boldsymbol{\psi}(\mathbf{x}))\eta(\mathbf{x}_2) | z \in \mathcal{X}\}.$$

Since $(\eta(\mathbf{X}_{i2}) - \mathbf{w}^\top \boldsymbol{\psi}(\mathbf{X}_i))\eta(X_{i2})$ is bounded, using Lemma C.2, we have

$$N(\mathcal{S}_h, L^2(\mathbb{P}_n), \epsilon) \leq \left(\frac{2\|K\|_{\text{TV}} A(L+1)B^2}{h\epsilon} \right)^4,$$

where A is parameter from Lemma C.2 and B is the upper bound for $\|\eta\|_\infty$ and $\|\boldsymbol{\psi}_{jk}\|_\infty$. We can upper bound the variance of the process as:

$$\sigma_p^2 := \sup_z \mathbb{E} K_h(X_{i1} - z)^2 (\eta(\mathbf{X}_{i2}) - \mathbf{w}^\top \boldsymbol{\psi}(\mathbf{X}_i))^2 \eta(X_{i2})^2 = \mathcal{O}(h^{-1}).$$

Therefore, from Lemma C.1,

$$\mathbb{E} \left[\sup_{z \in \mathcal{X}} U_n(z) \right] \leq C n^{-1/2} \sigma_p \sqrt{\log \frac{2\|K\|_{\text{TV}}(L+1)B^2}{h\sigma_p}} = C \frac{\sqrt{\log h^{-1}}}{\sqrt{nh}}. \quad (\text{B.3})$$

Since $\sup_z |K_h(x_1 - z)(\eta(\mathbf{x}_2) - \mathbf{w}^\top \boldsymbol{\psi}(\mathbf{x}))\eta(\mathbf{x}_2)| \leq U := Ch^{-1}$, by Lemma C.3, we have

$$\mathbb{P} \left(\sup_z U_n(z) \geq \mathbb{E} \left[\sup_z U_n(z) \right] + t \sqrt{2 \left(\sigma_p^2 + 2U \mathbb{E} \left[\sup_z U_n(z) \right] \right)} + \frac{2Ut^2}{3} \right) \leq \exp(-nt^2). \quad (\text{B.4})$$

Combining (B.3) and (B.4), with $t = \sqrt{\frac{\log n}{n}}$, then, with probability at least $1 - \frac{1}{n}$, we have

$$\begin{aligned} \sup_z U_n(z) &\leq \mathbb{E} \left[\sup_z U_n(z) \right] + \sqrt{2\sigma_p^2 + 4U \mathbb{E} \left[\sup_z U_n(z) \right]} \cdot \sqrt{\frac{\log n}{n}} + \frac{2U}{3} \cdot \frac{\log n}{n} \\ &\leq C \sqrt{\frac{\log n}{nh}}, \end{aligned}$$

where we used the inequality $\sqrt{xy} \leq 2^{-1}(x + y)$ and the assumption that $(nh)^{-1} \log n = \mathcal{O}(1)$. \square

Similarly, we can have two other corollaries:

Corollary B.5. With probability at least $1 - \frac{1}{n}$, we have

$$\min_{z \in \mathcal{X}} \frac{1}{n} \sum_{i=1}^n K_h(X_{i1} - z) (\eta(\mathbf{X}_{i2}) - \mathbf{w}^{*T} \boldsymbol{\psi}(\mathbf{X}_i)) \eta(X_{i2}) \geq b - C \sqrt{\frac{\log n}{nh}}.$$

Note that in Assumption 4.1, $\mathbb{E} K_h(X_1 - z) \eta(\mathbf{X}_2) (\eta(\mathbf{X}_2) - \mathbf{w}^{*T} \boldsymbol{\psi}(X)) \geq b$.

Corollary B.6. With probability at least $1 - \frac{2}{n}$, we have

$$\sup_{z \in \mathcal{X}} \left| \frac{1}{n} \sum_{i=1}^n t(\mathbf{X}_i, z) \eta(\mathbf{X}_{i2}) - \frac{1}{n} \sum_{i=1}^n t(\mathbf{X}'_i, z) \eta(\mathbf{X}'_{i2}) \right| \leq C \sqrt{\frac{\log n}{nh}},$$

where $t(\mathbf{X}_i, z) := K_h(X_{i1} - z) (\eta(\mathbf{X}_{i2}) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{X}_i))$ and $\{\mathbf{X}'_i\}$ is another dataset for estimating $\tilde{\mathbf{w}}$.

Now we can proof the Lemma 7.5.

Proof of Lemma 7.5. Combining Corollary B.2 and Corollary B.5, with probability $1 - \frac{2}{n}$,

$$\sup_{z \in \mathcal{X}} \left\| \frac{1}{n} \sum_{i=1}^n K_h(X_{i1} - z) (\eta(\mathbf{X}_{i2}) - \mathbf{w}^{*T} \boldsymbol{\psi}(\mathbf{X}_i)) \boldsymbol{\psi}(\mathbf{X}_i)^T \right\|_{2,\infty} \leq C \left(\frac{m + \sqrt{m \log(dnh^{-1})}}{\sqrt{nh}} \right).$$

Setting $\lambda = C \left(\frac{m + \sqrt{m \log(dnh^{-1})}}{\sqrt{nh}} \right)$, \mathbf{w}^* is a feasible solution for the optimization problem (7.7).

Since the loss function is $\|\mathbf{w}\|_1$, we can get that with probability at least $1 - \frac{2}{n}$, the solution $\tilde{\mathbf{w}}$ to problem (7.7) exists, and $\|\tilde{\mathbf{w}}\|_1 \leq \|\mathbf{w}^*\|_1 \leq L$. \square

Proof of Lemma 7.6. The loss function

$$L_z(\alpha, \gamma) = \frac{1}{2n} \sum_{i=1}^n K_n(X_{i1} - z) (Y_i - \alpha \cdot \eta(\mathbf{X}_{i2}) - \gamma^T \boldsymbol{\psi}(\mathbf{X}_i))^2$$

is a polynomial of degree 2. From the optimization problem (7.8),

$$\begin{aligned} & \nabla_\alpha L_z(\tilde{\alpha}, \hat{\gamma}) - \tilde{\mathbf{w}} \nabla_\gamma L_z(\tilde{\alpha}, \hat{\gamma}) \\ &= \nabla_\alpha L_z(\tilde{\alpha}, \gamma^*) + \nabla_{\alpha\gamma}^2 L_z(\tilde{\alpha}, \gamma^*) (\hat{\gamma} - \gamma^*) - \tilde{\mathbf{w}} \nabla_\gamma L_z(\tilde{\alpha}, \gamma^*) - \tilde{\mathbf{w}} \nabla_{\gamma\gamma}^2 L_z(\tilde{\alpha}, \gamma^*) (\hat{\gamma} - \gamma^*) \\ &= \nabla_\alpha L_z(\tilde{\alpha}, \gamma^*) - \tilde{\mathbf{w}} \nabla_\gamma L_z(\tilde{\alpha}, \gamma^*) + \left(\nabla_{\alpha\gamma}^2 L_z(\tilde{\alpha}, \gamma^*) - \tilde{\mathbf{w}} \nabla_{\gamma\gamma}^2 L_z(\tilde{\alpha}, \gamma^*) \right) (\hat{\gamma} - \gamma^*) = 0. \end{aligned} \quad (\text{B.5})$$

Furthermore,

$$\begin{aligned} \nabla_\alpha L_z(\tilde{\alpha}, \gamma^*) &= \nabla_\alpha L_z(\alpha^*, \gamma^*) + \nabla_{\alpha\alpha}^2 L_z(\alpha^*, \gamma^*) (\tilde{\alpha} - \alpha^*), \\ \tilde{\mathbf{w}} \nabla_\gamma L_z(\tilde{\alpha}, \gamma^*) &= \tilde{\mathbf{w}} \nabla_\gamma L_z(\alpha^*, \gamma^*) + \tilde{\mathbf{w}} \nabla_{\gamma\alpha}^2 L_z(\alpha^*, \gamma^*) (\tilde{\alpha} - \alpha^*). \end{aligned} \quad (\text{B.6})$$

Combining (B.5) and (B.6), we have

$$(\nabla_{\alpha\alpha}^2 L_z - \tilde{\mathbf{w}} \nabla_{\alpha\gamma}^2 L_z) (\tilde{\alpha} - \alpha^*) = (\nabla_{\alpha\gamma}^2 L_z - \tilde{\mathbf{w}} \nabla_{\gamma\gamma}^2 L_z) (\gamma^* - \hat{\gamma}) + \tilde{\mathbf{w}} \nabla_\gamma L_z(\alpha^*, \gamma^*) - \nabla_\alpha L_z(\alpha^*, \gamma^*).$$

Therefore,

$$\tilde{\alpha} - \alpha^* = \frac{(\nabla_{\alpha\gamma}^2 L_z - \tilde{\mathbf{w}}^T \nabla_{\gamma\gamma}^2 L_z)(\gamma^* - \hat{\gamma}) + \tilde{\mathbf{w}} \nabla_r L_z(\alpha^*, \gamma^*) - \nabla_{\alpha} L_z}{\nabla_{\alpha\alpha}^2 L_z - \tilde{\mathbf{w}} \nabla_{\alpha\gamma}^2 L_z}.$$

Note that since $\nabla_{\alpha\alpha}^2 L_z - \tilde{\mathbf{w}} \nabla_{\alpha\gamma}^2 L_z$ is the schur component of the nonsingular Hessian matrix Σ_z^* , it must be positive. \square

Then we give an Lemma to bound the denominator of $\tilde{\alpha} - \alpha^*$.

Lemma B.7. Under Assumption 4.1 and conditions in Theorem 4.2, with probability $1 - \frac{c}{n}$,

$$\min_z \nabla_{\alpha\alpha}^2 L_z - \tilde{\mathbf{w}} \nabla_{\alpha\gamma}^2 L_z \geq b/2.$$

Proof. In terms of Lemma 7.5, now we already have

$$\|\tilde{\mathbf{w}}\|_1 \leq L, \tag{B.7}$$

$$\sup_{z \in \mathcal{X}} \left\| \frac{1}{n} \sum_{i=1}^n t(\mathbf{X}'_i, z) \boldsymbol{\psi}(\mathbf{X}'_i)^T \right\|_{2,\infty} \leq C \left(\frac{m + \sqrt{m \log(dnh^{-1})}}{\sqrt{nh}} \right). \tag{B.8}$$

We are going to prove

$$\min_z \nabla_{\alpha\alpha}^2 L_z - \tilde{\mathbf{w}} \nabla_{\alpha\gamma}^2 L_z = \min_z \frac{1}{n} \sum_{i=1}^n K_h(X_{i1} - z) \eta(\mathbf{X}_{i2}) (\eta(\mathbf{X}_{i2}) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{X}_i)) \geq b/2.$$

Apply Lemma B.1 to (B.8), with probability $1 - c/n$,

$$\left\| \mathbb{E} K_h(X_1 - z) (\eta(\mathbf{X}_2) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{X})) \boldsymbol{\psi}(\mathbf{X}) \right\|_{2,\infty} \leq C \left(\frac{m + \sqrt{m \log(dnh^{-1})}}{\sqrt{nh}} \right).$$

Since

$$\begin{aligned} & \tilde{\mathbf{w}}^T \cdot \mathbb{E} K_h(X_1 - z) (\eta(\mathbf{X}_2) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{X})) \boldsymbol{\psi}(\mathbf{X}) + \mathbb{E} K_h(X_1 - z) \eta(\mathbf{X}_2) (\eta(\mathbf{X}_2) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{X})) \\ &= [-\tilde{\mathbf{w}}^T \mathbf{1}] \Sigma_z \begin{bmatrix} -\tilde{\mathbf{w}} \\ 1 \end{bmatrix} \geq b \sqrt{\|\tilde{\mathbf{w}}\|_2 + 1} \geq b. \end{aligned}$$

And

$$\begin{aligned} & \tilde{\mathbf{w}}^T \cdot \mathbb{E} K_h(X_1 - z) (\eta(\mathbf{X}_2) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{X})) \boldsymbol{\psi}(\mathbf{X}) \\ & \leq \|\tilde{\mathbf{w}}^T\|_{2,1} \cdot \left\| \mathbb{E} K_h(X_1 - z) (\eta(\mathbf{X}_2) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{X})) \boldsymbol{\psi}(\mathbf{X}) \right\|_{2,\infty} \\ & \leq \|\tilde{\mathbf{w}}\|_1 \cdot \left\| \mathbb{E} K_h(X_1 - z) (\eta(\mathbf{X}_2) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{X})) \boldsymbol{\psi}(\mathbf{X}) \right\|_{2,\infty} \\ & \leq C \left(\frac{m + \sqrt{m \log(dnh^{-1})}}{\sqrt{nh}} \right). \end{aligned}$$

Therefore, we have with probability $1 - c/n$,

$$\min_z \mathbb{E} K_h(X_1 - z) \eta(\mathbf{X}_2) (\eta(\mathbf{X}_2) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{X})) \geq b - C \left(\frac{m + \sqrt{m \log(dnh^{-1})}}{\sqrt{nh}} \right).$$

Apply Lemma B.4, under conditions in Theorem 4.2, we have with probability $1 - c/n$,

$$\min_z \frac{1}{n} \sum_{i=1}^n K_h(X_{i1} - z) \eta(\mathbf{X}_{i2}) (\eta(\mathbf{X}_{i2}) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{X}_i)) \geq b - C \left(\frac{m + \sqrt{m \log(dnh^{-1})}}{\sqrt{nh}} \right) - C \sqrt{\frac{\log n}{nh}} \geq b/2,$$

for sufficient large n .

□

Proof of Lemma 7.7. Recall that

$$I_1 = \frac{(\nabla_{\alpha\gamma}^2 L_z - \tilde{\mathbf{w}}^T \nabla_{\gamma\gamma}^2 L_z)(\gamma^* - \hat{\gamma})}{\nabla_{\alpha\alpha}^2 L_z - \tilde{\mathbf{w}} \nabla_{\alpha\gamma}^2 L_z}.$$

Under Assumption 4.1 $P(\|\hat{\gamma} - \gamma^*\|_{2,1} > Cr_e) \leq 1/n$, so we focus on upper bounding the numerator and lower bounding the denominator:

$$\begin{aligned} \|\nabla_{\alpha\gamma}^2 L_z - \tilde{\mathbf{w}}^T \nabla_{\gamma\gamma}^2 L_z\|_{2,\infty} &= \left\| \frac{1}{n} \sum_{i=1}^n t(\mathbf{X}_i, z) \boldsymbol{\psi}(\mathbf{X}_i)^T \right\|_{2,\infty}, \\ \nabla_{\alpha\alpha}^2 L_z - \tilde{\mathbf{w}} \nabla_{\alpha\gamma}^2 L_z &= \frac{1}{n} \sum_{i=1}^n t(\mathbf{X}_i, z) \eta(\mathbf{X}_{i2}). \end{aligned}$$

From the optimization problem (7.7) and Lemma 7.5, we have

$$\begin{aligned} \|\tilde{\mathbf{w}}\|_1 &\leq L, \\ \left\| \frac{1}{n} \sum_{i=1}^n t(\mathbf{X}_i', z) \boldsymbol{\psi}(\mathbf{X}_i')^T \right\|_{2,\infty} &\leq C \left(\frac{m + \sqrt{m \log(dnh^{-1})}}{\sqrt{nh}} \right). \end{aligned}$$

Combined with Corollary B.3 and Lemma B.7, we get

$$\begin{aligned} \left\| \frac{1}{n} \sum_{i=1}^n t(\mathbf{X}_i, z) \boldsymbol{\psi}(\mathbf{X}_i)^T \right\|_{2,\infty} &\leq C \left(\frac{m + \sqrt{m \log(dnh^{-1})}}{\sqrt{nh}} \right), \\ \frac{1}{n} \sum_{i=1}^n t(\mathbf{X}_i, z) \eta(\mathbf{X}_{i2}) &\geq b/2. \end{aligned}$$

Therefore, with probability $1 - \frac{7}{n}$, we have

$$I_1 = \frac{(\nabla_{\alpha\gamma}^2 L_z - \tilde{\mathbf{w}}^T \nabla_{\gamma\gamma}^2 L_z)(\gamma^* - \hat{\gamma})}{\nabla_{\alpha\alpha}^2 L_z - \tilde{\mathbf{w}} \nabla_{\alpha\gamma}^2 L_z} \leq C \cdot r_e \cdot \frac{m + \sqrt{m \log(dnh^{-1})}}{\sqrt{nh}}.$$

□

Proof of Lemma 7.8. We have the following representation

$$\begin{aligned} I_{21} &= \frac{\frac{1}{n} \sum_{i=1}^n t(\mathbf{X}_i, z) [(f^*(X_{i1}) - f^*(z)) \eta(\mathbf{X}_{i2}) + g(\mathbf{X}_i) - \gamma^{*T} \boldsymbol{\psi}(\mathbf{X}_i)]}{p(z)} \\ &= \frac{\frac{1}{n} \sum_{i=1}^n t(\mathbf{X}_i, z) (f^*(X_{i1}) - f^*(z)) \eta(\mathbf{X}_{i2})}{p(z)} + \frac{\frac{1}{n} \sum_{i=1}^n t(\mathbf{X}_i, z) [g(\mathbf{X}_i) - \gamma^{*T} \boldsymbol{\psi}(\mathbf{X}_i)]}{p(z)}. \end{aligned}$$

From Assumption 4.1, $|g(\mathbf{X}_i) - \gamma^{*T} \boldsymbol{\psi}(\mathbf{X}_i)| \leq Cr_a$. Since $t(X_i, z) = K_h(X_{i1} - z)(\eta(\mathbf{X}_{i2}) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{X}_i))$, it follows from Lemma A.1 that

$$\frac{1}{n} \sum_{i=1}^n t(\mathbf{X}_i, z) [g(\mathbf{X}_i) - \gamma^{*T} \boldsymbol{\psi}(\mathbf{X}_i)] \leq CBr_a \leq Cr_a.$$

From Corollary B.4, we have that $p(z) > b/4$. Therefore, with probability $1 - \frac{c}{n}$,

$$\frac{\frac{1}{n} \sum_{i=1}^n t(\mathbf{X}_i, z) [g(\mathbf{X}_i) - \gamma^{*T} \boldsymbol{\psi}(\mathbf{X}_i)]}{p(z)} \leq 4Cr_a/b = Cr_a. \quad (\text{B.9})$$

Next, we deal with the first term. We consider the process and function class of interest:

$$U_n(z) = \frac{1}{n} \sum_{i=1}^n t(\mathbf{X}_i, z) (f^*(X_{i1}) - f^*(z)) \eta(\mathbf{X}_{i2}) - \mathbb{E} t(\mathbf{X}_i, z) (f^*(X_{i1}) - f^*(z)) \eta(\mathbf{X}_{i2}),$$

$$\mathcal{S}_h = \{S_z(x_1, \mathbf{x}_2) = K_h(x_1 - z) (\eta(\mathbf{x}_2) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{x})) \eta(\mathbf{x}_2) (f^*(x_1) - f^*(z)) \mid z \in \mathcal{X}\}.$$

Note that \mathcal{S}_h is a subset of $\mathcal{F}_{h,1} \times \mathcal{F}_{h,2}$, where $\mathcal{F}_{h,1} := \{h^{-1}K(h^{-1}(x_1 - z))h^2(\mathbf{x}_2) \mid z \in \mathcal{X}\}$ and $\mathcal{F}_{h,2} := \{f(x_1) - f(z) \mid z \in \mathcal{X}\}$. Similar to the argument used to obtain (A.3), we prove

$$N(\mathcal{F}_{h,1}, L^2(\mathbb{P}_n), \epsilon) \leq \left(\frac{2\|K\|_{\text{TV}} AB^2(L+1)}{h\epsilon} \right)^4.$$

On the other hand, $f \in \mathcal{H}(2, L)$ on a bounded set \mathcal{X} and is hence Lipschitz. Therefore, we have $|(f(x_1) - f(z_1)) - (f(x_1) - f(z_2))| \leq \|f'\|_\infty |z_1 - z_2|$, and

$$N(\mathcal{F}_{h,2}, L^2(\mathbb{P}_n), \epsilon) \leq N(\mathcal{X}, d(\cdot, \cdot), \|f'\|_\infty^{-1} \epsilon) \leq \frac{\|f'\|_\infty \text{Diam}(\mathcal{X})}{\epsilon}.$$

By Lemma C.7 and Lemma C.8, we obtain

$$N(\mathcal{S}_h, L^2(\mathbb{P}_n), \epsilon) \leq \|f'\|_\infty \text{Diam}(\mathcal{X}) \left(\frac{2\|K\|_{\text{TV}} A(L+1)B^2}{h} \right)^4 \left(\frac{\|\mathcal{F}_{h,1}\|_\infty \vee \|\mathcal{F}_{h,2}\|_\infty}{\epsilon/2} \right)^5 = O(h^{-9} \epsilon^{-5}).$$

We bound the variance using the Taylor expansion as follows

$$\begin{aligned} \sigma_{\mathbb{P}}^2 &:= \sup_{z \in \mathcal{X}} \mathbb{E} \left[\left(K_h(X_1 - z) \eta(\mathbf{X}_2) (\eta(\mathbf{X}_2) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{X})) (f(X_1) - f(z)) \right)^2 \right] \\ &\leq \sup_{z \in \mathcal{X}} (L+1)B^2 h^{-1} \int K^2(h^{-1}(x_1 - z)) (f(x_1) - f(z))^2 p_{X_1}(x) dx \\ &= \sup_{z \in \mathcal{X}} (L+1)B^2 h^{-1} \int K^2(u) (f(z + uh) - f(z))^2 p_{X_1}(z + uh) du \\ &= \sup_{z \in \mathcal{X}} (L+1)B^2 h^{-1} \int K^2(u) (f'(z)uh + o(uh))^2 (p_{X_1}(z) + p'_{X_1}(z)uh + o(uh)) du \\ &= \sup_{z \in \mathcal{X}} \left[(f'(z))^2 p_{X_1}(z) \right] \cdot \int u^2 K^2(u) du \cdot (L+1)B^2 h + o(h) = \mathcal{O}(h). \end{aligned}$$

Next, we study the uniform upper bound of functions in \mathcal{S}_h . Depending on whether $h^{-1}(x_1 - z)$ lies in the support, we have two cases

- if $h^{-1}(x_1 - z) \notin \mathcal{X}$, then $K_h(x_1 - z)h^2(\mathbf{x}_2)(f(x_1) - f(z)) = 0$;
- if $h^{-1}(x_1 - z) \in \mathcal{X}$, then by mean value theorem,

$$\begin{aligned} & |K_h(x_1 - z)\eta(\mathbf{x}_2)(\eta(\mathbf{x}_2) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{x})) (f(x_1) - f(z))| \\ & \leq h^{-1} \|K\|_\infty \|f'\|_\infty B^2(L+1) \cdot \text{Diam}(\mathcal{X}) = \|K\|_\infty \|f'\|_\infty \text{Diam}(\mathcal{X}). \end{aligned}$$

Therefore $\sup_{s_z \in \mathcal{S}_h} \|s_z\|_\infty \leq U := \|K\|_\infty \|f'\|_\infty \text{Diam}(\mathcal{X}) = \mathcal{O}(1)$. By Lemma C.1, we then have

$$\mathbb{E} \left[\sup_{z \in \mathcal{X}} U_n(z) \right] \leq C n^{-1/2} h^{1/2} \sqrt{\log \frac{h^{-9/5}}{h^{1/2}}} = \mathcal{O}(\sqrt{n^{-1} h \log h^{-1}}). \quad (\text{B.10})$$

We bound the expectation as follows

$$\begin{aligned} & \mathbb{E}[K_h(X_1 - z)\eta(\mathbf{X}_2)(\eta(\mathbf{X}_2) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{X}_i))(f(X_1) - f(z))] \\ & = h^{-1} \int K(h^{-1}(x - z))(f(x) - f(z)) p_{X_1}(x) \cdot \mathbb{E}[\phi(\mathbf{X})|X_1 = x] dx \\ & = \int K(u) [f'(z)uh + f''(z)(uh)^2/2 + o(u^2 h^2)] [p_{X_1}(z) + p'_{X_1}(z)uh + o(uh)] \\ & \quad \cdot \left(\mathbb{E}[\phi(\mathbf{X})|X_1 = z] + uh \frac{d}{dz} \mathbb{E}[\phi(\mathbf{X})|X_1 = z] + o(uh) \right) du = \mathcal{O}(h^2), \end{aligned} \quad (\text{B.11})$$

where in the last step we used the fact that $K(\cdot)$ is an even function, $f \in \mathcal{H}(2, L)$ on a bounded set \mathcal{X} , and $\phi(\mathbf{X}) = \eta(\mathbf{X}_2)(\eta(\mathbf{X}_2) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{X}_i))$. Using Lemma C.3, we then have

$$\mathbb{P} \left(\sup_{z \in \mathcal{X}} U_n(z) - \mathbb{E} \left[\sup_{z \in \mathcal{X}} U_n(z) \right] \geq t \sqrt{2(\sigma_P^2 + 2U \mathbb{E} \left[\sup_{z \in \mathcal{X}} U_n(z) \right])} + \frac{2Ut^2}{3} \right) \leq \exp(-nt^2), \quad (\text{B.12})$$

where $\sup_{s_z \in \mathcal{S}_h} \|s_z\|_\infty \leq U := \|K\|_\infty \|f'\|_\infty \text{Diam}(\mathcal{X}) = \mathcal{O}(1)$. Combining (B.10), (B.11), and (B.12), with $t = \sqrt{\log n/n}$, we have with probability at least $1 - 1/n$,

$$\begin{aligned} & \sup_{z \in \mathcal{X}} \frac{1}{n} \sum_{i=1}^n K_h(X_{i1} - z)\eta(\mathbf{X}_2)(\eta(\mathbf{X}_2) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{X}_i))(f(X_{i1}) - f(z)) \\ & = \sup_{z \in \mathcal{X}} U_n(z) + \sup_{z \in \mathcal{X}} \mathbb{E}[K_h(X_1 - z)\eta(\mathbf{X}_2)(\eta(\mathbf{X}_2) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{X}_i))(f(X_1) - f(z))] \\ & = \mathcal{O}(h^2 + \sqrt{n^{-1} h \log n}), \end{aligned}$$

where in the last step we used the assumption that $(nh)^{-1} \log n = \mathcal{O}(1)$. Combined with $p(z) > b/4$, we have that, with probability at least $1 - c/n$,

$$\frac{\frac{1}{n} \sum_{i=1}^n t(\mathbf{X}_i, z)(f^*(X_{i1}) - f^*(z))\eta(\mathbf{X}_{i2})}{p(z)} \leq 4C \frac{h^2 + \sqrt{n^{-1} h \log n}}{b} \leq C(h^2 + \sqrt{n^{-1} h \log n}). \quad (\text{B.13})$$

Combining (B.13) and (B.9), we have

$$I_{21} \leq C(h^2 + \sqrt{n^{-1} h \log n} + r_a),$$

with probability $1 - \frac{c}{n}$, which completes the proof. \square

We need the following two technical lemmas to prove Lemma 7.9.

Lemma B.8. With probability $1 - \frac{c}{n}$,

$$|\hat{\sigma}^2 - \sigma^2| \leq Cr_n,$$

where $\hat{\sigma}^2 := n^{-1} \sum_{i=1}^n \hat{\varepsilon}_i^2 = n^{-1} \sum_{i=1}^n (Y_i - \hat{f}(X_{i1})\eta(\mathbf{X}_{i2}) - \hat{g}(X_{i1}, \mathbf{X}_{i2}))^2$.

Proof. We have the following decomposition

$$\begin{aligned} |\hat{\sigma}^2 - \sigma^2| &= \left| \frac{1}{n} \sum_{i=1}^n \hat{\varepsilon}_i^2 - \frac{1}{n} \sum_{i=1}^n \varepsilon_i^2 + \frac{1}{n} \sum_{i=1}^n \varepsilon_i^2 - \sigma^2 \right| \\ &\leq \underbrace{\frac{1}{n} \sum_{i=1}^n (\hat{\varepsilon}_i - \varepsilon_i)^2}_{\text{(I)}} + \underbrace{\left| \frac{2}{n} \sum_{i=1}^n (\hat{\varepsilon}_i - \varepsilon_i) \varepsilon_i \right|}_{\text{(II)}} + \underbrace{\left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i^2 - \sigma^2 \right|}_{\text{(III)}}. \end{aligned}$$

For all i , applying Theorem 4.1, the following holds with probability $1 - \frac{7}{n}$:

$$\begin{aligned} |\hat{\varepsilon}_i - \varepsilon_i| &= |Y_i - \varepsilon_i - \hat{f}(X_{i1})\eta(\mathbf{X}_{i2}) - \hat{g}(\mathbf{X}_i)| \\ &\leq |(f_*(X_{i1}) - \hat{f}(X_{i1}))\eta(\mathbf{X}_{i2})| + |g(\mathbf{X}_i) - \hat{g}(\mathbf{X}_i)| \\ &\leq |(f_*(X_{i1}) - \hat{f}(X_{i1}))\eta(\mathbf{X}_{i2})| + |g(\mathbf{X}_i) - \gamma^*\psi(\mathbf{X}_i)| + |\hat{\gamma}\psi(\mathbf{X}_i) - \gamma^*\psi(\mathbf{X}_i)| \\ &\leq Cr_n + Cr_a + C\sqrt{mr_e} \\ &\leq Cr_n, \end{aligned}$$

where $r_n := r_a + \sqrt{n^{-1}h \log n} + h^2 + (nh)^{-1/2} \log n + r_e\sqrt{m}$. Applying Bernstein's inequality (Bernstein, 1964), we have

$$\mathbb{P}\left(\left|\frac{1}{n} \sum_{i=1}^n \varepsilon_i^2 - \sigma^2\right| > C_1 \sqrt{\frac{\sigma^2 \log n}{n}}\right) \leq \frac{2}{n} \quad \text{and} \quad \mathbb{P}\left(\left|\frac{1}{n} \sum_{i=1}^n |\varepsilon_i| - \mathbb{E}|\varepsilon|\right| > C_2 \sqrt{\frac{\sigma^2 \log n}{n}}\right) \leq \frac{2}{n}.$$

Suppose n is large enough, so that $\mathbb{E}|\varepsilon| \geq \sqrt{\sigma^2 \log n/n}$. Then

$$\begin{aligned} &\mathbb{P}\left(\Pi > C(C_2 + 1)\mathbb{E}|\varepsilon|r_n\right) \\ &\leq \mathbb{P}\left(\max_{i \in [n]} |\hat{\varepsilon}_i - \varepsilon_i| > Cr_n\right) + \mathbb{P}\left(\frac{1}{n} \sum_{i=1}^n |\varepsilon_i| > \mathbb{E}|\varepsilon| + C_2 \sqrt{\frac{\sigma^2 \log n}{n}}\right) \leq \frac{11}{n}. \end{aligned}$$

Combining all the bounds, the proof is complete. \square

Lemma B.9. With probability $1 - \frac{c}{n}$, we have

$$\sup_z |\hat{q}^2(z) - q^2(z)| \leq C \sqrt{\frac{\log n}{nh^3}}.$$

Proof. Recall that

$$q^2(z) = \mathbb{E} \left[K_h(X_1 - z)^2 (\eta(\mathbf{X}_2) - \tilde{\mathbf{w}}(z)^T \boldsymbol{\psi}(\mathbf{X}))^2 \right], \quad (\text{B.14})$$

$$\hat{q}^2(z) = \frac{1}{n} \sum_{i=1}^n K_h(X_{i1} - z)^2 (\eta(\mathbf{X}_{i2}) - \tilde{\mathbf{w}}(z)^T \boldsymbol{\psi}(\mathbf{X}_i))^2. \quad (\text{B.15})$$

We consider the following process and the function class of interest:

$$\begin{aligned} U_n(z) &= \hat{q}^2(z) - q^2(z), \\ S_\eta(z) &= \{S(x_1, x_2) = K_h(x_1 - z)^2 (\eta(x_2) - \tilde{\mathbf{w}}(z)^T \boldsymbol{\psi}(\mathbf{x}))^2\}. \end{aligned}$$

Due to Lemma C.2, we have

$$N(S_\eta, L^2(\mathbb{P}_n), \epsilon) \leq \left(\frac{2\|K\|_{\text{TV}} A}{h^2 \epsilon} \right)^4.$$

Similarly, we can bound that:

$$\sigma_p^2 = \mathbb{E} \left[K_h(x_1 - z)^4 (\eta(x_2) - \tilde{\mathbf{w}}(z)^T \boldsymbol{\psi}(\mathbf{x}))^4 \right] \leq Ch^{-3} \mathbb{E} K_h(x_1 - z) = O(h^{-3}).$$

Thus, from Lemma C.1,

$$\mathbb{E}[\sup_z U_n(z)] \leq C\sigma_p n^{-1/2} \sqrt{\log \frac{2\|K\|_{\text{TV}} A}{h^2 \sigma_p}} = C\sqrt{\frac{\log(h^{-1})}{nh^3}}. \quad (\text{B.16})$$

Besides, the upper bound if $S_\eta(z)$ is $U = O(h^{-2})$, since $(\eta(x_2) - \tilde{\mathbf{w}}(z)^T \boldsymbol{\psi}(\mathbf{x}))^2$ can be bounded by constant $(L+1)^2 B^2$. By Lemma C.3, we have

$$\mathbb{P} \left(\sup_z U_n(z) \geq \mathbb{E} \left[\sup_z U_n(z) \right] + t \sqrt{2 \left(\sigma_p^2 + 2U \mathbb{E} \left[\sup_z U_n(z) \right] \right)} + \frac{2Ut^2}{3} \right) \leq \exp(-nt^2). \quad (\text{B.17})$$

Combining (B.16) and (B.17), with $t = \sqrt{\log n/n}$, we have, with probability at least $1 - 1/n$,

$$\begin{aligned} \sup_z |U_n(z)| &\leq \mathbb{E} \left[\sup_z U_n(z) \right] + \sqrt{\frac{\log n}{n}} \cdot \sqrt{2 \left(\sigma_p^2 + 2U \mathbb{E} \left[\sup_z U_n(z) \right] \right)} + \frac{2U}{3} \cdot \frac{\log n}{n} \\ &\leq C\sqrt{\frac{\log(h^{-1})}{nh^3}} + C\sqrt{\frac{\log n}{nh^3}} \leq C\sqrt{\frac{\log n}{nh^3}}, \end{aligned}$$

which completes the proof. \square

Proof of Lemma 7.9. We first establish a uniform lower bound on $q^2(z)$. From the optimization problem 7.7, we have

$$\nabla_{\alpha\alpha}^2 L_z - \mathbf{w}^T \nabla_{\alpha,\gamma}^2 L_z = \frac{1}{n} \sum_{i=1}^n K_h(X'_{i1} - z) (\eta(\mathbf{X}'_{i2}) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{X}'_i)) \eta(\mathbf{X}'_{i2}) \geq b - C\sqrt{\frac{\log n}{nh}}.$$

Note that we set $\lambda_2 = b - C\sqrt{\frac{\log n}{nh}}$. Applying Lemma B.4, we have

$$\mathbb{E}K_h(X_{i1} - z)(\eta(\mathbf{X}_{i2}) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{X}_i))\eta(\mathbf{X}_{i2}) \geq b - C\sqrt{\frac{\log n}{nh}}. \quad (\text{B.18})$$

Denote $\varphi(\mathbf{X}_i) := \eta(\mathbf{X}_{i2}) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{X}_i)$. We then have

$$\begin{aligned} q^2(z) &= \mathbb{E}K_h(X_{i1} - z)^2(\eta(\mathbf{X}_{i2}) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{X}_i))^2 \\ &= \int_{\mathbf{x} \in \mathcal{X}^d} \frac{1}{h^2} K\left(\frac{x_1 - z}{h}\right)^2 \varphi(\mathbf{x})^2 p_{X_1, \mathbf{X}_2}(x_1, \mathbf{x}_2) dx_1 d\mathbf{x}_2 \\ &= \frac{1}{h} \iint_{u \in \mathcal{X}, \mathbf{x}_2 \in \mathcal{X}^{d-1}} K(u)^2 \varphi(z + hu, \mathbf{x}_2)^2 p_{X_1, \mathbf{X}_2}(z + hu, \mathbf{x}_2) du d\mathbf{x}_2 \\ &\geq \frac{b}{h} \iint_{u \in \mathcal{X}, \mathbf{x}_2 \in \mathcal{X}^{d-1}} K(u)^2 \varphi(z + hu, \mathbf{x}_2)^2 du d\mathbf{x}_2 \\ &\geq \frac{b \left(\iint_{u \in \mathcal{X}, \mathbf{x}_2 \in \mathcal{X}^{d-1}} K(u) \varphi(z + hu, \mathbf{x}_2) du d\mathbf{x}_2 \right)^2}{h \iint_{u \in \mathcal{X}, \mathbf{x}_2 \in \mathcal{X}^{d-1}} 1 du d\mathbf{x}_2} \\ &\geq \frac{b}{h \|\mathcal{X}\|^d} \cdot \left(\iint_{u \in \mathcal{X}, \mathbf{x}_2 \in \mathcal{X}^{d-1}} K(u) \varphi(z + hu, \mathbf{x}_2) \eta(x_2) p_{X_1, \mathbf{X}_2}(z + hu, \mathbf{x}_2) / B^2 du d\mathbf{x}_2 \right)^2 \\ &= \frac{b}{h \|\mathcal{X}\|^d B^4} \cdot \left(\mathbb{E}K_h(X_{i1} - z)(\eta(\mathbf{X}_{i2}) - \tilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{X}_i))\eta(\mathbf{X}_{i2}) \right)^2 \\ &\geq \frac{b}{h \|\mathcal{X}\|^d B^4} \cdot \left(b - C\sqrt{\frac{\log n}{nh}} \right)^2 \\ &\geq C \cdot h^{-1}, \end{aligned}$$

where in the third equation we set $x_1 = z + hu$, in the first inequality we used $p_{X_1, \mathbf{X}_2}(z + hu, \mathbf{x}_2) \geq b$, in the second inequality we applied Cauchy inequality, in the third inequality we used $\eta(x_2)p_{X_1, \mathbf{X}_2}(z + hu, \mathbf{x}_2) \leq B^2$ and in the fourth inequality we used (B.18). From Lemma B.9, with probability $1 - \frac{c}{n}$,

$$\sup_z \left| \frac{\hat{q}^2(z)}{q^2(z)} - 1 \right| \leq C \left(\frac{\sqrt{\log n}}{\sqrt{nh}} \right).$$

Since $\left| \frac{\hat{q}^2(z)}{q^2(z)} - 1 \right| = \left| \left(\frac{\hat{q}(z)}{q(z)} - 1 \right) \left(\frac{\hat{q}(z)}{q(z)} + 1 \right) \right|$ and $\frac{\hat{q}(z)}{q(z)} + 1$ converges to a constant, we have

$$\sup_z \left| \frac{\hat{q}(z)}{q(z)} - 1 \right| \leq C \left(\frac{\sqrt{\log n}}{\sqrt{nh}} \right), \quad (\text{B.19})$$

with probability $1 - \frac{c}{n}$. Similarly, with probability $1 - \frac{7}{n}$,

$$\sup_z \left| \frac{\hat{\sigma}}{\sigma} - 1 \right| \leq Cr_n. \quad (\text{B.20})$$

Finally, since

$$\sup_{z \in \mathcal{X}} \left| \frac{\hat{\sigma} \hat{q}(z)}{\sigma q(z)} - 1 \right| \leq \sup_{z \in \mathcal{X}} \left| \frac{\hat{\sigma}}{\sigma} - 1 \right| + \sup_{z \in \mathcal{X}} \left| \frac{\hat{q}(z)}{q(z)} - 1 \right| \cdot \left| \frac{\hat{\sigma}}{\sigma} \right|,$$

with probability $1 - \frac{c}{n}$,

$$\sup_{z \in \mathcal{X}} \left| \frac{\widehat{\sigma} \widehat{q}(z)}{\sigma q(z)} - 1 \right| \leq C \left(\frac{\sqrt{\log n}}{\sqrt{nh}} + r_n \right),$$

which completes the proof. \square

Finally, we prove Lemma 7.10, Lemma 7.11 and Lemma 7.12. Recall that

$$\begin{aligned} \mathbb{G}_n(z) &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\varepsilon_i K_h(X_{i1} - z) (\eta(X_{i2}) - \widetilde{\mathbf{w}}(z)^T \boldsymbol{\psi}(\mathbf{X}_i))}{\sigma q(z)} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\varepsilon_i t(\mathbf{X}_i, z)}{\sigma q(z)}, \\ \mathbb{G}_n^{(1)}(z) &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \left[\frac{\varepsilon_i I(|\varepsilon_i| \leq b_n)}{\sigma} \cdot \frac{K_h(X_{i1} - z) (\eta(\mathbf{X}_{i2}) - \widetilde{\mathbf{w}}(z)^T \boldsymbol{\psi}(\mathbf{X}_i))}{q(z)} \right. \\ &\quad \left. - \mathbb{E} \left(\frac{\varepsilon_i I(|\varepsilon_i| \leq b_n)}{\sigma} \cdot \frac{K_h(X_{i1} - z) (\eta(\mathbf{X}_{i2}) - \widetilde{\mathbf{w}}(z)^T \boldsymbol{\psi}(\mathbf{X}_i))}{q(z)} \right) \right], \\ &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \left[\frac{\varepsilon_i I(|\varepsilon_i| \leq b_n) t(\mathbf{X}_i, z)}{\sigma q(z)} - \mathbb{E} \left(\frac{\varepsilon_i I(|\varepsilon_i| \leq b_n) t(\mathbf{X}_i, z)}{\sigma q(z)} \right) \right], \\ \mathbb{G}_n^{(2)}(z) &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\xi_i \varepsilon_i I(|\varepsilon_i| \leq b_n) \cdot K_h(X_{i1} - z) (\eta(\mathbf{X}_{i2}) - \widetilde{\mathbf{w}}^T \boldsymbol{\psi}(\mathbf{X}_i))}{\sigma q(z)} \\ &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\xi_i \varepsilon_i I(|\varepsilon_i| \leq b_n) t(\mathbf{X}_i, z)}{\sigma q(z)}, \\ \widehat{\mathbb{G}}_n(z) &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \xi_i \cdot \frac{\widehat{\varepsilon}_i K_h(X_{i1} - z) (\eta(\mathbf{X}_{i2}) - \widetilde{\mathbf{w}}(z)^T \boldsymbol{\psi}(\mathbf{X}_i))}{\widehat{\sigma} \widehat{q}(z)} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\xi_i \widehat{\varepsilon}_i t(\mathbf{X}_i, z)}{\widehat{\sigma} \widehat{q}(z)}, \end{aligned}$$

and

$$W_n = \sup_{z \in \mathcal{X}} \mathbb{G}_n(z), \quad W_n^{(1)} = \sup_{z \in \mathcal{X}} \mathbb{G}_n^{(1)}(z), \quad W_n^{(2)}(j_n) = \sup_{z \in \mathcal{X}} \mathbb{G}_n^{(2)}(z) \quad \text{and} \quad \widehat{W}_n(j_n) = \sup_{z \in \mathcal{X}} \widehat{\mathbb{G}}_n(z).$$

Proof of Lemma 7.10. Since ε_i are subgaussian, with probability $1 - \frac{2}{n}$,

$$\max_i |\varepsilon_i| \leq b_n.$$

Thus, with probability $1 - \frac{2}{n}$,

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\varepsilon_i t(\mathbf{X}_i, z)}{\sigma q(z)} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\varepsilon_i I(|\varepsilon_i| \leq b_n) t(\mathbf{X}_i, z)}{\sigma q(z)}. \quad (\text{B.21})$$

Therefore, we just have to bound

$$\sqrt{n} \mathbb{E} \left(\frac{\varepsilon_i I(|\varepsilon_i| \leq b_n) t(\mathbf{X}_i, z)}{\sigma q(z)} \right).$$

For $\mathbb{E}\varepsilon_i I(|\varepsilon_i| \leq b_n)$, we have

$$\begin{aligned}
|\mathbb{E}\varepsilon_i I(|\varepsilon_i| \leq b_n)| &= |\mathbb{E}\varepsilon - \mathbb{E}\varepsilon_i I(|\varepsilon_i| > b_n)| \\
&= |\mathbb{E}\varepsilon_i I(|\varepsilon_i| > b_n)| \\
&\leq \mathbb{E}|\varepsilon_i| \cdot I(|\varepsilon_i| > b_n) \\
&= \int_0^\infty \mathbb{P}(|\varepsilon_i| \cdot I(|\varepsilon_i| > b_n) \geq t) dt \\
&= b_n \cdot P(|\varepsilon_i| \geq b_n) + \int_{b_n}^\infty \mathbb{P}(|\varepsilon_i| \geq t) dt \\
&\leq C \frac{\log n}{n} + \int_{b_n}^\infty e^{-\frac{t^2}{2\sigma^2}} dt \leq C \frac{\log n}{n}.
\end{aligned}$$

Therefore

$$\begin{aligned}
\sqrt{n} \mathbb{E} \left(\frac{\varepsilon_i I(|\varepsilon_i| \leq b_n) t(\mathbf{X}_i, z)}{\sigma q(z)} \right) &\leq \sqrt{n} \mathbb{E} \left(\frac{|\varepsilon_i I(|\varepsilon_i| \leq b_n)| \cdot |t(\mathbf{X}_i, z)|}{\sigma q(z)} \right) \\
&\leq C \frac{\log n h^{1/2}}{\sigma n^{1/2}} \cdot \mathbb{E}(|t(\mathbf{X}_i, z)|) \\
&\leq C \frac{4(L+1)B^2 \log n h^{1/2}}{n^{1/2} \sigma b} \leq C \frac{\log n \cdot h^{1/2}}{n^{1/2}}. \tag{B.22}
\end{aligned}$$

Hence, combining (B.21) and (B.22), with probability $1 - \frac{2}{n}$,

$$|W_n - W_n^{(1)}| \leq C \frac{\log n \cdot h^{1/2}}{n^{1/2}} \leq C n^{-c},$$

which completes the proof. \square

Proof of Lemma 7.11. In this proof, our goal is to bound $W_n^{(1)} - W_n^{(2)}$. Since $W_n^{(1)} = \sup_{z \in \mathcal{X}} \mathbb{G}_n^{(1)}(z)$, $W_n^{(2)}(\mathbf{j}_n) = \sup_{z \in \mathcal{X}} \mathbb{G}_n^{(2)}(z)$ and $\mathbb{G}_n^{(1)}(z), \mathbb{G}_n^{(2)}(z)$ have the same denominators, we will focus on the numerators of them. We will define $\mathbb{F}_n^{(1)}$ as the numerator of $\mathbb{G}_n^{(1)}(z)$ and $\mathbb{F}_n^{(2)}$ as a term very close to the numerator of $\mathbb{G}_n^{(2)}(z)$. Firstly, we bound the different between $\mathbb{F}_n^{(1)}$ and $\mathbb{F}_n^{(2)}$. Then we bound the difference between $\mathbb{F}_n^{(2)}$ and the numerator of $\mathbb{G}_n^{(2)}(z)$. Finally, we give the bound of $W_n^{(1)} - W_n^{(2)}(\mathbf{j}_n)$.

Let

$$\begin{aligned}
\mathbb{F}_n^{(1)}(z) &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \left[\varepsilon_i I(|\varepsilon_i| \leq b_n) t(\mathbf{X}_i, z) - \mathbb{E}(\varepsilon_i I(|\varepsilon_i| \leq b_n) t(\mathbf{X}_i, z)) \right], \\
\mathbb{F}_n^{(2)}(z) &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \xi_i \cdot \left[\varepsilon_i I(|\varepsilon_i| \leq b_n) t(\mathbf{X}_i, z) - \mathbb{E}_n(\varepsilon_i I(|\varepsilon_i| \leq b_n) t(\mathbf{X}_i, z)) \right]
\end{aligned}$$

and

$$M_n^{(1)} = \sup_{z \in \mathcal{X}} \mathbb{F}_n^{(1)}(z) \quad \text{and} \quad M_n^{(2)}(\mathbf{j}_n) = \sup_{z \in \mathcal{X}} \mathbb{F}_n^{(2)}(z),$$

where $\mathbf{j}_n := \{(\mathbf{X}_1, \varepsilon_1), \dots, (\mathbf{X}_n, \varepsilon_n)\}$. Note that $\mathbb{F}_n^{(1)}(z)$ is a function of $z, \mathbf{X}_i, \varepsilon_i$, where $\mathbf{X}_i, \varepsilon_i$ are random variables. However, when analyzing $\mathbb{F}_n^{(2)}(z)$, we consider $\mathbf{X}_i, \varepsilon_i$ as fixed samples and the randomness comes from ξ_i . Let

$$\mathcal{S}_h = \{S_z(\epsilon, \mathbf{X}) = \epsilon I(|\epsilon| \leq b_n) t(\mathbf{X}, z) | z \in \mathcal{X}\}.$$

We are going to prove that \mathcal{S}_h is a VC-type class (see Definition C.4). First,

$$\|S_z(\mathbf{X}, \epsilon)\|_\infty \leq b_n \frac{1}{h} (L+1) B^2 = C \frac{\log n}{h} =: b.$$

According to Lemma C.2,

$$N(S_h, L_2(\mathbb{Q}), b\tau) \leq \left(\frac{2\|K\|_{tv} A}{\tau} \right)^4,$$

for any finite measure \mathbb{Q} . Thus, the parameters (in terms of Definition C.4, there are 4 parameters for VC-type class) for \mathcal{S}_h are

$$\begin{cases} v = 4 = \mathcal{O}(1) \\ a = 2\|K\|_{TV} A = \mathcal{O}(1) \\ b = \mathcal{O}(\log n \cdot h^{-1}) \\ \sigma^2 = \mathbb{E}S_z^2(\epsilon, \mathbf{X}) = \sigma^2 \cdot \frac{C}{h} = \mathcal{O}(h^{-1}) \end{cases},$$

and

$$K_h := Av(\log n \vee \log(ab/\sigma)) = \mathcal{O}(\log n).$$

Applying Lemma C.5, there exists a tight Gaussian random element $B = \{B(g) : g \in \mathcal{S}_h\}$ such that

$$\mathbb{E}[B(g_1)B(g_2)] = \mathbb{E}[g_1(\mathbf{X}_1, \varepsilon_1)g_2(\mathbf{X}_1, \varepsilon_1)] - \mathbb{E}[g_1(\mathbf{X}_1, \varepsilon_1)]\mathbb{E}[g_2(\mathbf{X}_1, \varepsilon_1)],$$

for all $g_1, g_2 \in \mathcal{S}_h$. Assume $W^0 := \sup_z B(S_z)$. We have for any $\gamma \in (0, 1)$,

$$\mathbb{P}\left(|M_n^{(1)} - W^0| > C \cdot \left(\frac{\log^2 n}{\gamma^{1/2} h n^{1/2}} + \frac{\log^{5/4} n}{\gamma^{1/2} n^{1/4} h^{3/4}} + \frac{\log n}{\gamma^{1/3} n^{1/6} h^{2/3}} \right)\right) \leq A' \left(\gamma + \frac{\log n}{n} \right),$$

where A' is an absolute constant. Under conditions in Theorem 4.2, since $h \asymp n^{-\delta_2}$, $\delta_2 < 1/4$, there exists c such that $\frac{1}{n^{1/6} h^{2/3}} \leq n^{-c}$. Hence, with $\gamma = n^{-3c/2}$, there exist constants c_1, c_2 such that

$$\mathbb{P}\left(|M_n^{(1)} - W^0| > c_1 n^{-c}\right) \leq c_2 n^{-c}. \quad (\text{B.23})$$

By applying Lemma C.6 to W^0 and $M_n^{(2)}(\mathbf{J}_n)$ we can find that $b^2 K_n \leq n\sigma^2$ is equivalent to $\frac{\log^3 n}{nh} \leq C$ and

$$\phi_n = C \cdot \frac{\log^{1/2} n}{h^{1/2} n^{1/2}} + C \cdot \frac{\log^{5/4} n}{n^{1/4} h^{3/4}} \quad \text{and} \quad \gamma_n(\delta) = C \cdot \frac{1}{\delta} \cdot \frac{\log^{5/4} n}{n^{1/4} h^{3/4}} + \frac{1}{n}.$$

Therefore, if $\frac{\log^3 n}{nh} \leq C$, there exists a $S_n \in (\mathcal{X}^d)^n \times [-b_n, b_n]^n$, such that $\mathbb{P}(\mathbf{J}_n \in S_n) \geq 1 - \frac{3}{n}$, and for a fixed $\mathbf{J}_n = \mathbf{j}_n \in S_n$, we have for all $\delta > 0$,

$$\mathbb{P}\left(|M_n^{(2)}(\mathbf{j}_n) - W^0| > C \cdot \frac{\log^{1/2} n}{h^{1/2} n^{1/2}} + C \cdot \frac{\log^{5/4} n}{n^{1/4} h^{3/4}} + \delta\right) \leq A'' \cdot \left(C \cdot \frac{1}{\delta} \cdot \frac{\log^{5/4} n}{n^{1/4} h^{3/4}} + \frac{1}{n} \right),$$

where A'' is an absolute constant. Similarly, under conditions in Theorem 4.2, since $h \asymp n^{-\delta_2}$, $\delta_2 < 1/4$, there also exists c_1, c_2, c , such that

$$\mathbb{P}\left(|M_n^{(2)}(\mathbf{j}_n) - W^0| > c_1 n^{-c}\right) \leq c_2 n^{-c}. \quad (\text{B.24})$$

Combining (B.23) and (B.24), we have for any $\mathbf{j}_n \in S_n$,

$$\mathbb{P}\left(|M_n^{(2)}(\mathbf{j}_n) - M_n^{(1)}| > c_1 n^{-c}\right) \leq c_2 n^{-c}.$$

Next, we are going to bound:

$$\left(\frac{1}{\sqrt{n}} \sum_{i=1}^n \xi_i\right) \cdot \mathbb{E}_n(\varepsilon_i I(|\varepsilon_i| \leq b_n) t(\mathbf{X}_i, z)).$$

Applying Lemma A.1, we have with probability $1 - \frac{1}{n}$,

$$\begin{aligned} \left|\mathbb{E}_n(\varepsilon_i I(|\varepsilon_i| \leq b_n) t(\mathbf{X}_i, z))\right| &= \left|\frac{1}{n} \sum_{i=1}^n \varepsilon_i I(|\varepsilon_i| \leq b_n) t(\mathbf{X}_i, z)\right| \\ &\leq \frac{1}{n} \sum_{i=1}^n b_n (L+1) B K_h(X_{i1} - z) \leq C \frac{\log^{3/2} n}{n^{1/2} h^{1/2}}. \end{aligned}$$

Since $\left(\frac{1}{\sqrt{n}} \sum_{i=1}^n \xi_i\right) \sim N(0, 1)$, with probability $1 - \frac{2}{n}$,

$$\left|\frac{1}{\sqrt{n}} \sum_{i=1}^n \xi_i\right| \leq C \log n.$$

Hence, with probability $1 - \frac{3}{n}$,

$$\left|\left(\frac{1}{\sqrt{n}} \sum_{i=1}^n \xi_i\right) \cdot \mathbb{E}_n(\varepsilon_i I(|\varepsilon_i| \leq b_n) t(\mathbf{X}_i, z))\right| \leq C \frac{\log^{5/2} n}{n^{1/2} h^{1/2}}.$$

Therefore, under conditions in Theorem 4.2, there exists a set S_n , such that $\mathbb{P}(\mathbf{J}_n \in S_n) \geq 1 - \frac{c}{n}$, and for a fixed $\mathbf{J}_n = \mathbf{j}_n \in S_n$, with probability $1 - \frac{c}{n} - c_2 n^{-c} = 1 - c'_2 n^{-c}$,

$$\begin{aligned} |W_n^{(1)} - W_n^{(2)}(\mathbf{j}_n)| &\leq \frac{4B}{\sigma b} \left(|M_n^{(1)} - M_n^{(2)}| + \sup_z \left| \mathbb{E}_n(\varepsilon_i I(|\varepsilon_i| \leq b_n) t(\mathbf{X}_i, z)) \right| \right) \\ &\leq C \cdot \left(c_1 n^{-c} + \frac{\log^{5/2} n}{n^{1/2} h^{1/2}} \right) \leq c'_1 n^{-c}. \end{aligned}$$

□

Proof of lemma 7.12. We introduce an intermediate process:

$$\mathbb{G}_n^{(3)}(z) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\xi_i \varepsilon_i I(|\varepsilon_i| \leq b_n) t(\mathbf{X}_i, z)}{\widehat{\sigma} \widehat{q}(z)}.$$

Let $W_n^{(3)}(\mathbf{j}_n) = \sup_{z \in \mathcal{X}} \mathbb{G}_n^{(3)}(\mathbf{j}_n)$, the big picture of this proof is that we firstly bound $W_n^{(2)}(\mathbf{j}_n) - W_n^{(3)}(\mathbf{j}_n)$ by some properties of Gaussian process and then we bound $W_n^{(3)}(\mathbf{j}_n) - \widehat{W}_n(\mathbf{j}_n)$ directly. Combining them, we get the bound for $W_n^{(2)}(\mathbf{j}_n) - \widehat{W}_n(\mathbf{j}_n)$.

Define

$$\begin{aligned} \Delta \mathbb{G}_n(z) &:= \widehat{\mathbb{G}}_n^{(3)}(z) - \mathbb{G}_n^{(2)}(z) \\ &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\xi_i \varepsilon_i I(|\varepsilon_i| \leq b_n) t(\mathbf{X}_i, z)}{\sigma q(z)} \left(\frac{\sigma q(z)}{\widehat{\sigma} \widehat{q}(z)} - 1 \right). \end{aligned}$$

Note that $\Delta \mathbb{G}_n(z)$ is a Gaussian process with mean zero and

$$\text{Var}(\Delta \mathbb{G}_n(z)) \leq \frac{1}{n} \sum_{i=1}^n \frac{b_n^2 t^2(\mathbf{X}_i, z)}{\sigma^2 q^2(z)} \left(\frac{\sigma q(z)}{\widehat{\sigma} \widehat{q}(z)} - 1 \right)^2 \leq C \cdot \log^2 n \cdot \left(r_n + \sqrt{\frac{\log n}{nh}} \right)^2 =: \sigma_\Delta, \quad (\text{B.25})$$

where in the last inequality we use Lemma 7.9, Lemma A.1, $q^2(z) \geq Ch^{-1}$, and $t(\mathbf{X}_i, z) \leq Ch^{-1}$. Under conditions in Theorem 4.2, there exists $c, \delta > 0$ such that

$$r_n + \sqrt{\frac{\log n}{nh}} < n^{-c-\delta}.$$

Therefore, $\sigma_\Delta < n^{-2c}$. From Lemma 7.9, for a large enough n , there exists a bounded set $S_0 \in \mathcal{X}^d \times R^d$, such that $P(\mathbf{J}_n = (\mathbf{X}_1^n, \varepsilon_1^n) \in S_0) \geq 1 - 3/n$ and, for $\mathbf{j}_n \in S_0$,

$$\left(\frac{\sigma q(z)}{\widehat{\sigma} \widehat{q}(z)} - 1 \right) \leq 2.$$

Consider the following function class:

$$\mathcal{K} := \left\{ S_z(\mathbf{x}, \varepsilon) = \frac{\varepsilon I(|\varepsilon| \leq b_n) t(\mathbf{x}, z)}{\sigma q(z)} \left(\frac{\sigma q(z)}{\widehat{\sigma} \widehat{q}(z)} - 1 \right) \middle| z \in \mathcal{X}, \mathbf{J}_n \in S_0 \right\}.$$

We are going to show that \mathcal{K} is a VC-type class. First,

$$\|S_z(\mathbf{x}, \varepsilon)\|_\infty \leq C \log nh^{-1/2} K(h^{-1}(x_1 - z)) \leq \log nh^{-1/2} =: b.$$

According to Lemma C.2,

$$N(\mathcal{K}, L_2(\mathbb{Q}), b\tau) \leq \left(\frac{2\|K\|_{tv} A}{\tau} \right)^4,$$

for any finite measure \mathbb{Q} . Thus the parameters for this VC-type class function are

$$\begin{cases} v = 4 = \mathcal{O}(1) \\ a = 2\|K\|_{TV} A = \mathcal{O}(1) \\ b = \mathcal{O}(\log n \cdot h^{-1/2}) \\ \sigma^2 = \mathbb{E} S_z^2(\mathbf{x}, \varepsilon) \leq C \mathbb{E} \frac{(\log n)^2 K_h(X_{i1} - z)^2}{q^2(z)} = \mathcal{O}(1). \end{cases}$$

In addition,

$$\mathbb{E} \left[\left(\Delta \mathbb{G}_n(z_1) - \Delta \mathbb{G}_n(z_2) \right)^2 \right] \leq \mathbb{E}_n \left[\left(S_{z_1}(\mathbf{X}_i, \varepsilon_i) - S_{z_2}(\mathbf{X}_i, \varepsilon_i) \right)^2 \right],$$

for all $z_1, z_2 \in \mathcal{X}$. The covering number for the index set \mathcal{X} with respect to the intrinsic semi-metric induced from the Gaussian process $\Delta \mathbb{G}_n(z)$ is bounded by the uniform covering number for the function class \mathcal{K} . Therefore, an application of Corollary 2.2.8 in [Vaart and Wellner \(1997\)](#) gives,

$$\begin{aligned} \mathbb{E} \left[\sup_{z \in \mathcal{X}} |\Delta \mathbb{G}_n(z)| \right] &\leq C \log^2 n \left(r_n + \sqrt{\frac{\log n}{nh}} \right)^2 \sqrt{(1+v) \log \left(\frac{ab}{\log^2 n (r_n + \sqrt{\log n / (nh)^{-1/2}})^2} \right)} \\ &\leq C \log^2 n \left(r_n + \sqrt{\frac{\log n}{nh}} \right)^2 \sqrt{\log \left(\frac{h^{-1/2}}{\log n (r_n + \sqrt{\log n / (nh)^{-1/2}})^2} \right)} \\ &\leq n^{-2c}. \end{aligned}$$

With Borell–TIS inequality [Borell \(1975\)](#), for all $u > 0$,

$$\mathbb{P} \left(\sup_{z \in \mathcal{X}} |\Delta \mathbb{G}_n(z)| \geq \mathbb{E} \left[\sup_{z \in \mathcal{X}} |\Delta \mathbb{G}_n(z)| \right] + u \right) \leq \exp \left(\frac{-u^2}{2\sigma_\Delta^2} \right).$$

Therefore, setting $u = \sqrt{\log n}/n^c$, with probability $1 - 2/n$,

$$\sup_{z \in \mathcal{X}} |\Delta \mathbb{G}_n(z)| \leq \mathbb{E} \left[\sup_{z \in \mathcal{X}} |\Delta \mathbb{G}_n(z)| \right] + u \leq n^{-c+\delta}.$$

Since $W_n^{(3)}(\mathbf{j}_n) = \sup_{z \in \mathcal{X}} \mathbb{G}_n^{(3)}(\mathbf{j}_n)$, we have

$$|W_n^{(2)}(\mathbf{j}_n) - W_n^{(3)}(\mathbf{j}_n)| \leq \sup_{z \in \mathcal{X}} |\Delta \mathbb{G}_n(z)|.$$

Hence, for any $\mathbf{j}_n \in S_0$,

$$\mathbb{P} \left(|W_n^{(2)}(\mathbf{j}_n) - W_n^{(3)}(\mathbf{j}_n)| > n^{-c+\delta} \right) \leq \frac{2}{n}. \quad (\text{B.26})$$

At last, since ε_i are subgaussian, with probability $1 - 3/n$,

$$\mathbb{G}_n^{(3)}(z) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\xi_i \varepsilon_i I(|\varepsilon_i| \leq b_n) t(\mathbf{X}_i, z)}{\widehat{\sigma} \widehat{q}(z)} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\xi_i \varepsilon_i t(\mathbf{X}_i, z)}{\widehat{\sigma} \widehat{q}(z)}.$$

Thus,

$$\mathbb{G}_n^{(3)}(z) - \widehat{\mathbb{G}}_n(z) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\xi_i (\varepsilon_i - \widehat{\varepsilon}_i) t(\mathbf{X}_i, z)}{\widehat{\sigma} \widehat{q}(z)}.$$

Since, with probability $1 - c/n$,

$$\max_i |\varepsilon_i - \widehat{\varepsilon}_i| \leq C r_n,$$

we have

$$\max_i |\xi_i| \leq C \log n,$$

with probability $1 - 2/n$. Therefore,

$$\begin{aligned} \sup_z |\mathbb{G}_n^{(3)}(z) - \widehat{\mathbb{G}}_n(z)| &\leq \sup_z \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\xi_i(\varepsilon_i - \widehat{\varepsilon}_i)t(\mathbf{X}_i, z)}{\widehat{\sigma}\widehat{q}(z)} \right| \\ &\leq \sup_z \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{C \log n \cdot r_n |t(\mathbf{X}_i, z)|}{\widehat{\sigma}\widehat{q}(z)} \leq C \log n \cdot r_n. \end{aligned}$$

and

$$\mathbb{P}\left(|W_n^{(3)}(\mathbf{j}_n) - \widehat{W}_n(\mathbf{j}_n)| > Cr_n \cdot \log n\right) \leq \frac{c}{n}. \quad (\text{B.27})$$

Combining (B.26) and (B.27), for any $\mathbf{j}_n \in S_0$, there exists c_1, c_2, c such that

$$\mathbb{P}\left(|W_n^{(2)}(\mathbf{j}_n) - \widehat{W}_n(\mathbf{j}_n)| > c_1 n^{-c}\right) \leq c_2 n^{-c}.$$

□

C Collection of Known Results on Empirical Processes

In this section, we provide some known results on empirical processes that are used in the technical proofs.

Lemma C.1 (Theorem 3.12, [Koltchinskii \(2011\)](#)). Assume that the functions in \mathcal{F} , defined on \mathcal{X} , are uniformly bounded by a constant U and $F(\cdot)$ is the envelope of \mathcal{F} , such that, $|f(x)| \leq F(x)$ for all $x \in \mathcal{X}$ and $f \in \mathcal{F}$. Define $\sigma_P^2 = \sup_{f \in \mathcal{F}} \mathbb{E}[f^2]$. Let X_1, \dots, X_n be i.i.d. copies of the random variable X . We denote the empirical measure as $\mathbb{P}_n = n^{-1} \sum_{i=1}^n \delta_{X_{i1}}$. If for some $A, V \geq 0$ and for all $\varepsilon > 0$ and $n \geq 1$, the covering entropy satisfies

$$N(\mathcal{F}, L^2(\mathbb{P}_n); \varepsilon) \leq \left(\frac{A \|F\|_{L^2(\mathbb{P}_n)}}{\varepsilon} \right)^V,$$

then for any i.i.d. subgaussian mean zero random variables $\varepsilon_1, \dots, \varepsilon_n$ there exists a universal constant C such that

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(X_{i1}) \right] \leq C \left[\sqrt{\frac{V}{n}} \sigma_P \sqrt{\log \frac{A \|F\|_{L^2(\mathbb{P})}}{\sigma_P}} + \frac{VU}{n} \log \frac{A \|F\|_{L^2(\mathbb{P})}}{\sigma_P} \right].$$

Furthermore, we also have

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (f(X_{i1}) - \mathbb{E}f(X)) \right] \leq C \left[\sqrt{\frac{V}{n}} \sigma_P \sqrt{\log \frac{A \|F\|_{L^2(\mathbb{P})}}{\sigma_P}} + \frac{VU}{n} \log \frac{A \|F\|_{L^2(\mathbb{P})}}{\sigma_P} \right].$$

Lemma C.2 (Lemma 3, [Giné and Nickl \(2009\)](#)). Let $K : \mathbb{R} \mapsto \mathbb{R}$ be a bounded variation function. Define the function class $\mathcal{F}_h = \{K((t - \cdot)/h) \mid t \in \mathbb{R}\}$. There exists $A < \infty$ such that for all probability measures Q on \mathbb{R} , we have

$$\sup_Q N(\mathcal{F}_h, L^2(Q), \epsilon) \leq \left(\frac{2\|K\|_{\text{TV}} A}{\epsilon} \right)^4, \text{ for any } \epsilon \in (0, 1).$$

Lemma C.3 ([Bousquet \(2002\)](#)). Let X_1, \dots, X_n be independent random variables and \mathcal{F} is a function class such that there exist η_n and τ_n^2 satisfying

$$\sup_{f \in \mathcal{F}} \|f\|_\infty \leq \eta_n \quad \text{and} \quad \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \text{Var}(f(X_{i1})) \leq \tau_n^2.$$

Define the random variable Z being the suprema of an empirical process

$$Z = \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n (f(X_{i1}) - \mathbb{E}f(X_{i1})) \right|.$$

Then for any $z > 0$, we have the following concentration inequality on the suprema

$$\mathbb{P} \left(Z \geq \mathbb{E}Z + z \sqrt{2(\tau_n^2 + 2\eta_n \mathbb{E}Z)} + 2z^2 \eta_n / 3 \right) \leq \exp(-nz^2).$$

Definition C.4 (VC-type class, [Chernozhukov et al. \(2014a\)](#)). Let \mathcal{G} be a class of measurable functions on a measurable space (S, \mathcal{S}) , and let $b > 0, a \geq e$ and $v \geq 1$ be some constants. Then the class \mathcal{G} is called $\text{VC}(b, a, v)$ type class if it is uniformly bounded in absolute value by b , and the covering numbers of \mathcal{G} satisfy

$$\sup_Q N(\mathcal{G}, L_2(Q), b\tau) \leq (a/\tau)^v, \quad 0 < \tau < 1,$$

where the supremum is taken over all finitely discrete probability measures Q on (S, \mathcal{S}) .

Lemma C.5 (Slepian–Stein type coupling for suprema of empirical processes, Theorem A.1 in [Chernozhukov et al. \(2014a\)](#)). Let X_1, \dots, X_n be i.i.d random variables taking values in a measurable space (S, \mathcal{S}) . Let \mathcal{G} be a pointwise-measurable $\text{VC}(b, a, v)$ type function class for some $b > 0, a \geq e$ and $v \geq 1$. Let $\sigma^2 > 0$ be any constant such that $\sup_{g \in \mathcal{G}} \mathbb{E}[g(X_1)^2] \leq \sigma^2 \leq b^2$. Define the empirical process

$$\mathbb{G}_n(g) := \frac{1}{\sqrt{n}} \sum_{i=1}^n (g(X_i) - \mathbb{E}[g(X_1)]), \quad g \in \mathcal{G},$$

and let

$$W_n := \|\mathbb{G}_n\|_{\mathcal{G}} := \sup_{g \in \mathcal{G}} |\mathbb{G}_n(g)|$$

denote the supremum of the empirical process. Let $B = \{B(g) \mid g \in \mathcal{G}\}$ be a tight Gaussian random element in $l^\infty(\mathcal{F})$ with mean zero and covariance function

$$\mathbb{E}[B(g_1) B(g_2)] = \mathbb{E}[g_1(X_1) g_2(X_1)] - \mathbb{E}[g_1(X_1)] \mathbb{E}[g_2(X_1)],$$

for all $g_1, g_2 \in \mathcal{G}$. It is well known that such a process exists under the VC type assumption; see [Vaart and Wellner \(1997\)](#), pages 100-101. For some sufficiently large absolute constant A , let

$$K_n := Av(\log n \vee \log(ab/\sigma)).$$

Then for every $\gamma \in (0, 1)$, one can construct a random variable W^0 on an enriched probability space such that: (i) $W^0 \stackrel{d}{=} \|B\|_{\mathcal{G}}$ and (ii)

$$\mathbb{P}\left(|W_n - W^0| > \frac{bK_n}{(\gamma n)^{1/2}} + \frac{(b\sigma)^{1/2}K_n^{3/4}}{\gamma^{1/2}n^{1/4}} + \frac{b^{1/3}\sigma^{2/3}K_n^{2/3}}{\gamma^{1/3}n^{1/6}}\right) \leq A' \left(\gamma + \frac{\log n}{n}\right),$$

where A' is an absolute constant.

Lemma C.6 (Slepian-Stein type coupling for suprema of conditional multiplier processes, Theorem A.2 in [Chernozhukov et al. \(2014a\)](#)). Let ξ_1, \dots, ξ_n be independent $N(0, 1)$ random variables independent of $X_1^n := \{X_1, \dots, X_n\}$, and let $\xi_1^n := \{\xi_1, \dots, \xi_n\}$. Define the Gaussian multiplier process

$$\tilde{\mathbb{G}}_n(g) := \tilde{\mathbb{G}}_n(X_1^n, \xi_1^n) := \frac{1}{\sqrt{n}} \sum_{i=1}^n \xi_i \left(g(X_i) - \mathbb{E}_n[g(X_i)] \right), \quad g \in \mathcal{G},$$

and for $x_1^n \in \mathcal{S}^n$. Let $\tilde{W}_n(x_1^n) := \|\tilde{\mathbb{G}}_n(x_1^n, \xi_1^n)\|_{\mathcal{G}}$ denote the supremum of this process calculated for a fixed $X_n^n = x_1^n$. In addition, let

$$\phi_n := \sqrt{\frac{\sigma^2 K_n}{n}} + \left(\frac{b^2 \sigma^2 K_n^3}{n} \right)^{1/4} \quad \text{and} \quad \gamma_n(\delta) := \frac{1}{\delta} \left(\frac{b^2 \sigma^2 K_n^3}{n} \right)^{1/4} + \frac{1}{n}.$$

Suppose that $b^2 K_n \leq n\sigma^2$. Then for every $\delta > 0$, there exist a set $S_{n,0} \in \mathcal{S}^n$ such that $\mathbb{P}(X_1^n \in S_{n,0}) \geq 1 - \frac{3}{n}$ and, for every $x_1^n \in S_{n,0}$, one can construct a random variable W^0 on an enriched probability space such that: $W^0 \stackrel{d}{=} \|B\|_{\mathcal{G}}$ and (ii)

$$\mathbb{P}\left(|\tilde{W}_n(x_1^n) - W^0| > (\phi_n + \delta)\right) \leq A'' \gamma_n(\delta),$$

where A'' is an absolute constant.

Lemma C.7. Let \mathcal{F}_1 and \mathcal{F}_2 be two function classes defined on \mathcal{X} , such that $\mathcal{F}_1 \subseteq \mathcal{F}_2$. It holds that

$$N(\mathcal{F}_1, L_2(Q), \epsilon) \leq N(\mathcal{F}_2, L_2(Q), \epsilon/2),$$

for any measure Q on \mathcal{X} .

Proof. For an ϵ -covering set N_2 of \mathcal{F}_2 such that $|N_2| = N(\mathcal{F}_2, L_2(Q), \epsilon/2)$. We define an ϵ -covering set N_1 for \mathcal{F}_1 as follows: For each $f_2 \in N_2$, we find a function f_1 in $\mathcal{F}_1 \cap \mathbb{B}(f_2, \epsilon/2)$, and all these f_1 forms the set N_1 . Next we prove that N_1 constructed in this way is an ϵ -covering of \mathcal{F}_1 .

By definition of N_2 , for each function f' in \mathcal{F}_1 we have a function $f'_2 \in N_2$ such that $\|f' - f'_2\|_{L^2(Q)} \leq \epsilon/2$. By definition of N_1 , for f'_2 there exists a function $f'_1 \in N_1 \subseteq \mathcal{F}_1$ such that $\|f'_2 - f'_1\|_{L^2(Q)} \leq \epsilon/2$. By triangle inequality,

$$\|f' - f'_1\|_{L^2(Q)} \leq \|f' - f'_2\|_{L^2(Q)} + \|f'_2 - f'_1\|_{L^2(Q)} \leq \epsilon,$$

and hence N_1 is an ϵ -covering of \mathcal{F}_1 . □

Lemma C.8. Let \mathcal{F}_1 and \mathcal{F}_2 be two function classes defined on \mathcal{X} satisfying

$$N(\mathcal{F}_1, L_2(Q), a_1\epsilon) \leq C_1\epsilon^{-v_1} \text{ and } N(\mathcal{F}_2, L_2(Q), a_2\epsilon) \leq C_2\epsilon^{-v_2}$$

for some $C_1, C_2, a_1, a_2, v_1, v_2 > 0$ and any $0 < \epsilon < 1$. Define $\|\mathcal{F}_\ell\|_\infty = \sup\{\|f\|_\infty, f \in \mathcal{F}_\ell\}$ for $\ell = 1, 2$ and $U = \|\mathcal{F}_1\|_\infty \vee \|\mathcal{F}_2\|_\infty$. For the function classes $\mathcal{F}_\times = \{f_1 f_2 \mid f_1 \in \mathcal{F}_1, f_2 \in \mathcal{F}_2\}$ and $\mathcal{F}_+ = \{f_1 + f_2 \mid f_1 \in \mathcal{F}_1, f_2 \in \mathcal{F}_2\}$, we have for any $\epsilon \in (0, 1)$,

$$\begin{aligned} N(\mathcal{F}_\times, L_2(Q), \epsilon) &\leq C_1 C_2 \left(\frac{2a_1 U}{\epsilon} \right)^{v_1} \left(\frac{2a_2 U}{\epsilon} \right)^{v_2}, \\ N(\mathcal{F}_+, L_2(Q), \epsilon) &\leq C_1 C_2 \left(\frac{2a_1 U}{\epsilon} \right)^{v_1} \left(\frac{2a_2 U}{\epsilon} \right)^{v_2}, \end{aligned}$$

for any measure Q on \mathcal{X} .

Proof. For any $\epsilon \in (0, 1)$, let $N_1 = \{f_{11}, \dots, f_{1N_1}\}$ and $N_2 = \{f_{21}, \dots, f_{2N_2}\}$ be the $\epsilon/(2U)$ -net of \mathcal{F}_1 and \mathcal{F}_2 respectively with

$$N_1 \leq C_1 \left(\frac{2a_1 U}{\epsilon} \right)^{v_1} \text{ and } N_2 \leq C_2 \left(\frac{2a_2 U}{\epsilon} \right)^{v_2}.$$

Define the set $N = \{f_{1j} f_{2k} \mid f_{1j} \in N_1, f_{2k} \in N_2\}$. We now show that N is an ϵ -net for \mathcal{F}_\times . For any $f_1 f_2 \in \mathcal{F}$, there exist two functions $f_{1j} \in N_1$ and $f_{2k} \in N_2$ such that $\|f_1 - f_{1j}\|_{L_2(Q)} \leq \epsilon/(2U)$ and $\|f_2 - f_{2k}\|_{L_2(Q)} \leq \epsilon/(2U)$. Moreover, we have $f_{1j} f_{2k} \in N$ and

$$\|f_1 f_2 - f_{1j} f_{2k}\|_{L_2(Q)} \leq \|\mathcal{F}_2\|_\infty \|f_1 - f_{1j}\|_{L_2(Q)} + \|\mathcal{F}_1\|_\infty \|f_2 - f_{2k}\|_{L_2(Q)} \leq \epsilon.$$

Therefore N is the ϵ -net for \mathcal{F}_\times . Similarly, we also have

$$\|(f_1 + f_2) - (f_{1j} + f_{2k})\|_{L_2(Q)} \leq \|f_1 - f_{1j}\|_{L_2(Q)} + \|f_2 - f_{2k}\|_{L_2(Q)} \leq \epsilon/U.$$

So $N' = \{f_{1j} + f_{2k} \mid f_{1j} \in N_1, f_{2k} \in N_2\}$ is the ϵ/U -net of \mathcal{F}_+ . We finally complete the proof by showing that

$$|N'| = |N| = N_1 N_2 \leq C_1 C_2 \left(\frac{2a_1 U}{\epsilon} \right)^{v_1} \left(\frac{2a_2 U}{\epsilon} \right)^{v_2}.$$

□