

FITTING A BIVARIATE ADDITIVE MODEL BY LOCAL POLYNOMIAL REGRESSION

BY JEAN D. OPSOMER¹ AND DAVID RUPPERT²

Iowa State University and Cornell University

While the additive model is a popular nonparametric regression method, many of its theoretical properties are not well understood, especially when the backfitting algorithm is used for computation of the estimators. This article explores those properties when the additive model is fitted by local polynomial regression. Sufficient conditions guaranteeing the asymptotic existence of unique estimators for the bivariate additive model are given. Asymptotic approximations to the bias and the variance of a homoscedastic bivariate additive model with local polynomial terms of odd and even degree are computed. This model is shown to have the same rate of convergence as that of univariate local polynomial regression.

1. Introduction. Nonparametric regression methods are a flexible and growing class of models in the statistician's toolbox. They allow researchers to evaluate data without having to postulate a shape for the relationship between the response variable and the covariate(s). Unfortunately, nonparametric regression methods become more cumbersome to implement when the number of covariates increases, and the ability to visually inspect estimated relationships is often lost when there are more than two covariates. An elegant solution to these problems is provided by the *additive model*, originally suggested by Friedman and Stuetzle (1981) and popularized by Hastie and Tibshirani (1990). The additive model assumes that the conditional expectation function of the dependent variable Y can be written as a sum of smooth terms in the covariates X_1, \dots, X_D :

$$(1) \quad E(Y|X = (x_1, \dots, x_D)) = m(x_1, \dots, x_D) = m_1(x_1) + \dots + m_D(x_D).$$

The backfitting algorithm proposed by Buja, Hastie and Tibshirani (1989) and the related fitting procedure in S-PLUS [see Chambers and Hastie (1992)] have made the additive model a popular choice for multivariate nonparametric fitting.

Compared to the development of practical applications, the understanding of the theoretical properties of the additive model has lagged. Stone (1985) showed in the case of additive splines that the optimal rate of convergence achievable by additive model estimators is independent of the number of covariates, and Burman (1990) proposed a cross-validation method for select-

Received November 1995; revised May 1996.

¹Research supported by NIEHS Training Grant ES07261 while at Cornell University.

²Research supported by NSA Grant MDA904-95-H-1025 and NSF Grant DMS-93-06196.

AMS 1991 *subject classifications*. Primary 62G07; secondary 62H99.

Key words and phrases. Additive model, local polynomial regression, optimal rates, existence, backfitting.

ing the number of knots. Other authors have proven existence results in the context of smoothing splines [Wahba (1986), Gu, Bates, Chen and Wahba (1989), and Chen (1993)] and interpolated regressograms [Härdle and Hall (1993)]. More recently, a paper by Linton and Nielsen (1995) describes a fitting procedure for bivariate additive models based on local linear regression and marginal integration, and they showed that their procedure also achieves the same $O_p(n^{-2/5})$ rate of convergence for the additive model as for the univariate local linear estimator. In the case of additive modeling through backfitting, the theoretical investigations are greatly complicated by the fact that the estimators are defined as the solution of an iterative algorithm. Only when two covariates are present are explicit expressions for the estimators available. As Linton and Nielsen (1995) note, however, “these expressions appear quite intractable.”

Buja, Hastie and Tibshirani (1989) provide sufficient conditions that guarantee the convergence of the backfitting algorithm or, equivalently, the existence of the estimators. These conditions are only generally satisfied for regression splines and parametric terms, but not by kernel regression or local polynomial regression. This is unfortunate, because local polynomial regression has recently been shown to possess many desirable theoretical and practical properties [e.g., Cleveland and Devlin (1988), Fan, Gasser, Gijbels, Brockmann and Engel (1993), and Ruppert and Wand (1994)] and its combination with backfitting has proven to be very popular for fitting additive models in S-PLUS.

In this article, we will explore two important theoretical issues concerning the bivariate additive model, in the context of backfitted estimators using local polynomial regression:

1. What are sufficient conditions guaranteeing convergence of backfitting?
2. What are the asymptotic properties of the estimators?

In another paper [Opsomer and Ruppert (1995)], these results are used to develop a fully automated plug-in bandwidth selection method.

The rest of the article will proceed as follows. In Section 2, the bivariate additive model estimators are defined. In Section 3, sufficient conditions for the existence of unique estimators are discussed. Section 4 derives conditional asymptotic bias and variance expressions when the local polynomials are of odd degree. Section 5 extends the results to local regression of even degree.

2. Definition of the estimators. Let $(X_1, Z_1, Y_1), \dots, (X_n, Z_n, Y_n)$ be a set of independent and identically distributed \mathbb{R}^3 -valued random variables. We assume the following model:

$$Y_i = \alpha + m_1(X_i) + m_2(Z_i) + \varepsilon_i,$$

where the ε_i are independent and identically distributed with mean 0 and variance σ^2 . To ensure identifiability of the functions m_1 and m_2 , we include the intercept α and assume $E(m_1(X_i)) = E(m_2(Z_i)) = 0$.

We introduce some notation. Let $\mathbf{Y} = (Y_1, \dots, Y_n)^T$ and similarly for \mathbf{X} and \mathbf{Z} , and we write the vectors of additive functions at the observation points as $\mathbf{m}_1 = (m_1(X_1), \dots, m_1(X_n))^T$, $\mathbf{m}_2 = (m_2(Z_1), \dots, m_2(Z_n))^T$. For any constant d , \mathbf{d} is the n -valued vector $(d, \dots, d)^T$. Let $\mathbf{s}_{1,x}^T, \mathbf{s}_{2,z}^T$ represent the *equivalent kernels* for the local polynomial regression at x and z . In the case of x , this equivalent kernel can be written as

$$\mathbf{s}_{1,x}^T = \mathbf{e}_1^T (\mathbf{X}_x^T \mathbf{W}_x \mathbf{X}_x)^{-1} \mathbf{X}_x^T \mathbf{W}_x,$$

where $\mathbf{e}_1^T = (1, 0)$,

$$\mathbf{W}_x = \text{diag} \left\{ \frac{1}{h_1} K \left(\frac{X_1 - x}{h_1} \right), \dots, \frac{1}{h_1} K \left(\frac{X_n - x}{h_1} \right) \right\}$$

for some kernel function K and bandwidth h_1 and

$$\mathbf{X}_x = \begin{bmatrix} 1 & (X_1 - x) & \cdots & (X_1 - x)^{p_1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & (X_n - x) & \cdots & (X_n - x)^{p_1} \end{bmatrix},$$

where p_1 is the order of the local polynomials for fitting m_1 [see Ruppert and Wand (1994)]. A similar expression holds for $\mathbf{s}_{2,z}^T$. Let \mathbf{S}_1 and \mathbf{S}_2 represent the smoother matrices whose rows are the equivalent kernels at the observations \mathbf{X} and \mathbf{Z} , respectively:

$$\mathbf{S}_1 = \begin{bmatrix} \mathbf{s}_{1,X_1}^T \\ \vdots \\ \mathbf{s}_{1,X_n}^T \end{bmatrix}, \quad \mathbf{S}_2 = \begin{bmatrix} \mathbf{s}_{2,Z_1}^T \\ \vdots \\ \mathbf{s}_{2,Z_n}^T \end{bmatrix}.$$

We define the vector of fitted values at the observation points as

$$\hat{\mathbf{m}} = \hat{\alpha} + \hat{\mathbf{m}}_1 + \hat{\mathbf{m}}_2,$$

in which $\hat{\alpha} = \bar{Y}$, and $\hat{\mathbf{m}}_1$ and $\hat{\mathbf{m}}_2$ are the solutions to the set of estimating equations

$$\begin{bmatrix} \mathbf{I} & \mathbf{S}_1^* \\ \mathbf{S}_2^* & \mathbf{I} \end{bmatrix} \begin{bmatrix} \hat{\mathbf{m}}_1 \\ \hat{\mathbf{m}}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{S}_1^* \\ \mathbf{S}_2^* \end{bmatrix} \mathbf{Y},$$

where $\mathbf{S}_1^* = (\mathbf{I} - \mathbf{1}\mathbf{1}^T/n)\mathbf{S}_1$ and similarly for \mathbf{S}_2^* . As discussed in Hastie and Tibshirani (1990), this adjustment of the smoothers, which they refer to as *centering*, is necessary to ensure uniqueness of the solutions to the estimating equations (if they exist), by requiring $\sum_{i=1}^n m_1(X_i) = \sum_{i=1}^n m_2(Z_i) = 0$. In practice, the estimating equations are solved using the backfitting algorithm, but in the bivariate case they also have the explicit solution

$$(2) \quad \begin{aligned} \hat{\mathbf{m}}_1 &= \{\mathbf{I} - (\mathbf{I} - \mathbf{S}_1^* \mathbf{S}_2^*)^{-1} (\mathbf{I} - \mathbf{S}_1^*)\} \mathbf{Y} \equiv \mathbf{W}_1 \mathbf{Y}, \\ \hat{\mathbf{m}}_2 &= \{\mathbf{I} - (\mathbf{I} - \mathbf{S}_2^* \mathbf{S}_1^*)^{-1} (\mathbf{I} - \mathbf{S}_2^*)\} \mathbf{Y} \equiv \mathbf{W}_2 \mathbf{Y}, \end{aligned}$$

provided the inverses exist. For \hat{m} , the explicit estimator expression is

$$(3) \quad \hat{m} = \{\mathbf{1}\mathbf{1}^T/n + 2\mathbf{I} - (\mathbf{I} - \mathbf{S}_1^*\mathbf{S}_2^*)^{-1}(\mathbf{I} - \mathbf{S}_1^*) - (\mathbf{I} - \mathbf{S}_2^*\mathbf{S}_1^*)^{-1}(\mathbf{I} - \mathbf{S}_2^*)\}\mathbf{Y} \equiv \mathbf{W}\mathbf{Y}.$$

3. Existence. We will specify a set of circumstances under which the estimators in (2) and (3) are guaranteed to exist. Let $f(x, z)$ represent the joint density of X_i and Z_i , with $f_X(x)$ and $f_Z(z)$ the corresponding marginal densities. For the kernel function K , we write the moments of K as $\mu_j(K) = \int u^j K(u) du$ for any j and let $R(K) = \int K(u)^2 du$.

One of the important issues in the theoretical derivations in this and following sections is whether or not an observation (X_i, Z_i) is close to the boundary of its domain. Using the notation of Ruppert and Wand (1994), we can formalize this by defining

$$D_{x, h_1} = \{t: (x + h_1 t) \in \text{supp}(f_X)\} \cap \text{supp}(K).$$

We then say that x is an *interior point* if and only if $D_{x, h_1} = \text{supp}(K)$. Otherwise, x is a *boundary point*. Another way to understand this distinction is depicted in Figure 1: for a given bandwidth value h_1 , x_1 is a boundary point, because there are values of x for which $K((x - x_1)/h_1) \neq 0$ outside of $\text{supp}(f_X) = [a, b]$, while x_2 is an interior point, because all values for which $K((x - x_2)/h_2) \neq 0$ are inside that support. Analogous definitions hold for D_{z, h_2} and z . We define the *boundary moments* of K with respect to x as

$$\mu_j(K, x) = \int_{D_{x, h_1}} u^j K(u) du,$$

where the dependency on h_1 will be suppressed for notational simplicity, and similarly for $R(K, x)$. Clearly, if x is an interior point, $\mu_j(K, x) = \mu_j(K)$ and $R(K, x) = R(K)$. Let \mathbf{N}_p represent the $(p + 1) \times (p + 1)$ matrix whose (i, j) th element is equal to $\mu_{i+j-2}(K)$ and $\mathbf{M}_p(u)$ be the same as \mathbf{N}_p , but with the first column replaced by $(1, u, \dots, u^p)^T$. As in Ruppert and Wand (1994), define the kernel

$$K_{(p)}(u) = \{|\mathbf{M}_p(u)|/|\mathbf{N}_p|\}K(u).$$

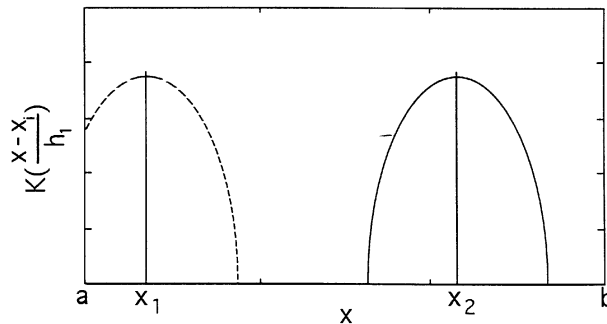


FIG. 1. Graphical representation of a boundary point (x_1) and an interior point (x_2) for a kernel function K and $\text{supp}(f_X) = [a, b]$.

For the boundary moments, we define the matrices $\mathbf{N}_p(x)$ and $\mathbf{M}_p(u, x)$ exactly as above, but with the $\mu_j(K)$ replaced by $\mu_j(K, x)$. We also define the boundary kernel

$$K_{(p)}(u, x) = \{|\mathbf{M}_p(u, x)|/|\mathbf{N}_p(x)|\}K(u)I_{(u \in D_{x, h_1})}.$$

If x is not at the boundary of $\text{supp}(f)$, $K_{(p)}(u, x) = K_{(p)}(u)$. Analogous definitions hold for the other covariate z .

We make the following assumptions.

ASSUMPTION 1. The kernel K is bounded and continuous, it has compact support and its first derivative has a finite number of sign changes over its support. Also, $\mu_j(K) = 0$ for all odd j and $\mu_{p_1+1}(K_{(p_1)}), \mu_{p_2+1}(K_{(p_2)}) \neq 0$.

ASSUMPTION 2. The densities f, f_X and f_Z are bounded and continuous, have compact support and their first derivatives have a finite number of sign changes over their supports. Also, $f_X(x) > 0, f_Z(z) > 0$ for all $(x, z) \in \text{supp}(f)$ and

$$(4) \quad \sup_{x, z} \left| \frac{f(x, z)}{f_X(x)f_Z(z)} - 1 \right| < 1.$$

ASSUMPTION 3. As $n \rightarrow \infty$, $h_1, h_2 \rightarrow 0$ and $nh_1/\log n, nh_2/\log n \rightarrow \infty$.

The following two lemmas show that, under these assumptions, the matrix inverses in the estimators (2) and (3) are well defined for local polynomials of any degree p_1 and p_2 “for sufficiently large n .” Strictly speaking, the lemmas only prove the existence of the estimator $\hat{\mathbf{m}}_1$, but it is clear that the results also hold for $\hat{\mathbf{m}}_2$ and $\hat{\mathbf{m}}$. The proofs are given in Appendix A.

LEMMA 3.1. *Under Assumptions 1–3, the following asymptotic approximations hold uniformly over all elements of the matrices:*

$$\begin{aligned} \mathbf{S}_1^* &= \mathbf{S}_1 - \mathbf{1}\mathbf{1}^T/n + o(\mathbf{1}\mathbf{1}^T/n) \quad a.s., \\ \mathbf{S}_1^*\mathbf{S}_2^* &= \mathbf{T}_{12}^* + o(\mathbf{1}\mathbf{1}^T/n) \quad a.s., \end{aligned}$$

where \mathbf{T}_{12}^* is a matrix whose (i, j) th element is

$$[\mathbf{T}_{12}^*]_{ij} = \frac{1}{n} \frac{f(X_i, Z_j)}{f_X(X_i)f_Z(Z_j)} - \frac{1}{n}.$$

LEMMA 3.2. *If Assumptions 1–3 hold, then $(\mathbf{I} - \mathbf{T}_{12}^*)$ is invertible for all n and*

$$P\{\text{there exists } N \text{ such that } (\mathbf{I} - \mathbf{S}_1^*\mathbf{S}_2^*) \text{ is invertible for all } n \geq N\} = 1.$$

When $(\mathbf{I} - \mathbf{S}_1^*\mathbf{S}_2^*)^{-1}$ exists,

$$\begin{aligned} (\mathbf{I} - \mathbf{S}_1^*\mathbf{S}_2^*)^{-1} &= (\mathbf{I} - \mathbf{T}_{12}^*)^{-1} + o(\mathbf{1}\mathbf{1}^T/n) \quad a.s. \\ &= \mathbf{I} + O(\mathbf{1}\mathbf{1}^T/n) \quad a.s. \end{aligned}$$

uniformly over all elements of the matrices.

REMARK 3.1. The restriction (4) in Assumption 2 ensures that $\|\mathbf{T}_{12}^*\|_r < 1$, where $\|\mathbf{A}\|_r$ denotes the *maximum row sum matrix norm* of the square matrix \mathbf{A} : $\|\mathbf{A}\|_r = \max_{1 \leq i \leq n} \sum_{j=1}^n |A_{ij}|$ [Horn and Johnson (1985)]. The main reason for selecting this norm was convenience, and it is clear that any other norm $\|\cdot\|$ which ensures that $(\mathbf{I} - \mathbf{T}_{12}^*)$ is invertible whenever $\|\mathbf{T}_{12}^*\| < 1$ would be equally appropriate. To assess the restrictiveness of (4), let us evaluate its effect on the bivariate normal distribution with censored support. For simplicity, assume that the mean of the distribution is at the center of the range for both covariates and that the standardized ranges (i.e., range divided by standard deviation) are the same in both dimensions. Let r represent this standardized range for both covariates. We can then rewrite (4) as

$$(5) \quad \sup_{-r/2 \leq x, z \leq r/2} \left| \frac{1}{\sqrt{1-\rho^2}} \exp\left(-\frac{\rho}{2(1-\rho^2)}(\rho x^2 - 2xz + \rho z^2)\right) - 1 \right| < 1.$$

Figure 2 displays the values of the correlation coefficient ρ for which (5) is satisfied as a function of the standardized range r . Clearly, there is a trade-off between the ranges of X_i and Z_i (in units of their standard deviations)

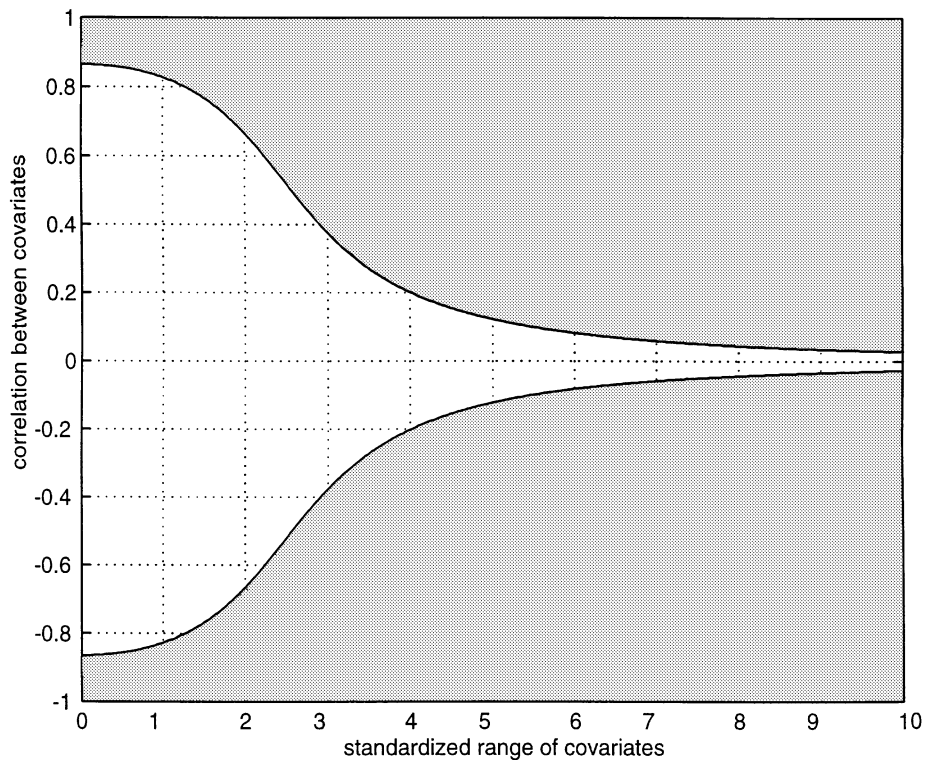


FIG. 2. Allowable values for the correlation between covariates following a censored normal distribution (unshaded region).

and the correlation. For instance, for a range of 3 standard deviation units, the correlation can take on values $-0.37 \leq \rho \leq 0.37$. While this may seem quite restrictive, it is important to realize that the constraint (4) is a *sufficient* condition for the existence of the estimators, not a necessary one. Note also that when the ranges of X_i and Z_i are very small relative to their standard deviation, that is, when the distribution is almost uniform, the amount of allowable correlation approaches $(-0.9, 0.9)$.

REMARK 3.2. Assumptions 1–3 are also somewhat stronger than ones usually made for local polynomial regression. Specifically, the kernel function and the density of the covariates have additional “smoothness” restrictions, and the maximum rate at which the bandwidths h_1, h_2 approach 0 is slightly slower. The purpose behind these restrictions is to allow us to use the uniform convergence results of Pollard (1984) in the proof of Lemma 3.1. They will not significantly affect the applicability of the results, since many commonly used kernels (including the Epanečnikov kernel) easily satisfy these conditions, and since the optimal rates of h_1, h_2 are fractional powers of n and therefore unaffected by the presence of the $\log n$ term in Assumption 3.

4. Conditional mean-squared error properties. As shown in Theorem 4.1 of Ruppert and Wand (1994), if the degree p of the local polynomial is even, the estimator has asymptotic bias of order $O_p(h^{p+2})$ in the interior, which is the same order as the estimator computed by local polynomial regression of (odd) degree $p + 1$. The asymptotic bias of the former estimator also contains an additional term. For these reasons, several authors [e.g., Fan and Gijbels (1994)] have argued that odd-degree local polynomials are preferable to even-degree ones. In this section, we will therefore restrict our attention to the case where p_1 and p_2 are odd. We will briefly discuss the situation where p_1 and p_2 are even in Section 5.

Let

$$D^p m_1 = \begin{bmatrix} \frac{d^p m_1(X_1)}{dx^p} \\ \vdots \\ \frac{d^p m_1(X_n)}{dx^p} \end{bmatrix}$$

and

$$E(m_1^{(p)}(X_i)|\mathbf{Z}) = \begin{bmatrix} E(m_1^{(p)}(X_1)|Z_1) \\ \vdots \\ E(m_1^{(p)}(X_n)|Z_n) \end{bmatrix}$$

and analogously for $D^p m_2$ and $E(m_2^{(p)}(Z_i)|\mathbf{X})$. Also, let $\mathbf{t}_i^T, \mathbf{v}_j$ represent the i th row and j th column of $(\mathbf{I} - \mathbf{T}_{12}^*)^{-1}$, respectively, and let \mathbf{e}_i^T represent the i th unit vector.

In addition to Assumptions 1–3 from Section 3, we also need the following assumption.

ASSUMPTION 4. The (p_1+1) th derivative of m_1 and the (p_2+1) th derivative of m_2 exist and are continuous and bounded.

In the theorem and corollaries that follow, we will only show the results for $\hat{\alpha}$, \hat{m}_1 and \hat{m} . It is clear that the results for \hat{m}_2 can be found by interchanging X_i and Z_i and the subscripts 1 and 2 in those for \hat{m}_1 .

THEOREM 4.1. *Suppose that Assumptions 1–4 hold. For the observation points (X_i, Z_i) , $i = 1, \dots, n$, the conditional bias and variance of $\hat{\alpha}$ and $\hat{m}_1(X_i)$ can be approximated by*

$$E(\hat{\alpha} - \alpha | \mathbf{X}, \mathbf{Z}) = \alpha + O_p\left(\frac{1}{\sqrt{n}}\right),$$

$$\begin{aligned} & E(\hat{m}_1(X_i) - m_1(X_i) | \mathbf{X}, \mathbf{Z}) \\ &= \frac{1}{(p_1+1)!} h_1^{p_1+1} \mu_{p_1+1}(K_{(p_1)}, X_i) m_1^{(p_1+1)}(X_i) \\ & \quad + \frac{1}{(p_1+1)!} h_1^{p_1+1} \mu_{p_1+1}(K_{(p_1)}) \left((\mathbf{t}_i^T - \mathbf{e}_i^T) D^{p_1+1} \mathbf{m}_1 - E(m_1^{(p_1+1)}(X_i)) \right) \\ & \quad - \frac{1}{(p_2+1)!} h_2^{p_2+1} \mu_{p_2+1}(K_{(p_2)}) \left(\mathbf{t}_i^T E(m_2^{(p_2+1)}(Z_i) | \mathbf{X}) - E(m_2^{(p_2+1)}(Z_i)) \right) \\ & \quad + O_p\left(\frac{1}{\sqrt{n}}\right) + o_p(h_1^{p_1+1} + h_2^{p_2+1}) \end{aligned}$$

and

$$\text{Var}(\hat{\alpha}) = \frac{\sigma^2}{n},$$

$$\text{Var}(\hat{m}_1(X_i) | \mathbf{X}, \mathbf{Z}) = \sigma^2 \frac{R(K_{(p_1)}, X_i)}{n h_1} f_X(X_i)^{-1} + o_p\left(\frac{1}{n h_1}\right).$$

The conditional bias and variance of $\hat{m}(X_i, Z_i)$ are

$$\begin{aligned} & E(\hat{m}(X_i, Z_i) - m(X_i, Z_i) | \mathbf{X}, \mathbf{Z}) \\ &= \frac{1}{(p_1+1)!} h_1^{p_1+1} \left(\mu_{p_1+1}(K_{(p_1)}, X_i) m_1^{(p_1+1)}(X_i) + \mu_{p_1+1}(K_{(p_1)}) \right. \\ & \quad \left. \times ((\mathbf{t}_i^T - \mathbf{e}_i^T) D^{p_1+1} \mathbf{m}_1 - \mathbf{v}_i^T E(m_1^{(p_1+1)}(X_i) | \mathbf{Z})) \right) \\ & \quad + \frac{1}{(p_2+1)!} h_2^{p_2+1} \left(\mu_{p_2+1}(K_{(p_2)}, Z_i) m_2^{(p_2+1)}(Z_i) + \mu_{p_2+1}(K_{(p_2)}) \right. \\ & \quad \left. \times ((\mathbf{t}_i^T - \mathbf{e}_i^T) D^{p_2+1} \mathbf{m}_2 - \mathbf{v}_i^T E(m_2^{(p_2+1)}(Z_i) | \mathbf{X})) \right) \\ & \quad + o_p(h_1^{p_1+1} + h_2^{p_2+1}) \end{aligned}$$

and

$$\begin{aligned} \text{Var}(\hat{m}(X_i, Z_i)|\mathbf{X}, \mathbf{Z}) &= \sigma^2 \left(\frac{R(K_{(p_1)}, X_i)}{nh_1} f_X(X_i)^{-1} + \frac{R(K_{(p_2)}, Z_i)}{nh_2} f_Z(Z_i)^{-1} \right) \\ &\quad + o_p \left(\frac{1}{nh_1} + \frac{1}{nh_2} \right). \end{aligned}$$

COROLLARY 4.1. *If the observation point (X_i, Z_i) lies in the interior of $\text{supp}(f)$, the conditional bias and variance of $\hat{m}_1(X_i)$ are approximated by*

$$\begin{aligned} E(\hat{m}_1(X_i) - m_1(X_i)|\mathbf{X}, \mathbf{Z}) &= \frac{1}{(p_1 + 1)!} h_1^{p_1+1} \mu_{p_1+1}(K_{(p_1)}) \left(\mathbf{t}_i^T D^{p_1+1} \mathbf{m}_1 - E(m_1^{(p_1+1)}(X_i)) \right) \\ &\quad - \frac{1}{(p_2 + 1)!} h_2^{p_2+1} \mu_{p_2+1}(K_{(p_2)}) \left(\mathbf{t}_i^T E(m_2^{(p_2+1)}(Z_i)|\mathbf{X}) - E(m_2^{(p_2+1)}(Z_i)) \right) \\ &\quad + O_p \left(\frac{1}{\sqrt{n}} \right) + o_p(h_1^{p_1+1} + h_2^{p_2+1}) \end{aligned}$$

and

$$\text{Var}(\hat{m}_1(X_i)|\mathbf{X}, \mathbf{Z}) = \sigma^2 \frac{R(K_{(p_1)})}{nh_1} f_X(X_i)^{-1} + o_p \left(\frac{1}{nh_1} \right).$$

The conditional bias and variance of $\hat{m}(X_i, Z_i)$ are

$$\begin{aligned} E(\hat{m}(X_i, Z_i) - m(X_i, Z_i)|\mathbf{X}, \mathbf{Z}) &= \frac{1}{(p_1 + 1)!} h_1^{p_1+1} \mu_{p_1+1}(K_{(p_1)}) \left(\mathbf{t}_i^T D^{p_1+1} \mathbf{m}_1 - \mathbf{v}_i^T E(m_1^{(p_1+1)}(X_i)|\mathbf{Z}) \right) \\ &\quad + \frac{1}{(p_2 + 1)!} h_2^{p_2+1} \mu_{p_2+1}(K_{(p_2)}) \left(\mathbf{v}_i^T D^{p_2+1} \mathbf{m}_2 - \mathbf{t}_i^T E(m_2^{(p_2+1)}(Z_i)|\mathbf{X}) \right) \\ &\quad + o_p(h_1^{p_1+1} + h_2^{p_2+1}) \end{aligned}$$

and

$$\begin{aligned} \text{Var}(\hat{m}(X_i, Z_i)|\mathbf{X}, \mathbf{Z}) &= \sigma^2 \left(\frac{R(K_{(p_1)})}{nh_1} f_X(X_i)^{-1} + \frac{R(K_{(p_2)})}{nh_2} f_Z(Z_i)^{-1} \right) \\ &\quad + o_p \left(\frac{1}{nh_1} + \frac{1}{nh_2} \right). \end{aligned}$$

A convenient error criterion that only uses the fitted values at the observation points is provided by the conditional mean averaged squared error (MASE), discussed by Härdle, Hall and Marron (1988). The MASE of m can

be written as

$$\begin{aligned} \text{MASE}(h_1, h_2|\mathbf{X}, \mathbf{Z}) &= \frac{1}{n} \sum_{i=1}^n E\{(\hat{m}(X_i, Z_i) - m(X_i, Z_i)|\mathbf{X}, \mathbf{Z})^2\} \\ &= \frac{1}{n} \sum_{i=1}^n \{E(\hat{m}(X_i, Z_i) - m(X_i, Z_i)|\mathbf{X}, \mathbf{Z})^2\} \\ &\quad + \frac{1}{n} \sum_{i=1}^n \text{Var}(\hat{m}(X_i, Z_i)|\mathbf{X}, \mathbf{Z}) \end{aligned}$$

and its asymptotic approximation is easily constructed from the previous results. To simplify the notation, let

$$\begin{aligned} \theta_{11}(r) &= \frac{1}{n} \sum_{i=1}^n (\mathbf{t}_i^T D^r \mathbf{m}_1 - \mathbf{v}_i^T E(m_1^{(r)}(X_i)|\mathbf{Z}))^2, \\ \theta_{22}(r) &= \frac{1}{n} \sum_{i=1}^n (\mathbf{v}_i^T D^r \mathbf{m}_2 - \mathbf{t}_i^T E(m_2^{(r)}(Z_i)|\mathbf{X}))^2 \end{aligned}$$

and

$$\theta_{12}(r, s) = \frac{1}{n} \sum_{i=1}^n (\mathbf{t}_i^T D^r \mathbf{m}_1 - \mathbf{v}_i^T E(m_1^{(r)}(X_i)|\mathbf{Z}))(\mathbf{v}_i^T D^s \mathbf{m}_2 - \mathbf{t}_i^T E(m_2^{(s)}(Z_i)|\mathbf{X})).$$

COROLLARY 4.2. *The conditional MASE for the bivariate additive model is approximated by*

$$\begin{aligned} &\text{MASE}(h_1, h_2|\mathbf{X}, \mathbf{Z}) \\ &= \left(\frac{\mu_{p_1+1}(K_{(p_1)})}{(p_1+1)!}\right)^2 h_1^{2p_1+2} \theta_{11}(p_1+1) + \left(\frac{\mu_{p_2+1}(K_{(p_2)})}{(p_2+1)!}\right)^2 h_2^{2p_2+2} \theta_{22}(p_2+1) \\ &\quad + \frac{\mu_{p_1+1}(K_{(p_1)})}{(p_1+1)!} \frac{\mu_{p_2+1}(K_{(p_2)})}{(p_2+1)!} h_1^{p_1+1} h_2^{p_2+1} \theta_{12}(p_1+1, p_2+1) \\ &\quad + \sigma^2 \frac{1}{n} \sum_{i=1}^n \left(\frac{R(K_{(p_1)})}{nh_1} f_X(X_i)^{-1} + \frac{R(K_{(p_2)})}{nh_2} f_Z(Z_i)^{-1}\right) \\ &\quad + o_p\left(h_1^{2p_1+2} + h_2^{2p_2+2} + \frac{1}{nh_1} + \frac{1}{nh_2}\right). \end{aligned}$$

If \mathbf{X} and \mathbf{Z} are independent, the preceding results can be simplified significantly. The expression for $\theta_{11}(r)$ becomes

$$\theta_{11}(r) = \frac{1}{n} \sum_{i=1}^n (m_1^{(r)}(X_i) - E(m_1^{(r)}(X_i)))^2$$

and similarly for $\theta_{22}(r)$. Since $\theta_{12}(r, s) = O_p(1/n)$ in this case, the term can be ignored.

COROLLARY 4.3. *If \mathbf{X} and \mathbf{Z} are independent, the conditional bias and variance of $\hat{\alpha}$ and $\hat{m}_1(X_i)$ in the interior of $\text{supp}(f)$ can be approximated by*

$$E(\hat{\alpha} - \alpha | \mathbf{X}, \mathbf{Z}) = \alpha + O_p\left(\frac{1}{\sqrt{n}}\right),$$

$$\begin{aligned} E(\hat{m}_1(X_i) - m_1(X_i) | \mathbf{X}, \mathbf{Z}) \\ &= \frac{1}{(p_1 + 1)!} h_1^{p_1+1} \mu_{p_1+1}(K_{(p_1)}) (m_1^{(p_1+1)}(X_i) - E(m_1^{(p_1+1)}(X_i))) \\ &\quad + O_p\left(\frac{1}{\sqrt{n}}\right) + o_p(h_1^{p_1+1} + h_2^{p_2+1}) \end{aligned}$$

and

$$\text{Var}(\hat{\alpha}) = \frac{\sigma^2}{n},$$

$$\text{Var}(\hat{m}_1(X_i) | \mathbf{X}, \mathbf{Z}) = \sigma^2 \frac{R(K_{(p_1)})}{nh_1} f_X(X_i)^{-1} + o_p\left(\frac{1}{nh_1}\right).$$

For $\hat{m}(X_i, Z_i)$,

$$\begin{aligned} E(\hat{m}(X_i, Z_i) - m(X_i, Z_i) | \mathbf{X}, \mathbf{Z}) \\ &= \frac{1}{(p_1 + 1)!} h_1^{p_1+1} \mu_{p_1+1}(K_{(p_1)}) (m_1^{(p_1+1)}(X_i) - E(m_1^{(p_1+1)}(X_i))) \\ &\quad + \frac{1}{(p_2 + 1)!} h_2^{p_2+1} \mu_{p_2+1}(K_{(p_2)}) (m_2^{(p_2+1)}(Z_i) - E(m_2^{(p_2+1)}(Z_i))) \\ &\quad + o_p(h_1^{p_1+1} + h_2^{p_2+1}), \end{aligned}$$

$$\begin{aligned} \text{Var}(\hat{m}(X_i, Z_i) | \mathbf{X}, \mathbf{Z}) &= \sigma^2 \left(\frac{R(K_{(p_1)})}{nh_1} f_X(X_i)^{-1} + \frac{R(K_{(p_2)})}{nh_2} f_Z(Z_i)^{-1} \right) \\ &\quad + o_p\left(\frac{1}{nh_1} + \frac{1}{nh_2}\right) \end{aligned}$$

and the conditional MASE is

$$\begin{aligned} \text{MASE}(h_1, h_2 | \mathbf{X}, \mathbf{Z}) \\ &= \left(\frac{\mu_{p_1+1}(K_{(p_1)})}{(p_1 + 1)!} \right)^2 h_1^{2p_1+2} \theta_{11}(p_1 + 1) + \left(\frac{\mu_{p_2+1}(K_{(p_2)})}{(p_2 + 1)!} \right)^2 h_2^{2p_2+2} \theta_{22}(p_2 + 1) \\ &\quad + \sigma^2 \frac{1}{n} \sum_{i=1}^n \left(\frac{R(K_{(p_1)})}{nh_1} f_X(X_i)^{-1} + \frac{R(K_{(p_2)})}{nh_2} f_Z(Z_i)^{-1} \right) \\ &\quad + o_p\left(h_1^{2p_1+2} + h_2^{2p_2+2} + \frac{1}{nh_1} + \frac{1}{nh_2}\right). \end{aligned}$$

REMARK 4.1. Theorem 4.1 shows that the additive model fitted by local polynomial regression of degree $p_1 = p_2 = p$ and computed through backfitting has the same $O_p(n^{-(p+1)/(2p+3)})$ rate of convergence as univariate local polynomial regression. However, there are a number of interesting differences between the asymptotic bias of the terms of the additive model and that of the (nonadditive) local polynomial regression estimator. In the case of \hat{m}_1 , the bias contains a term based on the curvature of m_1 ,

$$-\frac{1}{(p_1 + 1)!} \mu_{p_1+1}(K_{(p_1)}) h_1^{p_1+1} (\mathbf{t}_i^T E(m_1^{(p_1+1)}(X_i)|\mathbf{Z}) - E(m_1^{(p_1+1)}(X_i))),$$

another term based on the curvature of m_2 ,

$$-\frac{1}{(p_2 + 1)!} \mu_{p_2+1}(K_{(p_2)}) h_2^{p_2+1} (\mathbf{t}_i^T E(m_2^{p_2+1}(Z_i)|\mathbf{X}) - E(m_2^{p_2+1}(Z_i)))$$

and a third term caused only by the centering adjustment of \mathbf{S}_1^* , $O_p(1/\sqrt{n})$ (but note that the first two terms are also centered around their means). Only the first of these terms has an equivalence in univariate local polynomial regression. As shown in Corollary 4.3, the m_2 curvature term disappears if \mathbf{X} and \mathbf{Z} are independent, so that the asymptotic bias of \hat{m}_1 no longer depends on m_2 . The term $O_p(1/\sqrt{n})$ is not related to the dependence between \mathbf{X} and \mathbf{Z} . As shown in the proof of Theorem 4.1 (below), this centering bias of $\hat{m}_1(X_i)$ and $\hat{m}_2(Z_i)$ cancels out with the bias of $\hat{\alpha}$, so that this term does not appear in the bias of $\hat{m}(X_i, Z_i)$. Another difference with local polynomial regression is that the asymptotic bias at a point (X_i, Z_i) not only depends on the curvature of m_1 and m_2 at that point, but is a weighted average of the curvature at all the observation points, with the weights determined by the matrix $(\mathbf{I} - \mathbf{T}_{12}^*)^{-1}$. This difference again disappears when \mathbf{X} and \mathbf{Z} are independent.

REMARK 4.2. If \mathbf{X} and \mathbf{Z} are independent, Corollary 4.3 shows that \hat{m} has another interesting property. Suppose, for simplicity, that we are fitting the additive model by local polynomial regression of degree p . This additive model is unbiased as long as the unknown functions m_1 and m_2 are polynomials of degree less than or equal to $p + 1$. This differs from nonadditive local polynomial regression of degree p , which is only unbiased when the unknown function itself is of degree less than or equal to p . This effect is due to the centering adjustment, which replaces $m_1^{(p+1)}(X_i)$ in the nonadditive local polynomial regression bias by $(m_1^{(p+1)}(X_i) - E(m_1^{(p+1)}(X_i)))$. It is easy to see that the “centered” derivatives are indeed 0 for polynomials up to degree $p + 1$.

REMARK 4.3. Unlike the bias, the asymptotic variance terms are identical to those found in univariate local polynomial regression; that is, the asymptotic variance of the estimator of m_1 does not depend on the simultaneous estimation of m_2 . This somewhat surprising result is reminiscent of what happens in a two-way ANOVA design with an equal number of observations in each cell, where the effect for one factor can be estimated without adjusting for the other factor. The reason for this apparent independence is that any

product of the two smoothers $\mathbf{S}_1^* \mathbf{S}_2^*$ is of order $O_p(\mathbf{11}^T/n)$, as shown in Lemma 3.1, and hence is of smaller order than the leading variance terms. This result holds regardless of the dependence between \mathbf{X} and \mathbf{Z} , as long as restriction (4) is satisfied.

REMARK 4.4. We compare these asymptotic bias and variance results with those of Linton and Nielsen (1995) for the case $p_1 = p_2 = 1$. If we use f_Z as the weighting function q in their result, the theorem of Linton and Nielsen (1995) can be rewritten as

$$\begin{aligned} E(\hat{m}_1(X_i) - m_1(X_i)|\mathbf{X}, \mathbf{Z}) &= \frac{1}{2}\mu_2(K)h_1^2 m_1''(X_i) + \frac{1}{2}\mu_2(K)h_2^2 E(m_2''(Z_i)) + o_p(h_1^2 + h_2^2), \\ \text{Var}(\hat{m}_1(X_i)|\mathbf{X}, \mathbf{Z}) &= \sigma^2 R(K) \frac{1}{nh_1} E_Z(f_{X|Z}(X_i|Z)^{-1}) + o_p\left(\frac{1}{nh_1}\right). \end{aligned}$$

Both expressions ignore the boundary effects, so that Corollary 4.1 provides the relevant comparison. The rates of convergence for the bias and the variance are the same. The major differences in the bias are due to the effect in our results of the matrix $(\mathbf{I} - \mathbf{T}_{12}^*)^{-1}$, which makes the bias at X_i dependent on the curvature at all other observation points, and of the centering adjustment, which they do not account for. An interesting difference between the estimators occurs when \mathbf{X} and \mathbf{Z} are independent. It seems quite natural to expect from additive model estimators that, when the covariates are independent, the asymptotic bias for estimating one of the component functions does not depend on the behavior of the other function. As explained in Remark 4.1, the backfitting estimator indeed has this property, while the Linton–Nielsen estimator does not, as can readily be seen from the above bias expression. Unless the bias effects of the component functions happen to offset each other, this is likely to result in increased bias relative to the backfitting estimator. The comparison for the asymptotic variances is more straightforward, and, interestingly, the asymptotic variance of the backfitted estimators can be shown to be smaller than that of the “marginal integration” estimators, unless \mathbf{X} and \mathbf{Z} are independent. This is easily proven by noting that, in general,

$$E_Z(f_{X|Z}(X_i|Z)^{-1}) \geq f_X(X_i)^{-1}$$

by Jensen’s inequality for the function $h(x) = 1/x$. Since h is strictly convex, we get strict inequality unless $f_{Z|X}(X_i|z) = f_X(X_i)$ for almost all z . Thus, we get strict inequality for all X_i values in a set of positive probability unless \mathbf{X} and \mathbf{Z} are independent.

PROOF OF THEOREM 4.1. We first prove the theorem for the case $p_1 = p_2 = 1$. To simplify the notation, we suppress the fact that the bias and the variance we are approximating are conditional on \mathbf{X} and \mathbf{Z} .

For $\hat{\alpha}$, it is easy to see that

$$E(\hat{\alpha}) = \alpha + \bar{m}_1 + \bar{m}_2 = \alpha + O_p\left(\frac{1}{\sqrt{n}}\right).$$

In the case of $\hat{\mathbf{m}}_1$, we have

$$(6) \quad E(\hat{\mathbf{m}}_1) = (\mathbf{I} - (\mathbf{I} - \mathbf{S}_1^* \mathbf{S}_2^*)^{-1} (\mathbf{I} - \mathbf{S}_1^*)) (\boldsymbol{\alpha} + \mathbf{m}_1 + \mathbf{m}_2),$$

and we will apply the same Taylor expansion approximations as in Theorem 2.1 of Ruppert and Wand (1994). Let

$$\mathbf{Q}_{m_1}(x) = \begin{bmatrix} (X_1 - x)^2 \\ \vdots \\ (X_n - x)^2 \end{bmatrix} \frac{\partial^2 m_1(x)}{\partial x^2}$$

and

$$\mathbf{Q}_1 = \begin{bmatrix} \mathbf{s}_{1, X_1}^T \mathbf{Q}_{m_1}(X_1) \\ \vdots \\ \mathbf{s}_{1, X_n}^T \mathbf{Q}_{m_1}(X_n) \end{bmatrix},$$

and similarly for $\mathbf{Q}_{m_2}(z)$ and \mathbf{Q}_2 . Letting $\mathbf{h}_1^2 \equiv h_1^2 \mathbf{1}$, we can write

$$\mathbf{S}_2 \mathbf{m}_2 = \mathbf{m}_2 + \frac{1}{2} \mathbf{Q}_2 + o_p(\mathbf{h}_2^2)$$

and hence

$$(\mathbf{I} - \mathbf{S}_1^* \mathbf{S}_2^*)^{-1} (\mathbf{I} - \mathbf{S}_1^*) \mathbf{m}_2 = \mathbf{m}_2 + \frac{1}{2} (\mathbf{I} - \mathbf{S}_1^* \mathbf{S}_2^*)^{-1} \mathbf{S}_1^* \mathbf{Q}_2 + o_p(\mathbf{h}_2^2).$$

Similarly,

$$(\mathbf{I} - \mathbf{S}_1^* \mathbf{S}_2^*)^{-1} (\mathbf{I} - \mathbf{S}_1^*) \mathbf{m}_1 = \bar{\mathbf{m}}_1 - \frac{1}{2} (\mathbf{I} - \mathbf{S}_1^* \mathbf{S}_2^*)^{-1} \mathbf{Q}_1^* + o_p(\mathbf{h}_1^2),$$

where $\mathbf{Q}_1^* = (\mathbf{I} - \mathbf{1}\mathbf{1}^T/n) \mathbf{Q}_1$. Plugging these results into expression (6), we write the bias vector for $\hat{\mathbf{m}}_1$ as

$$(7) \quad E(\hat{\mathbf{m}}_1 - \mathbf{m}_1) = \frac{1}{2} (\mathbf{I} - \mathbf{S}_2^* \mathbf{S}_1^*)^{-1} (\mathbf{Q}_1^* - \mathbf{S}_1^* \mathbf{Q}_2) + O_p\left(\frac{1}{\sqrt{n}}\right) + o_p(\mathbf{h}_1^2 + \mathbf{h}_2^2).$$

The bias of $\hat{\mathbf{m}}_2$ and $\hat{\mathbf{m}}$ can be computed analogously. Note that the terms \bar{m}_1 and \bar{m}_2 will cancel in the expectation of $\hat{\mathbf{m}}$, so that the approximation term $O_p(\mathbf{1}/\sqrt{n})$ will not appear in that case.

The asymptotic bias in Theorem 2.2 of Ruppert and Wand (1994) can be rewritten as

$$\mathbf{s}_{1, x}^T \mathbf{Q}_{m_1}(x) = h_1^2 \mu_2(K, x) m_1''(x) + o_p(h_1^2).$$

By letting $\mathbf{M}_1 = \text{diag}\{\mu_2(K, X_1), \dots, \mu_2(K, X_n)\}$, we have therefore

$$\mathbf{Q}_1 = h_1^2 \mathbf{M}_1 \mathbf{D}^2 \mathbf{m}_1 + o_p(\mathbf{h}_1^2)$$

and

$$(8) \quad \mathbf{Q}_1^* = \mathbf{Q}_1 - \mu_2(K)h_1^2 E(\mathbf{m}_1''(X_i)) + o_p(\mathbf{h}_1^2).$$

Because

$$\int |(\mu_2(K, x) - \mu_2(K))m_1''(x)|f_X(x)dx = O_p(h_1),$$

the matrix of boundary moments \mathbf{M}_1 is replaced by $\mu_2(K)\mathbf{I}$ in the latter term. Similar expressions hold for \mathbf{Q}_2 and \mathbf{Q}_2^* . Next, using the fact that

$$[\mathbf{N}(x)^{-1}]_{11}\mu(K, x) + [\mathbf{N}(x)^{-1}]_{12}\mu_1(K, x) = 1,$$

one can compute

$$\mathbf{s}_{1,x}^T \mathbf{M}_2 D^2 \mathbf{m}_2 = \mu_2(K) E(m_2''(Z_i)|x) + o_p(1)$$

and hence

$$\mathbf{S}_1 \mathbf{Q}_2 = \mu_2(K)h_2^2 E(m_2''(Z_i)|\mathbf{X}) + o_p(\mathbf{h}_2^2)$$

and

$$(9) \quad \mathbf{S}_1^* \mathbf{Q}_2 = \mathbf{S}_1 \mathbf{Q}_2 - \mu_2(K)h_2^2 E(\mathbf{m}_2''(Z_i)) + o_p(\mathbf{h}_2^2).$$

Plugging in results (8) and (9), as well as Lemma 3.2 into the bias vector (7), we obtain

$$\begin{aligned} E(\hat{\mathbf{m}}_1 - \mathbf{m}_1) &= \frac{1}{2}h_1^2((\mathbf{I} - \mathbf{T}_{12}^*)^{-1}\mathbf{M}_1 D^2 \mathbf{m}_1 - \mu_2(K)E(m_1''(X_i))) \\ &\quad - \frac{1}{2}h_2^2\mu_2(K)((\mathbf{I} - \mathbf{T}_{12}^*)^{-1}E(m_2''(Z_i)|\mathbf{X}) - E(m_2''(Z_i))) \\ &\quad + O_p\left(\frac{1}{\sqrt{n}}\right) + o_p(\mathbf{h}_1^2 + \mathbf{h}_2^2) \\ &= \frac{1}{2}h_1^2(\mathbf{M}_1 D^2 \mathbf{m}_1 + \mu_2(K)((\mathbf{I} - \mathbf{T}_{12}^*)^{-1} - \mathbf{I})D^2 \mathbf{m}_1 - \mu_2(K)E(m_1''(X_i))) \\ &\quad - \frac{1}{2}h_2^2\mu_2(K)((\mathbf{I} - \mathbf{T}_{12}^*)^{-1}E(m_2''(Z_i)|\mathbf{X}) - E(m_2''(Z_i))) \\ &\quad + O_p\left(\frac{1}{\sqrt{n}}\right) + o_p(\mathbf{h}_1^2 + \mathbf{h}_2^2). \end{aligned}$$

The computations for the bias of $\hat{\mathbf{m}}$ are entirely analogous.

With \mathbf{W}_1 defined as in (2), the variance of $\hat{m}_1(X_i)$ is

$$\begin{aligned} \text{Var}(\hat{m}_1(X_i)) &= \sigma^2 \mathbf{e}_i^T \mathbf{W}_1 \mathbf{W}_1^T \mathbf{e}_i \\ (10) \quad &= \sigma^2 \{1 - 2\mathbf{e}_i^T (\mathbf{I} - \mathbf{S}_1^* \mathbf{S}_2^*)^{-1} (\mathbf{I} - \mathbf{S}_1^*) \mathbf{e}_i \\ &\quad + \mathbf{e}_i^T (\mathbf{I} - \mathbf{S}_1^* \mathbf{S}_2^*)^{-1} (\mathbf{I} - \mathbf{S}_1^*) (\mathbf{I} - \mathbf{S}_1^*)^T (\mathbf{I} - \mathbf{S}_1^* \mathbf{S}_2^*)^{-T} \mathbf{e}_i\}. \end{aligned}$$

Similar computations as in the proof of Lemma 3.1 lead to

$$\begin{aligned} [\mathbf{S}_1^*]_{ii} &= \frac{1}{nh_1} f_X(X_i)^{-1} [\mathbf{N}(X_i)^{-1}]_{11} K(0) - \frac{1}{n} + o_p\left(\frac{1}{n}\right) \\ &= \frac{1}{nh_1} f_X(X_i)^{-1} [\mathbf{N}(X_i)^{-1}]_{11} K(0) + o_p\left(\frac{1}{nh_1}\right) \end{aligned}$$

and show that, when we multiply \mathbf{S}_1^* by a matrix of order $O(\mathbf{11}^T/n)$, the resulting product is $o(\mathbf{11}^T/nh_1)$.

For the matrix $\mathbf{S}_1^* \mathbf{S}_1^{*T}$, note first that

$$\mathbf{S}_1^* \mathbf{S}_1^{*T} = \mathbf{S}_1 \mathbf{S}_1^T - \mathbf{11}^T/n + o_p(\mathbf{11}^T/n) = \mathbf{S}_1 \mathbf{S}_1^T + o_p\left(\frac{\mathbf{11}^T}{nh_1}\right).$$

Letting \approx denote equality up to order $(1 + o_p(1))$, the (i, j) th element of $\mathbf{S}_1 \mathbf{S}_1^T$ is

$$\begin{aligned} [\mathbf{S}_1 \mathbf{S}_1^T]_{ij} &\approx \frac{1}{n^2} \sum_{k=1}^n f_X(X_i)^{-1} f_X(X_j)^{-1} [\mathbf{N}(X_i)^{-1}]_{11} [\mathbf{N}(X_j)^{-1}]_{11} \\ &\quad \times K_{h_1}(X_k - X_i) K_{h_1}(X_k - X_j) \\ &+ \frac{1}{n^2} \sum_{k=1}^n \frac{1}{h_1} f_X(X_i)^{-1} f_X(X_j)^{-1} [\mathbf{N}(X_i)^{-1}]_{12} [\mathbf{N}(X_j)^{-1}]_{11} \\ &\quad \times K_{h_1}(X_k - X_i) (X_k - X_i) K_{h_1}(X_k - X_j) \\ &+ \frac{1}{n^2} \sum_{k=1}^n \frac{1}{h_1} f_X(X_i)^{-1} f_X(X_j)^{-1} [\mathbf{N}(X_i)^{-1}]_{11} [\mathbf{N}(X_j)^{-1}]_{12} \\ &\quad \times K_{h_1}(X_k - X_i) K_{h_1}(X_k - X_j) (X_k - X_j) \\ &+ \frac{1}{n^2} \sum_{k=1}^n \frac{1}{h_1^2} f_X(X_i)^{-1} f_X(X_j)^{-1} [\mathbf{N}(X_i)^{-1}]_{12} [\mathbf{N}(X_j)^{-1}]_{12} \\ &\quad \times K_{h_1}(X_k - X_i) (X_k - X_i) K_{h_1}(X_k - X_j) (X_k - X_j). \end{aligned}$$

When $i = j$,

$$[\mathbf{S}_1^* \mathbf{S}_1^{*T}]_{ii} = \frac{1}{nh_1} f_X(X_i)^{-1} R(K, X_i) + o_p\left(\frac{1}{nh_1}\right).$$

Let us write $H(u)$ for the function $uK(u)$. The elements of $\mathbf{S}_1^* \mathbf{S}_1^{*T}$ not on the diagonal are equal to

$$\begin{aligned} [\mathbf{S}_1^* \mathbf{S}_1^{*T}]_{ij} &\approx \frac{1}{nh_1} f_X(X_i)^{-1} [\mathbf{N}(X_i)^{-1}]_{11} [\mathbf{N}(X_j)^{-1}]_{11} K * K\left(\frac{X_j - X_i}{h_1}\right) \\ &+ \frac{1}{nh_1} f_X(X_i)^{-1} [\mathbf{N}(X_i)^{-1}]_{11} [\mathbf{N}(X_j)^{-1}]_{12} K * H\left(\frac{X_j - X_i}{h_1}\right) \end{aligned}$$

$$\begin{aligned}
& + \frac{1}{nh_1} f_X(X_i)^{-1} [\mathbf{N}(X_i)^{-1}]_{12} [\mathbf{N}(X_j)^{-1}]_{11} H * K \left(\frac{X_j - X_i}{h_1} \right) \\
& + \frac{1}{nh_1} f_X(X_i)^{-1} [\mathbf{N}(X_i)^{-1}]_{12} [\mathbf{N}(X_j)^{-1}]_{12} H * H \left(\frac{X_j - X_i}{h_1} \right) - \frac{1}{n},
\end{aligned}$$

so that $O_p(1/n) \sum_{j=1}^n [\mathbf{S}_1^{*T} \mathbf{S}_1^*]_{ij} = o_p(1/nh_1)$. Plugging these results and Lemma 3.2 into the terms of (10), we find

$$\text{Var}(\hat{m}_1(X_i)) = \sigma^2 \frac{R(K, X_i)}{nh_1} f_X(X_i)^{-1} + o_p\left(\frac{1}{nh_1}\right).$$

With \mathbf{W} as defined in (3), we can write the variance of $\hat{m}(X_i, Z_i)$ as

$$\begin{aligned}
(11) \quad \text{Var}(\hat{m}(X_i, Z_i)) &= \sigma^2 \mathbf{e}_i^T \mathbf{W} \mathbf{W}^T \mathbf{e}_i \\
&= \sigma^2 \{ 5/n + 4 - 4\mathbf{e}_i^T (\mathbf{I} - \mathbf{S}_1^* \mathbf{S}_2^*)^{-1} (\mathbf{I} - \mathbf{S}_1^*) \mathbf{e}_i \\
&\quad - 4\mathbf{e}_i^T (\mathbf{I} - \mathbf{S}_2^* \mathbf{S}_1^*)^{-1} (\mathbf{I} - \mathbf{S}_2^*) \mathbf{e}_i \\
&\quad - 2\mathbf{e}_i^T (\mathbf{I} - \mathbf{S}_1^* \mathbf{S}_2^*)^{-1} (\mathbf{I} - \mathbf{S}_1^*) \mathbf{1} \mathbf{1}^T / n \mathbf{e}_i \\
&\quad - 2\mathbf{e}_i^T (\mathbf{I} - \mathbf{S}_2^* \mathbf{S}_1^*)^{-1} (\mathbf{I} - \mathbf{S}_2^*) \mathbf{1} \mathbf{1}^T / n \mathbf{e}_i \\
&\quad + 2\mathbf{e}_i^T (\mathbf{I} - \mathbf{S}_1^* \mathbf{S}_2^*)^{-1} (\mathbf{I} - \mathbf{S}_1^*) (\mathbf{I} - \mathbf{S}_2^*)^T (\mathbf{I} - \mathbf{S}_2^* \mathbf{S}_1^*)^{-T} \mathbf{e}_i \\
&\quad + \mathbf{e}_i^T (\mathbf{I} - \mathbf{S}_1^* \mathbf{S}_2^*)^{-1} (\mathbf{I} - \mathbf{S}_1^*) (\mathbf{I} - \mathbf{S}_1^*)^T (\mathbf{I} - \mathbf{S}_1^* \mathbf{S}_2^*)^{-T} \mathbf{e}_i \\
&\quad + \mathbf{e}_i^T (\mathbf{I} - \mathbf{S}_2^* \mathbf{S}_1^*)^{-1} (\mathbf{I} - \mathbf{S}_2^*) (\mathbf{I} - \mathbf{S}_2^*)^T (\mathbf{I} - \mathbf{S}_2^* \mathbf{S}_1^*)^{-T} \mathbf{e}_i \}.
\end{aligned}$$

After computations entirely analogous to those in Lemma 3.1, we find that the (i, j) th element of $\mathbf{S}_1^* \mathbf{S}_2^{*T}$ is

$$[\mathbf{S}_1^* \mathbf{S}_2^{*T}]_{ij} = [\mathbf{T}_{12}^*]_{ij} + o\left(\frac{1}{n}\right) = o\left(\frac{1}{nh_1}\right) \quad \text{a.s.},$$

so that

$$\begin{aligned}
\text{Var}(\hat{m}(X_i, Z_i)) &= \sigma^2 \left\{ \frac{R(K, X_i)}{nh_1} f_X(X_i)^{-1} + \frac{R(K, Z_i)}{nh_2} f_Z(Z_i)^{-1} \right\} \\
&\quad + o_p\left(\frac{1}{nh_1} + \frac{1}{nh_2}\right).
\end{aligned}$$

The generalization to arbitrary odd p_1, p_2 is straightforward. Let $\mathbf{Q}_1^{(p_1+1)}$ and $\mathbf{Q}_2^{(p_2+1)}$ represent the higher order generalizations corresponding to \mathbf{Q}_1 and \mathbf{Q}_2 for the local linear case, and let $\mathbf{M}_{1, p_1} = \text{diag}\{\mu_{p_1+1}(K_{(p_1)}, X_1), \dots, \mu_{p_1+1}(K_{(p_1)}, X_n)\}$. To compute the bias of $\hat{\mathbf{m}}_1$, we note that (6)–(7) still hold after replacing the $\frac{1}{2}\mathbf{Q}_1$ and $\frac{1}{2}\mathbf{Q}_2$ by $(1/(p_1+1)!) \mathbf{Q}_1^{(p_1+1)}$ and $(1/(p_2+1)!) \mathbf{Q}_1^{(p_2+1)}$, respectively, and adjusting the orders of the approximations from h_1^2 and h_2^2 to

$h_1^{p_1+1}$ and $h_2^{p_2+1}$. Equations (8) and (9) can be generalized analogously. Therefore,

$$\begin{aligned} E(\hat{\mathbf{m}}_1 - \mathbf{m}_1) &= \frac{1}{(p_1 + 1)!} h_1^{p_1+1} (\mathbf{M}_{1, p_1} D^{p_1+1} \mathbf{m}_1 \\ &\quad + \mu_{p_1+1}(K_{(p_1)}) ((\mathbf{I} - \mathbf{T}_{12}^*)^{-1} - \mathbf{I}) D^{p_1+1} \mathbf{m}_1 \\ &\quad - \mu_{p_1+1}(K_{(p_1)}) E(m_1^{(p_1+1)}(X_i))) \\ &\quad - \frac{1}{(p_2 + 1)!} h_2^{p_2+1} \mu_{p_2+1}(K_{(p_2)}) ((\mathbf{I} - \mathbf{T}_{12}^*)^{-1} E(m_2^{(p_2+1)}(Z_i)|\mathbf{X}) \\ &\quad - E(m_2^{(p_2+1)}(Z_i))) \\ &\quad + O_p\left(\frac{1}{\sqrt{n}}\right) + o_p(\mathbf{h}_1^{p_1+1} + \mathbf{h}_2^{p_2+1}). \end{aligned}$$

The bias of $\hat{\mathbf{m}}$ is computed analogously.

For the variance, we note that (10) and (11) still hold. Using approximation (15), it is easy to compute that

$$\begin{aligned} [\mathbf{S}_1^*]_{ii} &= \frac{1}{nh_1} f_X(X_i)^{-1} [\mathbf{N}_{p_1}(X_i)^{-1}]_{11} K(0) + o_p\left(\frac{1}{nh_1}\right), \\ [\mathbf{S}_1^* \mathbf{S}_1^{*T}]_{ii} &= \frac{1}{nh_1} f_X(X_i)^{-1} R(K_{(p_1)}, X_i) + o_p\left(\frac{1}{nh_1}\right), \end{aligned}$$

with all other terms in (10) of order $o_p(1/nh_1)$, so that

$$\text{Var}(\hat{m}_1(X_i)) = \sigma^2 \frac{R(K_{(p_1)}, X_i)}{nh_1} f_X(X_i)^{-1} + o_p\left(\frac{1}{nh_1}\right).$$

A similar derivation leads to the desired expression for $\text{Var}(\hat{m}(X_i, Z_i))$. \square

5. Extension to local polynomials of even degree. We now consider the case when both of the degrees p_1 and p_2 are even. Let

$$C_{p_1}(z) = E(\mu_{p_1+1}(K_{(p_1)}, X_i) m_1^{(p_1+1)}(X_i) | z)$$

and similarly for $C_{p_2}(x)$. We define the matrices $\mathbf{F}_X = \text{diag}\{f'_X(X_1)/f_X(X_1), \dots, f'_X(X_n)/f_X(X_n)\}$, and similarly for \mathbf{F}_Z . \mathbf{M}_{1, p_1} is defined in the proof of Theorem 4.1, and \mathbf{M}_{2, p_2} is the corresponding matrix for the second covariate. We also replace Assumptions 3 and 4 by the following assumptions.

ASSUMPTION 3'. As $n \rightarrow \infty, h_1, h_2 \rightarrow \infty$ and $nh_1/\log n, nh_2/\log n, nh_1h_2/\log(n) \rightarrow \infty$.

ASSUMPTION 4'. The $(p_1 + 2)$ th derivative of m_1 and the $(p_2 + 2)$ th derivative of m_2 exist and are continuous and bounded.

We only state the asymptotic bias and variance of $\hat{m}_1(X_i)$. The proof of the following theorem can be found in Opsomer (1995).

THEOREM 5.1. *Assume that p_1 and p_2 are even and that Assumptions 1, 2, 3' and 4' hold. For the observation points (X_i, Z_i) , $i = 1, \dots, n$, the conditional bias and variance of $\hat{m}_1(X_i)$ can be approximated by*

$$\begin{aligned}
& E(\hat{m}_1(X_i) - m_1(X_i) | \mathbf{X}, \mathbf{Z}) \\
&= \frac{h_1^{p_1+1}}{(p_1+1)!} (\mathbf{t}_i^T \mathbf{M}_{1,p_1} D^{p_1+1} \mathbf{m}_1 - E(C_{p_1}(Z_i))) \\
&+ h_1^{p_1+2} \left(\mu_{p_1+2}(K_{(p_1)}, X_i) \left(\frac{f'_X(X_i)}{f_X(X_i)} \frac{m_1^{(p_1+1)}(X_i)}{(p_1+1)!} + \frac{m_1^{(p_1+2)}(X_i)}{(p_1+2)!} \right) \right. \\
&\quad + \mu_{p_1+2}(K_{(p_1)})(\mathbf{t}_i^T - \mathbf{e}_i^T) \left(\mathbf{F}_X \frac{D^{p_1+1} \mathbf{m}_1}{(p_1+1)!} + \frac{D^{p_1+2} \mathbf{m}_1}{(p_1+2)!} \right) \\
&\quad \left. - \mu_{p_1+2}(K_{(p_1)}) E \left(\frac{f'_X(X_i)}{f_X(X_i)} \frac{m_1^{(p_1+1)}(X_i)}{(p_1+1)!} + \frac{m_1^{(p_1+2)}(X_i)}{(p_1+2)!} \right) \right) \\
&- \frac{h_2^{p_2+1}}{(p_2+1)!} (\mathbf{t}_i^T C_{p_2}(\mathbf{X}) - E(C_{p_2}(X_i))) \\
&- h_2^{p_2+2} \mu_{p_2+2}(K_{(p_2)}) \left(\mathbf{t}_i^T E \left(\frac{f'_Z(Z_i)}{f_Z(Z_i)} \frac{m_2^{(p_2+1)}(Z_i)}{(p_2+1)!} + \frac{m_2^{(p_2+2)}(Z_i)}{(p_2+2)!} \middle| \mathbf{X} \right) \right. \\
&\quad \left. - E \left(\frac{f'_Z(Z_i)}{f_Z(Z_i)} \frac{m_2^{(p_2+1)}(Z_i)}{(p_2+1)!} + \frac{m_2^{(p_2+2)}(Z_i)}{(p_2+2)!} \right) \right) \\
&+ O_p \left(\frac{1}{\sqrt{n}} \right) + o_p(h_1^{p_1+2} + h_2^{p_2+2})
\end{aligned}$$

and

$$\text{Var}(\hat{m}_1(X_i) | \mathbf{X}, \mathbf{Z}) = \sigma^2 \frac{R(K_{(p_1)})}{nh_1} f_X(X_i)^{-1} + o_p \left(\frac{1}{nh_1} \right).$$

COROLLARY 5.1. *When \mathbf{X} and \mathbf{Z} are independent, the conditional bias of $\hat{m}_1(X_i)$ in the interior of $\text{supp}(f)$ is*

$$\begin{aligned}
& E(\hat{m}_1(X_i) - m_1(X_i) | \mathbf{X}, \mathbf{Z}) \\
&= -\frac{h_1^{p_1+1}}{(p_1+1)!} E(C_{p_1}(Z_i)) \\
&+ h_1^{p_1+2} \mu_{p_1+2}(K_{(p_1)}) \left(\left(\frac{f'_X(X_i)}{f_X(X_i)} \frac{m_1^{(p_1+1)}(X_i)}{(p_1+1)!} + \frac{m_1^{(p_1+2)}(X_i)}{(p_1+2)!} \right) \right. \\
&\quad \left. - E \left(\frac{f'_X(X_i)}{f_X(X_i)} \frac{m_1^{(p_1+1)}(X_i)}{(p_1+1)!} + \frac{m_1^{(p_1+2)}(X_i)}{(p_1+2)!} \right) \right) \\
&+ O_p \left(\frac{1}{\sqrt{n}} \right) + o_p(h_1^{p_1+2} + h_2^{p_2+2}).
\end{aligned}$$

REMARK 5.1. If X_i is in the interior, the bias of $\hat{m}_1(X_i)$ is of order $O_p(h_1^{p_1+2} + h_2^{p_2+2})$. To see this, note that $\mu_{p_1+1}(K_{(p_1)}) = 0$ for even p_1 , so that $\mathbf{t}_i^T \mathbf{M}_{1,p_1} D^{p_1+1} \mathbf{m}_1 = O_p(h_1)$ and $C_{p_1}(Z_i) = O_p(h_1)$ for interior X_i . This order is the same as that found in Theorem 4.1 of Ruppert and Wand (1994) for nonadditive local polynomial regression, but the bias in Theorem 5.1 contains several additional terms. If X_i is on the boundary, $\mathbf{t}_i^T \mathbf{M}_{1,p_1} D^{p_1+1} \mathbf{m}_1$ is the leading term, so that the bias is of order $O_p(h_1^{p_1+1} + h_2^{p_2+2})$. If \mathbf{X} and \mathbf{Z} are independent, most of the additional terms disappear, with the exception of those due to centering, as shown in Corollary 5.1.

APPENDIX A

Proofs of lemmas.

PROOF OF LEMMA 3.1. For simplicity, we prove the result for $p_1 = p_2 = 1$, and then show how it can be extended to arbitrary p_1 and p_2 . We first prove that the approximations in the lemma hold in probability. Let $\mathbf{A} \approx \mathbf{B}$ denote $\mathbf{A} = \mathbf{B}(1 + o_p(1))$ componentwise for any matrices \mathbf{A}, \mathbf{B} of the same dimension. Theorem 2.2 of Ruppert and Wand (1994) shows that, for any x ,

$$(12) \quad \left(\frac{1}{n} \mathbf{X}_x^T \mathbf{W}_x \mathbf{X}_x \right)^{-1} \approx f_X(x)^{-1} \mathbf{A}_1^{-1} \mathbf{N}(x)^{-1} \mathbf{A}_1^{-1},$$

where $\mathbf{A}_1 = \text{diag}\{1, h_1\}$. We can therefore write the (i, j) th element of \mathbf{S}_1 as

$$[\mathbf{S}_1]_{ij} \approx \frac{1}{n} f_X(X_i)^{-1} h_1^{-1} [\mathbf{N}(X_i)^{-1}]_{11} K\left(\frac{X_j - X_i}{h_1}\right) + \frac{1}{n} f_X(X_i)^{-1} h_1^{-1} [\mathbf{N}(X_i)^{-1}]_{12} K\left(\frac{X_j - X_i}{h_1}\right) \left(\frac{X_j - X_i}{h_1}\right).$$

Using standard results from density estimation, it is straightforward to compute that

$$(13) \quad [\mathbf{S}_1^*]_{ij} = [\mathbf{S}_1]_{ij} - \frac{1}{n} + o_p\left(\frac{1}{n}\right).$$

In the case of $\mathbf{S}_1 \mathbf{S}_2$, similar computations show that

$$[\mathbf{S}_1 \mathbf{S}_2]_{ij} \approx \frac{1}{n^2} f_X(X_i)^{-1} [\mathbf{N}(X_i)^{-1}]_{11} \sum_{k=1}^n f_Z(Z_k)^{-1} [\mathbf{N}(Z_k)^{-1}]_{11} \times K_{h_1}(X_k - X_i) K_{h_2}(Z_j - Z_k) + \frac{1}{n^2} f_X(X_i)^{-1} [\mathbf{N}(X_i)^{-1}]_{11} \sum_{k=1}^n f_Z(Z_k)^{-1} [\mathbf{N}(Z_k)^{-1}]_{12} K_{h_1}(X_k - X_i) \times \frac{1}{h_2} K_{h_2}(Z_j - Z_k)(Z_j - Z_k)$$

$$\begin{aligned}
& + \frac{1}{n^2} f_X(X_i)^{-1} [\mathbf{N}(X_i)^{-1}]_{12} \sum_{k=1}^n f_Z(Z_k)^{-1} [\mathbf{N}(Z_k)^{-1}]_{11} \frac{1}{h_1} K_{h_1}(X_k - X_i) \\
& \quad \times (X_k - X_i) K_{h_2}(Z_j - Z_k) \\
& + \frac{1}{n^2} f_X(X_i)^{-1} [\mathbf{N}(X_i)^{-1}]_{12} \sum_{k=1}^n f_Z(Z_k)^{-1} [\mathbf{N}(Z_k)^{-1}]_{12} \frac{1}{h_1} (X_k - X_i) \\
& \quad \times \frac{1}{h_2} K_{h_2}(Z_j - Z_k) (Z_j - Z_k) \\
& = \frac{1}{n} \frac{f(X_i, Z_j)}{f_X(X_i) f_Z(Z_j)} + o_p\left(\frac{1}{n}\right).
\end{aligned}$$

Using (13) and this result, we immediately find

$$(14) \quad [\mathbf{S}_1^* \mathbf{S}_2^*]_{ij} = [\mathbf{T}_{12}^*]_{ij} + o_p\left(\frac{1}{n}\right).$$

We now prove that (13) and (14) hold uniformly for all i, j and with $o_p(1)$ replaced by $o(1)$, using the theory from Chapter 2 of Pollard (1984) and two technical results which are proven in Appendix B. Expression (12) involves estimates of the moments of K up to order 2. By Assumptions 1 and 2 and Lemma B.1, the classes of graphs of the translation classes

$$\left\{ (\cdot - x)^t K\left(\frac{\cdot - x}{h_1}\right) : x \in \text{supp}(f_X) \right\} \quad \text{for } t = 0, 1, 2$$

have polynomial discrimination (see Appendix B for definitions of these terms). Theorem II.37 of Pollard (1994) therefore guarantees that the approximation (12) holds uniformly for all $x \in \text{supp}(f_X)$. For (13), the translation classes are

$$\left\{ f_X(\cdot)^{-1} [\mathbf{N}(\cdot)^{-1}]_{1t} K\left(\frac{\cdot - x}{h_1}\right) \left(\frac{\cdot - x}{h_1}\right)^{t-1} : x \in \text{supp}(f_X) \right\}$$

for $t = 1, 2$, and the same reasoning as above, combined with the fact that (12) holds uniformly, can be used to show that the approximation (13) holds uniformly over all i, j , proving the first part of the lemma.

Similarly, we consider (14). In this case, the translation classes are

$$\left\{ f_Z(\cdot - z) [\mathbf{N}(\cdot)^{-1}]_{1u} K\left(\frac{\cdot - z}{h_2}\right) \left(\frac{\cdot - z}{h_2}\right)^{u-1} K\left(\frac{\cdot - x}{h_1}\right) \left(\frac{\cdot - x}{h_1}\right)^{t-1} \right\}$$

for $t = 1, 2, u = 1, 2$. Since these translation classes are generated by the product of functions whose graphs have polynomial discrimination, Lemma 2.25 of Pollard (1984) and Lemma B.2 ensure that the assumptions of Theorem 2.37 of Pollard (1984) are satisfied. The approximation (14) therefore also holds uniformly over all i, j , completing the proof for the local linear case.

The generalization to arbitrary degree local polynomials is straightforward. The (i, j) th element of \mathbf{S}_1 is now

$$\begin{aligned} [\mathbf{S}_1]_{ij} &\approx \frac{1}{n} f_X(X_i)^{-1} \sum_{k=1}^{p_1+1} [\mathbf{N}_{p_1}(X_i)^{-1}]_{1k} K_{h_1}(X_j - X_i) \left(\frac{X_j - X_i}{h_1}\right)^{k-1} \\ (15) \qquad &\approx \frac{1}{nh_1} f_X(x)^{-1} K_{(p_1)}\left(\frac{X_j - x}{h_1}, x\right), \end{aligned}$$

so that again

$$[\mathbf{S}_1^*]_{ij} = [\mathbf{S}_1]_{ij} - \frac{1}{n} + o_p\left(\frac{1}{n}\right),$$

since $\mu(K_{(p)}, x) = 1$ for any x . As before, we can show that the classes of graphs

$$\left\{ K_{(p_1)}\left(\frac{\cdot - x}{h_1}\right) : x \in \text{supp}(f_X) \right\}$$

have polynomial discrimination for any p_1 using Pollard (1984) and the results in Appendix B, so that this approximation also holds uniformly over all i, j . The reasoning for $\mathbf{S}_1^* \mathbf{S}_2^*$ is completely analogous. \square

PROOF OF LEMMA 3.2. For any matrix \mathbf{A} , write $\rho(\mathbf{A})$ for the spectral radius of \mathbf{A} . Because of Assumption 2, we know that $\rho(\mathbf{T}_{12}^*) \leq \|\mathbf{T}_{12}^*\|_r < 1$ so that $(\mathbf{I} - \mathbf{T}_{12}^*)$ is invertible. By Lemma 3.1, it immediately follows that

$$P(\exists N: \rho(\mathbf{S}_1^* \mathbf{S}_2^*) < 1 \text{ if } n \geq N) = 1,$$

establishing the first part of the lemma.

Assume now that $(\mathbf{I} - \mathbf{S}_1^* \mathbf{S}_2^*)^{-1}$ exists. If we can show that $o(\mathbf{1}\mathbf{1}^T/n)(\mathbf{I} - \mathbf{T}_{12}^*)^{-1} = o(\mathbf{1}\mathbf{1}^T/n)$, then

$$(\mathbf{I} - \mathbf{S}_1^* \mathbf{S}_2^*)^{-1} = (\mathbf{I} - \mathbf{T}_{12}^*)^{-1} + o(\mathbf{1}\mathbf{1}^T/n) \quad \text{a.s.}$$

by the formula for the inverse of a sum of matrices [Horn and Johnson (1985), page 19]. Since $\rho(\mathbf{T}_{12}^*) < 1$, we can write

$$(\mathbf{I} - \mathbf{T}_{12}^*)^{-1} = \mathbf{I} + \sum_{p=1}^{\infty} \mathbf{T}_{12}^{*p}.$$

Using Assumption 2 again, it is easy to show by induction that, for all p , $\max_{i,j} \|[\mathbf{T}_{12}^{*p}]_{ij}\| \leq (1 - \varepsilon)^p/n$ for some $\varepsilon > 0$. Hence,

$$\max_{i,j} |(\mathbf{I} - \mathbf{T}_{12}^*)^{-1} - \mathbf{I}| \leq \frac{K}{n}$$

for all n and $o(\mathbf{1}\mathbf{1}^T/n)(\mathbf{I} - \mathbf{T}_{12}^*)^{-1} = o(\mathbf{1}\mathbf{1}^T/n)$ as desired. \square

APPENDIX B

Results from empirical process theory. We briefly state some of the definitions used in Pollard (1984). For a probability measure Q on S , a class of functions \mathcal{F} and $\varepsilon > 0$, define the *covering number* $N(\varepsilon, Q, \mathcal{F})$ as the smallest value of m for which there exist functions g_1, \dots, g_m (not necessarily in \mathcal{F}) such that $\min_j P|f - g_j| \leq \varepsilon$ for each f in \mathcal{F} . The *graph* of a real-valued function f on a set S , written $\text{gr}(f)$, is defined as the subset

$$\text{gr}(f) = \{(s, t): 0 \leq t \leq f(s) \text{ or } f(s) \leq t < 0\}$$

of $S \times \mathbb{R}$. A class \mathcal{D} of subsets of some space H has *polynomial discrimination* if there exists a polynomial $\rho(\cdot)$ such that, from every set of N points in H , the class picks out at most $\rho(N)$ distinct subsets.

We define a *translation class* of functions on g as the class $\{g(\cdot - x)\}$ for any function g on a set $S \subseteq \mathbb{R}$. It will be convenient to set $g(t) = 0$ for $t \notin S$, so that the domain of $g(t - x)$ does not depend on x . The set of graphs generated by the translation class on g will be written as \mathcal{S}_g . We also define a *monotonicity change* for a function g as (1) any point t_0 for which $g(t)$ changes from monotone increasing to monotone decreasing (or vice versa) in an interval $(t_0 - \varepsilon, t_0 + \varepsilon)$ for some $\varepsilon > 0$, and (2) any set of points (t_1, t_2) , $t_1 \neq t_2$, for which $g(t) = c$ for all $t \in (t_1, t_2)$ for some c and $g(t)$ changes from monotone increasing to monotone decreasing (or vice versa) in an interval $(t_1 - \varepsilon, t_2 + \varepsilon)$ for some $\varepsilon > 0$. We prove two lemmas.

LEMMA B.1. (i) *Suppose that the function $g(u)$ on $S \subseteq \mathbb{R}$ has a finite number of monotonicity changes. Then the set of graphs \mathcal{S}_g has polynomial discrimination in $S \times \mathbb{R}$.*

(ii) *If the function $h(u)$ has the same properties as g , then the following sets also have polynomial discrimination in $S \times \mathbb{R}$: \mathcal{S}_{gh} generated by the functions $\{g(\cdot - x)h(\cdot)\}$ and \mathcal{S}'_{gh} generated by $\{g(\cdot - x)h(\cdot - x)\}$.*

PROOF. Let us first look at the simplest possible case, where g does not change sign and has no monotonicity changes. We assume, without loss of generality, that $g \geq 0$ and is increasing. From a set containing only two points, \mathcal{S}_g can never pick out both singletons, since, for any $x_1 < x_2$, we have $g(s - x_1) \geq g(s - x_2)$ for all s , so that $\text{gr}(g(\cdot - x_2)) \subseteq \text{gr}(g(\cdot - x_1))$. By Lemma 2.17 of Pollard (1984), \mathcal{S}_g therefore has linear discrimination.

Suppose now that g is monotone. We can write

$$g = g^+ + g^- \equiv gI_{\{g>0\}} + gI_{\{g \leq 0\}}$$

so that

$$\text{gr}(g(\cdot - x)) = \text{gr}(g^+(\cdot - x)) \cup \text{gr}(g^-(\cdot - x)).$$

By Lemma 2.15 of Pollard (1984), we conclude that \mathcal{S}_g is also a polynomial class in this case.

Next, suppose that g has no sign changes but has exactly one monotonicity change. Without loss of generality, let $g \geq 0$. If the monotonicity changes from decreasing to increasing, then the graphs of g can be written as the union of graphs of two monotone functions. If the change is from increasing to decreasing, we can write it as the intersection of two graphs of monotone functions. Hence, Lemma 2.15 of Pollard (1984) can also be applied in this case.

Finally, consider any function g with a finite number of monotonicity changes. Begin by writing it as the sum of a positive and a nonnegative function, say g^+ and g^- . The graphs of each of these can be written as the finite union of graphs of functions with one monotonicity change as well as at most two graphs of monotone functions. Hence, for g^+ ,

$$\text{gr}(g^+(\cdot - x)) = \bigcup_{i=1}^{M_+} (\text{gr}(g_{i1}^+(\cdot - x)) \cap \text{gr}(g_{i2}^+(\cdot - x))),$$

where some of the sets $\text{gr}(g_{ij}^+(\cdot - x))$ can be empty sets for $i = 1$ and $i = M_+$. A similar expression holds for g^- . Using Lemma 2.15 of Pollard (1984) again proves result (i).

Result (ii) follows immediately from the fact that if the derivative of a function g has a finite number of sign changes, then so does the function itself. \square

LEMMA B.2. *Let \mathcal{S}, \mathcal{R} represent two classes of bounded, real-valued functions on S and T , respectively. Suppose there exists constant $A_g, A_r, w_g, w_r > 0$ and the covering numbers of \mathcal{S}, \mathcal{R} satisfy*

$$\begin{aligned} N(\varepsilon, P, \mathcal{S}) &\leq A_g \varepsilon^{-w_g} \quad \text{for } 0 < \varepsilon < 1, \\ N(\varepsilon, P, \mathcal{R}) &\leq A_r \varepsilon^{-w_r} \quad \text{for } 0 < \varepsilon < 1 \end{aligned}$$

for any probability measure P . Let $f(u, t) = g(u)r(t)$, the product function on $S \times T$. The covering numbers of the class $\mathcal{F} = \{f: f = g \times r, g \in \mathcal{S}, r \in \mathcal{R}\}$ satisfy

$$N(\varepsilon, P, \mathcal{F}) \leq A_f \varepsilon^{-w_f} \quad \text{for } 0 < \varepsilon < 1$$

for some A_f, w_f .

PROOF. Let G represent a function for which $|g| \leq G$ for all $g \in \mathcal{S}$, and R a similar function for \mathcal{R} . Let $M_g = \max |G|, M_r = \max |R|$. We need to find how many functions f_i are required for each $\varepsilon > 0$, so that

$$\min_{f_i} P|f - f_i| < \varepsilon$$

for all f , or, equivalently, how many pairs of functions (g_k, r_l) are required, so that

$$\min_{(g_k, r_l)} P|g \times r - g_k \times r_l| < \varepsilon$$

for all (g, r) . Since

$$P|g \times r - g_k \times r_l| \leq M_r P|g - g_k| + M_g P|r - r_l|,$$

if we require that

$$\begin{aligned} \min_{g_k} P|g - g_k| &< \frac{\varepsilon}{2M_r}, \\ \min_{r_l} P|r - r_l| &< \frac{\varepsilon}{2M_g}, \end{aligned}$$

the number of functions g_k, r_l required to achieve this are $A_g(\varepsilon/2M_r)^{-w_g}$ and $A_r(\varepsilon/2M_g)^{-w_r}$, respectively. Hence,

$$N(\varepsilon, P, \mathcal{F}) \leq A_f \varepsilon^{-w_f},$$

with $A_f = A_g A_r (2M_r)^{w_g} (2M_g)^{w_r}$ and $w_f = w_g + w_r$. \square

REFERENCES

- BUJA, A., HASTIE, T. J. and TIBSHIRANI, R. J. (1989). Linear smoothers and additive models. *Ann. Statist.* **17** 453–555.
- BURMAN, P. (1990). Estimation of generalized additive models. *J. Multivariate Anal.* **32** 230–255.
- CHAMBERS, J. M. and HASTIE, T. J., eds. (1992). *Statistical models in S*. Wadsworth and Brooks/Cole, Pacific Grove, CA.
- CHEN, Z. (1993). Fitting multivariate regression functions by interaction spline models. *J. Roy. Statist. Soc. Ser. B* **55** 473–491.
- CLEVELAND, W. and DEVLIN, S. (1988). Locally weighted regression: an approach to regression analysis by local fitting. *J. Amer. Statist. Assoc.* **83** 596–610.
- FAN, J., GASSER, T., GJJBELS, I., BROCKMANN, M. and ENGEL, J. (1993). Local polynomial fitting: a standard for nonparametric regression. Mimeo Series 2302, Institute of Statistics, Univ. North Carolina, Chapel Hill.
- FAN, J. and GJJBELS, I. (1994). Adaptive order polynomial fitting: bandwidth robustification and bias reduction. Unpublished manuscript.
- FRIEDMAN, J. H. and STUETZLE, W. (1981). Projection pursuit regression. *J. Amer. Statist. Assoc.* **76** 817–823.
- GU, C., BATES, D. M., CHEN, Z. and WAHBA, G. (1989). The computation of gev functions through householder tridiagonalization with application to the fitting of interaction spline models. *SIAM J. Matrix Anal. Appl.* **10** 457–480.
- HÄRDLE, W., HALL, P. and MARRON, J. S. (1988). How far are automatically chosen regression smoothing parameters from their optimum? *J. Amer. Statist. Assoc.* **83** 86–95.
- HÄRDLE, W. and HALL, P. (1993). On the backfitting algorithm for additive regression models. *Statist. Neerlandica* **47** 43–57.
- HASTIE, T. J. and TIBSHIRANI, R. J. (1990). *Generalized Additive Models*. Chapman and Hall, Washington, DC.
- HORN, R. A. and JOHNSON, C. A. (1985). *Matrix Analysis*. Cambridge Univ. Press.
- LINTON, O. and NIELSEN, J. P. (1995). A kernel method of estimating structured nonparametric regression based on marginal integration. *Biometrika* **82** 93–100.
- OPSOMER, J.-D. (1995). Optimal bandwidth selection for fitting an additive model by local polynomial regression. Ph.D. dissertation, Cornell Univ.
- OPSOMER, J.-D. and RUPPERT, D. (1995). A fully automated bandwidth selection method for fitting additive models. Preprint 95-32, Dept. Statistics, Iowa State Univ.
- POLLARD, D. (1984). *Convergence of Stochastic Processes*. Springer, New York.

- RUPPERT, D. and WAND, M. P. (1994). Multivariate locally weighted least squares regression. *Ann. Statist.* **22** 1346–1370.
- STONE, C. J. (1985). Additive regression and other nonparametric models. *Ann. Statist.* **13** 689–705.
- WAHBA, G. (1986). Partial and interaction splines for the semiparametric estimation of functions of several variables. In *Computer Science and Statistics: Proceedings of the 18th Symposium on the Interface* (T. J. Boardman, ed.) 75–80. Amer. Statist. Assoc., Alexandria, VA.

DEPARTMENT OF STATISTICS
SNEDECOR HALL
IOWA STATE UNIVERSITY
AMES, IOWA 50011-1210

COLLEGE OF ENGINEERING
RHODES HALL
CORNELL UNIVERSITY
ITHACA, NEW YORK 14853-3801