

# Weierstraß-Institut für Angewandte Analysis und Stochastik

im Forschungsverbund Berlin e.V.

Preprint

ISSN 0946 – 8633

## Component identification and estimation in nonlinear high-dimensional regression models by structural adaptation

Samarov, Alexander, Spokoiny, Vladimir and Vial, Celine

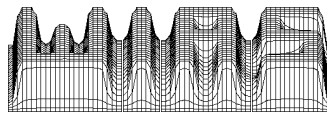
University of Massachusetts-Lowell  
and Massachusetts Institute of Technology.  
77 Massachusetts Avenue,  
Cambridge, MA 02139-4307, U.S.A.  
E-Mail: [samarov@mit.edu](mailto:samarov@mit.edu)

Weierstrass Institute  
and Humboldt University Berlin,  
Mohrenstr. 39, 10117 Berlin, Germany  
E-Mail: [spokoiny@wias-berlin.de](mailto:spokoiny@wias-berlin.de)  
URL: <http://www.wias-berlin.de/~spokoiny>

ENSAI, Campus Ker Lann,  
rue B. Pascal, 35170 Bruz, France  
E-Mail: [vial@ensai.fr](mailto:vial@ensai.fr)

No. 828

Berlin 2003



---

1991 *Mathematics Subject Classification.* 62H30, 62J02.

*Key words and phrases.* Structural adaptation, partially linear model, component analysis.

Edited by

Weierstraß-Institut für Angewandte Analysis und Stochastik (WIAS)

Mohrenstraße 39

D — 10117 Berlin

Germany

Fax: + 49 30 2044975

E-Mail: [preprint@wias-berlin.de](mailto:preprint@wias-berlin.de)

World Wide Web: <http://www.wias-berlin.de/>

## Abstract

This article proposes a new method of analysis of a partially linear model whose nonlinear component is completely unknown. The target of analysis is identification of the set of regressors which enter in a nonlinear way in the model function, and the complete estimation of the model including slope coefficients of the linear component and the link function of the nonlinear component. The procedure also allows for selecting the significant regression variables. As a by-product, we develop a test of linear hypothesis against a partially linear alternative, or, more generally, a test that the nonlinear component is  $M$ -dimensional for  $M = 0, 1, 2, \dots$ .

The approach proposed in this article is fully adaptive to the unknown model structure and applies under mild conditions on the model. The only important assumption is that the dimensionality of nonlinear component is relatively small. The theoretical results indicate that the procedure provides a prescribed level of the identification error and estimates the linear component with the accuracy of order  $n^{-1/2}$ . A numerical study demonstrates a very good performance of the method even for small or moderate sample sizes.

## 1 Introduction

We consider the model

$$y = f(\mathbf{x}) + \varepsilon, \quad f(\mathbf{x}) = \theta^\top \mathbf{x}_1 + G(\mathbf{x}_2), \quad (1.1)$$

where  $\mathbf{x}^\top = (\mathbf{x}_1^\top, \mathbf{x}_2^\top)$ ,  $\dim(\mathbf{x}_2) = M$ ,  $\dim(\mathbf{x}_1) = d - M$ , and  $M \ll d$ . Function  $G(\cdot)$ , vector of coefficients  $\theta$ , and the distribution of the noise  $\varepsilon$  are unknown. And most importantly, we do not know with respect to which  $d_1 = d - M$  variables  $\mathbf{x}_1$  the model is linear.

The model (1.1) naturally generalizes the linear model and are called a *partially linear* model. Such models can be used in analysis of high dimensional data when the assumption of linearity is too restrictive. They can also be used as a natural alternative to a linear model in the problem of testing the linearity assumption. A general case with a high dimensional nonlinear component makes the analysis complicated because of the “curse of dimensionality” problem. In this paper we consider the situation in which the nonlinear component is low dimensional, that is,  $M$  is relatively small.

Hristache, Juditsky and Spokoiny (2001) and Hristache, Juditsky, Polzehl and Spokoiny (2001) (referred to as HJS and HJPS, respectively, in the rest of the paper) proposed a new method of exploring a high-dimensional regression model with the help of a general structural adaptation approach. The aim of the present article is to apply this approach to the estimation and inference in the partially linear model (1.1). The analysis includes, in particular, estimation of the degree of nonlinearity  $M$ , identifying with respect to which  $d - M$  variables  $\boldsymbol{x}_1$  the model is linear or equivalently which  $M$  variables enter in  $f$  in a nonlinear fashion, estimation of the vector  $\boldsymbol{\theta}$  and of the nonlinear link function  $G$ .

It is important to note that the approach proposed here provides also a new method of selecting significant variables in nonparametric regression in case when the dimensionality of the nonlinear component is relatively small. More specifically, after selecting  $M$  (significant) nonlinear variables, one can further select variables among the linear ones using standard methods of linear regression analysis or by testing significance of linear variable slopes as suggested in Section 4.3 below.

As a by-product of our analysis, we develop a test of the hypothesis of linearity against a partially linear alternative, and, more generally, a test of the hypothesis that the dimensionality of the nonlinear component does not exceed the prescribed value  $M$ .

Following the work of Engle, Granger, Rice and Weiss (1986), much attention has been directed to estimating model (1.1). See, for example, Heckman (1986), Rice (1986), Chen (1988), Robinson (1988), Speckman (1988), Gao (1995), Schick (1996a,b), Bhattacharya and Zhao (1997), Mammen and Van der Geer (1997), Hamilton and Truong (1997), Eubank, Kambour, Kim, Klipple, Reese and Schimek (1998), Schimek (2000), Golubev and Härdle (2000, 2001). Further references and applications of partially linear models could be found in the recent book by Härdle, Liang, and Gao (2000). This literature addressed the problem of estimation of the parametric and nonparametric components of the model (1.1) under the assumption that the “nonlinear” variables  $\boldsymbol{x}_2$  are specified and, in fact, most papers assume that  $M = 1$ . Various estimators have been proposed which achieve root- $n$  rate or are semiparametrically efficient for estimating the parametric component  $\boldsymbol{\theta}$  as well as those which achieve the usual nonparametric rates for estimating  $G(\boldsymbol{x}_2)$ .

To our knowledge, the only paper which addressed the problem of selecting which variables  $\boldsymbol{x}_2$  enter nonlinearly in the model (1.1) was Chen and Chen (1991). That paper proposed a model-selection-type rule and showed that the probability of the correct identification by this method goes to one as the sample size goes to infinity. Härdle and Korostelev (1996) showed for the similar problem of selecting the significant variables in an additive model that the error of classification can be made exponentially small. In this paper we consider

another setup which seems to be more appealing for practical applications. Namely, we develop a nonlinear component identification method which guarantees a prescribed level of model misspecification uniformly over the class of models whose nonlinear component is separated away from the linear one by the squared distance of order  $n^{-1} \log n$  or larger. Our results are essentially nonasymptotic and apply for a small or moderate sample size.

Härdle, Spokoiny and Sperlich (2001) considered a similar problem of identifying the linear component for an additive model, using a wavelet (Haar) expansion of every additive component. The advantage of the structure adaptive procedure proposed here is that the additive structure is not required and is not used in the method.

There also exists a large literature on testing a parametric, in particular linear, regression model against nonparametric alternative. See, for example, Eubank and Spiegelman (1990), Eubank and Hart (1992), Ledwina (1994), Härdle and Mammen (1993), Fan (1996), Hart (1997), Stute (1997), Horowitz and Spokoiny (2001) and references therein. Our testing results are stated in the spirit of Spokoiny (2001) focusing on the minimal separation distance between the null and the alternative providing test consistency.

The paper is organized as follows. Section 2 contains the description of the structure adaptive estimation algorithm. Accuracy of estimation by the proposed method is described in Section 3. Further problems of identification of the nonlinear component and of estimation of slope coefficients of the linear component are discussed in Section 4. Section 5 presents results of a simulation study and an application to real data. Conclusion and some extensions of the method are presented in Section 6. The proofs are collected in the Appendix.

## 2 Structure adaptive procedure

This section explains the adaptive estimation procedure starting with a short heuristic discussion.

### 2.1 Preliminaries

The idea of structural adaptation from HJPS can be summarized as follows.

- (i) knowing the structural information helps better estimate the model function;
- (ii) a good pilot estimator of the model function helps recover some structural information about the model.

These two observations lead to the following iterative procedure: we start with a purely nonparametric estimator of the model function; then the above two steps (estimation of the model and estimation of the structure) are iterated several times increasing the amount of structural information and improving the quality of model estimation during iteration.

HJPS considered the problem of estimation for a multi-index model in which the regression function is of the form  $f(\mathbf{x}) = g(\theta_1^\top \mathbf{x}, \dots, \theta_M^\top \mathbf{x})$ , where  $\theta_1, \dots, \theta_M$  are unknown index vectors in  $\mathbb{R}^d$ . The partially linear model (1.1) can be regarded as a special case of the multi-index model with  $M + 1$  indices. Indeed,  $f(\mathbf{x})$  depends on  $\mathbf{x}$  only through  $\theta^\top \mathbf{x}_1$  and the coordinate vectors corresponding to the nonlinear component. So, one can formally apply the procedure from HJPS in the considered case. However, the special structure of the model (1.1) allows to considerably simplify the procedure and further analysis that justifies a separate treatment of the partially linear models.

Here the structure of the model (1.1) is described by the set  $\mathcal{J}$  of indices corresponding to the nonlinear component  $\mathbf{x}_2$ . An alternative description can be done by using the average gradient idea. Namely, if the function  $f(\mathbf{x})$  is linear with respect to the  $m$ th coordinate function  $x_m$ , then the partial derivative  $\partial f / \partial x_m$  is a constant, and therefore, the variance  $V_m$  of the  $m$ th partial derivative can be used to measure the degree of nonlinearity of the  $m$ th coordinate. Suppose that some information about the set  $\mathcal{J}$  or, equivalently, about the values  $V_m$  is available. Now we explain how this information can be used for improving the quality of estimation of the model function  $f$ . A local linear estimator of the function  $f$  and its gradient  $\nabla f$  at a point  $X_i$  is given by

$$\begin{pmatrix} \widehat{f}(X_i) \\ \widehat{\nabla f}(X_i) \end{pmatrix} = \left\{ \sum_{j=1}^n \begin{pmatrix} 1 \\ X_{ij} \end{pmatrix} \begin{pmatrix} 1 \\ X_{ij} \end{pmatrix}^\top K\left(\frac{|X_{ij}|^2}{b^2}\right) \right\}^{-1} \sum_{j=1}^n Y_j \begin{pmatrix} 1 \\ X_{ij} \end{pmatrix} K\left(\frac{|X_{ij}|^2}{b^2}\right)$$

where  $X_{ij} = X_j - X_i$ ,  $b$  is a *bandwidth* and  $K(\cdot)$  is a univariate kernel  $K$  on the positive semiaxis, supported on  $[0, 1]$ . The bandwidth  $b$  should be selected in a way that the ball with the radius  $b$  and the center the point of estimation  $X_i$  contains at least  $d + 1$  design points which for large value of  $d$  leads to a the bandwidth  $b$  of order one and to a huge estimation bias. This phenomenon is called the ‘‘curse of dimensionality’’. Observe now that the function  $f$  has anisotropic smoothness properties: smoothness of  $G$  in direction of the nonlinear component, and infinite smoothness (corresponding to a linear function) in other directions. This suggests to apply an anisotropic bandwidth for estimating the model function and its gradient. So, the ‘ideal’ estimator which utilizes the knowledge of the set  $\mathcal{J}$  can be defined by using the different bandwidths for different components of the vector  $\mathbf{x}$ . Let  $\mathbf{b} = \text{diag}(b_1, \dots, b_d)$  be a diagonal matrix with the diagonal entries

$b_1, \dots, b_d$ . Define the local linear estimator with the *anisotropic bandwidth*  $\mathbf{b}$  by

$$\begin{pmatrix} \widehat{f}(X_i) \\ \widehat{\nabla}f(X_i) \end{pmatrix} = \left\{ \sum_{j=1}^n \begin{pmatrix} 1 \\ X_{ij} \end{pmatrix} \begin{pmatrix} 1 \\ X_{ij} \end{pmatrix}^\top K_d(X_{ij}, \mathbf{b}) \right\}^{-1} \sum_{j=1}^n Y_j \begin{pmatrix} 1 \\ X_{ij} \end{pmatrix} K_d(X_{ij}, \mathbf{b}), \quad (2.1)$$

where  $K_d(\mathbf{x}, \mathbf{b}) = K(|\mathbf{b}^{-1}\mathbf{x}|^2)$ . Knowing that the value  $V_m$  is significantly positive (that is,  $m$  is presumably in  $\mathcal{J}$ ) leads to the choice of smaller bandwidth values  $b_m$  for such  $m$  and possibly larger bandwidths for the other regressors. This would help to avoid the ‘‘curse of dimensionality’’ problem if the dimension  $M$  of the nonlinear component is not too large, cf. Carroll, Fan, Gijbels and Wand (1997) or HJPS.

Next we explain how the structural information can be extracted from the pilot estimator (2.1) of the model function. Define for every coordinate  $x_m$  of  $\mathbf{x} \in \mathbb{R}^d$  a set of functions  $\psi_{1m}, \dots, \psi_{Lm}$  satisfying the conditions:

$$\sum_{i=1}^n \psi_{lm}(X_{i,m}) = 0, \quad n^{-1} \sum_{i=1}^n \psi_{lm}(X_{i,m}) \psi_{l'm}(X_{i,m}) = \delta_{ll'}.$$

In other words,  $\{\psi_{lm}, l = 1, \dots, L\}$  is a orthonormal set of functions with respect to the design of  $m$ th coordinate. Each of  $\psi_{lm}$  is also orthogonal to the constant function. The latter property implies that if  $f$  is linear with respect to  $x_m$ , then

$$\beta_{lm}^* = n^{-1} \sum_{i=1}^n \nabla f_m(X_i) \psi_{lm}(X_{i,m}) \equiv 0 \quad (2.2)$$

for every  $l = 1, \dots, L$ , where  $\nabla f_m(\mathbf{x}) = \partial f / \partial x_m(\mathbf{x})$ . Thus, the sum

$$v_m^* = \sum_{l=1}^L (\beta_{lm}^*)^2$$

can be used as the measure of nonlinearity of  $f$  with respect to  $x_m$ . Having estimated the gradient of  $f$  for all  $X_i$ , we can also estimate the coefficients  $\beta_{lm}$  with

$$\widehat{\beta}_{lm} = n^{-1} \sum_{i=1}^n \widehat{\nabla}f_m(X_i) \psi_{lm}(X_{i,m}) \quad (2.3)$$

and use the sum  $\widehat{v}_m = \widehat{\beta}_{1m}^2 + \dots + \widehat{\beta}_{Lm}^2$  as the estimated degree of nonlinearity of  $f$  with respect to the  $m$ th regression variable.

Next, the quantities  $\widehat{v}_m$  can be used to define new anisotropic bandwidth  $\mathbf{b}$  taking smaller bandwidths for the regressors with large  $\widehat{v}_m$ .

**Remark 2.1.** Similarly to HJPS, we use here the estimation method based on the Fourier expansion of the gradient  $\nabla f(\mathbf{x})$ . Alternatively, one can estimate  $V_m$  directly using the

average of  $|\nabla f_m(X_i)|^2$ . However, a detailed calculation (not given in the paper) shows that the procedure based on such a direct estimation of the quadratic functionals  $V_m$  leads to worse estimation results. At the same time, the loss of information from replacing  $V_m$  with  $v_m^*$  as a measure of nonlinearity is not significant if  $L$  is chosen sufficiently large, see more on the choice of  $L$  in Section 6.

## 2.2 Iterative procedure

We now present the description of the method. The procedure involves input parameters  $h_1, a_h, \rho_1, \rho_{\min}, a_\rho$  and  $\eta$ . The parameter of ellipticity  $\rho$  decreases geometrically from  $\rho_1$  to  $\rho_{\min}$  by the factor  $a_\rho < 1$  while the bandwidth  $h$  increases geometrically from  $h_1$  by the factor  $a_h > 1$  during iterations. The value  $\eta$  can be interpreted as the “memory parameter” of the procedure. The choice of these parameters, as well as of the set of basis functions  $\{\psi_{lm}\}$  will be discussed in Section 2.3. The algorithm reads as follows:

1. Select  $h_1$ . Set  $\widehat{v}_1^{(0)} = \dots = \widehat{v}_d^{(0)} = 0$ , and  $k = 1$ . Compute for  $i = 1, \dots, n$

$$\widehat{V}_i^{(0)} = \sum_{j=1}^n \begin{pmatrix} 1 \\ X_{ij} \end{pmatrix} \begin{pmatrix} 1 \\ X_{ij} \end{pmatrix}^\top, \quad \widehat{S}_i^{(0)} = \sum_{j=1}^n \begin{pmatrix} 1 \\ X_{ij} \end{pmatrix} Y_j,$$

where  $X_{ij} = X_j - X_i$ .

2. Compute

$$b_m^{(k)} = h_k \left( 1 + \rho_k^{-2} \widehat{v}_m^{(k-1)} \right)^{-1/2}, \quad m = 1, \dots, d. \quad (2.4)$$

Define  $\mathbf{b}^{(k)} = \text{diag}(b_1^{(k)}, \dots, b_d^{(k)})$ .

3. For every  $X_i$  compute

$$V_i^{(k)} = \eta V_i^{(k-1)} + (1 - \eta) \sum_{j=1}^n \begin{pmatrix} 1 \\ X_{ij} \end{pmatrix} \begin{pmatrix} 1 \\ X_{ij} \end{pmatrix}^\top K_d(X_{ij}, \mathbf{b}^{(k)}), \quad (2.5)$$

$$S_i^{(k)} = \eta S_i^{(k-1)} + (1 - \eta) \sum_{j=1}^n \begin{pmatrix} 1 \\ X_{ij} \end{pmatrix} Y_j K_d(X_{ij}, \mathbf{b}^{(k)}), \quad (2.6)$$

and

$$\begin{pmatrix} \widehat{f}^{(k)}(X_i) \\ \widehat{\nabla f}^{(k)}(X_i) \end{pmatrix} = \left( V_i^{(k)} \right)^{-1} S_i^{(k)}. \quad (2.7)$$

Compute  $\widehat{\varepsilon}_i^{(k)} = Y_i - \widehat{f}^{(k)}(X_i)$ .

4. Compute for  $m = 1, \dots, d$  and  $l = 1, \dots, L$

$$\widehat{\beta}_{lm}^{(k)} = n^{-1} \sum_{i=1}^n \widehat{\nabla f}_m^{(k)}(X_i) \psi_{lm}(X_{i,m}), \quad \widehat{v}_m^{(k)} = \sum_{l=1}^L \left| \widehat{\beta}_{lm}^{(k)} \right|^2. \quad (2.8)$$



If  $\widehat{v}_m^{(k)} > 1$ , then set  $\widehat{v}_m^{(k)} = 1$ .

5. Increase  $k$  by 1. Set  $\rho_{k+1} = a_\rho \rho_k$ ,  $h_{k+1} = a_h h_k$ . If  $\rho_{k+1} \geq \rho_{\min}$ , then set  $k = k + 1$  and continue with Step 2; otherwise terminate.

Denote by  $k_n$  the number of iterations and by  $\widehat{\beta}_{lm} = \widehat{\beta}_{lm}^{(k_n)}$  (resp.  $\widehat{v}_m = \widehat{v}_m^{(k_n)}$ ) the last step estimators of  $\beta_m^*$  (resp.  $v_m^*$ ). Similarly,  $\widehat{f}(X_i)$  denotes the last step estimator of  $f(X_i)$  and  $\widehat{\varepsilon}_i = \widehat{\varepsilon}_i^{(k_n)} = Y_i - \widehat{f}(X_i)$ .

**Remark 2.2.** At every step  $k$  of the algorithm the bandwidth  $b_m^{(k)}$  for the  $m$ th regressor is selected between  $h_k$  and  $h_k \rho_k$  depending on the value  $\widehat{v}_m^{(k-1)}$ . For the linear component, the values  $\widehat{v}_m^{(k-1)}$  should be small leading to a bandwidth about  $h_k$ , while for the nonlinear regressors with a large value  $v_m^*$ , the estimator  $\widehat{v}_m^{(k-1)}$  will be also large leading to a bandwidth about  $h_k \rho_k$ . During iteration the parameter  $h$  grows to  $h_{\text{final}} \geq 1$  while  $h_k \rho_k$  decreases to  $\rho_{\min}$  leading to the adaptive anisotropic bandwidth at the last step.

**Remark 2.3.** We cut  $\widehat{v}_m^{(k)}$  at one at step 4 in order to avoid too strong shrinkage in direction of  $m$ th regressor which may occur for too large values of  $\widehat{v}_m^{(k)}$ .

### 2.3 Choice of parameters

It is obvious that the quality of estimation by the proposed method strongly depends on the rule for changing the parameters  $h$  and  $\rho$ , and, in particular, on their values at the initial and final iteration. Some related discussion about this choice can be found in HJPS. The general idea is to ensure that the parameter  $h$  grows to one and  $h\rho$  decreases under the constraint that at every iteration  $k$  there exist enough design points in every or almost every local ellipsoidal neighborhoods  $E^{(k)}(X_i) = \{\mathbf{x} : |(\mathbf{b}^{(k)})^{-1}(\mathbf{x} - X_i)|^2 \leq 1\}$ .

Assuming that every  $b_m^{(k)}$  is close to the ‘ideal bandwidth’  $b_m^{*(k)} = h_k(1 + \rho_k^{-2} v_m^*)^{-1/2}$  we observe, that neighborhood  $E^{(k)}(X_i)$  is stretched at each iteration step by factor  $a_h$  in all directions and is shrunk by a factor about  $a_\rho$  in directions of the  $M$ -dimensional nonlinear component  $\mathcal{J}$  where  $a_h$  and  $a_\rho$  are parameters of the procedure. Therefore, the Lebesgue measure of every such neighborhood is changed each time by a factor about  $a_h^d a_\rho^M$ . This leads to the constraint  $a_h^d a_\rho^M \geq 1$ , cf. Assumption 3 in Section 3 below. Under the assumption of a random design with a positive density, this would result in an increase of the mean number of design points inside each  $E^{(k)}(X_i)$ . Our theoretical results will be stated for the choice  $h_1 \asymp n^{-1/\max\{4,d\}}$ ,  $h_{\max} \asymp 1$ ,  $\rho_1 = 1$ ,  $\rho_{\min} \asymp (n^{-1} \log n)^{1/3}$ , see Section 3 for more details. Similarly to HJPS, such a choice under the constraint  $a_h^d a_\rho^M > 1$  is possible only for  $M \leq 3$ .

We recommend to define for every  $m = 1, \dots, d$  the set of functions  $\psi_{lm}$ ,  $l = 1, \dots, L$

by orthogonalizing the set of polynomials  $x_m, x_m^2, \dots, x_m^L$  with respect to the design of the  $m$ th regressor under the constraint  $\sum_{i=1}^n \psi_{lm}(X_{i,m}) = 0$ . A model or variable dependent choice of the basis  $\{\psi_{lm}\}$  is possible as well. The “memory parameter”  $\eta$  used in (2.5) and (2.6) can be taken between 0.1 and 0.5. The number  $L$  can be taken between 5 and 10, see Section 6 for more discussion.

**Remark 2.4.** Similarly to HJS and HJPS we apply in our numerical study a slightly modified procedure. The only difference is in the definition of the estimated vectors  $\widehat{\beta}_{lm}$ . Namely we define

$$\widehat{\beta}_{lm}^{(k)} = \left( \sum_{i=1}^n w_i^{(k)} \right)^{-1} \sum_{i=1}^n w_i^{(k)} \widehat{\nabla} f_m^{(k)}(X_i) \psi_{lm}(X_{i,m}),$$

where  $w_i^{(k)}$  is square root of the smallest eigenvalue of the matrix  $V_i^{(k)}$ , that is,  $w_i^{(k)} = \lambda_{\min}^{1/2}(V_i^{(k)})$ . In addition, the basis functions  $\psi_{lm}$  should be modified as each step to satisfy the condition  $\sum_{i=1}^n w_i^{(k)} \psi_{lm}(X_i) = 0$ .

## 2.4 Estimation of the noise variance

The variance  $\sigma^2$  of the noise  $\varepsilon$  does not enter in the description of the method. However, it will be used for defining the stopping rule of the algorithm and the resampling procedure in Section 4. Here we briefly discuss how this variance can be estimated under the assumption of the noise homogeneity at every step of the algorithm.

A natural variance estimator can be defined on the base of residuals squared after each the step  $k$ :  $|\widehat{\sigma}^{(k)}|^2 = n^{-1} \sum_{i=1}^n |\widehat{\varepsilon}_i^{(k)}|^2$ . This simple crude estimator can be refined, see e.g. Gasser, Sroka and Jennen-Steinmetz (1986) or Spokoiny (2002) and reference therein. Namely, the residuals  $\widehat{\varepsilon}_i^{(k)}$  can obviously be represented in the form  $\widehat{\varepsilon}_i^{(k)} = \sum_{j=1}^n c_{ij}^{(k)} Y_j$  where  $c_{ij}^{(k)}$  are known coefficients. These coefficients are random and dependent on the  $Y_i$ ’s through the random bandwidths  $\widehat{b}_m^{(k)}$ . However, our theoretical results indicate that one can ignore this dependence and proceed as if the coefficients  $c_{ij}^{(k)}$  were deterministic and correspond to the “ideal” bandwidths  $b_m^{*(k)}$ .

Next, if the function  $f$  is sufficiently smooth, then the distribution of the residuals  $\widehat{\varepsilon}_i$  only weakly depends on this function  $f$  and can be effectively evaluated for  $f \equiv 0$ . In the last case,  $\mathbf{E}|\widehat{\varepsilon}_i^{(k)}|^2 = \sigma^2 \sum_{j=1}^n |c_{ij}^{(k)}|^2$  that leads to the estimator

$$|\widehat{\sigma}^{(k)}|^2 = \frac{1}{n} \sum_{i=1}^n \left( \sum_{j=1}^n |c_{ij}^{(k)}|^2 \right)^{-1} |\widehat{\varepsilon}_i^{(k)}|^2. \quad (2.9)$$

The properties of this estimator are briefly discussed in Section 3.4 below.

### 3 Accuracy of estimation

In this section we present the results about the accuracy of estimation of the functionals  $\beta_{lm}^*$  and  $v_m^*$  by the proposed iterative procedure.

#### 3.1 Assumptions

As in HJPS, it is useful to proceed with the renormalized link function. In what follows we consider the model

$$f(\mathbf{x}) = \mathbf{x}_1^\top \theta + g(\mathcal{R}^* \mathbf{x}_2) \quad (3.1)$$

where  $\mathcal{R}^*$  is the diagonal  $M \times M$  matrix with diagonal entries  $\sqrt{v_m^*}$ ,  $m \in \mathcal{J}$  and  $g$  is a nonlinear link function.

Our main results will be stated under the following assumptions.

**Assumption 1. (Kernel)** The kernel  $K(\cdot)$  is continuously differentiable decreasing function on  $\mathbb{R}_+$  with  $K(0) = 1$  and  $K(x) = 0$  for all  $|x| \geq 1$ .

**Assumption 2. (Errors)** The random variables  $\varepsilon_i$  in (1.1) are independent and normally distributed with zero mean and variance  $\sigma^2$ .

**Assumption 3. (Range of parameters  $h_k, \rho_k$ )** The parameters of the procedure fulfill  $\rho_1 = 1$ ,  $\rho_{\min} = (\sigma^2 n^{-1} L \log n)^{1/3}$ ,  $h_1 = C_0 n^{-\frac{1}{4vd}}$  with a constant  $C_0 \geq 1$ ,  $h_{\max} \geq 1$  and  $a_h^d a_\rho^M \geq 1$ .

**Assumption 4. (Link function)** The function  $g$  from (3.1) is twice differentiable with a bounded second derivative, so that, for some constant  $C_g$  and for all  $u, v \in \mathbb{R}^M$

$$|g(v) - g(u) - (v - u) \nabla g(u)| \leq C_g |u - v|^2;$$

Our last assumption concerns the design properties. In what follows we assume that the design is deterministic. That is,  $X_1, \dots, X_n$  are non-random points in  $\mathbb{R}^d$ . Note, however, that the case of a random design can be considered as well, supposing  $X_1, \dots, X_n$  independent and identically distributed random points in  $\mathbb{R}^d$  with a design density  $p(x)$ . Then all the results should be understood conditionally on the design.

In order for the procedure to work, we have to suppose that the design points  $(X_i)$  are “well diffused”, as a consequence, at  $k$ th iteration of the algorithm, all local gradient estimators from (2.7) corresponding to the anisotropic bandwidth  $\mathbf{b}^{(k)} = \text{diag}(b_1^{(k)}, \dots, b_d^{(k)})$  from (2.4) are well defined. The latter is equivalent to the condition that all the matrices  $V_i^{(k)}$  from (2.5) are non-singular. We also define for the  $k$ th iteration the “ideal

anisotropic bandwidth"  $\mathbf{b}^{*(k)}$  having the diagonal entries  $b_m^{*(k)} = (1 + \rho_k^{-2} v_m^*)^{-1/2} h_k$ . The closeness of  $\mathbf{b}^{(k)}$  to the "ideal bandwidth"  $\mathbf{b}^{*(k)}$  can be characterized by the values  $U_m^{(k)} = (b_m^{(k)}/b_m^{*(k)})^2 = (1 + \rho_k^{-2} v_m^*) / (1 + \rho_k^{-2} \widehat{v}_m^{(k-1)})$ ,  $m = 1, \dots, d$ . If  $\widehat{v}_m^{(k-1)} = v_m^*$ , then  $U_m^{(k)} = 1$ . The condition we need means that at the step  $k$  of the algorithm, for every anisotropic bandwidth  $\mathbf{b} = \text{diag}(b_1, \dots, b_d)$  close to  $\mathbf{b}^{*(k)}$  in the above sense, the design is regular within the elliptic neighborhood with the center at each point  $X_i$  and with the principal semiaxis  $b_m$ ,  $m = 1, \dots, d$ .

Define  $Z_{ij}^{(k)} = (\mathbf{b}^{*(k)})^{-1}(X_j - X_i)$  for  $i, j = 1, \dots, n$ . These vectors describe locations of the design points in the coordinate system shifted by  $X_i$  and rescaled by  $\mathbf{b}^{*(k)}$ . For a vector  $U = (U_1, \dots, U_d)^\top \in \mathbb{R}^d$  with  $U_m \geq 0$ , define  $D_U = \text{diag}(U_1, \dots, U_d)$ . Then, for  $\mathbf{b} = D_U^{-1/2} \mathbf{b}^{*(k)}$ , it holds  $K_d(X_{ij}, \mathbf{b}) = K((Z_{ij}^{(k)})^\top D_U Z_{ij}^{(k)})$ . Set

$$\begin{aligned} N_i^{(k)}(U) &= \sum_{j=1}^n K\left((Z_{ij}^{(k)})^\top D_U Z_{ij}^{(k)}\right), \quad i = 1, \dots, n, \\ \mathcal{V}_i^{(k)}(U) &= \sum_{j=1}^n \begin{pmatrix} 1 \\ Z_{ij}^{(k)} \end{pmatrix} \begin{pmatrix} 1 \\ Z_{ij}^{(k)} \end{pmatrix}^\top K\left((Z_{ij}^{(k)})^\top D_U Z_{ij}^{(k)}\right), \quad i = 1, \dots, n. \end{aligned}$$

In what follows  $\|A\|$  stands for the matrix norm associated with the Euclidean vector norm:  $\|A\| = \sup_\lambda |A\lambda|/|\lambda|$ .

**Assumption 5. (Design)** There exist constants  $C_V$ ,  $C_K$ ,  $C_{K'}$  and some  $\alpha \in (0, 1/2)$ , such that for all vectors  $U \in \mathbb{R}^d$  satisfying  $|U_m - 1| \leq \alpha$ ,  $m = 1, \dots, d$ , and for  $k \leq k_n$  the inverse matrices  $\mathcal{V}_i^{(k)}(U)^{-1}$  are well defined with

$$N_i^{(k)}(U) \left\| \mathcal{V}_i^{(k)}(U)^{-1} \right\| \leq C_V, \quad i = 1, \dots, n;$$

Moreover, for  $i, j = 1, \dots, n$ ,

$$\begin{aligned} \sum_{i=1}^n \frac{1}{N_i^{(k)}(U)} K\left((Z_{ij}^{(k)})^\top D_U Z_{ij}^{(k)}\right) &\leq C_K, \\ \sum_{i=1}^n \frac{1}{N_i^{(k)}(U)} \left| K'\left((Z_{ij}^{(k)})^\top D_U Z_{ij}^{(k)}\right) \right| &\leq C_{K'}, \\ \sum_{j=1}^n \frac{1}{N_i^{(k)}(U)} \left| K'\left((Z_{ij}^{(k)})^\top D_U Z_{ij}^{(k)}\right) \right| &\leq C_{K'}. \end{aligned}$$

where  $K'$  means the derivative of  $K$ .

**Remark 3.1.** As already mentioned in HJS and HJPS, in the case of random design with a continuous positive density one can fix some constants  $C_V$ ,  $C_K$  and  $C_{K'}$  (which depend on the dimension  $d$  and the design distribution) such that the bounds in Assumption 5 hold with probability which converges to 1 exponentially fast as  $n$  grows to infinity.

### 3.2 The first step estimator

Our first result describes the quality of the estimators  $\widehat{\beta}_{lm}^{(1)}$  obtained at the first step of the algorithm. These estimators correspond to usual nonparametric local linear estimation of the function  $f$  and its gradient. We also state the result about the accuracy of estimation of the values  $v_m^*$  at the first step.

Let  $\beta_m^*$  denote an  $L$ -vector with the components  $\beta_{lm}^*$ ,  $l = 1, \dots, L$ .

**Proposition 3.1.** *Let Assumptions 1 through 5 hold. For the first-iteration estimator  $\widehat{\beta}_m^{(1)}$  of the vector  $\beta_m^*$ ,  $m = 1, \dots, d$ , it holds:*

$$\widehat{\beta}_m^{(1)} - \beta_m^* = s_m h_1 + \frac{\eta_m}{h_1 \sqrt{n}},$$

where deterministic  $L$ -vectors  $s_m$  satisfy  $|s_m| \leq C_g C_V^{1/2} v_{(1)}^*$  with  $v_{(1)}^* = \max_{m=1, \dots, d} v_m^*$ , and the  $\eta_m$ 's are mean zero Gaussian random  $L$ -vectors with components  $\eta_{lm}$  such that  $E|\eta_{lm}|^2 \leq 2\sigma^2 C_V^2 C_K$ ,  $l = 1, \dots, L$ . Also, it holds

$$\mathbf{P} \left( \max_{m=1, \dots, d} \left| \widehat{\beta}_m^{(1)} - \beta_m^* \right| > \delta_1 \right) \leq \frac{1}{n}, \quad (3.2)$$

where

$$\delta_1 = C_g C_V^{1/2} v_{(1)}^* h_1 + \frac{\sqrt{2L} \sigma C_V C_K^{1/2} z_n}{h_1 \sqrt{n}}, \quad (3.3)$$

and  $z_n = (1 + 2 \log(nd) + 2 \log \log(nd))^{1/2}$ .

Moreover, for the first-iteration estimator  $\widehat{v}_m^{(1)}$ ,  $m = 1, \dots, d$ , it holds:

$$\mathbf{P} \left( \left| \widehat{v}_m^{(1)} - v_m^* \right| \leq \delta_1^2 + 2\delta_1 \tau_{m,1}, \forall m = 1, \dots, d \right) \geq 1 - \frac{1}{n}, \quad (3.4)$$

where  $\tau_{m,1} = \sqrt{v_m^*} (1 + v_m^*)^{-1/2} \leq \min\{1, \sqrt{v_m^*}\}$ .

### 3.3 The quality of the final estimators

Now we present the result which indicates how the accuracy of estimation can be improved by the iterative algorithm. As in HJS and HJPS, the quality of the final estimators depends on the 'direction'. This quality is of order  $n^{-1/2}$  for the linear component and is worse for the nonlinear component. This fact has a very simple explanation: estimation of a nonlinear component is a harder task than that of a linear one; hence, the worse accuracy. To express this fact, we introduce the scaling factor  $P_{\rho, m} = (1 + \rho^{-2} v_m^*)^{-1/2}$ , where  $\rho$  is a running parameter of the procedure. Note that  $P_{\rho, m} = 1$  for all linear regressors which have  $v_m^* = 0$ . If  $v_m^*$  is a positive constant, then  $P_{\rho, m} \asymp \rho$ . We will see that the estimation

error  $\widehat{\beta}_{lm} - \beta_{lm}^*$ , after being multiplied by  $P_{\rho,m}$ , can be bounded uniformly over  $l, m$  at every step of the algorithm. This implies, in particular, that the quality of estimation of the nonlinear component is about  $P_{\rho}^{-1} \asymp \rho^{-1}$  times worse than the quality for the linear one.

In the next theorem and in Theorem 4.1 below,  $\rho$  (resp.  $h$ ) denotes  $\rho_{k_n}$  (resp.  $h_{k_n}$ ) at the last iteration. Recall that  $h, \rho$  satisfy conditions  $h \geq 1$  and  $\rho = (\sigma^2 n^{-1} L \log n)^{1/3}$ .

**Theorem 3.1.** *Let Assumptions 1 through 5 hold. Then there exist a random set  $A$  with  $\mathbf{P}(A) \geq 1 - 3k_n/n$  and, for every  $m = 1, \dots, d$ , a Gaussian zero mean random vector  $\xi_m^* = (\xi_{1m}^*, \dots, \xi_{Lm}^*)^\top \in \mathbb{R}^L$  defined as a linear combination of the errors  $\varepsilon_i$  with deterministic coefficients, which depend on the ‘‘ideal’’ bandwidth  $\mathbf{b}^* = \mathbf{b}^{*(k_n)}$ , the design  $X_1, \dots, X_n$ , basis functions  $\psi_{lm}(\cdot)$ , and the kernel  $K$  only, and such that*

$$\mathbf{E}|\xi_{lm}^*|^2 \leq 2\sigma^2 C_V^2 C_K \quad l = 1, \dots, L, m = 1, \dots, d,$$

and on  $A$  it holds

$$\begin{aligned} \max_{m=1, \dots, d} \left| P_{\rho,m}(\widehat{\beta}_m - \beta_m^*) - n^{-1/2} \xi_m^* \right| &\leq C (\sigma^2 n^{-1} L \log n)^{2/3}, \\ \max_{m=1, \dots, d} \left| P_{\rho,m}(\widehat{\beta}_m - \beta_m^*) \right| &\leq \delta_n, \\ \max_{m=1, \dots, d} \left| P_{\rho,m}^2(\widehat{v}_m - v_m^*) \right| &\leq \delta_n^2 + 2\delta_n \tau_m, \end{aligned} \quad (3.5)$$

where  $C = C(d, M, C_g, C_V, C_K, C_{K'}, \bar{\psi})$ ,  $\bar{\psi} = \max_{i,l,m} |\psi_{lm}(X_i)|$ ,

$$\delta_n = \sqrt{2C_V^2 C_K \sigma^2 n^{-1} L z_n^2} + C (\sigma^2 n^{-1} L \log n)^{2/3} \quad (3.6)$$

and  $\tau_m = \rho \sqrt{v_m^*} (\rho^2 + v_m^*)^{-1/2} \leq \min\{\rho, \sqrt{v_m^*}\}$ . This implies that on  $A$  for every  $m \notin \mathcal{J}$ , with  $\omega_n = C (\sigma^2 n^{-1} L \log n)^{2/3}$ :

$$|\widehat{\beta}_m - n^{-1/2} \xi_m^*| \leq \omega_n, \quad |\widehat{\beta}_m| \leq \delta_n, \quad \left| |\widehat{\beta}_m|^2 - n^{-1} |\xi_m^*|^2 \right| \leq \omega_n^2 + 2\omega_n \delta_n. \quad (3.7)$$

**Remark 3.2.** The meaning of the random set  $A$  appearing in Theorem 3.1 can be understood as follows. The result of every iteration of the algorithm is random. With some probability it may happen that the estimation result at some step of the procedure does not follow the model structure. For instance, with some probability,  $\widehat{v}_m$  can be large even if  $v_m^* = 0$ . Our results indicate that the overall probability of such events is rather small and their complement is precisely the set  $A$  (of a dominating probability) on which the procedure ‘works’, that is, the iterative procedure leads to improvement of the quality of estimation at every iteration. The other results of Theorem 3.1 claim that on the set  $A$ , the adaptive estimators  $\widehat{\beta}_{lm}$  behave essentially as the ‘ideal’ estimators  $\widehat{\beta}_{lm}^*$  corresponding to the ‘ideal’ bandwidth  $\mathbf{b}^*$ . Since our further analysis is based on the final step estimators  $\widehat{\beta}_{lm}$ , all our results that follow will also be stated conditionally on this set  $A$ .

**Remark 3.3. (Origin of the constraint  $M \leq 3$ )** It follows from the proof of Theorem 3.1 that the bias of the ‘ideal’ estimators  $\widehat{\beta}_{lm}^*$  based on the local linear smoothing with the ‘ideal’ bandwidth  $\mathbf{b}^* = \mathbf{b}^{*(k_n)}$  is of the order  $(n^{-1} \log n)^{-2/3}$  only if the dimensionality  $M$  of the nonlinear component does not exceed 3. For  $M \geq 4$ , the model dependent bias of estimation is of order  $n^{-1/2}$  or larger while the stochastic component (which is model free) is of order  $n^{-1/2}$ . The same applies for the adaptive estimators  $\widehat{\beta}_{lm}$ . Therefore, the leading term in the estimation loss is model free only for  $M \leq 3$ , and the estimators  $\widehat{\beta}_{lm}$  do not achieve asymptotic normality at root-n rate for  $M \geq 4$ .

### 3.4 Variance estimation

The algorithm delivers an estimator  $\widehat{\sigma}^2$ , see (2.9), of the error variance  $\sigma^2$ . This estimator also utilizes the estimated structural information and improves upon the purely nonparametric variance estimators. Spokoiny (2002) has shown that in a general high dimensional regression model with  $d > 8$ , a root-n consistent estimation of the variance  $\sigma^2$  is impossible. Here the use of the structural assumption allows to relax this condition and to get a root-n accuracy for any  $d$ .

**Theorem 3.2.** *Let Assumptions 1 through 5 hold. There exists a constant  $C_\sigma$ , which depends on the constants entering in these assumptions only, such that*

$$\mathbf{P}(\sqrt{n}|\widehat{\sigma}^2 - \sigma^2| > C_\sigma \sigma^2 \lambda) \leq 2e^{-\lambda^2/4} + 3k_n/n.$$

## 4 Inference in a partially linear model

This section explains how the model (1.1) can be explored using our iterative procedure and results of Section 3. First we state the important separation result that will be used in the analysis below.

Let some integer  $\mathcal{M}$  be fixed. We put the estimated values  $\widehat{v}_m$  in the decreasing order  $\widehat{v}_{(1)} \geq \widehat{v}_{(2)} \geq \widehat{v}_{(3)} \dots$  and denote by  $\widehat{\mathcal{J}}_{\mathcal{M}}$  the index set corresponding to the  $\mathcal{M}$  largest values  $\widehat{v}_m$ . Theorem 3.1 implies the following separation result.

**Theorem 4.1.** *Let  $u_n = \delta_n/\rho < \sqrt{2} - 1$  with  $\rho = \rho_{k_n}$  and  $\delta_n$  from (3.6). Let  $r$  be some number satisfying  $r \geq 1$ . If  $v_m^* > (rs_r \delta_n)^2$  for all  $m \in \mathcal{J}$  where*

$$s_r = \frac{1 + \sqrt{1 + (r^2 + 1)(1 - u_n^2 - 2u_n)}}{1 - u_n^2 - 2u_n},$$

*then it holds on the random set  $A$  defined in Theorem 3.1  $\widehat{v}_m > r^2 \delta_n^2$  for  $m \in \mathcal{J}$  and  $\widehat{v}_m \leq \delta_n^2$  for  $m \notin \mathcal{J}$  and thus,  $\mathcal{J} \subseteq \widehat{\mathcal{J}}_{\mathcal{M}}$  for all  $\mathcal{M} \geq M$ .*

**Remark 4.1.** The result of Theorem 4.1 applied with  $r = 1$  yields the sufficient separation condition: if  $v_m^* > (s_1 \delta_n)^2$ , then, with a high probability,  $\hat{v}_m > \delta_n^2$  for  $m \in \mathcal{J}$  and  $\hat{v}_m \leq \delta_n^2$  for  $m \notin \mathcal{J}$ , and therefore  $\hat{\mathcal{J}}_M = \mathcal{J}$ . For application of this result to the resampling scheme below in this section, we introduced the factor  $r \geq 1$ , which ensures a qualified separation between linear and nonlinear component.

The value  $u_n = \delta_n / \rho$  is small at least if  $n$  is sufficiently large. Hence,  $s_r$  defined in Theorem 4.1 is bounded by a constant depending on  $r$  only and therefore, the threshold  $t^* = (rs_r \delta_n)^2$ , providing with a high probability a correct separation between linear and nonlinear components is of order  $\delta_n^2 \approx (n^{-1} \log n)$ . It can be easily seen that the separation with the prescribed level of the identification error is impossible if the separation distance square is smaller in order than  $n^{-1}$ . Therefore, the procedure provides a near optimal rate of separation within a log-factor.

#### 4.1 Testing the hypothesis about $M$

Here we discuss the problem how the estimators  $\hat{v}_m$  of  $v_m^*$  can be used for selecting the nonlinear component and for testing the hypothesis that the dimensionality  $M$  of the nonlinear component does not exceed the prescribed value  $\mathcal{M}$ . As special cases, for  $\mathcal{M} = 0$  we get the hypothesis that the original model is linear, and for  $\mathcal{M} = 1$ , the hypothesis that the nonlinear component is univariate. Then the effective nonlinear dimension of the model can be estimated by the minimal  $\mathcal{M}$  such that the hypothesis  $M \leq \mathcal{M}$  is not rejected.

The idea of the method is very simple: reject  $H_{\mathcal{M}} : M \leq \mathcal{M}$  if the value  $\hat{v}_{(\mathcal{M}+1)}$  is significantly positive. To formalize the procedure, we have to specify, for a given  $\alpha$ , the critical value  $t_\alpha$  such that the test has the significance level about  $\alpha$ . Suppose that the true model satisfies  $M \leq \mathcal{M}$  and that the values  $v_m^*$  for all  $m \in \mathcal{J}$  exceed the value  $t^* = (rs_r \delta_n)^2$  for some  $r \geq 1$ . Then Theorems 3.1 and 4.1 imply that

- under the null hypothesis  $M \leq \mathcal{M}$ , the index  $(\mathcal{M} + 1)$  corresponds with a high probability to a linear component;
- for  $m \notin \mathcal{J}$ , the distributions of the  $\hat{\beta}_{lm}$ 's and of  $\hat{v}_m$  only weakly depend on the model function  $f$ , see Remark 3.2;
- for every  $m \in \mathcal{J}$ , if  $v_m^*$  is separated from zero by distance of order  $\delta_n^2$ , then the same is true with a high probability for the estimator  $\hat{v}_m$ .

These observations suggest to apply the resampling scheme that mimics only the distri-



bution of the values  $\widehat{v}_{(1)}, \dots, \widehat{v}_{(\mathcal{M}+1)}$ . More precisely, we construct an artificial regression function  $\widetilde{f}_{\mathcal{M}}$  that has exactly  $\mathcal{M}$ -dimensional nonlinear component corresponding to  $m \in \widehat{\mathcal{J}}_{\mathcal{M}}$  and such that all the functionals of type  $\beta_{lm}^*$  constructed for this function  $\widetilde{f}_{\mathcal{M}}$  coincide with the  $\widehat{\beta}_{lm}$ 's, that is,

$$\frac{1}{n} \sum_{i=1}^n \frac{\partial \widetilde{f}_{\mathcal{M}}(X_i)}{\partial x_m} \psi_{lm}(X_{i,m}) = \begin{cases} \widehat{\beta}_{lm} & \text{if } m \in \widehat{\mathcal{J}}_{\mathcal{M}}, \\ 0 & \text{otherwise,} \end{cases} \quad l = 1, \dots, L.$$

The function  $\widetilde{f}_{\mathcal{M}}$  can be selected as a linear combination of the functions  $x_m^l/l$  for  $l = 2, \dots, L+1$ :

$$\widetilde{f}_{\mathcal{M}}(\mathbf{x}) = \sum_{m \in \widehat{\mathcal{J}}_{\mathcal{M}}} \sum_{l=2}^{L+1} c_{lm} x_m^l / l, \quad (4.1)$$

where the coefficients  $c_{lm}$  must fulfill

$$\Psi_m c_m = \widehat{\beta}_m, \quad m \in \widehat{\mathcal{J}}_{\mathcal{M}}. \quad (4.2)$$

Here  $c_m$  (resp.  $\widehat{\beta}_m$ ) denotes the vector in  $\mathbb{R}^L$  with the elements  $c_{lm}$  (resp.  $\widehat{\beta}_{lm}$ ) and  $\Psi_m$  is the  $L \times L$  matrix whose elements are the sums

$$\Psi_{m,l'l} = \frac{1}{n} \sum_{i=1}^n X_{i,m}^l \psi_{l'm}(X_{i,m}), \quad l, l' = 1, \dots, L.$$

We resample from the model

$$\widetilde{Y}_i = \widetilde{f}_{\mathcal{M}}(X_i) + \widehat{\sigma}_i \widetilde{\varepsilon}_i,$$

where  $\widetilde{\varepsilon}_i$  are i.i.d. standard normal. The variances  $\widehat{\sigma}_i^2$  either fulfill  $\widehat{\sigma}_i^2 = \widehat{\sigma}^2$  for the variance estimator  $\widehat{\sigma}^2$  from Section 2.4, or they simply are defined by  $\widehat{\sigma}_i^2 = \widehat{\varepsilon}_i^2$ . The first proposal suits well the case of a homogeneous noise, and the second one is similar to the wild bootstrap idea and should be used if the assumption of noise homogeneity is questionable.

The recommended estimator of the critical value can be computed by using the following simulation procedure:

1. For each  $i = 1, \dots, n$ , generate  $\widetilde{Y}_i = \widetilde{f}_{\mathcal{M}}(X_i) + \widehat{\sigma}_i \widetilde{\varepsilon}_i$ , where  $\widetilde{\varepsilon}_i$  is sampled randomly from the standard normal law.
2. Use the data set  $\{\widetilde{Y}_i, X_i : i = 1, \dots, n\}$  to estimate gradient projections  $\beta_{lm}^*$  with estimator (2.3) based on gradient estimator (2.1) with the last step bandwidth  $\mathbf{b} = \mathbf{b}^{(k_n)}$ . Denote the resulting estimator by  $\widetilde{\beta}_{lm}$ . Compute  $\widetilde{v}_m = \sum_{l=1}^L |\widetilde{\beta}_{lm}|^2$  for every  $m = 1, \dots, d$  and the statistic  $\widetilde{T}_{\mathcal{M}}$ , that is  $\widetilde{v}_{(\mathcal{M}+1)}$ .

3. Define  $t_\alpha$  as the  $(1-\alpha)$ -quantile of the empirical distribution of  $\tilde{T}_M$  that is obtained by repeating steps 1-2 many times.

**Theorem 4.2.** *Let Assumptions 1–5 hold and  $\min_{m \in \mathcal{J}} v_m^* \geq (rs_r \delta_n)^2$  with  $r = s_1$ . If  $\mathcal{M} = M$ , then*

$$\mathbf{P}(H_{\mathcal{M}} \text{ is rejected}) \leq \alpha + 3(k_n + 1)/n.$$

## 4.2 Identification of the nonlinear component

Here we describe how the effective nonlinear dimension  $M$  and the index set  $\mathcal{J}$  corresponding to the nonlinear component can be estimated using the above testing procedure. Let some positive  $\alpha < 1$  be fixed. Starting with  $\mathcal{M} = 0$ , we consider the model with  $\mathcal{M}$ -dimensional nonlinear component due to (4.1) and (4.2) and test the hypothesis  $M \leq \mathcal{M}$  at the level  $\alpha$  as described in Section 4.1. Terminate if the hypothesis  $M \leq \mathcal{M}$  is not rejected, otherwise increase  $\mathcal{M}$  by one. Finally we set  $\widehat{\mathcal{M}} =$  “the first nonrejected  $\mathcal{M}$ ” and  $\widehat{\mathcal{J}} = \widehat{\mathcal{J}}_{\widehat{\mathcal{M}}}$ .

**Theorem 4.3.** *Let Assumptions 1 through 5 hold and  $\min_{m \in \mathcal{J}} v_m^* \geq (rs_r \delta_n)^2$  with  $r = s_1$ . Then*

$$\mathbf{P}(\widehat{\mathcal{J}} \neq \mathcal{J}) \leq \alpha + 3(k_n + M)/n.$$

**Remark 4.2.** It can be easily checked that the results of Theorems 4.2 and 4.3 continue to hold even if the test level  $\alpha$  depends on  $n$  and goes to zero as  $n$  grows. In particular, one can take  $\alpha = n^{-a}$  with  $a < 1/2$ . With such a choice, our method leads to a consistent estimation of the set  $\mathcal{J}$ .

## 4.3 Estimation and inference for the linear component

The method described above allows to classify the regressors into linear and nonlinear. Moreover, the result of classification is correct with a dominating probability provided the sample size is large enough. The impact of every linear regression variable in the model function is characterized by the corresponding slope coefficient  $\theta_m$  from (1.1). Here we discuss how these slope coefficients can be estimated. We use again the observation that  $\partial f / \partial x_m \equiv \theta_m$  for every  $m \notin \mathcal{J}$ . Therefore, the sum

$$\widehat{\theta}_m = \frac{1}{n} \sum_{i=1}^n \widehat{\nabla} f_m(X_i) \tag{4.3}$$

is a reasonable estimator of  $\theta_m$ . Here  $\widehat{\nabla}f(X_i)$  is the gradient estimator obtained at the last step of the algorithm. Our next result claims that  $\widehat{\theta}_m$  from (4.3) estimates the true value  $\theta_m$  with the root-n accuracy and that it can be very well approximated by a Gaussian random variable. This result can be viewed as an application of Theorem 3.1 for  $m \notin \mathcal{J}$  and  $\psi_{lm} \equiv 1$ .

**Theorem 4.4.** *Let Assumptions 1 through 5 hold. Then for every  $m \notin \mathcal{J}$ , there exists a Gaussian zero mean random variable  $\gamma_m^*$  which is defined as a linear combination of the errors  $\varepsilon_i$  with deterministic coefficients, depending on the “ideal” bandwidth  $\mathbf{b}^* = \mathbf{b}^{*(k_n)}$ , the design  $X_1, \dots, X_n$ , the basis functions  $\psi_{lm}(\cdot)$  and the kernel  $K$  only, and such that*

$$\mathbf{E}|\gamma_m^*|^2 \leq 2\sigma^2 C_V^2 C_K,$$

and on the random set  $A$  from Theorem 3.1 with  $\mathbf{P}(A) \geq 1 - 3k_n/n$ , it holds

$$\max_{m \notin \mathcal{J}} \left| \widehat{\theta}_m - \theta_m - n^{-1/2} \gamma_m^* \right| \leq C_1 (\sigma^2 n^{-1} L \log n)^{2/3},$$

where  $C_1 = C_1(d, M, C_g, C_V, C_K, C_{K'}, \bar{\psi})$ .

**Remark 4.3.** The above estimator  $\widehat{\theta}_m$  can be slightly refined by explicitly using the estimated structural information about the model. Namely, an application of the anisotropic bandwidth  $\widehat{\mathbf{b}} = \text{diag}(\widehat{b}_1, \dots, \widehat{b}_d)$  with  $\widehat{b}_m$  from the last iteration for  $m \in \widehat{\mathcal{J}}$  and  $\widehat{b}_m = \infty$  for  $m \notin \widehat{\mathcal{J}}$  leads under condition of the correct identification to the classical partially linear estimator for the case with known  $\mathcal{J}$ , see e.g. Härdle, Liang and Gao (1999).

**Remark 4.4. (Selecting significant regressors)** The procedure in Section 4.2 can be also used for identifying the significant components. All the regressors entering in the nonlinear component are automatically significant. The linear regressors can be further analyzed for significance. Theorem 4.4 claims that the normalized estimation error  $\sqrt{n}(\widehat{\theta}_m - \theta_m)$  is asymptotically normal. Moreover, the asymptotic variance of  $\widehat{\theta}_m$  can be easily estimated. Indeed,  $\widehat{\theta}_m$  is a linear combination of the observations  $Y_i$  with the known coefficients  $c_{im}$ , that is,  $\widehat{\theta}_m = \sum_i c_{im} Y_m$ . Then  $\widehat{\sigma}_m^2 = \sigma^2 \sum_i c_{im}^2$  is an estimator of  $\text{Var}(\widehat{\theta}_m)$ . The search of significant regressors can be done by the rule  $|\widehat{\theta}_m|^2 \geq \lambda^2 \widehat{\sigma}_m^2$  for some  $\lambda > 0$ , see illustration of this procedure in Section 5. We skip the further discussion for the reasons of space.

## 5 Simulated and real data results

In this section we illustrate the performance of the proposed method on some simulated examples and give a real data application. With the simulated examples we aim to illustrate

how the performance of the proposed method depends on the sample size  $n$ , dimension  $d$  of the model, the dimensionality of the nonlinear component  $M$  and the noise variance  $\sigma^2$ . We especially focus on the component classification results: identification of the nonlinear variables and selection of the significant variables. We also demonstrate how the quality of estimation of the nonlinear components improves during iteration.

In our simulation study we apply the modified procedure (see Remark 2.4) with the following parameter setting:

$$\rho_1 = 1, \quad \rho_{\min} = n^{-1/3}, \quad a_\rho = e^{-1/6}, \quad \eta = 0.2, \quad a_h = a_\rho^{-1/2}.$$

The initial bandwidth  $h_1$  is selected from the condition  $\#\{i : M_h(X_i) \geq d + 1\} \geq n/2$ , where  $M_h(x)$  stand for the number of the design points  $X_i$  in the ball of radius  $h$  and center  $x$ . This condition ensures that for at least a half of the design points the local gradient estimator is well defined. This setting leads to the number of iterations  $k(n) \approx \frac{\log(\rho_1/\rho_{\min})}{\log a_\rho} \approx 2 \log n$ .

The procedure utilizes the kernel  $K(|x|^2) = (1 - |x|^2)_+^2$ . For every  $m \leq d$ , the basis system  $\{\psi_{1m}(x_m), \dots, \psi_{Lm}(x_m)\}$  is obtained using polynomials of  $x_m$  of degree from one to  $L$  satisfying the conditions  $\sum_{i=1}^n w_i \psi_{lm}(X_{i,m}) \psi_{l'm}(X_{i,m}) / \sum_{i=1}^n w_i = \delta_{ll'}$  and  $\sum_{i=1}^n \psi_{lm}(X_{i,m}) w_i = 0$  where  $w_i = w_i^{(k)} = \lambda_{\min}^{1/2}(V_i^{(k)})$  for  $k$ th iteration with  $k \geq 1$ . We apply  $L = 6$ .

In our simulation study we consider the model  $Y_i = \theta^\top X_i + g(X_{i,d-M+1}, \dots, X_{i,d}) + \varepsilon_i$  for  $M$  between 1 and 3. The vector  $\theta$  is taken of the form  $\theta = (1, 2, 3, 4, 0, \dots, 0)^\top$ . The link function  $g$  is  $g(u) = g_1(u) = e^u + e^{-u}$  for  $M = 1$ ,  $g(u_1, u_2) = g_1(u_1)g_1(u_2)$  for  $M = 2$  and  $g(u_1, u_2, u_3) = g_1(u_1)g_1(u_2)g_1(u_3)$  for  $M = 3$ . The dimension  $d$  is taken  $4 + M$  or larger. The errors  $\varepsilon_i$  are i.i.d. normal with parameters  $(0, \sigma^2)$  for  $\sigma^2 = 0.1$ . The design  $X_1, \dots, X_n$  is modelled randomly so that each  $X_i$  follows  $\text{Norm}(0.2, 0.8^2)$ -distribution, restricted to the  $[-1, 1]^d$ -cube. The experiments were done for sample size  $n = 100, 200, 400$ . The results displaying the quality of estimation by the iterative algorithm are summarized in Tables 1 for  $M = 1$  and in Table 2 for  $M = 2$ . We display the mean losses  $|\widehat{v}_m|$  for one linear regressor and  $|\widehat{v}_m^{(k)} - v_m^*|$  for nonlinear regressors where  $v_m^* = |\beta_m^*|^2$  and  $\beta_m^*$  is the vector with the components  $\beta_{lm}^* = \sum_{i=1}^n w_i^{(k)} \nabla f_m(X_{i,m}) \psi_{lm}(X_{i,m}) / \sum_{i=1}^n w_i^{(k)}$ .

It is interesting to observe that the quality of estimating the linear regressor  $x_1$  improves with growing dimension  $d$ .

In Table 2 we demonstrate in addition how the error of estimation depends on the noise variance  $\sigma^2$ . One can see that the estimation risk for the nonlinear components only

Table 1: Case  $M = 1$ : mean loss  $|\widehat{v}_m - v_m^*|$  for the nonlinear regressor for the first, second, fourth, and final iteration and final losses  $|\widehat{v}_1|$  for the first linear regressor. Results are obtained from  $N = 250$  simulations. The interquartile range of the losses is given in parentheses.

$d$	$n$	nonlinear regressor			linear regressor $x_1$
		1st	4th	final	final
5	100	0.9580	0.6656	0.3069	0.0139
		(0.1865)	(0.1546)	(0.2400)	(0.0113)
5	200	0.9395	0.7711	0.2424	0.0072
		(0.1378)	(0.1300)	(0.2024)	(0.0057)
6	200	0.9432	0.7207	0.1641	0.0018
		(0.1231)	(0.1067)	(0.1766)	(0.0016)
8	200	0.9362	0.6703	0.2232	0.0006
		(0.1253)	(0.1003)	(0.1797)	(0.0005)
10	100	0.9574	0.6743	0.5822	0.0005
		(0.2064)	(0.1526)	(0.2756)	(0.0004)
10	200	0.9406	0.6777	0.3690	0.0002
		(0.1522)	(0.1202)	(0.2213)	(0.0002)
10	400	0.9348	0.7217	0.2316	0.0001
		(0.0925)	(0.0838)	(0.1399)	(0.0001)

slightly increases with  $\sigma$  while it is essentially proportional to  $\sigma$  for the linear one. An explanation might be that the estimation error for the nonlinear components is mostly due to the nonparametric bias which disappears in the linear components during iteration process by structural adaptation.

Table 2: Case  $M = 2$ : mean loss  $|\widehat{v}_m - v_m^*|$  for the nonlinear regressors for the first, second, fourth, and final iteration and final losses  $|\widehat{v}_1|$  for the first linear regressor. Results are obtained from  $N = 250$  simulations. The interquartile range of the losses is given in parentheses.

$d$	$n$	$\sigma^2$	1st nonlinear regressor			2nd nonlinear regressor			linear regressor $x_1$
			1st	4th	final	1st	4th	final	final
6	200	0.1	4.6117	3.7349	0.4763	4.6337	3.7576	0.4473	0.0081
			(0.6646)	(0.5028)	(0.5211)	(0.6257)	(0.5402)	(0.4785)	(0.0063)
8	200	0.1	4.6397	3.4423	0.4244	4.5942	3.4085	0.4058	0.0025
			(0.6683)	(0.5431)	(0.4108)	(0.6646)	(0.4621)	(0.4607)	(0.0019)
10	100	0.1	4.6338	3.1450	0.7573	4.6862	3.1642	0.7089	0.0043
			(0.8840)	(0.7307)	(0.5302)	(1.0155)	(0.7312)	(0.4938)	(0.0032)
10	200	0.1	4.5537	3.2806	0.5812	4.5904	3.2917	0.5489	0.0011
			(0.7458)	(0.5065)	(0.3404)	(0.7649)	(0.5875)	(0.4014)	(0.0010)
10	400	0.1	4.5198	3.5276	0.4457	4.5584	3.5594	0.4319	0.0004
			(0.4850)	(0.3566)	(0.3121)	(0.4168)	(0.3562)	(0.3023)	(0.0003)
10	400	0.2	4.5198	3.5284	0.4403	4.5584	3.5602	0.4325	0.0007
			(0.4850)	(0.3483)	(0.3949)	(0.4167)	(0.3642)	(0.3948)	(0.0006)
10	400	0.4	4.5198	3.5297	0.4637	4.5584	3.5615	0.4666	0.0017
			(0.4850)	(0.3316)	(0.4891)	(0.4167)	(0.3727)	(0.5029)	(0.0013)

The next figure illustrates the result of Theorem 4.1 about separation between linear and nonlinear component. Let  $\mathcal{L}_N(\xi)$  denote the empirical distribution of the random variable  $\xi$  based on its sample of size  $N$ . A good separation between linear and nonlinear components means that the functions  $\mathcal{L}_N(\widehat{v}_m)$  for every  $m \in \mathcal{J}$  and  $1 - \mathcal{L}_N(\widehat{v}_m)$  for  $m \notin \mathcal{J}$  have non-overlapping support.

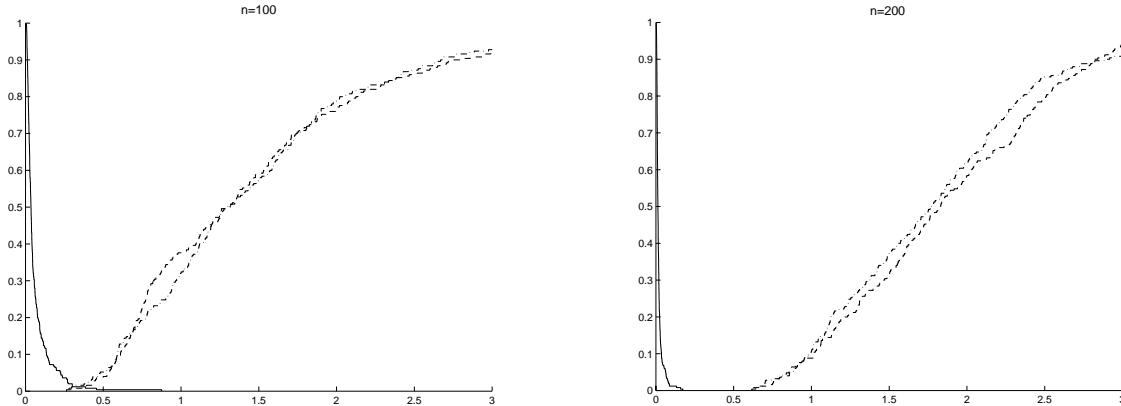


Figure 1: Case  $M = 2$ ,  $d = 6$  :  $\mathcal{L}_N(\widehat{v}_5)$ ,  $\mathcal{L}_N(\widehat{v}_6)$  (dotted lines) and  $1 - \mathcal{L}_N(\max_{m=1,\dots,4} \widehat{v}_m)$  (solid line) for  $n = 100, 200$  from 250 realizations.

We observe a very good separability for  $n = 100$  and a possibility of perfect separation for  $n = 200$ .

Next we illustrate how the quality of estimation of the linear component improves with the sample size. Figure 2 shows box-plots of the estimation errors  $n^{1/2} \|\widehat{\theta} - \theta^*\|$  of the linear component after the final iteration for  $d = 6$ ,  $M = 2$  and different sample sizes  $n$ .

Table 3 illustrates the performance of the test of the hypothesis  $M \leq \mathcal{M}$  and the quality of the classification rule from Sections 4.2 and 4.3 for different  $M$ ,  $d$  and  $n$ . In this table we present the fraction of wrong classifications for every of nonlinear regressors and for the whole model.

One can observe once again that the results (the fraction of wrong classifications) improve as the dimensionality  $d$  grows. This can be explained by the fact that the distribution of the test statistic used for classification will be more and more degenerated with growing dimension  $d$ .

Another observation is that for  $M = 3$ , the procedure requires some minimal sample size to start selecting all the three nonlinear components. For  $n = 100$  we obtain for almost all the cases  $\widehat{M} < M$ . For  $n = 200$  and  $d = 7$  we correctly classify in only about 30% cases but for  $d = 10$  the fraction of wrong classifying is already under control.

Figure 3 illustrates the quality of estimation of the noise variance  $\sigma^2$  by  $\widehat{\sigma}^2$  for one example

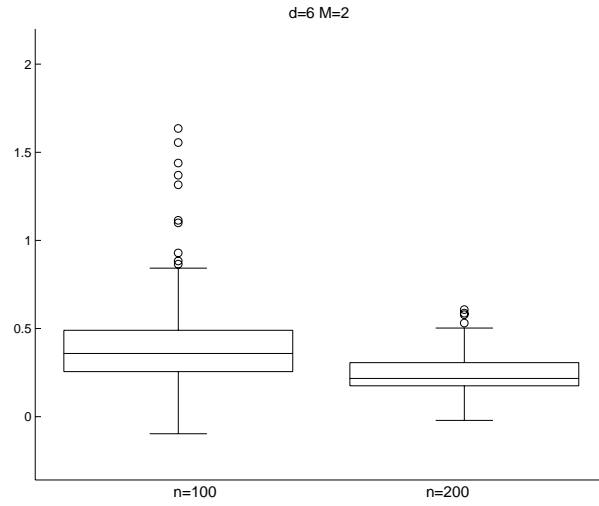


Figure 2: Case  $M = 2$  : box-plots of the estimation errors  $n^{1/2}\|\hat{\theta} - \theta^*\|$  of the linear component after the final iteration for  $d = 6$ . Results are obtained from  $N = 250$  simulations.

Table 3: Fraction of wrong classifications for every nonlinear regressor and for the whole model. Results are obtained from  $N = 250$  simulations and 500 bootstrap replications.

$M$	$d$	$n$	1st n.c.	2nd n.c.	3rd n.c.	$\{\hat{\mathcal{J}} \neq \mathcal{J}\}$
1	5	100	0.152	–	–	0.18
		200		–	–	
		400	0	–	–	
1	10	200	0	–	–	0.
		400	0	–	–	0.00
1	20	400	0	–	–	0.
2	6	100	0.268	0.308	–	0.38
		200	0.056	0.048	–	0.1
		400	0.004	0.004	–	0.024
2	10	200	0	0	–	0.008
		400	0	0	–	0.0
2	20	400	0	0	–	0
3	7	100	0.976	0.96	0.964	0.992
		200	0.62	0.656	0.656	0.748
		400	0.076	0.06	0.072	0.1
3	10	200	0.004	0	0.004	0.004
		400	0	0	0	0

with  $d = 6$ ,  $M = 2$  and different sample size  $n$ . The results are in agreement with the root-n consistency of the estimator  $\hat{\sigma}^2$ .

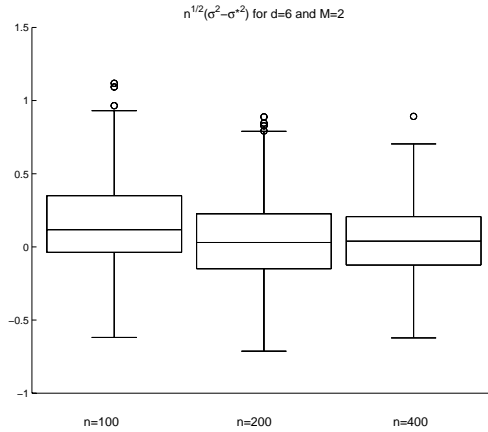


Figure 3: Box-plots of the estimation errors  $n^{1/2}||\hat{\sigma}^2 - \sigma^{*2}||$  for  $d = 6$ ,  $M = 2$  and different sample size  $n$ .

## 5.1 A real data example

This section presents an application of the procedure to a real data set. We consider the example from Sperlich (1998) and Härdle, Spokoiny and Sperlich (2001) where a subsample of the Socio-Economic Panel of Germany from 1992 was studied. The target of analysis is the weekly number of working hours,  $Y_i$ , of 607 women with job and living together with a partner. The following explanatory variables were used: the age of woman, between 25 and 60,  $X_1$ ; her earning per hour,  $X_2$ ; the prestige index of her kind of profession (Treiman prestige index),  $X_3$ ; the monthly rent or redemption for their apartment or house,  $X_4$ ; the monthly net income of their partner,  $X_5$ ; the number of years of education,  $X_6$ ; the unemployment rate at the particular tract they live in,  $X_7$ ; and the number of children younger than 16 years,  $X_8$ .

The estimates  $\hat{v}_m$  obtained by our estimation procedure are given in Table 4. We also got the estimate  $\hat{\sigma}^2 = 0.736$  for error variance  $\sigma^2$ .

Table 4: Estimates  $\hat{v}_m$  of  $v_m^*$ .

$\hat{v}_1$	$\hat{v}_2$	$\hat{v}_3$	$\hat{v}_4$	$\hat{v}_5$	$\hat{v}_6$	$\hat{v}_7$	$\hat{v}_8$
0.05259	0.00729	0.00441	0.00012	0.00060	0.00142	0.00015	0.00875

Next we identify the linear component starting with  $\mathcal{M} = 0$  as described in Section 4.2.



Table 5 gives the p-values  $PV_{\mathcal{M}}$  for each test  $H_{\mathcal{M}}$ , which are obtained during the bootstrap procedure, defined as:

$$PV_{\mathcal{M}} = \frac{1}{B} \sum_{b=1}^B 1_{\{\hat{v}_{(\mathcal{M}+1)}^{(b)} > \hat{v}_{(\mathcal{M}+1)}\}}$$

The first three hypotheses  $H_0$ ,  $H_1$  and  $H_2$  are rejected at 10% level, and there is clearly no rejection of  $H_3$ . So, for the considered model, the nonlinear dimension is estimated as three and the nonlinear variables are  $X_1$ ,  $X_2$ , and  $X_3$ . Our linear/nonlinear variable classification results coincide with those from Härdle, Sperlich and Spokoiny (2001), but with quite different p-values: in our results  $X_1$  (age) is the most nonlinear and  $X_2$  (earning per hour) is the least nonlinear variable among the three, while in Härdle, Sperlich and Spokoiny (2001) the situation is reversed. Note that while their identification was made under the assumption of additive model structure, our results are obtained for a general situation when such additive structure is not required.

Table 5: p-values for consecutive tests

$\mathcal{M}$	$\hat{v}_{(\mathcal{M}+1)}$	p-values
0	$\hat{v}_1$	0.003996
1	$\hat{v}_2$	0.086913
2	$\hat{v}_3$	0.01998
3	$\hat{v}_5$	0.47153

## 6 Conclusion and outlook

The paper has introduced a new method of exploring a partially linear model based on the idea of structural adaptation. The method applies under mild assumptions on the underlying regression function and the regression design. The procedure is fully adaptive and does not require any prior information. The results claim that the proposed procedure with a high probability correctly identifies the nonlinear component and estimates the linear component with the optimal rate  $n^{-1/2}$  provided that the dimension of the nonlinear component is not larger than 3. The simulation results demonstrate an excellent performance of the procedure for all considered situations. An important feature of the method is that it is very stable with respect to high dimensionality and for a non-regular design.

**Non-Gaussian or heterogeneous noise.** The method and results can be easily extended to models with homogeneous non-Gaussian noise satisfying some exponential moment con-

ditions. Another interesting issue is applicability of the method for a general heterogeneous or dependent noise, in particular, to time series models and financial data. We leave these extensions for further research.

**The case with  $M \geq 4$ .** The method continues to apply even if  $M \geq 4$  and iterations would lead to improvement of the bias. However, the bound for the bias of order  $(n^{-1} \log n)^{2/3}$  can be achieved only for  $M \leq 3$ . For larger  $M$ , the bias will be of order  $n^{-1/2}$  or bigger and the procedure does not provide root-n consistent estimation of the functionals  $\beta_{lm}^*$ . So, if the hypothesis  $M \leq 3$  is rejected, then we recommend to apply for the choice of  $M$  some model selection criteria like cross-validation or Mallows  $C_p$ .

**Data-driven choice of parameter  $L$ .** The method depends upon the parameter  $L$  describing the number of basis functions for every regressor. In the univariate case, either an  $n$ -dependent or data-driven choice of such a parameter is usually applied, see Hart (1997) or Spokoiny (2001) and references therein. An adaptive choice of  $L$  in the considered problem is an interesting question for further research.

**Semiparametrically efficient estimation of the linear component.** Due to the result of Theorem 4.4, the proposed estimator of the parameter  $\theta$  is root-n consistent and asymptotically normal. However, it is unlikely that this or the refined estimator of  $\theta$  from Section 4.3 is semiparametrically efficient in the sense of minimization of the asymptotic variance, see e.g. Bickel et al. (1998). A modification of the method leading to the semiparametrically efficient estimation of linear part will be discussed elsewhere.

**Estimation of the nonlinear component.** After the nonlinear component is identified, it can be estimated using the standard methods of nonparametric statistics. Actually, the algorithm gives an estimator of the whole function  $f$  and of the linear component, so that the nonlinear component can be extracted as well. This estimator corresponds to the local linear smoothing of the nonparametric  $M$ -dimensional function with the bandwidth about  $h\rho \approx \rho_{\min}$ , and may not achieve the best rate. To improve the quality of estimation, one can apply the classical cross-validation technique for selecting the bandwidth in the direction of the nonparametric component.

**Discrete and categorical data.** Note that the assumption of linearity is meaningful for discrete or categorical variables as well. It means that the influence of the corresponding regressor is independent of the other variables and therefore, at least in the binary case, can be modelled linearly. Moreover, the procedure easily applies for the situation with discrete data without any change.

## 7 Appendix

Here we collect the proofs of the main results. For the ease of exposition, we consider only the main procedure (without weights) and only the case of  $\eta = 0$ . The general case can be considered in the same way.

### 7.1 One-step improvement

Suppose that we are given some fixed numbers  $h$  and  $\rho$  (which mean the current values  $h_k$  and  $\rho_k$ ) and a vector  $\mathbf{v} = (v_1, \dots, v_d)^\top \in \mathbb{R}^d$  which can be viewed as an approximation of  $\mathbf{v}^* = (v_1^*, \dots, v_d^*)^\top$  obtained at the previous step of the algorithm. Set also

$$b_m = h (1 + \rho^{-2} v_m)^{-1/2}, \quad m = 1, \dots, d, \quad (7.1)$$

and define  $\mathbf{b} = \text{diag}(b_1, \dots, b_d)$ . Define also  $\widehat{f}(X_i)$ ,  $\widehat{\nabla} f(X_i)$  and  $\widehat{\beta}_{lm}$  by (2.1) and (2.3) for all  $l = 1, \dots, L$  and  $m = 1, \dots, d$  with the just defined bandwidth  $\mathbf{b}$ . We aim to evaluate the estimation errors  $\widehat{\beta}_{lm} - \beta_{lm}^*$ . To describe the results, we introduce the shrinking factors  $P_{\rho,m} = (1 + \rho^{-2} v_m^*)^{-1/2}$  and define

$$U_m = P_{\rho,m}^2 (1 + \rho^{-2} v_m) = (1 + \rho^{-2} v_m^*)^{-1} (1 + \rho^{-2} v_m)$$

and similarly  $U_m^* = P_{\rho,m}^2 (1 + \rho^{-2} v_m^*) = 1$ . Clearly the vector  $U = (U_1, \dots, U_d)^\top \in \mathbb{R}^d$  uniquely describes  $\mathbf{v}$ , so that we consider later in this section  $v = v(U)$  and similarly  $\widehat{\beta}_{lm} = \widehat{\beta}_{lm}(U)$  for the functionals  $\widehat{\beta}_{lm}$  in (2.3). Let  $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_d)^\top$  be a vector in  $\mathbb{R}^d$  with entries  $\alpha_m \in (0, 1)$ . Define

$$\mathcal{U}_{\boldsymbol{\alpha}} = \{U = (U_1, \dots, U_d)^\top \in \mathbb{R}^d : |U_m - 1| \leq \alpha_m, m = 1, \dots, d\}.$$

We also define  $\alpha^* = \max_{m=1, \dots, d} \alpha_m$ .

**Proposition 7.1.** *Let Assumptions 1 through 5 hold. Let  $\beta_{lm}(U) = \mathbf{E} \widehat{\beta}_{lm}(U)$ . Then*

$$\sup_{U \in \mathcal{U}_{\boldsymbol{\alpha}}} \sum_{m=1}^d \sum_{l=1}^L |P_{\rho,m} \{\beta_{lm}(U) - \beta_{lm}^*\}|^2 \leq \left( \frac{C_g C_V^{1/2}}{1 - \alpha^*} \rho^2 h \right)^2$$

and, for every  $l = 1, \dots, L$  and  $m = 1, \dots, d$ , there exists a zero mean Gaussian random variable  $\xi_{lm}$  defined as a linear combination of the errors  $\varepsilon_i$  with deterministic coefficients, which depend on  $\mathbf{v}^*$ , the design  $\{X_i\}$ , the basis functions  $\psi_{lm}(\cdot)$ , and the kernel  $K$  only, and such that

$$\max_{m,l} \mathbf{E} \xi_{lm}^2 \leq 2\sigma^2 C_V^2 C_K \quad (7.2)$$

and

$$\mathbf{P} \left( \max_{m,l} \sup_{U \in \mathcal{U}_\alpha} \left| P_{\rho,m} \{ \widehat{\beta}_{lm}(U) - \beta_{lm}(U) \} - \frac{\xi_{lm}}{h\sqrt{n}} \right| > \frac{\sigma \bar{\psi} C_{\alpha,n} |\alpha|}{h\sqrt{n}} \right) \leq \frac{2}{n},$$

where the maximum is taken over  $m = 1, \dots, d$  and  $l = 1, \dots, L$ ,  $\bar{\psi} = \max_{i,l,m} |\psi_{lm}(X_i)|$  and

$$C_{\alpha,n} = \left( \frac{\sqrt{2} C_V C_{K'}}{(1 - \alpha^*)^{3/2}} + \frac{2^{3/2} C_V^2 C_{K'} C_K}{(1 - \alpha^*)^{5/2}} \right) \left( 2 + \sqrt{2 \log(ndL) + d \log(4n)} \right).$$

Let  $\beta_m^*$  denote, as in Proposition 3.1, an  $L$ -vector with the components  $\beta_{lm}^*$  and  $\widehat{\beta}_m = \widehat{\beta}_m(U)$  its estimator with the components  $\widehat{\beta}_{lm}(U)$ .

**Corollary 7.1.** *Let  $z_n = (1 + 2 \log(nd) + 2 \log \log(nd))^{1/2}$  and*

$$\delta = \frac{C_g C_V^{1/2}}{1 - \alpha^*} h \rho^2 + \frac{\sqrt{2L} \sigma C_V C_K^{1/2} z_n}{h\sqrt{n}} + \frac{\sqrt{L} \sigma \bar{\psi} C_{\alpha,n} |\alpha|}{h\sqrt{n}}. \quad (7.3)$$

Then under the conditions of Proposition 7.1 it holds

$$\mathbf{P} \left( \max_{m=1,\dots,d} \sup_{U \in \mathcal{U}_\alpha} \left| P_{\rho,m} \left( \widehat{\beta}_m(U) - \beta_m^* \right) \right| > \delta \right) \leq 3/n.$$

The corollary helps bound the estimation error  $P_{\rho,m}^2 (\widehat{v}_m(U) - v_m^*)$ .

**Proposition 7.2.** *Under the conditions of Proposition 7.1,*

$$\mathbf{P} \left( \sup_{U \in \mathcal{U}_\alpha} \left| P_{\rho,m}^2 (\widehat{v}_m(U) - v_m^*) \right| \leq \delta^2 + 2\delta\tau_m \text{ for all } m = 1, \dots, d \right) \geq 1 - 3/n$$

where  $\tau_m = \rho \sqrt{v_m^*} (\rho^2 + v_m^*)^{-1/2} \leq \min \{ \rho, \sqrt{v_m^*} \}$ .

## 7.2 Proof of Proposition 7.1

Denote by  $P_\rho$  the diagonal  $d \times d$ -matrix with the diagonal entries  $P_{\rho,m}$ , that is,  $P_\rho = \text{diag}\{P_{\rho,1}, \dots, P_{\rho,d}\}$ . Similarly, for  $U = (U_1, \dots, U_d)^\top \in \mathbb{R}^d$ , define  $D_U = \text{diag}\{U_1, \dots, U_d\}$ . Next, for every  $i, j \leq n$ , define  $Z_{ij} = h^{-1} P_\rho^{-1} (X_j - X_i)$ ,  $K_{ij}(U) = K(Z_{ij}^\top D_U Z_{ij})$

$$\begin{aligned} \mathcal{V}_i(U) &= \sum_{j=1}^n \begin{pmatrix} 1 \\ Z_{ij} \end{pmatrix} \begin{pmatrix} 1 \\ Z_{ij} \end{pmatrix}^\top K_{ij}(U), \\ \widehat{\mathcal{S}}_i(U) &= h^{-1} \mathcal{V}_i(U)^{-1} \sum_{j=1}^n \begin{pmatrix} 1 \\ Z_{ij} \end{pmatrix} Y_j K_{ij}(U). \end{aligned}$$

It is easy to check that for the  $(m+1)$ th component  $\widehat{\mathcal{S}}_{i,m}(U)$  of  $\widehat{\mathcal{S}}_i(U)$  it holds  $\widehat{\mathcal{S}}_{i,m}(U) = P_{\rho,m} \widehat{\nabla} f_m(X_i)$  and hence,

$$P_{\rho,m} \widehat{\beta}_{lm}(U) = n^{-1} \sum_{i=1}^n \widehat{\mathcal{S}}_{i,m}(U) \psi_{lm}(X_{i,m}).$$

The model equation (1.1) implies  $\widehat{s}_i(U) = s_i(U) + \zeta_i(U)$  with

$$\begin{aligned} s_i(U) &= h^{-1} \mathcal{V}_i(U)^{-1} \sum_{j=1}^n \begin{pmatrix} 1 \\ Z_{ij} \end{pmatrix} f(X_j) K_{ij}(U), \\ \zeta_i(U) &= h^{-1} \mathcal{V}_i(U)^{-1} \sum_{j=1}^n \begin{pmatrix} 1 \\ Z_{ij} \end{pmatrix} \varepsilon_j K_{ij}(U). \end{aligned}$$

This yields, for each coordinate  $m = 1, \dots, d$ ,

$$\begin{aligned} P_{\rho,m} \{ \mathbf{E} \widehat{\beta}_{lm}(U) - \beta_{lm}^* \} &= \frac{1}{n} \sum_{i=1}^n \{ s_{i,m}(U) - P_{\rho,m} \nabla f_m(X_i) \} \psi_{lm}(X_{i,m}), \\ P_{\rho,m} \{ \widehat{\beta}_{lm}(U) - \mathbf{E} \widehat{\beta}_{lm} \} &= \frac{1}{n} \sum_{i=1}^n \zeta_{i,m}(U) \psi_{lm}(X_{i,m}). \end{aligned}$$

Clearly  $\zeta_{lm}(U) := n^{-1} \sum_{i=1}^n \zeta_{i,m}(U) \psi_{lm}(X_{i,m})$  is for every  $U$  a linear combination of the Gaussian errors  $\varepsilon_i$  and therefore it is also a Gaussian vector in  $\mathbb{R}^d$ .

Define  $\mathcal{E}_d$  is the projection from  $\mathbb{R}^{d+1}$  onto  $\mathbb{R}^d$  dropping the zero coordinate:  $\mathcal{E}_d(x_0, \dots, x_d)^\top = (x_1, \dots, x_d)^\top$ . It is easy to see that the following three statements imply the claimed result:

$$\sup_{U \in \mathcal{U}_\alpha} |\mathcal{E}_d s_i(U) - P_\rho \nabla f(X_i)| \leq \frac{C_g C_V^{1/2}}{1 - \alpha^*} h \rho^2, \quad i = 1, \dots, n, \quad (7.4)$$

$$\mathbf{P} \left( \max_{l,m} \sup_{U \in \mathcal{U}_\alpha} |\zeta_{lm}(U) - \zeta_{lm}(U^*)| > \frac{\sigma C_{\alpha,n} |\boldsymbol{\alpha}|}{h \sqrt{n}} \right) \leq 2/n, \quad (7.5)$$

$$\max_{l,m} \mathbf{E} |\zeta_{lm}(U^*)|^2 \leq \frac{2\sigma^2 C_V^2 C_K}{h^2 n}. \quad (7.6)$$

where the maximum is taken over  $l = 1, \dots, L$  and  $m = 1, \dots, d$ . Indeed, the last two statements of the proposition directly follows from (7.5) and (7.6) for  $\xi_{lm} = h\sqrt{n} \zeta_{lm}(U^*)$ .

Next, (7.4) implies

$$n^{-1} \sum_{i=1}^n \sum_{m=1}^d |s_{i,m}(U) - P_{\rho,m} \nabla f_m(X_i)|^2 \leq \left( \frac{C_g C_V^{1/2}}{1 - \alpha^*} h \rho^2 \right)^2.$$

Since the vectors  $\psi_{lm} \in \mathbb{R}^n$  are orthonormal for different  $l$ , it follows for the Bessel inequality for every  $m \leq d$

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n |s_{i,m}(U) - P_{\rho,m} \nabla f_m(X_i)|^2 &\geq \sum_{l=1}^L \left| \frac{1}{n} \sum_{i=1}^n (s_{i,m}(U) - P_{\rho,m} \nabla f_m(X_i)) \psi_{lm}(X_{i,m}) \right|^2 \\ &= \sum_{l=1}^L P_{\rho,m}^2 \left( \mathbf{E} \widehat{\beta}_{lm}(U) - \beta_{lm}^* \right)^2 \end{aligned}$$

and thus,

$$\sum_{m=1}^d \sum_{l=1}^L P_{\rho,m}^2 \left( \mathbf{E} \widehat{\beta}_{lm}(U) - \beta_{lm}^* \right)^2 \leq \left( \frac{C_g C_V^{1/2}}{1 - \alpha^*} h \rho^2 \right)^2.$$

To check the statements (7.4)–(7.6), the following lemma will be useful.

**Lemma 7.1.** *Let  $|U_m - 1| \leq \alpha_m < 1$  for all  $m = 1, \dots, d$ . Then for all  $i, j$ , the inequality  $|Z_{ij}^\top D_U Z_{ij}| \leq 1$  implies  $|Z_{ij}|^2 \leq 1/(1 - \alpha^*)$  and  $1 + |Z_{ij}|^2 \leq 2/(1 - \alpha^*)$ .*

*Proof.* Note that the inequalities  $Z_{ij}^\top D_U Z_{ij} \leq 1$  and  $|U_m - 1| \leq \alpha_m$  imply

$$\left| Z_{ij}^\top D_U Z_{ij} - |Z_{ij}|^2 \right| = \left| Z_{ij}^\top (D_U - I) Z_{ij} \right| \leq \alpha^* |Z_{ij}|^2$$

and thus,  $|Z_{ij}|^2 \leq (1 - \alpha^*)^{-1} Z_{ij}^\top D_U Z_{ij}$ .  $\square$

Now we check (7.4). Since

$$\begin{aligned} \begin{pmatrix} h^{-1} f(X_i) \\ P_\rho \nabla f(X_i) \end{pmatrix} &= \mathcal{V}_i(U)^{-1} \sum_{j=1}^n \begin{pmatrix} 1 \\ Z_{ij} \end{pmatrix} \begin{pmatrix} 1 \\ Z_{ij} \end{pmatrix}^\top \begin{pmatrix} h^{-1} f(X_i) \\ P_\rho \nabla f(X_i) \end{pmatrix} K_{ij}(U) \\ &= h^{-1} \mathcal{V}_i(U)^{-1} \sum_{j=1}^n \begin{pmatrix} 1 \\ Z_{ij} \end{pmatrix} \{ f(X_i) + X_{ij}^\top \nabla f(X_i) \} K_{ij}(U) \end{aligned}$$

it follows

$$\begin{aligned} s_i(U) - \begin{pmatrix} h^{-1} f(X_i) \\ P_\rho \nabla f(X_i) \end{pmatrix} &= \frac{1}{h} \mathcal{V}_i(U)^{-1} \sum_{j=1}^n \begin{pmatrix} 1 \\ Z_{ij} \end{pmatrix} \{ f(X_j) - f(X_i) - X_{ij}^\top \nabla f(X_i) \} K_{ij}(U) \\ &= \frac{1}{h} \mathcal{V}_i(U)^{-1} \sum_{j=1}^n \begin{pmatrix} 1 \\ Z_{ij} \end{pmatrix} r_{ij} K_{ij}(U) \end{aligned}$$

where in view of (3.1)

$$r_{ij} = g(\mathcal{R}^* X_{j,2}) - g(\mathcal{R}^* X_{i,2}) - (\mathcal{R}^* X_{j,2} - \mathcal{R}^* X_{i,2})^\top \nabla g(\mathcal{R}^* X_{i,2})$$

with  $\mathcal{R}^*$  being the diagonal  $M \times M$  matrix with diagonal entries  $\sqrt{v_m^*}$ ,  $m \in \mathcal{J}$ . It is clear that

$$\left| \sqrt{v_m^*} X_{j,m} - \sqrt{v_m^*} X_{i,m} \right|^2 = h^2 v_m^* (1 + \rho^{-2} v_m^*)^{-1} Z_{ij,m}^2 \leq h^2 \rho^2 Z_{ij,m}^2.$$

Therefore,

$$|\mathcal{R}^* X_{j,2} - \mathcal{R}^* X_{i,2}|^2 \leq h^2 \rho^2 |Z_{ij}|^2.$$

This yields by Lemma 7.1 and Assumption 4 for every pair  $(i, j)$  with  $Z_{ij}^\top D_U Z_{ij} \leq 1$ :

$$|r_{ij}| \leq C_g h^2 \rho^2 (1 - \alpha^*)^{-1}.$$

Using the Cauchy-Schwarz inequality and Assumptions 5 we bound

$$\begin{aligned}
|\mathcal{E}_d s_i(U) - P_\rho \nabla f(X_i)| &\leq h^{-1} \sup_{\lambda \in \mathbb{R}^{d+1}; |\lambda|=1} \left| \lambda^\top \mathcal{V}_i(U)^{-1} \sum_{j=1}^n \begin{pmatrix} 1 \\ Z_{ij} \end{pmatrix} r_{ij} K_{ij}(U) \right| \\
&\leq \sup_{|\lambda|=1} h^{-1} \left[ \sum_{j=1}^n \lambda^\top \mathcal{V}_i(U)^{-1} \begin{pmatrix} 1 \\ Z_{ij} \end{pmatrix} \begin{pmatrix} 1 \\ Z_{ij} \end{pmatrix}^\top K_{ij}(U) \mathcal{V}_i(U)^{-1} \lambda \sum_{j=1}^n r_{ij}^2 K_{ij}(U) \right]^{1/2} \\
&\leq \frac{C_g h \rho^2}{1 - \alpha^*} \sup_{|\lambda|=1} \left( \lambda^\top \mathcal{V}_i(U)^{-1} \lambda \sum_{j=1}^n K_{ij}(U) \right)^{1/2} \\
&\leq (1 - \alpha^*)^{-1} C_g h \rho^2 \|N_i(U) \mathcal{V}_i(U)^{-1}\|^{1/2} \leq (1 - \alpha^*)^{-1} C_g C_V^{1/2} h \rho^2
\end{aligned}$$

and (7.4) follows.

By definition every  $\zeta_{lm}(U)$  is a linear combination of the  $\varepsilon_i$ 's, that is, there are coefficients  $c_{i,lm}(U)$  such that

$$\zeta_{lm}(U) = \sum_{i=1}^n c_{i,lm}(U) \varepsilon_i.$$

The coefficients  $c_{i,lm}(U)$  depend on the design  $X_1, \dots, X_n$ , the basis function  $\psi_{lm}$ , the kernel  $K$  and the vector  $U$ . Moreover, these coefficients satisfy the following conditions:

**Lemma 7.2.** *For every  $l = 1, \dots, L$  and  $m = 1, \dots, d$*

- (i)  $\sum_{i=1}^n |c_{i,lm}(U^*)|^2 \leq \frac{2C_V^2 C_K}{h^2 n}$  ;
- (ii)  $\sup_{U \in \mathcal{U}_\alpha} \sum_{i=1}^n |c_{i,lm}(U)|^2 \leq \frac{2C_V^2 C_K}{(1 - \alpha^*) h^2 n}$  ;
- (iii)  $\sup_{U \in \mathcal{U}_\alpha} \left| \frac{dc_{i,lm}(U)}{dU} \right| \leq \frac{\kappa_\alpha}{nh}$ , where
$$\kappa_\alpha = \sqrt{2}(1 - \alpha^*)^{-3/2} C_V C_{K'} \bar{\psi} + 2^{3/2} (1 - \alpha^*)^{-5/2} C_V^2 C_K C_{K'} \bar{\psi}.$$

*Proof.* Define for  $i, j = 1, \dots, n$

$$N_i(U) = \sum_{j=1}^n K_{ij}(U), \quad v_{ij}(U) = \mathcal{V}_i(U)^{-1} \begin{pmatrix} 1 \\ Z_{ij} \end{pmatrix}.$$

It follows from Lemma 7.1 and Assumption 5 that  $U \in \mathcal{U}_\alpha$  implies for every  $i, j$  with  $Z_{ij}^\top D_U Z_{ij} \leq 1$

$$|N_i(U) v_{ij}(U)| \leq C_V (1 + |Z_{ij}|^2)^{1/2} \leq C_V \sqrt{2} (1 - \alpha^*)^{-1/2}. \quad (7.7)$$

Next, for a fixed  $m \leq d$ , denote by  $v_{ij,m}(U)$  the  $(m+1)$ th component of  $v_{ij}(U)$ . Then

$$\begin{aligned}\zeta_{lm}(U) &= \frac{1}{nh} \sum_{i=1}^n \psi_{lm}(X_{i,m}) \sum_{j=1}^n v_{ij,m}(U) K_{ij}(U) \varepsilon_j \\ &= \sum_{j=1}^n \left( \frac{1}{nh} \sum_{i=1}^n \psi_{lm}(X_{i,m}) v_{ij,m}(U) K_{ij}(U) \right) \varepsilon_j = \sum_{j=1}^n c_{j,lm}(U) \varepsilon_j.\end{aligned}$$

Clearly  $\mathbf{E}|\zeta_{lm}(U)|^2 = \sigma^2 \sum_{j=1}^n c_{j,lm}^2(U)$ . The Cauchy-Schwarz inequality, (7.7) and Assumption 5 imply

$$\begin{aligned}\sum_{j=1}^n c_{j,lm}^2(U) &= \frac{1}{n^2 h^2} \sum_{j=1}^n \left( \sum_{i=1}^n \psi_{lm}(X_{i,m}) v_{ij,m}(U) K_{ij}(U) \right)^2 \\ &\leq \frac{1}{n^2 h^2} \sum_{j=1}^n \left( \sum_{i=1}^n \psi_{lm}^2(X_{i,m}) v_{ij,m}(U) K_{ij}(U) \right) \left( \sum_{i=1}^n v_{ij,m}(U) K_{ij}(U) \right) \\ &\leq \frac{2C_V^2}{(1-\alpha^*)n^2 h^2} \sum_{j=1}^n \left( \sum_{i=1}^n \psi_{lm}^2(X_{i,m}) \frac{K_{ij}(U)}{N_i(U)} \right) \left( \sum_{i=1}^n \frac{K_{ij}(U)}{N_i(U)} \right) \\ &\leq \frac{2C_V^2 C_K}{(1-\alpha^*)n^2 h^2} \sum_{j=1}^n \sum_{i=1}^n \psi_{lm}^2(X_{i,m}) \frac{K_{ij}(U)}{N_i(U)} \\ &= \frac{2C_V^2 C_K}{(1-\alpha^*)n^2 h^2} \sum_{i=1}^n \psi_{lm}^2(X_{i,m}) = \frac{2C_V^2 C_K}{(1-\alpha^*)n h^2}.\end{aligned}$$

As a particular case, with  $D_U = D_{U^*} = I$  and  $\alpha^* = 0$ , this yields

$$\sum_{j=1}^n c_{j,lm}^2(U^*) \leq \frac{2C_V^2 C_K}{n h^2}$$

and the first two assertions of the lemma follows.

Now we bound the derivative of each coefficient  $c_{j,l,m}(U)$  with respect to  $U$ . For every pair  $i, j$  such that  $Z_{ij}^\top D_U Z_{ij} \leq 1$ , Lemma 7.1 implies

$$\left| \frac{d}{dU} K_{ij}(U) \right| = \left| K'(Z_{ij}^\top D_U Z_{ij}) \right| |Z_{ij}|^2 \leq (1-\alpha^*)^{-1} \left| K'(Z_{ij}^\top D_U Z_{ij}) \right|.$$

Let  $\mathbf{o}_1$  and  $\mathbf{o}_2$  be unit vectors in  $\mathbb{R}^{d+1}$ . Then for every  $m = 1, \dots, d$

$$\begin{aligned}\frac{\mathbf{o}_1^\top \partial \mathcal{V}_i(U)^{-1} \mathbf{o}_2}{\partial U_m} &= -\mathbf{o}_1^\top \mathcal{V}_i(U)^{-1} \left( \frac{\partial}{\partial U_m} \mathcal{V}_i(U) \right) \mathcal{V}_i(U)^{-1} \mathbf{o}_2 \\ &= -\mathbf{o}_1^\top \mathcal{V}_i(U)^{-1} \left( \sum_{j=1}^n \begin{pmatrix} 1 \\ Z_{ij} \end{pmatrix} \begin{pmatrix} 1 \\ Z_{ij} \end{pmatrix}^\top K'(Z_{ij}^\top D_U Z_{ij}) Z_{ij,m}^2 \right) \mathcal{V}_i(U)^{-1} \mathbf{o}_2.\end{aligned}$$

Lemma 7.1 and Assumption 5 yield

$$\left| \frac{\partial \mathbf{o}_1^\top \mathcal{V}_i(U)^{-1} \mathbf{o}_2}{\partial U_m} \right| \leq \frac{2C_V^2}{(1-\alpha^*)|N_i(U)|^2} \sum_{j=1}^n \left| K'(Z_{ij}^\top D_U Z_{ij}) \right| Z_{ij,m}^2.$$



Since  $v_{ij,m}(U) = (1 + |Z_{ij}|^2)^{1/2} \mathbf{e}_m^\top \mathcal{V}_i(U)^{-1} \mathbf{o}_2$  where  $\mathbf{e}_m$  denotes the  $m$ th coordinate vector in  $\mathbb{R}^{d+1}$  and  $\mathbf{o}_2 = (1 + |Z_{ij}|^2)^{-1/2} \begin{pmatrix} 1 \\ Z_{ij} \end{pmatrix}$ , it follows for every pair  $i, j$  such that  $Z_{ij}^\top D_U Z_{ij} \leq 1$ :

$$\begin{aligned} \left| \frac{dv_{ij,m}(U)}{dU} \right| &\leq (1 + |Z_{ij}|^2)^{1/2} \left( \sum_{m'=1}^d \left| \frac{\partial \mathbf{e}_m^\top \mathcal{V}_i(U)^{-1} \mathbf{o}_2}{\partial U_{m'}} \right|^2 \right)^{1/2} \\ &\leq \frac{2^{3/2} C_V^2}{(1 - \alpha^*)^{3/2} |N_i(U)|^2} \left[ \sum_{m'=1}^d \left( \sum_{j=1}^n |K'(Z_{ij}^\top D_U Z_{ij})| Z_{ij,m'}^2 \right)^2 \right]^{1/2} \\ &\leq \frac{2^{3/2} C_V^2}{(1 - \alpha^*)^{3/2} |N_i(U)|^2} \sum_{m'=1}^d \sum_{j=1}^n |K'(Z_{ij}^\top D_U Z_{ij})| Z_{ij,m'}^2 \\ &\leq \frac{2^{3/2} C_V^2 C_{K'}}{(1 - \alpha^*)^{3/2} |N_i(U)|^2} \sum_{j=1}^n |K'(Z_{ij}^\top D_U Z_{ij})| |Z_{ij}|^2 \leq \frac{2^{3/2} C_V^2 C_{K'}}{(1 - \alpha^*)^{5/2} |N_i(U)|}. \end{aligned}$$

Since

$$\frac{dc_{j,lm}(U)}{dU} = \frac{1}{nh} \sum_{i=1}^n v_{ij,m}(U) \psi_{lm}(X_{i,m}) \frac{dK_{ij}(U)}{dU} + \frac{1}{nh} \sum_{i=1}^n \frac{dv_{ij,m}(U)}{dU} K_{ij}(U) \psi_{lm}(X_{i,m}).$$

the use of (7.7) and Assumption 5 yields

$$\begin{aligned} \left| \frac{dc_{j,lm}(U)}{dU} \right| &\leq \frac{\sqrt{2} C_V \bar{\psi}_{lm}}{nh(1 - \alpha^*)^{3/2}} \sum_{i=1}^n \frac{|K'(Z_{ij}^\top D_U Z_{ij})|}{|N_i(U)|} + \frac{2^{3/2} C_V^2 C_{K'} \bar{\psi}_{lm}}{nh(1 - \alpha^*)^{5/2}} \sum_{i=1}^n \frac{K_{ij}(U)}{|N_i(U)|} \\ &\leq \frac{\sqrt{2} C_V C_{K'} \bar{\psi}_{lm}}{nh(1 - \alpha^*)^{3/2}} + \frac{2^{3/2} C_V^2 C_{K'} C_K \bar{\psi}_{lm}}{nh(1 - \alpha^*)^{5/2}} \end{aligned}$$

and assertion (iii) of the lemma follows.  $\square$

Since  $\mathbf{E}|\zeta_{lm}(U)|^2 = \sigma^2 \sum_{j=1}^n c_{j,lm}^2(U)$ , condition (7.6) follows from Lemma 7.2, (i).

The following lemma is a minor modification of Lemma 8 of HJS.

**Lemma 7.3.** *Let  $r$  be a positive number and let  $\Gamma$  be a finite set. Let functions  $a_{i,\gamma}(u)$  of  $u \in \mathbb{R}^d$  obey the conditions*

$$\sup_{\gamma \in \Gamma} \sup_{|u - u^*| \leq r} \left| \frac{d}{du} a_{i,\gamma}(u) \right| \leq \kappa, \quad i = 1, \dots, n. \quad (7.8)$$

If the  $\varepsilon_i$ 's are independent  $\mathcal{N}(0, \sigma^2)$ -distributed random variables, then

$$\mathbf{P} \left( \sup_{\gamma \in \Gamma} \sup_{|u - u^*| \leq r} \frac{1}{\sqrt{n}} \left| \sum_{i=1}^n \{a_{i,\gamma}(u) - a_{i,\gamma}(u^*)\} \varepsilon_i \right| > \sigma \kappa r t \right) \leq \frac{2}{n}$$

where  $t = 2 + \sqrt{2 \log(n|\Gamma|) + d \log(4n)}$  and  $|\Gamma|$  denotes the number of elements in  $\Gamma$ .

The claim (7.5) follows from Lemma 7.2, (ii) and (iii), by the preceding lemma applied with  $a_{i,\gamma}(u) = \sqrt{nc_{j,lm}}(U)$ ,  $\kappa = \frac{\kappa_\alpha}{h\sqrt{n}}$ ,  $\Gamma = \{(m,l) : m = 1, \dots, d, l = 1, \dots, L\}$ , and  $r = |\alpha|$ . This completes the proof of the proposition.

**Remark 7.1.** In the proof of Proposition 7.1 we defined the random variables  $\xi_{lm}$  as  $\xi_{lm} = \xi_{lm}(U^*)$ . One can easily check that the result of the proposition continues to hold with  $\xi_{lm}$  replaced by  $\xi_{lm}(U)$  for any  $U \in \mathcal{U}_\alpha$  and with the right hand-side of (7.2) and with the constant  $C_{\alpha,n}$  multiplied by  $(1 - \alpha^*)^{-1}$ . This fact is used in the proof Theorem 4.2.

### 7.3 Proof of Corollary 7.1

By Proposition 7.1

$$\sup_{U \in \mathcal{U}_\alpha} \max_{m=1, \dots, d} \left| P_{\rho,m} \left( \mathbf{E} \widehat{\beta}_m(U) - \beta_m^* \right) \right| \leq \frac{C_g C_V^{1/2}}{1 - \alpha^*} \rho^2 h$$

and on a random set of probability as least  $1 - 2/n$

$$\left| P_{\rho,m} \left( \widehat{\beta}_m(U) - \mathbf{E} \widehat{\beta}_m(U) \right) - \frac{\xi_m}{h\sqrt{n}} \right| \leq \frac{\sqrt{L} \sigma \bar{\psi} C_{\alpha,n} |\alpha|}{h\sqrt{n}}, \quad \forall m = 1, \dots, d,$$

where  $\xi_m \in \mathbb{R}^L$ ,  $m = 1, \dots, d$ , are Gaussian random vectors with components  $\xi_{lm}$  from Proposition 7.1.

By Lemma 7 in HJS,

$$\mathbf{P} \left( |\xi_m| \geq z_n \sqrt{\mathbf{E} |\xi_m|^2} \right) \leq 1/(nd).$$

In view of (7.2)  $\mathbf{E} |\xi_m|^2 \leq 2L\sigma^2 C_V^2 C_K$ , and the corollary follows.

### 7.4 Proof of Proposition 7.2

The definition of  $\tau_m$  implies

$$P_{\rho,m}^2 v_m^* = (1 + \rho^{-2} v_m^*)^{-1} v_m^* = \tau_m^2 \leq \min \{ \rho^2, v_m^* \}. \quad (7.9)$$

**Lemma 7.4.** *If  $P_{\rho,m} |\widehat{\beta}_m - \beta_m^*| < \delta$ , then  $P_{\rho,m}^2 |\widehat{v}_m - v_m^*| < \delta^2 + 2\delta \tau_m$ .*

*Proof.* Define the vector  $\widehat{u}_m \in \mathbb{R}^L$  (resp.  $u_m^*$ ) whose elements are  $P_{\rho,m} \widehat{\beta}_{lm}$  (resp.  $P_{\rho,m} \beta_{lm}^*$ ). Clearly  $P_{\rho,m}^2 \widehat{v}_m(U) = |\widehat{u}_m(U)|^2$  and by (7.9)  $P_{\rho,m}^2 v_m^* = |u_m^*|^2 \leq \tau_m^2$ . It is easy to check that

$$\left| |\widehat{u}_m|^2 - |u_m^*|^2 \right| \leq |\widehat{u}_m - u_m^*|^2 + 2|\widehat{u}_m - u_m^*| \cdot |u_m^*|, \quad (7.10)$$

and lemma follows.  $\square$

The proposition follows from Corollary 7.1 and Lemma 7.4.

## 7.5 Proof of Proposition 3.1

The proof of the first claim is a simplified version of the proof of Proposition 7.1: just set there  $P_{\rho,m} = 1$ , drop  $\sup_U$ , and repeat the proofs of (7.4) and (7.6). The factor  $v_{(1)}^*$  in  $\delta_1$  comes from  $\mathcal{R}^*$  in (3.1). Next, applying Lemma 7 of HJS one gets the claim (3.2). The claim (3.4) follows from (3.2) and Lemma 7.4 applied with  $\rho = 1$ .

## 7.6 Proof of Theorem 3.1

Let the numbers  $h_k$  and  $\rho_k$  be as in the algorithm description,  $k = 1, \dots, k_n$ . Define successively the values  $\delta_k$  and  $d$ -vectors  $\alpha_k$  with components  $\alpha_{k,m}$  as follows:  $\alpha_1 = 0$ ,  $\delta_1$  as in (3.3), and for  $k = 2, \dots, k_n$

$$\begin{aligned}\delta_k &= \frac{C_g C_V^{1/2}}{(1 - \alpha_k^*)} h_k \rho_k^2 + \frac{\sqrt{2L} \sigma C_V C_K^{1/2} z_n}{h_k \sqrt{n}} + \frac{\sqrt{L} \sigma \bar{\psi} C_{\alpha_k, n} |\alpha_k|}{h_k \sqrt{n}}, \\ \alpha_{k,m} &= \rho_k^{-2} (2\delta_{k-1} \tau_{k,m} + \delta_{k-1}^2), \quad m = 1, \dots, d\end{aligned}\tag{7.11}$$

with  $\alpha_k^* = \max_{m=1, \dots, d} \alpha_{k,m}$ ,  $\tau_{k,m} = \rho_k \sqrt{v_m^*} (\rho_k^2 + v_m^*)^{-1/2} \leq \min\{\rho_k, \sqrt{v_m^*}\}$ , and with  $\bar{\psi}$  defined in Proposition 7.1 and  $z_n$  in Corollary 7.1.

We will need the following two lemmas proofs of which require only minor modifications in the proofs of Lemmas 4 and 5 from HJS.

**Lemma 7.5.** *For  $n$  sufficiently large, the  $\alpha_k$ 's satisfy  $\max_{k \leq k_n} \alpha_k^* < 1/4$ . In addition, for the last iteration  $k_n$ , it holds*

$$\mu_n := \frac{C_g C_V^{1/2}}{(1 - \alpha_{k_n}^*)} h_{k_n} \rho_{k_n}^2 + \frac{\sqrt{L} \sigma \bar{\psi} C_{\alpha_{k_n}, n} |\alpha_{k_n}|}{h_{k_n} \sqrt{n}} \leq C (\sigma^2 n^{-1} L \log n)^{2/3}$$

and  $\delta_{k_n} \leq \delta_n$ , where  $\delta_n$  is defined in (3.6) and  $C$  means a generic constant depending on  $d$ ,  $M$  and the constants from Assumptions 1 through 5 only.

*Proof.* Note that  $\alpha_{k,m} \leq \delta_{k-1}^2 / \rho_k^2$  for all  $m \notin \mathcal{J}$  and  $\alpha_{k,m} \leq \delta_{k-1}^2 / \rho_k^2 + 2\delta_{k-1} / \rho_k$  for  $m \in \mathcal{J}$ . The first assertion of the lemma easily follows from the fact that  $h_k \rho_k$  decreases during iteration, cf. Lemma 4 of HJS.

Since the dimensionality of the nonlinear component is bounded by  $M$ , it follows

$$|\alpha_k|^2 \leq (d - M) \delta_{k-1}^4 / \rho_k^4 + M (\delta_{k-1}^2 / \rho_k^2 + 2\delta_{k-1} / \rho_k)^2.$$

Further, the inequality  $|\alpha_{k_n-1}| \leq C_1$  with some constant  $C_1$  depending on  $d$  and  $M$  only implies in view of  $h_{k_n-1} \geq 1/a_h$  and  $1 \leq \rho_{k_n-1} (\sigma^2 n^{-1} L \log n)^{-1/3} \leq 1/a_\rho$  that

$$\delta_{k_n-1} \leq C (\sigma^2 n^{-1} L \log n)^{1/2}, \quad |\alpha_{k_n}| \leq C (\sigma^2 n^{-1} L \log n)^{1/6}.$$

Substituting this bound in the formula for  $\mu_n$  yields by  $h_{k_n} \geq 1$  and  $\rho_{k_n} = (\sigma^2 n^{-1} L \log n)^{1/3}$  that  $\mu_n \leq C (\sigma^2 n^{-1} L \log n)^{2/3}$  and therefore

$$\delta_{k_n} \leq \sqrt{2} C_V C_K^{1/2} (\sigma^2 n^{-1} L z_n^2)^{1/2} + C (\sigma^2 n^{-1} L \log n)^{2/3}.$$

□

**Lemma 7.6.** *Let  $n$  be sufficiently large. There exist random sets  $\mathcal{A}_1 \supseteq \dots \supseteq \mathcal{A}_{k_n-1}$  such that  $\mathbf{P}(\mathcal{A}_k) \geq 1 - \frac{3k}{n}$  and it holds on  $\mathcal{A}_k$*

$$\max_{m=1, \dots, d} |P_{\rho_{k+1}, m}(\widehat{\beta}_m^{(k)} - \beta_m^*)| \leq \delta_k, \quad k = 1, \dots, k_n - 1.$$

*Proof.* We proceed by induction in  $k$ . First by (3.2) there exists a random set  $\mathcal{A}_1$  with  $\mathbf{P}(\mathcal{A}_1) \geq 1 - 1/n$  such that  $\max_{m=1, \dots, d} |\widehat{\beta}_1 - \beta_1^*| \leq \delta_1$  on  $\mathcal{A}_1$ . This obviously implies

$$\max_{m=1, \dots, d} |P_{\rho_2, m}(\widehat{\beta}_1 - \beta_1^*)| \leq \delta_1.$$

Suppose now that there is  $\mathcal{A}_{k-1}$  such that  $\mathbf{P}(\mathcal{A}_{k-1}) \geq 1 - \frac{3(k-1)}{n}$  and it holds on  $\mathcal{A}_{k-1}$ :

$$\max_{m=1, \dots, d} |P_{\rho_k, m}(\widehat{\beta}_m^{(k-1)} - \beta_m^*)| \leq \delta_{k-1}.$$

Then on  $\mathcal{A}_{k-1}$  by Lemma 7.4  $P_{\rho_k, m}^2 |\widehat{v}_m^{(k-1)} - v_m^*| < \delta_{k-1}^2 + 2\delta_{k-1} \tau_{k, m}$  simultaneously for all  $m = 1, \dots, d$ , and denoting  $U^{(k)}$  a  $d$ -vector with components  $U_m^{(k)} = P_{\rho_k, m}^2 (1 + \rho_k^{-2} \widehat{v}_m^{(k-1)})$ , one gets  $U^{(k)} \in \mathcal{U}_{\alpha_k}$ .

By Corollary 7.1 there exists another random set  $\mathcal{A}_k$  with  $\mathbf{P}(\mathcal{A}_k) \geq 1 - 3/n$  such that on  $\mathcal{A}_k$  it holds for every  $U \in \mathcal{U}_{\alpha_k}$

$$\max_{m=1, \dots, d} |P_{\rho_k, m}(\widehat{\beta}_m(U) - \beta_m^*)| \leq \delta_k,$$

so that, with  $\mathcal{A}_k = \mathcal{A}_{k-1} \cap \mathcal{A}_k$ , we obtain  $\mathbf{P}(\mathcal{A}_k) \geq 1 - 3k/n$  and it holds on  $\mathcal{A}_k$

$$\max_{m=1, \dots, d} |P_{\rho_k, m}(\widehat{\beta}_m^{(k)} - \beta_m^*)| \leq \delta_k.$$

and, since for every  $m$   $P_{\rho_{k+1}, m} \leq P_{\rho_k, m}$ , the assertion follows. □

Let now  $\mathcal{A}_{k_n-1}$  be the random set with  $\mathbf{P}(\mathcal{A}_{k_n-1}) \geq 1 - \frac{3k_n-3}{n}$  shown in Lemma 7.6 so that on this set

$$\max_{m=1, \dots, d} |P_{\rho_{k_n}, m}(\widehat{\beta}_m^{(k_n-1)} - \beta_m^*)| \leq \delta_{k_n-1},$$

and for the corresponding  $d$ -vector  $U^{(k_n)}$  with components  $U_m^{(k_n)} = P_{\rho_{k_n}, m}^2 (1 + \rho_{k_n}^{-2} \hat{v}_m^{(k_n-1)})$ , it holds  $U^{(k_n)} \in \mathcal{U}_{\alpha_{k_n}}$ .

Let then  $\xi_m$  be the Gaussian  $L$ -vector with the components  $\xi_{lm}$  from Proposition 7.1 applied with  $h = h_{k_n}$  and  $\rho = \rho_{k_n}$ . Due to this proposition, there exists a random set  $A_{k_n}$  with  $\mathbf{P}(A_{k_n}) \geq 1 - 2/n$ , so that on  $A_{k_n}$  it holds for all  $U \in \mathcal{U}_{\alpha_{k_n}}$ :

$$\max_{m=1, \dots, d} |P_{\rho_{k_n}, m}(\hat{\beta}_m(U) - \beta_m^*) - \frac{\xi_m}{h\sqrt{n}}| \leq \mu_n,$$

where  $\mu_n$  is defined in Lemma 7.5. This yields for the set  $\mathcal{A}_{k_n} = \mathcal{A}_{k_n-1} \cap A_{k_n}$  that  $\mathbf{P}(\mathcal{A}_{k_n}) \geq 1 - \frac{3k_n-1}{n}$  and the final estimator  $\hat{\beta}_m = \hat{\beta}_m^{(k_n)}$  satisfies on  $\mathcal{A}_{k_n}$ :

$$\max_{m=1, \dots, d} |P_{\rho_{k_n}}(\hat{\beta}_m - \beta_m^*) - n^{-1/2} \xi_m^*| \leq \mu_n$$

where  $\xi_m^* = h^{-1} \xi_m$ . In view of  $h = h_{k_n} \geq 1$

$$\mathbf{E}|\xi_{lm}^*|^2 = h^{-2} \mathbf{E}|\xi_{lm}|^2 \leq 2\sigma^2 C_V^2 C_K$$

and the first two claims in (3.5) follow from Lemma 7.5. The last claim in (3.5) follow by applying Lemma 7 of HJS and Lemma 7.4. The first two inequalities in (3.7) follow from (3.5) by setting  $P_{\rho, m} = 1$  and  $\beta_m^* = 0$ . The last one is proved similarly to Lemma 7.4.

## 7.7 Proof of Theorem 3.2

The proof can be done similarly to Spokoiny (2002) using the bound for the bias of estimation from the proof of Proposition 7.1. We omit the details to save the space.

## 7.8 Proof of Theorem 4.1

In view of Theorem 3.1 on the set  $A$ , it holds  $\hat{v}_m \leq \delta_n^2$  for all  $m \notin \mathcal{J}$ . Therefore, it suffices to show that on  $A$ , it holds  $\hat{v}_m > r^2 \delta_n^2$  for every  $m \in \mathcal{J}$ . Next, by Theorem 3.1 again, for  $m \in \mathcal{J}$

$$\hat{v}_m > v_m^* - P_\rho^{-2} (\delta_n^2 + 2\delta_n P_\rho v_m^*) = v_m^* - \delta_n^2 (1 + v_m^* \rho^{-2}) - 2\delta_n (1 + v_m^* \rho^{-2})^{1/2} v_m^*.$$

Define  $s^2 = v_m^* / \delta_n^2$  and  $u_n = \delta_n / \rho$ . Then, on  $A$ ,

$$\delta_n^{-2} \hat{v}_m > s^2 - 1 - s^2 u_n^2 - 2s(1 + s^2 u_n^2)^{1/2} \geq s^2(1 - u_n^2 - 2u_n) - 1 - 2s.$$

It is straightforward to check that the right hand-side of this inequality as a function of  $s$  is greater than  $r^2$  for all  $s \geq s_r$ . Therefore, on  $A$ ,  $\delta_n^{-2} \hat{v}_m > r^2$  for  $m \in \mathcal{J}$  as required.

## 7.9 Proof of Theorem 4.2

To simplify the exposition, we suppose that the resampling scheme of Section 4.1 utilizes the true variance  $\sigma^2$  instead of the estimated variance  $\hat{\sigma}^2$ . This assumption is easily justified by the result of Theorem 3.2 claiming root- $n$  consistent estimation of  $\sigma^2$  by  $\hat{\sigma}^2$ .

The idea of the proof is to show that the variable  $\hat{v}_{(M+1)}$  and the similarly defined variable  $\tilde{v}_{(M+1)}$  for the resampling model have approximately the same distribution. Let  $A$  be the random set from Theorem 3.1 with  $\mathbf{P}(A) \geq 1 - 3k_n/n$ . It is obviously sufficient to show that

$$\mathbf{P}(\hat{\mathcal{J}}_M \neq \mathcal{J} \mid A) \leq \alpha + 3/n.$$

We therefore suppose that the event  $A$  holds true. Then, under the assumptions of the theorem, the nonlinear component is correctly identified and all the bounds of Theorem 3.1 hold. Moreover, for every  $m \notin \mathcal{J}$ , the value  $n\hat{v}_m$  can be approximated by  $|\xi_m^*|^2$ , where the distribution of the vector  $\xi_m^*$  depend on the ‘ideal’ bandwidth  $\mathbf{b}^* = \mathbf{b}^{*(k_n)}$ , the kernel  $K$ , basis functions  $\psi_{lm}(\cdot)$ , and the design  $X_1, \dots, X_n$  only.

Next we consider the model we resample from. This artificial model has the same structure (i.e. the same linear and nonlinear components) and differs from the original one only by the parameters of the linear component (they are equal to zero in the resampling model) and by the nonlinear link function. More specifically, the estimators  $\hat{v}_m$  based on the original model are the ‘‘true’’ values for the resampling model and the last step bandwidth  $\mathbf{b} = \mathbf{b}^{(k_n)}$  is the ‘‘ideal’’ bandwidth for the resampling model. Since the resampling model fulfills all the conditions that we impose on the original model, Theorem 3.1 (or Proposition 7.1 with  $\alpha = 0$  and  $\mathbf{b} = \mathbf{b}^{(k_n)}$ ) continues to apply. This yields, in particular, that on a set  $\tilde{A}_M$  with  $\mathbf{P}(\tilde{A}_M) \geq 1 - 3/n$ , the nonlinear component of the resampling model will be correctly identified. Moreover, due to Remark 7.1, every variable  $n\tilde{v}_m$  with  $m \notin \mathcal{J}$  can be approximated by the squared norm of a Gaussian random vector with the same distribution as  $\xi_m^*$ . And thus, it is true for  $n\tilde{v}_{(M+1)}$ . This yields, in particular, that the  $(1 - \alpha)$ -quantile evaluated from the distribution of  $n\tilde{v}_{(M+1)}$  applies up to the approximation error to  $n\hat{v}_{(M+1)}$ . It follows from Theorem 3.1 that the error of approximation of  $n\hat{v}_m$  by  $|\xi_m^*|^2$  can be bounded by  $n(\omega_n^2 + \omega_n\delta_n) \leq C'n^{-1/6}(\log n)^{5/6}$  for some constant  $C'$ . Therefore, at least for sufficiently large  $n$ , the approximation error is small and the assertion of the theorem follows.

## 7.10 Proof of Theorem 4.3

Let  $A$  be the random set described in Theorem 3.1 with  $\mathbf{P}(A) \geq 1 - 3k_n/n$ . In view of Theorem 4.1, it is sufficient to prove that  $\mathbf{P}(\widehat{\mathcal{M}} \neq M \mid A) \leq \alpha + 3M/n$ .

On  $A$  it holds  $\widehat{v}_m \leq \delta_n^2$  for all  $m \notin \mathcal{J}$  and  $\widehat{v}_m > (r\delta_n)^2$  for all  $m \in \mathcal{J}$  and  $r = s_1$ . Thus  $\widehat{v}_{(\mathcal{M})} > (s_1\delta_n)^2$  for all  $\mathcal{M} \leq M$  and  $\widehat{v}_{(M+1)} \leq \delta_n^2$ . For every  $\mathcal{M} < M$ , we resample from the model having precisely  $\mathcal{M}$  nonlinear regressors with  $\widehat{v}_m$  being the ‘true’ measure of nonlinearity for every  $m \in \widehat{\mathcal{J}}_{\mathcal{M}}$ .

Application of Propositions 7.1 and 7.2 with  $\alpha = 0$  to this artificial models and again Theorem 4.1 with  $r = 1$  ensures that on a set  $\widetilde{A}_{\mathcal{M}}$  with  $\mathbf{P}(\widetilde{A}_{\mathcal{M}}) \geq 1 - 3/n$ , every  $\widetilde{v}_m$  for  $m \notin \mathcal{J}_{\mathcal{M}}$  fulfills  $\widetilde{v}_m \leq \delta_n^2$ . Hence,  $\widetilde{v}_{(M+1)} \leq \delta_n^2$  on  $\widetilde{A}_{\mathcal{M}}$  and the same holds for the  $1 - \alpha$  quantile of  $\widetilde{v}_{(M+1)}$  provided that  $\alpha > 3/n$ . Therefore, for every  $\mathcal{M} < M$ , the hypothesis  $M \leq \mathcal{M}$  will be rejected on the intersection  $A \cap \widetilde{A}_{\mathcal{M}}$ . This yields

$$\mathbf{P}(\widehat{\mathcal{M}} < M \mid A) \leq 3(M - 1)/n. \quad (7.12)$$

Next the definition of  $\widehat{\mathcal{M}}$  implies the inclusion

$$\{\widehat{\mathcal{M}} > M\} \subseteq \{\widehat{v}_{(M+1)} > t_{\alpha}(M)\},$$

where  $t_{\alpha}(M)$  is evaluated in the resampling procedure with  $\mathcal{M} = M$ . Applying now Theorem 4.2 we get, using also (7.12), the desired bound for  $\mathbf{P}(\widehat{\mathcal{M}} \neq M)$ , and the theorem follows.

## References

- [1] Bhattacharya, P.K. and Zhao, P.L. (1997). Semiparametric inference in a partial linear model. *Ann. Statist.*, **25**, 244–262.
- [2] Bickel, P.J., Klaassen, C.A.J., Ritov, Y., and Wellner, J. (1998). *Efficient and adaptive estimation for semiparametric models*. New York: Springer. xix, 560 p.
- [3] Carroll, R.J., Fan, J., Gijbels, I., and Wand, M.P. (1997). Generalized partially linear single-index models. *J. Amer. Statist. Assoc.*, **92**, no. 438, 477–489.
- [4] Chen, H. (1988). Convergence rates for parametric components in a partly linear model. *Ann. Statist.*, **16**, 136–146.
- [5] Chen, H. and Chen, K. W. (1991). Selection of the splined variables and convergence rates in a partial linear model. *Canadian J. Statist.*, **19**, 323–339.
- [6] Cuzick, J. (1992a). Semiparametric additive regression. *J. Royal Statist. Soc., Series B*, **54**, 831–843.
- [7] Cuzick, J. (1992b). Efficient estimates in semiparametric additive regression models with unknown error distribution. *Ann. Statist.*, **20**, 1129–1136.

- [8] Engle, R.F., Granger, C.W.J., Rice, J., and Weiss, A. (1986). Semiparametric estimates of the relation between weather and electricity sales. *J. Amer. Statist. Ass.*, **81**, 310–320.
- [9] Eubank, R.L. and Spiegelman, C.H. (1990). Testing the goodness of fit of a linear model via nonparametric regression techniques. *J. Amer. Stat. Ass.*, **85**, 387–392.
- [10] Eubank, R.L. and Hart, J.D. (1992). Testing goodness-of-fit in regression via order selection criteria. *Ann. Statist.*, **20**, no. 3, 1424–1425.
- [11] Eubank, R.L., Kambour, E.L., Kim, J.T., Klipple, K., Reese, C.S. and Schimek, M. (1998). Estimation in partially linear models. *Computational Statistics Data Analysis*, **29**, 27–34.
- [12] Fan, J. (1996). Test of significance based on wavelet thresholding and Neyman’s truncation. *J. Amer. Statist. Ass.*, **91**, 674–688.
- [13] Gao, J. T. (1995). Asymptotic theory for partly linear models. *Communications in Statistics, Theory Methods*, **24**, 1985–2010.
- [14] Gasser, T., Sroka, L., and Jennen-Steinmetz, C. (1986). Residual variance and residual pattern in nonlinear regression. *Biometrika*, **73**, 625–633.
- [15] Golubev, G. and Härdle, W. (2000). On the second order minimax estimation in a partly linear model. *Math. Methods Statist.*, **9**, 160–175.
- [16] Hamilton, S.A. and Truong, Y.K. (1997). Local linear estimation in partly linear models. *J. Multivariate Analysis*, **60**, 1–19.
- [17] Härdle, W. and Korostelev, A. (1996). Search of significant variables in nonparametric additive regression. *Biometrika*, **83**, 541–549.
- [18] Härdle, W., Liang, H., and Gao, J. (2000). *Partially linear models*. Physica Verlag and Springer Verlag, Germany.
- [19] Härdle, W. and Mammen, E. (1993). Comparing nonparametric versus parametric regression fits. *Ann. Statist.*, **4**, 1926–1947.
- [20] Härdle, W., Sperlich, S., and Spokoiny, V. (2001). Structural tests for additive regression. *J. Amer. Stat. Acc.*, **96**, no. 456, 1333–1347.
- [21] Hart, J. (1997). *Nonparametric smoothing and lack-of-fit tests*. Springer Verlag, New York.
- [22] Horowitz, J.L. and Spokoiny, V. (2001). An adaptive, rate-optimal test of a parametric mean regression model against a nonparametric alternative. *Econometrica*, **69**, 599–631.
- [23] Hristache, M., Juditsky, A., and Spokoiny, V. (2001). Direct estimation of the index coefficient in a single-index model. *Ann. Statist.*, **29**, no. 3, 595–623.
- [24] Hristache, M., Juditsky, A., Polzehl, J., and Spokoiny, V. (2001). Structure adaptive approach for dimension reduction. *Ann. Statist.*, **29**, no. 6, 1537–1566.
- [25] Heckman, N.E. (1986). Spline smoothing in partly linear models. *J. Royal Statist. Soc., Series B*, **48**, 244–248.
- [26] Ledwina, T. (1994). Data-driven version of Neyman’s smooth test of fit. *J. Amer. Stat. Ass.*, **89**, no. 427, 1000–1005.
- [27] Mammen, E. and van de Geer, S. (1997). Penalized estimation in partial linear models. *Ann. Statist.*, **25**, 1014–1035.
- [28] Rice, J. (1986). Convergence rates for partially splined models. *Statistics Probability Letters*, **4**, 203–208.



- [29] Robinson, P.M. (1988). Root-n-consistent semiparametric regression. *Econometrica*, **56**, 931–954.
- [30] Schick, A. (1996a). Weighted least squares estimates in partly linear regression models. *Statistics Probability Letters*, **27**, 281–287.
- [31] Schick, A. (1996b). Root-n consistent estimation in partly linear regression models. *Statistics Probability Letters*, **28**, 353–358.
- [32] Schimek, M. (2000). Estimation and inference in partially linear models with splines. *J. Statistical Planning Inference*, **91**, 525–540.
- [33] Speckman, P. (1988). Kernel smoothing in partial linear models. *J. R. Statist. Soc., Series B*, **50**, 413–436.
- [34] Sperlich, S. (1998). *Additive modelling and testing model specification*. Aachen, Shaker.
- [35] Spokoiny, V. (2001). Data driven testing the fit of linear models. *Math. Methods of Statistics*, **10**, no. 4, 465–497.
- [36] Spokoiny, V. (2002). Variance estimation for high-dimensional regression models. *J. of Multivariate Analysis*, pp. 1–23 (doi:10.1006/jmva.2001.2023).
- [37] Stute, W. (1997). Nonparametric model checks for regression. *Ann. Statist.*, **25**, no. 2, 613–641.