

## ASYMPTOTICS FOR DOUBLY FLEXIBLE LOGSPLINE RESPONSE MODELS<sup>1</sup>

BY CHARLES J. STONE

*University of California, Berkeley*

Consider a  $\mathcal{Y}$ -valued response variable having a density function  $f(\cdot|x)$  that depends on an  $\mathcal{X}$ -valued input variable  $x$ . It is assumed that  $\mathcal{X}$  and  $\mathcal{Y}$  are compact intervals and that  $f(\cdot|\cdot)$  is continuous and positive on  $\mathcal{X} \times \mathcal{Y}$ . Let  $F(\cdot|x)$  denote the distribution function of  $f(\cdot|x)$  and let  $Q(\cdot|x)$  denote its quantile function. A finite-parameter exponential family model based on tensor-product  $B$ -splines is constructed. Maximum likelihood estimation of the parameters of the model based on independent observations of the response variable at fixed settings of the input variable yields estimates of  $f(\cdot|\cdot)$ ,  $F(\cdot|\cdot)$  and  $Q(\cdot|\cdot)$ . Under mild conditions, if the number of parameters suitably tends to infinity as  $n \rightarrow \infty$ , these estimates have optimal rates of convergence. The asymptotic behavior of the corresponding confidence bounds is also investigated.

**1. Discussion of results.** Consider a  $\mathcal{Y}$ -valued random response variable  $Y$  having an unknown density function  $f(\cdot|x)$  that depends on an  $\mathcal{X}$ -valued input variable  $x$ ; here  $\mathcal{X}$  and  $\mathcal{Y}$  are intervals in  $\mathbb{R}$  having positive length. It is assumed that  $f(\cdot|\cdot)$  is continuous and positive on  $\mathcal{X} \times \mathcal{Y}$ . Let  $F(\cdot|x)$  and  $Q(\cdot|x)$  denote the distribution function and quantile function, respectively, corresponding to  $f(\cdot|x)$ . Let fixed inputs (design points)  $x_1, \dots, x_n \in \mathcal{X}$  be given and let  $Y_1, \dots, Y_n$  be independent random variables such that  $Y_i$  has density function  $f(\cdot|x_i)$  for  $1 \leq i \leq n$ ; here  $Y_1, \dots, Y_n$  are the response variables corresponding to the settings  $x_1, \dots, x_n$ , respectively, of the input variable. Observations on these response variables can be used for inference concerning  $f(\cdot|\cdot)$ ,  $F(\cdot|\cdot)$  and  $Q(\cdot|\cdot)$ .

The classical approach is to assume a fixed parametric model  $f(\cdot|\theta_1, \dots, \theta_K)$  for the density function of  $Y$  and consider fixed parametric forms  $\theta_k = h_k(x; \beta_k)$  for the dependence of  $\theta_1, \dots, \theta_K$  on  $x$ . Normal linear models and generalized linear models in which  $Y$  has a gamma distribution with known shape parameter are of this form [see McCullagh and Nelder (1983)].

A refinement of the classical approach is to assume that  $\theta_k = h_k(x)$  for  $1 \leq k \leq K$ , where  $h_1, \dots, h_K$  are unknown continuous functions on  $\mathcal{X}$ , approximate these functions by members of some flexible  $J$ -dimensional linear space  $\mathcal{H}$  such as a space of polynomial, trigonometric series or polynomial splines, and let  $J \rightarrow \infty$ , as  $n \rightarrow \infty$ .

---

Received September 1989; revised December 1990.

<sup>1</sup>Research supported in part by NSF Grants DMS-86-00409 and DMS-89-02016.

AMS 1980 subject classifications. Primary 62G05; secondary 62F12.

Key words and phrases. Input-response model, exponential families, B-splines, maximum likelihood, rates of convergence.

A further refinement is to choose a flexible  $K$ -dimensional linear space  $\mathcal{S}$  and a basis  $B_1, \dots, B_K$  of  $\mathcal{S}$ , approximate  $\log f(\cdot|x)$  by  $\theta_1 B_1 + \dots + \theta_K B_K - C(\theta_1, \dots, \theta_K)$ , where  $C(\theta_1, \dots, \theta_K)$  is the normalizing constant, approximate the dependence of  $\theta_1, \dots, \theta_K$  on  $x$  by members of some flexible  $J$ -dimensional linear space  $\mathcal{H}$  and let  $J, K \rightarrow \infty$  as  $n \rightarrow \infty$ . This *doubly-flexible* approach will be pursued in the present paper, with  $\mathcal{S}$  and  $\mathcal{H}$  being spaces of polynomial splines.

For theoretical purposes,  $\mathcal{X}$  and  $\mathcal{Y}$  are required to be compact subintervals of  $\mathbb{R}$ . Let  $\mathcal{S}$  be a standard linear space of spline functions of a given order  $q \geq 1$  on  $\mathcal{Y}$  having dimension  $K \geq 2$ . (The functions in  $\mathcal{S}$  are piecewise polynomials of degree  $q - 1$  or less. If  $q = 1$ , we choose them to be right-continuous on  $\mathcal{Y}$  and continuous at the right endpoint of  $\mathcal{Y}$ ; if  $q \geq 2$ , they are  $(q - 2)$ -times continuously differentiable on  $\mathcal{Y}$ .) Let  $B_1, \dots, B_K$  be a basis of  $\mathcal{S}$  consisting of  $B$ -splines [see deBoor (1978)]. Then  $B_1, \dots, B_K$  are nonnegative and sum to one on  $\mathcal{Y}$ .

Given  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_K)^t \in \Theta$ , set  $s(y; \boldsymbol{\theta}) = \theta_1 B_1(y) + \dots + \theta_K B_K(y)$  for  $y \in \mathcal{Y}$ ,  $C(\boldsymbol{\theta}) = \log \int \exp(s(y; \boldsymbol{\theta})) dy$  and  $f(y; \boldsymbol{\theta}) = \exp(s(y; \boldsymbol{\theta}) - C(\boldsymbol{\theta}))$  for  $y \in \mathcal{Y}$ . Also, set

$$\Theta = \{\boldsymbol{\theta} = (\theta_1, \dots, \theta_K)^t \in \mathbb{R}^K: \theta_1 + \dots + \theta_K = 0\}.$$

Then  $f(\cdot; \boldsymbol{\theta})$ ,  $\boldsymbol{\theta} \in \Theta$ , defines an identifiable exponential family; it is referred to as a *logspline model* since  $\log f(\cdot; \boldsymbol{\theta}) \in \mathcal{S}$ . The theory of such models was developed in Stone (1990), which is a precursor to the present paper. Barron and Sheu (1991) independently obtained results for logspline models as well as for similar models involving polynomials and trigonometric series.

Let  $\mathcal{H}$  be a standard linear space of spline functions of a given order (which is not necessarily the same as that of  $\mathcal{S}$ ) having dimension  $J \geq 1$  and let  $H_1, \dots, H_J$  be a basis of  $\mathcal{H}$  consisting of  $B$ -splines. Let  $\mathcal{B}$  denote the collection of  $J \times K$  matrices  $\boldsymbol{\beta} = (\beta_{jk})$  of real numbers  $\beta_{jk}$ ,  $1 \leq j \leq J$  and  $1 \leq k \leq K$ , such that  $\sum_k \beta_{jk} = 0$  for  $1 \leq j \leq J$ , which can be regarded as a  $[J(K - 1)]$ -dimensional subspace of  $\mathbb{R}^{JK}$ . Let  $\boldsymbol{\beta} \in \mathcal{B}$ . For  $1 \leq k \leq K$ , let  $h_k(\cdot; \boldsymbol{\beta})$  be the real-valued function on  $\mathcal{X}$  defined by  $h_k(x; \boldsymbol{\beta}) = \sum_j \beta_{jk} H_j(x)$  for  $x \in \mathcal{X}$ . Set  $\mathbf{h}(x; \boldsymbol{\beta}) = (h_1(x; \boldsymbol{\beta}), \dots, h_K(x; \boldsymbol{\beta}))^t$  for  $x \in \mathcal{X}$  and observe that  $\mathbf{h}(\cdot; \boldsymbol{\beta})$  is a  $\Theta$ -valued function on  $\mathcal{X}$ . Also, set  $f(y|x; \boldsymbol{\beta}) = f(y; \mathbf{h}(x; \boldsymbol{\beta})) = \exp(s(y; \mathbf{h}(x; \boldsymbol{\beta})) - C(\mathbf{h}(x; \boldsymbol{\beta})))$  for  $\boldsymbol{\beta} \in \mathcal{B}$ ,  $x \in \mathcal{X}$  and  $y \in \mathcal{Y}$ . Then  $f(\cdot|x; \boldsymbol{\beta})$  is a positive density function on  $\mathcal{Y}$  for  $\boldsymbol{\beta} \in \mathcal{B}$  and  $x \in \mathcal{X}$ . We refer to  $f(\cdot|x; \boldsymbol{\beta})$ ,  $x \in \mathcal{X}$  and  $\boldsymbol{\beta} \in \mathcal{B}$ , as a *logspline response model*.

The log-likelihood function  $l(\boldsymbol{\beta})$ ,  $\boldsymbol{\beta} \in \mathcal{B}$ , is defined by

$$l(\boldsymbol{\beta}) = \sum_i \log f(Y_i|x_i; \boldsymbol{\beta}) = \sum_i [s(Y_i; \mathbf{h}(x_i; \boldsymbol{\beta})) - C(\mathbf{h}(x_i; \boldsymbol{\beta}))], \quad \boldsymbol{\beta} \in \mathcal{B}.$$

Set  $c(\boldsymbol{\beta}) = \sum_i C(\mathbf{h}(x_i; \boldsymbol{\beta}))$  for  $\boldsymbol{\beta} \in \mathcal{B}$ . Then  $l(\boldsymbol{\beta}) = \sum_i s(Y_i; \mathbf{h}(x_i; \boldsymbol{\beta})) - c(\boldsymbol{\beta})$  for  $\boldsymbol{\beta} \in \mathcal{B}$ . The expected log-likelihood function  $\lambda(\boldsymbol{\beta})$ ,  $\boldsymbol{\beta} \in \mathcal{B}$ , is given by

$$\lambda(\boldsymbol{\beta}) = El(\boldsymbol{\beta}) = \sum_i \int s(y; \mathbf{h}(x_i; \boldsymbol{\beta})) f(y|x_i) dy - c(\boldsymbol{\beta}), \quad \boldsymbol{\beta} \in \mathcal{B}.$$

The functions  $l(\cdot)$ ,  $c(\cdot)$  and  $\lambda(\cdot)$  can also be viewed as functions on  $\mathbb{R}^{JK}$ . For  $\beta \in \mathcal{B}$ , let  $\mathbf{I}(\beta)$  denote the corresponding information matrix, which equals the Hessian matrix of  $c(\cdot)$  at  $\beta$  and is a positive semidefinite symmetric  $JK \times JK$  matrix. Thus  $c(\cdot)$  is a convex function on  $\mathcal{B}$  and  $l(\cdot)$  and  $\lambda(\cdot)$  are concave functions on  $\mathcal{B}$ .

It is assumed that  $\mathcal{H}$  is nonsingular relative to the design set: If  $h \in \mathcal{H}$  and  $h(x_1) = \cdots = h(x_n) = 0$ , then  $h = 0$  on  $\mathcal{X}$ . Then  $\tau' \mathbf{I}(\beta) \tau > 0$  for  $\beta \in \mathcal{B}$  and  $\tau \in \mathcal{B}$  with  $\tau \neq 0$ . Thus  $c(\cdot)$  is strictly convex and  $l(\cdot)$  and  $\lambda(\cdot)$  are strictly concave on  $\mathcal{B}$ .

Let  $\hat{\beta}$  denote the maximum likelihood estimate of  $\beta$ ; that is, the value of  $\beta \in \mathcal{B}$  that maximizes the log-likelihood function. Then  $\hat{\beta}$  may or may not exist. Under the nonsingularity assumption on  $\mathcal{H}$ , if  $\hat{\beta}$  exists, then it is unique. Given  $x \in \mathcal{X}$ , consider the maximum likelihood estimate  $\hat{f}(\cdot|x) = f(\cdot|x; \hat{\beta})$  of  $f(\cdot|x)$  and let  $\hat{F}(\cdot|x)$  and  $\hat{Q}(\cdot|x)$  denote the corresponding maximum likelihood estimates of  $F(\cdot|x)$  and  $Q(\cdot|x)$ .

Similarly,  $\lambda(\cdot)$  has at most one maximum on  $\mathcal{B}$ . It is easily seen that  $\lambda(\cdot)$  does have a maximum on  $\mathcal{B}$  and hence a unique maximum  $\beta^*$  on  $\mathcal{B}$ . Consider the function  $f^*(\cdot|\cdot)$  on  $\mathcal{X} \times \mathcal{Y}$  defined by  $f^*(y|x) = f(y|x; \beta^*)$  for  $x \in \mathcal{X}$  and  $y \in \mathcal{Y}$ . Let  $F^*(\cdot|x)$  and  $Q^*(\cdot|x)$  denote the distribution function and quantile function, respectively, corresponding to  $f^*(\cdot|x)$ .

The knot sequences defining  $\mathcal{S}$  and  $\mathcal{H}$  are allowed depending on  $n$ , but it is assumed that they are  $\sigma$ -quasiuniform in the sense of Schumaker [(1981), page 216]: The ratios of the differences between consecutive knots are bounded away from zero and infinity uniformly in  $n$ . We make the mild assumption on the design points that there is an  $M > 0$  (independent of  $n$ ) such that, for  $n$  sufficiently large ( $n \gg 1$ ),

$$(1) \quad M^{-1}n \int h^2(x) dx \leq \sum_i h^2(x_i) \leq Mn \int h^2(x) dx, \quad h \in \mathcal{H}.$$

[The nonsingularity assumption on  $\mathcal{H}$  is an immediate consequence of (1).]

Given a subinterval  $I$  of  $\mathcal{X}$ , let  $|I|$  denote the length of  $I$  and set  $N(I) = \#\{i: x_i \in I\}$ . Under (7) below, in light of the  $\sigma$ -quasiuniformity of the knot sequence defining  $\mathcal{H}$ , a sufficient condition for (1) is that for every  $\delta > 0$ , there is an  $M > 0$  such that for  $n \gg 1$ ,

$$(2) \quad \begin{aligned} M^{-1}n|I| &\leq N(I) \\ &\leq Mn|I| \quad \text{for every subinterval } I \text{ of } \mathcal{X} \text{ such that } |I| \geq n^{\delta-1}. \end{aligned}$$

Let  $\mathcal{T} = \mathcal{H} \otimes \mathcal{S}$  denote the tensor product of  $\mathcal{H}$  and  $\mathcal{S}$ ; that is, the linear space of real-valued functions on  $\mathcal{X} \times \mathcal{Y}$  spanned by functions of the form  $h \otimes s$  as  $h$  and  $s$  range over  $\mathcal{H}$  and  $\mathcal{S}$ , respectively; here  $(h \otimes s)(x, y) = h(x)s(y)$  for  $x \in \mathcal{X}$  and  $y \in \mathcal{Y}$ . Then  $\mathcal{T}$  has dimension  $JK$  and the functions  $H_j \otimes B_k$ ,  $1 \leq j \leq J$  and  $1 \leq k \leq K$ , form a basis of  $\mathcal{T}$ .

Given a real-valued function  $g(\cdot|\cdot)$  on  $\mathcal{X} \times \mathcal{Y}$ , set

$$\|g(\cdot|\cdot)\|_\infty = \sup_{x \in \mathcal{X}} \sup_{y \in \mathcal{Y}} g(y|x).$$

Also, set  $\delta_{\mathcal{J}} = \inf_{t \in \mathcal{J}} \|\log f(\cdot | \cdot) - t\|_{\infty}$ . Under the  $\sigma$ -quasiuniform condition on the knot sequences,  $\delta_{\mathcal{J}} \rightarrow 0$  as  $J, K \rightarrow \infty$ ; see Schumaker [(1981), Theorem 12.8] for this result and for an upper bound to  $\delta_{\mathcal{J}}$  in terms of the smoothness of  $\log f(\cdot | \cdot)$ .

Since  $f(\cdot | \cdot)$  is continuous and positive on the compact set  $\mathcal{X} \times \mathcal{Y}$ ,  $\log f(\cdot | \cdot)$  is bounded and continuous on this set. Under (1) and the  $\sigma$ -quasiuniform condition on the knot sequences, it was shown in Stone (1989) that

$$(3) \quad \|\log f(\cdot | \cdot) - \log f^*(\cdot | \cdot)\|_{\infty} = O(\delta_{\mathcal{J}}).$$

It follows from (3) that

$$(4) \quad \|f(\cdot | \cdot) - f^*(\cdot | \cdot)\|_{\infty} = O(\delta_{\mathcal{J}}),$$

$$(5) \quad \|F(\cdot | \cdot) - F^*(\cdot | \cdot)\|_{\infty} = O(\delta_{\mathcal{J}})$$

and

$$(6) \quad \|Q(\cdot | \cdot) - Q^*(\cdot | \cdot)\|_{\infty} = O(\delta_{\mathcal{J}}).$$

[In (6) the supremum is over  $p$  and  $x$  with  $0 < p < 1$  and  $x \in \mathcal{X}$ .]

From now on, it is assumed that

$$(7) \quad JK = o(n^{(1/2)-\varepsilon}) \quad \text{for some } \varepsilon \in (0, \tfrac{1}{2}).$$

This is slightly stronger than the assumption  $JK = o(\sqrt{n})$ , which arises in Portnoy (1986, 1988).

In Section 2 it will be shown that  $\hat{\beta}$  exists except on an event whose probability tends to zero with  $n$ . There the asymptotic behavior of  $\hat{\beta}$  will also be determined.

In Section 3, it will be shown that

$$(8) \quad \hat{f}(y|x) - f^*(y|x) = O_P(\sqrt{JK/n}),$$

$$(9) \quad \int [\hat{f}(y|x) - f^*(y|x)]^2 dy = O_P(JK/n),$$

$$(10) \quad \max_{x,y} |\hat{f}(y|x) - f^*(y|x)| = O_P(\sqrt{JK(\log JK)/n}),$$

$$(11) \quad \max_y |\hat{F}(y|x) - F^*(y|x)| = O_P(\sqrt{J/n})$$

and

$$(12) \quad \max_p |\hat{Q}(p|x) - Q^*(p|x)| = O_P(\sqrt{J/n}).$$

In (8),  $x \in \mathcal{X}$  and  $y \in \mathcal{Y}$  are fixed, while in (9), (11) and (12),  $x \in \mathcal{X}$  is fixed. The order of magnitude  $\sqrt{JK}$  in (8) is plausible: There are about  $n/(JK)$  trials per unknown parameter  $\beta_{jk}$ , so the asymptotic standard deviation of the  $\hat{\beta}_{jk}$ 's should be proportional to  $\sqrt{JK/n}$ . In light of the local support of the  $B$ -splines, the asymptotic standard deviation of  $\hat{f}(y|x)$  should have the same order of magnitude as that of the  $\hat{\beta}_{jk}$ 's.

Suppose that  $\delta_{\mathcal{G}} = O(J^{-p_1} + K^{-p_2})$ , where  $p_1 > \frac{1}{2}$  and  $p_2 > \frac{1}{2}$ . Set

$$p = \frac{2}{p_1^{-1} + p_2^{-1}}, \quad \gamma_1 = \frac{p}{p_1(2p+2)} = \frac{p_2}{p_1 + p_2 + 2p_1p_2}$$

and

$$\gamma_2 = \frac{p}{p_2(2p+2)} = \frac{p_1}{p_1 + p_2 + 2p_1p_2}.$$

Let  $a_n \sim b_n$  mean that  $a_n/b_n$  is bounded away from zero and infinity as  $n \rightarrow \infty$ .

Suppose that  $p > 1$ . Choose  $J$  and  $K$  such that  $J \sim n^{\gamma_1}$  and  $K \sim n^{\gamma_2}$ . Then  $\delta_{\mathcal{G}} = O(n^{-p/(2p+2)})$ . Also  $JK \sim n^{1/(p+1)}$ , so (7) holds. By (4) and (8),

$$(13) \quad \hat{f}(y|x) - f(y|x) = O_P(n^{-p/(2p+2)});$$

by (4) and (9),

$$(14) \quad \int [\hat{f}(y|x) - f(y|x)]^2 dy = O_P(n^{-2p/(2p+2)}) = O_P(n^{-p/(p+1)}).$$

Choose  $J$  and  $K$  such that  $J \sim (n/\log n)^{\gamma_1}$  and  $K \sim (n/\log n)^{\gamma_2}$ . Then

$$\delta_{\mathcal{G}} = O((n/\log n)^{-p/(2p+2)}).$$

Also  $JK \sim (n/\log n)^{1/(p+1)}$ , so (7) again holds. By (4) and (10),

$$(15) \quad \max_{x,y} |\hat{f}(y|x) - f(y|x)| = O_P((n/\log n)^{-p/(2p+2)}).$$

Suppose  $p_1 > \frac{1}{2}$  and  $p_2 > 2p_1/(2p_1 - 1)$  and set  $a = p_1/[p_2(2p_1 + 1)]$ . Choose  $J$  and  $K$  such that  $J \sim n^{(1/(2p_1+1))}$ ,  $K^{-1} = O(n^{-a})$  and (7) holds. Then  $\delta_{\mathcal{G}} = O(n^{-p_1/(2p_1+1)})$ . By (5) and (11),

$$(16) \quad \max_y |\hat{F}(y|x) - F(y|x)| = O_P(n^{-p_1/(2p_1+1)});$$

by (6) and (12),

$$(17) \quad \max_p |\hat{Q}(p|x) - Q(p|x)| = O_P(n^{-p_1/(2p_1-1)}).$$

Under reasonable further specifications, the rates of convergence in (13)–(17) are optimal [see Stone (1980, 1982) and Hasminskii and Ibragimov (1990)]. For fixed  $y \in \mathcal{Y}$ , the rate of convergence  $\hat{F}(y|x) - F(y|x) = O_P(n^{-p_1/(2p_1+1)})$  can be achieved by using a different estimate under the corresponding smoothness assumption on  $F(y|x)$  as a function of  $x$  without having to make any smoothness assumption on  $F(y|x)$  as a function of  $y$ . [Observe that  $F(y|x) = E(\text{ind}(Y \leq y)|X = x)$  and see Stone (1980).]

Let  $J, K \rightarrow \infty$  as  $n \rightarrow \infty$ ; let  $\tau^*$  be defined as  $f^*(y|x)$ ,  $F^*(y|x)$  with  $y$  in the interior of  $\mathcal{Y}$ , or  $Q^*(p|x)$  with  $0 < p < 1$ ; and let  $\hat{\tau}$  be defined as the corresponding maximum likelihood estimate  $\hat{f}(y|x)$ ,  $\hat{F}(y|x)$  or  $\hat{Q}(p|x)$ . Let  $\text{ASD}(\hat{\tau})$  and  $\text{SE}(\hat{\tau})$  denote the asymptotic standard deviation and standard error, respectively, of  $\hat{\tau}$ , as usually defined in terms of the information matrix

in large-sample parametric inference. Then, uniformly for  $x \in \mathcal{X}$ ,  $\text{SE}(\hat{\tau})/\text{ASD}(\hat{\tau}) = 1 + o_p(1)$  and the distributions of  $(\hat{\tau} - \tau^*)/\text{ASD}(\hat{\tau})$  and  $(\hat{\tau} - \tau^*)/\text{SE}(\hat{\tau})$  converge to the standard normal distribution as  $n \rightarrow \infty$ . These results will be verified in Section 4, where explicit formulas for the various asymptotic standard deviations and standard errors will be given.

According to the last result, for  $0 < \alpha < 1$ ,  $\hat{\tau} \pm z_{1-\alpha/2}\text{SE}(\hat{\tau})$  is an asymptotic  $100(1 - \alpha)\%$  confidence interval for  $\tau^*$ ; here  $z_{1-\alpha/2}$  is the  $(1 - \alpha/2)$ th quantile of the standard normal distribution. Such confidence intervals are useful in practice, but they must be interpreted with care. Under the additional, but dubious, assumption that  $(\tau - \tau^*)/\text{ASD}(\hat{\tau}) = o(1)$ , the confidence intervals for  $\tau^*$  can be interpreted as confidence intervals for  $\tau$  itself.

The arguments used in Sections 2–4, which are natural outgrowths of those developed in Stone (1985, 1986 and 1990), also apply when the fixed design is replaced by a random sample from the distribution of a random variable  $X$  having a density function that is bounded away from zero and infinity on  $\mathcal{X}$  [in which case, a suitable probabilistic version of (2) is easily verified]. Alternatively, the joint density function  $f(\cdot, \cdot)$  can be estimated by  $\hat{f}(x, y) = f(x, y; \hat{\beta}_1) = \exp[\sum_j \sum_k \hat{\beta}_{1jk} H_j(x) B_k(y) - c(\hat{\beta}_1)]$ , where  $\hat{\beta}_1 = (\hat{\beta}_{1jk})$  is the maximum likelihood estimate and  $c(\hat{\beta}_1)$  is the normalizing constant. The asymptotic behavior of this estimate follows from results in Barron and Sheu (1991) or from the extension of results in Stone (1990) given in Koo (1988). The corresponding estimate of the marginal density function of  $X$  is given by  $f(x; \hat{\beta}_1) = \int f(x, y; \hat{\beta}_1) dy$ . This leads to the alternative conditional density function estimate  $\hat{f}_1(y|x) = f(x, y; \hat{\beta}_1)/f(x; \hat{\beta}_1)$  which has the same form as the estimate  $f(y|x)$  defined above, but with an estimate  $\hat{\beta}_1$  that differs somewhat from  $\hat{\beta}$ . The alternative conditional density estimate inherits the accuracy of the corresponding estimate of the joint density function. [The preceding remarks in this paragraph were suggested by a reviewer. It should be pointed out that the alternative estimate of the conditional density function achieves the rates of convergence obtained in the present paper only under an auxiliary smoothness assumption on the marginal density function of  $X$ . In the related context of nonparametric regression, Fan (1990) refers to estimates of the regression function that require such an auxiliary assumption as being *inadmissible*: They are dominated by estimates that achieve the optimal rate of convergence without requiring the auxiliary smoothness assumption on the marginal density function of  $X$ .]

The numerical and practical aspects of logspline modelling were treated in Stone and Koo (1986) and Kooperberg and Stone (1991). The results to date clearly indicate that logspline modelling is competitive with other approaches such as kernel density estimation. A numerical investigation of logspline response modelling has yet to be carried out. Such a study would undoubtedly go beyond what is mathematically tractable. In particular,  $\mathcal{X}$  and  $\mathcal{Y}$  could be unbounded if linear restrictions were imposed on the tails of the various splines entering into the model. Also, it would be worthwhile to study stepwise selection of the basis functions of the model, as introduced in Smith (1982) and used successfully by Breiman and Peters (1988), Friedman and Silverman

(1989), Friedman (1991), Breiman (1989, 1991), Kooperberg and Stone (1991) and Jin (1990).

**2. Parameter estimation.** For  $\mathbf{b} = (b_{jk}) \in \mathcal{B}$ , let  $|\mathbf{b}|$  denote the non-negative square root of  $\sum_j \sum_k b_{jk}^2$ . In the next result,  $x \in \mathcal{X}$ ,  $y \in \mathcal{Y}$  and  $0 < p < 1$ ; the quantities  $j$  and  $k$  in (b) and the quantity  $j$  in (c) are allowed to depend on  $n$  in an arbitrary deterministic manner.

**THEOREM 1.** (a)  $\hat{\beta}$  exists except on an event whose probability tends to zero as  $n \rightarrow \infty$ .

$$(b) \hat{\beta}_{jk} - \beta_{jk}^* = O_P(\sqrt{JK/n}).$$

$$(c) (1/K) \sum_k (\hat{\beta}_{jk} - \beta_{jk}^*)^2 = O_P(JK/n).$$

$$(c) |\hat{\beta} - \beta^*|^2 = O_P(J^2 K^2 / n).$$

$$(e) \max_{j,k} |\hat{\beta}_{jk} - \beta_{jk}^*| = O_P(\sqrt{JK(\log JK)/n}).$$

The proof of Theorem 1 is divided into a number of lemmas. For  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_K)^t \in \Theta$ , let  $|\boldsymbol{\theta}|$  denote the nonnegative square root of  $\sum_k \theta_k^2$ . Also, let  $\|s\|_2$  and  $\|s\|_\infty$  be defined in the usual manner for functions  $s$  on  $\mathcal{Y}$ .

**LEMMA 1.** Let  $M > 0$ . Then there is an  $M_1 > 0$  such that if  $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \Theta$ ,  $\|s(\cdot; \boldsymbol{\theta}_1)\|_\infty \leq M$  and  $\|s(\cdot; \boldsymbol{\theta}_2)\|_\infty \leq M$ , then

$$[C(\boldsymbol{\theta}_2) - C(\boldsymbol{\theta}_1)]^2 \leq M_1 \|s(\cdot; \boldsymbol{\theta}_2) - s(\cdot; \boldsymbol{\theta}_1)\|_2^2.$$

**PROOF.** Since

$$C(\boldsymbol{\theta}_2) - C(\boldsymbol{\theta}_1) = \log \int e^{s(y; \boldsymbol{\theta}_2)} dy - \log \int e^{s(y; \boldsymbol{\theta}_1)} dy$$

and  $0 < \text{length}(\mathcal{Y}) < \infty$ , the desired result follows from the Schwarz inequality and elementary properties of the exponential and logarithm functions.  $\square$

**LEMMA 2.** Let  $M > 0$ . Then there is an  $M_1 > 0$  such that if  $\boldsymbol{\theta}, \boldsymbol{\theta}^* \in \Theta$ ,  $\|\log f(\cdot; \boldsymbol{\theta}^*)\|_\infty \leq M$  and  $\|s(\cdot; \boldsymbol{\theta}) - s(\cdot; \boldsymbol{\theta}^*)\|_\infty \leq M$ , then

$$(a) \|\log f(\cdot; \boldsymbol{\theta}) - \log f(\cdot; \boldsymbol{\theta}^*)\|_\infty \leq M_1$$

and

$$(b) M_1^{-1} K^{-1} |\boldsymbol{\theta} - \boldsymbol{\theta}^*|^2 \leq \|\log f(\cdot; \boldsymbol{\theta}) - \log f(\cdot; \boldsymbol{\theta}^*)\|_2^2 \leq M_1 K^{-1} |\boldsymbol{\theta} - \boldsymbol{\theta}^*|^2.$$

**PROOF.** By Lemma 3 of Stone (1990), there is an  $M_2 > 0$  such that  $|C(\boldsymbol{\theta}^*)| \leq M_2$ . Thus  $\|s(\cdot; \boldsymbol{\theta}^*)\|_\infty \leq M + M_2$  and hence  $\|s(\cdot; \boldsymbol{\theta})\|_\infty \leq 2M + M_2$ . Consequently, there is an  $M_3 > 0$  such that  $|C(\boldsymbol{\theta})| \leq M_3$  and hence  $\|\log f(\cdot; \boldsymbol{\theta}) - \log f(\cdot; \boldsymbol{\theta}^*)\|_\infty \leq M + M_2 + M_3$ , which yields (a). According to (12) of Stone (1986), there is an  $M_4 > 0$  such that  $M_4^{-1} K^{-1} |\boldsymbol{\theta} - \boldsymbol{\theta}^*|^2 \leq \|s(\cdot; \boldsymbol{\theta}) - s(\cdot; \boldsymbol{\theta}^*)\|_2^2 \leq M_4 K^{-1} |\boldsymbol{\theta} - \boldsymbol{\theta}^*|^2$ . Hence, by Lemma 1, there is an  $M_5 > 0$  such that

$$\|\log f(\cdot; \boldsymbol{\theta}) - \log f(\cdot; \boldsymbol{\theta}^*)\|_2^2 \leq M_5 K^{-1} |\boldsymbol{\theta} - \boldsymbol{\theta}^*|^2.$$

Observe that

$$\begin{aligned}\|\log f(\cdot; \boldsymbol{\theta}) - \log f(\cdot; \boldsymbol{\theta}^*)\|_2^2 &= \int \left( \sum_k (\theta_k - \theta_k^*) B_k(y) - [C(\boldsymbol{\theta}) - C(\boldsymbol{\theta}^*)] \right)^2 dy \\ &= \int \left( \sum_k \{ \theta_k - \theta_k^* - [C(\boldsymbol{\theta}) - C(\boldsymbol{\theta}^*)] \} B_k(y) \right)^2 dy.\end{aligned}$$

Thus, by (12) of Stone (1986), there is an  $M_6 > 0$  such that

$$\|\log f(\cdot; \boldsymbol{\theta}) - \log f(\cdot; \boldsymbol{\theta}^*)\|_2^2 \geq M_6^{-1} K^{-1} \sum_k \{ \theta_k - \theta_k^* - [C(\boldsymbol{\theta}) - C(\boldsymbol{\theta}^*)] \}^2.$$

Now  $\boldsymbol{\theta}, \boldsymbol{\theta}^* \in \Theta$ , so  $\sum_k (\theta_k - \theta_k^*) = 0$  and hence

$$\sum_k \{ \theta_k - \theta_k^* - [C(\boldsymbol{\theta}) - C(\boldsymbol{\theta}^*)] \}^2 \geq \sum_k (\theta_k - \theta_k^*)^2 = |\boldsymbol{\theta} - \boldsymbol{\theta}^*|^2.$$

Consequently,  $\|\log f(\cdot; \boldsymbol{\theta}) - \log f(\cdot; \boldsymbol{\theta}^*)\|_2^2 \geq M_6^{-1} K^{-1} |\boldsymbol{\theta} - \boldsymbol{\theta}^*|^2$ . Therefore, (b) is valid.  $\square$

LEMMA 3.  $\|s(\cdot; \mathbf{h}(\cdot; \boldsymbol{\beta}_2)) - s(\cdot; \mathbf{h}(\cdot; \boldsymbol{\beta}_1))\|_\infty \leq |\boldsymbol{\beta}_2 - \boldsymbol{\beta}_1|$  for  $\boldsymbol{\beta}_1, \boldsymbol{\beta}_2 \in \mathcal{B}$ .

PROOF. By the properties of  $B$ -splines,  $\sum_j \sum_k H_j^2(x) B_k^2(y) \leq \sum_j \sum_k H_j(x) B_k(y) = 1$ . Thus, by the Schwarz inequality,

$$\begin{aligned}[s(y; \mathbf{h}(x; \boldsymbol{\beta}_2)) - s(y; \mathbf{h}(x; \boldsymbol{\beta}_1))]^2 &= \left( \sum_j \sum_k (\beta_{2jk} - \beta_{1jk}) H_j(x) B_k(y) \right)^2 \\ &\leq |\boldsymbol{\beta}_2 - \boldsymbol{\beta}_1|^2,\end{aligned}$$

where  $\boldsymbol{\beta}_1 = (\beta_{1jk})$  and  $\boldsymbol{\beta}_2 = (\beta_{2jk})$ .  $\square$

The next result follows from (1), the  $\sigma$ -quasiuniformity of the knot sequence defining  $\mathcal{H}$  and (viii) of de Boor [(1978), page 155] [see the proof of (12) of Stone (1986)].

LEMMA 4. *There is an  $M > 0$  such that, for  $n \gg 1$ ,*

$$\begin{aligned}M^{-1} n J^{-1} |\boldsymbol{\beta}_2 - \boldsymbol{\beta}_1|^2 &\leq \sum_i |\mathbf{h}(x_i; \boldsymbol{\beta}_2) - \mathbf{h}(x_i; \boldsymbol{\beta}_1)|^2 \\ &\leq M n J^{-1} |\boldsymbol{\beta}_2 - \boldsymbol{\beta}_1|^2, \quad \boldsymbol{\beta}_1, \boldsymbol{\beta}_2 \in \mathcal{B}.\end{aligned}$$

Let  $\boldsymbol{\beta} \in \mathcal{B}$ . Then, by a direct computation,

$$(18) \quad \boldsymbol{\tau}' \mathbf{I}(\boldsymbol{\beta}) \boldsymbol{\tau} = \sum_i \int [s(y; \mathbf{h}(x_i; \boldsymbol{\tau})) - a_i]^2 f(y|x_i; \boldsymbol{\beta}) dy, \quad \boldsymbol{\tau} \in \mathcal{B},$$



where  $a_i = \int s(y; \mathbf{h}(x_i; \boldsymbol{\tau})) f(y|x_i; \boldsymbol{\beta}) dy$ . Let  $\boldsymbol{\beta} \in \mathcal{B}$ . Then  $(d/dt)\lambda(\boldsymbol{\beta}^* + t(\boldsymbol{\beta} - \boldsymbol{\beta}^*))|_{t=0} = 0$  and

$$\frac{d^2}{dt^2}\lambda(\boldsymbol{\beta}^* + t(\boldsymbol{\beta} - \boldsymbol{\beta}^*)) = -(\boldsymbol{\beta} - \boldsymbol{\beta}^*)^t \mathbf{I}(\boldsymbol{\beta}^* + t(\boldsymbol{\beta} - \boldsymbol{\beta}^*))(\boldsymbol{\beta} - \boldsymbol{\beta}^*).$$

Thus, by (18),

$$(19) \quad \lambda(\boldsymbol{\beta}) - \lambda(\boldsymbol{\beta}^*) = -\int_0^1 (1-t) \left( \sum_i \int [s(y; \mathbf{h}(x_i; \boldsymbol{\beta} - \boldsymbol{\beta}^*)) - a_i(t)]^2 \times f(y|x_i; \boldsymbol{\beta} + t(\boldsymbol{\beta} - \boldsymbol{\beta}^*)) dy \right) dt,$$

where  $a_i(t) = \int s(y; \mathbf{h}(x_i; \boldsymbol{\beta} - \boldsymbol{\beta}^*)) f(y|x_i; \boldsymbol{\beta} + t(\boldsymbol{\beta} - \boldsymbol{\beta}^*)) dy$  for  $1 \leq i \leq n$ .

Choose  $\varepsilon \in (0, \frac{1}{2})$  satisfying (7):  $JK = o(n^{(1/2)-\varepsilon})$ .

**LEMMA 5.** *There is a  $\delta > 0$  such that if  $n \gg 1$ ,  $\boldsymbol{\beta} \in \mathcal{B}$  and  $|\boldsymbol{\beta} - \boldsymbol{\beta}^*| = n^\varepsilon JK / \sqrt{n}$ , then  $\lambda(\boldsymbol{\beta}) - \lambda(\boldsymbol{\beta}^*) \leq -\delta n^{2\varepsilon} JK$ .*

**PROOF.** It follows from (3), (19), Lemma 2(a) and Lemma 3 [see the proof of Lemma 4 in Stone (1990)] that there is a  $\delta_1 > 0$  such that

$$\lambda(\boldsymbol{\beta}) - \lambda(\boldsymbol{\beta}^*) \leq -\delta_1 \sum_i \int [\log f(y|x_i; \boldsymbol{\beta}) - \log f(y|x_i; \boldsymbol{\beta}^*)]^2 dy.$$

By lemmas 2(b), 3 and 4, there is a  $\delta_2 > 0$  such that

$$\sum_i \int [\log f(y|x_i; \boldsymbol{\beta}) - \log f(y|x_i; \boldsymbol{\beta}^*)]^2 dy \geq \delta_2 n J^{-1} K^{-1} |\boldsymbol{\beta} - \boldsymbol{\beta}^*|^2.$$

Consequently, the desired result holds.  $\square$

**LEMMA 6.** *Let  $\delta > 0$ . Then there is a  $\delta_1 > 0$  such that if  $n \gg 1$ ,  $\boldsymbol{\beta} \in \mathcal{B}$  and  $|\boldsymbol{\beta} - \boldsymbol{\beta}^*| \leq n^\varepsilon JK / \sqrt{n}$ , then*

$$P\left(l(\boldsymbol{\beta}) - l(\boldsymbol{\beta}^*) - [\lambda(\boldsymbol{\beta}) - \lambda(\boldsymbol{\beta}^*)] \geq \frac{\delta}{2} n^{2\varepsilon} JK\right) \leq \exp(-\delta_1 n^{2\varepsilon} JK).$$

**PROOF.** Write  $l(\boldsymbol{\beta}) - l(\boldsymbol{\beta}^*) - [\lambda(\boldsymbol{\beta}) - \lambda(\boldsymbol{\beta}^*)] = \sum_i Z_i$ , where

$$\begin{aligned} Z_i &= \log f(Y_i|x_i; \boldsymbol{\beta}) - \log f(Y_i|x_i; \boldsymbol{\beta}^*) \\ &\quad - E[\log f(Y_i|x_i; \boldsymbol{\beta}) - \log f(Y_i|x_i; \boldsymbol{\beta}^*)]. \end{aligned}$$

It follows from (3), Lemma 2(a) and Lemma 3 that there is an  $M_1 > 0$  such

that  $P(|Z_i| \leq M_1) = 1$  for  $1 \leq i \leq n$ . Observe that, for  $1 \leq i \leq n$ ,  $EZ_i = 0$  and

$$\begin{aligned} \text{var}(Z_i) &\leq E\left\{\left[\log f(Y_i|x_i; \beta) - \log f(Y_i|x_i; \beta^*)\right]^2\right\} \\ &= \int \left[\log f(y|x_i; \beta) - \log f(y|x_i; \beta^*)\right]^2 f(y|x_i) dy. \end{aligned}$$

We now conclude from Lemmas 2(b), 3 and 4 and the boundedness of  $f(\cdot|\cdot)$  that there is an  $M_2 > 0$  such that  $\sum_i \text{var}(Z_i) \leq M_2 n^{2\epsilon} JK$ . The desired result now follows from Bernstein's inequality [see (2.13) in Hoeffding, (1963)].  $\square$

LEMMA 7. *Let  $\delta > 0$ . Then there is a  $\delta_1 > 0$  such that*

$$|l(\beta_2) - l(\beta_1) - [\lambda(\beta_2) - \lambda(\beta_1)]| \leq \frac{\delta}{2} n^{2\epsilon} JK$$

for  $n \gg 1$ ,  $\beta_1, \beta_2 \in \mathcal{B}$ ,  $|\beta_1 - \beta^*| = n^\epsilon JK / \sqrt{n}$ ,  $|\beta_2 - \beta^*| = n^\epsilon JK / \sqrt{n}$  and  $|\beta_2 - \beta_1| \leq \delta_1 n^{2\epsilon-1} JK$ .

PROOF. Observe that

$$|l(\beta_2) - l(\beta_1) - [\lambda(\beta_2) - \lambda(\beta_1)]| \leq 2n \|\log f(\cdot|\cdot; \beta_2) - \log f(\cdot|\cdot; \beta_1)\|_\infty.$$

By Lemmas 1 and 3, there is an  $M_1 > 0$  such that  $\|\log f(\cdot|\cdot; \beta_2) - \log f(\cdot|\cdot; \beta_1)\|_\infty \leq M_1 |\beta_2 - \beta_1|$ . Thus the desired result is valid.  $\square$

The diameter of a subset  $B$  of  $\mathcal{B}$  is defined as  $\sup\{|\beta_2 - \beta_1|: \beta_1, \beta_2 \in B\}$ . The next result is easily established by considering suitable inscribed and circumscribed  $JK$ -dimensional cubes.

LEMMA 8. *Let  $\delta_1 > 0$ . Then there is an  $M > 0$  such that, for  $n \gg 1$ ,*

$$\{\beta \in \mathcal{B}: |\beta - \beta^*| = n^\epsilon JK / \sqrt{n}\}$$

can be covered by  $\exp(MJK(\log n))$  subsets of  $\mathcal{B}$  each having diameter at most  $\delta_1 n^{2\epsilon-1} JK$ .

LEMMA 9. (a)  $\hat{\beta}$  exists except on an event whose probability tends to zero as  $n \rightarrow \infty$ .

(b)  $|\hat{\beta} - \beta^*| = O_p(n^\epsilon JK / \sqrt{n})$ .

PROOF. Set  $\mathcal{B}_1 = \{\beta \in \mathcal{B}: |\beta - \beta^*| \leq n^\epsilon JK / \sqrt{n}\}$ . Then  $\mathcal{B}_1$  is a compact set whose boundary relative to  $\mathcal{B}$  is  $\mathcal{B}_2 = \{\beta \in \mathcal{B}: |\beta - \beta^*| = n^\epsilon JK / \sqrt{n}\}$ . By Lemma 5 there is a  $\delta > 0$  such that  $\lambda(\beta) - \lambda(\beta^*) \leq -\delta n^{2\epsilon} JK$  for  $\beta \in \mathcal{B}_2$ . Thus it follows from Lemmas 6–8 that, except on an event whose probability tends to zero as  $n \rightarrow \infty$ ,  $l(\beta) < l(\beta^*)$  for  $\beta \in \mathcal{B}_2$ , so  $l(\cdot)$  has a local maximum in the interior of  $\mathcal{B}_1$  relative to  $\mathcal{B}$ . The desired results now follow from the strict concavity of  $l(\cdot)$  on  $\mathcal{B}$ .  $\square$

The next result follows from (4), (7) and Lemmas 2a, 3 and 9.

LEMMA 10. *There is a positive constant  $M$  such that, except on an event whose probability tends to zero as  $n \rightarrow \infty$ ,  $M^{-1} \leq f(\cdot | \cdot; \beta^* + t(\hat{\beta} - \beta^*)) \leq M$  for  $0 \leq t \leq 1$ .*

LEMMA 11. *There is an  $M > 0$  such that, for  $n \gg 1$ ,  $\beta \in \mathcal{B}$  and  $\tau \in \mathcal{B}$ ,*  
 $M^{-1}nJ^{-1}K^{-1}|\tau|^2 \min f(\cdot | \cdot; \beta) \leq \tau' \mathbf{I}(\beta) \tau \leq MnJ^{-1}K^{-1}|\tau|^2 \max f(\cdot | \cdot; \beta).$

PROOF. Set  $\min = \min f(\cdot | \cdot; \beta)$  and  $\max = \max f(\cdot | \cdot; \beta)$ . Using (1), (18) and (12) of Stone (1986), we see that there are positive numbers  $\delta_1$ ,  $\delta_2$  and  $\delta$  such that, for  $n \gg 1$ ,  $\beta \in \mathcal{B}$  and  $\tau \in \mathcal{B}$ ,

$$\begin{aligned} \tau' \mathbf{I}(\beta) \tau &\geq \delta_1(\min) \sum_i \int [s(y; \mathbf{h}(x_i; \tau)) - a_i]^2 dy \\ &= \delta_1(\min) \sum_i \int \left( \sum_k [h_k(x_i; \tau) - a_i] B_k(y) \right)^2 dy \\ &\geq \delta_2(\min) K^{-1} \sum_i \sum_k [h_k(x_i; \tau) - a_i]^2 \\ &= \delta_2(\min) K^{-1} \sum_i \sum_k \left( \sum_j (\tau_{jk} - a_i) H_j(x_i) \right)^2 \\ &\geq \delta_2(\min) K^{-1} \sum_i \sum_k \left( \sum_j \tau_{jk} H_j(x_i) \right)^2 \\ &\geq \delta(\min) nJ^{-1}K^{-1}|\tau|^2. \end{aligned}$$

Similarly, we conclude from (18) that  $\tau' \mathbf{I}(\beta) \tau \leq (\max) \sum_i \int [s(y; \mathbf{h}(x_i; \tau))]^2 dy$  and hence that there is an  $M > 0$  such that, for  $n \gg 1$ ,  $\beta \in \mathcal{B}$  and  $\tau \in \mathcal{B}$ ,  $\tau' \mathbf{I}(\beta) \tau \leq M(\max) nJ^{-1}K^{-1}|\tau|^2$ . Thus the desired result is valid.  $\square$

Let  $\mathbf{S}(\beta) \in \mathcal{B}$  denote the score at  $\beta$ ; that is, the  $JK$ -dimensional matrix the entry in row  $j$  and column  $k$  of which is

$$\frac{\partial l(\beta)}{\partial \beta_{ij}} = \sum_i H_j(x_i) \left( B_k(Y_i) - \frac{\partial C}{\partial \theta_k}(\mathbf{h}(x_i; \beta)) \right).$$

[In computing  $\partial C(\theta)/\partial \theta_k$ , we let  $\theta$  range over  $\mathbb{R}^K$ .] Then  $E\mathbf{S}(\beta^*) = \mathbf{0}$  and

$$\begin{aligned} E|\mathbf{S}(\beta^*)|^2 &= \sum_i \sum_j \sum_k H_j^2(x_i) \text{var}(B_k(Y_i)) \\ &\leq \sum_i \left( \sum_j H_j^2(x_i) \right) E \left( \sum_k B_k^2(Y_i) \right) \leq n. \end{aligned}$$

Consequently, the following result is valid.

LEMMA 12.  $|\mathbf{S}(\boldsymbol{\beta}^*)|^2 = O_p(n)$ .

The maximum likelihood equation  $\mathbf{S}(\hat{\boldsymbol{\beta}}) = \mathbf{0}$  for  $\hat{\boldsymbol{\beta}}$  can be written as

$$\int_0^1 \frac{d}{dt} S(\boldsymbol{\beta}^* + t(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)) dt = -\mathbf{S}(\boldsymbol{\beta}^*).$$

Thus it can be rewritten as  $\mathbf{D}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*) = \mathbf{S}(\boldsymbol{\beta}^*)$ , where  $\mathbf{D}$  is the  $JK \times JK$  matrix given by  $\mathbf{D} = \int_0^1 \mathbf{I}(\boldsymbol{\beta}^* + t(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)) dt$ .

LEMMA 13.  $|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*| = O_p(JK/\sqrt{n})$ .

PROOF. According to the maximum likelihood equation for  $\hat{\boldsymbol{\beta}}$ ,

$$(20) \quad (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)^t \mathbf{D}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*) = (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)^t \mathbf{S}(\boldsymbol{\beta}^*).$$

It follows from Lemma 12 that

$$(21) \quad (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)^t \mathbf{S}(\boldsymbol{\beta}^*) = O_p(|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*|\sqrt{n}).$$

According to Lemmas 10 and 11, there is an  $M > 0$  such that, except on an event whose probability tends to zero as  $n \rightarrow \infty$ ,

$$(22) \quad (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)^t \mathbf{D}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*) \geq M^{-1} n J^{-1} K^{-1} |\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*|^2.$$

We conclude from (20)–(22) that  $n J^{-1} K^{-1} |\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*|^2 = O_p(|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*|\sqrt{n})$ , which yields the desired result.  $\square$

Let  $\text{VC}(\mathbf{S}(\boldsymbol{\beta}^*))$  denote the variance–covariance matrix of  $\mathbf{S}(\boldsymbol{\beta}^*)$ .

LEMMA 14. *There is an  $M > 0$  such that, for  $n \gg 1$  and  $\boldsymbol{\tau} \in \mathcal{B}$ ,*

$$M^{-1} n J^{-1} K^{-1} |\boldsymbol{\tau}|^2 \leq \boldsymbol{\tau}^t \text{VC}(\mathbf{S}(\boldsymbol{\beta}^*)) \boldsymbol{\tau} \leq M n J^{-1} K^{-1} |\boldsymbol{\tau}|^2.$$

PROOF. Since

$$(23) \quad \boldsymbol{\tau}^t \text{VC}(\mathbf{S}(\boldsymbol{\beta}^*)) \boldsymbol{\tau} = \sum_i \int [s(y; \mathbf{h}(x_i; \boldsymbol{\tau})) - a_i]^2 f(y|x_i) dy, \quad \boldsymbol{\tau} \in \mathcal{B},$$

where  $a_i = \int s(y; \mathbf{h}(x_i; \boldsymbol{\tau})) f(y|x_i) dy$  for  $1 \leq i \leq n$ , the desired result follows from the argument used to prove Lemma 11.  $\square$

Let  $\boldsymbol{\beta} \in \mathcal{B}$ . Then there is a positive semidefinite symmetric  $JK \times JK$  matrix  $[\mathbf{I}(\boldsymbol{\beta})]^-$  having range  $\mathcal{B}$  such that  $\mathbf{I}(\boldsymbol{\beta})[\mathbf{I}(\boldsymbol{\beta})]^- \boldsymbol{\tau} = [\mathbf{I}(\boldsymbol{\beta})]^- \mathbf{I}(\boldsymbol{\beta}) \boldsymbol{\tau} = \boldsymbol{\tau}$  for  $\boldsymbol{\tau} \in \mathcal{B}$ . [Consider the orthogonal diagonalization of  $\mathbf{I}(\boldsymbol{\beta})$ .] The matrix  $[\mathbf{I}(\boldsymbol{\beta})]^-$  is referred to as the generalized inverse of  $\mathbf{I}(\boldsymbol{\beta})$ .

Let  $\hat{\phi} \in \mathcal{B}$  be the approximation to  $\hat{\beta} - \beta^*$  defined by  $\mathbf{I}(\beta^*)\hat{\phi} = \mathbf{S}(\beta^*)$ . Then  $\hat{\phi} = [\mathbf{I}(\beta^*)]^{-1}\mathbf{S}(\beta^*)$ ,  $E\hat{\phi} = \mathbf{0}$  and

$$(24) \quad \text{VC}(\hat{\phi}) = [\mathbf{I}(\beta^*)]^{-1} \text{VC}(\mathbf{S}(\beta^*)) [\mathbf{I}(\beta^*)]^{-1}.$$

The next result follows from (24) and Lemmas 10, 11 and 14. (Consider a symmetric square root of  $[\mathbf{I}(\beta^*)]^{-1}$ .)

LEMMA 15. *There is an  $M > 0$  such that if  $n \gg 1$ , then*

$$(25) \quad M^{-1}n^{-1}JK|\tau|^2 \leq \tau^t [\mathbf{I}(\beta^*)]^{-1} \tau \leq Mn^{-1}JK|\tau|^2, \quad \tau \in \mathcal{B},$$

and

$$(26) \quad M^{-1}n^{-1}JK|\tau|^2 \leq \text{var}(\tau^t \hat{\phi}) \leq Mn^{-1}JK|\tau|^2, \quad \tau \in \mathcal{B}.$$

Given  $y \in \mathcal{Y}$  and  $\beta \in \mathcal{B}$ , let  $\mathbf{G}(y|x; \beta) \in \mathcal{B}$  denote the gradient of  $\log f(y|x; \cdot)$  at  $\beta$ : the  $J \times K$  matrix the entry in row  $j$  and column  $k$  of which is

$$H_j(x) \left( B_k(y) - \frac{\partial C}{\partial \theta_k}(\mathbf{h}(x; \beta)) \right).$$

It follows from Lemma 3 of Stone (1990) and Lemma 10 that

$$(27) \quad \max_k \left| \frac{\partial C}{\partial \theta_k}(\mathbf{h}(x; \beta^*)) \right| = O(K^{-1}).$$

Thus there is an  $M > 0$  such that

$$(28) \quad |\mathbf{G}(y|x; \beta^*)| \leq M, \quad x \in \mathcal{X} \text{ and } y \in \mathcal{Y}.$$

Observe that  $\mathbf{S}(\beta) = \sum_i \mathbf{G}(Y_i|x_i; \beta)$  for  $\beta \in \mathcal{B}$  and hence that

$$(29) \quad \hat{\phi}_{jk} = \sum_i ([\mathbf{I}(\beta^*)]^{-1} \mathbf{G}(Y_i|x_i; \beta^*))_{jk}.$$

The quantities  $j$  and  $k$  in (a) of the next result and the quantity  $j$  in (b) and (d) are allowed to depend on  $n$  in an arbitrary deterministic manner.

LEMMA 16. (a)  $\hat{\phi}_{jk} = O_P(\sqrt{JK/n})$ .

(b)  $(1/K) \sum_k \hat{\phi}_{jk}^2 = O_P(JK/n)$ .

(c)  $|\hat{\phi}|^2 = O_P(J^2 K^2/n)$ .

(d)  $\max_k |\hat{\phi}_{jk}| = O_P(\sqrt{JK(\log K)/n})$ .

(e)  $\max_{j,k} |\hat{\phi}_{jk}| = O_P(\sqrt{JK(\log JK)/n})$ .

PROOF. Now  $\max_{j,k} E\hat{\phi}_{jk}^2 = \max_{j,k} \text{var}(\hat{\phi}_{jk}) = O(JK/n)$  by (26), so (a)–(c) hold. By (25) and (28), there is an  $M > 0$  such that, for  $n \gg 1$ ,

$$\max_{x,y} |\tau^t [\mathbf{I}(\beta^*)]^{-1} \mathbf{G}(y|x; \beta^*)| \leq Mn^{-1}JK|\tau|, \quad \tau \in \mathcal{B},$$

and hence

$$\max_{x,y} \max_{j,k} |([\mathbf{I}(\boldsymbol{\beta}^*)]^{-1} \mathbf{G}(y|x; \boldsymbol{\beta}^*))_{jk}| \leq MJK/n.$$

Parts (d) and (e) now follow from (7), (29) and Bernstein's inequality.  $\square$

LEMMA 17.  $|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^* - \hat{\boldsymbol{\varphi}}|^2 = O_P(n^{-2} J^3 K^2 (\log JK)).$

PROOF. It follows from the maximum likelihood equation that

$$\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^* = \hat{\boldsymbol{\varphi}} - [\mathbf{I}(\boldsymbol{\beta}^*)]^{-1} [\mathbf{D} - \mathbf{I}(\boldsymbol{\beta}^*)](\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*).$$

By (25),

$$|[\mathbf{I}(\boldsymbol{\beta}^*)]^{-1} [\mathbf{D} - \mathbf{I}(\boldsymbol{\beta}^*)](\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)|^2 = O_P(n^{-2} (JK)^2 |[\mathbf{D} - \mathbf{I}(\boldsymbol{\beta}^*)](\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)|^2).$$

The entry in row  $(j, k)$  and column  $(j', k')$  of  $\mathbf{D} - \mathbf{I}(\boldsymbol{\beta}^*)$  can be written as

$$\sum_{j''} \sum_{k''} A_{jkj'k'j''k''} (\hat{\beta}_{j''k''} - \beta_{j''k''}^*),$$

where

$$\begin{aligned} A_{jkj'k'j''k''} &= \sum_i \int_0^1 (1-t) H_j(x_i) H_{j'}(x_i) H_{j''}(x_i) \\ &\quad \times \frac{\partial^3 C}{\partial \theta_k \partial \theta_{k'} \partial \theta_{k''}} (\mathbf{h}(x_i; \boldsymbol{\beta}^* + t(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*))). \end{aligned}$$

Thus the entry in row  $(j, k)$  of  $[\mathbf{D} - \mathbf{I}(\boldsymbol{\beta}^*)](\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)$  is

$$\sum_{j'} \sum_{k'} \sum_{j''} \sum_{k''} A_{jkj'k'j''k''} (\hat{\beta}_{j'k'} - \beta_{j'k'}^*) (\hat{\beta}_{j''k''} - \beta_{j''k''}^*),$$

which is dominated in absolute value by

$$\max_{j,k} (\hat{\beta}_{jk} - \beta_{jk}^*)^2 \sum_{j'} \sum_{k'} \sum_{j''} \sum_{k''} |A_{jkj'k'j''k''}|.$$

There is a positive integer  $J_0$  such that  $A_{jkj'k'j''k''} = 0$  unless  $|j' - j| \leq J_0$  and  $|j'' - j| \leq J_0$ . Thus, by (8) of Stone (1989), there is an  $M_1 > 0$  such that

$$\sum_{j'} \sum_{j''} |A_{jkj'k'j''k''}| \leq M_1 J^{-1} n \sup_{x \in \mathcal{X}} \max_{0 \leq t \leq 1} \left| \frac{\partial^3 C}{\partial \theta_k \partial \theta_{k'} \partial \theta_{k''}} (\mathbf{h}(x; \boldsymbol{\beta}^* + t(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)) \right|.$$

Consequently [see the proof of Lemma 15 of Stone (1990)],

$$\sum_j \sum_k \left( \sum_{j'} \sum_{k'} \sum_{j''} \sum_{k''} |A_{jkj'k'j''k''}| \right)^2 = O_P(nK^{-1})$$

and hence  $|\mathbf{D} - I(\boldsymbol{\beta}^*)|(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)|^2 = O_P(nK^{-1} \max_{j,k} (\hat{\beta}_{jk} - \beta_{jk}^*)^4)$ . Therefore,

$$|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^* - \boldsymbol{\varphi}|^2 = O_P\left(n^{-1}J^2K \max_{j,k} (\hat{\beta}_{jk} - \beta_{jk}^*)^4\right).$$

We now conclude from Lemma 16(e) that

$$\max_{j,k} (\hat{\beta}_{jk} - \beta_{jk}^*)^2 = O_P\left(n^{-1}JK(\log JK) + n^{-1}J^2K \max_{j,k} (\hat{\beta}_{jk} - \beta_{jk}^*)^4\right).$$

Thus  $\max_{j,k} (\hat{\beta}_{jk} - \beta_{jk}^*)^2 = O_P(n^{-1}JK(\log JK))$  by (7) and Lemma 9(b), so

$$|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^* - \hat{\boldsymbol{\varphi}}|^2 = O_P(n^{-2}J^3K^2(\log JK)). \quad \square$$

Theorem 1(a) coincides with Lemma 9(a) and the remaining parts of Theorem 1 follow from (7) and Lemmas 16 and 17.

**3. Functional estimation.** In this section, (8)–(12) will be verified. To help the reader follow the details, we first indicate the proof of (10). It follows from (3) that  $\|\log f^*(\cdot|\cdot)\|_\infty = O(1)$ . Thus, to prove (10), it suffices to show that

$$(30) \quad \max_{x,y} |\log f(y|x; \hat{\boldsymbol{\beta}}) - \log f(y|x; \boldsymbol{\beta}^*)| = O_P(\sqrt{JK(\log JK)/n}).$$

Let  $\nabla C(\boldsymbol{\theta})$  denote the gradient vector of  $C(\cdot)$  at  $\boldsymbol{\theta}$ , whose  $k$ th entry  $\partial C(\boldsymbol{\theta})/\partial \theta_k$  is computed as  $\boldsymbol{\theta}$  ranges over  $\mathbb{R}^K$ . Observe that

$$(31) \quad \begin{aligned} & \log f(y|x; \hat{\boldsymbol{\beta}}) - \log f(y|x; \boldsymbol{\beta}^*) \\ &= [\mathbf{G}(y|x; \boldsymbol{\beta}^*)]^t \hat{\boldsymbol{\varphi}} + [\mathbf{G}(y|x; \boldsymbol{\beta}^*)]^t (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^* - \hat{\boldsymbol{\varphi}}) - R, \end{aligned}$$

where

$$(32) \quad R = C(\mathbf{h}(x; \hat{\boldsymbol{\beta}})) - C(\mathbf{h}(x; \boldsymbol{\beta}^*)) - [\nabla C(\mathbf{h}(x; \boldsymbol{\beta}^*))]^t \mathbf{h}(x; \hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*).$$

According to Lemma 18(c) below,

$$(33) \quad \max_{x,y} |[\mathbf{G}(y|x; \boldsymbol{\beta}^*)]^t \hat{\boldsymbol{\varphi}}| = O_P(\sqrt{JK(\log JK)/n});$$

according to (7), (28) and Lemma 17,

$$(34) \quad \max_{x,y} |[\mathbf{G}(y|x; \boldsymbol{\beta}^*)]^t (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^* - \hat{\boldsymbol{\varphi}})| = o_P(\sqrt{JK/n});$$

and, according to (7) and Lemma 19 below,

$$(35) \quad \max_x |R| = o_P(\sqrt{JK/n}).$$

The desired result (30) follows from (31)–(35).

- LEMMA 18. (a)  $[\mathbf{G}(y|x; \boldsymbol{\beta}^*)]^t \hat{\boldsymbol{\varphi}} = O_P(\sqrt{JK/n})$ .  
 (b)  $\int |[\mathbf{G}(y|x; \boldsymbol{\beta}^*)]^t \hat{\boldsymbol{\varphi}}|^2 dy = O_P(JK/n)$ .  
 (c)  $\max_{x,y} |[\mathbf{G}(y|x; \boldsymbol{\beta}^*)]^t \hat{\boldsymbol{\varphi}}| = O_P(\sqrt{JK(\log JK)/n})$ .

PROOF. Part (a) follows from (7) and Lemma 17. In order to verify (b), choose  $x \in \mathcal{X}$ . There are at most  $J_0$  values of  $j$  such that  $H_j(x) > 0$ . For any such  $j$ ,

$$\begin{aligned} \left( \frac{1}{K} \sum_k \left| \hat{\beta}_{jk} - \beta_{jk}^* - \hat{\phi}_{jk} \right| \right)^2 &\leq \frac{1}{K} \sum_k \left( \hat{\beta}_{jk} - \beta_{jk}^* - \hat{\phi}_{jk} \right)^2 \\ &\leq \frac{1}{K} \sum_j \sum_k \left( \hat{\beta}_{jk} - \beta_{jk}^* - \hat{\phi}_{jk} \right)^2. \end{aligned}$$

Thus, by (7) and Lemma 17,

$$\max_j \frac{1}{K} \sum_k \left| \hat{\beta}_{jk} - \beta_{jk}^* - \hat{\phi}_{jk} \right| = O_P(n^{-1} J^{3/2} K^{1/2} \sqrt{\log JK}) = o_P(\sqrt{J/n}),$$

so (b) holds. Part (c) follows from (27) and Lemma 16(a, b); (d) follows from (27), Lemma 16(b) and (12) of Stone (1986); (e) follows from (27) and Lemma 16(e); and (f) follows from (7) and (e).  $\square$

LEMMA 19. *The following is valid:*

$$\begin{aligned} \max_x \left| C(\mathbf{h}(x; \hat{\beta})) - C(\mathbf{h}(x; \beta^*)) - [\nabla C(\mathbf{h}(x; \beta^*))]^t \mathbf{h}(x; \hat{\beta} - \beta^*) \right| \\ = O_P(JK(\log JK)/n). \end{aligned}$$

PROOF. Observe that

$$\begin{aligned} C(\mathbf{h}(x; \hat{\beta})) - C(\mathbf{h}(x; \beta^*)) \\ = [\nabla C(\mathbf{h}(x; \beta^*))]^t \mathbf{h}(x; \hat{\beta} - \beta^*) \\ + \int_0^1 (1-t) \left( \int [s(y; \mathbf{h}(x; \hat{\beta} - \beta^*)) - a(t)]^2 \right. \\ \left. \times f(y|x; \beta^* + t(\hat{\beta} - \beta^*)) dy \right) dt, \end{aligned}$$

where  $a(t) = \int s(y; \mathbf{h}(x; \hat{\beta} - \beta^*)) f(y|x; \beta^* + t(\hat{\beta} - \beta^*)) dy$  [see (18)]. The desired result now follows from Lemma 10, Theorem 1(e), and (12) of Stone (1986).  $\square$

LEMMA 20. (a)  $\max_{x,y} |\log f(y|x; \hat{\beta}) - \log f(y|x; \beta^*) - [\mathbf{G}(y|x; \beta^*)]^t \hat{\phi}| = o_P(\sqrt{JK/n})$ .

$$\begin{aligned} \max_{x,y} \left| F(y|x; \hat{\beta}) - F(y|x; \beta^*) \right. \\ \left. - \int_{y' \leq y} f(y'|x; \beta^*) [\mathbf{G}(y'|x; \beta^*)]^t \hat{\phi} dy' \right| = o_P(\sqrt{J/n}). \end{aligned}$$

(b)



PROOF. Part (a) follows from (31), (32), (34) and (35). Since

$$F(y|x; \hat{\beta}) - F(y|x; \beta^*) = \int_0^y (e^{\log f(y'|x; \hat{\beta}) - \log f(y'|x; \beta^*)} - 1) f(y'|x; \beta^*) dy',$$

(b) follows from (a) and Lemma 18(c).  $\square$

LEMMA 21.

$$\max_y \left| \int_{y' \leq y} f(y'|x; \beta^*) [\mathbf{G}(y'|x; \beta^*)]^t \hat{\phi} dy' \right| = O_P(\sqrt{J/n}).$$

PROOF. Let  $\tau$  be the member of  $\mathcal{B}$ , the entry in row  $j$  and column  $k$  of which is

$$H_j(x) \frac{\partial C}{\partial \theta_k}(\mathbf{h}(x; \beta^*)).$$

Then  $E(\tau^t \hat{\phi}) = 0$  since  $E\hat{\phi} = \mathbf{0}$  and  $\text{var}(\tau^t \hat{\phi}) = O(J/n)$  by (26) since  $|\tau|^2 = O(K^{-1})$  by (27), so  $\tau^t \hat{\phi} = O_P(\sqrt{J/n})$ . Thus to prove the desired result it suffices to verify that

$$(36) \quad \max_y \left| \sum_j H_j(x) \sum_k \hat{\phi}_{jk} \int_{y' \leq y} f(y'|x; \beta^*) B_k(y') dy' \right| = O_P(\sqrt{J/n}).$$

For any given value of  $y$ , all but a bounded number of terms  $\int_{y' \leq y} f(y'|x; \beta^*) B_k(y') dy'$  are equal to  $\int f(y'|x; \beta^*) B_k(y') dy'$  or to zero. By (4), (7) and Lemma 16(d), the total contribution of the bounded number of exceptional terms is  $O_P(K^{-1} \sqrt{JK}(\log K)/n) = O_P(\sqrt{J/n})$ . Let the  $B$ -splines  $B_1, \dots, B_k$  be ordered according to the right endpoints of their supports. In order to verify (36), it suffices to show that

$$(37) \quad \max_{k'} \left| \sum_{k \leq k'} \hat{\phi}_{jk} \int f(y|x; \beta^*) B_k(y) dy \right| = O_P(\sqrt{J/n}).$$

Let  $\mathcal{K}$  be a subset of consecutive integers in  $\{1, \dots, K\}$  and let  $K'$  denote the cardinality of  $\mathcal{K}$ . Let  $\tau$  denote the  $J \times K$  matrix having entry  $\int f(y|x; \beta^*) B_k(y) dy$  in row  $j$  and column  $k$  for  $k \in \mathcal{K}$  and all other entries equal to zero. Then

$$(38) \quad |\tau|^2 = O_P(K'/K^2).$$

Since

$$\text{var} \left( \sum_{k \in \mathcal{K}} \hat{\phi}_{jk} \int f(y|x; \beta^*) B_k(y) dy \right) = \tau^t \text{VC}(\hat{\phi}) \tau,$$

it follows from (26) and (38) that

$$(39) \quad \text{var} \left( \sum_{k \in \mathcal{K}} \hat{\phi}_{jk} \int f(y|x; \beta^*) B_k(y) dy \right) = O \left( \frac{JK'}{nK} \right).$$

Observe next that

$$(40) \quad \sum_{k \in \mathcal{K}} \hat{\phi}_{jk} \int f(y|x; \beta^*) B_k(y) dy = \sum_i Z_{\mathcal{K}i},$$

where

$$(41) \quad Z_{\mathcal{K}i} = \left( \sum_{k \in \mathcal{K}} ([\mathbf{I}(\beta^*)]^{-1} \mathbf{G}(Y_i|x_i; \beta^*))_{jk} \right) \int f(y|x; \beta^*) B_k(y) dy,$$

$$1 \leq i \leq n,$$

are independent random variables whose sum has mean zero. By (4), (25) and (28),

$$(42) \quad |Z_{\mathcal{K}i}| \leq b = O(n^{-1}J\sqrt{K}).$$

It follows from (39)–(42) and Bernstein's inequality that there is a  $\delta > 0$  such that

$$(43) \quad P \left( \left| \sum_{k \in \mathcal{K}} \hat{\phi}_{jk} \int f(y|x; \beta^*) B_k(y) dy \right| \geq A\sqrt{J/n} (K'/K)^a \right) \leq 2 \left\{ \exp \left[ -\delta A\sqrt{n/(JK)} (K'/K)^a \right] + \exp \left[ -\delta A^2 (K/K')^{1-2a} \right] \right\}$$

for  $A > 0$  and  $0 < a < \frac{1}{2}$ .

Set  $R = \min[r: 2^r \geq K]$ . For  $0 \leq r \leq R$ , let  $\mathcal{M}_r$  denote the collection of all sets of integers of the form  $\{(m-1)2^r + 1, \dots, m2^r\}$ , where  $1 \leq m \leq K/2^r$  and note that the cardinality of  $\mathcal{M}_r$  is at most  $K/2^r$  and that each set in  $\mathcal{M}_r$  has cardinality  $2^r$ . It follows from (7) and (43) that, for any  $\alpha > 0$ ,  $A$  can be chosen sufficiently large so that, for  $0 < a < \frac{1}{4}$  and  $n \gg 1$ ,

$$(44) \quad P \left( \left| \sum_{k \in \mathcal{K}} \hat{\phi}_{jk} \int f(y|x; \beta^*) B_k(y) dy \right| \geq A\sqrt{J/n} (K'/K)^a \right. \\ \left. \text{for some } \mathcal{K} \in \mathcal{M}_0 \cup \dots \cup \mathcal{M}_R \right) \leq \alpha.$$

For  $1 \leq k' \leq K$ ,  $\{1, \dots, k'\}$  can be written as a disjoint union of sets  $\mathcal{K} \in \mathcal{M}_0 \cup \dots \cup \mathcal{M}_R$  such that for  $0 \leq r \leq R$ , there is at most one such  $\mathcal{K} \in \mathcal{M}_r$ . Thus it follows from (44) that (37) holds.  $\square$

Equations (8), (9), (11) and (12) follow from (4), (7) and Lemmas 18, 20 and 21.

**4. Asymptotic normality.** In this section the asymptotic normality of  $(\hat{\tau} - \tau^*)/\text{ASD}(\hat{\tau})$  and  $(\hat{\tau} - \tau^*)/\text{SE}(\hat{\tau})$  will be established, where  $\hat{\tau}$  is  $\hat{f}(y|x)$ ,  $\hat{F}(y|x)$  or  $\hat{Q}(p|x)$ .

The next result follows from (1), (4), (18), (23) and (12) of Stone (1986).

**LEMMA 22.** *There is an  $M > 0$  such that if  $n \gg 1$ , then*

$$|\tau' \text{VC}(\mathbf{S}(\beta^*))\tau - \tau' \mathbf{I}(\beta^*)\tau| \leq MnJ^{-1}K^{-1}\delta_{\mathcal{F}}|\tau|^2, \quad \tau \in \mathcal{B}.$$

According to (24), (25), (28) and Lemma 22, there is an  $M > 0$  such that if  $n \gg 1$ , then

$$(45) \quad \left| \text{var}([\mathbf{G}(y|x; \beta^*)]^t \hat{\phi}) - [\mathbf{G}(y|x; \beta^*)]^t [\mathbf{I}(\beta^*)]^{-1} \mathbf{G}(y|x; \beta^*) \right| \leq Mn^{-1}JK\delta_{\mathcal{J}}, \quad x \in \mathcal{X} \text{ and } y \in \mathcal{Y}.$$

Throughout the remainder of the section, it is assumed that  $J, K \rightarrow \infty$  as  $n \rightarrow \infty$ . Under this assumption,  $\delta_{\mathcal{J}} \rightarrow 0$  as  $n \rightarrow \infty$ . Also, it follows from (27) that there is an  $M > 0$  such that if  $n \gg 1$ , then

$$(46) \quad M^{-1} \leq |\mathbf{G}(y|x; \beta^*)| \leq M, \quad x \in \mathcal{X} \text{ and } y \in \mathcal{Y}.$$

It follows from (26) and (46) that if  $n \gg 1$ , then

$$(47) \quad M^{-1}n^{-1}JK \leq \text{var}([\mathbf{G}(y|x; \beta^*)]^t \hat{\phi}) \leq Mn^{-1}JK, \quad x \in \mathcal{X} \text{ and } y \in \mathcal{Y}.$$

LEMMA 23. *Uniformly for  $x \in \mathcal{X}$  and  $y \in \mathcal{Y}$ , the distribution of*

$$\frac{[\mathbf{G}(y|x; \beta^*)]^t \hat{\phi}}{\text{SD}([\mathbf{G}(y|x; \beta^*)]^t \hat{\phi})}$$

*converges to the standard normal distribution as  $n \rightarrow \infty$ .*

PROOF. Observe that  $[\mathbf{G}(y|x; \beta^*)]^t \hat{\phi} = \sum_i Z_i$ , where

$$Z_i = [\mathbf{G}(y|x; \beta^*)]^t [\mathbf{I}(\beta^*)]^{-1} \mathbf{G}(Y_i|x_i; \beta^*), \quad 1 \leq i \leq n.$$

The random variables  $Z_1, \dots, Z_n$  are independent and their sum has mean zero. Moreover, by (25) and (28), there is an  $M > 0$  such that  $|Z_i| \leq Mn^{-1}JK$  for  $1 \leq i \leq n$ . The desired result now follows from (7), (47) and the central limit theorem [see the corollary on page 201 of Chung (1974)].  $\square$

Sets  $\mathbf{G}^*(y|x) = \mathbf{G}(y|x; \beta^*)$  and  $\hat{\mathbf{G}}(y|x) = \mathbf{G}(y|x; \hat{\beta})$ . Then

$$\text{ASD}(\hat{f}(y|x)) = f^*(y|x) \left\{ [\mathbf{G}^*(y|x)]^t [\mathbf{I}(\beta^*)]^{-1} \mathbf{G}^*(y|x) \right\}^{1/2},$$

$$\text{SE}(\hat{f}(y|x)) = \hat{f}(y|x) \left\{ [\hat{\mathbf{G}}(y|x)]^t [\mathbf{I}(\hat{\beta})]^{-1} \hat{\mathbf{G}}(y|x) \right\}^{1/2},$$

$$\begin{aligned} \text{ASD}(\hat{F}(y|x)) &= \left[ \left( \int_{y' \leq y} f^*(y'|x) \mathbf{G}^*(y'|x) dy' \right)^t [\mathbf{I}(\beta^*)]^{-1} \right. \\ &\quad \left. \times \int_{y' \leq y} f^*(y'|x) \mathbf{G}^*(y'|x) dy' \right]^{1/2}, \end{aligned}$$

$$\text{SE}(\hat{F}(y|x)) = \left[ \left( \int_{y' \leq y} \hat{f}(y'|x) \hat{\mathbf{G}}(y'|x) dy' \right)^t [\mathbf{I}(\hat{\beta})]^{-1} \int_{y' \leq y} \hat{f}(y'|x) \hat{\mathbf{G}}(y'|x) dy' \right]^{1/2},$$

$$\text{ASD}(\hat{Q}(p|x)) = \frac{\text{ASD}(\hat{F}(y|x))}{f^*(y|x)} \Big|_{y=Q^*(p|x)}$$

and

$$\text{SE}(\hat{Q}(p|x)) = \frac{\text{SE}(\hat{F}(y|x))}{\hat{f}(y|x)} \bigg|_{y=\hat{Q}(p|x)}.$$

It follows from (47) and Lemmas 20(a) and 23 that, uniformly for  $x \in \mathcal{X}$  and  $y \in \mathcal{Y}$ , the distribution of

$$\frac{\log f(y|x; \hat{\beta}) - \log f(y|x; \beta^*)}{\text{SD}(\mathbf{G}(y|x; \hat{\beta}))}$$

converges to the standard normal distribution as  $n \rightarrow \infty$ . It now follows easily from (45) and (47) that the distribution of

$$\frac{f(y|x; \hat{\beta}) - f(y|x; \beta^*)}{\text{ASD}(f(y|x; \hat{\beta}))}$$

converges to the standard normal distribution as  $n \rightarrow \infty$  uniformly for  $x \in \mathcal{X}$  and  $y \in \mathcal{Y}$ .

The next result follows from (1), (4), (10), (18) and (12) of Stone (1986).

LEMMA 24. *Uniformly for  $\tau \in \mathcal{B}$ ,*

$$\left| [\mathbf{I}(\hat{\beta}) - \mathbf{I}(\beta^*)] \tau \right|^2 = O_p(n J^{-1} K^{-1} (\log JK)) |\tau|^2.$$

Since

$$\left\{ [\mathbf{I}(\hat{\beta})]^{-1} [\mathbf{I}(\beta^*)]^{-1} \right\} \tau = [\mathbf{I}(\beta^*)]^{-1} [\mathbf{I}(\beta^*) - \mathbf{I}(\hat{\beta})] [\mathbf{I}(\hat{\beta})]^{-1} \tau, \quad \tau \in \mathcal{B},$$

the next result follows from Lemmas 10, 11 and 24.

LEMMA 25. *Uniformly for  $\tau \in \mathcal{B}$ ,*

$$\left| \left\{ [\mathbf{I}(\hat{\beta})]^{-1} - [\mathbf{I}(\beta^*)]^{-1} \right\} \tau \right|^2 = O_p(n^{-3} (JK)^3 (\log JK)) |\tau|^2.$$

LEMMA 26.

$$\max_{x,y} \left| \mathbf{G}(y|x; \hat{\beta}) - \mathbf{G}(y|x; \beta^*) \right|^2 = O_p(n^{-1} JK^{-1} (\log JK)).$$

PROOF. Observe that  $\mathbf{G}(y|x; \hat{\beta}) - \mathbf{G}(y|x; \beta^*)$  is the  $J \times K$  matrix the entry in row  $j$  and column  $k$  of which is  $H_j(x)/B_k(y)[f(y|x; \beta^*) - f(y|x; \hat{\beta})] dy$ . The desired result now follows easily from (10).  $\square$

It follows easily from (7), (10), (45)–(47) and Lemmas 25 and 26 that, uniformly for  $x \in \mathcal{X}$  and  $y \in \mathcal{Y}$ ,

$$\frac{\text{SE}(f(y|x; \hat{\beta}))}{\text{ASD}(f(y|x; \hat{\beta}))} = 1 + o_p(1) \quad \text{as } n \rightarrow \infty.$$

This completes the proof of asymptotic normality for  $\hat{\tau} = f(y|x; \hat{\beta})$ .

Let  $y$  be in the interior of  $\mathcal{Y}$ .

LEMMA 27. *There is an  $M > 0$  such that if  $n \gg 1$ , then*

$$M^{-1}K^{-1} \leq \left| \int_{y' \leq y} f(y'|x; \beta^*) \mathbf{G}(y'|x; \beta^*) dy' \right|^2 \leq MK^{-1}, \quad x \in \mathcal{X}.$$

PROOF. The entry in row  $j$  and column  $k$  of  $\int_{y' \leq y} f(y'|x; \beta^*) \mathbf{G}(y'|x; \beta^*) dy'$  is

$$\begin{aligned} & \int_{y' \leq y} \left( H_j(x) B_k(y') - H_j(x) \frac{\partial C}{\partial \theta_k}(\mathbf{h}(x; \beta^*)) \right) f(y'|x; \beta^*) dy' \\ &= H_j(x) \left( \int_{y' \leq y} B_k(y') f(y'|x; \beta^*) dy' \right. \\ & \quad \left. - \int_{y' \leq y} f(y'|x; \beta^*) dy' \int B_k(y') f(y'|x; \beta^*) dy' \right). \end{aligned}$$

The upper bound in the conclusion of the lemma now obviously holds, so to complete the proof it suffices to show that there is an  $M > 0$  such that if  $n \gg 1$ , then

$$\begin{aligned} M^{-1}K^{-1} &\leq \sum_k \left( \int_{y' \leq y} B_k(y') f(y'|x; \beta^*) dy' \right. \\ & \quad \left. - \int_{y' \leq y} f(y'|x; \beta^*) dy' \int B_k(y') f(y'|x; \beta^*) dy' \right)^2 \end{aligned}$$

for  $x \in \mathcal{X}$ . This lower bound is easily established by noting that if the support of  $B_k(\cdot)$  is to the right of  $y$ , then  $\int_{y' \leq y} B_k(y') f(y'|x; \beta^*) dy' = 0$  for  $x \in \mathcal{X}$ .  $\square$

By (26) and Lemma 27, there is an  $M > 0$  such that if  $n \gg 1$ , then

$$\begin{aligned} (48) \quad M^{-1}n^{-1}J &\leq \text{var} \left( \int_{y' \leq y} f(y'|x; \beta^*) [\mathbf{G}(y'|x; \beta^*)]^t \hat{\phi} dy' \right) \\ &\leq Mn^{-1}J, \quad x \in \mathcal{X}. \end{aligned}$$

LEMMA 28. *Uniformly for  $x \in \mathcal{X}$ , the distribution of*

$$\frac{\int_{y' \leq y} f(y'|x; \beta^*) [\mathbf{G}(y'|x; \beta^*)]^t \hat{\phi} dy'}{\text{SD} \left( \int_{y' \leq y} f(y'|x; \beta^*) [\mathbf{G}(y'|x; \beta^*)]^t \hat{\phi} dy' \right)}$$

*converges to the standard normal distribution as  $n \rightarrow \infty$ .*

PROOF. Observe that  $\int_{y' \leq y} f(y'|x; \beta^*) [\mathbf{G}(y'|x; \beta^*)]^t \hat{\phi} dy' = \sum_i Z_i$ , where

$$Z_i = \int_{y' \leq y} f(y'|x; \beta^*) [\mathbf{G}(y'|x; \beta^*)]^t [\mathbf{I}(\beta^*)]^{-1} \mathbf{G}(Y_i|x_i; \beta^*) dy', \quad 1 \leq i \leq n.$$

By (25), (28) and Lemma 27, there is an  $M > 0$  such that  $|Z_i| \leq Mn^{-1}J\sqrt{K}$  for  $1 \leq i \leq n$ . The remainder of the proof is as in that of Lemma 23, with (48) used instead of (47).  $\square$

The next result follows from (24), (25), (48) and Lemmas 22 and 27.

LEMMA 29. *Uniformly for  $x \in \mathcal{X}$ ,*

$$\frac{\text{ASD}(F(y|x; \hat{\beta}))}{\text{SD}\left(\int_{y' \leq y} f(y'|x; \beta^*) [\mathbf{G}(y'|x; \beta^*)]^t \hat{\phi} dy\right)} \rightarrow 1 \quad \text{as } n \rightarrow \infty.$$

It follows from (48) and Lemmas 20(b), 28 and 29 that, uniformly for  $x \in \mathcal{X}$ , the distribution of

$$\frac{F(y|x; \hat{\beta}) - F(y|x; \beta^*)}{\text{ASD}(F(y|x; \hat{\beta}))}$$

converge to the standard normal distribution as  $n \rightarrow \infty$ . It follows from (4), (7), (10), (25), and (48) and Lemmas 25–27 and 29 that, uniformly for  $x \in \mathcal{X}$ ,

$$\frac{\text{SE}(F(y|x; \hat{\beta}))}{\text{ASD}(F(y|x; \hat{\beta}))} = 1 + o_p(1) \quad \text{as } n \rightarrow \infty.$$

This completes the proof of asymptotic normality for  $\hat{\tau} = F(y|x; \hat{\beta})$ .

Observe that  $F^*(Q^*(p|x)|x) = p$  and  $\hat{F}(\hat{Q}(p|x)|x) = p$ . Thus

$$\begin{aligned} [\hat{F}(y|x) - F^*(y|x)]_{y=Q^*(p|x)} &= -[\hat{F}(\hat{Q}(p|x)|x) - \hat{F}(Q^*(p|x)|x)] \\ &= \int_{Q^*(p|x)}^{\hat{Q}(p|x)} \hat{f}(y|x) dy. \end{aligned}$$

We now conclude from (4) and (10) that, uniformly for  $x \in \mathcal{X}$ ,

$$\hat{Q}(p|x) - Q^*(p|x) = [1 + o_p(1)] \frac{[\hat{F}(y|x) - F^*(y|x)]}{f^*(y|x)} \Big|_{y=Q^*(p|x)}.$$

The argument used to establish asymptotic normality for  $\hat{\tau} = \hat{F}(y|x)$  applies with  $y = Q^*(p|x)$  (even though this value of  $y$  depends on  $n$ ). Thus asymptotic normality for  $\hat{\tau} = \hat{Q}(p|x)$  is valid.

## REFERENCES

- BARRON, A. R. and SHEU, C.-H. (1991). Approximation of density functions by sequences of exponential families. *Ann. Statist.* **19** 1347–1369.  
 BREIMAN, L. (1989). Fitting additive models to regression data. Technical Report 209, Dept. Statistics, Univ. California, Berkeley.

- BREIMAN, L. (1991). The  $\Pi$ -method for estimating multivariate functions (with discussion). *Technometrics* **33** 125–162.
- BREIMAN, L. and PETERS, S. (1988). Comparing automatic smoothers. Technical Report 161, Dept. Statistics, Univ. California, Berkeley.
- CHUNG, K. L. (1974). *A Course in Probability Theory*, 2nd ed. Academic, New York.
- DE BOOR, C. (1978). *A Practical Guide to Splines*. Springer, New York.
- FAN, J. (1990). A remedy to regression estimators and nonparametric minimax efficiency. Technical Report, Dept. Statistics, Univ. North Carolina, Chapel Hill.
- FRIEDMAN, J. H. (1991). Multivariate adaptive regression splines (with discussion). *Ann. Statist.* **19** 1–141.
- FRIEDMAN, J. H. and SILVERMAN, B. W. (1989). Flexible parsimonious smoothing and additive modeling (with discussion). *Technometrics* **31** 3–39.
- HASMINSKI, R. and IBRAGIMOV, I. A. (1990). On density estimation in the view of Kolmogorov's ideas in approximation theory. *Ann. Statist.* **18** 999–1010.
- HOEFFDING, W. (1963). Probability inequalities for sums of bounded random variables. *J. Amer. Statist. Assoc.* **58** 13–30.
- JIN, K. (1990). Empirical smoothing parameter selection in adaptive estimation, Ph. D. dissertation, Dept. Statistics, Univ. California, Berkeley.
- KOO, J.-Y. (1988). Tensor product splines in the estimation of regression, exponential response functions and multivariate densities. Ph. D. dissertation. Dept. Statistics, Univ. California, Berkeley.
- KOOPERBERG, C. and STONE, C. J. (1991). A study of logspline density estimation. *Comput. Statist. Data Anal.* To appear.
- MCCULLAGH, P. and NELDER, J. A. (1983). *Generalized Linear Models*. Chapman and Hall, London.
- PORTNOY, S. (1986). On the central limit in  $\mathbb{R}^p$  when  $p \rightarrow \infty$ . *Probab. Theory Related Fields* **73** 571–583.
- PORTNOY, S. (1988). Asymptotic behavior of likelihood methods for exponential families when the number of parameters tends to infinity. *Ann. Statist.* **16** 356–366.
- SCHUMAKER, L. L. (1981). *Spline Functions: Basic Theory*. Wiley, New York.
- SMITH, P. L. (1982). Curve fitting and modeling with splines using statistical variable selection techniques. NASA Report CR-166034. NASA, Langley Research Center, Hampton, Va.
- 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. (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, Boston.
- STONE, C. J. (1990). Large-sample inference for log-spline models. *Ann. Statist.* **18** 717–741.
- STONE, C. J. and KOO, C.-Y. (1986). Logspline density estimation. *Contemp. Math.* **59** 1–15.

DEPARTMENT OF STATISTICS  
UNIVERSITY OF CALIFORNIA  
BERKELEY, CALIFORNIA 94720