# LOCAL POLYNOMIAL REGRESSION: OPTIMAL KERNELS AND ASYMPTOTIC MINIMAX EFFICIENCY*

JIANQING FAN[1]**, THEO GASSER[2], IRÈNE GIJBELS[3]***,
MICHAEL BROCKMANN[4]* AND JOACHIM ENGEL[4]*

[1]Department of Statistics, University of North Carolina, Chapel Hill, NC 27599-3260, U.S.A.
[2]Biostatistics Department, ISPM, University of Zürich, CH-8006 Zürich, Switzerland
[3]Institut de Statistique, Université Catholique de Louvain, Voie du Roman Pays, 20,
B-1348 Louvain-la-Neuve, Belgium
[4]Institute of Applied Mathematics, University of Heidelberg,
Im Neuenheimer Feld 294, D-6900 Heidelberg, Germany

**Abstract.** We consider local polynomial fitting for estimating a regression function and its derivatives nonparametrically. This method possesses many nice features, among which automatic adaptation to the boundary and adaptation to various designs. A first contribution of this paper is the derivation of an optimal kernel for local polynomial regression, revealing that there is a universal optimal weighting scheme. Fan (1993, *Ann. Statist.*, **21**, 196–216) showed that the univariate local linear regression estimator is the best linear smoother, meaning that it attains the asymptotic linear minimax risk. Moreover, this smoother has high minimax risk. We show that this property also holds for the multivariate local linear regression estimator. In the univariate case we investigate minimax efficiency of local polynomial regression estimators, and find that the asymptotic minimax efficiency for commonly-used orders of fit is 100% among the class of all linear smoothers. Further, we quantify the loss in efficiency when going beyond this class.

*Key words and phrases*:   Curve estimation, local polynomials, minimax efficiency, minimax risk, multivariate curve estimation, nonparametric regression, universal optimal weighting scheme.

## 1.  Introduction

   In parametric regression the form of the regression function is imposed by the model whereas in nonparametric regression this form is determined by available data. Even when parametric modeling is the ultimate goal, nonparametric methods can prove useful for an exploratory analysis and for checking and improving functional models. Details on various nonparametric regression techniques can be found in the monographs of Eubank (1988), Müller (1988), Härdle (1990), Wahba (1990), Green and Silverman (1994) and Wand and Jones (1995), among others.
   For response variables $Y_1, \ldots, Y_n \in \mathbb{R}$, there are explanatory variables $X_1, \ldots, X_n \in \mathbb{R}^d$, $(X_1, Y_1), \ldots, (X_n, Y_n)$ being independent and identically distributed random variables for the random design model with regression function $m$ given by

$$m(x) = \mathrm{E}(Y \mid X = x), \tag{1.1}$$

and conditional variance function $\sigma^2(x) = \mathrm{Var}(Y \mid X = x)$. Non-random $X_1, \ldots, X_n$ are assumed for the following fixed design model:

$$Y_i = m(X_i) + \varepsilon_i \tag{1.2}$$

where the $\varepsilon_i$ are independent and satisfy $\mathrm{E}(\varepsilon_i) = 0$ and $\mathrm{Var}(\varepsilon_i) = \sigma^2(X_i)$. While the $X$ come from a $d$ dimensional probability density with support $\mathrm{Supp}(f_X)$ in model (1.1), a "regular design" is assumed for (1.2), leading formally also to a design density $f_X$.
   Of interest is the estimation of the regression function $m$ and its derivatives. Most of the literature deals with the one-dimensional case $d = 1$. Among linear estimators there are smoothing splines (Wahba (1990)), kernel estimators of the evaluation type (Nadaraya (1964) and Watson (1964)) and of the convolution type (Gasser and Müller (1984)) and local polynomial fitting (for early references see for example Cleveland and Loader (1996)). Spline smoothers are close to kernel estimators (see Silverman (1984)) with an advantage in terms of minimax properties for the latter (see Jennen-Steinmetz and Gasser (1988)). The relative merits of evaluation and convolution weights for kernel estimation have been discussed by Chu and Marron (1991), Fan (1992), Hastie and Loader (1993) and discussions therein. In summary, the evaluation weights lead to an undesirable form of the bias (Gasser and Engel (1990) and Fan (1992)), while convolution weights pay a price in variance for random designs. Fan (1992) evidenced that the local linear fit overcomes these drawbacks. Moreover, a detailed efficiency study by Fan (1993) revealed that local linear regression estimators (for $d = 1$) achieve full asymptotic minimax efficiency among all linear estimators (for the convolution kernel estimators this holds for fixed design only). Further, the asymptotic minimax efficiency remains at 89.4% among all estimators. Another nice feature of the local linear regression estimator is that the bias at the boundary stays automatically of the same order as in the interior, without use of specific boundary kernels. See Lejeune (1985) and Fan and Gijbels (1992).

In this paper we focus on the more general local polynomial fitting in the univariate ($d = 1$) and the multivariate ($d > 1$) case. Recently, Ruppert and Wand (1994) established asymptotic expressions for the conditional bias and variance of local polynomial regression estimators, and showed that these estimators also adapt automatically to the boundary. The main contribution of the present paper is to investigate the minimax efficiency of local polynomial fits in one and more dimensions. In a first stage, while establishing the framework for this study, we derive a universal optimal weighting scheme under local polynomial fitting.

The paper is organized as follows. Section 2 presents the universal optimal weighting scheme in one dimension, which will be used later on in the efficiency study. Section 3 generalizes the results on the optimal weight function to multivariate local linear fitting. Section 4 contains the main contribution of this paper, i.e. the study on minimax efficiency of local polynomial regression estimators in one and more dimensions.

## 2. Universal optimal weighting scheme in one dimension

A local polynomial regression at $x_0$ is computed by minimizing

$$(2.1) \qquad \sum_{i=1}^{n} \left\{ Y_i - \sum_{j=0}^{p} b_j(x_0)(X_i - x_0)^j \right\}^2 K\left(\frac{X_i - x_0}{h_n}\right),$$

where $K(\cdot)$ denotes a weight function and $h_n$ is a smoothing parameter or bandwidth. Denote by $\hat{b}_j(x_0)$ ($j = 0, \ldots, p$) the solution of the least squares problem (2.1). By Taylor's formula $\hat{m}_\nu(x_0) = \nu! \hat{b}_\nu(x_0)$ is an estimator for $m^{(\nu)}(x_0)$. Putting

$$X = \begin{pmatrix} 1 & (X_1 - x_0) & \cdots & (X_1 - x_0)^p \\ \vdots & \vdots & & \vdots \\ 1 & (X_n - x_0) & \cdots & (X_n - x_0)^p \end{pmatrix},$$

$$Y = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}, \qquad \hat{b}(x_0) = \begin{pmatrix} \hat{b}_0(x_0) \\ \vdots \\ \hat{b}_p(x_0) \end{pmatrix}$$

and

$$W = \text{diag}\left\{ K\left(\frac{X_i - x_0}{h_n}\right) \right\},$$

the $n \times n$ diagonal matrix of weights, the solution to the least squares problem (2.1) can be written as

$$(2.2) \quad \hat{b}(x_0) = (X^T W X)^{-1} X^T W Y$$

$$= \begin{pmatrix} S_{n,0}(x_0) & S_{n,1}(x_0) & \cdots & S_{n,p}(x_0) \\ S_{n,1}(x_0) & S_{n,2}(x_0) & \cdots & S_{n,p+1}(x_0) \\ \vdots & \vdots & \ddots & \vdots \\ S_{n,p}(x_0) & S_{n,p+1}(x_0) & \cdots & S_{n,2p}(x_0) \end{pmatrix}^{-1} \begin{pmatrix} T_{n,0}(x_0) \\ T_{n,1}(x_0) \\ \vdots \\ T_{n,p}(x_0) \end{pmatrix}$$

$$= S_n^{-1} T_n,$$

where

$$S_{n,j}(x_0) = \sum_{i=1}^{n} K\left(\frac{X_i - x_0}{h_n}\right)(X_i - x_0)^j, \qquad j = 0, 1, \ldots, 2p,$$

$$T_{n,j}(x_0) = \sum_{i=1}^{n} K\left(\frac{X_i - x_0}{h_n}\right)(X_i - x_0)^j Y_i, \qquad j = 0, 1, \ldots, p.$$

Hence

$$(2.3) \qquad \hat{b}_\nu(x_0) = e_\nu^T \hat{b}(x_0) = \sum_{i=1}^{n} W_\nu^n\left(\frac{X_i - x_0}{h_n}\right) Y_i$$

where $e_\nu = (0, \ldots, 0, 1, 0, \ldots, 0)^T$ with 1 at the $(\nu+1)^{th}$ position and the weight function $W_\nu^n(t) = c_\nu^T S_n^{-1}\{1, th_n, \ldots, (th_n)^p\}^T K(t)$.

Expression (2.3) reveals that the estimator $\hat{b}_\nu(x_0)$ is very much like a conventional kernel estimator except that the kernel $W_\nu^n$ is defined in terms of the design points $X_i$ and the location point $x_0$. The weights in (2.3) satisfy the following discrete orthogonality relation:

$$(2.4) \qquad \sum_{i=1}^{n}(X_i - x_0)^q W_\nu^n\left(\frac{X_i - x_0}{h_n}\right) = \delta_{\nu,q} \qquad 0 \leq \nu, q \leq p,$$

which leads to zero finite bias for polynomials up to order $p$. Such moment conditions and the respective zero bias are satisfied only *asymptotically* for convolution kernel estimators.

Calculation of the local polynomial regression estimators reduces to computing the quantities $S_{n,j}(x_0)$ and $T_{n,j}(x_0)$ followed by calculating (2.2). Fast computation algorithms for nonparametric curve estimators, in particular local polynomials, are provided in Fan and Marron (1994) and Seifert *et al.* (1994).

The following lemma provides a representation of the local polynomial regression estimator in terms of an equivalent kernel estimator. Its proof is simple and is omitted.

LEMMA 2.1. *Assume that the design density $f_X(\cdot)$ is continuous and positive at $x_0$. Then*

$$(2.5) \qquad \hat{b}_\nu(x_0) = \frac{1}{nh_n^{\nu+1} f_X(x_0)} \sum_{i=1}^{n} K_\nu^*\left(\frac{X_i - x_0}{h_n}\right) Y_i\{1 + o_P(1)\}$$

*where*

$$(2.6) \qquad K_\nu^*(t) = e_\nu^T S^{-1}(1, t, \ldots, t^p)^T K(t) = \sum_{l=0}^{p} S^{\nu l} t^l K(t)$$

*with $S = \{\int t^{j+l} K(t)dt\}_{0 \leq j,l \leq p}$ and $S^{-1} = (S^{jl})_{0 \leq j,l \leq p}$.*

We refer to $K_\nu^*$ as the *equivalent kernel.* Since $\int u^q K_\nu^*(u) du = e_\nu^T S^{-1} S e_q$ it follows that the equivalent kernel satisfies the following moment conditions:

$$(2.7) \qquad \int u^q K_\nu^*(u) du = \delta_{\nu,q} \qquad 0 \le \nu, q \le p,$$

which are an asymptotic version of the discrete moment conditions presented in (2.4). In fact the equivalent kernel $K_\nu^*$ is, up to normalizing constants, a kernel of order $(\nu, p+1)$ as defined by Gasser *et al.* (1985). Equivalent kernels have been previously used for analyzing polynomial fitting by Lejeune (1985) and Müller (1987) in a slightly different way. See also Ruppert and Wand (1994).

From the above considerations it is easy to derive the asymptotic expressions for the conditional bias and variance of the estimator $\hat{m}_\nu(x_0)$. As obtained by Ruppert and Wand (1994), these are given by

$$E\{\hat{m}_\nu(x_0) \mid X_1, \ldots, X_n\} - m^{(\nu)}(x_0)$$
$$= \left\{ \int t^{p+1} K_\nu^*(t) dt \right\} \frac{\nu! m^{(p+1)}(x_0)}{(p+1)!} h_n^{p+1-\nu} \{1 + o_P(1)\},$$
$$\mathrm{Var}\{\hat{m}_\nu(x_0) \mid X_1, \ldots, X_n\} = \frac{\nu!^2 \sigma^2(x_0)}{n h_n^{2\nu+1} f_X(x_0)} \int K_\nu^{*2}(t) dt \{1 + o_P(1)\}.$$

*Remark.* For polynomial fitting with a symmetric kernel, it is preferable to choose $p$ of the order $p = \nu + 1, \nu + 3, \ldots$ (see e.g. Ruppert and Wand (1994) and Fan and Gijbels (1995$b$)). Contrary to expectation, the use of a lower order polynomial $p^* = \nu, \nu + 2, \ldots$ with one parameter less does not lead to a smaller asymptotic variance. However, an additional lower order bias arises depending on $m^{(p^*+1)}$ and also on $f_X$ and $f_X'$. A classical example is the Nadaraya-Watson estimator obtained by local constant fitting. Henceforth, we assume $(p - \nu)$ odd. See Fan and Gijbels (1995$b$) for a detailed discussion on the choice of the degree of the polynomial, and for an adaptive procedure for choosing the degree of the polynomial.

Minimization of the asymptotically Mean Squared Error (MSE) leads to the asymptotically optimal local bandwidth

$$h_{\mathrm{opt}}(x_0) = \left[ \frac{(p+1)!^2 (2\nu+1) \int K_\nu^{*2}(t) dt \sigma^2(x_0)}{2n(p+1-\nu)\{\int t^{p+1} K_\nu^*(t) dt\}^2 \{m^{(p+1)}(x_0)\}^2 f_X(x_0)} \right]^{1/(2p+3)}.$$

Minimizing the asymptotically Mean Integrated Squared Error (MISE) results in the asymptotically optimal global bandwidth

$$h_{\mathrm{opt}} = \left[ \frac{(p+1)!^2 (2\nu+1) \int K_\nu^{*2}(t) dt \int \sigma^2(x)/f_X(x) w(x) dx}{2n(p+1-\nu)\{\int t^{p+1} K_\nu^*(t) dt\}^2 \int \{m^{(p+1)}(x)\}^2 w(x) dx} \right]^{1/(2p+3)}$$

for some weight function $w \ge 0$. It is understood that the denominators do not vanish. Effective estimators of $h_{\mathrm{opt}}(x_0)$ and $h_{\mathrm{opt}}$ can be found in Fan and Gijbels (1995$a$) and Ruppert, Sheather and Wand (1995).

Suppose that $\mathrm{Supp}(f_X) = [0,1]$. A point at the left boundary is of the form $x_0 = ch_n$ where $c \geq 0$, and the right boundary points can be defined similarly. The moments are now defined by $s_{j,c} = \int_{-c}^{\infty} u^j K(u) du$ (and $s_{j,c} = \int_{-\infty}^{c} u^j K(u) du$ for the right boundary point $x_0 = 1 - ch_n$), leading to an equivalent boundary kernel (compare with (2.6))

$$(2.8) \qquad K_{\nu,c}^*(t) = e_\nu^T S_c^{-1}(1,t,\ldots,t^p)^T K(t) \quad \text{with} \quad S_c = (s_{j+l,c})_{0 \leq j,l \leq p}.$$

This equivalent kernel differs from the one in (2.6) only in the matrix $S$, and satisfies the boundary moment conditions of Gasser et al. (1985). This reflects clearly that the polynomial method adapts automatically to the boundary, as shown in Fan and Gijbels (1992) for the local linear regression and extended by Ruppert and Wand (1994) to the general case.

Next, the question arises which weight function should be used for different choices of $\nu$ and $p$. The asymptotically MSE and MISE, with optimal choice of the bandwidth, depends on the weight function through

$$(2.9) \qquad T_\nu(K) \equiv \left| \int t^{p+1} K_\nu^*(t) dt \right|^{2\nu+1} \left\{ \int K_\nu^{*2}(t) dt \right\}^{p+1-\nu}.$$

"Optimal" kernels $K_\nu^*$, minimizing the right-hand side of (2.9), have been derived by Gasser et al. (1985) and Granovsky and Müller (1991), postulating a minimal number of sign changes to avoid degeneracy. The following theorem provides a simple solution in the context of polynomial fitting, and reveals a *universal optimal weighting scheme* (i.e. the solution is independent of $p$ and $\nu$) for the interior.

THEOREM 2.1. *The Epanechnikov weight function $K(z) = 3/4(1 - z^2)_+$ is the optimal kernel in the sense that it minimizes $T_\nu(K)$ over all non-negative symmetric functions $K$. It also induces kernels $K_\nu^*$ which are optimal in the sense of Gasser et al. (1985). Similarly, the minimum variance kernel minimizing $T_\nu^*(K) = \int K_\nu^{*2}(t) dt$ is the uniform kernel $1/2 I_{\{|z| \leq 1\}}$.*

The proof of this result is given in Appendix A.1. The optimal kernel $K_\nu^*$, of order $(\nu, p+1)$, is given by (see Gasser et al. (1985))

$$(2.10) \qquad K_{\nu,p}^{\mathrm{opt}}(z) = \sum_{j=0}^{p+1} \lambda_j z^j,$$

where

$$(2.11) \qquad \lambda_j = \begin{cases} 0 \\ \qquad \text{if} \quad j+p+1 \text{ odd} \\ \dfrac{(-1)^{(j+\nu)/2} C_p(p+1-\nu)(p+1+j)!}{j!(j+\nu+1)2^{2p+3}\left(\dfrac{p+1-j}{2}\right)!\left(\dfrac{p+1+j}{2}\right)!} \\ \qquad \text{if} \quad j+p+1 \text{ even,} \end{cases}$$

with $C_p = (p + \nu + 2)!/\{(\frac{p+1+\nu}{2})!(\frac{p+1-\nu}{2})!\}$. Moreover, we have the following explicit formulae for its $(p+1)^{th}$ moment and $L_2$-norm:

(2.12)
$$\left| \int t^{p+1} K_{\nu,p}^{opt}(t)dt \right| = \frac{C_p\{(p+1)!\}^2}{(2p+3)!},$$
$$\int \{K_{\nu,p}^{opt}(t)\}^2 dt = \frac{C_p^2(p+1-\nu)^2}{(2\nu+1)(2p+3)(p+\nu+2)2^{2p+2}}.$$

Can a similar kind of universal optimal weighting scheme be given at the boundary? For example, if one wants to estimate $m(x_0)$, with $x_0 = 1 - ch_n$ which kernel function should be used? In this case the MSE optimality criterion can be defined as (2.9) but with $K_\nu^*$ replaced by $K_{\nu,c}^*$. Denote the resulting criterion by $T_{p,\nu}^c(K)$. When using a local polynomial fit with the Uniform kernel, the weighting scheme still possesses the minimum variance property. The behavior is thus optimal in this sense for the interior as well as for the boundary. This follows immediately from a characterization of minimum variance kernels given by Müller (1991).

The situation is not as simple when taking MSE instead of variance as a criterion for optimality. In the context of kernel estimation there is so far no convincing solution for such optimal boundary kernels available except at the point $c = 0$. This makes it also difficult to judge the quality of the Epanechnikov weight function, when fitting local polynomials at the boundary. For the most left boundary point with $c = 0$, Cheng et al. (1993) shows that the kernel $K(z) = (1 - z)I_{[0,1]}(z)$ is the optimal one, independent of $p$ and $\nu$. (For the most right boundary point, the optimal kernel is $K(z) = (1+z)I_{[-1,0]}(z)$.) For other boundary points the solution is not yet obtained. One difficulty is that it is not clear how to define an appropriate target function. For example, the point $x_n = ch_n$ does not correspond to the same point when two different kernel functions are used, each using the optimal bandwidth.

Table 1 shows the value of $T_{\nu,p}^c(K)^{1/(2p+3)}$, which is the constant factor that depends on $K$ in the MSE expression. Recall that $K(z) = (1 - z)I_{[0,1]}(z)$ is the optimal kernel at the most left boundary point $c = 0$, whereas the Epanechnikov kernel is the optimal one at the least left boundary point $c = 1$—interior point. Note that for the Gaussian kernel, the point $x_n = ch_n$ with $c = 1$ is still a boundary point, whereas for the other kernels, $c = 1$ corresponds to interior points. The triangular kernel used in Table 1 is defined by

$$K(z) = (1 - |z|)_+,$$

where the subscript '+' refers to taking the positive part. Note that this kernel is identical to the optimal kernels at the boundary points 0 and 1. Note also that at the most left boundary point $c = 0$, the Epanechnikov kernel performs very close to the optimal boundary kernel $K(z) = (1 - z)I_{[0,1]}(z)$. This is a justification to use the Epanechnikov kernel even at boundary regions.

It is known that the choice of the kernel function $K$ is not very important for the performance of the resulting estimators, both theoretically and empirically.

Table 1    Values of the target function $T^c_{\nu,p}(K)^{1/(2p+3)}$ for several kernels $K$.

| $\nu$ | $p$ | Gaussian | Triangular | Uniform | Epanechnikov | Biweight |
|---|---|---|---|---|---|---|
| | | | | $c = 0$ | | |
| 0 | 1 | 2.0731 | 1.1817 | 1.2167 | 1.1856 | 1.1830 |
| 0 | 3 | 3.9236 | 2.0996 | 2.1387 | 2.1041 | 2.1012 |
| 1 | 2 | 5.2638 | 3.4588 | 3.6082 | 3.4713 | 3.4630 |
| 1 | 4 | 19.9751 | 11.7790 | 12.2056 | 11.8194 | 11.7939 |
| 2 | 3 | 10.8991 | 7.7564 | 8.0845 | 7.7810 | 7.7664 |
| | | | | $c = 0.5$ | | |
| 0 | 1 | 0.8010 | 0.6095 | 0.6294 | 0.6155 | 0.6055 |
| 0 | 3 | 1.4926 | 0.7771 | 0.6284 | 0.7712 | 0.7957 |
| 1 | 2 | 2.3042 | 0.9103 | 0.8433 | 0.9060 | 0.9161 |
| 1 | 4 | 3.5218 | 0.6282 | 1.4678 | 0.7734 | 1.1863 |
| 2 | 3 | 5.2022 | 1.1657 | 0.8313 | 1.1343 | 1.2651 |
| | | | | $c = 1$ | | |
| 0 | 1 | 0.8834 | 0.5942 | 0.6084 | 0.5908 | 0.5923 |
| 0 | 3 | 1.3077 | 0.7909 | 0.8020 | 0.7873 | 0.7893 |
| 1 | 2 | 0.7293 | 0.8671 | 0.9021 | 0.8647 | 0.8692 |
| 1 | 4 | 2.8039 | 1.4774 | 1.5257 | 1.4724 | 1.4808 |
| 2 | 5 | 1.7734 | 1.2976 | 1.3474 | 1.2927 | 1.3010 |

However, since the Epanechnikov kernel is optimal in minimizing MSE and MISE at interior points and is nearly optimal at the most boundary point, we recommend to use this kernel function. From Table 1, one can see that the Biweight kernel function or the Triweight kernel function (not reported in Table 1), which are respectively of the form

$$K_2(z) = \text{const}\{(1 - z^2)_+\}^2 \quad \text{and} \quad K_3(z) = \text{const}\{(1 - z^2)_+\}^3$$

perform very closely to the Epanechnikov kernel. Therefore, these kernels are also recommendable. The structure of this kind of kernels also enables one to implement fast computing algorithms, which is another reason for this recommendation.

## 3.    Local linear fitting in higher dimensions

We consider only multivariate local linear fitting ($p = 1$, $\nu = 0$) which is the case of most practical interest, since the sparsity of data in higher dimensions becomes more of a problem for higher order polynomials ("curse of dimensionality"). In principle, the methodology generalizes to higher order polynomials using a multi-indices notation. In this section we derive the optimal weighting scheme for multivariate local linear fitting. Minimax properties are discussed in Section 4.

Let $K$ be a function defined in $\mathbb{R}^d$ and let $x = (x_1, \ldots, x_d)$—a point in the $d$-dimensional space—be an interior point of the support of the density of the $d$-variate vector $(X_1, \ldots, X_d)$. The observations are given by $(X_1, Y_1), \ldots, (X_n, Y_n)$,

with $X_j = (X_{j1}, \ldots, X_{jd})$, $j = 1, \ldots, n$. Denote by

$$K_B(u) = \frac{1}{|B|} K(B^{-1}u),$$

where $B$ is a nonsingular matrix, called bandwidth matrix. Further let

$$W = \mathrm{diag}\{K_B(X_1 - x), \ldots, K_B(X_n - x)\}$$

be the diagonal matrix of weights. Then, the local linear regression is to find $\hat{b}_0(x)$ and $\hat{b}_1(x)$ to minimize

$$\sum_{i=1}^n \{Y_i - b_0(x) - b_1^T(x)(X_i - x)\}^2 K_B(X_i - x),$$

and the local linear regression estimator is $\hat{m}(x) = \hat{b}_0(x)$. Asymptotic analysis (see Theorem 2.1 of Ruppert and Wand (1994)) yields the following expression for the conditional mean squared error:

$$\begin{aligned}
\mathrm{E}[\{\hat{m}(x) - m(x)\}^2 \mid X_1, \ldots, X_n] \\
= \left[ \frac{1}{4} \left( \mathrm{tr} \left\{ H(x) BB^T \int K(u) uu^T du \right\} \right)^2 \right. \\
\left. + \frac{1}{n|B|} \int K^2(u) du \frac{\sigma^2(x)}{f_X(x)} \right] \{1 + o_P(1)\},
\end{aligned}$$

where $H(x)$ stands for the Hessian matrix of $m$ at $x$.

Without loss of generality $K$ satisfies

$$\int u_i u_j K(u) du = \delta_{ij} \mu_2(K)$$

where $\mu_2(K)$ is a positive constant. Differentiating the conditional mean squared error with respect to the matrix $BB^T$ leads to the necessary condition for a local minimum

$$\frac{1}{2} \mu_2^2(K) \mathrm{tr}(HBB^T) H - \frac{\nu_0(K)\sigma^2(x)}{2nf_X(x)|B|}(BB^T)^{-1} = 0$$

(see Rao (1973), p. 72), where $\nu_0(K) = \int K^2(u) du$. If the Hessian matrix $H(x)$ is positive or negative definite this equation has a unique solution

(3.1)
$$BB^T = \left( \frac{\nu_0(K)\sigma^2(x)|H^*|^{1/2}}{\mu_2^2(K)ndf_X(x)} \right)^{2/(d+4)} (H^*)^{-1}$$

with

$$H^* = \begin{cases} H & \text{for positive definite } H \\ -H & \text{for negative definite } H \end{cases}$$

which constitutes a minimum. The optimal bandwidth matrix $B$ itself can be chosen as any matrix satisfying the last equation. As can be seen from (3.2) below, the MSE does not depend on the particular choice of $B$. Relation (3.1) leads to insight into the problem of multivariate bandwidth choice: when performing an eigenvalue decomposition of $H^*$, one gets first an optimal rotation of the coordinate system, aligning according to the Hessian matrix. The respective eigenvalues lead to a scaling of the weight function $K$ in direction of the new axes, similar to the analogous problem of bandwidth choice in one dimension. The case of an indefinite $H$ comprises features which are characteristic for dimensions higher than one. It is then possible to choose the bandwidths in different directions appropriately such that the above leading term of the bias vanishes. In such regions with lower order bias, it would be necessary to perform a higher order analysis. In addition, zero eigenvalues of the Hessian matrix lead to problems similar to the linear case in one dimension. A detailed discussion is beyond the scope of this paper.

Substituting the optimal choice for the bandwidth matrix into the asymptotic expression for the MSE, we get

$$(3.2) \quad \text{AMSE} = \frac{d+4}{4} d^{4/(d+4)} \{\nu_0^2(K)\mu_2^d(K)\}^{2/(d+4)} \{\frac{\sigma^4(x)}{n^2 f_X^2(x)}|H^+(x)|\}^{2/(d+4)}.$$

The above formula leads to the formalization of optimal weight functions: find $K$ such that

$$(3.3) \quad \nu_0^2(K)\mu_2^d(K) = \left\{\int K^2(u)du\right\}^2 \int u_1^2 K(u)du \cdots \int u_d^2 K(u)du$$

is minimized subject to

$$\int K(u)du = 1, \quad \int uK(u)du = 0, \quad K \geq 0, \quad \int u_i u_j K(u)du = \delta_{ij}\mu_2(K).$$

The solution given in the next theorem is derived in Appendix A.2.

THEOREM 3.1. *The optimal weight function is the spherical Epanechnikov kernel*

$$K_0(u) = \frac{d(d+2)}{2S_d}(1 - u_1^2 - \cdots - u_d^2)_+,$$

*where $S_d = 2\pi^{d/2}/\Gamma(d/2)$ denotes the area of the surface of the d-dimensional unit ball.*

Thus the theorem also provides an answer to the open problem about the optimal support of the weight function in higher dimensions. The solution is contrary to a suggestion made by Epanechnikov (1969) to use the product of the one-dimensional kernel of the form $0.75(1 - u_i^2)_+$. With the optimal weight function, we can easily obtain that

$$\mu_2(K_0) = \frac{1}{d+4} \quad \text{and} \quad \nu_0(K_0) = \frac{2d(d+2)}{(d+4)S_d}.$$

These two moments are useful in bandwidth selection and for risk calculation.

## 4.   Minimax efficiency

In this section minimax efficiency of univariate and multivariate local poly-
nomial estimators is studied, generalizing results obtained by Fan (1993) for the
special case $d = 1$, $p = 1$ and $\nu = 0$. We only focus on estimation of $m^{(\nu)}(x_0)$
for interior points $x_0$. It is shown that the asymptotic minimax efficiency for
commonly-used orders of fit is 100% among all linear estimators and that a rela-
tively small loss in efficiency occurs when allowing nonlinear estimators as well. For
an illuminating account on recent developments of minimax theory, see Donoho
and Liu (1991), Donoho (1994) and Donoho *et al.* (1995) and the references therein,
where attention is focused on density estimation and white noise models.

### 4.1   *Estimating $m^{(\nu)}$ in one dimension*

Most commonly-used function estimators are linear, i.e. admit a represen-
tation of the form $\sum_{i=1}^{n} w_i(x, X_1, \ldots, X_n) Y_i$. In this section, previous minimax
results for univariate local linear fitting (see Fan (1993)) are extended to higher
order polynomials and to derivative estimation.

Without loss of generality, our goal is to estimate the functional $S_\nu(m) =
m^{(\nu)}(0)$, where 0 represents an arbitrary interior point. Consider

$$(4.1) \qquad \mathcal{C}_{p+1} = \left\{ m : \left| m(z) - \sum_{j=0}^{p} \frac{m^{(j)}(0)}{j!} z^j \right| \le C \frac{|z|^{p+1}}{(p+1)!} \right\},$$

which includes all regression functions whose $(p+1)^{th}$ derivative is bounded by $C$.

Further basic conditions are needed in order to obtain the minimax results:

CONDITION A.
(a)  $\sigma(\cdot)$ is continuous at the point 0,
(b)  $f_X(\cdot)$ is continuous at the point 0 with $f_X(0) > 0$,
(c)  $p - \nu$ is odd.

The *linear minimax risk* is defined as

$$(4.2) \qquad R_{\nu,L}(n, \mathcal{C}_{p+1}) = \inf_{S_\nu \text{ linear}} \sup_{m \in \mathcal{C}_{p+1}} E[\{\hat{S}_\nu - m^{(\nu)}(0)\}^2 \mid X_1, \ldots, X_n],$$

and the *minimax risk* is

$$(4.3) \qquad R_\nu(n, \mathcal{C}_{p+1}) = \inf_{\hat{T}_\nu} \sup_{m \in \mathcal{C}_{p+1}} E[\{\hat{T}_\nu - m^{(\nu)}(0)\}^2 \mid X_1, \ldots, X_n].$$

The latter involves all possible estimators including nonparametric estimators with
data-driven choice of smoothing parameters.

As has been shown by Donoho and Liu (1991) and Fan (1993) the modulus
of continuity

$$(4.4) \qquad \omega_\nu(\varepsilon) = \sup\{|m_1^{(\nu)}(0) - m_0^{(\nu)}(0)| : m_0, m_1 \in \mathcal{C}_{p+1}, \|m_1 - m_0\| = \varepsilon\}$$

is a key tool for deriving lower bounds for the minimax risk.

Denote the optimal kernel of order $(\nu, p+1)$ by $K_{\nu,p}^{\mathrm{opt}}(x)$ (see (2.10)). Recall that, as shown in Theorem 2.1, $K_{\nu,p}^{\mathrm{opt}}(x) = K_{0\nu}^{*}(x)$—the equivalent kernel of the Epanechnikov weight function. Let

$$(4.5) \qquad r = 2(p+1-\nu)/(2p+3), \qquad \text{and} \qquad s = (2\nu+1)/(2p+3),$$

and denote

$$(4.6) \qquad \tau_{\nu,p} = r^{-r} s^{-s} \left\{ \frac{C}{(p+1)!} \right\}^{2s} \left\{ \frac{\sigma^2(0)}{n f_X(0)} \right\}^r.$$

THEOREM 4.1.    *The linear minimax risk is bounded from above and below by*

$$(4.7) \qquad B_{\nu,p}\{1 + o_P(1)\} \geq R_{\nu,L}(n, C_{p+1}) \geq b_{\nu,p}\{1 + o_P(1)\},$$

*where*

$$B_{\nu,p} = \tau_{\nu,p} \left( \int |t^{p+1} K_{\nu,p}^{\mathrm{opt}}(t)| dt \right)^{2s} \left( \int \{K_{\nu,p}^{\mathrm{opt}}(t)\}^2 dt \right)^r$$

*and*

$$b_{\nu,p} = \tau_{\nu,p} \left( \left| \int t^{p+1} K_{\nu,p}^{\mathrm{opt}}(t) dt \right| \right)^{2s} \left( \int \{K_{\nu,p}^{\mathrm{opt}}(t)\}^2 dt \right)^r.$$

Using expressions in (2.12), it can explicitly be calculated that

$$(4.8) \qquad b_{\nu,p} = \left( \frac{2p+3}{2\nu+1} \right) \left\{ \frac{(p+\nu+2)!}{\left( \dfrac{p+1+\nu}{2} \right)! \left( \dfrac{p+1-\nu}{2} \right)!} \right\}^2$$

$$\times \left\{ \frac{r}{(p+\nu+2)4^{p+2}} \right\}^r \left\{ \frac{(p+1)!C}{(2p+3)!} \right\}^{2s} \left\{ \frac{\sigma^2(0)}{n f_X(0)} \right\}^r.$$

How far apart are the lower bound and upper bound in the linear minimax risk? To get an impression of this we present in Table 2 the ratio of the square roots of $b_{\nu,p}$ and $B_{\nu,p}$ (the MSE lower bound and upper bound) which is given by

$$(4.9) \qquad \theta_{\nu,p} = \left( \frac{|\int t^{p+1} K_{\nu,p}^{\mathrm{opt}}(t) dt|}{\int |t^{p+1} K_{\nu,p}^{\mathrm{opt}}(t)| dt} \right)^s.$$

It is easy to see that this ratio is equal to one for the commonly-used cases: $p = 1$, $\nu = 0$ and $p = 2$, $\nu = 1$. We conjecture that the sharp risk is obtained by replacing $|K_{\nu,p}^{\mathrm{opt}}|$ in the definition of $B_{\nu,p}$ by $|K_{\nu,p}^{\mathrm{opt}*}|$ which minimizes

$$T_{\nu,p}^{*}(K) = \left\{ \int |t^{p+1} K_{\nu}^{*}(t)| dt \right\}^{2\nu+1} \left\{ \int K_{\nu}^{*2}(t) dt \right\}^{p+1-\nu}$$

Table 2. Ratio $\theta_{\nu,p}$ of the minimax lower and upper bounds.

| $p$ | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| $p - \nu$ | | | | | |
| 1 | 1 | 1 | .9902 | .9124 | .7620 |
| 3 | | | .9478 | .8126 | .6611 |
| 5 | | | | | .8983 |

(compare with (2.9)). However, the solution to this problem does not usually admit an explicit formula and has a so complicated shape that it is only of theoretical interest. See Sacks and Ylvisaker (1981).

Let $\hat{m}_\nu^*(0) = \nu!\hat{b}_\nu(0)$ be the estimator resulting from a local polynomial fit of order $p$ with the Epanechnikov kernel $K_0$ and with the bandwidth

$$(4.10) \qquad h_n = \left[ \frac{(2\nu + 1)(p + 1)!^2 \int \{K_{\nu,p}^{\mathrm{opt}}(t)\}^2 dt \sigma^2(0)}{2(p + 1 - \nu)\{\int |t^{p+1} K_{\nu,p}^{\mathrm{opt}}(t)| dt\}^2 f_X(0) C^2 n} \right]^{1/(2p+3)}$$

It is worthwhile noting that these choices of the kernel function and the bandwidth reduce to the choices provided in Fan (1993) for the special case that $p = 1$ and $\nu = 0$. With the above choices of the kernel and the bandwidth the local polynomial estimator has a high linear minimax efficiency, as is established in the next theorem.

THEOREM 4.2. *The local polynomial estimator $\hat{m}_\nu^*(0)$ has high linear minimax efficiency for estimating $m^{(\nu)}(0)$ in the sense that*

$$(4.11) \qquad \frac{R_{\nu,L}(n, \mathcal{C}_{p+1})}{\sup_{m \in \mathcal{C}_{p+1}} \mathrm{E}[\{\hat{m}_\nu^*(0) - m^{(\nu)}(0)\}^2 \mid X_1, \ldots, X_n]} \geq \theta_{\nu,p}^2 \{1 + o_P(1)\}.$$

Thus, for commonly-used local polynomial fitting with $p = 1$, $\nu = 0$ or $p = 2$, $\nu = 1$, the estimator $\hat{m}_\nu^*(0)$ is the best linear estimator. For other values of $p$ and $\nu$ (unless the value is too large to be practical), the local polynomial fit is nearly the best, as evidenced by Table 2. In fact, if our conjecture after Theorem 4.1 is true, the efficiency of $\hat{m}_\nu^*(0)$ is far higher than $\theta_{\nu,p}$.

The behavior of the minimax risk $R_\nu(n, \mathcal{C}_{p+1})$ over the class of regression functions $\mathcal{C}_{p+1}$ is established in statement (4.12) below. The second statement in the theorem shows that the local polynomial regression estimator $\hat{m}_r^*(0)$ comes asymptotically fairly close to the minimax risk. This estimator is efficient in the rate, and nearly efficient in the constant factor.

THEOREM 4.3. *The minimax risk (4.3) is asymptotically bounded by*

$$(4.12) \qquad B_{\nu,p} \geq R_\nu(n, \mathcal{C}_{p+1}) \geq (0.894)^2 b_{\nu,p},$$

*with $b_{\nu,p}$ and $B_{\nu,p}$ as in (4.7). Moreover, the local polynomial regression estimator $\hat{m}_{\nu}^*(0)$ has an asymptotic efficiency of at least $89.4\,\theta_{\nu,p}\%$ among all estimators:*

$$(4.13) \qquad \frac{R_{\nu}(n, \mathcal{C}_{p+1})}{\sup_{m \in \mathcal{C}_{p+1}} \mathrm{E}[\{\hat{m}_{\nu}^*(x_0) - m^{(\nu)}(x_0)\}^2 \mid X_1, \ldots, X_n]}$$
$$\geq (0.894)^2 \theta_{\nu,p}^2 + o_P(1),$$

*with $\theta_{\nu,p}$ as in (4.9).*

The minimax theory provided in this section is an additional justification of the intuitively appealing local polynomial approximation method. Other justifications such as graphical representation and finite sample simulations can be found in Fan (1992), Hastie and Loader (1993) and discussions therein.

## 4.2  *Local linear fitting in higher dimensions*

In this section we will show that the multivariate local linear regression estimator, as defined in Section 3, possesses the same minimax properties as the univariate local linear regression estimator: with an appropriate choice of the bandwidth matrix and the kernel function, the multivariate local linear regression estimator achieves asymptotically the linear minimax risk, and comes asymptotically fairly close to the minimax risk. To accomplish this, we assume that the regression function is in the class.

$$\mathcal{C}_2 = \left\{ m : |m(z) - (m'(x))(z - x)^T| \leq \frac{1}{2}(z - x)C(z - x)^T \right\},$$

where $C$ is a positive definite $(d \times d)$-matrix. Intuitively, this class includes regression functions whose Hessian matrix is bounded by $C$. We use $R_{0,L}(n, \mathcal{C}_2)$ and $R_0(n, \mathcal{C}_2)$ to denote respectively the minimax linear risk and the minimax risk, defined similarly to (4.2) and (4.3). Without loss of generality, we assume that $x = 0$.

THEOREM 4.4.  *Suppose that Condition A (a) and (b) hold. Then, the local linear fit with spherically symmetric Epanechnikov weight function is a best linear estimator and has minimax efficiency at least 89.4% in the sense similar to (4.13) in Theorem 4.3. Moreover, the linear minimax risk is given by*

$$(4.14) \qquad\qquad R_{0,L}(n, \mathcal{C}_2) = r_d\{1 + o_P(1)\},$$

*where*

$$(4.15) \quad r_d = \frac{d}{4}\left(\frac{2}{S_d}\right)^{4/(d+4)}(d+2)^{4/(d+4)}(d+4)^{-d/(d+4)}\left\{\frac{\sigma^4(0)}{n^2 f_X^2(0)}|C|\right\}^{2/(d+4)}.$$

*Further, the minimax risk $R_0(n, \mathcal{C}_2)$ is asymptotically bounded by*

$$(4.16) \qquad\qquad r_d \geq R_0(n, \mathcal{C}_2) \geq (0.894)^2 r_d.$$

In the univariate case (i.e. $d = 1$), the quantity $r_d$ reduces to the factor in Fan (1993).

## Acknowledgement

## Appendix—Proofs of the results

### A.1  Proof of Theorem 2.1

In order to use the same normalization constant as Gasser *et al.* (1985), we define the equivalent kernel by (compare with (2.6))

$$W_\nu^*(x) = (-1)^\nu \nu! e_\nu^T S^{-1} (1, x, \ldots, x^p)^T K_0(x).$$

The following properties hold for $W_\nu^*$:

1.  $W_\nu^*$ is a $(\nu, p+1)$ kernel in the sense of Gasser *et al.* (1985).
2.  Obviously, $W_\nu^*(-1) = W_\nu^*(1) = 0$.
3.  Because $\nu + p$ is odd, the last element of the $\nu$-th row of $S^{-1}$ is zero and therefore by (2.6), $W_\nu^*$ is a polynomial of degree $p+1$ on $[-1, 1]$.

Theorem C in Granovsky and Müller (1991) states that these properties characterize the optimal kernel. □

### A.2  Proof of Theorem 3.1

Let

(A.1)  $$\varphi(K) = \left\{ \int K^2(u) du \right\}^2 \int u_1^2 K(u) du \cdots \int u_d^2 K(u) du.$$

We had assumed without loss of generality that

(A.2)  $$\int u_1^2 K(u) du = \cdots = \int u_d^2 K(u) du,$$

which leads to

(A.3)  $$\varphi(K) = \left\{ \int K^2(u) du \right\}^2 \left\{ \frac{1}{d} \int (u_1^2 + \cdots + u_d^2) K(u) du \right\}^d.$$

Property (A.3) allows us to find the optimal $K$ through the following minimization problem:

Minimize $\int K^2(u) du$ subject to

$$\int K(u) du = 1, \quad K \geq 0, \quad \int u K(u) du = 0,$$

(A.4)  $$\int u_i u_j K(u) du = 0 \quad \text{when } i \neq j$$

$$\int (u_1^2 + \cdots + u_d^2) K(u) du = \int (u_1^2 + \cdots + u_d^2) K_0(u) du.$$

Now for any nonnegative kernel $K \geq 0$, let $\delta - K - K_0$.

$$\int \|u\|^2 \delta(u) du = 0, \qquad \int \delta(u) du = 0.$$

Hence, we find that

$$\int \delta(u) K_0(u) du = \int_{\{\|u\|^2 \leq 1\}} \delta(u)(1 - \|u\|^2) du$$

$$= -\int_{\{\|u\|^2 > 1\}} \delta(u)(1 - \|u\|^2) du$$

$$= \int_{\{\|u\|^2 > 1\}} K(u)(\|u\|^2 - 1) du \geq 0,$$

and therefore,

$$\int K^2(u) du = \int K_0^2(u) du + 2 \int K_0(u) \delta(u) du + \int \delta^2(u) du \geq \int K_0^2(u) du,$$

which proves that $K_0$ is the optimal kernel. $\square$

### A.3    Proofs of Theorems 4.1–4.3

### A.3.1    Upper bound

Let $\hat{m}^{(\nu)}(0)$ be the estimator resulting from a local polynomial fit of order $p$ with the Epanechnikov weight function $K_0$ and bandwidth $h_n$. Let $K_{0\nu}^* = K_{\nu,p}^{\mathrm{opt}}$ be its equivalent kernel, as defined in (2.6). Then,

$$\sup_{m \in \mathcal{C}_{p+1}} \mathrm{E}[\{\hat{m}_\nu(0) - m^{(\nu)}(0)\}^2 \mid X_1, \ldots, X_n]$$

$$\leq \left( \int |t^{p+1} K_{\nu,p}^{\mathrm{opt}}(t)| dt \frac{C}{(p+1)!} \right)^2 h_n^{2(p+1-\nu)}$$

$$+ \int \{K_{\nu,p}^{\mathrm{opt}}(t)\}^2 dt \frac{\sigma^2(0)}{n f_X(0)} h_n^{-(2\nu+1)}$$

$$\equiv A_1 h_n^{2(p+1-\nu)} + A_2 h_n^{-(2\nu+1)}.$$

The bandwidth that minimizes the above quantity is given by

$$h_n = \left\{ \frac{(2\nu+1)A_2}{2(p+1-\nu)A_1} \right\}^{1/(2p+3)},$$

which is exactly the bandwidth defined in (4.10). With this choice of $h_n$ we obtain via simple algebra that

(A.5)    $$\sup_{m \in \mathcal{C}_{p+1}} \mathrm{E}[\{\hat{m}_\nu(0) - m^{(\nu)}(0)\}^2 \mid X_1, \ldots, X_n]$$

$$\leq A_1^s A_2^r \left[ \left\{ \frac{2\nu+1}{2(p+1-\nu)} \right\}^r + \left\{ \frac{2\nu+1}{2(p+1-\nu)} \right\}^{-s} \right]$$

$$= A_1^s A_2^r (2p+3)(2\nu+1)^{-s} [2(p+1-\nu)]^{-r}$$

$$= B_{\nu,p},$$

with $B_{\nu,p}$ as in Theorem 4.1. This establishes the upper bound for the linear minimax risk.

A.3.2    *Lower bound*

Let us evaluate the modulus of continuity defined by (4.4). Take an $f \in \mathcal{C}$ and let

(A.6)                    $m_1(x) = \delta^{p+1} f(x/\delta), \qquad m_0(x) = -m_1(x),$

where $\delta$ is a positive constant to be determined later on. Clearly, $m_0, m_1 \in \mathcal{C}$.

Now, selecting $\delta$ such that

$$\|m_1 - m_0\|^2 = 4\delta^{2p+3}\|f\|^2 = \varepsilon^2,$$

which is equivalent to taking

$$\delta = \left(\frac{\varepsilon^2}{4\|f\|^2}\right)^{1/(2p+3)},$$

we obtain that

(A.7)     $\omega_\nu(\varepsilon) \geq |m_1^{(\nu)}(0) - m_0^{(\nu)}(0)| = 2|f^{(\nu)}(0)|\left(\frac{\varepsilon^2}{4\|f\|^2}\right)^{(p+1-\nu)/(2p+3)}.$

Recall the optimal kernel of order $(\nu, p+1)$ given in (2.10).

Defining

(A.8)                    $g(x) = \begin{cases} K_{\nu,p}^{\mathrm{opt}}(x) & \text{if } |x| \leq 1 \\ 0 & \text{otherwise,} \end{cases}$

we need the following lemma.

LEMMA A.1.    *The function $g(\cdot)$ satisfies*

(A.9)            $\left| g(x) - \sum_{j=0}^{p} g^{(j)}(0)\frac{x^j}{j!} \right| \leq \left| \frac{g^{(p+1)}(0)}{(p+1)!}x^{p+1} \right|.$

PROOF.    It is obvious that (A.9) holds whenever $|x| \leq 1$. When $|x| > 1$, $g(x) = 0$ and (A.9) becomes

(A.10)                    $\left| \sum_{j=0}^{p} \lambda_j x^j \right| \leq |\lambda_{p+1} x^{p+1}|.$

Since the polynomial above is either even or odd (see (2.11)), we only need to check (A.10) for $x > 1$. Note that the polynomial (2.10) has $k$ roots on $[-1, 1]$ (see Lemmas 2 and 3 of Gasser *et al.* (1985)). Hence $K_{\nu,p}^{\mathrm{opt}}(x)$ can not change its sign when $x > 1$. Thus

$$\mathrm{sgn}(\lambda_{p+1})K_{\nu,p}^{\mathrm{opt}}(x) > 0,$$

namely

$$|\lambda_{p+1}||x^{p+1} > -\operatorname{sgn}(\lambda_{p+1})\sum_{j-0}^{p}\lambda_j x^j \quad\text{when}\quad x > 1.$$

Using a similar argument we get,

$$-\operatorname{sgn}(\lambda_{p+1})\sum_{j=0}^{p}\lambda_j x^j > 0 \quad\text{when}\quad x > 1.$$

Combining the above two statements, we obtain (A.10). $\square$

We are now ready to establish a lower bound for the linear minimax risk. Take,

$$(A.11) \qquad f(x) = g(ax) \quad\text{with}\quad a = \left(\frac{C}{(p+1)!\lambda_{p+1}}\right)^{1/(p+1)},$$

where $g$ is defined by (A.8). Then, by Lemma A.1, $f \in \mathcal{C}$. It follows from (A.11) that

$$(A.12) \qquad \|f\|^2 = \frac{\|g\|^2}{a} = \frac{\|K_{\nu,p}^{\mathrm{opt}}\|^2}{a}, \qquad f^{(\nu)}(0) = a^{\nu}\nu!\lambda_\nu.$$

Substituting (A.12) into (A.7), we obtain

$$(A.13) \qquad \omega_\nu(\varepsilon) \ge 2\nu!|\lambda_\nu|\left(\frac{C}{(p+1)!|\lambda_{p+1}|}\right)^{s}\|K_{\nu,p}^{\mathrm{opt}}\|^{-r}\left(\frac{\varepsilon^2}{4}\right)^{r/2}.$$

Applying Theorem 6 of Fan (1993), we find that the minimax bound of the best linear procedure is

$$(A.14) \quad R_{\nu,L}(n,\mathcal{C}_{p+1})$$
$$\ge r^r s^s\left[\nu!\lambda_\nu\left(\frac{C}{(p+1)!|\lambda_{p+1}|}\right)^{s}\|K_{\nu,p}^{\mathrm{opt}}\|^{-r}\left(\frac{\upsilon^2(0)}{nf_X(0)}\right)^{r/2}\right]^2$$
$$= b_{\nu,p},$$

using (2.10) and (2.12), with $b_{\nu,p}$ as defined in (4.8). This leads to a lower bound for the linear minimax risk.

PROOF OF THEOREM 4.1.   Statement (4.7) follows immediately from the upper and lower bound given in respectively (A.5) and (A.14). $\square$

PROOF OF THEOREM 4.2.   The maximum risk of $m^{(\nu)}(0)$ is given by (A.5). The result follows from Theorem 4.1. $\square$

PROOF OF THEOREM 4.3.   Theorem 4.3 is an immediate consequence of Theorem 4.1, (A.13) and an application of Theorem 6 of Fan (1993). $\square$

## A.4   *Proof of Theorem 4.1*

It can easily been shown that the maximum risk of the local linear regression smoother is bounded by the left-hand side of (4.15). See (3.2) for a similar expression. Therefore, the minimax risks are bounded by the left-hand side of (4.15). To establish the lower bound, we apply Theorem 6 of Fan (1993). To this end, let the modulus of continuity be

$$\omega(\varepsilon) = \sup\{|m_1(0) - m_0(0)| : m_0, m_1 \in \mathcal{C}_2, \|m_1 - m_0\| = \varepsilon\},$$

where $\| \cdot \|$ is the $L_2$-norm. Without loss of generality, we assume that $C$ is a diagonal matrix given by $C = \text{diag}\{\lambda_1, \ldots, \lambda_d\}$. Take

$$m_0(x) = \frac{\delta^2}{2}\left(1 - \delta^{-2}\sum_{j=1}^{d}\lambda_j x_j^2\right)_+^2.$$

Then, $m \in \mathcal{C}_2$. It can easily be computed that

$$\|m_0\|^2 = \frac{2\delta^{4+d}S_d}{|C|^{1/2}d(d+2)(d+4)}.$$

Setting the above expression to $\varepsilon^2/4$ leads to

$$\delta = \left\{\frac{|C|^{1/2}d(d+2)(d+4)\varepsilon^2}{8S_d}\right\}^{1/(d+4)}.$$

Now, taking the pair $m_1 = -m_0$ and $m_0$, we have that $\|m_1(0) - m_0(0)\| = \varepsilon$. Therefore,

$$\omega(\varepsilon) \geq |m_1(0) - m_0(0)| = \delta^2 = \left\{\frac{|C|^{1/2}d(d+2)(d+4)}{8S_d}\right\}^{2/(d+4)}\varepsilon^{4/(d+4)}.$$

Now applying Theorem 6 of Fan (1993) with $p = 4/(d+4)$ and $q = 1 - p$, we obtain

$$(A.15) \qquad R_{0,L}(n, \mathcal{C}_2) \geq \frac{p^p q^q}{4}\omega^2\left(2\sqrt{\frac{\sigma^2(0)}{nf_X(0)}}\right)(1 + o(1))$$

$$= \frac{d}{4}\left(\frac{2}{S_d}\right)^{4/(d+4)}(d+2)^{4/(d+4)}(d+4)^{-d/(d+4)}$$

$$\times \left\{\frac{\sigma^4(0)}{n^2 f_X^2(0)}|C|\right\}^{2/(d+4)}\{1 + o(1)\}.$$

The fact that the lower and upper bound are the same leads to the conclusions on linear minimax risk. As for the nonlinear minimax risk $R_{0,L}(n, \mathcal{C}_2)$, the conclusion follows directly from Theorem 6 of Fan (1993) together with (A.15). $\square$

# References

Cheng, M. Y., Fan, J. and Marron, J. S. (1993). Minimax efficiency of local polynomial fit estimators at boundaries, Institute of Statistics Mimeo Series #2098, University of North Carolina, Chapel Hill.

Chu, C. K. and Marron, J. S. (1991). Choosing a kernel regression estimator, *Statist. Sci.*, **6**, 404–433.

Cleveland, W. S. and Loader, C. (1996). Smoothing by local regression: principles and methods, *Computational Statistics*, **11** (to appear).

Donoho, D. L. (1994). Statistical estimation and optimal recovery, *Ann. Statist.*, **22**, 238–270.

Donoho, D. L. and Liu, R. C. (1991). Geometrizing rate of convergence III, *Ann. Statist.*, **19**, 668–701.

Donoho, D. L., Johnstone, I. M., Kerkyacharian, G. and Picard, D. (1995). Wavelet shrinkage: asymptopia?, *J. Roy. Statist. Soc. Ser. B*, **57**, 301–369.

Epanechnikov, V. A. (1969). Nonparametric estimation of a multidimensional probability density, *Theory Probab. Appl.*, **13**, 153–158.

Eubank, R. L. (1988). *Spline Smoothing and Nonparametric Regression*, Marcel Dekker, New York.

Fan, J. (1992). Design-adaptive nonparametric regression, *J. Amer. Statist. Assoc.*, **87**, 998–1004.

Fan, J. (1993). Local linear regression smoothers and their minimax efficiency, *Ann. Statist.*, **21**, 196–216.

Fan, J. and Gijbels, I. (1992). Variable bandwidth and local linear regression smoothers, *Ann. Statist.*, **20**, 2008–2036.

Fan, J. and Gijbels, I. (1995a). Data-driven bandwidth selection in local polynomial fitting: variable bandwidth and spatial adaptation, *J. Roy. Statist. Soc. Ser. B*, **57**, 371–394.

Fan, J. and Gijbels, I. (1995b). Adaptive order polynomial fitting: bandwidth robustification and bias reduction, *Journal of Computational and Graphical Statistics*, **4**, 213–227.

Fan, J. and Marron, J. S. (1994). Fast implementations of nonparametric curve estimators, *Journal of Computational and Graphical Statistics*, **3**, 35–56.

Gasser, T. and Engel, J. (1990). The choice of weights in kernel regression estimation, *Biometrika*, **77**, 377–381.

Gasser, T. and Müller, H.-G. (1984). Estimating regression functions and their derivatives by the kernel method, *Scand. J. Statist.*, **11**, 171–185.

Gasser, T., Müller, H.-G. and Mammitzsch, V. (1985). Kernels for nonparametric curve estimation, *J. Roy. Statist. Soc. Ser. B*, **47**, 238–252.

Granovsky, B. L. and Müller, H.-G. (1991). Optimizing kernel methods: a unifying variational principle, *Internat. Statist. Rev.*, **59**, 373–388.

Green, P. J. and Silverman, B. W. (1994). *Nonparametric Regression and Generalized Linear Models: a Roughness Penalty Approach*, Chapman and Hall, London.

Härdle, W. (1990). *Applied Nonparametric Regression*, Cambridge University Press, Boston.

Hastie, T. J. and Loader, C. (1993). Local regression: automatic kernel carpentry (with discussion), *Statist. Sci.*, **8**, 120–143.

Jennen-Steinmetz, C. and Gasser, T. (1988). A unifying approach to nonparametric regression estimation, *J. Amer. Statist. Assoc.*, **83**, 1084–1089.

Lejeune, M. (1985). Estimation non-paramétrique par noyaux: régression polynomiale mobile, *Rev. Statist. Appl.*, **33**, 43–68.

Müller, H.-G. (1987). Weighted local regression and kernel methods for nonparametric curve fitting, *J. Amer. Statist. Assoc.*, **82**, 231–238.

Müller, H.-G. (1988). *Nonparametric Analysis of Longitudinal Data*, Springer, Berlin.

Müller, H.-G. (1991). Smooth optimum kernel estimators near endpoints, *Biometrika*, **78**, 521–530.

Nadaraya, E. A. (1964). On estimating regression, *Theory Probab. Appl.*, **9**, 141–142.

Rao, C. R. (1973). *Linear Statistical Inference and Its Applications*, Wiley, New York.

Ruppert, D. and Wand, M. P. (1994). Multivariate weighted least squares regression, *Ann. Statist.*, **22**, 1346–1370.

Ruppert, D., Sheather, S. J. and Wand, M. P. (1995). An effective bandwidth selector for local least squares regression, *J. Amer. Statist. Assoc.*, **90**, 1257–1270.

Sacks, J. and Ylvisaker, D. (1981). Asymptotically optimum kernels for density estimation at a point, *Ann. Statist.*, **9**, 334–346.

Seifert, B., Brockmann, M., Engel, J. and Gasser, T. (1994). Fast algorithms for nonparametric curve estimation, *Journal of Computational and Graphical Statistics*, **3**, 192–213.

Silverman, B. W. (1984). Spline smoothing: the equivalent variable kernel method, *Ann. Statist.*, **12**, 898–916.

Wahba, G. (1990). *Spline Models for Observational Data*, SIAM, Philadelphia.

Wand, M. P. and Jones, M. C. (1995). *Kernel Smoothing*, Chapman and Hall, London.

Watson, G. S. (1964). Smooth regression analysis, *Sankhyā Ser. A*, **26**, 359–372.