Nonparametric identification and estimation of a generalized additive model with a flexible additive structure and unknown link *

Songnian Chen[†] Nianqing Liu[‡] Jian Zhang[§] Yahong Zhou[‡]

September, 2023

Abstract

This paper proposes a nonparametric approach to identify and estimate the generalized additive model with a flexible additive structure and with possibly discrete variables when the link function is unknown. Our approach allows for a flexible additive structure which provides the applied researchers the flexibility to specify their model according to economic theory or practical experience. Motivated by the concerns from empirical research, our method also allows for multiple discrete variables in the covariates. By transforming our model into a generalized additive model with univariate component functions, our identification and estimation hence follows a procedure adapted from the case with univariate components. The estimators converge to normal distributions in large sample with a one-dimensional convergence rate for the link function and a d_k -dimensional convergence rate for the link function and a k.

Key words: Generalized additive model, flexible additive structure, nonparametric regression, kernel estimation

JEL classification: C13, C14, C21

1 Introduction

Flexibility in model specification is one of the key features pursued by the applied economists when they use a nonparametric method in applications, since different applications need different model specifications which are given by economic theory or practical experience. In addition, there are potentially multiple discrete covariates in many economic datasets. It is hence an empirical concern to handle discrete covariates appropriately in the estimation procedures of economic applications.

In this paper, we address the above two concerns in estimating the generalized additive model with an unknown link function as

$$H(x) = G\left(\sum_{k=1}^{K} f_k(x^k)\right),\tag{1}$$

^{*}We thank the co-editor, Liangjun Su, and three anonymous referees for comments that have greatly improved the article. We are very grateful to Chunrong Ai, Xiaohong Chen, Yanqin Fan, LungFei Lee and Arthur Lewbel for their constructive comments and suggestions. We also thank Liang Chen, Timothy Christensen, Qingfeng Liu, Ryo Okui, Quang Vuong, Hanghui Zhang, Xinyu Zhang, and seminar/conference participants at HKUST, SUFE, Xiamen University, 2019 Asian Meeting of Econometric Society, 2019 China Meeting of Econometric Society, 2019 Young Econometrician Asia Pacific (YEAP) for helpful comments. The usual disclaimer applies. Zequn Jin provided capable research assistance.

⁺(corresponding author) School of Economics, Zhejiang University, Hangzhou, China, email: snchen2022@zju.edu.cn.

[‡]School of Economics, The Shanghai University of Finance and Economics (SUFE), Shanghai, China.

[§]School of Economics, Nankai University, Tianjin, China.

where $H(\cdot)$ is a function that can be consistently estimated (such as nonparametric regression), but the link function $G(\cdot)$ and the component functions $f_k(\cdot)$'s are unknown with $x \equiv (x^1, \ldots, x^K) \in \mathbb{R}^d$ and $x^k \in \mathbb{R}^{d_k}$. The first concern is addressed by allowing a flexible grouping of the covariates in the subvectors x^k for $k = 1, \ldots, K$. In this way, a researcher can group the covariates according to economic theory or practical experience , instead of having to restrict one or all of sub-vectors to be univariate. The second concern is addressed by allowing discrete covariates in the estimation procedure. The parametric version of this functional restriction has been implemented in many economic applications, including the very popular specification of constant elasticity of substitution (CES) in the estimation of production function. See, e.g., Kmenta (1967); Hodges (1969); Paraskevopoulos (1979); Antras (2004); Klump, McAdam, and Willman (2007); Berkowitz, Ma, and Nishioka (2017), among others.

To identify the model primitives of $G(\cdot)$ and $f_k(\cdot)$'s, we transform the model (1) by a known mapping into a new model with the link function $G(\cdot)$ and some univariate component functions $\tilde{f}_k(\cdot)$'s. We then identify the new model by applying some existing identification approach to the generalized additive model with univariate components. Closely following the identification strategy, we propose a three-step procedure to estimate the link $G(\cdot)$ and the original components $f_k(\cdot)$. The consistency and asymptotic normality is then established for the estimator of the link $G(\cdot)$ at a one-dimensional convergence rate and for the estimator of the component $f_k(\cdot)$ at a d_k -dimensional convergence rate.

Our paper contributes to the estimation of generalized additive model. With a known link function and only univariate component functions, Chen, Härdle, Linton, and Severance-Lossin (1996), Linton and Härdle (1996), Horowitz and Mammen (2004), and Ma (2012), among others, estimated the univariate components at a one-dimensional convergence rate. Their estimators hence have no curse of dimensionality. With an unknown link and only univariate components, Horowitz (2001), Horowitz and Mammen (2007, 2011), and Lin, Pan, Lv, and Zhang (2018), among others, recovered the univariate components still at a one-dimensional convergence rate and hence avoided the curse of dimensionality. Jacho-Chávez, Lewbel, and Linton (2010, JLL hereafter) generalized the framework with only univariate components (and an unknown link) to allow multivariate components, as long as one component function is univariate. Our paper further generalizes the model to allow for a flexible specification of additivity, and the existence of a univariate component function is not needed. In a related area, Lewbel, Lu, and Su (2015) provided a nonparametric test of whether the monotonic transformation structure is correctly specified. With a weaker notion of separability, Pinkse (2001) developed the estimators of $\tilde{f}_1(\cdot), \ldots, \tilde{f}_K(\cdot)$ in a nonparametric regression with weak separability as $E(Y|X = x, Z = z) = \tilde{G}(x, \tilde{f}_1(z^1), \dots, \tilde{f}_K(z^K))$ where \tilde{G} is monotone in $\tilde{f}_1, \dots, \tilde{f}_K$, and furthermore all of $\tilde{f}_1(\cdot), \ldots, \tilde{f}_K(\cdot)$ are monotone in their respective first arguments. He showed that the functions $\tilde{f}_1(\cdot), \ldots, \tilde{f}_K(\cdot)$ can be identified up to a monotonic transformation. The generalized additive model are in general identified up to location and sign-scale normalizations.¹ Our paper is most relevant to Horowitz (2001) and JLL in this research line. To clarify our contributions relative to them, consider model (1) with K = 2 and $d_1, d_2 \ge 2$. Horowitz (2001) identified such a model by further imposing an additive structure on both $f_1(\cdot)$ and $f_2(\cdot)$ as $f_1(x^1) = \sum_{k=1}^{d_1} f_{1k}(x_k^1)$ and $f_2(x^2) = \sum_{k=1}^{d_2} f_{2k}(x_k^2)$. Although such an extra additive structure reduces the dimensionality of

¹Other related papers include Ma and Song (2015) who estimated the unknown link function of varying index coefficient models (VICM) by the means of B-splines, as well as Kohler and Krzyżak (2017), and Schmidt-Hieber (2020). The latter two articles estimated nonparametric regression by deep neural network (DNN) methods, and have natural links to the generalized additive model.

this problem to 1, it is vulnerable to misspecification error. The economic theory might rule out any additional additive structure on the components of $f_1(\cdot)$ and $f_2(\cdot)$. JLL identified this model by imposing an additive structure on one of $f_1(\cdot)$ and $f_2(\cdot)$ as $f_1(x^1) = f_{11}(x_1^1) + f_{12}(x_2^1, \ldots, x_{d_1}^1)$ or $f_2(x^2) = f_{21}(x_1^2) + f_{22}(x_2^2, \ldots, x_{d_2}^2)$. The extra additive structure imposed by JLL is weaker than the one of Horowitz (2001), but their identification requires a large image/support condition (see Condition I2.(iv) of their Assumption I) which substantially restricts its applicability in real empirical applications. Their identification strategy also rules out discrete elements in x^1 and x^2 (see condition I1 of their Assumption I), and hence further restricts their applicability in real applications.² In contrast, we identify such a model without imposing any extra additive structure or any large image/support condition. Our identification approach also allows for discrete elements in x^1 and x^2 .

Our paper also contributes to the research line of the identification of model primitives by exploiting the monotonicity restrictions on nonparametric functions. One of our key identification steps exploits the monotonicity of the unknown link $G(\cdot)$ to transform the original model into a new model with univariate components. Our identification is hence established by this connection between our model with a flexible grouping and the transformed model with univariate components. The identification of latter has been well studied. The monotonicity of transformation function has been employed to identify the model primitives of different variants of transformation model by, e.g., Khan (2001), Chen (2002, 2010a,b, 2012), and Chen and Zhang (2020). Moreover, the monotonicity of nonparametric function on latent random variable has been used to identify the non-separable models by, e.g., Chesher (2003) and Matzkin (2003). In the auction literature, the monotonicity of bidding strategy helps to identify the value distribution by, e.g., Guerre, Perrigne, and Vuong (2000, 2009), Athey and Haile (2002), Li and Zheng (2009), Marmer and Shneyerov (2012), Gentry and Li (2014), and Li and Liu (2018). The monotonicity of strategies is also used to identify discrete games by, e.g., Tang (2010); De Paula and Tang (2012); Grieco (2014); Liu, Vuong, and Xu (2017), with a notable exception of Lewbel and Tang (2015). To test whether monotonicity restrictions hold, Hoderlein, Su, White, and Yang (2016) provided a testing procedure in the structural model without strategic interaction; while Liu and Vuong (2020) proposed nonparametric tests for monotonicity of strategies in the games of incomplete information.³

The rest of this paper is organized in the following way. Section 2 presents our generalized additive model with two component functions. It also lays out our strategy to identify the link function $G(\cdot)$ and the component functions $f_k(\cdot)$ for k = 1, 2. In Section 3, we propose a nonparametric estimation procedure closely following the identification strategy. Section 4 then establishes the large-sample properties of our estimators. In Section 5, a simulation is used to demonstrate the finite sample performance of our nonparametric estimators. Section 6 briefly discusses how to extend our framework to cases with discrete covariates and more than two components. The paper is concluded in Section 7. An appendix collects the proofs of our theorems. The online Supplemental Material (SM) collects some notations and technical lemmas (as well as their proofs).

²Note that discrete regressors are still not allowed to enter any component functions in their extension to handle discrete regressors (see Section 6 of Jacho-Chávez, Lewbel, and Linton, 2010).

³While maintaining the monotonicity restriction on bidding strategies, Liu and Luo (2017) proposed a nonparametric inference procedure to compare the valuation distributions in first price auctions.

2 The model and identification

We consider the generalized additive model with an unknown link function as follows,

$$H(x) = G\left(\sum_{k=1}^{K} f_k(x^k)\right),$$
(2)

where $H(\cdot)$ is a function which can be identified directly by the joint distribution of observables and can therefore be consistently estimated, such as the mean regression function $E(Y|X = \cdot)$ or the quantile regression function $Q_{Y|X}(\tau_0|\cdot)$ for a given τ_0 , and $x = (x^1, \dots, x^K)$ such that $x^k \in \mathbb{R}^{d_k}$ for $d_k \ge 1$. The parameter of interest includes the unknown link function $G(\cdot)$ and the component functions $f_k(\cdot)$ for $k = 1, \dots, K$.

For convenience of discussion, we focus on the case of two component functions in the link, i.e. the model is simplified as

$$H(x) = G(f_1(x^1) + f_2(x^2)),$$
(M)

where the unknown link function $G(\cdot)$ is monotonic, $x \equiv (x^1, x^2) \in \mathbb{R}^d$ and $x^k \in \mathbb{R}^{d_k}$ for k = 1, 2. We will return to the general case with more than two components in Section 6.2. Clearly, $d = d_1 + d_2$. Throughout the paper, let $X \equiv (X^1, X^2)$ be a random vector in \mathbb{R}^d with X^k denoting a random vector in \mathbb{R}^{d_k} , and $x \equiv (x^1, x^2)$ be its realized value with $x^k \in \mathbb{R}^{d_k}$, for k = 1, 2. In addition, let $p_V(\cdot)$ (or $p_{V^s|V^t}(\cdot|v^t)$) denote the probability density function of any given random vector/variable V (or the conditional density function of V^s given $V^t = v^t$).

In this paper, we aim to provide the identification and estimation of $G(\cdot)$ and $f_k(\cdot)$'s in such a model under reasonably weak restrictions motivated by empirical concerns. Specifically, we allow for a flexible division of (x^1, x^2) according to economic theory or practical experience,⁴ and discrete variables in x^1 and/or x^2 . The latter is motivated by the presence of discrete variables in many economic datasets. For presentation purpose, we first consider the case of (x^1, x^2) to only have continuous variables. We then return to the case with discrete variables in Section 6.1.

We obtain the nonparametric identification of (M) in three steps. In the first step, we transform it into a new generalized additive model with univariate components. The new model has the same link function as (M). In the second step, the transformed model is identified by a strategy adapted from Horowitz (2001). The original component functions are identified in the third step by applying the inverse of step-one transformation.

We first transform the original model (M) into a new model with univariate components. Such a transformation is given by the following theorem.

Theorem 1. Under a strictly monotonic link function $G(\cdot)$, the generalized additive model (M) can be transformed equivalently to

$$\mathcal{H}(z) = G(\tilde{f}_1(z^1) + \tilde{f}_2(z^2)), \tag{M'}$$

where $\mathcal{H}(z) = E[H(X)|\zeta_1(X^1) = z^1, \zeta_2(X^2) = z^2]$, the inverse of $\tilde{f}_k(\cdot)$ is $\tilde{f}_k^{-1}(s) = \int G(s + f_{-k}(x^{-k})) \cdot w_{-k}(x^{-k}) dx^{-k}$, and $\zeta_k(x^k) = \int H(x) \cdot w_{-k}(x^{-k}) dx^{-k}$ with freely chosen non-negative weight functions $w_k(\cdot)$ for k = 1, 2 where -k denotes the index other than k in $\{1, 2\}$.

⁴For example, let $(x^1, x^2) = (x_1, x_2, x_3, x_4)$. Our model allows all possible divisions, such as $x^1 = (x_1, x_2), x^2 = (x_3, x_4)$ or $x^1 = x_1, x^2 = (x_2, x_3, x_4)$.

Theorem 1 transforms the original model (M) into a new model (M') which is easier to analyze for two reasons. First, the new function $\mathcal{H}(\cdot)$ can be identified, since the function $\mathcal{H}(\cdot)$ and hence its weighted integrals $\zeta_k(\cdot)$'s are identified. Second, both of the new components $\tilde{f}_1(\cdot)$ and $\tilde{f}_2(\cdot)$ are univariate. Moreover, the functions $\tilde{f}_k(\cdot)$'s and their inverses are monotonic when the link $G(\cdot)$ is monotonic. To simplify the notation, hereafter let $Z = (Z^1, Z^2)$ with $Z^k = \zeta_k(X^k)$, and $z = (z^1, z^2)$ with $z^k \in \mathbb{R}$ for k = 1, 2.

Before proceeding with the identification of new model (M'), we give the identifying assumptions as follows,

Assumption I (Identification condition). (*i*) Location normalization: $\tilde{f}_1(z_0^1) = \tilde{f}_2(z_0^2) = 0$ for some interior point (z_0^1, z_0^2) in the support of Z;

(ii) Scale normalization: $\int w_3(z^1) / \tilde{f}'_1(z^1) dz^1 = 1$ where $w_3(\cdot)$ is some non-negative weight function by choice; (iii) Monotonicity: the link function $G(\cdot)$ is strictly monotonic.

Parts (i) and (ii) of Assumption I specify the location and scale normalizations needed for the identification. Similar normalizations have been adopted by the literature (see, e.g., Horowitz, 2001) to identify the generalized additive model. Note that our identification strategy still works (with minor change) if the location normalization is relaxed to $\tilde{f}_k(z_0^k) = \tilde{f}_{k0}$ with some known constant $\tilde{f}_{k0} \in \mathbb{R}$ for k = 1, 2. We can also adopt other location and scale normalizations, such as the ones of JLL. Part (iii) imposes a monotonicity condition on the link function $G(\cdot)$. Such a monotonicity condition is used to guarantee the existence of $\tilde{f}_k(\cdot)$'s and their inverses.

In the second step, we turn to identify the new model (M'). Such an identification is achieved in two stages by applying a strategy adapted from Horowitz (2001). In the first stage, we identify the transformed components $\tilde{f}_1(\cdot)$ and $\tilde{f}_2(\cdot)$. In the second stage, the unknown link $G(\cdot)$ is identified.

We now turn to the identification of transformed components $\tilde{f}_k(\cdot)$'s. Let $\mathcal{H}(z) = E[H(X)|Z = z]$, and $\partial_k g(z) = \partial g(z) / \partial z_k$ for any multivariate function g(z). The identification idea comes from the following two basic equations:

$$\partial_1 \mathcal{H}(z) = G'(\tilde{f}_1(z^1) + \tilde{f}_2(z^2)) \cdot \tilde{f}'_1(z^1)$$
(3)

$$\partial_2 \mathcal{H}(z) = G'(\tilde{f}_1(z^1) + \tilde{f}_2(z^2)) \cdot \tilde{f}'_2(z^2)$$
(4)

Let (4) be divided by (3), we obtain

$$\frac{\tilde{f}_2'(z^2)}{\tilde{f}_1'(z^1)} = \frac{\partial_2 \mathcal{H}(z)}{\partial_1 \mathcal{H}(z)}.$$
(5)

We next multiply both sides by $w_3(z^1)$ and integrate (i) by z^1 on the whole support of $Z^1 \equiv \zeta_1(X^1)$ and (ii) by z^2 from z_0^2 to z^2 , and get

$$\tilde{f}_2(z^2) = \int_{z_0^2}^{z^2} \tilde{f}_2'(z^2) dz^2 \cdot \int w_3(z^1) / \tilde{f}_1'(z^1) dz^1 = \int_{z_0^2}^{z^2} \int \frac{\partial_2 \mathcal{H}(z)}{\partial_1 \mathcal{H}(z)} \cdot w_3(z^1) dz^1 dz^2,$$
(C2)

where the first equality comes from the location and scale normalizations imposed by Assumption I. *The second transformed component function* $\tilde{f}_2(\cdot)$ *is hence identified by* (C2).

The identification of $\tilde{f}_1(\cdot)$ follows a similar strategy. Specifically, we apply the former strategy to

 $\tilde{f}_1'(z^1)/\tilde{f}_2'(z^2) = [\partial_1 \mathcal{H}(z)]/[\partial_2 \mathcal{H}(z)]$ and then have

(1)

$$\tilde{f}_{1}(z^{1}) = \left[\int_{z_{0}^{1}}^{z^{1}} \int \frac{\partial_{1} \mathcal{H}(z)}{\partial_{2} \mathcal{H}(z)} \cdot w_{4}(z^{2}) dz^{2} dz^{1} \right] / \left[\int w_{4}(z^{2}) / \tilde{f}_{2}'(z^{2}) dz^{2} \right],$$
(6)

which can identify the first transformed component $\tilde{f}_1(\cdot)$ if the denominator $\int w_4(z^2)/\tilde{f}_2'(z^2)dz^2$ can be identified. This is achieved by the scale normalization and (5) as

$$\frac{1}{\int w_4(z^2)/\tilde{f}_2'(z^2)dz^2} = \frac{\int \frac{w_3(z^2)}{\tilde{f}_1'(z^1)}dz^1}{\int \frac{w_4(z^2)}{\tilde{f}_2'(z^2)}dz^2} = \int \frac{w_3(z^1)}{\int \frac{\tilde{f}_1'(z^1)}{\tilde{f}_2'(z^2)} \cdot w_4(z^2)dz^2}dz^1 = \int \frac{w_3(z^1)}{\int \frac{\partial_1\mathcal{H}(z)}{\partial_2\mathcal{H}(z)} \cdot w_4(z^2)dz^2}dz^1,$$

which introduces an expression to identify the first transformed component $\tilde{f}_1(\cdot)$ as follows

$$\tilde{f}_1(z^1) = c \cdot \int_{z_0^1}^{z^1} \int \frac{\partial_1 \mathcal{H}(z)}{\partial_2 \mathcal{H}(z)} \cdot w_4(z^2) dz^2 dz^1, \tag{C1}$$

where $c = \int \omega_3(z^1) \cdot \left[\int \left[\partial_1 \mathcal{H}(z) / \partial_2 \mathcal{H}(z) \right] \cdot \omega_4(z^2) dz^2 \right]^{-1} dz^1$. Consequently, the first transformed component function $\tilde{f}_1(\cdot)$ is identified by (C1).

After identifying the transformed components $\tilde{f}_k(\cdot)$'s, we now investigate the identification of the unknown link $G(\cdot)$. The function $T(z) = \tilde{f}_1(z^1) + \tilde{f}_2(z^2)$ is identified once the transformed components $\tilde{f}_1(\cdot)$ and $\tilde{f}_2(\cdot)$ are identified. *The unknown link function* $G(\cdot)$ *is then identified* by the nonparametric regression of H(X) on T(Z), namely E[H(X)|T(Z)], due to the following result

$$E[H(X)|T(Z) = \tau] = E[\mathcal{H}(Z)|T(Z) = \tau] = G(\tau), \tag{L}$$

where the first equality comes from the fact that, given T(Z), the conditional expectation of H(X)and $\mathcal{H}(Z) = E[H(X)|Z]$ are the same by the law of iterated expectation; and the second equality holds due to the restriction given by (M'). In particular, when H(X) is a nonparametric regression E(Y|X), by the law of iterated expectation, the identification equation (L) for the link $G(\cdot)$ can be further simplified as

$$G(\tau) = E[Y|T(Z) = \tau]. \tag{L'}$$

In the final step, we use the inverse of step-one transformation to identify the original components $f_1(\cdot)$ and $f_2(\cdot)$. Notice that the original link $G(\cdot)$ has already been identified in step two. This is accomplished by the following mapping from the inverse of step-one transformation,

$$f_k(x^k) = \tilde{f}_k(\zeta_k(x^k)), \text{ for } k = 1, 2,$$
(7)

which can be derived by replacing *s* with $f_k(x^k)$ in the expressions of $\tilde{f}_k^{-1}(\cdot)$ of Theorem 1 and exploring the equality of (M). Both of the original component functions $f_k(\cdot)$ for k = 1, 2 are then identified, since $\zeta_k(\cdot)$'s are identified functions by their definitions in Theorem 1, and $\tilde{f}_k(\cdot)$'s have been identified in step two.

We summarize the above discussion on the identification of the link function $G(\cdot)$ and the original

component functions $f_k(\cdot)$'s in the following theorem whose proof is omitted.

Theorem 2. Let Assumption I hold. Given the expressions in (C1), (C2), and (L) are well defined, the link function $G(\cdot)$ is identified by (L), and the original component functions are identified by (7) where the transformed component functions $\tilde{f}_k(\cdot)$'s are given by (C1) and (C2) and the weighted integrals $\zeta_k(\cdot)$'s are defined by Theorem 1 for k = 1, 2. In particular, when H(x) = E(Y|X = x), the link function $G(\cdot)$ is identified by a simplified expression as (L').

Theorem 2 identifies the link $G(\cdot)$ and the original components $f_k(\cdot)$'s for k = 1, 2 by applying Horowitz (2001)'s strategy to the transformed model (M') in Theorem 1. In addition, Theorem 2 establishes the identification of model primitives when there are only two components within the link. Such an identification strategy can be easily extended to the case of more than two components (i.e. K > 2). We will briefly discuss such an extension in Section 6.2.

Remark 1: Theorem 2 shows that the link $G(\cdot)$ and components $f_k(\cdot)$'s are identified under each chosen set of weights $w_k(\cdot)$, k = 1, ..., 4. The choice of weights $w_k(\cdot)$ affect the efficiency of estimating $G(\cdot)$ and $f_k(\cdot)$'s.

3 Estimation

This section only considers estimating the parameter of interest in the case of nonparametric (mean) regression for $H(\cdot)$, i.e. H(x) = E(Y|X = x). We leave other cases of $H(\cdot)$ (such as the case of quantile regression) for future research. In the case of nonparametric regression, note that

$$E[Y|Z = z] = \mathcal{H}(z) = G(\tilde{f}_1(z^1) + \tilde{f}_2(z^2))$$
(8)

by the law of iterated expectation. That is, our estimation problem is essentially the same as Horowitz (2001)'s one, in which all component functions are univariate, if the true $\zeta_1(\cdot)$ and $\zeta_2(\cdot)$ were used in our estimation. Consequently, we propose a three-step estimation procedure by the kernel method to recover the parameter of interest, namely the link function $G(\cdot)$ and the component functions $f_k(\cdot)$ for k = 1, 2. We leave other nonparametric alternatives such as the sieve method proposed by, e.g., Ai and Chen (2003) and Chen (2007), for future research. In the first step, the nonparametric regression $H(\cdot)$ and its partial integrals $\zeta_k(\cdot)$'s are recovered by the local polynomial method, then the partial derivatives $\partial_k \mathcal{H}(\cdot)$ for k = 1, 2 are estimated by another local polynomial regression of Y on two generated regressors $\widehat{Z}^1 = \widehat{\zeta}_1(X^1)$ and $\widehat{Z}^2 = \widehat{\zeta}_2(X^2)$. In step two, the (transformed) components $\tilde{f}_k(\cdot)$'s are estimated through the expressions of (C1) and (C2) by replacing $\partial_k \mathcal{H}(\cdot)$ for k = 1, 2 with their step-one estimates, and the link $G(\cdot)$ is recovered through the local polynomial regression according to (L'). In the third step, the original components $f_k(\cdot)$ for k = 1, 2 are then recovered according to (7) by replacing $\tilde{f}_k(\cdot)$ and $\zeta_k(\cdot)$ with their nonparametric estimates.

Such a kernel estimation approach has several attractive features. First, the estimation strategy closely follows the identification idea laid out in Section 2. In particular, it transforms the estimation problem with multivariate components to the one with univariate components. The latter has been well studied in the literature. Second, it can group (X^1, X^2) in a flexible way. This flexibility can be important to adopt the generalized additive model in real empirical applications, since many applications may specify some or even all component functions to be multivariate. Third, we only need one continuous variable in X^1 and X^2 , namely, the other covariates in X^1 and X^2 can be all

discrete. For presentation purpose, we consider the case of (X^1, X^2) to only have continuous variables here. We will return to the case with discrete variables in Section 6.1.

Specifically, our estimation approach proceeds in three steps as follows.

Step 1. Estimation of $\partial_k \mathcal{H}(\cdot)$. We first use a local *r*th-order polynomial method to estimate $H(x) = E[Y|X^1 = x^1, X^2 = x^2]^5$. We use a leave-one-out estimator $\hat{H}_{-j}(x)$, namely, the intercept of

$$\widehat{\alpha} = \arg\min_{\alpha} \sum_{i \neq j} \left(Y_i - \sum_{0 \le |\mathbf{k}| \le r} \alpha_{\mathbf{k}} (X_i - x)^{\mathbf{k}} \right)^2 K\left(\frac{X_i - x}{h_H}\right),$$

where $\mathbf{k} = (k_1, k_2, \dots, k_d)$ is a *d*-tuple of integers, $|\mathbf{k}| = k_1 + k_2 + \dots + k_d$, $(X_i - x)^{\mathbf{k}} = (X_i^1 - x^1)^{k_1} \times (X_i^2 - x^2)^{k_2} \times \dots \times (X_i^d - x^d)^{k_d}$, and $K(x_1, \dots, x_d) = \prod_{\ell=1}^d k(x_\ell)$ with $k(\cdot)$ being a univariate kernel function (i.e. a multiplicative kernel is used in the multivariate case). More details of local polynomial regression could be found in the Appendix S.1. The generated regressors are estimated by $\hat{\zeta}_1(X_i^1) = (1/n) \cdot \sum_{j=1}^n \hat{H}_{-j}(X_i^1, X_j^2)$ and $\hat{\zeta}_2(X_i^2) = (1/n) \cdot \sum_{j=1}^n \hat{H}_{-j}(X_j^1, X_i^2)$ with the weights $w_k(\cdot)$ to be the marginal densities of X^k on \mathcal{S}_{X^k} , namely $w_k(\cdot) = p_{X^k}(\cdot)$, for k = 1, 2.

Finally, the partial derivatives $\partial_k \mathcal{H}(\cdot)$ can then be recovered by another local *r*th-order polynomial estimation, i.e. the slope coefficients of

$$\begin{split} \widehat{\beta} = & \arg\min_{\beta} \sum_{i=1}^{n} \left(Y_{i} - \sum_{0 \le k_{1} + k_{2} \le r} \beta_{k_{1},k_{2}} (\widehat{\zeta}_{1} \left(X_{i}^{1} \right) - z^{1})^{k_{1}} (\widehat{\zeta}_{2} \left(X_{i}^{2} \right) - z^{2})^{k_{2}} \right)^{2} \\ & \quad \cdot k \Big(\frac{\widehat{\zeta}_{1} \left(X_{i}^{1} \right) - z^{1}}{h_{\mathcal{H}}} \Big) k \Big(\frac{\widehat{\zeta}_{2} \left(X_{i}^{2} \right) - z^{2}}{h_{\mathcal{H}}} \Big). \end{split}$$

Denote the derivative estimators by $\partial_k \hat{\mathcal{H}}(z)$ for k = 1, 2.

Step 2. Estimation of the transformed model. The transformed component functions $\tilde{f}_k(\cdot)$'s are estimated by the sample analogue of (C1) and (C2) as follows:

$$\widehat{f}_1(z^1) = \widehat{c} \int_{z_0^1}^{z^1} \int \frac{\partial_1 \widehat{\mathcal{H}}(z)}{\partial_2 \widehat{\mathcal{H}}(z)} \omega_4(z^2) dz^2 dz^1, \quad \widehat{f}_2(z^2) = \int_{z_0^2}^{z^2} \int \frac{\partial_2 \widehat{\mathcal{H}}(z)}{\partial_1 \widehat{\mathcal{H}}(z)} \omega_3(z^1) dz^1 dz^2,$$

where $\hat{c} = \int \omega_3(z^1) \left[\int \left[\partial_1 \hat{\mathcal{H}}(z) / \partial_2 \hat{\mathcal{H}}(z) \right] \cdot \omega_4(z^2) dz^2 \right]^{-1} dz^1.$

The link function $G(\cdot)$ is then estimated by the intercept of

$$\widehat{\gamma} = \arg\min_{\gamma} \sum_{i=1}^{n} \left(Y_i - \sum_{0 \le k \le r} \gamma_k (\widehat{f}_1(\widehat{Z}_i^1) + \widehat{f}_2(\widehat{Z}_i^2) - \tau)^k \right)^2 k \left(\frac{\widehat{f}_1(\widehat{Z}_i^1) + \widehat{f}_2(\widehat{Z}_i^2) - \tau}{h_G} \right),$$

where $\widehat{Z}_i^k = \widehat{\zeta}_k(X_i^k)$ for k = 1, 2.

Step 3. Estimation of the original component functions $f_k(\cdot)$ **'s.** Lastly, the original component functions $f_1(\cdot)$ and $f_2(\cdot)$ are estimated by

$$\widehat{f}_k(x^k) = \widehat{\widetilde{f}}_k(\widehat{\zeta}_k(x^k)), \text{ for } k = 1, 2.$$

Three remarks are in order. First, our estimators essentially have similar asymptotic properties to

⁵Here, *r* is also the smoothness of unknown functions and densities. See Assumption A.3.

the Horowitz (2001)'s estimators if the true partial integrals $\zeta_k(\cdot)$'s were used so that the first step is not needed. Second, we use local polynomial regressions instead of local constant ones to address the boundary bias issue (see also Fan and Gijbels (1992)). Third, the step-two estimation of the link $G(\cdot)$ can be viewed as a result of estimating it by a sample analogue of (L'). It can also be viewed as a result of recovering $G(\cdot)$ by a sample analogue of a moment condition of $G(\tau) = E[Y|f_1(X^1) + f_2(X^2) = \tau]$ which comes from (M) and the law of iterated expectation.

4 Large Sample Properties

In this section, we study the large sample properties of the estimators proposed in Section 3. Let $d_1 \ge d_2$ only for presentation purposes.⁶ We first state the assumptions under which the large sample properties of our estimators are established. Let $int(\Theta)$ denote the interior of any given set Θ . Let S_W be the support of a random vector/variable W, and S_G be defined as $\{\tau : \tau = f_1(x^1) + f_2(x^2) \text{ for some } (x^1, x^2) \in S_{(X^1, X^2)} \}$.

Assumption A.1 (DGP). $\{(Y_i, X_i)\}_{i=1}^n$ is an *i.i.d.* sample from the distribution of (Y, X) which satisfies (M) and (i) $E(|Y|^{4+s}|X = x) \leq C$ for some finite C, positive s, and all $x \in S_X$; (ii) Var(Y|X = x) is continuous in x.

Assumption A.2 (distribution of X). The random vector X satisfies (i) S_X is compact; (ii) the distribution of X is absolutely continuous with respect to Lebesgue measure and has density of $p_X(\cdot) > 0$ in the interior of S_X ; (iii) there exist some compact intervals $\mathcal{I}_1 \subset int(S_{Z^1})$, $\mathcal{I}_2 \subset int(S_{Z^2})$ and some $\underline{c} > 0$ such that (a) $\tilde{f}'_k(z^k) \ge \underline{c}$ for all $z^k \in \mathcal{I}_k$ and k = 1, 2, (b) $P(X : Z^k \in \mathcal{I}_k, k = 1, 2) > 0$, (c) $z_0^k \in \mathcal{I}_k$ where z_0^k is defined in Assumption I for k = 1, 2, (d) $|G'(\cdot)| \ge \underline{c}$ on S_G .

Assumption A.3 (smoothness of *G*, f_k and p_X). (*i*) The link function $G(\cdot)$ is (r + 1) times continuously differentiable. (*ii*) The component functions $f_k(\cdot)$ for k = 1, 2 and density $p_X(\cdot)$ are (r + 1) times differentiable with respect to any mixture of its arguments with uniformly continuous derivatives on their supports S_{X^k} and S_X .

Assumption A.4 (weights). (i) For k = 1, 2, the weight function $w_k(\cdot) = p_{X^k}(\cdot)$. (ii) For k = 3, 4, the weight function $w_k(\cdot)$ is non-negative and bounded with support $S_{w_k} \subset \mathcal{I}_{k-2}$ such that $w_k(\cdot)$ has (r+1)-th continuous derivatives on S_{w_k} with $\int w_k(z^{k-2})dz^{k-2} = 1$.

Assumption A.5 (kernel). The univariate kernel function $k(\cdot)$ is symmetric, bounded, and continuously differentiable on its support [-1,1] For any $d' \ge 1$ and a kernel function $K(\cdot)$ on $[-1,1]^{d'}$, there is $K(s_1,\ldots,s_{d'}) = \prod_{j=1}^{d'} k(s_j)$. Let $H_{\mathbf{j}}(u) = u^{\mathbf{j}}K(u)$ for all integers $\mathbf{j} = (j_1, j_2, \cdots, j_d)$ and $u \in \mathbb{R}^d$. Then $H_{\mathbf{j}}(u)$ is Lipschitz continuous on $[-1,1]^d$ for all \mathbf{j} with $0 \le |\mathbf{j}| \le 2r + 1$.

Assumption A.6 (bandwidth). As $n \to \infty$, the bandwidth sequences h_H , h_H , and h_G go to zero and satisfy:

$$\begin{array}{ll} (i) \ nh_{H}^{d+r+1}/log(n) \to \infty, & nh_{\mathcal{H}}^{6}/log(n) \to \infty, & nh_{G}^{3}/log(n) \to \infty, \\ (ii) \ h_{H}^{d_{2}}/h_{\mathcal{H}} \to 0, & log(n)^{2}/[nh_{H}^{d_{1}/2+r+1}h_{\mathcal{H}}^{3}] \to \gamma_{1}, & n \cdot h_{H}^{d_{2}} \cdot h_{\mathcal{H}}^{2r} \to \tilde{\gamma}_{1}, \\ (iii) \ h_{H}^{r+1}/h_{G} \to 0, & n \cdot h_{H}^{d_{1}} \cdot h_{G}^{2} \to \infty, & nh_{G}^{2r+3} \to \gamma_{2}, & n \cdot h_{H}^{2r+2} \cdot h_{G} \to \gamma_{3}, & nh_{H}^{d_{2}+2r+2} \to \tilde{\gamma}_{2}, \end{array}$$

⁶If $d_1 < d_2$, we can define $\overline{x}^1 = x^2$ and $\overline{x}^2 = x^1$. It then follows that $\overline{d}_1 > \overline{d}_2$ where \overline{d}_k denotes the dimension of \overline{x}^k for k = 1, 2. We then study the new model of $\overline{H}(\overline{x}^1, \overline{x}^2) = G(\overline{f}_1(\overline{x}^1) + \overline{f}_2(\overline{x}^2))$ where $\overline{H}(\overline{x}^1, \overline{x}^2) = H(x^1, x^2)$, $\overline{f}_1(\overline{x}^1) = f_2(x^2)$, and $\overline{f}_2(\overline{x}^2) = f_1(x^1)$.

(*iv*) $h_{\mathcal{H}}^r/h_G \to 0$, $h_G/h_{\mathcal{H}} \to \delta_G$, $n \cdot h_{\mathcal{H}}^{2r} \cdot h_G \to \gamma_4$,

where $\gamma_1, \ldots, \gamma_4, \tilde{\gamma}_1, \tilde{\gamma}_2$, and δ_G are some non-negative constants.

Assumption A.1 describes the model and Data Generating Process (DGP). Assumption A.2 (i) and (ii) give some regularity conditions on the support and density function of the random vector X. With the normalization conditions in Assumption I, Assumption A.2 (iii) provides sufficient conditions to identify the component functions $f_k(\cdot)$'s and the link function $G(\cdot)$.

Assumption A.3 contains some smoothness conditions on the link function $G(\cdot)$, the component functions $f_k(\cdot)$'s, and the density function $p_X(\cdot)$. They require those functions having a smoothness of (r + 1) in order to make our Taylor-series expansions to proper orders. In addition, they imply that the transformed component functions $\tilde{f}_k(\cdot)$'s also have (r + 1) derivatives which are uniformly continuous on their supports.

Assumption A.4 describes the condition on the weight functions $w_k(\cdot)$ for k = 1, ..., 4. For k = 1, 2, it uses the marginal density of X^k on S_{X^k} as the weight $w_k(\cdot)$ to estimate the partial integrations $\zeta_k(\cdot)$ in step one of our estimation approach laid out in Section 3. Other weights for $w_1(\cdot)$ and $w_2(\cdot)$ can also be used. For k = 3, 4, it requires the weight function $w_k(\cdot)$ to be (r + 1) times continuously differentiable on its support.

Assumption A.5 gives the restrictions on the univariate kernel function $k(\cdot)$ which builds all multivariate kernel functions throughout this paper in a multiplicative way. This assumption is also used in other local polynomial literature. See, e.g., Kong, Linton, and Xia (2010) and JLL. This assumption is utilized to derive the uniform asymptotic representation of local polynomial estimators.

Assumption A.6 specifies the conditions on the choices of bandwidths used in our kernel estimation. These conditions permit various combinations of bandwidths h_H , h_H , and h_G . For example, they are satisfied when $h_H \in (n^{-1/(r+1+d)}, n^{-(r+1)/[r \cdot (2r+3)]})$, and $h_H = h_G = n^{-(r+1)/[r \cdot (2r+3)]}$ for large enough r. They ensure that the remainder terms are negligible in each stage of our estimation. In particular, conditions (ii)-(iv) control the contributions from the previous estimation steps to the asymptotic variances of $\hat{f}_k(\cdot)$ and $\hat{G}(\cdot)$ for k = 1, 2.

We now present the asymptotic results of our estimators of the component functions $f_k(\cdot)$ and the link function $G(\cdot)$ for k = 1, 2. We first consider the estimation of original component functions $f_k(\cdot)$ for k = 1, 2. Our third theorem gives the asymptotic properties of the estimators $\hat{f}_k(\cdot)$ for k = 1, 2.

Theorem 3. Suppose that Assumptions I, A.1-A.6 hold. Then, for every k = 1, 2, as $n \to \infty$: (i) $\sup_{x^k \in S_{X^k}} |\widehat{f}_k(x^k) - f_k(x^k)| \to 0$ in probability, and (ii) for any $x^k \in S_{X^k}$, $\sqrt{nh_H^{d_k}} (\widehat{f}_k(x^k) - f_k(x^k) - B_{nf_k}(x^k)) \xrightarrow{d} N(0, \sigma_k^2(x^k))$ where $B_{nf_k}(x^k)$ and $\sigma_k^2(x^k)$ are given by (S.1.1) and (S.1.2), respectively.

Theorem 3 establishes the uniform consistency and asymptotic normality of our estimators of original component functions $f_k(\cdot)$ for k = 1, 2. It shows that the only contributions from previous estimation steps are in the resulting biases of $\hat{f}_k(\cdot)$ for k = 1, 2 in the final step. The variances of previous steps do not contribute into the variances of $\hat{f}_k(\cdot)$, namely the asymptotic variances of $\hat{f}_k(\cdot)$ do not enter the ones of $\hat{f}_k(\cdot)$. In particular, since the estimator can be represented as $\hat{f}_k(\cdot) = \hat{f}_k(\hat{\zeta}_k(\cdot))$, the asymptotic bias term $B_{nf_k}(x^k)$ consists of two parts. The first part $h_{\mathcal{H}}^r \mathfrak{B}_k(\zeta_k(x^k))$ is the bias of the infeasible estimator $\check{f}_k(\cdot)$ of $f_k(\cdot)$ if the (unobserved) true $\zeta_k(\cdot)$'s were used in all three steps. Specifically, the infeasible estimator $\hat{f}_k(\cdot)$ of the transformed component function $\tilde{f}_k(\cdot)$ is obtained by using the true $\zeta_k(\cdot)$'s, instead of their estimators $\hat{\zeta}_k(\cdot)$'s, to recover $\partial_k \mathcal{H}(\cdot)$ in the first step. The second part

 $h_{H}^{r} \cdot [\tilde{f}_{k}^{\prime}(\zeta_{k}(x^{k}))D_{k}(x^{k}) + \tilde{\mathcal{B}}_{k}(\zeta_{k}(x^{k}))]$ is the additional bias brought by using the estimators $\hat{\zeta}_{k}(\cdot)$'s, instead of the true functions $\zeta_{k}(\cdot)$'s, in all three steps.

Two additional remarks are in order. First, the asymptotic bias terms $B_{nf_k}(x^k)$ for k = 1, 2 is controllable in general when we use bandwidths satisfying Assumption A.6, i.e. $\limsup_{n\to\infty} \sqrt{nh_H^{d_k}}B_{nf_k}(x^k) < \infty$ holds. Second, there are two ways to consistently estimate the asymptotic variances $\sigma_k^2(x^k)$ for k = 1, 2. The first way exploits the expression of $\sigma_k^2(x^k)$ and replaces its population terms with their nonparametric consistent estimators. The other way is to estimate $\sigma_k^2(x^k)$ by adapting the bootstrap method for nonparametric regression. See, e.g., Härdle and Bowman (1988); Hall (1992); Hall and Horowitz (2013), among others.

We next consider the estimation of link function $G(\cdot)$. Our next theorem summarizes the large sample properties of our link estimator $\widehat{G}(\cdot)$. Let S_G be the compact set $\{\tau : \tau = f_1(x^1) + f_2(x^2) \text{ for some } (x^1, x^2) \in S_X\}$ where S_X is the support of X.

Theorem 4. Let Assumptions I, A.1-A.6 hold. Then as $n \to \infty$: (i) $\sup_{\tau \in S_G} |\widehat{G}(\tau) - G(\tau)| \to 0$ in probability, and (ii) for any $\tau \in S_G$, $\sqrt{nh_G} \cdot (\widehat{G}(\tau) - G(\tau) - B_{nG}(\tau)) \xrightarrow{d} N(0, \sigma_G^2(\tau))$ where $B_{nG}(\tau)$ and $\sigma_G^2(\tau)$ are defined by (S.1.3) and (S.1.4), respectively.

Theorem 4 shows the uniform convergence and asymptotic normality of our kernel estimator of link $G(\cdot)$. Several remarks are in order. First, the asymptotic bias $B_{nG}(\tau)$ consists of three terms. The first term $h_G^{r+1}e'_{1G}\{S_r^G\}^{-1}S_r^{G,r+1}G_{r+1}(\tau)^7$ comes from the infeasible estimation of link $G(\cdot)$ when the (unobserved) true $\zeta_k(\cdot)$ and $\tilde{f}_k(\cdot)$ for k = 1, 2, instead of their estimators, were used in the second step to recover $G(\cdot)$. It is a bias term of a standard nonparametric regression. The other two terms are the additional biases caused by using the feasible estimators $\hat{\zeta}_k(\cdot)$ and $\hat{f}_k(\cdot)$ for k = 1, 2, instead of their true functions, in the second step to estimate $G(\cdot)$. Second, similar to the case of $\hat{f}_k(\cdot)$'s, the asymptotic bias is controllable under Assumption A.6. Third, our asymptotic variance $\sigma_G^2(\tau)$ can be estimated through replacing its population quantities with their consistent estimators.

5 A Simulation Study

This section demonstrates the finite sample performance of our estimator by some Monte Carlo experiments. We adopt the following data generating process with the sample sizes of 400 and 800, each replicated 200 times:

$$Y = 1\{f_1(X^1) + f_2(X_1^2, X_2^2) - U > 0\},\$$

where the regressors X^1 , X_1^2 , and X_2^2 are independent truncated normal on [-3, 3] with mean 0 and standard deviation of 2, and the error term U is independent of all regressors and distributed according to standard normal N(0, 1). The true link and component functions are specified as

$$G(\tau) = \Phi(\tau), \quad f_1(x^1) = x^1, \quad f_2(x_1^2, x_2^2) = x_1^2 \cdot x_2^2,$$

where $\Phi(\cdot)$ is the distribution function of standard normal.

Two remarks are in order. First, under this specification, the partial integrals are $\zeta_1(x^1) = E[\Phi(x^1 + X_1^2 \cdot X_2^2)]$ and $\zeta_2(x_1^2, x_2^2) = E[\Phi(X^1 + x_1^2 \cdot x_2^2)]$, and the transformed components are the correspondent

⁷See (S.1.3) in the Appendix S.1.

inverse functions with $\tilde{f}_1(\zeta_1(x^1)) = x^1$ and $\tilde{f}_2(\zeta_2(x_1^2, x_2^2)) = x_1^2 \cdot x_2^2$. Second, the location normalization then requires $z_0^1 = \zeta_1(0)$ and $z_0^2 = \zeta_2(0,0)$ since $\tilde{f}_1(\zeta_1(0)) = 0$ and $\tilde{f}_2(\zeta_2(0,0)) = 0$. The symmetry of distributions of X^1 and X^2 implies that $\zeta_1(0) = \zeta_2(0,0) = \Phi(0) = 0.5$ which is used in the simulation. The scale normalization holds in the model with a constant weight function $w_3(z_1) = (\int_{0.3}^{0.7} [\tilde{f}'_1(z_1)]^{-1} dz_1)^{-1} \cdot 1\{0.3 \le z_1 \le 0.7\}.$

We next provide the implementation details of our estimation method. Let $\hat{\sigma}(W)$ denote the standard error of a given random variable W. To estimate $f_1(\cdot)$, $f_2(\cdot, \cdot)$, and $G(\cdot)$, we use local linear regressions with a second-order Gaussian kernel and the bandwidths of $h_H = \min \{\hat{\sigma}(X^1), \hat{\sigma}(X_1^2), \hat{\sigma}(X_2^2)\} \cdot n^{-1/7}$, $h_{\mathcal{H}} = \min \{\hat{\sigma}(\hat{Z}^1), \hat{\sigma}(\hat{Z}^2)\} \cdot n^{-1/8}$, and $h_G = \hat{\sigma}(\hat{f}_1(X^1) + \hat{f}_2(X_1^2, X_2^2)) \cdot n^{-1/5}$ following the simplified Silverman's rule of thumb (Silverman, 1986; Hansen, 2009). The weight function $w_4(\cdot)$ is chosen according to $w_4(z^2) = \frac{5}{3} \cdot 1\{0.2 \le z^2 \le 0.8\}$. Meanwhile, we replicate the estimators of JLL (a.k.a. "JLL estimators" in our paper) to do a side-by-side comparison. The details are given as follows. In the estimation of JLL, we also choose the second-order Gaussian kernel to do local linear regressions in all stages, use linear extrapolation to extend the integrand function when we do numerical integration and apply the silverman's rule of thumb to pick the bandwidths. To compute the integrals in our and JLL's estimators, we adopt the midpoint rule to calculate them numerically.

We now show the performance of our estimators and JLL estimators of $f_1(\cdot)$, $f_2(\cdot, \cdot)$ and $G(\cdot)$ to demonstrate how well our estimation procedure can recover the component and link functions at different locations. In particular, we report the bias (Bias), the standard deviation (SD), and the root mean square error (RMSE) for all estimators. Table 1 summarizes the simulation results for the estimation of components $f_1(\cdot)$, Table 2 is for $f_2(\cdot, \cdot)$ and Table 3 for the estimation of link $G(\cdot)$. We report in Tables 1-3 the simulation results for ours and JLL estimators at different points in the interior of the support of each function, where the left sections display the results for our estimators and right sections for JLL.

			ours				JLL			
п	x^1	Bias	SD	RMSE		Bias	SD	RMSE		
	-1	0.109	0.162	0.194		0.141	0.249	0.286		
400	0	-0.005	0.121	0.121		0.018	0.246	0.246		
	1	-0.104	0.161	0.191		-0.148	0.272	0.309		
	-1	0.095	0.115	0.149		0.129	0.220	0.254		
800	0	0.002	0.101	0.100		-0.007	0.165	0.164		
	1	-0.089	0.126	0.154		-0.122	0.185	0.221		

Table 1: Simulation results for the estimation of component function $f_1(x^1)$

Tables 1-2 show the estimation of components $f_1(\cdot)$ and $f_2(\cdot, \cdot)$, respectively. Table 1 shows the performance of our component estimator $\hat{f}_1(x^1)$ for $x^1 = -1, 0, 1$. They show that our estimator $\hat{f}_1(\cdot)$ performs reasonably well even under the moderate sample size of 400. When the sample size increases from 400 to 800, the RMSEs of $\hat{f}_1(\cdot)$ decline significantly. Moreover, the estimation biases are relatively small under both sample sizes of 400 and 800. Table 2 reports the estimation results for $f_2(x_1^2, x_2^2)$ for all $x_1^2 = -1, 0, 1$ and $x_2^2 = -1, 1$. We first look at the case of $x_2^2 = -1$ which is shown in the upper sections of table 2. The biases are relatively small under both n = 400 and n = 800. In addition, the

			ours				JLL				
п	x_{1}^{2}	x_{2}^{2}	Bias	SD	RMSE		Bias	SD	RMSE		
	-1	-1	-0.115	0.237	0.263		-0.107	0.346	0.361		
400	0	-1	-0.007	0.199	0.199		0.017	0.253	0.253		
	1	-1	0.083	0.239	0.252		0.084	0.315	0.326		
	-1	1	0.065	0.251	0.258		0.095	0.288	0.303		
	0	1	-0.008	0.201	0.201		-0.006	0.248	0.248		
	1	1	-0.078	0.242	0.254		-0.097	0.301	0.315		
	-1	-1	-0.103	0.168 0.197			-0.061	0.257	0.263		
800	0	-1	-0.005	0.151	51 0.151		-0.028	0.191	0.193		
	1	-1	0.086	0.182	0.201		0.030	0.245	0.247		
	-1	1	0.069	0.182	0.194		0.037	0.266	0.268		
	0	1	-0.022	0.152	0.153		0.023	0.210	0.210		
	1	1	-0.095	0.176	0.200		-0.056	0.267	0.273		

Table 2: Simulation results for the estimation of component function $f_2(x_1^2, x_2^2)$

Table 3: Simulation results for the estimation of link function $G(\tau)$

		ours			JLL			
п	τ	Bias	SD	RMSE	Bias	SD	RMSE	
	-3	0.000	0.006	0.006	0.013	0.039	0.041	
400	0	0.000	0.062	0.061	0.002	0.049	0.049	
	3	0.001	0.010	0.010	-0.016	0.045	0.047	
	-3	0.001	0.005	0.005	0.010	0.027	0.028	
800	0	0.002	0.043	0.042	0.001	0.041	0.041	
	3	0.000	0.004	0.004	-0.007	0.026	0.026	

decrease of RMSEs is significant when the sample size increases from 400 to 800. Our estimation of the two-dimensional function $f_2(\cdot, \cdot)$ hence performs reasonably well. We then look at the case of $x_2^2 = 1$ shown in the lower sections of table 2. Similar to the case of $x_2^2 = -1$, it confirms that (i) the biases are relatively satisfactory under both sample sizes of 400 and 800; (ii) our estimator becomes closer to its true value as the sample size increases.

Table 3 gives the performance of our link estimator $\hat{G}(\tau)$ for $\tau = -3, 0, 3$. In general our link estimator performs relatively well, although it is given by a nonparametric regression with a regressor generated by a two-step nonparametric estimation. The biases are reasonably small under n = 400 and n = 800. In addition, the RMSEs decrease when the sample size increases from 400 to 800.

Tables 1 - 3 also compare our results with JLL estimators. We can see that (i) our estimators have smaller variances and RMSEs than JLL estimators with two exceptions in the estimation of $f_1(\cdot)$, $f_2(\cdot)$ and $G(\cdot)$ and (ii) our estimators have biases in a magnitude similar to JLL. Thus, our estimators perform well in finite sample even though we do not require the existence of an univariate component like JLL.

				ours			JLL			Pinkse (2001)		
x^1	x_{1}^{2}	x_{2}^{2}	Bias	SD	RMSE	Bias	SD	RMSE	Bias	SD	RMSE	
	n = 400											
-1	-1	-1	-0.004	0.095	0.095	0.012	0.109	0.109	-0.012	0.105	0.106	
0	-1	-1	-0.046	0.076	0.089	-0.081	0.102	0.130	-0.075	0.071	0.103	
1	-1	-1	-0.024	0.036	0.043	-0.073	0.082	0.110	-0.049	0.034	0.060	
-1	0	-1	0.039	0.073	0.083	0.101	0.095	0.138	0.052	0.062	0.081	
0	0	-1	-0.004	0.084	0.084	0.014	0.106	0.106	0.003	0.079	0.079	
1	0	-1	-0.045	0.070	0.083	-0.092	0.100	0.135	-0.057	0.061	0.083	
-1	1	-1	0.020	0.039	0.044	0.068	0.065	0.094	0.050	0.035	0.061	
0	1	-1	0.035	0.075	0.082	0.093	0.097	0.134	0.078	0.073	0.107	
1	1	-1	-0.010	0.093	0.094	-0.010	0.108	0.108	0.008	0.107	0.107	
-1	-1	1	0.019	0.038	0.042	0.068	0.063	0.093	0.050	0.033	0.060	
0	-1	1	0.031	0.072	0.078	0.095	0.088	0.129	0.078	0.075	0.109	
1	-1	1	-0.016	0.098	0.099	-0.007	0.092	0.092	0.007	0.100	0.100	
-1	0	1	0.040	0.072	0.082	0.096	0.089	0.131	0.054	0.060	0.081	
0	0	1	-0.004	0.087	0.087	0.008	0.096	0.096	0.004	0.078	0.078	
1	0	1	-0.046	0.074	0.087	-0.095	0.095	0.134	-0.055	0.058	0.080	
-1	1	1	0.011	0.092	0.093	0.011	0.101	0.101	-0.016	0.115	0.116	
0	1	1	-0.036	0.079	0.086	-0.082	0.092	0.122	-0.078	0.075	0.108	
1	1	1	-0.021	0.038	0.043	-0.075	0.079	0.108	0.051	0.036	0.062	
						n = 800)					
-1	-1	-1	0.001	0.070	0.070	0.020	0.092	0.094	0.000	0.088	0.088	
0	-1	-1	-0.034	0.059	0.068	-0.066	0.081	0.105	-0.065	0.058	0.088	
1	-1	-1	-0.018	0.032	0.037	-0.048	0.048	0.048	0.042	0.026	0.049	
-1	0	-1	0.036	0.060	0.069	0.072	0.085	0.111	0.056	0.049	0.074	
0	0	-1	-0.001	0.065	0.064	-0.008	0.090	0.090	0.003	0.065	0.065	
1	0	-1	-0.033	0.058	0.066	-0.083	0.075	0.112	-0.048	0.048	0.068	
-1	1	-1	0.018	0.032	0.037	0.050	0.066	0.083	0.044	0.026	0.052	
0	1	-1	0.034	0.059	0.068	0.057	0.078	0.096	0.068	0.060	0.091	
1	1	-1	0.001	0.075	0.074	-0.026	0.082	0.085	0.005	0.081	0.081	
-1	-1	1	0.016	0.028	0.032	0.051	0.064	0.082	0.046	0.028	0.054	
0	-1	1	0.031	0.061	0.069	0.061	0.078	0.099	0.071	0.061	0.093	
1	-1	1	-0.005	0.077	0.077	-0.021	0.086	0.078	0.008	0.085	0.085	
-1	0	1	0.031	0.055	0.063	0.085	0.090	0.123	0.057	0.045	0.072	
0	0	1	-0.007	0.068	0.068	0.006	0.089	0.089	0.005	0.066	0.066	
1	0	1	-0.038	0.059	0.070	-0.072	0.073	0.103	-0.047	0.049	0.067	
-1	1	1	0.003	0.073	0.072	0.025	0.098	0.101	-0.004	0.088	0.088	
0	1	1	-0.033	0.061	0.069	-0.064	0.084	0.106	-0.068	0.059	0.090	
1	1	1	-0.018	0.032	0.032	-0.049	0.052	0.072	-0.043	0.025	0.050	

Table 4: Simulation results for the estimation of original regression function $H(x^1, x_1^2, x_2^2)$

We also compare our method with Pinkse (2001) in finite sample before concluding our simulation section. In the context of Pinkse (2001), the above model can be represented as

$$H(x^1, x^2) = \tilde{G}(x^1, \tilde{f}_2(x^2))$$

where $H(x^1, x_1^2, x_2^2) = \Phi(x^1 + x_1^2 x_2^2)$, $\tilde{f}_2(x^2) = x_1^2 x_2^2$ and $\tilde{G}(x^1, t) = \Phi(x^1 + t)$. The estimation of $H(\cdot)$ is the object of comparison here. To implement his approach, we apply local linear method for the first-step estimation and a weighted local constant regression for the second-step estimation. The second-step estimation closely follows the definition of his estimator. We also use the second order Gaussian kernel and bandwidths following the rule of thumb. The weight is chosen according to the simulation study of Pinkse (2001). We report the simulation results of his third estimator, namely S_{π} , here.

Table 4 shows the simulation results for ours, JLL and Pinkse's estimators of the overall function $H(x^1, x_1^2, x_2^2)$ under sample sizes 400 and 800. The comparison shows that (i) the RMSEs of our estimator decline significantly when the sample size increases from 400 to 800; (ii) our estimator has smaller variances and RMSEs than both of JLL and Pinkse (2001)'s estimators in most cases; (iii) our biases are comparable to the best ones between JLL and Pinkse (2001).

6 Extensions

6.1 Discrete covariates

We now turn to the case with discrete covariates in (x^1, x^2) . Let $X^k = (X^k_d, X^k_c)$ with discrete regressors $X^k_d \in \mathbb{R}^{a_k}$ ($a_k \ge 1$) and continuous regressors $X^k_c \in \mathbb{R}^{b_k}$ ($b_k \ge 1$) for k = 1, 2.

With mixed data of discrete and continuous regressors, our transformation and identification results, namely Theorems 1 and 2, still hold under proper choices of weight functions $w_k(\cdot)$ for k = 1, ..., 4 and proper definition of integration with respect to discrete variables. We follow Li and Racine (2007) to accommodate both discrete and continuous regressors in our estimation. We mainly need to modify the kernel regression estimators of $\hat{\zeta}_k(x^k)$ for k = 1, 2 and $\hat{H}_{-j}(x)$ in Step 1 (outlined in Section 3) as follows:

$$\widehat{\zeta}_1(x^1) = \frac{1}{n} \sum_{j=1}^n \widehat{H}_{-j}(x^1, X_j^2), \quad \widehat{\zeta}_2(x^2) = \frac{1}{n} \sum_{j=1}^n \widehat{H}_{-j}(X_j^1, x^2),$$

where $\hat{H}_{-i}(x)$ comes from the intercept of a leave-one-out local polynomial estimation

$$\min_{\alpha} \sum_{i \neq j} \left(Y_i - \sum_{0 \leq |\mathbf{k}| \leq r} \alpha_{\mathbf{k}} (X_{ci} - x_c)^{\mathbf{k}} \right)^2 K_{h_H, \lambda}(x, X_i),$$

where $\mathbf{k} = (k_1, k_2, \dots, k_{b_1+b_2})$, $K_{h_H,\lambda}(x, X) = \prod_{\ell=1}^{b_1} k \left(\frac{x_{\ell\ell}^1 - X_{\ell\ell}^1}{h_H} \right) \cdot \prod_{\ell=1}^{b_2} k \left(\frac{x_{\ell\ell}^2 - X_{\ell\ell}^2}{h_H} \right) \cdot \prod_{\ell=1}^{a_1} \lambda_{1\ell}^{N_{\ell}^1(x,X)} \cdot \prod_{\ell=1}^{a_2} \lambda_{2\ell}^{N_{\ell}^2(x,X)}$ and $N_{\ell}^k(x, X) = 1\{X_{d\ell}^k \neq x_{d\ell}^k\}$. Here, we use a multiplicative kernel function for the multivariate regressors. A univariate kernel of $k(\cdot)$ is adopted for the continuous regressors, and another univariate kernel of $l(X_{d\ell}^k, x_{d\ell}^k, \lambda_{k\ell}) = \lambda_{k\ell}^{1\{X_{d\ell}^k \neq x_{d\ell}^k\}}$ is employed for the discrete (and unordered) regressors with a bandwidth $\lambda_{k\ell} \in [0, 1]$.⁸

⁸If the discrete regressors are ordered, then a univariate kernel of $l(X_{d\ell}^k, x_{d\ell}^k, \lambda_{k\ell}) = \lambda_{k\ell}^{|X_{d\ell}^k - x_{d\ell}^k|}$ can be applied in this case. See

With the above adaption in our estimation (to accommodate the mixed data of discrete and continuous regressors), we can obtain the large sample properties of $\hat{f}_k(\cdot)$ for k = 1, 2 and $\hat{G}(\cdot)$ similar to those summarized by Theorems 3 and 4. In particular, the asymptotic variance of $\hat{f}_k(\cdot)$ has an order of $O(1/(nh_H^{b_k}))$ instead of $O(1/(nh_H^{b_k}))$ where $b_k < d_k$.

6.2 Multiple component functions

We next briefly discuss the extension of our method from the baseline model with two components to the case with more than two components.

Let K > 2. For any k = 2, ..., K, let $H_k(x^1, x^k) = \int H(x) \cdot p_{\tilde{X}^{-k}}(\tilde{x}^{-k}) d\tilde{x}^{-k}$ where \tilde{X}^{-k} is obtained by excluding X^1 and X^k from X, and \tilde{x}^{-k} is obtained by excluding x^1 and x^k from x. This constructed $H_k(x^1, x^k)$ is identified if the original H(x) is identified. We can transform the original model (2) with K components into the following new model with two components as

$$H_k(x^1, x^k) = G_k(f_1(x^1) + f_k(x^k)),$$
(9)

where $G_k(\tau) = \int G(\tau - f_k(x^k) + \sum_{\ell=2}^{K} f_\ell(x^\ell)) \cdot p_{\tilde{X}^{-k}}(\tilde{x}^{-k}) d\tilde{x}^{-k}$ is monotonic if the original link function $G(\cdot)$ is monotonic.

Our previous idea can be applied directly to the new model (9) to identify $f_1(\cdot)$ and $f_k(\cdot)$. Specifically, for any k = 2, ..., K, we use an idea similar to Theorem 1 to transform the new model (9) into the following model with two univariate components:

$$\mathcal{H}_k(z^1, z^k) = G_k(\tilde{f}_1(z^1) + \tilde{f}_k(z^k)), \tag{10}$$

where $\mathcal{H}_k(z^1, z^k) = E[H_k(X^1, X^k)|\zeta_1(X^1) = z^1, \zeta_k(X^k) = z^k]$, the inverse of $\tilde{f}_\ell(\cdot)$ is $\tilde{f}_\ell^{-1}(s) = \int G_k(s + f_{-\ell}(x^{-\ell})) \cdot w_{-\ell}(x^{-\ell}) dx^{-\ell}$, and $\zeta_\ell(x^\ell) = \int H_k(x^1, x^k) \cdot w_{-\ell}(x^{-\ell}) dx^{-\ell}$ with freely chosen weight functions $w_{-\ell}(\cdot)$ for $\ell = 1, k$ where $x^{-\ell}$ is x^k if $\ell = 1$ and is x^1 if $\ell = k$. The transformed components $\tilde{f}_1(\cdot)$ and $\tilde{f}_k(\cdot)$ can then be identified by (C1) and (C2), respectively, where $\mathcal{H}(\cdot)$ is replaced by $\mathcal{H}_k(\cdot)$. The original components are identified as $f_k(x^k) = \tilde{f}_k(\zeta_k(x^k))$ for all $k = 1, \ldots, K$. Once all of $f_k(\cdot), k = 1, \ldots, K$, are identified, the original link $G(\cdot)$ is identified by $G(\tau) = E[H(X)|\sum_{k=1}^K f_k(X^k) = \tau]$.

Similar to the case with two components (i.e. K = 2), we can closely follow the above identification strategy to estimate the link $G(\cdot)$ and the components $f_{\ell}(\cdot)$ in three steps for $\ell = 1, ..., K$. Let k = 2, ..., K. In the first step, we estimate the transformed function $\mathcal{H}_k(z^1, z^k)$ by the nonparametric sample analogue of its definition as $\widehat{E}[H_k(X^1, X^k)|\widehat{\zeta}_1(X^1) = z^1, \widehat{\zeta}_k(X^k) = z^k]$ where $H_k(x^1, x^k) = \int H(x) \cdot p_{\tilde{X}^{-k}}(\tilde{x}^{-k}) d\tilde{x}^{-k}$ and $\widehat{\zeta}_{\ell}(X^{\ell})$'s are also given by the sample analogues of $\zeta_{\ell}(X^{\ell}) = \int H_k(X^1, X^k) \cdot w_{-\ell}(X^{-\ell}) dX^{-\ell}$ for $\ell = 1, k$. Given the first-step estimator $\widehat{\mathcal{H}}_k(z^1, z^k)$, the second step follows Horowitz (2001)'s estimation procedure to estimate the transformed components $\tilde{f}_1(\cdot)$ and $\tilde{f}_k(\cdot)$ according to (C1) and (C2), respectively, with $\mathcal{H}(\cdot)$ replaced by $\mathcal{H}_k(\cdot)$ in the transformed model (10). Moreover, the link $G(\cdot)$ is recovered by $\widehat{G}(\tau) = \widehat{E}[Y|\sum_{\ell=1}^K \widehat{f}_\ell(\widehat{\zeta}_\ell(X^\ell)) = \tau]$. In step three, the original components are then recovered by $\widehat{f}_\ell(\cdot) = \widehat{f}_\ell(\widehat{\zeta}_\ell(\cdot))$ for $\ell = 1, \ldots, K$. Note that we will obtain (K - 1) estimates for the first component $f_1(\cdot)$. We hence aggregate them by their average to estimate $f_1(\cdot)$.

Li and Racine (2007) for more details.

7 Conclusion

In this paper, we consider estimating the generalized additive model with a flexible grouping and unknown link. To identify the model primitives, we transform the model into a new model with univariate components. We then identify the new model by applying the existing strategy for the generalized additive model with univariate components. Closely following the identification strategy, we propose a three-step procedure to estimate the link and original components. The consistency and asymptotic normality are then established for the link estimator at a one-dimensional convergence rate and for the component estimators at the convergence rates corresponding to their own dimensions.

This paper adopts a multi-step kernel method to estimate the component and link functions in the generalized additive model with a flexible additive structure and unknown link. Hahn, Liao, and Ridder (2018) studied nonparametric two-step sieve M estimation in a general class of semi/nonparametric models. As sieve method is convenient to implement in practice, it is interesting to use a multi-step sieve method to estimate the component and link functions in our framework. This is an interesting topic for future research.

Appendix

Appendix A proves the theorems given in the text. Appendix S.1 of Supplementary Material (SM) introduces some notations for the convenience of discussion in the text and proofs. All of technical lemmas are stated and shown in the Appendix S.2 of SM.

A Proofs of Theorems

A.1 Proof of Theorem 1

Proof. By definition, for k = 1, 2, we have

$$\zeta_k(x^k) = \int H(x) w_{-k}(x^{-k}) dx^{-k} = \int G(f_1(x^1) + f_2(x^2)) w_{-k}(x^{-k}) dx^{-k} = \delta_k(f_k(x^k)), \quad (A.1)$$

where the second equality comes from the model restriction (M). Here, the dependence of $\delta_k(\cdot)$ on the function $f_{-k}(\cdot)$ is abbreviated for simplicity of notation. It is easy to verify that $\delta_k(\cdot)$ is strictly monotonic and hence has an inverse function $\delta_k^{-1}(\cdot)$ if $G(\cdot)$ is strictly monotonic. Thus $f_k(x^k) = \delta_k^{-1}(\zeta_k(x^k))$. Because $\mathcal{H}(z) = E[H(X)|\zeta_1(X^1) = z^1, \zeta_2(X^2) = z^2]$ by definition, it follows that

$$\begin{aligned} \mathcal{H}(z) &= E \big[G \big(f_1(X^1) + f_2(X^2) \big) \Big| \zeta_1(X^1) = z^1, \zeta_2(X^2) = z^2 \big] \\ &= E \big[G \big(\delta_1^{-1}(\zeta_1(X^1)) + \delta_2^{-1}(\zeta_2(X^2)) \big) \Big| \zeta_1(X^1) = z^1, \zeta_2(X^2) = z^2 \big] \\ &= G \big(\delta_1^{-1}(z^1) + \delta_2^{-1}(z^2) \big). \end{aligned}$$

The desired conclusion is therefore established by letting $\tilde{f}_k(z^k) = \delta_k^{-1}(z^k)$ for k = 1, 2.

A.2 Proof of Theorem 3

Proof. Only the case for k = 2 is proved. The proof for k = 1 is similar. The definition of $\hat{f}_2(x^2)$ gives the following decomposition:

$$\widehat{f}_{2}(x^{2}) - f_{2}(x^{2}) = [\widehat{f}_{2}(\widehat{\zeta}_{2}(x^{2})) - \widetilde{f}_{2}(\widehat{\zeta}_{2}(x^{2}))] + [\widetilde{f}_{2}(\widehat{\zeta}_{2}(x^{2})) - \widetilde{f}_{2}(\zeta_{2}(x^{2}))], \quad (A.2)$$

where both terms on the right hand side of equality converge to 0 uniformly over $x^2 \in S_{X^2}$ in probability by Lemmas S.3 and S.6. Part (i) is hence established. The rest of proof is to show part (ii). The first term on the right hand side of (A.2) can be simplified as

$$\widehat{\tilde{f}}_{2}(\widehat{\zeta}_{2}(x^{2})) - \widetilde{f}_{2}(\widehat{\zeta}_{2}(x^{2})) = \widehat{\tilde{f}}_{2}(\zeta_{2}(x^{2})) - \widetilde{f}_{2}(\zeta_{2}(x^{2})) + O_{p}(\xi_{H2}(\xi_{\mathcal{H}}' + \xi_{H1})),$$
(A.3)

uniformly over x^2 as $n \to \infty$, where the third (remaining) term on the right hand side is due to $\int_{\zeta_2(x^2)}^{\widehat{\zeta}_2(x^2)} \int \left[\frac{\partial_2 \widehat{\mathcal{H}}(z)}{\partial_1 \widehat{\mathcal{H}}(z)} - \frac{\partial_2 \mathcal{H}(z)}{\partial_1 \mathcal{H}(z)}\right] w_3(z^1) dz^1 dz^2 = O_p\left(\xi_{H2}(\xi'_{\mathcal{H}} + \xi_{H1})\right)$ which is derived by applying a Taylor expansion similar to (S.2.17) on the (unweighted) integrand and Lemmas S.3 and S.5. Take a Taylor expansion to the second term on the right hand side of (A.2) to obtain

$$\tilde{f}_{2}(\hat{\zeta}_{2}(x^{2})) - \tilde{f}_{2}(\zeta_{2}(x^{2})) = \tilde{f}_{2}'(\zeta_{2}(x^{2}))(\hat{\zeta}_{2}(x^{2}) - \zeta_{2}(x^{2})) + O_{p}(\xi_{H2}^{2}),$$

uniformly over x^2 as $n \to \infty$. Consequently, with bandwidths satisfying Assumption A.6, the asymptotic representations of $\hat{f}_2(\cdot)$ given by Lemma S.6 and $\hat{\zeta}_2(\cdot)$ given by Lemma S.3 imply that

$$\widehat{f}_{2}(x^{2}) - f_{2}(x^{2}) = \widetilde{f}_{2}'(\zeta_{2}(x^{2})) \cdot J_{n2}(x^{2}) + \widetilde{\mathfrak{J}}_{n2}(\zeta_{2}(x^{2})) - E[\widetilde{\mathfrak{J}}_{n2}(\zeta_{2}(x^{2}))] + h_{H}^{r+1}[\widetilde{f}_{2}'(\zeta_{2}(x^{2}))D_{2}(x^{2}) + \widetilde{\mathcal{B}}_{2}(\zeta_{2}(x^{2}))] + h_{\mathcal{H}}^{r}\mathfrak{B}_{2}(\zeta_{2}(x^{2})) + o_{p}(h_{\mathcal{H}}^{r} + h_{H}^{r+1}),$$
(A.4)

uniformly over x^2 as $n \to \infty$. The asymptotic normality of part (ii) then follows by applying the Lindeberg-Feller central limit theorem (see Theorem 7.2.1 of Chung, 2001) to (A.4). The asymptotic bias is an immediate consequence of (A.4), and the asymptotic variance $\operatorname{Var}\left(\sqrt{nh_H^{d_2}} \cdot \left[\tilde{f}'_2(\zeta_2(x^2))J_{n2}(x^2) + \tilde{\mathfrak{J}}_{n2}(\zeta_2(x^2))\right]\right) = \sigma_2^2(x^2) + o(1)$ is obtained by a calculation similar to the one of asymptotic variance of a kernel density estimator. This completes the whole proof.

A.3 Proof of Theorem 4

Proof. For any i = 1, ..., n, let $T = T(x) = f_1(X^1) + f_2(X^2)$, $T_i = T(X_i) = f_1(X_i^1) + f_2(X_i^2)$, $\hat{T}_i = \hat{T}(x_i) = \hat{f}_{1i}(X_i^1) + \hat{f}_{2i}(X_i^2)$, and $p_T(\cdot)$ be the probability density function of T, where $\hat{f}_{ki}(\cdot)$ is the estimator of $f_k(\cdot)$ leaving observation i out for k = 1, 2. Since we have $\sup_{x \in S_X} |\hat{T}(x) - T(x)|$ similar to Lemma S.4, we can derive the asymptotic representation for any $\tau \in S_G$,

$$\begin{split} \widehat{G}(\tau) &- G(\tau) \\ = \frac{1}{nh_G} \sum_{i=1}^n e_{1G}' \mathcal{S}_{n,r}^G(\tau)^{-1} k \Big(\frac{T_i - \tau}{h_G} \Big) \Big\{ Y_i - \mu_G(T_i - \tau)' \beta_G(\tau) \Big\} \mu_G \Big(\frac{T_i - \tau}{h_G} \Big) + \frac{1}{nh_G^2} \sum_{i=1}^n e_{1G}' \mathcal{S}_{n,r}^G(\tau)^{-1} \cdot \\ & \left[\Big(\frac{\partial}{\partial u} t_G(u, Y_i; \tau) k(u) + t_G(u, Y_i; \tau) k'(u) \Big) \Big|_{u = \frac{T_i - \tau}{h_G}} \right] \cdot (\widehat{T}(X_i) - T_i) + o_p \Big(h_G^{r+1} + \sqrt{\log(n)/(nh_G)} \Big) \\ \end{split}$$

$$=:\Gamma_{1n}(\tau) + \Gamma_{2n}(\tau) + o_p \left(h_G^{r+1} + \sqrt{\log(n)/(nh_G)} \right)$$
(A.5)

uniformly over $\tau \in S_G$ as $n \to \infty$, where $t_G(u, Y_i; \tau) = \mu_G(u) (Y_i - \mu(u)' B_{h_G} \beta_G(\tau))$. The first term $\Gamma_{1n}(\tau)$ is the uniform Bahadur representation for local polynomial regression in Kong, Linton, and Xia (2010). The second term $\Gamma_{2n}(\tau)$ represents the error caused by using generated regressor \hat{T}_i . Thus, we have the uniform convergence of $\sup_{\tau \in S_G} |\hat{G}(\tau) - G(\tau)|$ and thus part (i) is proved. Similar to Lemma S.5, $\Gamma_{1n}(\tau)$ can be decomposed into a bias leading term and a stochastic leading term, i.e.⁹

$$\Gamma_{1n}(\tau) = \frac{1}{nh_G} \sum_{i=1}^n e_{1G}' \{S_r^G\}^{-1} \frac{Y_i - G(T_i)}{p_T(\tau)} \mu_G \Big(\frac{T_i - \tau}{h_G}\Big) k\Big(\frac{T_i - \tau}{h_G}\Big) + B_0(\tau) + R_{Gn},$$
(A.6)

where $R_{Gn} = o_p (h_G^{r+1} + \sqrt{1/(nh_G)})$, and $B_0(\tau) = e'_{1G} \{S_r^G\}^{-1} S_r^{G,r+1} G_{r+1}(\tau) \cdot h_G^{r+1}$. As for $\Gamma_{2n}(\tau)$, we can further decompose as under Assumption A.6,

$$\Gamma_{2n}(\tau) = \frac{1}{nh_G^2} \sum_{i=1}^n e_{1G}^{\prime} \{\mathcal{S}_r^G\}^{-1} p_T(\tau)^{-1} B_1(T_i, Y_i, \tau) \cdot \left(E[\widehat{T}(X_i)|X_i] - T_i \right) + \frac{1}{nh_G^2} \sum_{i=1}^n e_{1G}^{\prime} \{\mathcal{S}_r^G\}^{-1} \cdot p_T(\tau)^{-1} B_1(T_i, Y_i, \tau) \cdot \left(\widehat{T}(X_i) - E[\widehat{T}(X_i)|X_i] \right) + o_p \left(h_G^{r+1} + h_H^{r+1} + h_H^r + \sqrt{1/(nh_G)} \right)$$
$$=: \Gamma_{21n}(\tau) + \Gamma_{22n}(\tau) + o_p \left(h_G^{r+1} + h_H^{r+1} + h_H^r + \sqrt{1/(nh_G)} \right), \tag{A.7}$$

where $B_1(T_i, Y_i, \tau) = \left(\frac{\partial}{\partial u} t_G(u, Y_i; \tau) k(u) + t_G(u, Y_i; \tau) k'(u)\right) \Big|_{u = \frac{T_i - \tau}{h_G}}, E[\widehat{T}(X_i)|X_i] - T_i = \sum_{k=1}^2 B_{nfk}(X_i^k) = \sum_{k=1}^2 \left\{ h_{\mathcal{H}}^r \mathfrak{B}_k(\zeta_k(X_i^k)) + h_{\mathcal{H}}^{r+1}[\widetilde{f}'_k(\zeta_k(X_i^k))D_k(X_i^k) + \widetilde{\mathcal{B}}_k(\zeta_k(X_i^k))] \right\} \text{ and } \widehat{T}(X_i) - E[\widehat{T}(X_i)|X_i] = \sum_{k=1}^2 \left(\widetilde{f}'_k(\zeta_k(X_i^k)) \cdot J_{nk}(X_i^k) + \widetilde{\mathfrak{J}}_{nk}(\zeta_k(X_i^k)) - E[\widetilde{\mathfrak{J}}_{nk}(\zeta_k(X_i^k))|X_i] \right).$ $\Gamma_{21n}(\tau)$ is the additional bias due to the generated regressor $\widehat{T}(X_i)$. Similar to the arguments in (S.2.13) of Lemma S.5, we get

$$\Gamma_{21n}(\tau) = -e_{1G}'\{\mathcal{S}_{r}^{G}\}^{-1}p_{T}(\tau)^{-1}\frac{1}{h_{G}}\int\mu_{G}(\frac{T-\tau}{h_{G}})k(\frac{T-\tau}{h_{G}})G'(T)\cdot E\left[\sum_{k=1}^{2}B_{nfk}(X_{i}^{k})\big|T_{i}=T\right]p_{T}(T)dT + R_{0n}$$
$$= -G'(\tau)\cdot E\left[\sum_{k=1}^{2}\left\{h_{\mathcal{H}}^{r}\mathfrak{B}_{k}(\zeta_{k}(x^{k})) + h_{H}^{r+1}[\tilde{f}_{k}'(\zeta_{k}(x^{k}))D_{k}(x^{k}) + \tilde{\mathcal{B}}_{k}(\zeta_{k}(x^{k}))]\right\}\Big|T_{I}=\tau\right] + R_{0n}, \quad (A.8)$$

where $R_{0n} = o_p(h_H^{r+1} + h_H^r)$, and the last equality is due to (i) change of variables, (ii) Taylor expansion, and (iii) the fact that $e'_{1G} \{S_r^G\}^{-1} \int \mu_G(u)k(u)du = e'_{1G}e_{1G} = 1$.

Next consider $\Gamma_{22n}(\tau)$. It represents the additional stochastic term induced by $\widehat{T}(X_i)$. Similar to Lemma S.5, under Assumption A.6, $\Gamma_{22n}(\tau)$ can be written as

$$\Gamma_{22n}(\tau) = -e_{1G}'\{\mathcal{S}_{r}^{G}\}^{-1}p_{T}(\tau)^{-1}\frac{1}{nh_{G}}\sum_{i=1}^{n}\mu_{G}\left(\frac{T_{i}-\tau}{h_{G}}\right)k\left(\frac{T_{i}-\tau}{h_{G}}\right)G'(T_{i})$$

$$\cdot\sum_{k=1}^{2}\left(\tilde{f}_{k}'(\zeta_{k}(X_{i}^{k}))\cdot J_{nk}(X_{i}^{k}) + \tilde{\mathfrak{J}}_{nk}(\zeta_{k}(X_{i}^{k})) - E[\tilde{\mathfrak{J}}_{nk}(\zeta_{k}(X_{i}^{k}))|X_{i}]\right) + R_{1n}, \quad (A.9)$$

where $R_{1n} = o_p (h_G^{r+1} + h_H^{r+1} + h_H^r + \sqrt{1/(nh_G)})$. Follow the U-Statistics arguments similar to Lemma

⁹Here, We derive a weaker, point-wise representation rather than the uniform representation in Lemma S.5.

8 of Horowitz (1998)¹⁰, (A.9) can be represented as

$$\begin{split} \Gamma_{22n}(\tau) = & e_{1G}^{\prime}\{\mathcal{S}_{r}^{G}\}^{-1} (\int \mu_{G}(u)k(u)du) G^{\prime}(\tau) \sum_{k=1}^{2} \{\Gamma_{22n,k}(\tau) - E[\Gamma_{22n,k}(\tau)] + \widetilde{\Gamma}_{22n,k}(\tau) - E[\widetilde{\Gamma}_{22n,k}(\tau)] \} + R_{1n} \\ = & G^{\prime}(\tau) \sum_{k=1}^{2} \Gamma_{22n,k}(\tau) + G^{\prime}(\tau) \sum_{k=1}^{2} \widetilde{\Gamma}_{22n,k}(\tau) + R_{1n} \end{split}$$

for all $\tau \in S_G$, where $\Gamma_{22n,k}(\tau) = \frac{1}{nh_{\mathcal{H}}} \sum_{i=1}^n c^{2-k} \omega_{5-k}(Z_i^{-k}) q_k(\mathcal{Z}_{ki}^0)' e'_d \tilde{S}_r^{-1} V_k^{\tilde{\mu}} \left(\frac{Z_i^k - Z_0^k}{h_{\mathcal{H}}}\right) \frac{Y_i - \mathcal{H}(Z_i)}{p_Z(\mathcal{Z}_{ki}^0)} \mathcal{K}_k \left(\frac{z_0^k - Z_i^k}{h_{\mathcal{H}}}\right),$ $\tilde{\Gamma}_{22n,k}(\tau) = -\frac{1}{n} \sum_{i=1}^n \tilde{f}'_k (\zeta_k(X_i^k)) \frac{p_{X^1|T}(X_i^k|\tau)}{p_{X^k|X-k}(X_i^k|X_i^{-k})} (Y_i - \mathcal{H}(X_i)), \text{ and } E[\Gamma_{22n,k}(\tau)] = E[\tilde{\Gamma}_{22n,k}(\tau)] = 0.$ $\Gamma_{22n,k}(\tau)$ is the stochastic term due to the estimation of $\tilde{f}_k(\cdot)$, i.e. $\tilde{\mathfrak{J}}_{nk}(\zeta_k(X_i^k))$, and has a order of $O_p(1/\sqrt{nh_{\mathcal{H}}})$. $\tilde{\Gamma}_{22n,k}(\tau)$ is induced by the estimation of $\zeta_k(X_i^k)$, i.e. $\tilde{f}'_k(\zeta_k(X_i^k)) \cdot J_{nk}(X_i^k)$ with a

of $O_p(1/\sqrt{nh_{\mathcal{H}}})$. $\tilde{\Gamma}_{22n,k}(\tau)$ is induced by the estimation of $\zeta_k(X_i^k)$, i.e. $f'_k(\zeta_k(X_i^k)) \cdot J_{nk}(X_i^k)$ with a variance of order $O(1/\sqrt{n})$. Therefore, $\tilde{\Gamma}_{22n,k}(\tau)$ is of smaller order than $\Gamma_{22n,k}(\tau)$ and we conclude that

$$\Gamma_{22n}(\tau) = G'(\tau) \sum_{k=1}^{2} \Gamma_{22n,k}(\tau) + R_{1n}.$$
(A.10)

By combining the bias leading terms of (A.6) and (A.8), the asymptotic bias of $\hat{G}(\tau)$ can be established. By the stochastic parts of (A.6) and (A.10), the asymptotic normality and correspondent variance follow from Lindeberg-Feller central limit theorem. This complete the whole proof.

¹⁰Note that we use our Lemma S.1 instead of Horowitz (1998)'s Lemma 5 to characterize the projection error of U-statistics.

References

- AI, C., AND X. CHEN (2003): "Efficient estimation of models with conditional moment restrictions containing unknown functions," *Econometrica*, 71(6), 1795–1843.
- ANTRAS, P. (2004): "Is the US aggregate production function Cobb-Douglas? New estimates of the elasticity of substitution," *Contributions in Macroeconomics*, 4(1).
- ATHEY, S., AND P. A. HAILE (2002): "Identification of standard auction models," *Econometrica*, 70(6), 2107–2140.
- BERKOWITZ, D., H. MA, AND S. NISHIOKA (2017): "Recasting the iron rice bowl: The reform of China's state-owned enterprises," *Review of Economics and Statistics*, 99(4), 735–747.
- CHEN, R., W. HÄRDLE, O. B. LINTON, AND E. SEVERANCE-LOSSIN (1996): "Nonparametric estimation of additive separable regression models," in *Statistical Theory and Computational Aspects of Smoothing*, pp. 247–265. Springer.
- CHEN, S. (2002): "Rank estimation of transformation models," Econometrica, 70(4), 1683–1697.
- (2010a): "An integrated maximum score estimator for a generalized censored quantile regression model," *Journal of Econometrics*, 155(1), 90–98.
- —— (2010b): "Root-N-consistent estimation of fixed-effect panel data transformation models with censoring," *Journal of econometrics*, 159(1), 222–234.
- (2012): "Distribution-free estimation of the Box-Cox regression model with censoring," *Econometric Theory*, pp. 680–695.
- CHEN, S., AND H. ZHANG (2020): "Root-N-prediction of generalized heteroscedastic transformation regression models," *Journal of Econometrics*, 215(2), 305–340.
- CHEN, X. (2007): "Large sample sieve estimation of semi-nonparametric models," Handbook of econometrics, 6, 5549–5632.
- CHESHER, A. (2003): "Identification in nonseparable models," Econometrica, 71(5), 1405–1441.
- CHUNG, K. L. (2001): A course in probability theory. Academic press.
- DE PAULA, A., AND X. TANG (2012): "Inference of signs of interaction effects in simultaneous games with incomplete information," *Econometrica*, 80(1), 143–172.
- FAN, J., AND I. GIJBELS (1992): "Variable Bandwidth and Local Linear Regression Smoothers," *The Annals of Statistics*, 20(4), 2008–2036.
- GENTRY, M., AND T. LI (2014): "Identification in auctions with selective entry," *Econometrica*, 82(1), 315–344.
- GRIECO, P. L. (2014): "Discrete games with flexible information structures: An application to local grocery markets," *The RAND Journal of Economics*, 45(2), 303–340.

- GUERRE, E., I. PERRIGNE, AND Q. VUONG (2000): "Optimal nonparametric estimation of first-price auctions," *Econometrica*, 68(3), 525–574.
- —— (2009): "Nonparametric identification of risk aversion in first-price auctions under exclusion restrictions," *Econometrica*, 77(4), 1193–1227.
- HAHN, J., Z. LIAO, AND G. RIDDER (2018): "Nonparametric two-step sieve M estimation and inference," *Econometric Theory*, 34(6), 1281–1324.
- HALL, P. (1992): "On bootstrap confidence intervals in nonparametric regression," *The Annals of Statistics*, pp. 695–711.
- HALL, P., AND J. HOROWITZ (2013): "A simple bootstrap method for constructing nonparametric confidence bands for functions," *Annals of Statistics*, 41(4), 1892–1921.
- HANSEN, B. E. (2009): Lecture notes on nonparametrics. University of Wisconsin-Madison.
- HÄRDLE, W., AND A. W. BOWMAN (1988): "Bootstrapping in nonparametric regression: Local adaptive smoothing and confidence bands," *Journal of the American Statistical Association*, 83(401), 102–110.
- HODERLEIN, S., L. SU, H. WHITE, AND T. T. YANG (2016): "Testing for monotonicity in unobservables under unconfoundedness," *Journal of Econometrics*, 193(1), 183–202.
- HODGES, D. J. (1969): "A note on estimation of Cobb-Douglas and CES production function models," *Econometrica*, pp. 721–725.
- HOROWITZ, J. (1998): "Nonparametric estimation of a generalized additive model with an unknown link function," Working paper, University of Iowa.
- HOROWITZ, J. L. (2001): "Nonparametric estimation of a generalized additive model with an unknown link function," *Econometrica*, 69(2), 499–513.
- HOROWITZ, J. L., AND E. MAMMEN (2004): "Nonparametric estimation of an additive model with a link function," *Annals of Statistics*, 32(6), 2412–2443.
- —— (2007): "Rate-optimal estimation for a general class of nonparametric regression models with unknown link functions," *Annals of Statistics*, 35(6), 2589–2619.
- —— (2011): "Oracle-efficient nonparametric estimation of an additive model with an unknown link function," *Econometric Theory*, pp. 582–608.
- JACHO-CHÁVEZ, D., A. LEWBEL, AND O. LINTON (2010): "Identification and nonparametric estimation of a transformed additively separable model," *Journal of Econometrics*, 156(2), 392–407.
- KHAN, S. (2001): "Two-stage rank estimation of quantile index models," *Journal of Econometrics*, 100(2), 319–355.
- KLUMP, R., P. MCADAM, AND A. WILLMAN (2007): "Factor substitution and factor-augmenting technical progress in the United States: a normalized supply-side system approach," *The Review of Economics and Statistics*, 89(1), 183–192.

- KMENTA, J. (1967): "On estimation of the CES production function," *International Economic Review*, 8(2), 180–189.
- KOHLER, M., AND A. KRZYŻAK (2017): "Nonparametric Regression Based on Hierarchical Interaction Models," *IEEE Transactions on Information Theory*, 63(3), 1620–1630.
- KONG, E., O. LINTON, AND Y. XIA (2010): "Uniform bahadur representation for local polynomial estimates of m-regression and its application to the additive model," *Econometric Theory*, 26(5), 1529–1524.
- LEWBEL, A., X. LU, AND L. SU (2015): "Specification testing for transformation models with an application to generalized accelerated failure-time models," *Journal of Econometrics*, 184(1), 81–96.
- LEWBEL, A., AND X. TANG (2015): "Identification and estimation of games with incomplete information using excluded regressors," *Journal of Econometrics*, 189(1), 229–244.
- LI, H., AND N. LIU (2018): "Nonparamametric identification and estimation of double auctions with bargaining," Working paper, Shanghai University of Finance and Economics.
- LI, Q., AND J. S. RACINE (2007): *Nonparametric econometrics: theory and practice*. Princeton University Press.
- LI, T., AND X. ZHENG (2009): "Entry and competition effects in first-price auctions: theory and evidence from procurement auctions," *The Review of Economic Studies*, 76(4), 1397–1429.
- LIN, H., L. PAN, S. LV, AND W. ZHANG (2018): "Efficient estimation and computation for the generalised additive models with unknown link function," *Journal of Econometrics*, 202(2), 230–244.
- LINTON, O., AND W. HÄRDLE (1996): "Estimation of additive regression models with known links," *Biometrika*, 83(3), 529–540.
- LIU, N., AND Y. LUO (2017): "A Nonparametric Test for Comparing Valuation Distributions in First-Price Auctions," *International Economic Review*, 58(3), 857–888.
- LIU, N., AND Q. VUONG (2020): "Nonparametric tests for monotonicity of strategies in games of incomplete information," Working paper, New York University.
- LIU, N., Q. VUONG, AND H. XU (2017): "Rationalization and identification of binary games with correlated types," *Journal of Econometrics*, 201(2), 249–268.
- MA, S. (2012): "Two-step spline estimating equations for generalized additive partially linear models with large cluster sizes," *The Annals of Statistics*, 40(6), 2943–2972.
- MA, S., AND P. X.-K. SONG (2015): "Varying index coefficient models," *Journal of the American Statistical Association*, 110(509), 341–356.
- MARMER, V., AND A. SHNEYEROV (2012): "Quantile-based nonparametric inference for first-price auctions," *Journal of Econometrics*, 167(2), 345–357.
- MASRY, E. (1996): "Multivariate local polynomial regression for time series: uniform strong consistency and rates," *Journal of Time Series Analysis*, 17(6), 571–599.

- MATZKIN, R. L. (2003): "Nonparametric estimation of nonadditive random functions," *Econometrica*, 71(5), 1339–1375.
- PARASKEVOPOULOS, C. C. (1979): "Alternative estimates of the elasticity of substitution: An intermetropolitan CES production function analysis of US manufacturing industries, 1958-1972," *The Review of Economics and Statistics*, pp. 439–442.
- PINKSE, J. (2001): "Nonparametric regression estimation using weak separability," Working paper, University of British Columbia.
- POLLARD, D. (1984): Convergence of stochastic processes. Springer-Verlag.
- POWELL, J. L., J. H. STOCK, AND T. M. STOKER (1989): "Semiparametric estimation of index coefficients," *Econometrica*, pp. 1403–1430.
- SCHMIDT-HIEBER, J. (2020): "Nonparametric regression using deep neural networks with ReLU activation function,".
- SILVERMAN, B. W. (1986): Density Estimation for Statistics and Data Analysis. Chapman and Hall/CRC.
- TANG, X. (2010): "Estimating simultaneous games with incomplete information under median restrictions," *Economics Letters*, 108(3), 273–276.

Supplementary material to "Nonparametric identification and estimation of a generalized additive model with a flexible additive structure and unknown link"

Songnian Chen, Nianqing Liu, Jian Zhang, and Yahong Zhou

Abstract

The supplementary material includes two appendices. Appendix S.1 introduces some notations for the convenience of discussion in the text and proofs. Appendix S.2 states and proves some technical lemmas needed to show the main theorems in the text.

S.1 Notations

S.1.1 Local polynomial Regression

For *r*-th order local polynomial regression of Y_i on X_i , let $\mathbf{j} = (j_1, j_2, \cdots, j_d)$ be an arbitrary *d*-tuple of integers, denote $|\mathbf{j}| = j_1 + j_2 + \cdots + j_d$, $\mathbf{j}! = j_1! \times j_2! \times \cdots \times j_d!$, $x^{\mathbf{j}} = (x^1)^{j_1} \times (x^2)^{j_2} \times \cdots \times (x^d)^{j_d}$, $D^{\mathbf{j}}H(x) = \frac{\partial^{|\mathbf{j}|}H(x)}{\partial (x^1)^{j_1}\partial (x^2)^{j_2}\cdots\partial (x^d)^{j_d}}$, and $\sum_{0 \le |\mathbf{j}| \le r} = \sum_{k=0}^r \sum_{j_1+j_2+\cdots+j_d=k}$. The total number of *d*-tuples with $|\mathbf{j}| = s$ is $M_s = \binom{r+d-1}{d-1}$. We arrange these tuples in an ascending lexicographical order style as in Masry (1996)¹¹. The correspondent position of each tuple forms a one-to-one map which is called π_s , i.e. $\pi_s(1) = (s, 0, 0, \cdots, 0), \dots, \pi_s(M_s) = (0, 0, \cdots, 0, s)$. Denote a vector-value function $\mu(\cdot)$ for

 π_s , i.e. $\pi_s(1) = (s, 0, 0, \dots, 0), \dots, \pi_s(M_s) = (0, 0, \dots, 0, s)$. Denote a vector-value function $\mu(\cdot)$ for an arbitrary entry $x \in \mathbb{R}^d$ such that $\mu_s(x)$ is a $M_s \times 1$ vector with *l*-entry given by $[\mu_s(x)]_l = x^{\pi_s(l)}$. and we stack these vectors and define a $N_r \times 1$ vector as $\mu(x) = [\mu_0(x), \mu_1(x), \dots, \mu_r(x)]'$, where $N_r = M_0 + M_1 + \dots + M_r$. Also, we denote $M_s \times 1$ vectors $H_s(x)$ ($s = 0, 1, \dots, r+1$) to store H(x) and its derivatives (up to (r+1)-th order) such that the *l*-entry of $\alpha_s(x)$ equals to $[H_s(x)]_l = \frac{1}{\pi_s(l)!}D^{\pi_s(l)}H(x)$, and $\alpha(x)$ stacks $\alpha_s(x)$ ($s = 0, 1, \dots, r$.) as $\alpha(x) = [H_0(x), H_1(x), \dots, H_r(x)]'$, then $\mu(y - x)'\alpha(x)$ is the *r*-th order Taylor expansion of H(y) at x. Let $S_{n,p,q}(x)$ and $S_{p,q}$ be $M_p \times M_q$ matrices with (l, k)-element given by $[S_{n,p,q}(x)]_{l,k} = \int u^{\pi_p(l) + \pi_q(k)}K(u)p_X(x + h_Hu)du$ and $[S_{p,q}]_{l,k} = \int u^{\pi_p(l) + \pi_q(k)}K(u)du$, where $u = (u_1, u_2, \dots, u_d)$, $K(u) = K_1(u^1)K_2(u^2)$ with $u^1 = (u_1, \dots, u_d_1)$ and $u^2 = (u_{d_1+1}, \dots, u_d)$, and $p_X(\cdot)$ is the probability density function of X. Define $N_r \times N_r$ matrices $S_{n,r}(x)$ and S_r as

$$S_{n,r}(x) = \begin{pmatrix} S_{n,0,0}(x) & S_{n,0,1}(x) & \cdots & S_{n,0,r}(x) \\ S_{n,1,0}(x) & S_{n,1,1}(x) & \cdots & S_{n,1,r}(x) \\ \vdots & \vdots & \ddots & \vdots \\ S_{n,r,0}(x) & S_{n,r,1}(x) & \cdots & S_{n,r,r}(x) \end{pmatrix}, S_r = \begin{pmatrix} S_{0,0} & S_{0,1} & \cdots & S_{0,r} \\ S_{1,0} & S_{1,1} & \cdots & S_{1,r} \\ \vdots & \vdots & \ddots & \vdots \\ S_{r,0} & S_{r,1} & \cdots & S_{r,r} \end{pmatrix},$$

and $N_r \times M_{r+1}$ matrices $S_{n,r}^{r+1}(x)$ and S_r^{r+1} as $S_{n,r}^{r+1}(x) = (S_{n,0,r+1}(x)', S_{n,1,r+1}(x)', \cdots, S_{n,r,r+1}(x)')'$ and $S_r^{r+1} = (S_{0,r+1}', S_{1,r+1}', \cdots, S_{r,r+1}')'$. For *r*-th order local polynomial regression of Y_i on $Z_i = \zeta(X_i)$, similarly, for each 2-tuple $\tilde{\mathbf{j}} = (j_1, j_2)$, we can define summation, factorial operation, multiplication and partial derivatives. In the same style as M_s , N_r , $\pi_s(\cdot)$, we can define \tilde{M}_s , \tilde{N}_r , and the lexicographical

¹¹The highest priority of the order is based on j_1 , second we order by j_2 , so on and so forth, finally we order by j_d .

order map $\tau_s(\cdot)$. Similar to $\mu(\cdot)$, $S_{n,p,q}(x)$, $S_{p,q}$, $S_{n,r}(x)$, and S_r , we can define $\tilde{\mu}(\cdot)$, $\tilde{S}_{n,p,q}(z)$, $\tilde{S}_{p,q}$, $\tilde{S}_{n,r}(z)$, and \tilde{S}_r with $z = (z^1, z^2)$. Let $S_{n,p,q}(z, \zeta)$ be a $\tilde{M}_p \times \tilde{M}_q$ matrix with (l,k)-element defined by $[S_{n,p,q}(z,\zeta)]_{l,k} = \frac{1}{nh_{\mathcal{H}}^2} \sum_{i=1}^n \left(\frac{\zeta(X_i)-z}{h_{\mathcal{H}}}\right)^{\tau_p(l)+\tau_q(k)} \tilde{K}\left(\frac{\zeta(X_i)-z}{h_{\mathcal{H}}}\right)$, and $Q_{n,p,0}(z,\zeta)$ be a $\tilde{M}_p \times 1$ vector with k-th entry given by $[Q_{n,p,0}(z,\zeta)]_k = \frac{1}{nh_{\mathcal{H}}^2} \sum_{i=1}^n Y_i \left(\frac{\zeta(X_i)-z}{h_{\mathcal{H}}}\right)^{\tau_p(k)} \tilde{K}\left(\frac{\zeta(X_i)-z}{h_{\mathcal{H}}}\right)$, where $z = (z^1, z^2)$, $\zeta = (\zeta^1, \zeta^2)$, $\zeta(X_i) = (\zeta_1(X_i^1), \zeta_2(X_i^2))$, and $\tilde{K}(u) = k_1(u^1)k_2(u^2)$. Also, we define the kernel derivatives $\partial_k \tilde{K}(u) = k'_k(u^k)k_{-k}(u^{-k})$ for k = 1, 2. By stacking $S_{n,p,q}(z,\zeta)$ and $Q_{n,p}(z,\zeta)$, we define a $\tilde{N}_r \times \tilde{N}_r$ matrix $S_{n,r}(z,\zeta)$ and a $\tilde{N}_r \times 1$ vector $Q_{n,r}(z,\zeta)$ as

$$\mathcal{S}_{n,r}(z,\zeta) = \begin{pmatrix} \mathcal{S}_{n,0,0}(z,\zeta) & \mathcal{S}_{n,0,1}(z,\zeta) & \cdots & \mathcal{S}_{n,0,r}(z,\zeta) \\ \mathcal{S}_{n,1,0}(z,\zeta) & \mathcal{S}_{n,1,1}(z,\zeta) & \cdots & \mathcal{S}_{n,1,r}(z,\zeta) \\ \vdots & \vdots & \ddots & \vdots \\ \mathcal{S}_{n,r,0}(z,\zeta) & \mathcal{S}_{n,r,1}(z,\zeta) & \cdots & \mathcal{S}_{n,r,r}(z,\zeta) \end{pmatrix}, \quad \mathcal{Q}_{n,r}(z,\zeta) = \begin{bmatrix} \mathcal{Q}_{n,0,0}(z,\zeta) \\ \mathcal{Q}_{n,1,0}(z,\zeta) \\ \vdots \\ \mathcal{Q}_{n,r,0}(z,\zeta) \end{bmatrix}.$$

Then infeasible local polynomial estimator is $\tilde{\beta}(z) = B_{\mathcal{H}}^{-1} S_{n,r}(z,\zeta)^{-1} Q_{n,r}(z,\zeta)$ with unknown parameter $\zeta(\cdot)$, and correspondent feasible estimator is $\hat{\beta}(z) = B_{\mathcal{H}}^{-1} S_{n,r}(z,\hat{\zeta})^{-1} Q_{n,r}(z,\hat{\zeta})$, where $B_{\mathcal{H}}$ is a $\tilde{N}_r \times \tilde{N}_r$ diagonal matrix with diagonal vector $D_h = [D_{h,0}, D_{h,1}, \cdots, D_{h,r}]'$ and $D_{h,s} = (h_{\mathcal{H}}^{|\tau(k)|})_{k=1,2,\dots,\tilde{M}_s}$. In order to represent the first-order derivatives of $\mathcal{H}(z)$ by $\beta(z)$, we introduce a $\tilde{N}_r \times 2$ vector

In order to represent the first-order derivatives of $\mathcal{H}(z)$ by $\beta(z)$, we introduce a $N_r \times 2$ vector given by $e_d = \begin{pmatrix} 0 & 1 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 1 & 0 & \cdots & 0 \end{pmatrix}'$, then $(\partial_1 \mathcal{H}(z), \partial_2 \mathcal{H}(z))' = e'_d \beta(z)$. For *r*-th local polynomial regression of Y_i on $T_i = f_1(X_i^1) + f_2(X_i^2)$. Similar to $\mu(\cdot)$, $S_{n,r}(x)$, and S_r , we can also define $\mu_G(\cdot)$, $S_{n,r}^G(\tau)$, and S_r^G .

S.1.2 Convergence rates, bias and variance terms

For k = 1, 2,

$$B_{nf_{k}}(x^{k}) = h_{\mathcal{H}}^{r} \mathfrak{B}_{k}(\zeta_{k}(x^{k})) + h_{H}^{r+1} [\tilde{f}_{k}^{r}(\zeta_{k}(x^{k}))D_{k}(x^{k}) + \tilde{\mathcal{B}}_{k}(\zeta_{k}(x^{k}))],$$
(S.1.1)

$$\sigma_k^2(x^k) = \tilde{f}'_k(\zeta_k(x^k))^2 \cdot \left[\int \left(e_1' S_r^{-1} V_k^\mu(t) \right)^2 K_k(t)^2 dt \right] \int \frac{E[(Y - H(x))^- |X = x]}{p_{X^k | X^{-k}}(x^k | x^{-k})^2} \cdot p_X(x) dx^{-k}, \quad (S.1.2)$$

$$B_k c(\tau) = h^{r+1} e_X' c \left\{ S_k^G \right\}^{-1} S_k^{G,r+1} C_{r+1}(\tau)$$

$$B_{nG}(\tau) = h_G^{r+1} e_{1G}^{r} \{S_r^{\sigma}\}^{-1} S_r^{\sigma,r+1} G_{r+1}(\tau) - h_H^{r+1} G'(\tau) \sum_{k=1}^2 E\left[\tilde{f}_k'(\zeta_k(X^k)) D_k(X^k) + \tilde{\mathcal{B}}_k(\zeta_k(X^k)) \middle| T = \tau\right] - h_H^r G'(\tau) \sum_{k=1}^2 E\left[\mathfrak{B}_k(\zeta_k(X^k)) \middle| T = \tau\right],$$
(S.1.3)

$$\begin{split} \sigma_{G}^{2}(\tau) &= \frac{\operatorname{Var}(Y|T=\tau)}{p_{T}(\tau)} \int \left(e_{1G}^{\prime} \{S_{r}^{G}\}^{-1} \mu_{G}(u) \right)^{2} k(t)^{2} dt + \delta_{G} \cdot \sigma_{G2}^{2}(\tau), \end{split} \tag{S.1.4} \\ \sigma_{G2}^{2}(\tau) &= G^{\prime}(\tau)^{2} \sum_{k=1}^{2} \int c^{2(2-k)} \omega_{5-k} (Z_{i}^{-k})^{2} \\ &\cdot \left\{ \int \left(q_{k}(\mathcal{Z}_{ki}^{0})^{\prime} e_{d}^{\prime} \tilde{S}_{r}^{-1} V_{k}^{\tilde{\mu}}(t) \right)^{2} \mathcal{K}_{k}(t)^{2} dt \right\} \cdot \frac{\operatorname{Var}(Y|Z=\mathcal{Z}_{ki}^{0})}{p_{Z}(\mathcal{Z}_{ki}^{0})} \, dZ_{i}^{-k}, \end{split}$$

(S.1.5)

$$J_{nk}(x^{k}) = \frac{1}{nh_{H}^{d_{k}}} \sum_{i=1}^{n} K_{k}\left(\frac{X_{i}^{k} - x^{k}}{h_{H}}\right) \cdot \frac{Y_{i} - H(x_{i})}{p_{X^{k}|X^{-k}}(x^{k}|X_{i}^{-k})} e_{1}^{\prime} S_{r}^{-1} V_{k}^{\mu}\left(\frac{X_{i}^{k} - x^{k}}{h_{H}}\right),$$
(S.1.6)

$$\begin{aligned} D_{k}(x^{k}) &= e_{1}^{\prime}S_{r}^{-1}S_{r}^{r+1}\int H_{r+1}(x^{k},x^{-k})p_{X^{-k}}(x^{-k})dx^{-k}, \end{aligned} \tag{S.1.7} \\ \tilde{\mathfrak{J}}_{nk}(z^{k}) &= \frac{c^{2-k}}{nh_{\mathcal{H}}}\sum_{i=1}^{n}\omega_{5-k}(Z_{i}^{-k})\left[q_{k}(\mathcal{Z}_{ki})'e_{d}'\tilde{\mathfrak{S}}_{r}^{-1}\frac{Y_{i}-H(X_{i})}{p_{Z}(\mathcal{Z}_{ki})}V_{k}^{\tilde{\mu}}\Big(\frac{\zeta_{k}(X_{i}^{k})-z^{k}}{h_{\mathcal{H}}}\Big)\mathcal{K}_{k}\Big(\frac{z^{k}-\zeta_{k}(X_{i}^{k})}{h_{\mathcal{H}}}\Big) \\ &-q_{k}(\mathcal{Z}_{ki}^{0})'e_{d}'\tilde{\mathfrak{S}}_{r}^{-1}\frac{Y_{i}-H(X_{i})}{p_{Z}(\mathcal{Z}_{ki}^{0})}V_{k}^{\tilde{\mu}}\Big(\frac{\zeta_{k}(X_{i}^{k})-z_{0}^{k}}{h_{\mathcal{H}}}\Big)\mathcal{K}_{k}\Big(\frac{z_{0}^{k}-\zeta_{k}(X_{i}^{k})}{h_{\mathcal{H}}}\Big)\right] \end{aligned}$$

$$\begin{split} \mathfrak{B}_{k}(z^{k}) &= c^{2-k} \int_{z_{0}^{k}}^{z^{k}} \int q_{k}(\nu)' \mathfrak{D}(\nu) \omega_{5-k}(\nu^{-k}) d\nu^{-k} d\nu^{k}, \\ \tilde{\mathcal{B}}_{k}(z^{k}) &= c^{2-k} \int_{z_{0}^{k}}^{z^{k}} \int q_{k}(\nu)' \mathcal{D}(\nu) \omega_{5-k}(\nu^{-k}) d\nu^{-k} d\nu^{k}, \\ \mathfrak{D}(z) &= e'_{d} \tilde{S}_{r}^{-1} \tilde{S}_{r}^{r+1} \mathcal{H}_{r+1}(z), \\ \mathcal{D}(z) &= \tilde{e}_{1} \mathcal{D}_{1}(z) + \tilde{e}_{2} \mathcal{D}_{2}(z), \end{split}$$

where $c = \int \omega_3(z^1) \cdot \left[\int [\partial_1 \mathcal{H}(z)/\partial_2 \mathcal{H}(z)] \cdot \omega_4(z^2) dz^2 \right]^{-1} dz^1$, $V_k^{\mu}(u^k) = \int \mu(u^k, t^{-k}) K_{-k}(t^{-k}) dt^{-k}$, $V_k^{\tilde{\mu}}(\tilde{u}^k) = \int \tilde{\mu}(\tilde{u}^k, \tilde{t}^{-k}) k_{-k}(\tilde{t}^{-k}) d\tilde{t}^{-k}$, $T = f_1(X^1) + f_2(X^2)$, $\mathcal{X}_{1i} = (x^1, X_i^2)$, $\mathcal{X}_{2i} = (X_i^1, x^2)$, $\mathcal{Z}_{1i} = (z^1, Z_i^2)$, $\mathcal{Z}_{2i} = (Z_i^1, z^2)$, $\mathcal{Z}_{1i}^0 = (z_0^1, Z_i^2)$, $\mathcal{Z}_{2i}^0 = (Z_i^1, z_0^2)$, $\mathcal{K}_K(u^k) = \int_{-\infty}^{u^k} k_k(t^k) dt^k$, x_s^k denotes the s-th element of x^k , $e_1 = (1, 0, 0, \dots, 0)'$ is a $N_r \times 1$ vector, $e_{1G} = (1, 0, 0, \dots, 0)'$ is a $(r + 1) \times 1$ vector, $\tilde{e}_1 = (1, 0)'$, $\tilde{e}_2 = (0, 1)'$, $q_2(\nu) = \left[- \frac{\partial_2 \mathcal{H}(\nu)}{[\partial_1 \mathcal{H}(\nu)]^2}, \frac{1}{\partial_1 \mathcal{H}(\nu)} \right]'$, $q_1(\nu) = \left[\frac{1}{\partial_2 \mathcal{H}(\nu)}, -\frac{\partial_1 \mathcal{H}(\nu)}{[\partial_2 \mathcal{H}(\nu)]^2} \right]'$, and

$$\mathcal{D}_k(z) = -p_Z(z)^{-1} \frac{\partial}{\partial z^k} \Biggl\{ \sum_{\ell=1}^2 \frac{\partial}{\partial z^\ell} \mathcal{H}(z) \int D_\ell(x^\ell) p_{X^\ell | Z}(x^\ell | z) dx^\ell \cdot p_Z(z) \Biggr\},$$

where $D_k(x^k)$ are given by (S.1.7), respectively.

Furthermore, let $\xi_H = h_H^{r+1} + \sqrt{\log(n)/(nh_H^d)}$, $\xi_H = h_H^{r+1} + \sqrt{\log(n)/(nh_H^2)}$, $\xi'_H = h_H^r + \sqrt{\log(n)/(nh_H^2)}$, $\xi'_H = h_H^r + \sqrt{\log(n)/(nh_H^d)}$, and $\xi_{Hk} = h_H^{r+1} + \sqrt{\log(n)/(nh_H^d)}$ for k = 1, 2. Let S_Z be a compact set range of $\{(z^1, z^2) : z^1 = \zeta_1(x^1) \text{ and } z^2 = \zeta_2(x^2) \text{ for some } (x^1, x^2) \in S_X\}$, and S_{Z^k} be a compact set range of $\{z^k : z^k = \zeta_k(x^k) \text{ for some } x^k \in S_{X^k}\}$ for k = 1, 2.

S.2 Technical Lemmas

We state and show in this section the lemmas used to prove the theorems in the text.

S.2.1 Lemma S.1

Lemma S.1 modifies Lemma 3.1 of Powell, Stock, and Stoker (1989) and Lemma 5 of Horowitz (1998). It provides sufficient conditions for approximation error of U-statistic projection other than $o_p(1/\sqrt{n})$. In particular, it degenerates to the case of Lemma 3.1 of Powell, Stock, and Stoker (1989) when $\lambda_n = n$. Denote $U_n = 2 \cdot [n(n-1)]^{-1} \sum_{i=1}^{n-1} \sum_{j=i+1}^n q_n(W_i, W_j)$ and $\hat{U}_n = E[q_n(W_1, W_2)] + (2/n) \sum_{i=1}^n (E[q_n(W_i, W_j)|W_i] - E[q_n(W_1, W_2)])$. As a matter of fact, Lemma S.1 can be further modified as $\hat{U}_n - U_n = O_p \left[n^{-1} \cdot \sqrt{E[q_n(W_1, W_2)^2]} \right]$ by its proof.

Lemma S.1. Suppose that $\{W_i\}_{i=1}^n$ is a sequence of independently and identically distributed random variables

or vectors. Let $q_n(\cdot, \cdot)$ be a symmetric function, and λ_n be a sequence of positive scalars. If $E[q_n(W_1, W_2)^2] = o(\lambda_n)$, then $\hat{U}_n - U_n = o_p [\sqrt{\lambda_n}/n]$.

Proof. Follow the same idea as the proof of Lemma 3.1 of Powell, Stock, and Stoker (1989) to get $E(\hat{U}_n - U_n)^2 = O\left[n^{-2} \cdot E\left[q_n(W_1, W_2)^2\right]\right]$. Thus $(n^2/\lambda_n) \cdot E(\hat{U}_n - U_n)^2 = O\left[(n^2/\lambda_n) \cdot n^{-2} \cdot E\left[q_n(W_1, W_2)^2\right]\right] = o(1)$. The desired conclusion therefore holds by Markov's inequality. \Box

S.2.2 Lemma S.2

Lemma S.2 finds the uniform convergence rate and asymptotic representation of the nonparametric regression estimator $\hat{H}(\cdot)$. We give a proof of Lemma S.2 for completeness.

Lemma S.2. Let Assumptions A.1-A.5 hold, and the bandwidth h_H satisfy (i) $h_H \rightarrow 0$ and (ii) $\log(n)/(nh_H^d) \rightarrow 0$ as $n \rightarrow \infty$. Then

$$\sup_{x\in\mathcal{S}_X}\left|\widehat{H}(x)-H(x)\right|=O\big(\xi_H\big)$$

in probability as $n \to \infty$ *. Moreover, the asymptotic representation of* $\hat{H}(x) - H(x)$ *is given by*

$$\begin{aligned} \hat{H}(x) &- H(x) \\ = \frac{1}{nh_{H}^{d}} \sum_{i=1}^{n} \left(Y_{i} - \mu(X_{i} - x)'\alpha(x) \right) \left\{ e_{1}'S_{n,r}(x)^{-1}\mu\left(\frac{X_{i} - x}{h_{H}}\right) \right\} \\ & \cdot K_{1}\left(\frac{x^{1} - X_{i}^{1}}{h_{H}}\right) K_{2}\left(\frac{x^{2} - X_{i}^{2}}{h_{H}}\right) + O(\xi_{H}^{2}) \end{aligned}$$

as $n \to \infty$ in probability uniformly over $x \in S_X$, where $e_1 = (1, 0, \dots, 0)'$ is a $N_r \times 1$ vector, $\mu(X_i - x)'\alpha(x)$ represents the r-th order Taylor expansion of $H(X_i)$ at $X_i = x$. $S_{n,r}(x)$, $\mu(\cdot)$ and $\alpha(x)$ are defined in Appendix *S*.1.1.

Proof. The first part can be established by an argument similar to the proof of Theorem 6 of Masry (1996). Its proof is hence omitted here. According to the uniform bahadur representation in Remark 1 of Theorem 3.2 in Kong, Linton, and Xia (2010),

$$\begin{aligned} \widehat{H}(x) - H(x) &= \frac{1}{nh_{H}^{d}} e_{1}^{\prime} S_{n,r}(x)^{-1} B_{H}^{-1} \sum_{i=1}^{n} K\left(\frac{x - X_{i}}{h_{H}}\right) \left(Y_{i} - \mu(X_{i} - x)^{\prime} \alpha(x)\right) \mu(X_{i} - x) \\ &+ O\left(\frac{\log n}{nh_{H}^{d}}\right) \\ &= \frac{1}{nh_{H}^{d}} e_{1}^{\prime} S_{n,r}(x)^{-1} B_{H}^{-1} \sum_{i=1}^{n} K\left(\frac{x - X_{i}}{h_{H}}\right) \left(Y_{i} - \mu(X_{i} - x)^{\prime} \alpha(x)\right) \mu(X_{i} - x) \\ &+ O\left(\xi_{H}^{2}\right) \end{aligned}$$

as $n \to \infty$ in probability uniformly over $x \in S_X$, where $e_1 = (1, 0, \dots, 0)'(N_r - 1 \text{ copies of } 0)$, B_H is the diagonal matrix with diagonal vector $b_H = (b'_{H,s})'_{s=0,1,\dots,r}$ and $b_{H,s} = (h_H^{|\pi_s(k)|})_{k=1,2,\dots,M_s}$. By simplifying this equation, it establishes the second part and hence completes the whole proof.

S.2.3 Lemma S.3

Lemma S.3 shows the large sample properties of the estimators of partial integrations $\zeta_k(\cdot)$ for k = 1, 2. It establishes the uniform convergence rate and asymptotic representation of the estimators $\hat{\zeta}_k(\cdot)$'s for k = 1, 2. In particular, the asymptotic representation decomposes the difference between the estimator and true value of $\zeta_k(\cdot)$ (i.e. $\hat{\zeta}_k - \zeta_k$) into a weighted sum of i.i.d. quantities (with a mean of 0) and a bias term $h_H^r D_k(x^k)$ for k = 1, 2 up to some higher order error.

Lemma S.3. Let Assumptions A.1-A.6 hold. Then for any $k = 1, 2, as n \to \infty$, (i) $\sup_{x^k \in S_{X^k}} |\widehat{\zeta}_k(x^k) - \zeta_k(x^k)| = O(\xi_{Hk})$ in probability with $\xi_{Hk} = h_H^{r+1} + \sqrt{\log(n)/(nh_H^{d_k})}$. (ii) Moreover, for any $x^k \in S_{X^k}$, $\widehat{\zeta}_k(x^k) - \zeta_k(x^k)$ can be written as

$$\widehat{\zeta}_k(x^k) - \zeta_k(x^k) = J_{nk}(x^k) - E[J_{nk}(x^k)] + h_H^{r+1} \cdot D_k(x^k) + o_p(h_H^{r+1}),$$

where $J_{nk}(x^k)$ and $D_k(x^k)$ are defined respectively by (S.1.6) and (S.1.7).

Proof. Only $\zeta_1(\cdot)$ part is shown here. The $\zeta_2(\cdot)$ part can be shown similarly. Let $W = (Y, X^1, X^2)$. Apply Lemma 1 of Horowitz (1998) (or Theorem 2.37 of Pollard (1984)) and Lemma S.2 to obtain

$$\begin{split} \tilde{\zeta}_{1}(x^{1}) &- \zeta_{1}(x^{1}) \\ &= \frac{1}{n} \sum_{j=1}^{n} \left[\widehat{H}_{-j}(x^{1}, X_{j}^{2}) - H(x^{1}, X_{j}^{2}) \right] + o\left(\frac{\log(n)}{\sqrt{n}}\right) \\ &= \frac{1}{n(n-1)} \sum_{j=1}^{n} \sum_{i \neq j} \frac{1}{h_{H}^{d}} \left(Y_{i} - \mu(X_{i}^{1} - x^{1}, X_{i}^{2} - X_{j}^{2})' \alpha(x^{1}, X_{j}^{2}) \right) K_{1} \left(\frac{x^{1} - X_{i}^{1}}{h_{H}} \right) K_{2} \left(\frac{X_{j}^{2} - X_{i}^{2}}{h_{H}} \right) \\ &\cdot \left(e_{1}' S_{n,r}(x^{1}, X_{j}^{2})^{-1} \mu\left(\frac{X_{i}^{1} - x^{1}}{h_{H}}, \frac{X_{i}^{2} - X_{j}^{2}}{h_{H}} \right) \right) + O(\xi_{H}^{2}) + o\left(\frac{\log(n)}{\sqrt{n}} \right) \\ &= \frac{1}{n(n-1)} \sum_{j=1}^{n} \sum_{i \neq j} \psi_{1}(W_{i}, W_{j}) + \frac{1}{n(n-1)} \sum_{j=1}^{n} \sum_{i \neq j} \psi_{2}(W_{i}, W_{j}) + O(\xi_{H}^{2}) + o\left(\frac{\log(n)}{\sqrt{n}} \right) \\ &= \frac{1}{n} \sum_{i=1}^{n} E[\psi_{1}(W_{i}, W_{j})|W_{i}] + \frac{1}{n(n-1)} \sum_{i=1}^{n} \sum_{j \neq i} \left(\psi_{1}(W_{i}, W_{j}) - E[\psi_{1}(W_{i}, W_{j})|W_{i}] \right) \\ &+ \frac{1}{n} \sum_{j=1}^{n} E[\psi_{2}(W_{i}, W_{j})|W_{j}] + \frac{1}{n(n-1)} \sum_{i=1}^{n} \sum_{j \neq i} \left(\psi_{2}(W_{i}, W_{j}) - E[\psi_{2}(W_{i}, W_{j})|W_{j}] \right) \\ &+ O(\xi_{H}^{2}) + o\left(\frac{\log(n)}{\sqrt{n}} \right) \\ &=: T_{1n} + T_{2n} + T_{3n} + T_{4n} + O(\xi_{H}^{2}) + o\left(\frac{\log(n)}{\sqrt{n}} \right) \end{split}$$
(S.2.1)

as $n \to \infty$ in probability uniformly over $x^1 \in S_{X^1}$, where $\hat{H}_{-j}(x^1, X_j^2)$ is a leave-one-out local polynomial estimator and

$$\psi_1(W_i, W_j) = \frac{1}{h_H^d} \Big(Y_i - H(X_i) \Big) K_1 \Big(\frac{x^1 - X_i^1}{h_H} \Big) K_2 \Big(\frac{X_j^2 - X_i^2}{h_H} \Big) \\ \cdot \Big(e_1' S_{n,r}(x^1, X_j^2)^{-1} \mu \Big(\frac{X_i^1 - x^1}{h_H}, \frac{X_i^2 - X_j^2}{h_H} \Big) \Big),$$

$$\begin{split} \psi_2\big(W_i, W_j\big) &= \frac{1}{h_H^d} \Big(H(X_i) - \mu \big(X_i^1 - x^1, X_i^2 - X_j^2\big)' \alpha(x^1, X_j^2) \Big) K_1\Big(\frac{x^1 - X_i^1}{h_H}\Big) K_2\Big(\frac{X_j^2 - X_i^2}{h_H}\Big) \\ &\cdot \Big(e_1' S_{n,r}(x^1, X_j^2)^{-1} \mu \big(\frac{X_i^1 - x^1}{h_H}, \frac{X_i^2 - X_j^2}{h_H}\big) \Big). \end{split}$$

The rest of proof establishes the asymptotic representation of T_{1n} , T_{2n} , T_{3n} , and T_{4n} . It is accomplished in four steps. The asymptotic representation of T_{1n} characterizes the stochastic leading term, and T_{3n} characterizes the leading bias term.

Step 1. For T_{1n} ,

$$\begin{split} E[\psi_{1}(W_{i},W_{j})|W_{i}] \\ &= \frac{1}{h_{H}^{d_{1}}}K_{1}\Big(\frac{x^{1}-X_{i}^{1}}{h_{H}}\Big)\big(Y_{i}-H(X_{i})\big) \\ &\quad \cdot \int \Big(e_{1}'S_{n,r}(x^{1},X_{j}^{2})^{-1}\mu\Big(\frac{X_{i}^{1}-x^{1}}{h_{H}},\frac{X_{i}^{2}-X_{j}^{2}}{h_{H}}\Big)\Big)\frac{1}{h_{H}^{d_{2}}}K_{2}\Big(\frac{X_{j}^{2}-X_{i}^{2}}{h_{H}}\Big)p_{X^{2}}(X_{j}^{2})dX_{j}^{2} \\ &= \frac{1}{h_{H}^{d_{1}}}K_{1}\Big(\frac{x^{1}-X_{i}^{1}}{h_{H}}\Big)\big(Y_{i}-H(X_{i})\big)e_{1}'\big(S_{r}^{-1}+O(h_{H})\big) \\ &\quad \cdot \int \mu\Big(\frac{X_{i}^{1}-x^{1}}{h_{H}},\frac{X_{i}^{2}-X_{j}^{2}}{h_{H}}\Big)\frac{1}{h_{H}^{d_{1}}}K_{2}\Big(\frac{X_{j}^{2}-X_{i}^{2}}{h_{H}}\Big)\frac{p_{X^{2}}(X_{j}^{2})}{p_{X}(x^{1},X_{j}^{2})}dX_{j}^{2} \\ &= \frac{1}{h_{H}^{d_{1}}}K_{1}\Big(\frac{x^{1}-X_{i}^{1}}{h_{H}}\Big)\frac{Y_{i}-H(X_{i})}{p_{X^{1}|X^{2}}(x^{1}|X_{i}^{2})}e_{1}'S_{r}^{-1}V_{1}^{\mu}\Big(\frac{x^{1}-X_{i}^{1}}{h_{H}}\Big)+O(h_{H}\xi_{H1}), \end{split}$$

where the last second equation is given by change of variable and first order Taylor expansion, and the last equation is based on the proof of Theorem 6 in Masry (1996). Thus, T_{1n} can be written as

$$T_{1n} = \frac{1}{nh_H^{d_1}} \sum_{i=1}^n K_1\left(\frac{x^1 - X_i^1}{h_H}\right) \frac{Y_i - H(X_i)}{p_{X^1|X^2}(x^1|X_i^2)} e_1' S_r^{-1} V_1^{\mu}\left(\frac{x^1 - X_i^1}{h_H}\right) + O(h_H \xi_{H1}).$$

Step 2. For T_{2n} , it can be decomposed as

$$T_{2n} = \frac{1}{n(n-1)} \sum_{i=1}^{n} \sum_{j \neq i} \left(\psi_1(W_i, W_j) - E[\psi_1(W_i, W_j) | W_i] \right)$$

= $\frac{1}{nh_H^d} \sum_{i=1}^{n} K_1\left(\frac{x^1 - X_i^1}{h_H}\right) \frac{1}{n-1} \sum_{j \neq i} \left(\tilde{\psi}_1(W_i, W_j) - E[\tilde{\psi}_1(W_i, W_j) | W_i] \right)$
= $o\left(\frac{\log(n)^2}{nh_H^{d/2}}\right),$

where $\tilde{\psi}_1(W_i, W_j) = \frac{1}{h_H^{d_2}} K_2\left(\frac{X_j^2 - X_i^2}{h_H}\right) \left(Y_i - H(x_i)\right) \left(e_1' S_{n,p}(x^1, X_j^2)^{-1} \mu\left(\frac{X_i^1 - x^1}{h_H}, \frac{X_i^2 - X_j^2}{h_H}\right)\right)$, and the last equality is obtained by applying Lemma 1 of Horowitz (1998) (or Theorem 2.37 of Pollard (1984)).¹²

 $^{^{12}\}mathrm{A}$ similar argument is used by Horowitz (1998) to establish its (C.5).

Step 3. For T_{3n} , the summand $E[\psi(W_i, W_j)|W_j]$ can be simplified as

$$\begin{split} E[\psi_{2}(W_{i},W_{j})|W_{j}] \\ = E\left[\frac{1}{h_{H}^{d}}\Big(H(X_{i})-\mu(X_{i}^{1}-x^{1},X_{i}^{2}-X_{j}^{2})'\alpha(x^{1},X_{j}^{2})\Big)K_{1}\Big(\frac{x^{1}-X_{i}^{1}}{h_{H}}\Big)K_{2}\Big(\frac{X_{j}^{2}-X_{i}^{2}}{h_{H}}\Big) \\ \cdot \Big(e_{1}'S_{n,r}(x^{1},X_{j}^{2})^{-1}\mu\Big(\frac{X_{i}^{1}-x^{1}}{h_{H}},\frac{X_{i}^{2}-X_{j}^{2}}{h_{H}}\Big)\Big)\Big|W_{j}\right] \\ = e_{1}'S_{n,r}(x^{1},X_{j}^{2})^{-1}\int\frac{1}{h_{H}^{d}}\Big(H(X_{i})-\mu(X_{i}^{1}-x^{1},X_{i}^{2}-X_{j}^{2})'\alpha(x^{1},X_{j}^{2})\Big) \\ \cdot K_{1}\Big(\frac{x^{1}-X_{i}^{1}}{h_{H}}\Big)K_{2}\Big(\frac{X_{j}^{2}-X_{i}^{2}}{h_{H}}\Big)\mu\Big(\frac{X_{i}^{1}-x^{1}}{h_{H}},\frac{X_{i}^{2}-X_{j}^{2}}{h_{H}}\Big)p_{X}(X_{i})dX_{i} \\ = e_{1}'S_{n,r}(x^{1},X_{j}^{2})^{-1}\sum_{|\mathbf{s}|=r+1}\frac{1}{\mathbf{s}!}D^{\mathbf{s}}H(x^{1},X_{j}^{2})\int u^{\mathbf{s}}\mu(u)p_{X}(X_{j}+h_{H}u)du\cdot h_{H}^{r+1}+o(h_{H}^{r+1}) \\ = e_{1}'S_{n,r}(x^{1},X_{j}^{2})^{-1}\Big(h_{H}^{r+1}S_{n,r}^{r+1}(x^{1},X_{j}^{2})H_{r+1}(x^{1},X_{j}^{2})+o_{p}(h_{H}^{r+1})\Big) \\ = h_{H}^{r+1}e_{1}'S_{r}^{-1}S_{r}^{r+1}H_{r+1}(x^{1},X_{j}^{2})+o_{p}(h_{H}^{r+1}), \end{split}$$
(S.2.2)

where the last second equality is derived by change of variable in the integration and Taylor expansion, and the last equality is due to the approximations $S_{n,r}(x)^{-1} = \{p_X(x)\}^{-1}S_r^{-1} + O(h_H)$ and $S_{n,r}^{r+1}(x) = p_X(x)S_r^{r+1} + O(h_H)$ in the proof of Proposition 3.1 in Kong, Linton, and Xia (2010).¹³ Thus, the weighted sum can be represented as

$$T_{3n} = \frac{1}{n} \sum_{j=1}^{n} E[\psi_2(W_i, W_j) | W_j]$$

= $h_H^{r+1} e_1' S_r^{-1} S_r^{r+1} \left(\frac{1}{n} \sum_{j=1}^{n} H_{r+1}(x^1, X_j^2) p_X(x^1, X_j^2) \right) + o_p(h_H^{r+1})$
= $h_H^{r+1} e_1' S_r^{-1} S_r^{r+1} E[H_{r+1}(x^1, X^2)] + o_p(h_H^{r+1})$

Step 4. For T_{4n} , note that

$$\begin{split} & E\left[\left(\frac{1}{n-1}\sum_{j\neq i}\left[\psi_{2}(W_{i},W_{j})-E\left[\psi_{2}(W_{i},W_{j})|W_{i}\right]\right)^{2}\right]\\ &=\frac{1}{n-1}E\left[\left(\psi_{2}(W_{2},W_{1})-E\left[\psi_{2}(W_{2},W_{1})|W_{1}\right]\right)^{2}\right]\\ &\leq\frac{1}{n-1}E\left[\left(\psi_{2}(W_{2},W_{1})\right)^{2}\right]\\ &=O\left(\frac{h_{H}^{2r+2}}{n}\right), \end{split}$$

where the last equality is obtained by Taylor expansion similar to Step 3. By applying Lemma 1 of

¹³When *r* is even, $e'_1 S_r^{-1} S_r^{r+1} = 0$ and thus the first term on the right hand side of the last equality (S.2.2) vanishes. In this case, the bias term is of order $O(h^{r+2})$ if we further assume that all functions and densities are (r + 2) continuously differentiable.

Horowitz (1998) (or Theorem 2.37 of Pollard (1984)), we derive

$$T_{4n} = \frac{1}{n(n-1)} \sum_{i=1}^{n} \sum_{j \neq i} \left[\psi_2(W_i, W_j) - E\left[\psi_2(W_i, W_j) | W_j\right] \right]$$
$$= \frac{1}{n} \sum_{i=1}^{n} \frac{1}{n-1} \sum_{j \neq i} \left[\psi_2(W_i, W_j) - E\left[\psi_2(W_i, W_j) | W_i\right] \right]$$
$$= o\left(h_H^{r+1} \frac{\log(n)}{n}\right).$$

With the bandwidths satisfying Assumption A.6, combining steps 1-4 yields

$$\begin{split} \widehat{\zeta}_{1}(x^{1}) &- \zeta_{1}(x^{1}) \\ = & T_{1n} + T_{2n} + T_{3n} + T_{4n} + O(\xi_{H}^{2}) + o\left(\frac{\log(n)}{\sqrt{n}}\right) \\ &= \frac{1}{nh_{H}^{d_{1}}} \sum_{i=1}^{n} K_{1}\left(\frac{x^{1} - X_{i}^{1}}{h_{H}}\right) \cdot \frac{Y_{i} - H(X_{i})}{p_{X^{1}|X^{2}}(x^{1}|X_{i}^{2})} e_{1}^{\prime} S_{r}^{-1} V_{1}^{\mu}\left(\frac{x^{1} - X_{i}^{1}}{h_{H}}\right) \\ &+ e_{1}^{\prime} S_{r}^{-1} S_{r}^{r+1} E\left[H_{r+1}(x^{1}, X^{2})\right] h_{H}^{r+1} \\ &+ O(\xi_{H}^{2} + h_{H}\xi_{H1}) + o\left(\frac{\log(n)}{\sqrt{n}} + \frac{\log(n)^{2}}{nh_{H}^{d/2}} + h_{H}^{r+1} + h_{H}^{r+1}\frac{\log(n)}{n}\right) \\ &= J_{n1}(x^{1}) + D_{1}(x^{1})h_{H}^{r+1} + o(h_{H}^{r+1}) \end{split}$$

in probability as $n \to \infty$ uniformly over $x^1 \in S_{X^1}$, where the first term on the right hand side of last equality is given by the definition of $J_{n1}(x^1)$ in (S.1.6), and $E[J_{n1}(x^1)] = 0.^{14}$ The asymptotic representation of $\hat{\zeta}_1(x^1) - \zeta_1(x^1)$ is hence established.

Following an idea similar to the proof of Theorem 6 of Masry (1996), we have

$$\sup_{x^{1}\in\mathcal{S}_{X^{1}}}\left|\frac{1}{nh_{H}^{d_{1}}}\sum_{i=1}^{n}K_{1}\left(\frac{x^{1}-X_{i}^{1}}{h_{H}}\right)\cdot\frac{Y_{i}-H(X_{i})}{p_{X^{1}|X^{2}}(x^{1}|X_{i}^{2})}e_{1}'S_{p}^{-1}V^{\mu}\left(\frac{x^{1}-X_{i}^{1}}{h_{H}}\right)\right|=O\left(\sqrt{\frac{\log(n)}{nh_{H}^{d_{1}}}}\right)$$
(S.2.3)

in probability as $n \to \infty$. Based on the asymptotic representation, this implies that $\sup_{x^1 \in S_{X^1}} |\hat{\zeta}_1(x^1) - \zeta_1(x^1)| = O(h_H^{r+1} + \sqrt{\frac{\log(n)}{nh_H^{d_1}}}) = O(\xi_{H1})$ in probability as $n \to \infty$. This completes the proof. \Box

S.2.4 Lemma S.4

Lemma S.4 characterize the uniform convergence rate and asymptotic representation of the Local linear estimators. There are three terms in the leading part (excluding all higher order remainders) of the difference $\hat{\mathcal{H}}(z) - \mathcal{H}(z)$. The first term in the asymptotic representation is the oracle term with true $\zeta_1(\cdot)$ and $\zeta_2(\cdot)$. The second and third terms represent the error by estimating $\zeta_1(\cdot)$ and $\zeta_2(\cdot)$, respectively. Note that $\xi_{H1} \geq \xi_{H2} > 0$ due to $d_1 \geq d_2$. This implies that $O(\xi_{H1} + \xi_{H2}) = O(\xi_{H1})$.

 $[\]frac{1}{1^{4}\text{It is easy to obtain } E[J_{n1}(x^{1})] = E\left[\frac{1}{nh_{H}^{d_{1}}}K_{1}\left(\frac{x^{1}-X_{i}^{1}}{h_{H}}\right) \cdot \frac{1}{p_{X^{1}|X^{2}}(x^{1}|X_{i}^{2})}e_{1}'S_{r}^{-1}V_{1}^{\mu}\left(\frac{x^{1}-X_{i}^{1}}{h_{H}}\right)E[Y_{i}-H(X_{i})|X_{i}]\right] = 0 \text{ by Law of iterative expectations.}$

Lemma S.4. Suppose that Assumptions A.1-A.6 hold. Then

$$\sup_{z \in \mathcal{S}_Z} |B_{\mathcal{H}}(\widehat{\beta}(z) - \beta(z))| = O(\xi_{\mathcal{H}} + \xi_{H1})$$

in probability as $n \to \infty$. Moreover, the asymptotic representation of $\hat{\beta}(z) - \beta(z)$ is given by

$$\begin{split} & B_{\mathcal{H}}\big(\widehat{\beta}(z) - \beta(z)\big) \\ = & \frac{1}{nh_{\mathcal{H}}^2} \sum_{i=1}^n \widetilde{S}_{n,r}(z)^{-1} k\big(\frac{z^1 - \zeta_1(X_i^1)}{h_{\mathcal{H}}}\big) k\big(\frac{z^2 - \zeta_2(X_i^2)}{h_{\mathcal{H}}}\big) \Big\{Y_i - \widetilde{\mu}\big(\zeta(X_i) - z\big)'\beta(z)\Big\} \mu\Big(\frac{\zeta(X_i) - z}{h_{\mathcal{H}}}\Big) \\ & + \frac{1}{nh_{\mathcal{H}}^3} \sum_{i=1}^n \widetilde{S}_{n,r}(z)^{-1} \bigg[\Big(\frac{\partial}{\partial u^1} t(u, Y_i; z)\widetilde{K}(u) + t(u, Y_i; z)\partial_1\widetilde{K}(u)\Big)\Big|_{u = \frac{\zeta(X_i) - z}{h_{\mathcal{H}}}}\bigg] \big(\widehat{\zeta}_1(X_i^1) - \zeta_1(X_i^1)\big) \\ & + \frac{1}{nh_{\mathcal{H}}^3} \sum_{i=1}^n \widetilde{S}_{n,r}(z)^{-1} \bigg[\Big(\frac{\partial}{\partial u^2} t(u, Y_i; z)\widetilde{K}(u) + t(u, Y_i; z)\partial_2\widetilde{K}(u)\Big)\Big|_{u = \frac{\zeta(X_i) - z}{h_{\mathcal{H}}}}\bigg] \big(\widehat{\zeta}_2(X_i^2) - \zeta_2(X_i^2)\big) \\ & + O\big(\zeta_{\mathcal{H}}^2 + \zeta_{\mathcal{H}1}^2\big) \end{split}$$

as $n \to \infty$ in probability uniformly over $z \in S_Z$, where $\hat{\beta}(z)$ is the r-th order local polynomial estimator of true value $\beta(z)$, $u = (u^1, u^2)$, and $t(u, Y_i; z) = \tilde{\mu}(u)(Y_i - \tilde{\mu}(u)'B_H\beta(z))$. B_H , $\tilde{S}_{n,r}(z)$ and $\tilde{\mu}(u)$ are defined in Appendix S.1.1.

Proof. Note that

$$B_{\mathcal{H}}(\widehat{\beta}(z) - \beta(z)) = B_{\mathcal{H}}(\widehat{\beta}(z) - \widetilde{\beta}(z)) + B_{\mathcal{H}}(\widetilde{\beta}(z) - \beta(z)).$$

First we consider $B_{\mathcal{H}}(\widehat{\beta}(z) - \widetilde{\beta}(z))$.

$$B_{\mathcal{H}}(\widehat{\beta}(z) - \widetilde{\beta}(z)) = [\mathcal{S}_{n,r}(z,\widehat{\zeta})^{-1} - \mathcal{S}_{n,r}(z,\zeta)^{-1}]\mathcal{Q}_{n,r}(z,\zeta) + \mathcal{S}_{n,r}(z,\zeta)^{-1}[\mathcal{Q}_{n,r}(z,\widehat{\zeta}) - \mathcal{Q}_{n}(z,\zeta)] + [\mathcal{S}_{n,r}(z,\widehat{\zeta})^{-1} - \mathcal{S}_{n,r}(z,\zeta)^{-1}] \cdot [\mathcal{Q}_{n,r}(z,\widehat{\zeta}) - \mathcal{Q}_{n,r}(z,\zeta)].$$
(S.2.4)

 $S_{n,r}(\cdot)$ and $Q_{n,r}(\cdot)$ are defined in Appendix S.1.1. As for $Q_{n,r}(z, \hat{\zeta}) - Q_{n,r}(z, \zeta)$, we apply Taylor expansion. By Lemma 1 of Horowitz (1998) (or Theorem 2.37 of Pollard (1984)) and our Lemma S.3,

$$\mathcal{Q}_{n,r}(z,\hat{\zeta}) - \mathcal{Q}_{n,r}(z,\zeta) = \frac{1}{nh_{\mathcal{H}}^3} \sum_{i=1}^n \left\{ D_{\zeta} \mathcal{Q}_{in}^1(z,\zeta) \big(\hat{\zeta}_1(X_i^1) - \zeta_1(X_i^1) \big) + D_{\zeta} \mathcal{Q}_{in}^2(z,\zeta) \big(\hat{\zeta}_2(X_i^2) - \zeta_2(X_i^2) \big) \right\} + O(\xi_{H1}^2)$$
(S.2.5)

and $Q_{n,r}(z,\hat{\zeta}) - Q_{n,r}(z,\zeta) = O(\xi_{H1})$ in probability as $n \to \infty$ uniformly over $z \in S_Z$, where $D_{\zeta}Q_{in}^1(z,\zeta)$ is a $\tilde{N}_r \times 1$ vector with

$$\begin{split} & \left[D_{\zeta} \mathcal{Q}_{in}^{1}(z,\zeta) \right]_{l} = \\ & \left\{ \begin{split} & Y_{i} \big(\frac{\zeta_{2}(X_{i}^{2}) - z^{2}}{h_{\mathcal{H}}} \big)^{r_{2}} \partial_{1} \tilde{K} \big(\frac{z - \zeta(X_{i})}{h_{\mathcal{H}}} \big), & r_{1} = 0 \\ & Y_{i} \big(\frac{\zeta_{2}(X_{i}^{2}) - z^{2}}{h_{\mathcal{H}}} \big)^{r_{2}} \Big[\partial_{1} \tilde{K} \big(\frac{z - \zeta(X_{i})}{h_{\mathcal{H}}} \big) \big(\frac{\zeta_{1}(X_{i}^{1}) - z^{1}}{h_{\mathcal{H}}} \big)^{r_{1}} + r_{1} \tilde{K} \big(\frac{z - \zeta(X_{i})}{h_{\mathcal{H}}} \big) \big(\frac{\zeta_{1}(X_{i}^{1}) - z^{1}}{h_{\mathcal{H}}} \big)^{r_{1} - 1} \Big], & r_{1} \ge 1 \end{split} \right.$$

where $\tilde{r} = (r_1, r_2)$ is the correspondent power numbers of the *l*-th entry of $Q_{n,r}(z, \zeta)$, i.e.

$$\left[\mathcal{Q}_{n,r}(z,\zeta)\right]_{l}=\frac{1}{nh_{\mathcal{H}}^{2}}\sum_{i=1}^{n}\left(\frac{z-\zeta(X_{i})}{h_{\mathcal{H}}}\right)^{\tilde{r}}\tilde{K}\left(\frac{z-\zeta(X_{i})}{h_{\mathcal{H}}}\right).$$

Similarly, we can define $D_{\zeta} Q_{in}^2(z, \zeta)$.

As for $S_n(z, \hat{\zeta})^{-1} - S_n(z, \zeta)^{-1}$, Similarly, we can derive that

$$S_{n}(z,\hat{\zeta}) - S_{n}(z,\zeta) = \frac{1}{nh_{\mathcal{H}}^{3}} \sum_{i=1}^{n} \left\{ D_{\zeta}S_{in}^{1}(z,\zeta) \left(\widehat{\zeta}_{1}(X_{i}^{1}) - \zeta_{1}(X_{i}^{1}) \right) + D_{\zeta}S_{in}^{2}(z,\zeta) \left(\widehat{\zeta}_{2}(X_{i}) - \zeta_{2}(X_{i}) \right) \right\} + O(\xi_{H1}^{2}),$$

and $S_n(z, \hat{\zeta}) - S_n(z, \zeta) = O(\xi_{H1})$ as $n \to \infty$ in probability uniformly over $z \in S_Z$, where matrix $D_{\zeta}S_{in}^v(z,\zeta)$ (v = 1,2) satisfies that its (l,k)-entry $(l,k = 1,2,...,N_r)$ is

$$\begin{split} &\left[D_{\zeta}\mathcal{S}_{in}^{1}(z,\zeta)\right]_{lk} = \\ &\left\{ \begin{pmatrix} \frac{\zeta_{2}(X_{i}^{2})-z^{2}}{h_{\mathcal{H}}} \end{pmatrix}^{r_{2}} \partial_{1}\tilde{K} \begin{pmatrix} \frac{z-\zeta(X_{i})}{h_{\mathcal{H}}} \end{pmatrix}, & r_{1} = 0 \\ & \begin{pmatrix} \frac{\zeta_{2}(X_{i}^{2})-z^{2}}{h_{\mathcal{H}}} \end{pmatrix}^{r_{2}} \begin{bmatrix} \partial_{1}\tilde{K} \begin{pmatrix} \frac{z-\zeta(X_{i})}{h_{\mathcal{H}}} \end{pmatrix} \begin{pmatrix} \frac{\zeta_{1}(X_{i}^{1})-z^{1}}{h_{\mathcal{H}}} \end{pmatrix}^{r_{1}} + r_{1}\tilde{K} \begin{pmatrix} \frac{z-\zeta(X_{i})}{h_{\mathcal{H}}} \end{pmatrix} \begin{pmatrix} \frac{\zeta_{1}(X_{i}^{1})-z^{1}}{h_{\mathcal{H}}} \end{pmatrix}^{r_{1}-1} \end{bmatrix}, & r_{1} \geq 1 \end{split} \right], \end{split}$$

where $\tilde{r} = (r_1, r_2)$ is the power number of the (l, k)-element of $S_n(z, \zeta)$. Similarly, we can define $D_{\zeta}S_{in}^2(z,\zeta)$. Similar to the arguments in the proof of Theorem 3.2 in Kong, Linton, and Xia (2010), we have $\sup_{z \in S_Z} |S_n(z,\zeta) - \tilde{S}_{n,r}(z)| = O(\xi_{\mathcal{H}})$ as $n \to \infty$ in probability. Thus, the triangular inequality implies that

$$S_n(z,\widehat{\zeta}) - \tilde{S}_{n,r}(z) = O(\xi_{H1} + \xi_{\mathcal{H}})$$

as $n \to \infty$ in probability uniformly over $z \in S_Z$. Therefore, we can derive that

$$S_{n}(z,\hat{\zeta})^{-1} - S_{n}(z,\zeta)^{-1}$$

$$= -S_{n}(z,\hat{\zeta})^{-1} \Big(S_{n}(z,\hat{\zeta}) - S_{n}(z,\zeta) \Big) S_{n}(z,\zeta)^{-1}$$

$$= -\frac{1}{nh_{\mathcal{H}}^{3}} \sum_{i=1}^{n} \tilde{S}_{n,r}(z)^{-1} \Big\{ D_{\zeta} S_{in}^{1}(z,\zeta) \big(\hat{\zeta}_{1}(X_{i}^{1}) - \zeta_{1}(X_{i}^{1}) \big) + D_{\zeta} S_{in}^{2}(z,\zeta) \big(\hat{\zeta}_{2}(X_{i}^{2}) - \zeta_{2}(X_{i}^{2}) \big) \Big\} \tilde{S}_{n,r}(z)^{-1}$$

$$+ O\big(\xi_{\mathcal{H}} \cdot \xi_{H1} + \xi_{H1}^{2} \big)$$
(S.2.6)

and $S_n(z, \hat{\zeta})^{-1} - S_n(z, \zeta)^{-1} = O(\xi_{H1})$ as $n \to \infty$ in probability uniformly over $z \in S_Z$. Also, we have $Q_n(z, \zeta) = S_n(z, \zeta) B_H \tilde{\beta}(z)$. By Theorem 6 in Masry (1996), $\sup_{z \in S_Z} |B_H(\tilde{\beta}(z) - \beta(z))| = O(\xi_H)$ in probability as $n \to \infty$. Therefore, we have

$$\mathcal{Q}_n(z,\zeta) - \tilde{S}_{n,r}(z)B_{\mathcal{H}}\beta(z) = O(\xi_{\mathcal{H}})$$
(S.2.7)

as $n \to \infty$ in probability uniformly over $z \in S_Z$. According to (S.2.5), (S.2.6) and (S.2.7), (S.2.4) can be rewritten as

$$B_{\mathcal{H}}(\widehat{\beta}(z) - \widetilde{\beta}(z))$$

$$= [\mathcal{S}_{n}(z,\widehat{\zeta})^{-1} - \mathcal{S}_{n}(z,\zeta)^{-1}](\widetilde{S}_{n,r}(z)B_{\mathcal{H}}\beta(z) + O(\xi_{\mathcal{H}})) + (\widetilde{S}_{n,r}(z)^{-1} + O(\xi_{\mathcal{H}}))[\mathcal{Q}_{n}(z,\widehat{\zeta}) - \mathcal{Q}_{n}(z,\zeta)]$$

$$+ O(\xi_{\mathcal{H}1}^{2})$$

$$= \frac{1}{nh_{\mathcal{H}}^{3}} \sum_{i=1}^{n} \tilde{S}_{n,r}(z)^{-1} \left[-D_{\zeta} \mathcal{S}_{in}^{1}(z,\zeta) B_{\mathcal{H}} \beta(z) + D_{\zeta} \mathcal{Q}_{in}^{1}(z,\zeta) \right] \left(\hat{\zeta}_{1}(X_{i}^{1}) - \zeta_{1}(X_{i}^{1}) \right) \\ + \frac{1}{nh_{\mathcal{H}}^{3}} \sum_{i=1}^{n} \tilde{S}_{n,r}(z)^{-1} \left[-D_{\zeta} \mathcal{S}_{in}^{2}(z,\zeta) B_{\mathcal{H}} \beta(z) + D_{\zeta} \mathcal{Q}_{in}^{2}(z,\zeta) \right] \left(\hat{\zeta}_{2}(X_{i}^{2}) - \zeta_{2}(X_{i}^{2}) \right) \\ + O(\xi_{\mathcal{H}} \cdot \xi_{H1} + \xi_{H1}^{2}) \\ = \frac{1}{nh_{\mathcal{H}}^{3}} \sum_{i=1}^{n} \tilde{S}_{n,r}(z)^{-1} \left[\left(\frac{\partial}{\partial u^{1}} t(u, Y_{i}; z) \tilde{K}(u) + t(u, Y_{i}; z) \partial_{1} \tilde{K}(u) \right) \right|_{u = \frac{\zeta(X_{i}) - z}{h_{\mathcal{H}}}} \right] \left(\hat{\zeta}_{1}(X_{i}^{1}) - \zeta_{1}(X_{i}^{1}) \right) \\ = \frac{1}{nh_{\mathcal{H}}^{3}} \sum_{i=1}^{n} \tilde{S}_{n,r}(z)^{-1} \left[\left(\frac{\partial}{\partial u^{2}} t(u, Y_{i}; z) \tilde{K}(u) + t(u, Y_{i}; z) \partial_{2} \tilde{K}(u) \right) \right|_{u = \frac{\zeta(X_{i}) - z}{h_{\mathcal{H}}}} \right] \left(\hat{\zeta}_{2}(X_{i}^{2}) - \zeta_{2}(X_{i}^{2}) \right) \\ + O(\xi_{\mathcal{H}}^{2} + \xi_{H1}^{2})$$
(S.2.8)

and $B_{\mathcal{H}}(\hat{\beta}(z) - \tilde{\beta}(z)) = O(\xi_{H1})$ in probability as $n \to \infty$ uniformly over $z \in S_Z$. Second we consider $B_{\mathcal{H}}(\tilde{\beta}(z) - \beta(z))$. The asymptotic linear representation is a direct application of Theorem 3.2 in Kong, Linton, and Xia (2010), that is,

$$B_{\mathcal{H}}(\tilde{\beta}(z) - \beta(z)) = \frac{1}{nh_{\mathcal{H}}^2} \sum_{i=1}^n \tilde{S}_{n,r}(z)^{-1} k \Big(\frac{z^1 - \zeta_1(X_i^1)}{h_{\mathcal{H}}}\Big) k \Big(\frac{z^2 - \zeta_2(X_i^2)}{h_{\mathcal{H}}}\Big) \Big\{ Y_i - \tilde{\mu} \big(\zeta(X_i) - z\big)' \beta(z) \Big\} \mu \Big(\frac{\zeta(X_i) - z}{h_{\mathcal{H}}}\Big) + O(\xi_{\mathcal{H}}^2)$$
(S.2.9)

and $B_{\mathcal{H}}(\hat{\beta}(z) - \beta(z)) = O(\xi_{\mathcal{H}})$ in probability as $n \to \infty$ uniformly over $z \in S_Z$. Finally, the desired representation of $B_{\mathcal{H}}(\hat{\beta}(\cdot) - \beta(\cdot))$ can then be established by (S.2.8) and (S.2.9).

S.2.5 Lemma S.5

 $\partial_k \hat{\mathcal{H}}(z)$ is the *r*-th order local polynomial estimator of first derivatives $\partial_k \mathcal{H}(z)$ (k = 1, 2) based on a data $\{Y_i, \hat{\zeta}_1(X_i), \hat{\zeta}_2(X_i)\}_{i=1}^n$, while $\partial_k \tilde{\mathcal{H}}(z)$ is the infeasible version with a data $\{Y_i, \zeta_1(X_i), \zeta_2(X_i)\}_{i=1}^n$. Lemma S.5 studies the asymptotic properties of $\partial_k \hat{\mathcal{H}}(z)$. It shows the uniform convergence and asymptotic representation of such statistics. Particularly, the first two terms in the asymptotic representation come from the (asymptotic) representation of infeasible estimator $\partial_k \tilde{\mathcal{H}}(z)$, while the third term is the additional bias appearing in the difference between feasible and infeasible estimators, namely $\partial_k \hat{\mathcal{H}}(z) - \partial_k \tilde{\mathcal{H}}(z)$.

Lemma S.5. Suppose that Assumptions A.1-A.6 hold. Then for k = 1, 2, (i) $\sup_{z \in S_Z} |\partial_k \hat{\mathcal{H}}(z) - \partial_k \mathcal{H}(z)| = O(\xi'_{\mathcal{H}} + \xi_{H1})$ in probability as $n \to \infty$; (ii) $\partial_k \hat{\mathcal{H}}(z) - \partial_k \mathcal{H}(z)$ has an asymptotic representation as

$$\begin{bmatrix} \widehat{\partial_1 \mathcal{H}}(z) \\ \widehat{\partial_2 \mathcal{H}}(z) \end{bmatrix} - \begin{bmatrix} \partial_1 \mathcal{H}(z) \\ \partial_2 \mathcal{H}(z) \end{bmatrix} = \frac{1}{nh_{\mathcal{H}}^3} \sum_{i=1}^n e'_d \tilde{S}_{n,r}(z)^{-1} \tilde{K} \Big(\frac{z - \zeta(X_i)}{h_{\mathcal{H}}} \Big) \Big(Y_i - H(X_i) \Big) \mu \Big(\frac{\zeta(X_i) - z}{h_{\mathcal{H}}} \Big) + \mathfrak{D}(z) h_{\mathcal{H}}^r + \mathcal{D}(z) h_{\mathcal{H}}^{r+1} + \mathcal{D}(z) h_{\mathcal{H}}^{r+1} + \mathcal{D}(z) h_{\mathcal{H}}^r + \mathcal{D}(z) h_{\mathcal{H}}^r$$

in probability as $n \to \infty$ uniformly over $z \in S_Z$.

Proof. Note that

$$\begin{split} \begin{bmatrix} \widehat{\partial_1 \mathcal{H}}(z) \\ \widehat{\partial_2 \mathcal{H}}(z) \end{bmatrix} &- \begin{bmatrix} \partial_1 \mathcal{H}(z) \\ \partial_2 \mathcal{H}(z) \end{bmatrix} = e'_d \big(\widehat{\beta}(z) - \beta(z) \big) = e'_d B_{\mathcal{H}}^{-1} \cdot B_{\mathcal{H}} \big(\widehat{\beta}(z) - \beta(z) \big) \\ &= \frac{1}{h_{\mathcal{H}}} e'_d \cdot B_{\mathcal{H}} \big(\widehat{\beta}(z) - \beta(z) \big). \end{split}$$

Thus, part (i) is trivial by Lemma S.5. To find the asymptotic representation, we just need to further decompose $\hat{\beta}(z) - \beta(z)$. According to Lemma S.4, we just need to derive the asymptotic representation of following three parts,

$$\begin{split} A_{1n}(z) &= \frac{1}{nh_{\mathcal{H}}^3} \sum_{i=1}^n e'_d \tilde{S}_{n,r}(z)^{-1} \tilde{K} \Big(\frac{z - \zeta(X_i)}{h_{\mathcal{H}}} \Big) \Big\{ Y_i - \tilde{\mu} \big(z - \zeta(X_i) \big)' \beta(z) \Big\} \mu \Big(\frac{z - \zeta(X_i)}{h_{\mathcal{H}}} \Big); \\ A_{2n}(z) &= \frac{1}{nh_{\mathcal{H}}^4} \sum_{i=1}^n e'_d \tilde{S}_{n,r}(z)^{-1} \left[\Big(\frac{\partial}{\partial u^1} t(u, Y_i; z) \tilde{K}(u) + t(u, Y_i; z) \partial_1 \tilde{K}(u) \Big) \Big|_{u = \frac{\zeta(X_i) - z}{h_{\mathcal{H}}}} \right] (\widehat{\zeta}_1(X_i^1) - \zeta_1(X_i^1)); \\ A_{3n}(z) &= \frac{1}{nh_{\mathcal{H}}^4} \sum_{i=1}^n e'_d \tilde{S}_{n,r}(z)^{-1} \left[\Big(\frac{\partial}{\partial u^2} t(u, Y_i; z) \tilde{K}(u) + t(u, Y_i; z) \partial_2 \tilde{K}(u) \Big) \Big|_{u = \frac{\zeta(X_i) - z}{h_{\mathcal{H}}}} \right] (\widehat{\zeta}_2(X_i^2) - \zeta_2(X_i^2)). \end{split}$$

First we consider $A_{1n}(z)$,

$$A_{1n}(z) = \frac{1}{nh_{\mathcal{H}}^3} \sum_{i=1}^n e'_d \tilde{S}_{n,r}(z)^{-1} \tilde{K} \Big(\frac{z - \zeta(X_i)}{h_{\mathcal{H}}} \Big) \Big\{ \mathcal{H} \big(\zeta(X_i) \big) - \tilde{\mu} \big(z - \zeta(X_i) \big)' \beta(z) \Big\} \mu \Big(\frac{z - \zeta(X_i)}{h_{\mathcal{H}}} \Big) \\ + \frac{1}{nh_{\mathcal{H}}^3} \sum_{i=1}^n e'_d \tilde{S}_{n,r}(z)^{-1} \tilde{K} \Big(\frac{z - \zeta(X_i)}{h_{\mathcal{H}}} \Big) \Big\{ Y_i - \mathcal{H} \big(\zeta(X_i) \big) \Big\} \mu \Big(\frac{z - \zeta(X_i)}{h_{\mathcal{H}}} \Big) \\ =: A_{11n}(z) + A_{12n}(z).$$
(S.2.10)

As for $A_{11n}(z)$, note that

$$\begin{split} & E\left[A_{11n}(z)\right] \\ &= \frac{1}{h_{\mathcal{H}}^3} e'_d B_{\mathcal{H}}^{-1} \tilde{S}_{n,r}(z)^{-1} \int \tilde{K}\left(\frac{z-Z}{h_{\mathcal{H}}}\right) \Big\{ \mathcal{H}(Z) - \tilde{\mu}(z-Z)'\beta(z) \Big\} \mu\left(\frac{z-Z}{h_{\mathcal{H}}}\right) p_Z(Z) dZ \\ &= e'_d \tilde{S}_{n,r}(z)^{-1} \sum_{|\mathbf{j}|=r+1} \frac{1}{\mathbf{j}!} D^{\mathbf{j}} \mathcal{H}(z) \int u^{\mathbf{j}} \tilde{\mu}(u) \tilde{K}(u) p_Z(z+h_{\mathcal{H}}u) du \cdot h_{\mathcal{H}}^r + o(h^r) \\ &= e'_d \tilde{S}_{n,r}(z)^{-1} \tilde{S}_{n,r}^{r+1}(z) \mathcal{H}_{r+1}(z) \cdot h_{\mathcal{H}}^r + o(h_{\mathcal{H}}^r) \\ &= e'_d \tilde{S}_r^{-1} \tilde{S}_r^{r+1} \mathcal{H}_{r+1}(z) \cdot h_{\mathcal{H}}^r + o(h_{\mathcal{H}}^r), \end{split}$$

where the second equality is derived by change of variables and Taylor expansion, and the last equality is due to the approximations $\tilde{S}_{n,r}(z)^{-1} = \tilde{S}_r^{-1}p_Z(z)^{-1} + O(h_H)$ and $\tilde{S}_{n,r}^{r+1}(z) = \tilde{S}_r^{r+1}p_Z(z) + O(h_H)$ in the proof of Proposition 3.1 in Kong, Linton, and Xia (2010). Also, by Lemma 1 of Horowitz (1998) (or Theorem 2.37 of Pollard (1984)), we derive

$$A_{11n}(z) = E[A_{11n}(z)] + (A_{11n}(z) - E[A_{11n}(z)])$$

$$=e'_d \tilde{S}_r^{-1} \tilde{S}_r^{r+1} \mathcal{H}_{r+1}(z) \cdot h_{\mathcal{H}}^r + o(h_{\mathcal{H}}^r) + o\left(h_{\mathcal{H}}^r \frac{\log(n)}{n^{1/2}}\right)$$
$$=e'_d \tilde{S}_r^{-1} \tilde{S}_r^{r+1} \mathcal{H}_{r+1}(z) \cdot h_{\mathcal{H}}^r + o(h_{\mathcal{H}}^r)$$

Therefore by (S.2.10), we derive

$$A_{1n}(z) = \frac{1}{nh_{\mathcal{H}}^3} \sum_{i=1}^n e'_d \tilde{S}_{n,r}(z)^{-1} \tilde{K} \Big(\frac{z - \zeta(X_i)}{h_{\mathcal{H}}} \Big) \Big(Y_i - H(X_i) \Big) \mu \Big(\frac{z - \zeta(X_i)}{h_{\mathcal{H}}} \Big) + e'_d \tilde{S}_r^{-1} \tilde{S}_r^{r+1} \mathcal{H}_{r+1}(z) \cdot h_{\mathcal{H}}^r + o(h_{\mathcal{H}}^r).$$
(S.2.11)

Second we consider $A_{2n}(z)$, plugging the asymptotic representation of $\hat{\zeta}_1(X_i^1) - \zeta_1(X_i^1)$ given by Lemma S.3 into $A_{2n}(z)$, under Assumption A.6,

$$\begin{aligned} A_{2n}(z) \\ = e'_{d}\tilde{S}_{n,r}(z)^{-1} \left(\frac{1}{nh_{\mathcal{H}}^{4}} \sum_{i=1}^{n} \left[\left(\frac{\partial}{\partial u^{1}} t(u,Y_{i};z)\tilde{K}(u) + t(u,Y_{i};z)\partial_{1}\tilde{K}(u) \right) \Big|_{u = \frac{\zeta(X_{i}) - z}{h_{\mathcal{H}}}} \right] D_{1}(X_{i}^{1}) \cdot h_{H}^{r+1} \\ &+ \frac{1}{nh_{\mathcal{H}}^{4}} \sum_{i=1}^{n} \left[\left(\frac{\partial}{\partial u^{1}} t(u,Y_{i};z)\tilde{K}(u) + t(u,Y_{i};z)\partial_{1}\tilde{K}(u) \right) \Big|_{u = \frac{\zeta(X_{i}) - z}{h_{\mathcal{H}}}} \right] J_{n1}(X_{i}^{1}) \right) \\ &+ o(h_{H}^{r+1}) \\ =: e'_{d}\tilde{S}_{n,r}(z)^{-1} \left(A_{21n}(z) + A_{22n}(z) \right) + o(h_{H}^{r+1}). \end{aligned}$$
(S.2.12)

As for $A_{21n}(z)$, note that by product rule of derivatives

$$\begin{split} & \left(\frac{\partial}{\partial u^{1}}t(u,Y_{i};z)\tilde{K}(u)+t(u,Y_{i};z)\partial_{1}\tilde{K}(u)\right)\Big|_{u=\frac{\zeta(X_{i})-z}{h_{\mathcal{H}}}} \\ &= \frac{\partial}{\partial u^{1}}\left\{t(u,Y_{i};z)\tilde{K}(u)\right\}\Big|_{u=\frac{\zeta(X_{i})-z}{h_{\mathcal{H}}}} \\ &= \left(Y_{i}-\mathcal{H}(\zeta(X_{i}))\right)\left(\frac{\partial}{\partial u^{1}}\left\{\tilde{\mu}(u)\tilde{K}(u)\right\}\right)\Big|_{u=\frac{\zeta(X_{i})-z}{h_{\mathcal{H}}}} \\ &+ \frac{\partial}{\partial u^{1}}\left\{t(u,\mathcal{H}(z+h_{\mathcal{H}}u);z)\tilde{K}(u)\right\}\Big|_{u=\frac{\zeta(X_{i})-z}{h_{\mathcal{H}}}} \\ &- \tilde{\mu}(u)\tilde{K}(u)\Big|_{u=\frac{\zeta(X_{i})-z}{h_{\mathcal{H}}}} \cdot \frac{\partial}{\partial z^{1}}\mathcal{H}(z)\Big|_{z=\zeta(X_{i})} \cdot h_{\mathcal{H}}, \end{split}$$

where $t(u, Y_i; z)\tilde{K}(u) = \tilde{\mu}(u)(Y_i - \tilde{\mu}(u)'B_{\mathcal{H}}\beta(z))\tilde{K}(u)$. Thus, we can further decompose $A_{21n}(z)$ as

$$\begin{aligned} A_{21n}(z) \\ = & \frac{1}{nh_{\mathcal{H}}^4} \sum_{i=1}^n \left(Y_i - \mathcal{H}(\zeta(X_i)) \right) D_1(X_i^1) \left(\frac{\partial}{\partial u^1} \left\{ \tilde{\mu}(u) \tilde{K}(u) \right\} \right) \Big|_{u = \frac{\zeta(X_i) - z}{h_{\mathcal{H}}}} \cdot h_H^{r+1} \\ & + \frac{1}{nh_{\mathcal{H}}^4} \sum_{i=1}^n \frac{\partial}{\partial u^1} \left\{ t(u, \mathcal{H}(z + h_{\mathcal{H}}u); z) \tilde{K}(u) \right\} \Big|_{u = \frac{\zeta(X_i) - z}{h_{\mathcal{H}}}} \cdot D_1(X_i^1) \cdot h_H^{r+1} \end{aligned}$$

$$-\frac{1}{nh_{\mathcal{H}}^3}\sum_{i=1}^n \tilde{\mu}(u)\tilde{K}(u)\Big|_{u=\frac{\zeta(X_i)-z}{h_{\mathcal{H}}}} \cdot \frac{\partial}{\partial z^1}\mathcal{H}(z)\Big|_{z=\zeta(X_i)} \cdot D_1(X_i^1) \cdot h_{\mathcal{H}}^{r+1}$$
$$=:A_{211n}(z) + A_{212n}(z) + A_{213n}(z).$$

By Lemma 1 of Horowitz (1998) (or Theorem 2.37 of Pollard (1984)),

$$A_{211n}(z) = o\Big(\frac{h_{H}^{r+1}log(n)}{n^{1/2}h_{\mathcal{H}}^{3}}\Big),$$

and

$$\begin{split} &A_{212n}(z) \\ =& E[A_{212n}(z)] + \left(A_{212n}(z) - E[A_{212n}(z)]\right) \\ =& \frac{h_{H}^{r+1}}{h_{\mathcal{H}}^{4}} \int \frac{\partial}{\partial u^{1}} \Big\{ t(u, \mathcal{H}(z+h_{\mathcal{H}}u); z)\tilde{K}(u) \Big\} \Big|_{u=\frac{\zeta(X_{i})-z}{h_{\mathcal{H}}}} E[D_{1}(X_{i}^{1})|\zeta(X_{i}) = Z] p_{Z}(Z) dZ + o \Big(h_{H}^{r+1} \frac{h_{\mathcal{H}}^{(r-5)/2} log(n)}{n^{1/2}}\Big) \\ =& -\frac{h_{H}^{r+1}}{h_{\mathcal{H}}^{2}} \int \Big\{ t(u, \mathcal{H}(z+h_{\mathcal{H}}u); z)\tilde{K}(u) \Big\} \frac{\partial}{\partial u^{1}} E[D_{1}(X_{i}^{1})|\zeta(X_{i}) = z+h_{\mathcal{H}}u] p_{Z}(z+h_{\mathcal{H}}u) du + o \Big(\frac{h_{H}^{r+1}h_{\mathcal{H}}^{(r-5)/2} log(n)}{n^{1/2}}\Big) \\ =& O(h_{H}^{r+1} \cdot h_{\mathcal{H}}^{r}) + o \Big(\frac{h_{H}^{r+1}h_{\mathcal{H}}^{(r-5)/2} log(n)}{n^{1/2}}\Big), \end{split}$$

where the last second equality is derived by change of variable and integration by parts, and the last equality is due to Taylor expansion. As for $A_{213n}(z)$, similar to $A_{212n}(z)$, by Lemma 1 of Horowitz (1998) (or Theorem 2.37 of Pollard (1984)), we have

$$\begin{split} &A_{213n}(z) \\ =& E[A_{213n}(z)] + \left(A_{213n}(z) - E[A_{213n}(z)]\right) \\ =& -\frac{h_{H}^{r+1}}{h_{\mathcal{H}}^{3}} \int \tilde{\mu}(u)\tilde{K}(u)\big|_{u=\frac{Z-z}{h_{\mathcal{H}}}} \cdot \frac{\partial}{\partial z^{1}}\mathcal{H}(z)\Big|_{z=Z}E[D_{1}(X_{i}^{1})|\zeta(X_{i}) = Z]p_{Z}(Z)dZ + o\Big(\frac{h_{H}^{r+1}log(n)}{n^{1/2}h_{\mathcal{H}}^{2}}\Big) \\ =& -V_{r}^{\tilde{\mu}}\frac{\partial}{\partial z^{1}}\mathcal{H}(z)E[D_{1}(X_{i}^{1})|\zeta(X_{i}) = z]p_{Z}(z) \cdot \frac{h_{H}^{r+1}}{h_{\mathcal{H}}} \\ & -V_{r}^{\tilde{\mu}}(1)\frac{\partial}{\partial z^{1}}\left\{\frac{\partial}{\partial z^{1}}\mathcal{H}(z)E[D_{1}(X_{i}^{1})|\zeta(X_{i}) = z]p_{Z}(z)\right\} \cdot h_{H}^{r+1} \\ & -V_{r}^{\tilde{\mu}}(2)\frac{\partial}{\partial z^{2}}\left\{\frac{\partial}{\partial z^{1}}\mathcal{H}(z)E[D_{1}(X_{i}^{1})|\zeta(X_{i}) = z]p_{Z}(z)\right\} \cdot h_{H}^{r+1} \\ & + o\Big(h_{H}^{r+1} + \frac{h_{H}^{r+1}log(n)}{n^{1/2}h_{\mathcal{H}}^{2}}\Big), \end{split}$$

where $V_r^{\tilde{\mu}} = \int \tilde{\mu}(u)\tilde{K}(u)du$, $V_r^{\tilde{\mu}}(1) = \int u^1\tilde{\mu}(u)\tilde{K}(u)du$, and $V_r^{\tilde{\mu}}(2) = \int u^2\tilde{\mu}(u)\tilde{K}(u)du$. Therefore, by adding up $A_{211n}(z)$, $A_{212n}(z)$, and $A_{213n}(z)$, we derive

$$e_d'\tilde{S}_{n,r}(z)^{-1}A_{21n}$$

$$=e'_{d}(\tilde{S}_{r}^{-1}p_{Z}(z)^{-1}+O(h_{\mathcal{H}}))\cdot\left(-V_{r}^{\tilde{\mu}}\frac{\partial}{\partial z^{1}}\mathcal{H}(z)E[D_{1}(X_{i}^{1})|\zeta(X_{i})=z]p_{Z}(z)\cdot h_{H}^{r+1}\right)$$

$$-V_{r}^{\tilde{\mu}}(1)p_{Z}(z)^{-1}\frac{\partial}{\partial z^{1}}\left\{\frac{\partial}{\partial z^{1}}\mathcal{H}(z)E[D_{1}(X_{i}^{1})|\zeta(X_{i})=z]p_{Z}(z)\right\}\cdot h_{H}^{r+1}$$

$$-V_{r}^{\tilde{\mu}}(2)p_{Z}(z)^{-1}\frac{\partial}{\partial z^{2}}\left\{\frac{\partial}{\partial z^{1}}\mathcal{H}(z)E[D_{1}(X_{i}^{1})|\zeta(X_{i})=z]p_{Z}(z)\right\}\cdot h_{H}^{r+1}\right)$$

$$+o\left(h_{H}^{r+1}+\frac{h_{H}^{r+1}log(n)}{n^{1/2}h_{\mathcal{H}}^{3}}\right)$$

$$=-\tilde{e}_{1}p_{Z}(z)^{-1}\frac{\partial}{\partial z^{1}}\left\{\frac{\partial}{\partial z^{1}}\mathcal{H}(z)E[D_{1}(X_{i}^{1})|\zeta(X_{i})=z]p_{Z}(z)\right\}\cdot h_{H}^{r+1}$$

$$-\tilde{e}_{2}p_{Z}(z)^{-1}\frac{\partial}{\partial z^{2}}\left\{\frac{\partial}{\partial z^{1}}\mathcal{H}(z)E[D_{1}(X_{i}^{1})|\zeta(X_{i})=z]p_{Z}(z)\right\}\cdot h_{H}^{r+1}$$

$$+o\left(h_{H}^{r+1}+\frac{h_{H}^{r+1}log(n)}{n^{1/2}h_{\mathcal{H}}^{3}}\right),$$
(S.2.13)

where the last equality is due to the facts that $\tilde{S}_{r}^{-1}V_{r}^{\tilde{\mu}} = e_{1}, \ \tilde{S}_{r}^{-1}V_{r}^{\tilde{\mu}}(1) = (0, 1, 0, \dots, 0)'$, and $\tilde{S}_{r}^{-1}V_{r}^{\tilde{\mu}}(1) = (0, 0, 1, \dots, 0)'$. As for $A_{22n}(z)$, note that by (S.2.3) in Lemma S.3,

$$\sup_{x^1 \in \mathcal{S}_{X^1}} \left| J_{n1}(x^1) \right| = O\left(\sqrt{\frac{\log(n)}{nh_H^{d_1}}}\right)$$

in probability as $n \to \infty$, and $E[J_{n1}(X_i^1)|X_i^1 = x^1] = 0$. Thus by Lemma 1 of Horowitz (1998) (or Theorem 2.37 of Pollard (1984)),

$$A_{22n}(z) = \frac{1}{nh_{\mathcal{H}}^4} \sum_{i=1}^n \left[\left(\frac{\partial}{\partial u^1} t(u, Y_i; z) \tilde{K}(u) + t(u, Y_i; z) \partial_1 \tilde{K}(u) \right) \Big|_{u = \frac{z - \zeta(X_i)}{h_{\mathcal{H}}}} \right] J_{n1}(X_i^1) \right]$$
$$= o \left(\frac{\log(n)^{3/2}}{nh_{\mathcal{H}}^3 h_{\mathcal{H}}^{d_{1/2}}} \right).$$
(S.2.14)

Plugging (S.2.13) and (S.2.14) into (S.2.12), we get

$$A_{2n}(z) = -\tilde{e}_1 p_Z(z)^{-1} \frac{\partial}{\partial z^1} \left\{ \frac{\partial}{\partial z^1} \mathcal{H}(z) E[D_1(X_i^1) | \zeta(X_i) = z] p_Z(z) \right\} \cdot h_H^{r+1} - \tilde{e}_2 p_Z(z)^{-1} \frac{\partial}{\partial z^2} \left\{ \frac{\partial}{\partial z^1} \mathcal{H}(z) E[D_1(X_i^1) | \zeta(X_i) = z] p_Z(z) \right\} \cdot h_H^{r+1} + o \left(h_H^{r+1} + \frac{h_H^{r+1} log(n)}{n^{1/2} h_{\mathcal{H}}^3} + \frac{h_H^{r+1} h_{\mathcal{H}}^{(r-5)/2} log(n)}{n^{1/2}} + \frac{log(n)^{3/2}}{n h_{\mathcal{H}}^3 h_H^{d_1/2}} \right)$$
(S.2.15)

Similarly, we have the decomposition for $A_{3n}(z)$. By adding up the representations of $A_{1n}(z)$, $A_{2n}(z)$, and $A_{3n}(z)$, The desired conclusion therefore follows from Assumption A.6.

S.2.6 Lemma S.6

Let $\tilde{f}_k(\cdot)$ be an infeasible estimator of $\tilde{f}_k(\cdot)$ with an (infeasible) data of $\{Y_i, \zeta_1(X_i^1), \zeta_2(X_i^2)\}_{i=1}^n$, while $\hat{f}_k(\cdot)$ be a feasible estimator of $\tilde{f}_k(\cdot)$ with a data of $\{Y_i, \hat{\zeta}_1(X_i^1), \hat{\zeta}_2(X_i^2)\}_{i=1}^n$. Lemma S.6 establishes the uniform convergence rate and asymptotic representation of the feasible estimator of transformed component function $\tilde{f}_k(\cdot)$ for k = 1, 2. In particular, the first three terms in the asymptotic representation come from the (asymptotic) representation of infeasible estimator $\check{f}_k(z^k)$, while the fourth term is the additional bias appearing in the difference between feasible and infeasible estimators, namely $\hat{f}_k(z^k) - \check{f}_k(z^k)$.

Lemma S.6. If Assumptions A.1-A.6 hold, then for $k = 1, 2, (i) \hat{f}_k(z^k) - \tilde{f}_k(z^k) = \tilde{\mathfrak{J}}_{nk}(z^k) - E[\tilde{\mathfrak{J}}_{nk}(z^k)] + h_{\mathcal{H}}^{r+1}\tilde{\mathcal{B}}_k(z^k) + h_{\mathcal{H}}^{r+1}\tilde{\mathcal{B}}_k(z^k) + o_p(h_{\mathcal{H}}^r + h_{\mathcal{H}}^{r+1}), and (ii) \hat{f}_k(z^k) - \tilde{f}_k(z^k) = O_p(h_{\mathcal{H}}^r + \sqrt{\frac{\log(n)}{nh_{\mathcal{H}}}} + h_{\mathcal{H}}^{r+1}) as n \to \infty$ uniformly over $z^k \in \mathcal{S}_{Z^k}$.

Proof. Only the case for k = 2 is proved. The proof for k = 1 is similar. The definition of $\hat{f}_2(\cdot)$ yields

$$\widehat{f}_{2}(z^{2}) - \widetilde{f}_{2}(z^{2}) = \int_{z_{0}^{2}}^{z^{2}} \int \left[\frac{\partial_{2}\widehat{\mathcal{H}}(\nu)}{\partial_{1}\widehat{\mathcal{H}}(\nu)} - \frac{\partial_{2}\mathcal{H}(\nu)}{\partial_{1}\mathcal{H}(\nu)}\right] \omega_{3}(\nu^{1})d\nu^{1}d\nu^{2}$$
(S.2.16)

By applying Taylor expansion to the integrand,

$$\frac{\partial_{2}\widehat{\mathcal{H}}(\nu)}{\partial_{1}\widehat{\mathcal{H}}(\nu)} - \frac{\partial_{2}\mathcal{H}(\nu)}{\partial_{1}\mathcal{H}(\nu)} = \frac{\partial_{2}\widehat{\mathcal{H}}(\nu)}{\partial_{1}\mathcal{H}(\nu)} - \frac{\partial_{2}\mathcal{H}(\nu)}{[\partial_{1}\mathcal{H}(\nu)]^{2}} [\partial_{1}\widehat{\mathcal{H}}(\nu) - \partial_{1}\mathcal{H}(\nu)] + O(\xi_{H1}^{2} + [\xi_{\mathcal{H}}']^{2}) = q_{2}(\nu)' \cdot \left(\left[\widehat{\partial_{1}\widehat{\mathcal{H}}}(\nu) \\ \widehat{\partial_{2}\widehat{\mathcal{H}}}(\nu) \right] - \left[\frac{\partial_{1}\mathcal{H}(\nu)}{\partial_{2}\mathcal{H}(\nu)} \right] \right) + O(\xi_{H1}^{2} + [\xi_{\mathcal{H}}']^{2}) \tag{S.2.17}$$

in probability as $n \to \infty$ uniformly over $\nu \in S_Z$, where $q_2(\nu) = \left[-\frac{\partial_2 \mathcal{H}(\nu)}{[\partial_1 \mathcal{H}(\nu)]^2}, \frac{1}{\partial_1 \mathcal{H}(\nu)} \right]'$. By Lemma S.5 and plugging the representations of $\widehat{\partial_k \mathcal{H}}(\nu) - \partial_k \mathcal{H}(\nu)$ into (S.2.17), we derive

$$\frac{\partial_{2}\widehat{\mathcal{H}}(\nu)}{\partial_{1}\widehat{\mathcal{H}}(\nu)} - \frac{\partial_{2}\mathcal{H}(\nu)}{\partial_{1}\mathcal{H}(\nu)} = \frac{1}{nh_{\mathcal{H}}^{3}}q_{2}(\nu)'\sum_{i=1}^{n}e_{d}'\widetilde{S}_{n,r}(\nu)^{-1}\widetilde{K}\left(\frac{\nu-\zeta(X_{i})}{h_{\mathcal{H}}}\right)\left(Y_{i}-H(X_{i})\right)\mu\left(\frac{\zeta(X_{i})-\nu}{h_{\mathcal{H}}}\right)
+ q_{2}(\nu)'\mathfrak{D}(\nu)h_{\mathcal{H}}^{r} + q_{2}(\nu)'\mathcal{D}(\nu)h^{r+1}
+ o\left(h_{\mathcal{H}}^{r}+h_{H}^{r+1}\right) + O\left(\frac{\log(n)}{nh_{\mathcal{H}}^{4}} + \frac{\log(n)}{nh_{H}^{d}}\right)$$
(S.2.18)

in probability as $n \to \infty$ uniformly over $z \in S_Z$. Therefore by integrating (S.2.18) and Assumption A.6,

$$\begin{aligned} \hat{f}_{2}(z^{2}) &- \tilde{f}_{2}(z^{2}) \\ &= \int_{z_{0}^{2}}^{z^{2}} \int \left[\frac{\partial_{2} \hat{\mathcal{H}}(\nu)}{\partial_{1} \hat{\mathcal{H}}(\nu)} - \frac{\partial_{2} \mathcal{H}(\nu)}{\partial_{1} \mathcal{H}(\nu)} \right] \omega_{3}(\nu^{1}) d\nu^{1} d\nu^{2} \end{aligned}$$

$$=\frac{1}{nh_{\mathcal{H}}^{3}}\sum_{i=1}^{n}\int_{z_{0}^{2}}^{z^{2}}\int q_{2}(\nu)'e_{d}'\tilde{S}_{n,r}(\nu)^{-1}\tilde{K}(\frac{\nu-\zeta(X_{i})}{h_{\mathcal{H}}})(Y_{i}-H(X_{i}))\mu(\frac{\zeta(X_{i})-\nu}{h_{\mathcal{H}}})\omega_{3}(\nu^{1})d\nu^{1}d\nu^{2} +\mathfrak{B}_{2}(z^{2})h_{\mathcal{H}}^{r}+\tilde{B}_{2}(z^{2})h^{r+1}+o(h_{\mathcal{H}}^{r}+h_{\mathcal{H}}^{r+1}).$$
(S.2.19)

The rest of the proof is to analyse the first term of the right hand side in (S.2.19). A change of variables and a Taylor expansion show that

$$\begin{split} \frac{1}{nh_{\mathcal{H}}^{3}} \sum_{i=1}^{n} \int_{z_{0}^{2}}^{z_{0}^{2}} \int q_{2}(v)' e_{d}' \tilde{S}_{n,r}(v)^{-1} \tilde{K} \Big(\frac{v - \zeta(X_{i})}{h_{\mathcal{H}}} \Big) \big(Y_{i} - H(X_{i})\big) \mu \Big(\frac{\zeta(X_{i}) - v}{h_{\mathcal{H}}} \Big) \omega_{3}(v^{1}) dv^{1} dv^{2} \\ = \frac{1}{nh_{\mathcal{H}}^{2}} \sum_{i=1}^{n} \int_{z_{0}^{2}}^{z_{0}^{2}} q_{2}(\zeta_{1}(X_{i}^{1}), v^{2})' e_{d}' \tilde{S}_{r}^{-1} k_{2} \Big(\frac{v^{2} - \zeta_{2}(X_{i}^{2})}{h_{\mathcal{H}}} \Big) \frac{Y_{i} - H(X_{i})}{p_{Z}(\zeta_{1}(X_{i}^{1}), v^{2})} V_{2}^{\tilde{\mu}} \Big(\frac{\zeta_{2}(X_{i}^{2}) - v^{2}}{h_{\mathcal{H}}} \Big) \omega_{3}(\zeta_{1}(X_{i}^{1})) dv^{2} \\ + \frac{1}{nh_{\mathcal{H}}^{3}} \sum_{i=1}^{n} \int_{z_{0}^{2}}^{z_{0}^{2}} \int q_{2}(v)' e_{d}' \Big(\tilde{S}_{n,r}(v)^{-1} - \{p_{Z}(v)\}^{-1} \tilde{S}_{r}^{-1} \Big) \\ \cdot \tilde{K} \Big(\frac{v - \zeta(X_{i})}{h_{\mathcal{H}}} \Big) \big(Y_{i} - H(X_{i})\big) \mu \Big(\frac{\zeta(X_{i}) - v}{h_{\mathcal{H}}} \Big) \omega_{3}(v^{1}) dv^{1} dv^{2} \\ + \frac{1}{n} \sum_{i=1}^{n} \int_{(z_{0}^{2} - \zeta_{2}(X_{i}))/h_{\mathcal{H}}}^{(z^{2} - \zeta_{2}(X_{i}^{2}) + h_{\mathcal{H}}u)} \frac{\partial}{\partial z^{1}} \left[q_{2}(z^{1}, \zeta_{2}(X_{i}^{2}) + h_{\mathcal{H}}u)' e_{d}' \tilde{S}_{r}^{-1} \\ \cdot \frac{(Y_{i} - H(X_{i}))\omega_{3}(z^{1})}{p_{Z}(z^{1}, \zeta_{2}(X_{i}^{2}) + h_{\mathcal{H}}u)} \right] \right|_{z^{1} = \zeta_{1}(X_{i}^{1})} \left\{ \int u^{1}\tilde{\mu}(u)k_{1}(u^{1}) du^{1} \right\} k_{2}(u^{2}) du^{2} + o\left(\frac{1}{n}\right) \\ =: Q_{1n}(z^{2}) + Q_{2n}(z^{2}) + Q_{3n}(z^{2}) + o(h_{H}^{r+1}), \end{split}$$

$$(S.2.20)$$

where $o(1/n) = o(h_H^{r+1})$ is due to Assumption A.6. By Lemma 1 in Horowitz (1998) (or Theorem 2.37 in Pollard (1984)) and Assumption A.6, $Q_{2n}(z^2) = o(h_H^{r+1})$ and $Q_{3n}(z^2) = o(h_H^{r+1})$. As for $Q_{1n}(z^2)$, an integration by parts implies that

$$\begin{aligned} Q_{1n}(z^{2}) \\ = & \frac{1}{nh_{\mathcal{H}}} \sum_{i=1}^{n} \left[q_{2}(\zeta_{1}(X_{i}^{1}), z^{2})' e_{d}' \tilde{S}_{r}^{-1} \frac{Y_{i} - H(X_{i})}{p_{Z}(\zeta_{1}(X_{i}^{1}), z^{2})} V_{2}^{\tilde{\mu}} \Big(\frac{\zeta_{2}(X_{i}^{2}) - z^{2}}{h_{\mathcal{H}}} \Big) \omega_{3}(\zeta_{1}(X_{i}^{1})) \mathcal{K}_{2} \Big(\frac{z^{2} - \zeta_{2}(X_{i}^{2})}{h_{\mathcal{H}}} \Big) \\ & - q_{2}(\zeta_{1}(X_{i}^{1}), z_{0}^{2})' e_{d}' \tilde{S}_{r}^{-1} \frac{Y_{i} - H(X_{i})}{p_{Z}(\zeta_{1}(X_{i}^{1}), z_{0}^{2})} V_{2}^{\tilde{\mu}} \Big(\frac{\zeta_{2}(X_{i}^{2}) - z_{0}^{2}}{h_{\mathcal{H}}} \Big) \omega_{3}(\zeta_{1}(X_{i}^{1})) \mathcal{K}_{2} \Big(\frac{z_{0}^{2} - \zeta_{2}(X_{i}^{2})}{h_{\mathcal{H}}} \Big) \Big] \\ & - \frac{1}{n} \sum_{i=1}^{n} \int_{(z_{0}^{2} - \zeta_{2}(X_{i}))/h_{\mathcal{H}}}^{(z^{2} - \zeta_{2}(X_{i}))/h_{\mathcal{H}}} \frac{\partial}{\partial z^{2}} \left[q_{2}(\zeta_{1}(X_{i}^{1}), z^{2})' e_{d}' \tilde{S}_{r}^{-1} \\ & \cdot \frac{(Y_{i} - H(X_{i}))\omega_{3}(\zeta_{1}(X_{i}^{1}))}{p_{Z}(\zeta_{1}(X_{i}^{1}), z^{2})} V_{2}^{\tilde{\mu}}(u) \right] \Big|_{z^{2} = \zeta_{2}(X_{i}^{2}) + h_{\mathcal{H}}u} \mathcal{K}_{2}(u^{2}) du^{2} \\ &= \tilde{\mathfrak{J}}_{n2}(z^{2}) + Q_{12n}(z^{2}), \end{aligned}$$

$$(S.2.21)$$

where $\mathcal{K}_2(u) = \int_{-\infty}^{u} k_2(t) dt$. By Lemma 1 in Horowitz (1998) (or Theorem 2.37 in Pollard (1984)) and Assumption A.6, $Q_{12n}(z^2) = o(h_H^{r+1})$ in probability as $n \to \infty$ uniformly over $z^2 \in S_{Z^2}$. Rearranging (S.2.19), (S.2.20), and (S.2.21), then part (i) is proved. Also, an argument similar to the proof of Theorem 6 in Masry (1996) shows that $\sup_{z^2 \in S_{Z^2}} |\tilde{\mathfrak{J}}_{n2}(z^2)| = O_p\left(\sqrt{\frac{\log(n)}{nh_{\mathcal{H}}}}\right)$. Thus, part (ii) follows from Assumption A.6.