

# **$L_2$ Rate of Convergence for Interaction Spline Regression**

**By**

**Charles J. Stone**

**Technical Report No. 268**

**August 1990**

**This research was supported in part by NSF Grant DMS-8902016.**

**Department of Statistics  
University of California  
Berkeley, California**

# $L_2$ RATE OF CONVERGENCE FOR INTERACTION SPLINE REGRESSION<sup>1</sup>

BY CHARLES J. STONE

*University of California, Berkeley*

Let  $X_1, \dots, X_N, Y$  be random variables with  $E(Y^2) < \infty$  and let  $\mu(\cdot)$  denote the regression function of  $Y$  on  $\mathbf{X} = (X_1, \dots, X_N)$ . Let  $\mathcal{H}$  be the space of functions containing the constant and main-effect components and, possibly, some interaction components and let  $d$  be the maximum number of variables involved in any such component. Let  $\mu^*(\cdot)$  be the best approximation in  $\mathcal{H}$  to  $\mu(\cdot)$ :  $E[(\mu(\mathbf{X}) - \mu^*(\mathbf{X}))^2] \leq E[(\mu(\mathbf{X}) - h(\mathbf{X}))^2]$  for  $h \in \mathcal{H}$ . Let  $p$  be a suitably defined lower bound to the smoothness of  $\mu^*(\cdot)$ . Consider a random sample of size  $n$  from the joint distribution of  $\mathbf{X}$  and  $Y$ . Let  $\mathcal{G}$  be a suitably defined finite-dimensional subspace of  $\mathcal{H}$  consisting of splines. The dimension of  $\mathcal{G}$  is allowed to tend to infinity along with  $n$ . Let  $\hat{\mu}(\cdot)$  denote the least-squares estimate in  $\mathcal{G}$  of  $\mu^*(\cdot)$  based on the random sample. Under suitable conditions, the  $L_2$  rate of convergence of  $\hat{\mu}(\cdot)$  to  $\mu^*(\cdot)$  is  $n^{-p/(2p+d)}$  and optimal.

---

<sup>1</sup>This research was supported in part by National Science Foundation Grant DMS-8902016.

*AMS 1980 subject classifications.* Primary 62G20; secondary 62G05.

*Key words and phrases.* Functional inference, optimal rate of convergence, least-squares estimation,  $B$ -splines.

**1. Introduction.** Consider random variables  $X_1, \dots, X_N, Y$  and let  $\mu(\cdot)$  denote the regression function of  $Y$  on  $\mathbf{X} = (X_1, \dots, X_N)$ , so that  $\mu(\mathbf{x}) = E(Y|\mathbf{X} = \mathbf{x})$ . Let  $p$  denote the number of derivatives of the regression function (precise definitions will be given in Section 2). Then, under suitable conditions, the optimal  $L_2$  rate of convergence on compacts for estimating the regression function based on a random sample of size  $n$  from the joint distribution of  $\mathbf{X}$  and  $Y$  is  $n^{-p/(2p+N)}$  (see Stone, 1982). One formulation of the "curse of dimensionality" is that, for fixed  $p$ ,  $p/(2p+N)$  is close to zero when  $N$  is large.

To get a faster rate of convergence, we could assume that  $\mu(x_1, \dots, x_N)$  is an additive function of  $x_1, \dots, x_N$  or, more realistically, replace the goal of estimating the regression function by that of estimating the best additive approximation  $\mu^*(\cdot)$  to this function. The optimal rate of convergence is now  $n^{-p/(2p+1)}$  (see Stone, 1985). In effect, additivity reduces the dimensionality of the estimation problem from  $N$  to 1.

More generally, we could allow some low-order interaction components into the model. It is natural to conjecture that, under suitable conditions, the optimal rate of convergence should now be  $n^{-p/(2p+d)}$  with  $d$  being the maximum number of variables in any component. Suppose, for example, that  $N = 3$  and let  $\mu^*(\cdot)$  be the best approximation to the regression function of the form

$$\mu^*(x_1, x_2, x_3) = \mu_0^* + \mu_1^*(x_1) + \mu_2^*(x_2) + \mu_3^*(x_3) + \mu_{12}^*(x_1, x_2) + \mu_{13}^*(x_1, x_3),$$

which is a hierarchical model that includes the constant effect, all three main effects and two of the three two-factor interactions. The optimal rate of convergence should now be  $n^{-p/(2p+2)}$ . One main purpose of the present paper is to verify the stated conjecture.

In order to achieve the rate of convergence  $n^{-p/(2p+d)}$ , we will construct finite-dimensional spline spaces that can be used both to approximate  $\mu^*(\cdot)$  and, in conjunction with the method of least squares, to estimate this function from the sample data. Thereby, we will accomplish the second main purpose of this paper, which is to lend further theoretical support to the use of finite-dimensional spline spaces in functional inference.

2. **Statement of results.** Consider random variables  $X_1, \dots, X_N, Y$ , where  $X_1, \dots, X_N$  are  $[0, 1]$ -valued and  $Y$  has finite mean. Then  $\mathbf{X} = (X_1, \dots, X_N)$  ranges over  $C = [0, 1]^N$ . It is supposed that the following condition is satisfied.

CONDITION 1.  $\mathbf{X}$  has a density  $f$  that is bounded away from zero and infinity on  $C$ .

Let  $M_1$  and  $M_2$  be positive numbers such that  $M_1^{-1} \leq f \leq M_2$  on  $C$ . Then  $M_1, M_2 \geq 1$ .

Set

$$\langle h_1, h_2 \rangle = E[h_1(\mathbf{X})h_2(\mathbf{X})] = \int_C h_1(\mathbf{x})h_2(\mathbf{x})f(\mathbf{x})d\mathbf{x}$$

and

$$\|h\|^2 = \langle h, h \rangle = E[h^2(\mathbf{X})] = \int_C h^2(\mathbf{x})f(\mathbf{x})d\mathbf{x}$$

for square-integrable functions  $h_1, h_2, h$  on  $C$ . Two such functions are regarded as being equal if they differ only on a set of Lebesgue measure zero. Let  $\mu(\cdot)$  denote the regression function of  $Y$  on  $\mathbf{X}$ , which is defined by  $\mu(\mathbf{x}) = E(Y | \mathbf{X} = \mathbf{x})$  for  $\mathbf{x} \in C$ .

Given a subset  $v$  of  $\{1, \dots, N\}$ , let  $\mathcal{H}_v$  denote the collection of all square-integrable functions  $h$  on  $C$  that depend only on the coordinates  $x_l, l \in v$ , of  $\mathbf{x} = (x_1, \dots, x_N)$ . We refer to  $\#(v) - 1$  as the *interaction order* of  $\mathcal{H}_v$ . (If  $v$  is the empty set  $\emptyset$  or, equivalently, if  $\#(v) = 0$ , then  $\mathcal{H}_v$  is the space  $\mathcal{C}$  of constant functions on  $C$ .)

Let  $\mathcal{N}$  be a collection of subsets of  $\{1, \dots, N\}$  and set

$$\mathcal{H} = \sum_v \mathcal{H}_v = \sum_{v \in \mathcal{N}} \mathcal{H}_v = \left\{ \sum_v h_v : h_v \in \mathcal{H}_v \text{ for } v \in \mathcal{N} \right\}$$

and  $d = \max_{v \in \mathcal{N}} \#(v)$ . Then  $d-1$  is the maximum interaction order of the components of  $\mathcal{H}$ . Observe that  $d = 1$  if and only if every function in  $\mathcal{H}$  is additive. It is assumed that  $\mathcal{H}$  is *hierarchical*: if  $v$  is in  $\mathcal{N}$  and  $\eta$  is a subset of  $v$ , then  $\eta$  is in  $\mathcal{N}$ . Set

$$\mathcal{H}_v^0 = \{h \in \mathcal{H}_v : h \perp \mathcal{H}_\eta \text{ for every proper subset } \eta \text{ of } v\}, \quad v \in \mathcal{N}.$$

(Here  $h \perp \mathcal{H}_\eta$  means that  $\langle h, k \rangle = 0$  for  $k \in \mathcal{H}_\eta$ .) Then (under Condition 1)

$$\mathcal{H} = \bigoplus_v \mathcal{H}_v^0,$$

each  $h \in \mathcal{H}$  can be written uniquely in the form  $h = \sum_v h_v$ , where  $h_v \in \mathcal{H}_v^0$  for  $v \in \mathcal{N}$ ; clearly,  $h_\emptyset = E[h(\mathbf{X})]$ . (See Section 3 for the proof). We refer to  $\mathcal{H}_v^0, v \in \mathcal{N}$ , as the components of  $\mathcal{H}$ , to  $\mathcal{H}_\emptyset = \mathcal{C}$  as the constant component, to  $\mathcal{H}_v$  with  $\#(v) = 1$  as a main-

effect component, and to  $\mathcal{H}_v$  with  $\#(v) \geq 2$  as an interaction component. There is a unique best approximation  $\mu^*(\cdot)$  in  $\mathcal{H}$  to  $\mu(\cdot)$ :

$$E[(\mu(X) - \mu^*(X))^2] = \min_{h \in \mathcal{H}} E[(\mu(X) - h(X))^2].$$

(This follows from Lemma 1 in Section 3 by a standard completeness argument in the context of Hilbert space.) We can write  $\mu^*(\cdot) = \sum_v \mu_v^*(\cdot)$  for uniquely determined  $\mu_v^*(\cdot) \in \mathcal{H}_v$ ,  $v \in \mathcal{N}$ ; clearly  $\mu_{\emptyset}^*(\cdot) = E\mu^*(X) = E\mu(X) = EY$ . Observe that  $\mu^*(\cdot) = \mu(\cdot)$  if and only if  $\mu(\cdot) \in \mathcal{H}$ .

Let  $(X_1, Y_1), \dots, (X_n, Y_n)$  be a random sample of size  $n$  from the joint distribution of  $X$  and  $Y$  and set  $\bar{Y} = (Y_1 + \dots + Y_n)/n$ . It follows from Condition 1 that  $X_1, \dots, X_n$  are distinct (with probability one). Set  $C' = \{X_1, \dots, X_n\}$ . Consider the space of all real-valued functions whose domain includes  $C'$ . Let  $\langle \cdot, \cdot \rangle_n$  denote the semi-inner product on this space defined by

$$\langle h_1, h_2 \rangle_n = \frac{1}{n} \sum_i h_1(X_i) h_2(X_i)$$

and let  $\|\cdot\|_n$  denote the corresponding seminorm:  $\|h\|_n^2 = \langle h, h \rangle_n$ . Then  $\|1\|_n^2 = 1$ .

Let  $K = K_n$  denote a positive integer and let  $I_k$ ,  $1 \leq k \leq K$ , denote the subintervals of  $[0, 1]$  defined by  $I_k = [(k-1)/K, k/K)$  for  $1 \leq k < K$  and  $I_K = [1-1/K, 1]$  for  $k = K$ . Let  $m$  and  $q$  be fixed integers such that  $m \geq 0$  and  $m > q$ . Let  $\mathcal{S} = \mathcal{S}_n$  denote the collection of functions  $s$  on  $[0, 1]$  such that

(i) the restriction of  $s$  to  $I_k$  is a polynomial of degree  $m$  (or less) for  $1 \leq k \leq K$ ;

and, if  $q \geq 0$ ,

(ii)  $s$  is  $q$ -times continuously differentiable on  $[0, 1]$ .

A function satisfying (i) is called a piecewise polynomial; if  $m = 0$ , it is piecewise constant. A function satisfying (i) and (ii) is called a spline. Typically, splines are considered with  $q = m-1$  and then called linear, quadratic or cubic splines according as  $m = 1, 2$ , or  $3$ . Let  $B_j$ ,  $1 \leq j \leq J$ , denote the usual basis of  $\mathcal{S}$  consisting of  $B$ -splines (see de Boor, 1978). Then, in particular,  $B_j \geq 0$  on  $[0, 1]$  for  $1 \leq j \leq J$  and  $\sum B_j = 1$  on  $[0, 1]$ . Observe that  $K \leq J \leq (m+1)K$ .

Given a subset  $v$  of  $\{1, \dots, N\}$ , let  $\mathcal{S}_v$  denote the corresponding *interaction spline*

space, defined as the span of all functions  $g$  on  $C$  of the form

$$g(\mathbf{x}) = \prod_{l \in \nu} g_l(x_l), \quad \text{where } \mathbf{x} = (x_1, \dots, x_N) \text{ and } g_l \in \mathcal{S} \text{ for } l \in \nu.$$

Then  $\mathcal{G}_\nu$  has dimension  $J^{\#(\nu)}$ , where  $\#(\nu)$  is the number of integers in  $\nu$ . (In particular,  $\mathcal{G}_\emptyset = \mathcal{S}$  has dimension  $J^0 = 1$ .) Set  $\mathcal{G} = \sum_\nu \mathcal{G}_\nu = \{\sum_\nu g_\nu : g_\nu \in \mathcal{G}_\nu \text{ for } \nu \in \mathcal{N}\}$ . Also, set

$$\mathcal{G}_\nu^0 = \{g \in \mathcal{G}_\nu : g \perp_n \mathcal{G}_\eta \text{ for every proper subset } \eta \text{ of } \nu\}, \quad \nu \in \mathcal{N}.$$

(Here  $g \perp_n \mathcal{G}_\eta$  means that  $\langle g, h \rangle_n = 0$  for  $h \in \mathcal{G}_\eta$ .) Then  $\mathcal{G} = \sum_\nu \mathcal{G}_\nu^0$

The space  $\mathcal{G}$  is said to be *singular* (with respect to  $C'$ ) if there is a nonzero function  $g \in \mathcal{G}$  such that  $g = 0$  on  $C'$ ; otherwise,  $\mathcal{G}$  is said to be *nonsingular*. Suppose  $\mathcal{G}$  is nonsingular. Then  $\langle \cdot, \cdot \rangle_n$  is an inner product on  $\mathcal{G}$  and  $\| \cdot \|$  is a norm on  $\mathcal{G}$ ; that is,  $\|g\|_{2n} > 0$  for every nonzero function  $g \in \mathcal{G}$ . Moreover (see Lemma 2 in Section 3),  $\mathcal{G} = \bigoplus_\nu \mathcal{G}_\nu^0$ ; each  $g \in \mathcal{G}$  can be written uniquely in the form  $g = \sum_\nu g_\nu$  where  $g_\nu \in \mathcal{G}_\nu^0$  for  $\nu \in \mathcal{N}$ ; clearly,  $g_\emptyset = n^{-1} \langle 1, g \rangle$ .

Set  $d_1 = \max\{\#(\eta \cup \nu) : \eta, \nu \in \mathcal{N}\}$ . Then  $d \leq d_1 \leq 2d$ .

CONDITION 2.  $J^{d_1} = o(n^{1-\delta})$  for some  $\delta > 0$ .

The next result follows from Lemmas 4 and 5 in Section 3.

THEOREM 1. Suppose Conditions 1 and 2 hold. Then  $P(\mathcal{G} \text{ is singular}) = o(1)$ .

Let  $Y(\cdot)$  be defined by  $Y(\mathbf{X}_i) = Y_i$  for  $1 \leq i \leq n$ . Let  $\hat{\mu}(\cdot) = \sum_{\nu \in \mathcal{N}} \hat{\mu}_\nu(\cdot)$ , where  $\hat{\mu}_\nu(\cdot) \in \mathcal{G}_\nu^0$  for  $\nu \in \mathcal{N}$ , minimize  $\|Y(\cdot) - g\|_n^2 = n^{-1} \sum_1^n [Y_i - g(\mathbf{X}_i)]^2$ ,  $g \in \mathcal{G}$ . Then  $\hat{\mu}(\cdot)$  is the least-squares fit in  $\mathcal{G}$  to the sample data and  $\hat{\mu}_\emptyset(\cdot) = \bar{Y}$ . We think of  $\hat{\mu}(\cdot)$  as an estimate of  $\mu^*(\cdot)$  and of  $\hat{\mu}_\nu(\cdot)$  as an estimate of  $\mu_\nu^*(\cdot)$  for  $\nu \in \mathcal{N}$ . If  $\mathcal{G}$  is nonsingular, then  $\hat{\mu}(\cdot)$  and  $\hat{g}_\nu(\cdot)$ ,  $\nu \in \mathcal{N}$ , are uniquely determined.

CONDITION 3. The function  $E(Y^2 | \mathbf{X} = \mathbf{x})$ ,  $\mathbf{x} \in C$ , is bounded.

Given the positive number  $b_n$  and the random variable  $Z_n$  for  $n \geq n_0$ ,  $Z_n = O_p(b_n)$  means that  $\lim_{c \rightarrow \infty} \limsup_n P(|Z_n| > cb_n) = 0$ .

THEOREM 2. Suppose Conditions 1–3 hold. Then

$$\sup_{\mathbf{x} \in C} \text{var}(\hat{\mu}_\nu(\mathbf{x}) | \mathbf{X}_1, \dots, \mathbf{X}_n) = O_p(J^d/n), \quad \nu \in \mathcal{N},$$

so

$$\sup_{\mathbf{x} \in C} \text{var}(\hat{\mu}(\mathbf{x}) | \mathbf{X}_1, \dots, \mathbf{X}_n) = O_p(J^d/n).$$

Let  $0 < \beta \leq 1$ . A function  $h$  on  $C$  is said to satisfy a Hölder condition with exponent  $\beta$  if there is a positive number  $M$  such that  $|h(\mathbf{x}) - h(\mathbf{x}_0)| \leq M|\mathbf{x} - \mathbf{x}_0|^\beta$  for  $\mathbf{x}_0, \mathbf{x} \in C$ ; here  $|\mathbf{x}|$  is the Euclidean norm  $(x_1^2 + \cdots + x_N^2)^{1/2}$  of  $\mathbf{x} = (x_1, \dots, x_N)$ . Given an  $N$ -tuple  $\alpha = (\alpha_1, \dots, \alpha_N)$  of nonnegative integers, set  $[\alpha] = \alpha_1 + \cdots + \alpha_N$  and let  $D^\alpha$  denote the differentiable operator defined by

$$D^\alpha = \frac{\partial^{[\alpha]}}{\partial x_1^{\alpha_1} \cdots \partial x_N^{\alpha_N}}.$$

Set  $p = m + \beta$ . When the following condition is satisfied,  $p$  can be thought of as a lower bound to the smoothness of  $g^*$ .

CONDITION 4. For  $v \in \mathcal{N}$  and  $[\alpha] = m$ , the function  $\mu_v^*(\cdot)$  on  $C$  is  $m$ -times continuously differentiable and  $D^\alpha \mu_v^*(\cdot)$  satisfies a Hölder condition with exponent  $\beta$ .

THEOREM 3. Suppose Conditions 1–4 hold. Then

$$\|E(\hat{\mu}_v(\cdot) | \mathbf{X}_1, \dots, \mathbf{X}_n) - \mu_v^*(\cdot)\| = O_p\left[J^{-p} + \sqrt{J^d/n}\right], \quad v \in \mathcal{N},$$

so

$$\|E(\hat{\mu}(\cdot) | \mathbf{X}_1, \dots, \mathbf{X}_n) - \mu^*(\cdot)\| = O_p\left[J^{-p} + \sqrt{J^d/n}\right],$$

Theorems 2 and 3, which will be proven in Section 3, have the following consequence.

COROLLARY 1. Suppose Conditions 1–4 hold. Then

$$\|\hat{\mu}_v(\cdot) - \mu_v^*(\cdot)\| = O_p\left[J^{-p} + \sqrt{J^d/n}\right], \quad v \in \mathcal{N},$$

so

$$\|\hat{\mu}(\cdot) - \mu^*(\cdot)\| = O_p\left[J^{-p} + \sqrt{J^d/n}\right].$$

Given positive numbers  $a_n$  and  $b_n$  for  $n \geq n_0$ ,  $a_n \sim b_n$  means that  $a_n/b_n$  is bounded away from zero and infinity. Set  $\gamma = 1/(2p+d)$  and  $r = p/(2p+d)$ . Observe that if Condition 2 holds with  $J \sim n^\gamma$ , then  $p > (d_1 - d)/2$ . The next result follows from Corollary 1.

COROLLARY 2. Suppose Conditions 1–4 hold and that  $J \sim n^\gamma$ . Then

$$\|\hat{\mu}_v(\cdot) - \mu_v^*(\cdot)\| = O_p(n^{-r}), \quad v \in \mathcal{N},$$

so

$$\|\hat{\mu}(\cdot) - \mu^*(\cdot)\| = O_p(n^{-r}).$$

The  $L_2$  rate of convergence in Corollary 2 does not depend on  $N$ . It is clear from Stone (1982) with  $d = N$  that this rate is optimal. When  $d = N$ , it is possible to use the tensor-product extension of de Boor (1976) referred to in the proof of Lemma 9 below to obtain the pointwise and  $L_\infty$  rates of convergence of  $\hat{\mu}(\cdot)$  to  $\mu^*(\cdot)$  (see Koo, 1988).

Chen (1989) has obtained results along the lines of those of the present paper with penalized least-squares estimation with  $X_1, \dots, X_n$  replaced by deterministic design points  $\mathbf{x}_1, \dots, \mathbf{x}_n$ . For mathematical convenience, however, he imposes the severe restriction on the design points that they form a (suitably regular) balanced complete factorial design. Under this restriction, his results for  $d < N$  are new and imply those for  $d < N$ . He also assumes that  $\mu(\cdot) \in \mathcal{H}$  and (essentially) requires that  $p \geq dm$  for some positive integer  $m$  with  $2m > N$ , which is much more restrictive than the requirement  $p > (d_1 - d)/2$  for Corollary 2. (In a private communication, Chen stated that the condition  $2m > N$  in his paper can be replaced by the condition  $4m > d$ .)

When  $d = 1$ , the results in this section were obtained by Stone (1985). When  $d = 1$  and  $N = 1$ , similar results were obtained by Agarwal and Studden (1980) in the context of suitably regular deterministic designs. The results for additive regression ( $d = 1$ ) have been extended to robust additive regression by Mo (1990a, 1990b).

There is a growing literature on the methodological aspects of finite-dimensional additive and interactive spline modelling. See Stone and Koo (1985), Friedman and Silverman (1989), Breiman (1989) and, especially, Friedman (1991). For the corresponding methodology based on penalized least-squares estimation, see Barry (1983, 1986) and Wahba (1986). For other related additive methodologies, see Buja, Hastie and Tibshirani (1989) and Hastie and Tibshirani (1990).



3. **Proofs.** The arguments in this section were partly suggested by those in de Boor (1976) and Stone (1985).

LEMMA 1. Suppose Condition 1 holds, set  $\delta_1 = 1 - \sqrt{1 - M_1^{-2} M_2^{-2}} \in (0, 1]$ , and let  $h_v \in \mathcal{H}_v^0$  for  $v \in \mathcal{N}$ . Then

$$(3.1) \quad E \left[ \left[ \sum_v h_v(\mathbf{X}) \right]^2 \right] \geq \delta_1^{\#(\mathcal{N})-1} \sum_v E[h_v^2(\mathbf{X})].$$

PROOF. Recall that  $M_1, M_2 \geq 1$ . We will verify (3.1) by induction on  $\#(\mathcal{N})$ . Observe first that it is trivially true when  $\#(\mathcal{N}) = 1$ . Suppose  $\#(\mathcal{N}) \geq 2$  and that (3.1) holds whenever  $\mathcal{N}$  is replaced by  $\mathcal{N}'$  with  $\#(\mathcal{N}') < \#(\mathcal{N})$ . Choose a "maximal"  $\eta \in \mathcal{N}$  (that is, such that  $\eta$  is not a proper subset of any set  $v$  in  $\mathcal{N}$ ). We first verify that

$$(3.2) \quad E \left[ \left[ \sum_v h_v(\mathbf{X}) \right]^2 \right] \geq M_1^{-2} M_2^{-2} E[(h_\eta(\mathbf{X}))^2].$$

If  $\#(\eta) = N$ , then (3.2) follows immediately from the definition of  $\mathcal{H}_\eta^0$ . Suppose, instead, that  $1 \leq \#(\eta) \leq N-1$ . We can write  $\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_2)$ , where  $\mathbf{X}_1$  consists of  $X_l$ ,  $l \in \eta$ , in some order and  $\mathbf{X}_2$  consists of  $X_l$ ,  $l \notin \eta$ , in some order. Then  $\mathbf{X}_1$  is  $C_1$ -valued and  $\mathbf{X}_2$  is  $C_2$ -valued, where  $C_1 = [0, 1]^{N-\#(\eta)}$  and  $C_2 = [0, 1]^{\#(\eta)}$ . Let  $f_{\mathbf{X}_1}$  denote the density of  $\mathbf{X}_1$ ,  $f_{\mathbf{X}_2}$  the density of  $\mathbf{X}_2$  and  $f_{\mathbf{X}_2|\mathbf{X}_1}$  the conditional density of  $\mathbf{X}_2$  given  $\mathbf{X}_1$ . Then  $f_{\mathbf{X}_1}$  and  $f_{\mathbf{X}_2}$  are bounded above by  $M_2$ , so

$$f_{\mathbf{X}_2|\mathbf{X}_1}(\mathbf{x}_2|\mathbf{x}_1) \geq M_1^{-1} M_2^{-2} f_{\mathbf{X}_2}(\mathbf{x}_2), \quad \mathbf{x}_1 \in C_1 \text{ and } \mathbf{x}_2 \in C_2.$$

We can write  $h_\eta(\mathbf{x})$  as  $h_\eta(\mathbf{x}_2)$  for  $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)$ . Since  $f_{\mathbf{X}_1}$  is bounded below by  $M_1^{-1}$ , we conclude from the definition of  $\mathcal{H}_\eta^0$  that

$$\begin{aligned} E \left[ \left[ \sum_v h_v(\mathbf{X}) \right]^2 \right] &= \int_{C_1} f_{\mathbf{X}_1}(\mathbf{x}_1) d\mathbf{x}_1 \int_{C_2} \left[ h_\eta(\mathbf{x}_2) + \sum_{v \neq \eta} h_v(\mathbf{x}_1, \mathbf{x}_2) \right]^2 f_{\mathbf{X}_2|\mathbf{X}_1}(\mathbf{x}_2|\mathbf{x}_1) d\mathbf{x}_2 \\ &\geq M_1^{-2} M_2^{-2} \inf_{\mathbf{x}_1 \in C_1} \int_{C_2} \left[ h_\eta(\mathbf{x}_2) + \sum_{v \neq \eta} h_v(\mathbf{x}_1, \mathbf{x}_2) \right]^2 f_{\mathbf{X}_2}(\mathbf{x}_2) d\mathbf{x}_2 \\ &= M_1^{-2} M_2^{-2} \inf_{\mathbf{x}_1 \in C_1} E \left[ \left[ h_\eta(\mathbf{X}_2) + \sum_{v \neq \eta} h_v(\mathbf{x}_1, \mathbf{X}_2) \right]^2 \right] \\ &\geq M_1^{-2} M_2^{-2} E[h_\eta^2(\mathbf{X})] \end{aligned}$$

and hence that (3.2) again holds.

It follows from (3.2) that

$$(3.3) \quad E \left[ \left[ h_\eta(\mathbf{X}) - \beta \sum_{\mathbf{v} \neq \eta} h_{\mathbf{v}}(\mathbf{X}) \right]^2 \right] \geq M_1^{-2} M_2^{-2} E[h_\eta^2(\mathbf{X})], \quad \beta \in \mathbb{R}.$$

We conclude from (3.3) that

$$\left[ E \left[ h_\eta(\mathbf{X}) \sum_{\mathbf{v} \neq \eta} h_{\mathbf{v}}(\mathbf{X}) \right] \right]^2 \leq (1 - M_1^{-2} M_2^{-2}) E[h_\eta^2(\mathbf{X})] E \left[ \left[ \sum_{\mathbf{v} \neq \eta} h_{\mathbf{v}}(\mathbf{X}) \right]^2 \right]$$

and hence, by the induction hypothesis, that

$$\begin{aligned} E \left[ \left[ \sum_{\mathbf{v}} h_{\mathbf{v}}(\mathbf{X}) \right]^2 \right] &\geq \left[ 1 - \sqrt{1 - M_1^{-2} M_2^{-2}} \right] \left\{ E[h_\eta^2(\mathbf{X})] + E \left[ \left[ \sum_{\mathbf{v} \neq \eta} h_{\mathbf{v}}(\mathbf{X}) \right]^2 \right] \right\} \\ &\geq \left[ 1 - \sqrt{1 - M_1^{-2} M_2^{-2}} \right] \left\{ E[h_\eta^2(\mathbf{X})] + \left[ 1 - \sqrt{1 - M_1^{-2} M_2^{-2}} \right]^{\#(\mathcal{N})-2} \sum_{\mathbf{v} \neq \eta} E[h_{\mathbf{v}}^2(\mathbf{X})] \right\} \\ &\geq \left[ 1 - \sqrt{1 - M_1^{-2} M_2^{-2}} \right]^{\#(\mathcal{N})-1} \sum_{\mathbf{v}} E[h_{\mathbf{v}}^2(\mathbf{X})]. \end{aligned}$$

Therefore (3.1) holds for  $\mathcal{N}$ .  $\square$

**LEMMA 2.** Suppose  $\mathcal{G}$  is nonsingular,  $g_{\mathbf{v}} \in \mathcal{G}_{\mathbf{v}}^0$  for  $\mathbf{v} \in \mathcal{N}$  and  $\sum_{\mathbf{v}} g_{\mathbf{v}} = 0$ . Then  $g_{\mathbf{v}} = 0$  for  $\mathbf{v} \in \mathcal{N}$ .

**PROOF.** It suffices to show that if  $\mathbf{v}$  is maximal, then  $g_{\mathbf{v}} = 0$ . To this end, as an application of Lemma 1, we can write  $g_{\mathbf{v}} = \sum^{(\mathbf{v})} \tilde{g}_\eta$ , where  $\tilde{g}_\eta \in \mathcal{G}_\eta$  for  $\eta$  a proper subset of  $\mathbf{v}$  and  $\sum^{(\mathbf{v})}$  denotes summation over proper subsets of  $\mathbf{v}$ . Then

$$\|g_{\mathbf{v}}\|_n^2 = \langle g_{\mathbf{v}}, \sum^{(\mathbf{v})} \tilde{g}_\eta \rangle_n = 0$$

and hence  $g_{\mathbf{v}} = 0$ .  $\square$

Write  $\mathbf{X}_i = (X_{i1}, \dots, X_{iN})$  for  $1 \leq i \leq n$ .

**LEMMA 3.** Suppose Condition 1 holds and let  $t > 0$ . Then, except on an event having probability at most  $2(m+1)^N \exp(-2nt^2)$ , the inequalities

$$\left| \frac{1}{n} \sum_i \prod_l p_{1l}(X_{il}) p_{2l}(X_{il}) - E \left[ \prod_l p_{1l}(X_l) p_{2l}(X_l) \right] \right| \leq c_m^{2N} M_1 \sqrt{E \left[ \prod_l p_{1l}^2(X_{1l}) \right] E \left[ \prod_l p_{2l}^2(X_{2l}) \right]}$$

hold simultaneously for all polynomials  $p_{11}, \dots, p_{1N}, p_{21}, \dots, p_{2N}$  of degree  $m$ . Here

$c_m$  is a positive number that depends only on  $m$ .

PROOF. By an elementary compactness argument, there is a positive number  $c_m$  such that if  $p$  is a polynomial of degree  $m$ , then

$$(3.4) \quad \left[ \sum_0^m \frac{|p^{(k)}(0)|}{k!} \right]^2 \leq c_m^2 \int_0^1 p^2(x) dx.$$

It follows from Hoeffding's inequality (Theorem 1 of Hoeffding, 1963) that, except on an event having probability at most  $2(m+1)^{2N} \exp(-2nt^2)$ , the inequalities

$$(3.5) \quad \left| \frac{1}{n} \sum_i X_{i1}^{k_{11}} \cdots X_{iN}^{k_{1N}} X_{i1}^{k_{21}} \cdots X_{iN}^{k_{2N}} - E[X_1^{k_{11}} \cdots X_N^{k_{1N}} X_1^{k_{21}} \cdots X_N^{k_{2N}}] \right| < t$$

hold simultaneously for all choices in  $\{0, \dots, m\}$  of  $k_{11}, \dots, k_{1N}, k_{21}, \dots, k_{2N}$ . It follows from (3.4) and (3.5) that

$$\left| \frac{1}{n} \sum_i \prod_l p_{1l}(X_{il}) p_{2l}(X_{il}) - E[\prod_l p_{1l}(X_l) p_{2l}(X_l)] \right| \leq t c_m^{2N} \sqrt{\prod_l \int_0^1 p_{1l}^2(x) dx \prod_l \int_0^1 p_{2l}^2(x) dx}.$$

Since

$$E\left[\prod_l p_{1l}^2(X_l)\right] = \int_C \prod_l p_{1l}^2(x_l) f(\mathbf{x}) d\mathbf{x} \geq M_1^{-1} \int_C \prod_l p_{1l}^2(x_l) d\mathbf{x} \geq M_1^{-1} \prod_l \int_0^1 p_{1l}^2(x) dx$$

and, similarly,

$$E\left[\prod_l p_{2l}^2(X_l)\right] = M_1^{-1} \prod_l \int_0^1 p_{2l}^2(x) dx,$$

the desired result holds.  $\square$

LEMMA 4. Suppose Conditions 1 and 2 hold and let  $\varepsilon > 0$ . Then, except on an event whose probability tends to zero with  $n$ ,

$$(3.6) \quad |\langle g_1, g_2 \rangle_n - E[g_1(\mathbf{X})g_2(\mathbf{X})]| \leq \varepsilon \sqrt{E[g_1^2(\mathbf{X})]E[g_2^2(\mathbf{X})]},$$

$$g_1, g_2 \in \mathcal{G}_{\eta \cup \nu} \text{ for some } \eta, \nu \in \mathcal{N}.$$

PROOF. It suffices to verify the desired result when  $q = -1$  and  $d = N$ . Then  $d_1 = N$ ,  $\mathcal{G}$  is the span of all functions  $g$  on  $C$  of the form

$$g(\mathbf{x}) = \prod_{l=1}^N g_l(x_l), \quad \text{where } \mathbf{x} = (x_1, \dots, x_N) \text{ and } g_l \in \mathcal{O} \text{ for } 1 \leq l \leq N,$$

and (3.6) simplifies to

$$|\langle g_1, g_2 \rangle_n - E[g_1(\mathbf{X})g_2(\mathbf{X})]| \leq \varepsilon \sqrt{E[g_1^2(\mathbf{X})]E[g_2^2(\mathbf{X})]}, \quad g_1, g_2 \in \mathcal{G}.$$

Given  $k_1, \dots, k_N \in \{1, \dots, K\}$ , set  $\mathbf{k} = (k_1, \dots, k_N)$  and

$$I_{\mathbf{k}} = \{\mathbf{x} = (x_1, \dots, x_N): x_l \in I_{k_l} \text{ for } 1 \leq l \leq N\}.$$

Let  $\text{ind}(\cdot \in I_{\mathbf{k}})$  denote the indicator function of  $I_{\mathbf{k}}$ , which is defined by  $\text{ind}(\mathbf{x} \in I_{\mathbf{k}}) = 1$  for  $\mathbf{x} \in I_{\mathbf{k}}$  and  $\text{ind}(\mathbf{x} \in I_{\mathbf{k}}) = 0$  for  $\mathbf{x} \notin I_{\mathbf{k}}$ . Given polynomials  $p_{k_1}, \dots, p_{k_N}$  of degree  $m$ , define the multivariate polynomial  $p_{\mathbf{k}}$  by

$$p_{\mathbf{k}}(\mathbf{x}) = p_{k_1}(x_1) \cdots p_{k_N}(x_N), \quad \mathbf{x} = (x_1, \dots, x_N).$$

Then every function  $g$  in  $\mathcal{G}$  can be written in the form

$$g(\mathbf{x}) = \sum_{\mathbf{k}} p_{\mathbf{k}}(\mathbf{x}) \text{ind}(\mathbf{x} \in I_{\mathbf{k}}), \quad \mathbf{x} \in C.$$

Thus we can write

$$g_1(\mathbf{x}) = \sum_{\mathbf{k}} p_{1\mathbf{k}}(\mathbf{x}) \text{ind}(\mathbf{x} \in I_{\mathbf{k}}) \quad \text{and} \quad g_2(\mathbf{x}) = \sum_{\mathbf{k}} p_{2\mathbf{k}}(\mathbf{x}) \text{ind}(\mathbf{x} \in I_{\mathbf{k}}), \quad \mathbf{x} \in C.$$

Observe that

$$g_1(\mathbf{x})g_2(\mathbf{x}) = \sum_{\mathbf{k}} p_{1\mathbf{k}}(\mathbf{x})p_{2\mathbf{k}}(\mathbf{x}) \text{ind}(\mathbf{x} \in I_{\mathbf{k}}), \quad \mathbf{x} \in C.$$

Hence

$$\begin{aligned} E[g_1(\mathbf{X})g_2(\mathbf{X})] &= \sum_{\mathbf{k}} P(\mathbf{X} \in I_{\mathbf{k}}) E(p_{1\mathbf{k}}(\mathbf{X})p_{2\mathbf{k}}(\mathbf{X}) | \mathbf{X} \in I_{\mathbf{k}}), \\ E[g_1^2(\mathbf{X})] &= \sum_{\mathbf{k}} P(\mathbf{X} \in I_{\mathbf{k}}) E(p_{1\mathbf{k}}^2(\mathbf{X}) | \mathbf{X} \in I_{\mathbf{k}}), \end{aligned}$$

and

$$E[g_2^2(\mathbf{X})] = \sum_{\mathbf{k}} P(\mathbf{X} \in I_{\mathbf{k}}) E(p_{2\mathbf{k}}^2(\mathbf{X}) | \mathbf{X} \in I_{\mathbf{k}}).$$

Set  $\mathcal{J}_{\mathbf{k}} = \{i: 1 \leq i \leq n \text{ and } \mathbf{X}_i \in I_{\mathbf{k}}\}$ . Then

$$E_n[g_1(\mathbf{X})g_2(\mathbf{X})] = \sum_{\mathbf{k}} P_n(\mathbf{X} \in I_{\mathbf{k}}) E_n(p_{1\mathbf{k}}(\mathbf{X})p_{2\mathbf{k}}(\mathbf{X}) | \mathbf{X} \in I_{\mathbf{k}}),$$

where

$$\begin{aligned} E_n[g_1(\mathbf{X})g_2(\mathbf{X})] &= \langle g_1, g_2 \rangle_n = \frac{1}{n} \sum_i g_1(\mathbf{X}_i)g_2(\mathbf{X}_i), \\ P_n(\mathbf{X} \in I_{\mathbf{k}}) &= \frac{1}{n} \#(\mathcal{J}_{\mathbf{k}}), \end{aligned}$$

and

$$E_n(p_{1\mathbf{k}}(\mathbf{X})p_{2\mathbf{k}}(\mathbf{X}) | \mathbf{X} \in I_{\mathbf{k}}) = \frac{1}{\#(\mathcal{J}_{\mathbf{k}})} \sum_{i \in \mathcal{J}_{\mathbf{k}}} p_{1\mathbf{k}}(\mathbf{X}_i)p_{2\mathbf{k}}(\mathbf{X}_i).$$

Choose  $\varepsilon_1 \in (0, 1)$  such that  $\varepsilon_1^2 + 2\varepsilon_1 \leq \varepsilon$ . It follows from Conditions 1 and 2 and Bernstein's inequality (see Theorem 3 of Hoeffding, 1963) that, except on an event whose

probability tends to zero with  $n$ ,  $|P_n(X \in I_{\mathbf{k}}) - P(X \in I_{\mathbf{k}})| \leq \varepsilon_1 P(X \in I_{\mathbf{k}})$  for all  $\mathbf{k}$  and hence

$$\frac{(1-\varepsilon_1)}{M_1 K^N} \leq P_n(X \in I_{\mathbf{k}}) \leq \frac{(1+\varepsilon_1)M_2}{K^N} \quad \text{for all } \mathbf{k}.$$

By Condition 2,  $K^N = o(n^{1-\delta})$  for some  $\delta > 0$ . Thus there are positive numbers  $M_3$  and  $\delta$  such that, except on an event whose probability tends to zero with  $n$ ,  $\#(\mathcal{J}_{\mathbf{k}}) \geq M_3^{-1} n^\delta$  for all  $\mathbf{k}$ . We conclude from Lemma 3 that, except on an event whose probability tends to zero with  $n$ ,

$$\begin{aligned} |E_n(p_{1\mathbf{k}}(\mathbf{X})p_{2\mathbf{k}}(\mathbf{X}) | \mathbf{X} \in I_{\mathbf{k}}) - E(p_{1\mathbf{k}}(\mathbf{X})p_{2\mathbf{k}}(\mathbf{X}) | \mathbf{X} \in I_{\mathbf{k}})| \\ \leq \varepsilon_1 \sqrt{E(p_{1\mathbf{k}}^2(\mathbf{X}) | \mathbf{X} \in I_{\mathbf{k}})E(p_{2\mathbf{k}}^2(\mathbf{X}) | \mathbf{X} \in I_{\mathbf{k}})} \end{aligned}$$

for all  $\mathbf{k}$  and all choices of  $p_{1\mathbf{k}}$  and  $p_{2\mathbf{k}}$ . Consequently, except on an event whose probability tends to zero with  $n$ ,

$$\begin{aligned} |\langle g_1, g_2 \rangle_n - E[g_1(\mathbf{X})g_2(\mathbf{X})]| &\leq \varepsilon_1 \sum_{\mathbf{k}} P(\mathbf{X} \in I_{\mathbf{k}}) |E(g_1(\mathbf{X})g_2(\mathbf{X}) | \mathbf{X} \in I_{\mathbf{k}})| \\ &\quad + \varepsilon_1 (1+\varepsilon_1) \sum_{\mathbf{k}} P(\mathbf{X} \in I_{\mathbf{k}}) \sqrt{E(g_1^2(\mathbf{X}) | \mathbf{X} \in I_{\mathbf{k}})E(g_2^2(\mathbf{X}) | \mathbf{X} \in I_{\mathbf{k}})} \\ &\leq (\varepsilon_1^2 + 2\varepsilon_1) \sum_{\mathbf{k}} P(\mathbf{X} \in I_{\mathbf{k}}) \sqrt{E(g_1^2(\mathbf{X}) | \mathbf{X} \in I_{\mathbf{k}})E(g_2^2(\mathbf{X}) | \mathbf{X} \in I_{\mathbf{k}})} \\ &= (\varepsilon_1^2 + 2\varepsilon_1) \sum_{\mathbf{k}} \sqrt{E[g_1^2(\mathbf{X}) \text{ind}(\mathbf{X} \in I_{\mathbf{k}})]E[g_2^2(\mathbf{X}) \text{ind}(\mathbf{X} \in I_{\mathbf{k}})]} \\ &\leq \varepsilon \sqrt{E[g_1^2(\mathbf{X})]E[g_2^2(\mathbf{X})]}, \quad g_1, g_2 \in \mathcal{G}. \quad \square \end{aligned}$$

LEMMA 5. Suppose Conditions 1 and 2 hold and let  $0 < \delta_2 < \delta_1$ . Then, except on an event whose probability tends to zero with  $n$ ,

$$\|\sum_{\mathbf{v}} g_{\mathbf{v}}\|_n^2 \geq \delta_2^{\#(\mathcal{N})-1} \sum_{\mathbf{v}} \|g_{\mathbf{v}}\|_n^2, \quad g_{\mathbf{v}} \in \mathcal{G}_{\mathbf{v}}^0 \text{ for } \mathbf{v} \in \mathcal{N}.$$

PROOF. Applying Lemma 4 with  $0 < \#(\mathcal{N})\varepsilon < \delta_1^{\#(\mathcal{N})-1}$ , we see that, except on an event whose probability tends to zero with  $n$ , (3.6) holds and hence

$$\|g_{\mathbf{v}}\|_n^2 \leq (1+\varepsilon)E[g_{\mathbf{v}}^2(\mathbf{X})], \quad \mathbf{v} \in \mathcal{N} \text{ and } g_{\mathbf{v}} \in \mathcal{G}_{\mathbf{v}}^0$$

so

$$\sum_{\mathbf{v}} \|g_{\mathbf{v}}\|_n^2 \leq (1+\varepsilon) \sum_{\mathbf{v}} E[g_{\mathbf{v}}^2(\mathbf{X})].$$

Therefore, by Lemmas 1 and 4, except on an event whose probability tends to zero with  $n$ ,

$$\begin{aligned} \|\sum_{\mathbf{v}} g_{\mathbf{v}}\|_n^2 &\geq E[(\sum_{\mathbf{v}} g_{\mathbf{v}}(\mathbf{X}))^2] - \#(\mathcal{N})\varepsilon \sum_{\mathbf{v}} E[g_{\mathbf{v}}^2(\mathbf{X})] \\ &\geq (\delta_1^{\#(\mathcal{N})-1} - \#(\mathcal{N})\varepsilon) \sum_{\mathbf{v}} E[g_{\mathbf{v}}^2(\mathbf{X})] \\ &\geq \frac{\delta_1^{\#(\mathcal{N})-1} - \#(\mathcal{N})\varepsilon}{1+\varepsilon} - \sum_{\mathbf{v}} \|g_{\mathbf{v}}\|_n^2. \end{aligned}$$

Since  $\varepsilon$  can be made arbitrarily small, the desired result holds.  $\square$

Set  $\mathcal{J}_{\emptyset} = \{0\}$  and  $B_{\emptyset 0} = 1$ . For  $\mathbf{v} \in \mathcal{N}$  with  $\mathbf{v} \neq \emptyset$ , let  $\mathcal{J}_{\mathbf{v}}$  denote the collection of ordered  $\#(\mathbf{v})$ -tuples  $j_l$ ,  $l \in \mathbf{v}$ , with  $j_l \in \{1, \dots, J\}$  for  $l \in \mathbf{v}$ . Then  $\#(\mathcal{J}_{\mathbf{v}}) = J^{\#(\mathbf{v})}$ . For  $\mathbf{j} \in \mathcal{J}_{\mathbf{v}}$ , let  $B_{\mathbf{v}\mathbf{j}}$  denote the function on  $C$  given by

$$B_{\mathbf{v}\mathbf{j}}(\mathbf{x}) = \prod_{l \in \mathbf{v}} B_{j_l}(x_l), \quad \mathbf{x} = (x_1, \dots, x_N).$$

Then, for  $\mathbf{v} \in \mathcal{N}$ , the functions  $B_{\mathbf{v}\mathbf{j}}$ ,  $\mathbf{j} \in \mathcal{J}_{\mathbf{v}}$ , which are nonnegative and have sum one, form a basis of  $\mathcal{G}_{\mathbf{v}}$ .

Suppose  $\mathcal{G}$  is nonsingular and let  $g \in \mathcal{G}$ . Then  $g = \sum_{\mathbf{v}} g_{\mathbf{v}}$ , where  $g_{\mathbf{v}} \in \mathcal{G}_{\mathbf{v}}^0$ ,  $\mathbf{v} \in \mathcal{N}$ , are uniquely determined. Moreover,  $g_{\mathbf{v}} = \sum_{\mathbf{j}} b_{\mathbf{v}\mathbf{j}} B_{\mathbf{v}\mathbf{j}}$  for  $\mathbf{v} \in \mathcal{N}$ , where the  $b_{\mathbf{v}\mathbf{j}}$ 's are uniquely determined. Let  $\mathbf{v}$  and  $\mathbf{j}$  be fixed. Let  $G_{\mathbf{v}\mathbf{j}} \in \mathcal{G}$  denote the representer of the linear functional  $g \mapsto b_{\mathbf{v}\mathbf{j}}$  on  $\mathcal{G}$  relative to the inner product  $\langle \cdot, \cdot \rangle_n$ , so that  $b_{\mathbf{v}\mathbf{j}} = \langle G_{\mathbf{v}\mathbf{j}}, g \rangle_n$ . Now  $G_{\mathbf{v}\mathbf{j}} = \sum_{\mathbf{v}'} G_{\mathbf{v}\mathbf{j}\mathbf{v}'}$ , where  $G_{\mathbf{v}\mathbf{j}\mathbf{v}'} \in \mathcal{G}_{\mathbf{v}'}^0$  for  $\mathbf{v}' \in \mathcal{N}$ . Thus  $G_{\mathbf{v}\mathbf{j}\mathbf{v}'} = \sum_{\mathbf{j}' \in \mathcal{J}_{\mathbf{v}'}} \gamma_{\mathbf{v}\mathbf{j}\mathbf{v}'\mathbf{j}'} B_{\mathbf{v}'\mathbf{j}'}$  for  $\mathbf{v}' \in \mathcal{N}$ , where the  $\gamma_{\mathbf{v}\mathbf{j}\mathbf{v}'\mathbf{j}'}$ 's are uniquely determined. Observe that

$$(3.7) \quad \langle G_{\mathbf{v}\mathbf{j}}, G_{\mathbf{v}'\mathbf{j}'} \rangle_n = \gamma_{\mathbf{v}\mathbf{j}\mathbf{v}'\mathbf{j}'}, \quad \mathbf{v}, \mathbf{v}' \in \mathcal{N}, \mathbf{j} \in \mathcal{J}_{\mathbf{v}} \text{ and } \mathbf{j}' \in \mathcal{J}_{\mathbf{v}'}.$$

Observe also that, for  $\mathbf{v} \in \mathcal{N}$  and  $\mathbf{j} \in \mathcal{J}_{\mathbf{v}}$ ,  $\langle G_{\mathbf{v}\mathbf{j}}, B_{\mathbf{v}\mathbf{j}} \rangle_n = 1$  and hence  $0 < \|G_{\mathbf{v}\mathbf{j}}\|_n^2 = \gamma_{\mathbf{v}\mathbf{j}\mathbf{v}\mathbf{j}}$ .

LEMMA 6. Suppose Conditions 1 and 2 hold. Then there is a positive number  $M_3$ , which does not depend on  $J$ , such that, except on an event whose probability tends to zero with  $n$ ,

$$(3.8) \quad \|\sum_{\mathbf{v}} \sum_{\mathbf{j}} b_{\mathbf{v}\mathbf{j}} B_{\mathbf{v}\mathbf{j}}\|_n^2 \geq M_3^{-1} J^{-d} \sum_{\mathbf{v}} \sum_{\mathbf{j}} b_{\mathbf{v}\mathbf{j}}^2 \quad \text{if } \sum_{\mathbf{j}} b_{\mathbf{v}\mathbf{j}} B_{\mathbf{v}\mathbf{j}} \in \mathcal{G}_{\mathbf{v}}^0 \text{ for } \mathbf{v} \in \mathcal{N}.$$

PROOF. It follows from the basic properties of  $B$ -splines and repeated use of (viii) on page 155 of de Boor (1978) that, for some positive number  $M_4$ ,

$$\int_C [\sum_{\mathbf{j}} b_{\mathbf{v}\mathbf{j}} B_{\mathbf{v}\mathbf{j}}(\mathbf{x})]^2 d\mathbf{x} \geq 2M_4^{-1} J^{-\#(\mathbf{v})} \sum_{\mathbf{j}} b_{\mathbf{v}\mathbf{j}}^2$$

for all choices of  $v \in \mathcal{N}$  and  $b_{vj} \in \mathbb{R}$  and  $j \in \mathcal{J}_v$ . Thus, by Condition 1 and Lemma 4, except on an event whose probability tends to zero with  $n$ ,

$$\|\sum_j b_{vj} B_{vj}\|_n^2 \geq M_4^{-1} J^{\#(v)} \sum_j b_{vj}^2$$

for all such choices. The desired result now follows from Lemma 5.  $\square$

LEMMA 7. *Suppose Conditions 1 and 2 hold. Then, except on an event whose probability tends to zero with  $n$ ,*

$$(3.9) \quad \sum_{v'} \sum_j \gamma_{vjv'j'}^2 \leq M_3^2 J^{2d}, \quad v \in \mathcal{N} \text{ and } j \in \mathcal{J}_v$$

PROOF. Suppose that (3.8) holds and let  $v \in \mathcal{N}$  and  $j \in \mathcal{J}_v$ . Then

$$M_3^{-1} J^{-d} \gamma_{vjvj}^2 \leq M_3^{-1} J^{-d} \sum_{v'} \sum_j \gamma_{vjv'j'}^2 \leq \|G_{vj}\|_n^2 = \gamma_{vjvj},$$

so  $\gamma_{vjvj} \leq M_3 J^d$  and therefore (3.9) is valid. We now obtain the desired result from Lemma 6.  $\square$

LEMMA 8. *Suppose Conditions 1–3 hold. Then, except on an event whose probability tends to zero with  $n$ ,*

$$\max_v \max_{j \in \mathcal{J}_v} \text{var}(\hat{\beta}_{vj} | X_1, \dots, X_n) = O_p(J^d/n).$$

PROOF. Set  $\sigma^2(\mathbf{x}) = \text{var}(Y | X = \mathbf{x})$ ,  $\mathbf{x} \in C$ . It follows from Condition 3 that  $\sigma^2(\cdot)$  has a finite upper bound  $M_4$  on  $C$ .

Suppose that  $\mathcal{G}$  is nonsingular. Let  $Q$  denote orthogonal projection onto  $\mathcal{G}$  relative to the inner product  $\langle \cdot, \cdot \rangle_n$ . Then  $\langle g, Qh \rangle_n = \langle g, h \rangle_n$  for all real-valued functions  $h$  whose domain includes  $C'$  and all  $g \in \mathcal{G}$ . Given such a function  $h$ , write  $Qh$  in the form

$$Qh = \sum_v \sum_j b_{vj} B_{vj}, \quad \text{where } \sum_j b_{vj} B_{vj} \in \mathcal{G}_v^0 \text{ for } v \in \mathcal{N}.$$

Then  $b_{vj} = \langle G_{vj}, Qh \rangle_n = \langle G_{vj}, h \rangle_n$  and hence

$$b_{vj} = \sum_{v'} \sum_j \gamma_{vjv'j'} \langle B_{vj}, h \rangle_n, \quad v \in \mathcal{N} \text{ and } j \in \mathcal{J}_v$$

The least-squares estimate  $\hat{\mu}(\cdot)$  can be written as

$$\hat{\mu}(\cdot) = QY(\cdot) = \sum_v \sum_j \hat{\beta}_{vj} B_{vj}, \quad \text{where } \sum_j \hat{\beta}_{vj} B_{vj} \in \mathcal{G}_v^0 \text{ for } v \in \mathcal{N}.$$

Thus

$$\hat{\beta}_{vj} = \langle G_{vj}, Y(\cdot) \rangle_n = \sum_{v'} \sum_j \gamma_{vjv'j'} \langle B_{vj}, Y(\cdot) \rangle_n, \quad v \in \mathcal{N} \text{ and } j \in \mathcal{J}_v$$

Consequently,

$$\begin{aligned}\text{var}(\hat{\beta}_{vj} | \mathbf{X}_1, \dots, \mathbf{X}_n) &= \frac{1}{n^2} \sum_i \sigma^2(\mathbf{X}_i) \left[ \sum_{v'} \sum_j \gamma_{vjv'j'} B_{v'j'}(\mathbf{X}_i) \right]^2 \\ &\leq M_4 n^{-1} \|G_{vj}\|_n^2 \\ &= M_4 n^{-1} \gamma_{vj} \gamma_{vj}.\end{aligned}$$

The desired result now follows from Theorem 1 and Lemma 7.  $\square$

Theorem 2 follows from Lemma 8.

LEMMA 9. *Suppose Conditions 1–3 hold and that  $\mu^*(\cdot) = 0$ . Then*

$$\|E(\hat{\mu}_v(\cdot) | \mathbf{X}_1, \dots, \mathbf{X}_n)\|_n = O_P\left[\sqrt{J^d/n}\right], \quad v \in \mathcal{N}.$$

PROOF. Choose  $v \in \mathcal{N}$  and recall that  $B_{vj}, j \in \mathcal{J}_v$ , form a basis of  $\mathcal{G}_v$ . Let  $g \in \mathcal{G}_v$ . Then  $g = \sum_j b_j^{(v)} B_{vj}$ , where the  $b_j^{(v)}$ 's are uniquely determined. Suppose  $\mathcal{G}$  is nonsingular. Let  $G_j^{(v)}$  denote the representer of the linear functional  $g \mapsto b_j^{(v)}$  on  $\mathcal{G}_v$  relative to the inner product  $\langle \cdot, \cdot \rangle_n$ , so that  $b_j^{(v)} = \langle G_j^{(v)}, g \rangle$ . Then  $G_j^{(v)} = \sum_{j'} \gamma_{jj'}^{(v)} B_{vj'}$ , where the  $\gamma_{jj'}^{(v)}$ 's are uniquely determined. (Alternatively,  $(\gamma_{jj'}^{(v)})$  is the inverse of the Gram matrix  $(\langle B_{vj}, B_{vj'} \rangle)$ .) Let  $\mu^{(v)}(\cdot)$  denote the orthogonal projection of  $\mu(\cdot)$  onto  $\mathcal{G}_v$  (relative to the inner product  $\langle \cdot, \cdot \rangle_n$ ). Then  $\mu^{(v)}(\cdot) = \sum_j \beta_j^{(v)} B_{vj}$ , where

$$\beta_j^{(v)} = \sum_{j'} \gamma_{jj'}^{(v)} \langle B_{vj'}, \mu(\cdot) \rangle_n, \quad j \in \mathcal{J}_v.$$

Now

$$\|\mu^{(v)}(\cdot)\|_n^2 = \left\| \sum_j \beta_j^{(v)} B_{vj} \right\|_n^2 = \sum_j \sum_{j'} \beta_j^{(v)} \beta_{j'}^{(v)} \langle B_{vj}, B_{vj'} \rangle_n.$$

By Conditions 1 and 2, Bernstein's inequality applied to the binomial distribution, and the basic properties of  $B$ -splines,

$$\|\mu^{(v)}(\cdot)\|_n^2 = O_P(J^{\#(v)} \sum_j (\beta_j^{(v)})^2).$$

It follows from Conditions 1 and 2 by an extension of arguments in de Boor (1976) and Stone (1989) that there are numbers  $M_4 \in (0, \infty)$  and  $c \in (0, 1)$  (both independent of  $J$ ) such that, except on an event whose probability tends to zero with  $n$ ,

$$|\gamma_{jj'}^{(v)}| \leq M_4 J^{\#(v)} c^{|j'-j|}, \quad j, j' \in \mathcal{J}_v$$

Consequently,

$$\sum_j (\beta_j^{(v)})^2 = O_P\left[J^{2\#(v)} \sum_j \left[ \sum_{j'} c^{|j'-j|} |\langle B_{vj'}, \mu(\cdot) \rangle_n| \right]^2\right] = O_P(J^{2\#(v)} \sum_j (\langle B_{vj}, \mu(\cdot) \rangle_n)^2).$$



Since  $\mu^*(\cdot) = 0$ , we see that  $E(\langle B_{vj}, \mu(\cdot) \rangle_n) = E(B_{vj}(\mathbf{X})\mu(\mathbf{X})) = 0$  for  $j \in \mathcal{J}_v$ . Moreover,

$$\max_j \text{var}(\langle B_{vj}, \mu(\cdot) \rangle_n) = n^{-1} \max_j \text{var}(B_{vj}(\mathbf{X})\mu(\mathbf{X})) = O(n^{-1}J^{-\#(v)}).$$

Thus  $E[\sum_j (\langle B_{vj}, \mu(\cdot) \rangle_n)^2] = O(1/n)$  and hence  $\sum_j (\langle B_{vj}, \mu(\cdot) \rangle_n)^2 = O_P(1/n)$ . Consequently,  $\sum_j (\beta^{(v)})^2 = O_P(J^{2\#(v)}/n)$  and therefore

$$\|\mu^{(v)}(\cdot)\|_n^2 = O_P(J^{\#(v)}/n) = O_P(J^d/n), \quad v \in \mathcal{N}.$$

Let  $\mu_v^0(\cdot)$  denote the orthogonal projection of  $\mu(\cdot)$  onto  $\mathcal{G}_v^0$ , which equals the orthogonal projection of  $\mu^{(v)}(\cdot)$  onto  $\mathcal{G}_v^0$ . Then  $\|\mu_v^0(\cdot)\|_n^2 \leq \|\mu^{(v)}(\cdot)\|_n^2$  and hence

$$\|\mu_v^0(\cdot)\|_n^2 = O_P(J^d/n), \quad v \in \mathcal{N}.$$

Observe that  $E(\hat{\mu}(\cdot) | \mathbf{X}_1, \dots, \mathbf{X}_n)$  is the orthogonal projection of  $\mu(\cdot)$  onto  $\mathcal{G}$ . We can write this orthogonal projection as  $\sum_v \mu_v(\cdot)$ , where  $\mu_v(\cdot) \in \mathcal{G}_v^0$  for  $v \in \mathcal{N}$ . Now  $\mu_v^0(\cdot)$  is the orthogonal projection of  $\sum_v \mu_v(\cdot)$  onto  $\mathcal{G}_v^0$  for  $v \in \mathcal{N}$ , so

$$\begin{aligned} \|\sum_v \mu_v(\cdot)\|_n^2 &= \sum_v \langle \mu_v(\cdot), \sum_v \mu_v(\cdot) \rangle_n \\ &= \sum_v \langle \mu_v(\cdot), \mu_v^0(\cdot) \rangle_n \\ &\leq \sum_v \|\mu_v(\cdot)\|_n \|\mu_v^0(\cdot)\|_n \\ &\leq (\max_v \|\mu_v(\cdot)\|_n) \sum_v \|\mu_v^0(\cdot)\|_n. \end{aligned}$$

We conclude from Lemma 5 that

$$\|E(\hat{\mu}(\cdot) | \mathbf{X}_1, \dots, \mathbf{X}_n)\|_n^2 = \|\sum_v \mu_v(\cdot)\|_n^2 = O_P(\sum_v \|\mu_v^0(\cdot)\|_n^2) = O_P(J^d/n).$$

The desired result now follows by another application of Lemma 5.  $\square$

**LEMMA 10.** *Suppose Conditions 1–4 hold and that  $\mu^*(\cdot) = \mu(\cdot)$ . Then*

$$\|E(\hat{\mu}_v(\cdot) | \mathbf{X}_1, \dots, \mathbf{X}_n) - \mu_v^*(\cdot)\|_n^2 = O_P(J^{-2p} + J^{d-1}/n), \quad v \in \mathcal{N}.$$

**PROOF.** It follows from Condition 4 (see Theorem 12.8 of Schumaker, 1981) that there is a positive number  $M_4$  not depending on  $n$  or  $J$  such that, for  $v \in \mathcal{N}$ , there is a function  $g_v \in \mathcal{G}_v$  with  $\|g_v - \mu_v^*(\cdot)\|_\infty \leq M_4 J^{-p}$ ; here  $\|h\|_\infty = \sup_{\mathbf{x} \in C} |h(\mathbf{x})|$  is the  $L_\infty$  norm of a function  $h$  on  $C$ . Choose  $v \in \mathcal{N}$  and let  $\eta$  be a proper subset of  $v$ . Then  $E[B_{\eta j}(\mathbf{X})\mu_v^*(\mathbf{X})] = 0$  for  $j \in \mathcal{J}_\eta$  and hence

$$\max_j |E[B_{\eta j}(\mathbf{X})g_v(\mathbf{X})]| = O(J^{-\#(\eta)-p}).$$

Moreover,

$$\max_j \text{var}(B_{\eta j}(\mathbf{X})g_{\mathbf{v}}(\mathbf{X})) = O(J^{-\#(\eta)}).$$

Suppose  $\mathcal{G}$  is nonsingular. Let  $g_{\mathbf{v}\eta}^0$  denote the orthogonal projection of  $g_{\mathbf{v}}$  onto  $\mathcal{G}_{\eta}^0$  (relative to the inner product  $\langle \cdot, \cdot \rangle_n$ ). Arguing as in the proof of Lemma 9, we get that

$$\|g_{\mathbf{v}\eta}^0\|_n^2 = O_P(J^{-2p} + J^{\#(\eta)}/n) = O_P(J^{-2p} + J^{d-1}/n).$$

Write  $g_{\mathbf{v}} = \sum_{\eta \subset \mathbf{v}} g_{\mathbf{v}\eta}$ , where  $g_{\mathbf{v}\eta} \in \mathcal{G}_{\eta}^0$  for  $\eta \subset \mathbf{v}$ . Let  $\sum^{(\mathbf{v})}$  denote summation over all proper subsets of  $\mathbf{v}$ . Then  $g_{\mathbf{v}\eta}^0$  is the orthogonal projection of  $\sum^{(\mathbf{v})} g_{\mathbf{v}\eta}$  onto  $\mathcal{G}_{\eta}^0$ . We conclude, again by arguing as in the proof of Lemma 9, that

$$\|g_{\mathbf{v}} - g_{\mathbf{v}\mathbf{v}}\|_n^2 = \|\sum^{(\mathbf{v})} g_{\mathbf{v}\eta}\|_n^2 = O_P(J^{-p} + J^{d-1}/n).$$

Replacing  $g_{\mathbf{v}}$  by  $g_{\mathbf{v}\mathbf{v}}$  if necessary, we see that, for  $\mathbf{v} \in \mathcal{N}$ , there is a function  $g_{\mathbf{v}} \in \mathcal{G}_{\mathbf{v}}^0$  such that  $\|g_{\mathbf{v}} - \mu_{\mathbf{v}}^*(\cdot)\|_n^2 = O_P(J^{-2p} + J^{d-1}/n)$  and hence  $\|\sum_{\mathbf{v}} g_{\mathbf{v}} - \mu^*(\cdot)\|_n^2 = O_P(J^{-2p} + J^{d-1}/n)$ .

Write the orthogonal projection  $E(\hat{\mu}(\cdot) | \mathbf{X}_1, \dots, \mathbf{X}_n)$  of  $\mu(\cdot) = \mu^*(\cdot)$  onto  $\mathcal{G}$  as  $\sum_{\mathbf{v}} \mu_{\mathbf{v}}(\cdot)$ , where  $\mu_{\mathbf{v}}(\cdot) = E(\hat{\mu}_{\mathbf{v}}(\cdot) | \mathbf{X}_1, \dots, \mathbf{X}_n) \in \mathcal{G}_{\mathbf{v}}^0$  for  $\mathbf{v} \in \mathcal{N}$ . Observe that

$$\|\sum_{\mathbf{v}} \mu_{\mathbf{v}}(\cdot) - \mu^*(\cdot)\|_n^2 \leq \|\sum_{\mathbf{v}} g_{\mathbf{v}} - \mu^*(\cdot)\|_n^2.$$

Thus

$$\|\sum_{\mathbf{v}} \mu_{\mathbf{v}}(\cdot) - \mu^*(\cdot)\|_n^2 = O_P(J^{-2p} + J^{d-1}/n)$$

and hence

$$\|\sum_{\mathbf{v}} \mu_{\mathbf{v}}(\cdot) - \sum_{\mathbf{v}} g_{\mathbf{v}}\|_n^2 = O_P(J^{-2p} + J^{d-1}/n).$$

We conclude from Lemma 5 that

$$\|\mu_{\mathbf{v}}(\cdot) - g_{\mathbf{v}}\|_n^2 = O_P(J^{-2p} + J^{d-1}/n) \quad \mathbf{v} \in \mathcal{N},$$

and therefore that

$$\|\mu_{\mathbf{v}}(\cdot) - \mu_{\mathbf{v}}^*(\cdot)\|_n^2 = O_P(J^{-2p} + J^{d-1}/n) \quad \mathbf{v} \in \mathcal{N}. \quad \square$$

LEMMA 11. Suppose Conditions 1–4 hold. There is a positive number  $M_4$  not depending on  $n$  or  $J$  such that, except on an event whose probability tends to zero with  $n$ ,

$$\|g - \mu_{\mathbf{v}}^*(\cdot)\|^2 \leq M_4(\|g - \mu_{\mathbf{v}}^*(\cdot)\|_n^2 + J^{-2p}), \quad \mathbf{v} \in \mathcal{N} \text{ and } g \in \mathcal{G}_{\mathbf{v}}.$$

PROOF. Given  $\mathbf{v} \in \mathcal{N}$ , set  $h = \mu_{\mathbf{v}}^*(\cdot)$  and let  $g \in \mathcal{G}_{\mathbf{v}}$ . Then (see the proof of Lemma 4)  $g$  can be written in the form

$$g(\mathbf{x}) = \sum_{\mathbf{k}} p_{\mathbf{k}}(\mathbf{x}) \text{ind}(\mathbf{x} \in I_{\mathbf{k}}), \quad \mathbf{x} \in C.$$

By Condition 4, there is a function  $g_1$  of the same form such that  $\|g_1 - h\|_\infty \leq M_5 J^{-p}$ ,  $M_5$  being a positive number that does not depend on  $n$  of  $J$ . Then  $\|g_1 - h\| \leq M_5 J^{-p}$  and  $\|g_1 - h\|_n \leq M_5 J^{-p}$ , so

$$\|g - h\|^2 \leq 2\|g - g_1\|^2 + 2M_5^2 J^{-2p} \quad \text{and} \quad \|g - g_1\|_n^2 \leq 2\|g - h\|_n^2 + 2M_5^2 J^{-2p}.$$

It follows from Lemma 4 that, except on an event whose probability tends to zero with  $n$ ,  $\|g - g_1\|^2 \leq 2\|g - g_1\|_n^2$  and hence

$$\|g - h\|^2 \leq 4\|g - g_1\|_n^2 + 2M_5^2 J^{-2p} \leq 8\|g - h\|_n^2 + 10M_5^2 J^{-2p}. \quad \square$$

**PROOF OF THEOREM 3.** It follows from Lemma 9 applied to the regression function  $\mu(\cdot) - \mu^*(\cdot)$  and Lemma 10 applied to the regression function  $\mu^*(\cdot)$  that

$$\|E(\hat{\mu}_v(\cdot) | \mathbf{X}_1, \dots, \mathbf{X}_n) - \mu_v^*(\cdot)\|_n^2 = O_p(J^{-2p} + J^d/n), \quad v \in \mathcal{N}.$$

We conclude from Lemma 11 that

$$\|E(\hat{\mu}_v(\cdot) | \mathbf{X}_1, \dots, \mathbf{X}_n) - \mu_v^*(\cdot)\|^2 = O_p(J^{-2p} + J^d/n), \quad v \in \mathcal{N}. \quad \square$$

## REFERENCES

- AGARWAL, G. G. and STUDDEN, W. J. (1980). Asymptotic integrated mean square error using least squares and bias minimizing spline. *Ann. Statist.* 8 1307–1325.
- BARRY, D. (1983). *Nonparametric Bayesian Regression*. Ph. D. Dissertation, Department of Statistics, Yale University.
- BARRY, D. (1986). Nonparametric Bayesian regression. *Ann. Statist.* 14 934–953.
- BREIMAN, L. (1989). Fitting additive models to regression data. Technical Report No. 209, Department of Statistics, University of California at Berkeley.
- BUJA, A., HASTIE, T. and TIBSHIRANI, R. (1989). Linear smoothers and additive models (with discussion). *Ann. Statist.* 17 453–555.
- DE BOOR, C. (1976). A bound on the  $L_\infty$ -norm of  $L_2$ -approximation by splines in terms of a global mesh ratio. *Math. Comp.* 30 765–771.
- DE BOOR, C. (1978). *A Practical Guide to Splines*. Springer–Verlag, New York.
- CHEN, Z. (1989). Interaction spline models and their convergence rates. Manuscript.
- FRIEDMAN, J. H. (1991). Multivariate Adaptive Regression Splines. *Ann. Statist.*, to appear.
- FRIEDMAN, J. H. and SILVERMAN, B. W. (1989). Flexible parsimonious smoothing and additive modeling (with discussion). *Technometrics* 31 3–39.
- HASTIE, T. J. and TIBSHIRANI, R. J. (1990). *Generalized Additive Models*. Chapman and Hall, New York.
- HOEFFDING, W. (1963). Probability inequalities for sums of bounded random variables. *J. Amer. Statist. Assoc.* 58 13–30.
- KOO, J.-Y. (1988) *Tensor product splines in the estimation of regression, exponential response functions and multivariate densities*. Ph. D. dissertation, Department of Statistics, University of California at Berkeley.
- SCHUMAKER, L. L. (1981). *Spline Functions: Basic Theory*. Wiley, New York.
- STONE, C. J. (1980). Optimal rates of convergence for nonparametric estimators. *Ann. Statist.* 8 1348–1360.

- STONE, C. J. (1982). Optimal global rates of convergence for nonparametric regression. *Ann. Statist.* 10 1040–1053.
- STONE, C. J. (1985). Additive regression and other nonparametric models. *Ann. Statist.* 13 689–705.
- STONE, C. J. (1986). The dimensionality reduction principle for generalized additive models. *Ann. Statist.* 14 590–606.
- STONE, C. J. and KOO, C.-Y. (1986) Additive splines in statistics. In *1985 Statistical Computing Section Proceedings of the American Statistical Association* 45–48. American Statistical Association, Washington.
- STONE, C. J. (1989). Uniform error bounds involving logspline models. In *Probability, Statistics and Mathematics: Papers in Honor of Samuel Karlin* (T. W. Anderson, K. B. Athreya, and D. L. Iglehart, eds.) 335–355. Academic Press, Boston.
- 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*. (Boardman, T. J., ed.) 75–80. Amer. Statist. Assoc, Washington, D. C.

## REFERENCES

- AGARWAL, G. G. and STUDDEN, W. J. (1980). Asymptotic integrated mean square error using least squares and bias minimizing spline. *Ann. Statist.* **8** 1307–1325.
- BARRY, D. (1983). *Nonparametric Bayesian Regression*. Ph. D. Dissertation, Department of Statistics, Yale University.
- BARRY, D. (1986). Nonparametric Bayesian regression. *Ann. Statist.* **14** 934–953.
- BREIMAN, L. (1989). Fitting additive models to regression data. Technical Report No. 209, Department of Statistics, University of California at Berkeley.
- BUJA, A., HASTIE, T. and TIBSHIRANI, R. (1989). Linear smoothers and additive models (with discussion). *Ann. Statist.* **17** 453–555.
- DE BOOR, C. (1976). A bound on the  $L_\infty$ -norm of  $L_2$ -approximation by splines in terms of a global mesh ratio. *Math. Comp.* **30** 765–771.
- DE BOOR, C. (1978). *A Practical Guide to Splines*. Springer–Verlag, New York.
- CHEN, Z. (1989). Interaction spline models and their convergence rates. Manuscript.
- FRIEDMAN, J. H. (1991). Multivariate Adaptive Regression Splines. *Ann. Statist.*, to appear.
- FRIEDMAN, J. H. and SILVERMAN, B. W. (1989). Flexible parsimonious smoothing and additive modeling (with discussion). *Technometrics* **31** 3–39.
- HASTIE, T. J. and TIBSHIRANI, R. J. (1990). *Generalized Additive Models*. Chapman and Hall, New York.
- HOEFFDING, W. (1963). Probability inequalities for sums of bounded random variables. *J. Amer. Statist. Assoc.* **58** 13–30.
- KOO, J.-Y. (1988). *Tensor product splines in the estimation of regression, esponential response functions and multivariate densities*. Ph. D. dissertation, Department of Statistics, University of California at Berkeley.
- MO, M. (1990a). Robust additive regression I: population aspect. Manuscript.
- MO, M. (1990b). Robust additive regression II: finite sample approximations. Manuscript.
- SCHUMAKER, L. L. (1981). *Spline Functions: Basic Theory*. Wiley, New York.

- STONE, C. J. (1980). Optimal rates of convergence for nonparametric estimators. *Ann. Statist.* 8 1348–1360.
- STONE, C. J. (1982). Optimal global rates of convergence for nonparametric regression. *Ann. Statist.* 10 1040–1053.
- STONE, C. J. (1985). Additive regression and other nonparametric models. *Ann. Statist.* 13 689–705.
- STONE, C. J. (1986). The dimensionality reduction principle for generalized additive models. *Ann. Statist.* 14 590–606.
- STONE, C. J. and KOO, C.-Y. (1986) Additive splines in statistics. In *1985 Statistical Computing Section Proceedings of the American Statistical Association* 45–48. American Statistical Association, Washington.
- STONE, C. J. (1989). Uniform error bounds involving logspline models. In *Probability, Statistics and Mathematics: Papers in Honor of Samuel Karlin* (T. W. Anderson, K. B. Athreya, and D. L. Iglehart, eds.) 335–355. Academic Press, Boston.
- 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*. (Boardman, T. J., ed.) 75–80. Amer. Statist. Assoc, Washington, D. C.