

# Weierstraß–Institut für Angewandte Analysis und Stochastik

im Forschungsverbund Berlin e.V.

Preprint

ISSN 0946 – 8633

## Structure adaptive approach for dimension reduction

Marian Hristache<sup>1</sup>, Anatoli Juditsky<sup>2</sup>, Jörg Polzehl<sup>3</sup>, Vladimir Spokoiny<sup>3</sup>

submitted: March 28, 2000

<sup>1</sup> ENSAI

Campus Ker Lann rue B. Pascal  
35170 Bruz France  
E-Mail: hristach@@ensai.fr

<sup>2</sup> LMC

Domaine Universitaire B.P.53  
38041 Grenoble Cedex 9 France  
E-Mail: anatoli.iouditski@@inrialpes.fr

<sup>3</sup> Weierstrass Institute for Applied Analysis  
and Stochastics, Mohrenstraße 39

D – 10117 Berlin, Germany

E-Mail: polzehl@wias-berlin.de

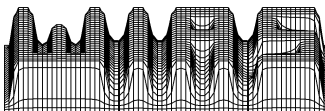
E-Mail: spokoiny@wias-berlin.de

<http://www.wias-berlin.de/~polzehl>

<http://www.wias-berlin.de/~spokoiny>

Preprint No. 569

Berlin 2000



---

1991 *Mathematics Subject Classification.* Primary 62G05; Secondary 62H40, 62G20.

*Key words and phrases.* dimension-reduction, multi-index model, index space, average derivative estimation, structural adaptation.

The correspondence should be sent to Vladimir Spokoiny.

Edited by

Weierstraß-Institut für Angewandte Analysis und Stochastik (WIAS)

Mohrenstraße 39

D — 10117 Berlin

Germany

Fax: + 49 30 2044975

E-Mail (X.400): c=de;a=d400-gw;p=WIAS-BERLIN;s=preprint

E-Mail (Internet): preprint@wias-berlin.de

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

## Abstract

We propose a new method of effective dimension reduction for a multi-index model which is based on iterative improvement of the family of average derivative estimates. The procedure is computationally straightforward and does not require any prior information about the structure of the underlying model. We show that in the case when the effective dimension  $m$  of the index space does not exceed 3, this space can be estimated with the rate  $n^{-1/2}$  under rather mild assumptions on the model.

## 1 Introduction

Suppose that the observations  $(Y_i, X_i)$ ,  $i = 1, \dots, n$ , are generated by the regression model

$$Y_i = f(X_i) + \varepsilon_i \quad (1.1)$$

where  $Y_i$  is a scalar response variables,  $X_i \in [-1, 1]^d$  are  $d$ -dimensional explanatory variables,  $\varepsilon_i$  are random errors and  $f(\cdot)$  is an unknown  $d$ -dimensional function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ .

We assume that  $f(x)$  has the specific structure:

$$f(x) = g_0(Tx). \quad (1.2)$$

Here  $g_0(\cdot)$  is an unknown  $m$ -dimensional *link* function and  $T$  is a linear orthogonal mapping from the high-dimensional space  $\mathbb{R}^d$  onto the space  $\mathbb{R}^m$  with an essentially smaller dimension  $m$ , satisfying the condition  $TT^\top = I_m$ , where  $T^\top$  stands for the transpose of  $T$ . In the statistical literature relations as in (1.1) and (1.2) are referred to as *multi-index regression* models. Model (1.2) is a rather general expression of the hypothesis that all the information about  $f(x)$  is “concentrated” in a low-dimensional projection  $Tx$ . If we adopt such a model, our intention can be both to find the *effective dimension*  $m$  and to describe the *index space*  $\mathcal{I} = \text{Im } T^\top$  which is also referred to as the *effective dimension space* or the space of *effective dimension reduction* in Li (1991, 1992) and Cook (1998). In the present paper we propose an algorithm to estimate the index

space when the effective dimension  $m$  is known *a priori*. Some extensions are discussed in Section 6.

Note first that the representation (1.2) is not unique. For instance, if  $O_m$  is an orthogonal transform in  $\mathbb{R}^m$ , then the function  $f$  can be rewritten in the form  $f(x) = g_1(T_1x)$  with  $g_1(z) = g_0(O_m z)$  and  $T_1 = O_m^\top T$ . Nevertheless, the index space  $\mathcal{I}$  is defined uniquely by (1.2) and it contains very important information about the model. As soon as the operator  $T$  which maps  $\mathbb{R}^d$  onto  $\mathbb{R}^m$  is fixed, the link function  $g_0$  can be estimated in a nonparametric way.

Various methods for dimension reduction have been proposed in the literature. Classical theory of *principle component analysis* considers mostly the case of multiple *linear* regression. Brillinger (1983) extended the method to the so called “generalized linear model” with normally distributed regressors. The underlying idea is to make some data transformation and then to proceed as if the model were linear. Under a similar assumption on the distribution of regressors, Li (1991) offered the so called “sliced inverse regression” approach. A modification of this method (principle Hessian directions) is explored in Li (1992) and Cook (1998). Samarov (1993) discussed an approach relying on average derivative estimation of some linear functionals of the gradient of the regression function  $f$ . However, the conditions for this method to work appear to be quite restrictive in application to real data. The main problem here is that, for large  $d$ , the data in the high dimensional space  $\mathbb{R}^d$  is very sparse (the so called “curse of dimensionality” problem).

Our approach can be seen as an iterative improvement of the average derivative estimator and can be used under weak assumptions on the model. The proposed procedure can be regarded as an extension of the method developed in Hristache, Juditsky and Spokoiny (1998) for the single-index model to the multi-index situation. In the sequel the latter paper is referred to as HJS98.

The paper is organized as follows: in the next section we discuss the heuristics behind the proposed approach. Then in Section 3 the estimation procedure is presented. The performance of the method is tested for some simulated datasets in Section 4. The theoretical setting is given and asymptotic properties of the algorithm are studied in Section 5. Section 6 shortly summarizes main results and discusses possible extensions and open problems. Finally, the proofs are collected in the appendix.

## 2 Basic ideas

Since the gradient  $F(X_i) = \nabla f(X_i)$  of the regression function  $f$  at every point  $X_i$  belongs to the index space  $\mathcal{I}$ , it seems quite natural to apply the principle component analysis for estimating this space: one can compute the matrix  $\mathcal{M}^* = \sum_{i=1}^n F(X_i)F^\top(X_i)$  and then use the eigenvalue decomposition of  $\mathcal{M}$ ,  $\mathcal{M}^* = O_d^\top \Lambda O_d$ . Here  $O_d$  is an orthonormal matrix and  $\Lambda$  is a diagonal matrix with decreasing eigenvalues. These matrices deliver a valuable important information about model (1.2): the first  $m$  columns of  $U$  (i.e. the first  $m$  eigenvectors of  $\mathcal{M}^*$ ) provide an orthonormal basis of the index space  $\mathcal{I}$ ; the corresponding eigenvalues show how fast the function  $f$  varies in each direction. In particular, the first eigenvector of  $\mathcal{M}^*$  is the direction in which  $f$  varies most (cf. Samarov (1993)). This leads to the natural idea, to first estimate  $\mathcal{M}^*$  from the data  $Y_1, \dots, Y_n$  and then to recover the index space  $\mathcal{I}$  using this estimate. Note that the matrix  $\mathcal{M}^*$  is a quadratic functional of the gradient of the regression function  $f$ . There is a number of papers on estimation of such functionals in the framework of nonparametric regression. Various estimation algorithms and results on their optimality can be found in Ibragimov, Nemirovskii and Khasmiskii (1986), Donoho and Nussbaum (1990), Fan (1991). The estimators in Samarov (1993) and Doksum and Samarov (1995) are based on kernel estimators of the regression function  $f$ , Huang and Fan (1998) applied the local polynomial fit, the procedure from Ibragimov, Nemirovskii and Khasmiskii (1986) is based on the Fourier expansion of the gradient  $F$  of the function  $f$ . Let us see how this latter idea applies to our problem.

Suppose that we are given a collection  $\{\psi_\ell, \ell = 1, \dots, L\}$  of functions  $\psi_\ell : \mathbb{R}^d \rightarrow \mathbb{R}$  which satisfy

$$\sum_{i=1}^n \psi_\ell(X_i)\psi_{\ell'}(X_i) = \delta_{\ell\ell'}$$

where  $\delta_{\ell\ell} = 1$  and  $\delta_{\ell\ell'} = 0$  for  $\ell \neq \ell'$ . Now, let  $\beta_\ell^*$ ,

$$\beta_\ell^* = \sum_{i=1}^n F(X_i)\psi_\ell(X_i), \tag{2.1}$$

be the  $\ell$ -th Fourier coefficient of  $F$  with respect to the basis system  $\{\psi_\ell\}$ . Note that each  $d$ -vector  $\beta_\ell^*$  is a linear functional of the gradient and hence belongs to  $\mathcal{I}$ . Thus if the dimension of the space spanned by  $\beta_1^*, \dots, \beta_L^*$  equals  $m$ , this set of vectors completely characterizes the index space  $\mathcal{I}$ , and one can identify the space  $\mathcal{I}$  by looking for the first  $m$  principal components of the set  $\beta_1, \dots, \beta_L$ .

In order to estimate  $\mathcal{M}^*$ , one can first construct an estimate  $\widehat{\beta}_\ell$  of each Fourier coefficient

$\beta_\ell^*$ , e.g.

$$\widehat{\beta}_\ell = \sum_{i=1}^n \widehat{F}(X_i) \psi_\ell(X_i) \quad (2.2)$$

on the basis of a pilot estimate  $\widehat{F}$  of the gradient, and then compose the estimate

$$\widehat{\mathcal{M}}_L = \sum_{\ell=1}^L \widehat{\beta}_\ell \widehat{\beta}_\ell^\top$$

of  $\mathcal{M}^*$ . Note that in order to ensure  $\widehat{\mathcal{M}}_L$  to be a consistent estimate of the matrix  $\mathcal{M}^*$  the number  $L$  of basis functions  $\psi_\ell$  should be taken growing with  $n$ . Otherwise  $\widehat{\mathcal{M}}_L$  estimates the matrix  $\mathcal{M}_L^*$  with

$$\mathcal{M}_L^* = \sum_{\ell=1}^L \beta_\ell \beta_\ell^\top.$$

On the other hand, recall that it is the index space  $\mathcal{I}$  we are interested in, and not the estimation of  $\mathcal{M}^*$ . It would be sufficient for our purposes to point out a fixed (possibly small) number of “test functions”  $\psi_\ell$  such that  $\text{rank}(\mathcal{M}_L^*) = m$  and the value  $\|\mathcal{M}^* - \mathcal{M}_L^*\|$  is not too large. The choice of a proper set of test functions  $\psi_\ell$ ,  $\ell = 1, \dots, L$  is a very sensitive issue of the proposed approach. Some heuristic ideas about it are discussed in more details in Section 3.4.

## 2.1 Equivalent representation

As we have already noticed, the model representation (1.2) is not unique. It is more convenient for our purposes to work with another one, which is distinctly defined by the set of test functions  $\psi_\ell$ ,  $\ell = 1, \dots, L$  and the regression function  $f$ .

Let us denote  $\mathcal{B}^*$  the  $d \times L$  matrix with the columns  $\beta_\ell^*$ ,  $\ell = 1, \dots, L$ , where the vectors  $\beta_\ell^*$  are as in (2.1). Obviously, each vector  $\beta_\ell^*$  belongs to  $\mathcal{I}$  and hence  $\text{rank}(\mathcal{B}^*) \leq m$ . We additionally suppose that  $\text{rank}(\mathcal{B}^*) = m$  which means that this matrix completely describes the index space  $\mathcal{I}$ .

Let  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$  be the ordered set of eigenvalues of the symmetric  $d \times d$ -matrix  $\mathcal{M}_L^* = \mathcal{B}^* (\mathcal{B}^*)^\top$ . Since  $\text{rank}(\mathcal{M}_L^*) = m$ , only the first  $m$  of them are positive and the remainings are equal to zero. Without loss of generality we assume that

$$\lambda_1 > \lambda_2 > \dots > \lambda_m \quad (2.3)$$

which ensures that the corresponding eigenvectors of unit length  $e_1, \dots, e_m$  are uniquely defined (up to a sign). These vectors belong to the index space  $\mathcal{I}$  and can be used as a

natural basis in it. We also denote  $\theta_k = \sqrt{\lambda_k} e_k$ ,  $k = 1, \dots, m$ . Since  $\lambda_k = 0$  for  $k > m$ , it also holds  $\theta_k = 0$  for those  $k$ .

We now represent the model (1.1), (1.2) in the form

$$f(x) = g(\theta_1^\top x, \dots, \theta_m^\top x) \quad (2.4)$$

where the new link function  $g$  is uniquely defined as soon as the vectors  $\theta_1, \dots, \theta_m$  are fixed. Usually a similar representation with vectors  $e_k = \theta_k/|\theta_k|$  in place of  $\theta_k$  is used:

$$f(x) = g_1(e_1^\top x, \dots, e_m^\top x). \quad (2.5)$$

However, the value  $\lambda_k$  characterizes a variability of the function  $f$  in the direction  $e_k$ . Thus the function  $g_1$  in (2.5) inherits the inhomogeneity of  $f$  in different directions. The benefit of using (2.4) is that the corresponding link function  $g$  is homogeneous w.r.t. its variables.

Let  $\mathcal{R}^*$  be a  $m \times d$ -matrix such that its transpose  $(\mathcal{R}^*)^\top = (\theta_1, \dots, \theta_m)$  has vectors  $\theta_1, \dots, \theta_m$  as columns. Then (2.4) can be rewritten as  $f(x) = g(\mathcal{R}^* x)$ . The matrix  $\mathcal{R}^*$  maps  $\mathbb{R}^d$  onto  $\mathbb{R}^m$  and determines the required effective dimension space. In what follows we refer to  $\mathcal{R}^*$  as the *effective dimension reduction matrix*, or simply the *e.d.r.*

The following lemma offers an explicit representation of the matrix  $\mathcal{R}^*$  via the orthogonal decomposition of the symmetric  $L \times L$ -matrix  $(\mathcal{B}^*)^\top \mathcal{B}^*$ .

**Lemma 2.1** *Let  $(\mathcal{B}^*)^\top \mathcal{B}^* = O \Lambda_L O^\top$  be the orthogonal decomposition of  $(\mathcal{B}^*)^\top \mathcal{B}^*$  where  $O$  is an orthogonal  $L \times L$ -matrix and  $\Lambda_L$  is a diagonal matrix with non-increasing eigenvalues  $\lambda'_1 \geq \lambda'_2 \geq \dots \geq \lambda'_L$ . Let also  $O_m$  be the block of the first  $m$  columns of  $O$ . Then  $\lambda'_k = \lambda_k$  for  $k \leq d$  and*

$$\mathcal{R}^* = (\mathcal{B}^* O_m)^\top. \quad (2.6)$$

Due to this lemma, the model (2.4) can be now rewritten in the form

$$f(x) = g(\mathcal{R}^* x) = g\left((\mathcal{B}^* O_m)^\top x\right) \quad (2.7)$$

which is used in the sequel.

## 2.2 Gradient estimation

Next we discuss the problem of estimating each linear functional  $\beta_\ell^*$  using a nonparametric estimate  $\widehat{F}$  of the gradient  $F$ , see (2.2). A standard way to estimate both  $f(X_i)$

and  $F(X_i)$  is to apply the local linear least squares approach:

$$\begin{pmatrix} \hat{f}(X_i) \\ \hat{F}(X_i) \end{pmatrix} = \underset{c \in \mathbb{R}, b \in \mathbb{R}^d}{\operatorname{arginf}} \sum_{j=1}^n \left[ Y_j - c - b^\top (X_j - X_i) \right]^2 K \left( \frac{|X_j - X_i|^2}{h^2} \right), \quad (2.8)$$

where a kernel  $K(\cdot)$  is positive and supported on  $[0, 1]$ , so that the weights of all points  $X_j$  outside a spherical neighborhood  $U_h(X_i)$  of diameter  $h$  around  $X_i$  vanish. The solution to this quadratic optimization problem can be represented as

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

where  $X_{ij} = X_j - X_i$ . As many other nonparametric estimates, the estimate (2.8) suffers from the data sparseness for large  $d$ . This phenomenon is often referred to as *curse of dimensionality*. Indeed, one has to select the bandwidth  $h$  in a way to provide at least  $d + 1$  design points in every (or almost every) spherical neighborhood  $U_h(X_i)$ . For the case of a random design with a positive density, this implies that a bandwidth  $h$  of order  $n^{-1/d}$  or even larger should be taken. For large  $d$  this leads to a very poor rate  $n^{-1/d}$  in estimation of  $F$ , and the same applies to the estimation of the vectors  $\beta_\ell^*$  (see Proposition 5.1 below).

At the same time, suppose for a moment that we know the mapping  $T : \mathbb{R}^d \rightarrow \mathbb{R}^m$ . Then we could use this information for estimating the  $m$ -dimensional link function  $g_0$  and its gradient  $\nabla g_0$ . This also provides an estimate of the gradient  $F(x) = T^\top \nabla g_0(Tx)$  of much better accuracy, which corresponds to an  $m$ -dimensional nonparametric problem on the “true” index space, instead of the original  $d$ -dimensional nonparametric estimate  $\hat{F}(x)$ . More specifically, a function  $f(x)$  of the form (2.7) remains constant when  $x$  varies in any direction orthogonal to the  $m$ -dimensional subspace  $\mathcal{I}$ . The above considerations leads to another estimate:

$$\begin{aligned} \begin{pmatrix} \hat{f}(X_i) \\ \hat{F}(X_i) \end{pmatrix} &= \underset{c \in \mathbb{R}, b \in \mathbb{R}^d}{\operatorname{arginf}} \sum_{j=1}^n \left[ Y_j - c - b^\top (X_j - X_i) \right]^2 K \left( \frac{|T(X_j - X_i)|^2}{h^2} \right) \\ &= \left\{ \sum_{j=1}^n \begin{pmatrix} 1 \\ X_{ij} \end{pmatrix} \begin{pmatrix} 1 \\ X_{ij} \end{pmatrix}^\top K \left( \frac{|TX_{ij}|^2}{h^2} \right) \right\}^{-1} \sum_{j=1}^n Y_j \begin{pmatrix} 1 \\ X_{ij} \end{pmatrix} K \left( \frac{|TX_{ij}|^2}{h^2} \right) \end{aligned}$$

The latter estimate of  $F(X_i)$  is based on averaging over a narrow cylinder  $\{x : |T(x - X_i)| \leq h\}$ , centered at  $X_i$ , which spans  $\mathcal{I}^\perp$ . This allows to use an essentially smaller bandwidth  $h$  and still have enough design points in every such neighborhood. On the other hand, the smaller bandwidth would decrease drastically the bias of estimation. Unfortunately this “ideal” estimate cannot be implemented in practice since it requires



the explicit knowledge of the target index space  $\mathcal{I}$ . A natural idea is to substitute the mapping  $T$  by its pilot estimate. This leads to the following *structural adaptation* approach. We proceed iteratively starting with the estimates  $\widehat{\beta}_\ell = \sum_{i=1}^n \widehat{F}(X_i) \psi_\ell(X_i)$ ,  $\ell = 1, \dots, L$  based on the fully nonparametric gradient estimate  $\widehat{F}$  with some  $h = h_1$ , see (2.8). Although this estimate is very rough, it contains some information about the structure of the model function  $f$  and, in particular, about the mapping  $T$ : all vectors  $\widehat{\beta}_\ell$  up to the estimation error, belong to the index space  $\mathcal{I}$ . This information can be used for producing another, more careful estimate of the gradient function and hence, of the vectors  $\beta_\ell^*$ . More precisely, let  $\widehat{\mathcal{B}}_1$  be the matrix composed from the vectors  $\widehat{\beta}_\ell$ ,  $\ell = 1, \dots, L$ . We define the gradient estimate  $\widehat{F}_2(X_i)$  at  $X_i$  by a local linear fit using the elliptic neighborhood  $\{x : |S_2(x - X_i)| \leq h_2\}$ , with  $S_2 = (I + \rho_2^{-2} \widehat{\mathcal{B}}_1 \widehat{\mathcal{B}}_1^\top)^{-1/2}$  for some  $\rho_2 < 1$  and  $h_2 > h_1$  (instead of the spherical windows  $\{x : |x - X_i| \leq h_1\}$ ). In other words, we shrink the original windows in all the directions  $\widehat{\beta}_\ell$  (since  $\rho_2 < 1$ ) and stretch them in all the orthogonal directions (since  $h_2 > h_1$ ):

$$\begin{aligned} \begin{pmatrix} \widehat{f}_2(X_i) \\ \widehat{F}_2(X_i) \end{pmatrix} &= \underset{c \in \mathbb{R}, b \in \mathbb{R}^d}{\operatorname{arginf}} \sum_{j=1}^n \left[ Y_j - c - b^\top (X_j - X_i) \right]^2 K \left( \frac{|S_2(X_j - X_i)|^2}{h_2^2} \right) \\ &= \left\{ \sum_{j=1}^n \begin{pmatrix} 1 \\ X_{ij} \end{pmatrix} \begin{pmatrix} 1 \\ X_{ij} \end{pmatrix}^\top K \left( \frac{|S_2 X_{ij}|^2}{h_2^2} \right) \right\}^{-1} \sum_{j=1}^n Y_j \begin{pmatrix} 1 \\ X_{ij} \end{pmatrix} K \left( \frac{|S_2 X_{ij}|^2}{h_2^2} \right). \end{aligned}$$

This leads to the estimates  $\widehat{\beta}_{2,\ell} = \frac{1}{n} \sum_{i=1}^n \widehat{F}_2(X_i) \psi_\ell(X_i)$  of  $\beta_\ell^*$  producing the matrix  $\widehat{\mathcal{B}}_2$ . We continue this way each time compressing the averaging windows in the direction of the current estimate  $\widehat{\mathcal{B}}_k$  and expanding them in orthogonal directions.

The results presented below show that this procedure allows to estimate the index space  $\mathcal{I}$  at the rate  $n^{-1/2}$  provided that  $m < 4$ .

### 3 Estimation algorithm

We now present the description of the method. The whole estimation procedure (we refer to it as Algorithm 1) is carried out in two basic steps: estimation of the vectors  $\beta_\ell^*$  and estimation the e.d.r. matrix  $\mathcal{R}^*$ . Below we discuss each step separately.

#### 3.1 Estimation of $\beta_\ell^*$ 's

The procedure involves input parameters  $h_1 < h_{\max}$  and  $\rho_{\min} < \rho_1$ , so that  $\rho$  decreases geometrically from  $\rho_1$  to  $\rho_{\min}$  by the factor  $a_\rho$  and  $h$  increases geometrically from  $h_1$  to  $h_{\max}$  by the factor  $a_h$  during iterations. The choice of these parameters as well as

the set of basis functions  $\{\psi_\ell\}$  will be discussed in the next section. The algorithm reads as follows:

**1 Initialization:** specify parameters  $\rho_1, \rho_{\min}, a_\rho, h_1, h_{\max}, a_h$  and the set of functions  $\{\psi_\ell\}$ ; set  $k = 1, \widehat{\mathcal{B}}_0 = 0$ ;

**2 Compute**  $S_k = \left(I + \rho_k^{-2} \widehat{\mathcal{B}}_{k-1} \widehat{\mathcal{B}}_{k-1}^\top\right)^{1/2}$ ;

**3 For every**  $i = 1, \dots, n$ , compute  $\widehat{F}_k(X_i)$  from the expression:

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

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

**4 Compute the vectors**  $\widehat{\beta}_{k,\ell} = \frac{1}{n} \sum_{i=1}^n \widehat{F}_k(X_i) \psi_\ell(X_i)$ ,  $\ell = 1, \dots, L$  and compose the matrix  $\widehat{\mathcal{B}}_k$  with columns  $\widehat{\beta}_{k,1}, \dots, \widehat{\beta}_{k,L}$ ;

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

By  $k(n)$  we denote the total number of iterations. The estimates  $\widehat{\beta}_{k(n),\ell}$  from the last iteration are used as the final estimates of  $\beta_\ell^*$ .

### 3.2 Computing the effective dimension reduction matrix

Let  $\widehat{\mathcal{B}}$  be an estimate of the matrix  $\mathcal{B}^*$  obtained by the previously described iterative procedure. We will see (Theorem 5.3) that this matrix estimates the target matrix  $\mathcal{B}^*$  with a reasonable accuracy but it is typically of the rank  $d$  and hence, it does not provide any dimension reduction. We estimate the effective dimension reduction matrix  $\mathcal{R}^*$  using the singular value decomposition of  $\widehat{\mathcal{B}}$  in place of  $\mathcal{B}^*$ , cf. (2.6). Namely, the product  $\widehat{\mathcal{B}}^\top \widehat{\mathcal{B}}$ , being symmetric and non-negative, can be represented in the form  $\widehat{\mathcal{B}}^\top \widehat{\mathcal{B}} = \widehat{O} \widehat{\Lambda} \widehat{O}^\top$  with the orthogonal  $L \times L$ -matrix  $\widehat{O}$  and the diagonal matrix  $\widehat{\Lambda}$ :  $\widehat{\Lambda} = \text{diag}\{\widehat{\lambda}_1, \dots, \widehat{\lambda}_L\}$  with non-increasing eigenvalues  $\widehat{\lambda}_1 \geq \dots \geq \widehat{\lambda}_L \geq 0$  (the squared singular values of  $\widehat{\mathcal{B}}$ ). The estimate  $\mathcal{R}_m$  of the true e.d.r. matrix  $\mathcal{R}^*$  from (2.6) is defined by

$$\mathcal{R}_m = (\widehat{\mathcal{B}} \widehat{O}_m)^\top \quad (3.1)$$

where  $\widehat{O}_m$  is the submatrix of  $\widehat{O}$  composed of its first  $m$  columns.

### 3.3 Choice of parameters of the algorithm

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. The values  $\rho_k$  decrease from  $\rho_1$  to  $\rho_{\min}$  while  $h_k$  increase during iteration from  $h_1$  to  $h_{\max}$ . The value  $h_1$  is to be selected in such a way that for every (or almost every) point  $X_i$ , the estimate  $\widehat{F}(X_i)$  is well defined. A necessary (and usually sufficient) condition is that every ball  $\{x : |x - X_i| \leq h_1\}$  contains at least  $d + 1$  design points; see the modified procedure in the next section for more discussion. Concerning the values  $h$  and  $\rho$  at the last iteration  $k(n)$ , the optimization of the risk of this estimate (see Corollary 5.1 in Section 5.3) leads to the following recommendation: the value  $h_{k(n)}$  for the last iteration should be possibly large that is, about  $\sqrt{d}$ , and then the value  $\rho_{\min}$  should be selected possibly small, but still providing enough design points in every or almost every local ellipsoidal neighborhoods  $E_k(X_i) = \{x : |S_k(x - X_i)| \leq h_k\}$ . For the case with  $m \leq 3$ , we propose the following rule of thumb:

$$\begin{aligned} \rho_1 &= 1, & \rho_{\min} &= n^{-1/3}, & a_\rho &= e^{-1/6}, \\ h_1 &= C_0 n^{-\frac{1}{4\sqrt{d}}}, & h_{\max} &= C_0, & a_h &= e^{\frac{1}{2(4\sqrt{d})}}, \end{aligned} \quad (3.2)$$

where  $C_0 \geq 1$  is to be defined depending on the design, see the modified procedure for a proposal.

The proposed rule leads to  $k(n) \approx 6 \log(\rho_1/\rho_{\min}) \approx 2 \log n$  iterations and provides that  $h_{k(n)} \approx h_{\max}$ . Note also that assuming the structure of the matrix  $\widehat{B}_{k-1} \widehat{B}_{k-1}^\top$  to follow the structure of the target matrix  $\mathcal{M}^*$ , neighborhood  $E_k(X_i)$  is stretched at each iteration step by factor  $a_h$  in all directions and is shrunk by factor  $a_\rho$  in directions of the  $m$ -dimensional index space  $\mathcal{I}$ . Therefore, the Lebesgue measure of every such neighborhood is changed each time by the factor  $e^{\frac{d}{2(4\sqrt{d})} - \frac{m}{6}}$  which is larger or equal to 1 for all  $m \leq 3$  and  $d > m$ . Under the assumption of a random design with a positive density, this would lead to an increase of the mean number of design points inside each  $E_k(X_i)$ .

### 3.4 Choice of functions $\psi_\ell$

As we have mentioned already, the choice of test functions  $\{\psi_\ell, \ell = 1, \dots, L\}$  is of primary importance for the practical efficiency of the proposed procedure. The main constraint on the set  $\{\psi_\ell\}$  is that the matrix  $\mathcal{B}^*$  is of the same rank as  $T$  and that the function  $g$  from the equivalent representation (2.7) is sufficiently smooth, see Assumption 3 below. It can be easily shown that the “ideal” choice of the set  $\{\psi_\ell\}$  can be obtained by orthogonalization of the components  $F_j = \partial f / \partial x_j$ ,  $j = 1, \dots, d$  of the

gradient  $F$ . This “ideal” collection of functions  $\psi_\ell$  would contain only  $m$  elements. Of course, this choice cannot be realized since it involves the unknown regression function  $f$ .

Note next that the functions (vectors)  $\psi_1, \dots, \psi_L$  form an orthonormal system in  $\mathbb{R}^n$  and  $\beta_\ell^*$  is the scalar product of the gradient  $F$  and the basis function  $\psi_\ell$ . The sum

$$F_L = \sum_{\ell=1}^L \beta_\ell^* \psi_\ell$$

is the projection of the gradient  $F$  on the linear subspace in  $\mathbb{R}^n$  spanned by  $\{\psi_\ell\}$ . One can easily check that  $\mathcal{M}_L^* = \sum_{i=1}^n F_L(X_i) F_L(X_i)^\top$ . Thus to prevent the loss of information due to the substitution of  $\mathcal{M}$  for  $\mathcal{M}_L$ , the set  $\{\psi_\ell\}$  should be selected rich enough. Our proposal is to define  $\{\psi_\ell\}$  by orthogonalizing the set of all polynomials  $x_{\ell_1} \dots x_{\ell_q}$  of the coordinate functions for some  $q \geq 1$  and all  $1 \leq \ell_1 \leq \dots \leq \ell_q$ .

A suitable alternative, especially for large  $d$ , is a basis system constructed by orthogonalizing a fully nonparametric estimate of the gradient.

## 4 Implementation and simulated results

In this section we illustrate the performance of the proposed algorithm on some simulated examples. First we discuss a slight modification of the procedure which allows to relax design assumptions.

### 4.1 Modified procedure

In Algorithm 1, at each step, we use a linear combination of the estimated gradient vectors  $\widehat{F}(X_i)$  as the estimate of the vector  $\beta_\ell^*$ . To guarantee some useful properties of this procedure, the estimates  $\widehat{F}(X_i)$  should be well defined, which in turn requires some local regularity of the design in the corresponding neighborhood of the point  $X_i$ , see Assumption 4 in Section 5. If such a condition is not satisfied even at a few points, then the corresponding gradient estimates would have a very large standard deviation which may deteriorate the quality of the index estimates  $\widehat{\beta}_\ell$ . We can avoid this problem by weighting each summand in the expression for  $\widehat{\beta}_{k,\ell}$  with some coefficients which express the degree of local regularity of the design. This leads to the modified procedure which is presented below.

**1 Initialization:** specify parameters  $\rho_1, \rho_{\min}, a_\rho, h_1, h_{\max}, a_h, C_w$  and the set

of functions  $\{\psi_\ell\}$ ; Define  $\bar{w}$  as the square root of the minimal eigenvalue of the matrix  $\bar{\mathcal{V}}$  with

$$\bar{\mathcal{V}} = \frac{1}{\mathbf{E}K(\zeta^\top \zeta)} \mathbf{E} \begin{pmatrix} 1 \\ \zeta \end{pmatrix} \begin{pmatrix} 1 \\ \zeta \end{pmatrix}^\top K(\zeta^\top \zeta)$$

where  $\zeta$  is random and uniformly distributed over the ball  $B_1 = \{x \in \mathbb{R}^d : |x| \leq 1\}$ :  $\bar{w}^2 = \lambda_{\min}(\bar{\mathcal{V}})$ ; set  $k = 1$ ,  $\hat{\mathcal{B}}_0 = 0$ ;

**2** Compute  $\hat{\mathcal{M}}_k = \hat{\mathcal{B}}_{k-1} \hat{\mathcal{B}}_{k-1}^\top$ . If  $\|\hat{\mathcal{M}}_k\| > 1$ , then normalize it by its maximal eigenvalue:  $\hat{\mathcal{M}}_k := \hat{\mathcal{M}}_k / \|\hat{\mathcal{M}}_k\|$ ; Set  $S_k = \left(I + \rho_k^{-2} \hat{\mathcal{M}}_k\right)^{1/2}$ ;

**3** For every  $i = 1, \dots, n$ , compute the matrix  $\hat{\mathcal{V}}_k(X_i)$  with

$$\hat{\mathcal{V}}_k(X_i) = \frac{1}{\sum_{j=1}^n K(W_{ij,k}^\top W_{ij,k})} \sum_{j=1}^n \begin{pmatrix} 1 \\ W_{ij,k} \end{pmatrix} \begin{pmatrix} 1 \\ W_{ij,k} \end{pmatrix}^\top K(W_{ij,k}^\top W_{ij,k})$$

where  $W_{ij,k} = h_k^{-1} S_k(X_j - X_i)$  and define  $w_i$  as the square root of the minimal eigenvalue of  $\hat{\mathcal{V}}_k(X_i)$ :  $w_i^2 = \lambda_{\min}(\hat{\mathcal{V}}_k(X_i))$ ;

**4** If the condition

$$\frac{1}{n} \sum_{i=1}^n w_i \geq C_w \bar{w}$$

is not fulfilled, then increase  $h_k$  by the factor  $a_h$ , that is,  $h_k := a_h h_k$ . If  $h_k > h_{\max}$ , then terminate, otherwise repeat from Step 3;

**5** For every  $i = 1, \dots, n$ , compute  $\hat{F}_k(X_i)$ :

$$\begin{pmatrix} \hat{f}_k(X_i) \\ \hat{F}_k(X_i) \end{pmatrix} = \left\{ \sum_{j=1}^n \begin{pmatrix} 1 \\ X_{ij} \end{pmatrix} \begin{pmatrix} 1 \\ X_{ij} \end{pmatrix}^\top K\left(\frac{|S_k X_{ij}|^2}{h_k^2}\right) \right\}^{-1} \sum_{j=1}^n Y_j \begin{pmatrix} 1 \\ X_{ij} \end{pmatrix} K\left(\frac{|S_k X_{ij}|^2}{h_k^2}\right);$$

**6** For every  $\ell = 1, \dots, L$ , compute the vector  $\hat{\beta}_{k,\ell}$

$$\hat{\beta}_{k,\ell} = \left( \sum_{i=1}^n w_i \right)^{-1} \sum_{i=1}^n \hat{F}_k(X_i) \psi_\ell(X_i) w_i$$

with the previously obtained  $w_i$ 's. Compose the matrix  $\hat{\mathcal{B}}_k$  with columns  $\hat{\beta}_{k,\ell}$ ,  $\ell = 1, \dots, L$ .

**7** Set  $\rho_{k+1} = a_\rho \rho_k$ , and  $h_{k+1} = a_h h_k$ . If  $\rho_{k+1} \geq \rho_{\min}$ , then set  $k = k + 1$  and continue with Step 2.

The last iteration estimate  $\hat{\mathcal{B}} = \hat{\mathcal{B}}_{k(n)}$  will be used for the dimension-reduction step.

## 4.2 Simulated datasets

In our simulation study we apply the modified procedure with the following parameter setting:

$$\begin{aligned} \rho_1 &= 1, & \rho_{\min} &= n^{-1/3}, & a_\rho &= e^{-1/6}, \\ h_1 &= n^{-\frac{1}{4\sqrt{d}}}, & h_{\max} &= 2\sqrt{d}, & a_h &= e^{\frac{1}{2(4\sqrt{d})}}. \end{aligned}$$

We also set  $C_w = 2^{-1/2}$ . In case of high dimensionality, i.e.  $d > 20$  a smaller value of  $C_w$  was necessary to guarantee the existence of valid bandwidths  $h_k$ . The basis system  $\{\psi_\ell\}$  is obtained by orthogonalization of the set of functions  $\{1, x_j, x_j x_k, j, k = 1, \dots, d\}$ . This setting leads to the number of iterations  $k(n) \approx \frac{\log(\rho_1/\rho_{\min})}{\log a_\rho} = 2 \log n$ .

The performance of the method is illustrated by means of the following examples. We consider the model  $Y_i = g(X_i^\top \theta_1, \dots, X_i^\top \theta_m)$  for  $m = 1$  and  $m = 2$ . The design  $X_1, \dots, X_n$  is modelled randomly with independent components so that every component of  $(X_i + 1)/2$  follows  $B(1, \tau)$ -distribution. The parameter  $\tau$  controls the skewness of the beta-distribution with  $\tau = 1$  corresponding to the uniform design. We also set

$$m = 1: g(u) = u \sin(\sqrt{5}u) \text{ and } \theta = (1, 2, 0, \dots, 0)^\top / \sqrt{5}.$$

$$\begin{aligned} m = 2: g(u_1, u_2) &= (u_1^3 + u_2)(u_1 - u_2^3) \text{ and } \theta_1 = (1, 1, 1, 0, \dots, 0)^\top / \sqrt{3}, \\ \theta_2 &= (1, -1, 0, \dots, 0)^\top / \sqrt{2}. \end{aligned}$$

The first situation corresponds essentially to example 8.2 from Li (1992). The procedure utilizes the biweight kernel  $K(|x|^2) = (1 - |x|^2)_+^2$ . The quality of estimation is measured using the criterion  $\|\mathcal{R}^*(I - \mathcal{P}_m)\|_2$ , where  $\mathcal{P}_m$  is the projector on the estimated index space  $\widehat{\mathcal{L}}$ , see Section 5.2 for more details.

Our objective is to illustrate the following features of the procedure:

- how the quality of estimation improves during iteration;
- dependence on the sample size  $n$  and the dimensionality  $d$ .
- how the results depend on skewness of the design and on the error variance  $\sigma^2$ .

We also compare the performance of our iterative procedure to a one step estimate with an “ideal” choice of the bandwidth, minimizing the criterion for the situation at hand. This bandwidth was selected in a separate small simulation from a grid of bandwidths, complying to the condition in step 4 of the procedure.

Figure 1 illustrates the quality of estimation of the index space for  $m = 1$ ,  $d = 10$ ,  $n = 400$  and  $\sigma = .1$ , providing the “best” view obtained by a one step estimate (left) and the view gained from our procedure (right).

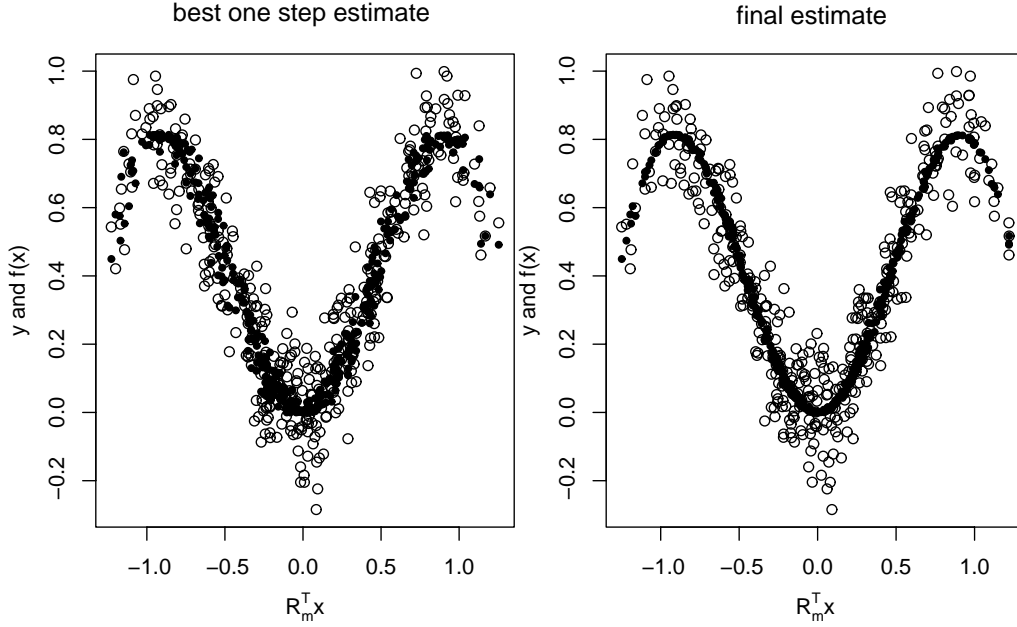


Figure 1: Best view for a one step estimate (left) and view from the last iteration (right) for  $g(u) = u \sin(\sqrt{5}u)$ ,  $m = 1$ ,  $d = 10$ ,  $n = 400$  and  $\sigma = .1$ . Values of  $y$  and  $f(x)$  are indicated by  $\circ$  and  $\bullet$  respectively.

The simulation results for different dimensionality  $d$  and sample size  $n$  are given in Tables 1 and 2. All simulation results show a considerable gain using the proposed iterative method. This gain increases drastically as the dimensionality  $d$  grows. The results from Table 2 for  $d = 10$  and different  $\sigma$ -values clearly illustrate the bias-variance trade-off. For the first step estimate as well as for the “best” one-step estimate the bias dominates and the quality of estimation only weakly depends on the noise variance while for our procedure the bias is essentially reduced during iteration and the final quality of estimation is proportional to the standard deviation  $\sigma$ . We also observe that the procedure performs stable in case of moderate error variance and design asymmetry.

The box-plots in Figure 2 provide some information about the distribution of the criterion  $\sqrt{n}\|\mathcal{R}^*(I - \hat{\mathcal{P}}_m)\|_2/\|\mathcal{R}^*\|_2$  for the “best” one step estimate and after the first, second, fourth, eighth and final iteration for  $d = 10$ ,  $m = 2$  and different sample size  $n$ . Results displayed are obtained from  $N = 250$  simulations. The results confirm the root- $n$  consistence of the final estimate as claimed by Theorem 5.1 from Section 5. Note that the losses even being multiplied by  $\sqrt{n}$  are still slightly improved with growing  $n$ .

Table 1: Case  $m = 1$ : mean loss  $\|\mathcal{R}^*(I - \widehat{\mathcal{P}}_m)\|_2 / \|\mathcal{R}^*\|_2$  for the “best” one step estimate and the first, second, fourth, eighth and final iteration. Results are obtained from  $N = 250$  simulations ( $N = 100$  in case of  $d > 10$ ). The interquartile range of the losses is given in parentheses.

$d$	$n$	$\sigma$	$\tau$	best	1st	2nd	4th	8th	final
3	200	0.1	1	0.0442 (0.032)	0.0508 (0.038)	0.0419 (0.031)	0.0359 (0.026)	0.0271 (0.019)	0.0236 (0.014)
4	200	0.1	1	0.0558 (0.034)	0.0606 (0.033)	0.0484 (0.024)	0.0417 (0.025)	0.0339 (0.02)	0.0309 (0.018)
6	200	0.1	1	0.0807 (0.036)	0.0829 (0.034)	0.0631 (0.024)	0.0536 (0.024)	0.0437 (0.02)	0.0389 (0.018)
10	100	0.1	1	0.343 (0.14)	0.341 (0.14)	0.208 (0.083)	0.146 (0.067)	0.105 (0.047)	0.0903 (0.04)
10	200	0.1	1	0.172 (0.066)	0.173 (0.065)	0.109 (0.036)	0.0854 (0.026)	0.0646 (0.02)	0.0537 (0.017)
10	400	0.1	1	0.101 (0.029)	0.103 (0.031)	0.0698 (0.024)	0.0573 (0.019)	0.0438 (0.015)	0.0369 (0.012)
10	800	0.1	1	0.0619 (0.019)	0.0642 (0.019)	0.0479 (0.015)	0.0409 (0.013)	0.032 (0.011)	0.0271 (0.0084)

Table 2: Case  $m = 2$ : mean loss  $\|\mathcal{R}^*(I - \widehat{\mathcal{P}}_m)\|_2 / \|\mathcal{R}^*\|_2$  for the “best” one step estimate and the first, second, fourth, eighth and final iteration. Results are obtained from  $N = 250$  simulations ( $N = 100$  in case of  $d > 10$ ). The interquartile range of the losses is given in parentheses.

$d$	$n$	$\sigma$	$\tau$	best	1st	2nd	4th	8th	final
3	200	0.1	1	0.0248 (0.018)	0.0248 (0.018)	0.0186 (0.013)	0.017 (0.013)	0.0154 (0.012)	0.0148 (0.01)
4	200	0.1	1	0.0445 (0.024)	0.0445 (0.024)	0.0306 (0.017)	0.0265 (0.013)	0.0219 (0.011)	0.0195 (0.011)
6	200	0.1	1	0.0925 (0.035)	0.0933 (0.034)	0.0732 (0.029)	0.0653 (0.028)	0.0583 (0.022)	0.0477 (0.02)
10	100	0.1	1	0.361 (0.092)	0.361 (0.092)	0.25 (0.081)	0.209 (0.072)	0.179 (0.065)	0.153 (0.06)
10	200	0.1	1	0.203 (0.047)	0.203 (0.047)	0.132 (0.037)	0.112 (0.033)	0.086 (0.031)	0.0559 (0.017)
10	400	0.1	1	0.123 (0.031)	0.124 (0.032)	0.0739 (0.021)	0.0605 (0.019)	0.0425 (0.012)	0.0269 (0.0077)
10	800	0.1	1	0.0749 (0.019)	0.0749 (0.019)	0.048 (0.013)	0.0403 (0.011)	0.0285 (0.0073)	0.0155 (0.0042)
20	800	0.1	1	0.189 (0.026)	0.191 (0.026)	0.127 (0.024)	0.107 (0.021)	0.0743 (0.017)	0.0281 (0.0076)
50	800	0.1	1	0.647 (0.18)	0.654 (0.14)	0.389 (0.058)	0.313 (0.042)	0.229 (0.039)	0.0729 (0.0095)
10	400	0.05	1	0.123 (0.032)	0.122 (0.032)	0.0706 (0.02)	0.0572 (0.018)	0.0379 (0.013)	0.017 (0.0051)
10	400	0.2	1	0.13 (0.031)	0.132 (0.032)	0.087 (0.022)	0.0731 (0.019)	0.0573 (0.015)	0.0506 (0.015)
10	400	0.1	0.75	0.118 (0.029)	0.119 (0.031)	0.0813 (0.023)	0.0758 (0.024)	0.0591 (0.025)	0.0234 (0.0074)
10	400	0.1	1.5	0.12 (0.032)	0.12 (0.034)	0.0918 (0.026)	0.0902 (0.027)	0.0826 (0.033)	0.0382 (0.014)



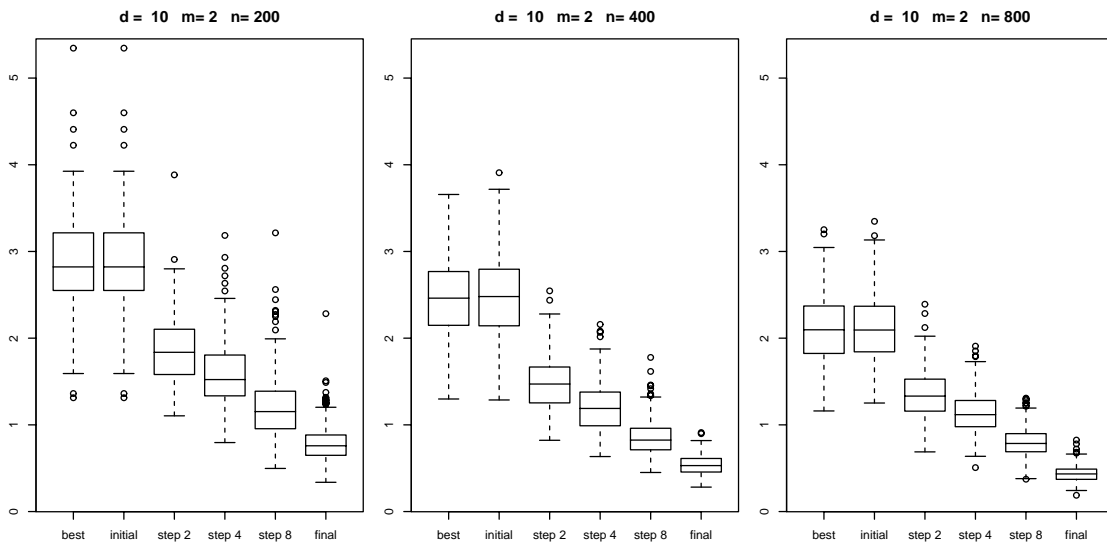


Figure 2: Simulation results in terms of  $\sqrt{n}\|\mathcal{R}^*(I - \widehat{\mathcal{P}}_m)\|_2/\|\mathcal{R}^*\|_2$  for  $m = 2$ ,  $d = 10$  and  $n = 200, 400, 800$  for the “best” one step estimate and the first, second, fourth, eighth and final iteration

## 5 Main results

In this section we present some results describing the properties of the previously introduced procedure.

### 5.1 Assumptions

We consider the following assumptions:

**Assumption 1. (Kernel)** The kernel  $K(\cdot)$  is continuously differentiable, monotonously 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. (Link function)** The function  $g$  from (2.7) is two times differentiable with a bounded second derivative, so that, for some constants  $C_g$  and for all  $u, v \in \mathbb{R}^m$ , it holds

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

Our last assumption concerns the design properties. In what follows we assume deterministic design, 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  $X_1, \dots, X_n$  i.i.d. random points in  $\mathbb{R}^d$  with a design density  $p(x)$ . Then all the result should be understood conditionally on the design.

In order Algorithm 1 to work, we have to suppose that the design points  $(X_i)$  are “well diffused” and, as a consequence, all the matrices  $V_k(X_i)$  are well defined.

The estimation procedure utilizes the matrices  $S_k$  with  $S_k^2 = I + \rho_k^{-2} \widehat{\mathcal{B}}_{k-1} \widehat{\mathcal{B}}_{k-1}^\top$  where  $\widehat{\mathcal{B}}_{k-1}$  is the estimate of the matrix  $\mathcal{B}^*$  constructed at the preceding iteration step. We also introduce an ‘ideal’ matrix  $S_k^* = (I + \rho_k^{-2} \mathcal{B}^* (\mathcal{B}^*)^\top)^{1/2}$  and define the matrix

$$U_k = (S_k^*)^{-1} S_k^2 (S_k^*)^{-1}.$$

This matrix  $U_k$  characterizes the accuracy of estimating the matrix  $\mathcal{B}^*$  by  $\widehat{\mathcal{B}}_{k-1}$ . If  $\widehat{\mathcal{B}}_{k-1} = \mathcal{B}^*$ , then  $U_k = I$ . We shall see that these matrices  $U_k$  are typically close to  $I$ . Define now, given a matrix  $U$  and  $k \leq k(n)$

$$\begin{aligned} Z_{ij,k} &= h_k^{-1} S_k^* (X_j - X_i), \quad i, j = 1, \dots, n, \\ N_{i,k}(U) &= \sum_{j=1}^n K(Z_{ij,k}^\top U Z_{ij,k}), \quad i = 1, \dots, n, \\ \mathcal{V}_{i,k}(U) &= \sum_{j=1}^n \begin{pmatrix} 1 \\ Z_{ij,k} \end{pmatrix} \begin{pmatrix} 1 \\ Z_{ij,k} \end{pmatrix}^\top K(Z_{ij,k}^\top U Z_{ij,k}), \quad i = 1, \dots, n. \end{aligned}$$

Our design assumption means in particular that the  $(d+1) \times (d+1)$ -matrices  $\mathcal{V}_{i,k}(U)$  are well defined for all  $U$  close to  $I$  and for all  $i \leq n$ .

We use below the notation  $\|A\|$  for the sup-norm of  $A$ :  $\|A\| = \sup_\lambda |A\lambda|/|\lambda|$ .

**Assumption 4. (Design)** There exist constants  $C_V$ ,  $C_K$ ,  $C_{K'}$  and some  $\alpha > 0$ , such that for all matrices  $U$  satisfying  $\|U - I\| \leq \alpha$  and for all  $k \leq k(n)$  the following conditions hold:

(1) the inverse matrices  $\mathcal{V}_{i,k}(U)^{-1}$  are well defined and

$$N_{i,k}(U) \|\mathcal{V}_{i,k}(U)^{-1}\| \leq C_V, \quad i = 1, \dots, n;$$

(2) For  $j = 1, \dots, n$

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

Here  $K'$  means the derivative of the kernel  $K$ .

**Remark 5.1** One can easily checked that for the case of a random design with a continuous positive density, one can fix some constant  $C_V$ ,  $C_K$  and  $C_{K'}$  depending on the dimension  $d$  and design density only and such that the conditions from Assumption 4 are fulfilled with a high probability converging exponentially fast to 1 as  $n$  grows.

In what follows by  $C, C_1, C_2, \dots$  we denote generic constants depending on  $d, C_g, C_V, C_K, C_{K'}, \psi_\ell, L$  and  $\sigma$  only.

## 5.2 Loss of information caused by estimated e.d.r.

An important characteristic of the estimated e.d.r.  $\mathcal{R}_m$  is the loss of information caused by this reduction. Due to the representation (2.7), the information contained in a unit vector  $v \in \mathbb{R}^d$  can be measured by the value  $|\mathcal{R}^*v|$ . A loss of information occurs if  $|\mathcal{R}^*v| > 0$  but  $|\mathcal{R}_m v| = 0$ . Let  $\Pi^*$  be the projector in  $\mathbb{R}^d$  onto the true index space  $\mathcal{I}$  and similarly,  $\mathcal{P}_m$  denote the projector in  $\mathbb{R}^d$  onto the estimated index space  $\widehat{\mathcal{I}}$  corresponding to the e.d.r.  $\mathcal{R}_m$ , that is  $\widehat{\mathcal{I}} = \text{Im } \mathcal{R}_m^\top$ . Then the total loss of information by e.d.r.  $\mathcal{R}_m$  can be measured by the value

$$\|\mathcal{R}^*(I - \mathcal{P}_m)\|_2$$

where  $\|A\|_2$  means the Euclidean norm of the matrix  $A$ , that is,  $\|A\|_2^2 = \text{tr } AA^\top = \text{tr } A^\top A$ . In the sequel we use the following obvious inequalities:  $\|A\| \leq \|A\|_2 \leq \sqrt{m}\|A\|$  where  $m$  is the rank of  $A$ .

The next result claims that the loss of information caused by the e.d.r.  $\mathcal{R}_m$  is of order  $n^{-1/2}$ .

**Theorem 5.1** *Let  $\mathcal{R}_m$  be defined by (3.1). For  $m \leq 3$ , there exists a sequence  $\varkappa_n \rightarrow 0$  as  $n \rightarrow \infty$  such that under Assumptions 1 through 4, it holds for sufficiently large  $n$  and every  $z \geq 1$ :*

$$\begin{aligned} \mathbf{P} \left( \|\mathcal{R}_m(I - \Pi^*)\|_2 > \frac{2zH_1}{\sqrt{n}} + Cz_n^2n^{-2/3} \right) &< ze^{-(z^2-1)/2} + \frac{3k(n)}{n}, \\ \mathbf{P} \left( \|\mathcal{R}^*(I - \mathcal{P}_m)\|_2 > \frac{2zH_1}{\sqrt{n(1-\varkappa_n)}} + Cz_n^2n^{-2/3} \right) &< ze^{-(z^2-1)/2} + \frac{3k(n)}{n}, \end{aligned}$$

with  $z_n = (1 + 2 \log n + 2 \log \log n)^{1/2}$  and

$$\begin{aligned} H_1 &= \sqrt{2} \sigma C_V C_K \bar{\psi} \sqrt{L}, \\ \bar{\psi} &= \max_{i=1, \dots, n} \max_{\ell=1, \dots, L} |\psi_\ell(X_i)|. \end{aligned} \tag{5.1}$$

### 5.3 Estimation of the index space

By construction,  $\mathcal{R}^*$  is an orthogonal mapping from  $\mathbb{R}^d$  to  $\mathbb{R}^m$ , that is,  $\mathcal{R}^*(\mathcal{R}^*)^\top$  is diagonal  $m \times m$ -matrix with the diagonal elements  $\lambda_1, \dots, \lambda_m$ . Moreover, the product  $\Pi^* = (\mathcal{R}^*)^\top (\mathcal{R}^*(\mathcal{R}^*)^\top)^{-1} \mathcal{R}^*$  is the projector in  $\mathbb{R}^d$  onto the corresponding index space  $\mathcal{I}$ . Similarly  $\mathcal{P}_m = \mathcal{R}_m(\mathcal{R}_m \mathcal{R}_m^\top)^{-1} \mathcal{R}_m^\top$  is the projector onto the estimated e.d.r. space. Thus the quality of the identification of the true index space can be qualified with the error of estimating  $\Pi^*$  with  $\mathcal{P}_m$ . We encounter the following identifiability problem: if, for instance, the last eigenvalue  $\lambda_m$  is (close to) zero, then the corresponding eigenvector  $e_m$  is not uniquely defined. The next result states that if the eigenvalue  $\lambda_m$  is separated away from zero, the estimated projector  $\mathcal{P}_m$  recovers  $\Pi^*$  at the rate  $n^{-1/2}$ .

**Theorem 5.2** *Let  $m \leq 3$  and Assumptions 1 through 4 hold. For  $n$  sufficiently large,*

$$P \left( \|\Pi^* - \mathcal{P}_m\|_2 > \frac{2\sqrt{2}\lambda_m^{-1/2} z H_1}{\sqrt{n(1-\varkappa_n)}} + C z_n^2 n^{-2/3} \right) \leq z e^{-(z^2-1)/2} + \frac{3k(n)}{n}$$

with  $\varkappa_n$  and  $H_1$  from Theorem 5.1.

### 5.4 Estimation of the matrix $\mathcal{B}^*$

In this section we present some results describing the quality of estimating the vectors  $\beta_\ell^*$  by the proposed estimation procedure. The first result describes the accuracy of the first step estimate, and the next result describes the quality of the final estimate.

#### 5.4.1 The first-step approximation

Let  $\widehat{\beta}_{1,\ell}$ ,  $\ell = 1, \dots, n$  be the family of the estimates obtained at the first step of the iterative procedure with  $\rho_1 = 1$ ,  $S_1 = I$  and some  $h_1$ .

**Proposition 5.1** *Under Assumptions 1 through 4, it holds for every  $\ell \leq L$*

$$\widehat{\beta}_{1,\ell} - \beta_\ell^* = C_{1,\ell} h_1 + \frac{\xi_{1,\ell}}{h_1 \sqrt{n}}$$

where  $C_{1,\ell}$  is a constant and  $\xi_{1,\ell}$  is a zero mean normal random vector in  $\mathbb{R}^d$  satisfying

$$\begin{aligned} C_{1,\ell} &\leq \sqrt{2} C_g C_V \overline{\psi}_\ell, \\ \mathbf{E}|\xi_{1,\ell}|^2 &\leq 2\sigma^2 C_V^2 C_K^2 \overline{\psi}_\ell^2. \end{aligned}$$

**Remark 5.2** The optimization of the risk of the first step estimate under the constraint  $h_1 \geq \text{Const. } h^{-1/d}$  leads to the following rule for the choice of  $h_1$ :  $h_1 = \text{Const. } n^{-\frac{1}{4\sqrt{d}}}$ . Hence, we get the accuracy for  $\widehat{\beta}_{1,\ell}$ :

$$|\widehat{\beta}_{1,\ell} - \beta_\ell^*| \leq \text{Const. } n^{-(\frac{1}{4} \wedge \frac{1}{d})}.$$

#### 5.4.2 Accuracy of the final estimate

Let  $\widehat{\beta}_\ell$ 's be the estimates of  $\beta_\ell^*$ 's obtained at the last iteration,  $\ell = 1, \dots, L$ . As previously,  $\widehat{\mathcal{B}}$  denotes the matrix composed by the vectors  $\widehat{\beta}_\ell$ . It turns out that the quality of estimation delivered by  $\widehat{\mathcal{B}}$  is not homogeneous w.r.t. to the orientation in the space  $\mathbb{R}^d$ . This heterogeneity is caused by application of elliptic windows for estimating the gradient vectors  $F(X_i)$ . To mimic this property, we introduce for every  $k \leq k(n)$  an operator ( $d \times d$ -matrix)  $P_{\rho_k}^* = (I + \rho_k^{-2} \mathcal{B}^* (\mathcal{B}^*)^\top)^{-1/2} = (S_k^*)^{-1}$  which, roughly speaking, multiply by the factor  $\rho_k$  within the index space  $\mathcal{I}$  while, being restricted to the orthogonal subspace  $\mathcal{I}^\perp$ , it coincides with the identity mapping.

**Theorem 5.3** *Let  $m \leq 3$  and Assumptions 1 through 4 hold. There exist a Gaussian zero mean random  $d \times L$ -matrix  $\xi^* \in \mathbb{R}^{dL}$  such that, with  $\rho = \rho_{k(n)}$  and  $n$  large enough*

$$\mathbf{P} \left( \left\| P_\rho^* (\widehat{\mathcal{B}} - \mathcal{B}^*) - \frac{\xi^*}{\sqrt{n}} \right\|_2 > C_1 z_n^2 n^{-2/3} \right) \leq \frac{3k(n) - 1}{n}$$

and

$$\mathbf{E} \|\xi^*\|_2^2 \leq 2\sigma^2 \overline{\psi}^2 L C_V^2 C_K^2 = H_1^2.$$

**Corollary 5.1** *Under the conditions of Theorem 5.3, for every  $z \geq 1$*

$$\mathbf{P} \left( \left\| P_\rho^* (\widehat{\mathcal{B}} - \mathcal{B}^*) \right\|_2 > \frac{z H_1}{\sqrt{n}} + C_1 z_n^2 n^{-2/3} \right) \leq z e^{-(z^2-1)/2} + \frac{3k(n) - 1}{n}.$$

## 6 Conclusions and outlook

We introduce a new method of dimension reduction based on the idea of structural adaptation. The method applies for a very broad class of regression models 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 delivers the optimal rate  $n^{-1/2}$  of estimating the index space provided that the effective dimensionality of the model is not large 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 works stable in high dimensional situations and for a non-regular design.

The procedure can be easily extended to the situation with a multivariate response variable  $Y \in \mathbb{R}^p$  with  $p > 1$ . The underlying multi-index assumption remains of the same functional form:  $\mathbf{E}(Y | X) = f(x) = g(X^\top \theta_1, \dots, X^\top \theta_m)$  where  $g$  is a vector function on  $\mathbb{R}^m$  with values in  $\mathbb{R}^p$ . This means that the gradient  $F_j = \nabla f_j$  of each component  $f_j$  of  $f$  belongs to the index space spanned by vectors  $\theta_1, \dots, \theta_m$  and one can utilize the same ideas as previously for estimating the index space  $\mathcal{I}$ . The only difference is that the basis functions  $\{\psi_\ell\}$  should also be vectors in  $\mathbb{R}^p$ . A reasonable example corresponds to the procedure which estimates for every component  $f_j$ ,  $j = 1, \dots, p$ , of the regression function  $f \in \mathbb{R}^p$  the vectors  $\beta_{1,j}^*, \dots, \beta_{L,j}^*$  with

$$\beta_{\ell,j}^* = \sum_{i=1}^n F_j(X_i) \psi_\ell(X_i), \quad \ell = 1, \dots, L,$$

and the same  $\psi_\ell$ 's and then utilizes the total collections of the vectors  $\{\widehat{\beta}_{\ell,j}\}$  with  $\ell = 1, \dots, L$  and  $j = 1, \dots, p$  for estimating the index space  $\mathcal{I}$ .

One more open question corresponds to the case of the unknown effective dimension  $m$ . This immediately leads to the following two problems: estimation of  $m$  and testing a  $m$ -index hypothesis. An important feature of the proposed iterative procedure is that it does not rely on the specific value of  $m$ . One can therefore expect that the matrix  $\widehat{\mathbf{B}}$  coming from the last step of the algorithm, can be used for answering the above mentioned problems.

Another interesting issue arises when considering multiple time series and especially financial data. We regard such extensions as topics for further research.

## 7 Appendix A: Proofs

Here we collect the proofs of the assertions formulated previously.

### 7.1 Proof of Lemma 2.1

Let  $o_k$  denote the  $k$ -th column of  $O$ . Then  $\theta_k = \mathcal{B}^* o_k$  is the eigenvector of  $\mathcal{M}_L^*$  with the eigenvalue  $\lambda'_k$ ,  $k \leq d$ . Indeed, with e.g.  $k = 1$ ,

$$\mathcal{M}_L^* \theta_1 = \mathcal{B}^* (\mathcal{B}^*)^\top \mathcal{B}^* o_1 = \mathcal{B}^* O \Lambda_L O^\top o_1 = \lambda'_1 \mathcal{B}^* o_1 = \lambda'_1 \theta_1. \quad (7.1)$$

Here we have used that the matrix  $O$  is orthogonal and hence,  $O^\top o_1 = (1, 0, \dots, 0)^\top$  and  $O\Lambda_L O^\top o_1 = O(\lambda'_1, 0, \dots, 0)^\top = \lambda'_1 o_1$ . Under condition (2.3) this implies that  $\lambda'_1 = \lambda_1$  and  $\theta_1$  is a multiple of  $e_1$ , and similarly for other  $k$ . In addition,

$$\theta_k^\top \theta_{k'} = o_k(\mathcal{B}^*)^\top \mathcal{B}^* o_{k'} = o_k O\Lambda_L O^\top o_{k'} = \lambda_k \delta_{kk'}, \quad k, k' = 1, \dots, L,$$

that is, the vectors  $\theta_k$  are orthogonal to each other and satisfy  $|\theta_k|^2 = \lambda_k$ ,  $k = 1, \dots, L$  and the assertion follows.

We now present the proofs of Proposition 5.1 and Theorems 5.1 through 5.3. All these results are based on the following technical assertion describing an improvement of the estimate  $\widehat{\mathcal{B}}$  at each iteration step.

## 7.2 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 fixed  $d \times L$ -matrix  $B$  which can be viewed as an approximation  $\widehat{\mathcal{B}}_{k-1}$  of  $\mathcal{B}^*$  obtained at the previous step. Set also

$$\begin{aligned} S_B &= \left( I + \rho^{-2} B B^\top \right)^{1/2}, \\ V_B(X_i) &= \sum_{j=1}^n \begin{pmatrix} 1 \\ X_{ij} \end{pmatrix} \begin{pmatrix} 1 \\ X_{ij} \end{pmatrix}^\top K \left( \frac{|S_B X_{ij}|^2}{h^2} \right) \\ \begin{pmatrix} \widehat{f}_B(X_i) \\ \widehat{F}_B(X_i) \end{pmatrix} &= V_B(X_i)^{-1} \sum_{j=1}^n Y_j \begin{pmatrix} 1 \\ X_{ij} \end{pmatrix} K \left( \frac{|S_B X_{ij}|^2}{h^2} \right) \end{aligned} \quad (7.2)$$

$$\widehat{\beta}_{B,\ell} = \frac{1}{n} \sum_{i=1}^n \widehat{F}_B(X_i) \psi_\ell(X_i) \quad (7.3)$$

where, recall,  $X_{ij} = X_j - X_i$ , and define the matrix  $\widehat{\mathcal{B}}_B$  with columns  $\widehat{\beta}_{B,\ell}$ ,  $\ell = 1, \dots, L$ . We aim to evaluate the estimation errors  $\widehat{\mathcal{B}}_B - \mathcal{B}^*$ . To describe the results, we introduce the matrix (linear operator)  $P_\rho^* = (I + \rho^{-2} \mathcal{B}^* (\mathcal{B}^*)^\top)^{-1/2}$ . Define also for some positive  $\delta < \rho/4$ , the set  $\mathfrak{B}_{\delta,\rho}$  by

$$\mathfrak{B}_{\delta,\rho} = \left\{ B : \|P_\rho^*(B - \mathcal{B}^*)\|_2 \leq \delta \right\}.$$

**Proposition 7.1** *Let Assumptions 1 through 4 hold. Then there exists Gaussian random  $d \times L$ -matrix  $\xi$  such that it holds with  $\alpha = 2\delta/\rho + \delta^2/\rho^2$*

$$\mathbf{P} \left( \sup_{B \in \mathfrak{B}_{\delta,\rho}} \left\| P_\rho^*(\widehat{\mathcal{B}}_B - \mathcal{B}^*) - \frac{\xi}{h\sqrt{n}} \right\|_2 > \frac{\sqrt{2} C_g C_V \bar{\psi} \sqrt{L}}{(1-\alpha)^{3/2}} h \rho^2 + \frac{\sigma \bar{\psi} \sqrt{L} C_{\alpha,n} \alpha}{h\sqrt{n}} \right) \leq 2/n$$

where

$$C_{\alpha,n} = \frac{1}{2} \left( \frac{\sqrt{2} C_V C_{K'}}{(1-\alpha)^2} + \frac{2\sqrt{2} C_V^2 C_{K'} C_K}{(1-\alpha)^3} \right) \left( 2 + \sqrt{(3+dL) \log(4n)} \right) \quad (7.4)$$

and

$$\mathbf{E} \|\xi\|_2^2 \leq 2\sigma^2 C_V^2 C_K^2 \bar{\psi}^2 L. \quad (7.5)$$

Before prove this statement, we present one straightforward corollary.

**Corollary 7.1** *It holds under Assumptions 1 through 4 for every  $z \geq 1$*

$$\mathbf{P} \left( \sup_{B \in \mathfrak{B}_{\delta,\rho}} \|P_\rho^*(\hat{\mathcal{B}}_B - \mathcal{B}^*)\|_2 > \bar{\psi} \sqrt{L} \left( \frac{\sqrt{2} C_g C_V h \rho^2}{(1-\alpha)^{3/2}} + \frac{z\sqrt{2} \sigma C_V C_K}{h\sqrt{n}} + \frac{\sigma C_{\alpha,n}\alpha}{h\sqrt{n}} \right) \right) \leq z e^{-(z^2-1)/2} + 2/n.$$

Indeed, the Gaussian vector  $\xi \in \mathbb{R}^{dL}$  fulfills with every  $z \geq 1$

$$\mathbf{P} \left( \|\xi\|_2 \geq z \sqrt{\mathbf{E} \|\xi\|_2^2} \right) \leq z e^{-(z^2-1)/2}$$

see Lemma 9 in HJS98, and the assertion follows from Proposition 7.1.

**Proof of Proposition 7.1:** We follow the line of the proof of Proposition 2 in HJS98 and focus here only on the essential points omitting technical details.

It is useful to define

$$u = \rho^{-1} P_\rho^* B, \quad U = P_\rho^* \left( I + \rho^{-2} B B^\top \right) P_\rho^* = (P_\rho^*)^2 + u u^\top$$

and similarly

$$u^* = \rho^{-1} P_\rho^* \mathcal{B}^*, \quad U^* = P_\rho^* \left( I + \rho^{-2} \mathcal{B}^* (\mathcal{B}^*)^\top \right) P_\rho^* = I$$

so that  $u, u^*$  are  $d \times L$ -matrices and  $U, U^*$  are  $d \times d$  symmetric matrices. Clearly  $B = \mathcal{B}^*$  implies  $U = I$  and the condition  $\|B - \mathcal{B}^*\|_2 \leq \delta$  implies  $\|u - u^*\|_2 \leq \delta/\rho$ , that is, the inclusion  $B \in \mathfrak{B}_{\delta,\rho}$  is equivalent to  $u \in \{u : \|u - u^*\|_2 \leq \delta/\rho\}$ . Due to Lemma 7.7 from Appendix B it also follows  $\|U - U^*\| = \|u u^\top - u^* (u^*)^\top\| \leq \alpha = 2\delta/\rho + \delta^2/\rho^2$  for all such  $u$ .

Next, for every  $i, j \leq n$ , define

$$\begin{aligned} Z_{ij} &= h^{-1} (P_\rho^*)^{-1} (X_j - X_i), \\ \mathcal{V}_i(U) &= \sum_{j=1}^n \begin{pmatrix} 1 \\ Z_{ij} \end{pmatrix} \begin{pmatrix} 1 \\ Z_{ij} \end{pmatrix}^\top K(Z_{ij}^\top U Z_{ij}), \\ \hat{s}_i(U) &= h^{-1} \mathcal{V}_i(U)^{-1} \sum_{j=1}^n \begin{pmatrix} 1 \\ Z_{ij} \end{pmatrix} Y_j K(Z_{ij}^\top U Z_{ij}). \end{aligned}$$



It is easy to check that  $\widehat{s}_i(U) = \begin{pmatrix} h^{-1} \widehat{f}_B(X_i) \\ P_\rho^* \widehat{F}_B(X_i) \end{pmatrix}$  and hence,

$$P_\rho^* \widehat{\beta}_{B,\ell} = \mathcal{E}_d n^{-1} \sum_{i=1}^n \widehat{s}_i(U) \psi_\ell(X_i)$$

where  $\mathcal{E}_d$  denotes the projector from  $\mathbb{R}^{d+1}$  onto  $\mathbb{R}^d$  keeping the last  $d$  coordinates.

The model equation (1.2) implies

$$\widehat{s}_i(U) = s_i(U) + \zeta_i(U)$$

with

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

so that

$$P_\rho^* (\widehat{\beta}_{B,\ell} - \beta_\ell^*) = \frac{1}{n} \sum_{i=1}^n \{ \mathcal{E}_d s_i(U) - P_\rho^* F(X_i) \} \psi_\ell(X_i) + \mathcal{E}_d n^{-1} \sum_{i=1}^n \zeta_i(U) \psi_\ell(X_i).$$

Clearly  $\xi_\ell(U) = \mathcal{E}_d n^{-1} \sum_{i=1}^n \zeta_i(U) \psi_\ell(X_i)$  is for every  $U$  a linear combination of the Gaussian errors  $\varepsilon_i$  and therefore it is also a Gaussian vector in  $\mathbb{R}^d$ . We define  $\xi(U)$  as  $d \times L$  matrix with columns  $\xi_\ell(U)$  and set  $\xi = \xi(U^*)$ . It is easy to see that the following three statements imply the desirable result:

$$\sup_{u: \|u-u^*\|_2 \leq \delta/\rho} |\mathcal{E}_d s_i(U) - P_\rho^* F(X_i)| \leq \frac{\sqrt{2} C_g C_V}{(1-\alpha)^{3/2}} h \rho^2, \quad i = 1, \dots, n, \quad (7.6)$$

$$\mathbf{P} \left( \sup_{u: \|u-u^*\|_2 \leq \delta/\rho} \|\xi(U) - \xi(U^*)\|_2 > \frac{\sigma C_{\alpha,n} \alpha}{h \sqrt{n}} \right) \leq 2/n \quad (7.7)$$

with  $U = (P_\rho^*)^2 + uu^\top$  and  $U^* = I$ , and for all  $\ell = 1, \dots, L$

$$\mathbf{E} |\xi_\ell(U^*)|^2 \leq \frac{2\sigma^2 C_V^2 C_K^2 \overline{\psi}_\ell^2}{h^2 n}. \quad (7.8)$$

To check these statements, the following lemma will be useful.

**Lemma 7.1** *Let  $\|U - I\| \leq \alpha < 1$ . Then for all  $i, j$  with  $Z_{ij}^\top U Z_{ij} \leq 1$ , it holds  $|Z_{ij}|^2 \leq (1 - \alpha)^{-1}$ .*

**Proof.** Note that the inequalities  $Z_{ij}^\top U Z_{ij} \leq 1$  and  $\|U - I\| \leq \alpha$  imply

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

and hence  $|Z_{ij}|^2 \leq (1 - \alpha)^{-1} Z_{ij}^\top U Z_{ij}$ . ■

First we evaluate the “bias” term  $\mathcal{E}_d s_i(U) - P_\rho^* F(X_i)$ . Since

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

it holds

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

where in view of (2.7)

$$r_{ij} = g(\mathcal{R}^* X_j) - g(\mathcal{R}^* X_i) - (\mathcal{R}^* X_j - \mathcal{R}^* X_i)^\top g'(\mathcal{R}^* X_i).$$

The use of  $P_\rho^* \mathcal{B}^* (\mathcal{B}^*)^\top P_\rho^* = \rho^2 (I - (P_\rho^*)^2)$  and  $\|I - (P_\rho^*)^2\| \leq 1$  provide

$$\begin{aligned} |(\mathcal{B}^*)^\top X_j - (\mathcal{B}^*)^\top X_i|^2 &= (X_j - X_i)^\top \mathcal{B}^* (\mathcal{B}^*)^\top (X_j - X_i) \\ &= ((P_\rho^*)^{-1} (X_j - X_i))^\top P_\rho^* \mathcal{B}^* (\mathcal{B}^*)^\top P_\rho^* (P_\rho^*)^{-1} (X_j - X_i) \\ &= h^2 \rho^2 Z_{ij}^\top (I - (P_\rho^*)^2) Z_{ij} \\ &\leq h^2 \rho^2 |Z_{ij}|^2 \end{aligned}$$

which also implies

$$|\mathcal{R}^* X_j - \mathcal{R}^* X_i| = |(\mathcal{B}^* O_m)^\top X_j - (\mathcal{B}^* O_m)^\top X_i| \leq h^2 \rho^2 |Z_{ij}|^2.$$

This yields by Lemma 7.1 and Assumption 3 for every pair  $(i, j)$  with  $Z_{ij}^\top U Z_{ij} \leq 1$ :

$$|r_{ij}| \leq \frac{C_g h^2 \rho^2}{1 - \alpha}, \quad 1 + |Z_{ij}|^2 \leq 1 + \frac{1}{1 - \alpha} \leq \frac{2}{1 - \alpha}$$

and using Assumptions 4 we bound

$$\begin{aligned}
|\mathcal{E}_d s_i(U) - P_\rho^* F(X_i)| &\leq h^{-1} \left| \mathcal{V}_i(U)^{-1} \sum_{j=1}^n \left( \frac{1}{Z_{ij}} \right) r_{ij} K(Z_{ij}^\top U Z_{ij}) \right| \\
&\leq \frac{C_g h \rho^2}{1 - \alpha} \|\mathcal{V}_i(U)\|^{-1} \left| \sum_{j=1}^n (1 + |Z_{ij}|^2)^{1/2} K(Z_{ij}^\top U Z_{ij}) \right| \\
&\leq \sqrt{2}(1 - \alpha)^{-3/2} C_g C_V h \rho^2
\end{aligned}$$

and (7.6) follows.

Further we study the stochastic components  $\xi_\ell(U)$ . It follows directly from the definition that there are vector coefficients  $c_{i,\ell}(U)$  such that

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

We now apply the following two technical results from HJS98, see Lemma 3, 10 there for a particular case with  $L = 1$  and  $\psi_\ell \equiv 1$ . Extension to general  $L$  and  $\psi_\ell$ 's is straightforward.

**Lemma 7.2** *It holds*

$$(i) \quad \sum_{i=1}^n |c_{i,\ell}(U^*)|^2 \leq \frac{2C_V^2 C_K^2 \bar{\psi}_\ell^2}{h^2 n};$$

$$(ii) \quad \sup_{U: \|U-I\| \leq \alpha} \sum_{i=1}^n |c_{i,\ell}(U)|^2 \leq \frac{2C_V^2 C_K^2 \bar{\psi}_\ell^2}{(1 - \alpha)h^2 n};$$

(iii) *For every unit vector  $e \in \mathbb{R}^d$*

$$\sup_{U: \|U-I\| \leq \alpha} \left\| \frac{d}{dU} e^\top c_{i,\ell}(U) \right\| \leq \frac{\kappa_\alpha \bar{\psi}_\ell}{nh}$$

with

$$\kappa_\alpha = \sqrt{2}(1 - \alpha)^{-3/2} C_V C_{K'} + 2\sqrt{2}(1 - \alpha)^{-5/2} C_V^2 C_{K'} C_K.$$

(iv) *For every unit vector  $e \in \mathbb{R}^d$*

$$\sup_{u: \|u-u^*\|_2 \leq \delta/\rho} \left\| \frac{d}{du} e^\top c_{i,\ell}(U) \right\| \leq \frac{\kappa'_\alpha \bar{\psi}_\ell}{nh}$$

with  $U = U_u = (P_\rho^*)^2 + uu^\top$  and

$$\kappa'_\alpha = \kappa_\alpha (1 - \alpha)^{-1/2} = \sqrt{2}(1 - \alpha)^{-2} C_V C_{K'} + 2\sqrt{2}(1 - \alpha)^{-3} C_V^2 C_{K'} C_K.$$

**Lemma 7.3** *Let  $r \geq 0$  and let vector-functions  $a_i(u)$  with  $u \in \mathbb{R}^p$  obey the conditions*

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

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

$$\mathbf{P} \left( \sup_{|u-u^*| \leq r} \frac{1}{\sqrt{n}} \left| \sum_{i=1}^n \{a_i(u) - a_i(u^*)\} \varepsilon_i \right| > \sigma \kappa r \left( 2 + \sqrt{(3+p) \log(4n)} \right) \right) \leq \frac{2}{n}.$$

Lemma 7.2, (i) implies (7.8). The statement (7.7) follows from Lemma 7.2, (iv), and Lemma 7.3 applied to the matrix  $\xi(U) \in \mathbb{R}^{dL}$  with columns  $\xi_\ell(U)$  and with  $U = U_u = (P_\rho^*)^2 + uu^\top$ , for details see again HJS98.

### 7.3 Proof of Theorem 5.3

To be able to apply Proposition 7.1 to the estimates  $\widehat{\beta}_{k,\ell}$  at step  $k$ , we need that the matrix  $B = \widehat{\mathcal{B}}_{k-1}$  coming as the result of the preceding iteration belongs to the set  $\mathfrak{B}_{\rho,\delta}$  with  $\rho = \rho_k$  and some  $\delta < \rho/4$ . Since the matrix  $\widehat{\mathcal{B}}_{k-1}$  is random, we have to check that the probability of the event  $\{\widehat{\mathcal{B}}_{k-1} \in \mathfrak{B}_{\rho_k,\delta}\} = \{B : \|P_\rho^*(B - \mathcal{B}^*)\|_2 \leq \rho\}$  is sufficiently large. Further we show that this property is fulfilled if  $n$  is large enough.

Let the numbers  $h_k$  and  $\rho_k$  be shown in the algorithm description,  $k = 1, \dots, k(n)$ . Define successively values  $\delta_k$  and  $\alpha_k$ ,  $k = 1, \dots, k(n)$  by  $\alpha_1 = 0$  and

$$\begin{aligned} \delta_k &= \bar{\psi} \sqrt{L} \left( \frac{\sqrt{2} C_g C_V}{(1 - \alpha_k)^{3/2}} h_k \rho_k^2 + \frac{\sqrt{2} \sigma C_V C_K z_n}{h_k \sqrt{n}} + \frac{\sigma C_{\alpha_k, n} \alpha_k}{2 h_k \sqrt{n}} \right), \\ \alpha_{k+1} &= \rho_{k+1}^{-2} (2\delta_k \rho_k + \delta_k^2) \end{aligned}$$

where  $z_n = (1 + 2 \log n + 2 \log \log n)^{1/2}$ .

**Lemma 7.4** *For  $m \leq 3$  and  $n$  sufficiently large, the values  $\alpha_k$ 's fulfill  $\max_{k \leq k(n)} \alpha_k < 1/4$ . In addition, for the last iteration  $k(n)$ , it holds*

$$\mu_n := \bar{\psi} \sqrt{L} \left( \frac{\sqrt{2} C_g C_V}{(1 - \alpha_{k(n)})^{3/2}} h_{k(n)} \rho_{k(n)}^2 + \frac{\sigma C_{\alpha_{k(n)}, n} \alpha_{k(n)}}{h_{k(n)} \sqrt{n}} \right) \leq C_1 z_n^2 n^{-2/3}.$$

**Proof.** See Lemma 5 in HJS98. ■

Next, successive application of the results of Propositions 7.1 and Corollary 7.1 with  $z_n = (1 + 2 \log n + 2 \log \log n)^{1/2}$  leads to the following

**Lemma 7.5** *Let  $n$  be sufficiently large. There exists random sets  $\mathcal{A}_1 \supseteq \dots \supseteq \mathcal{A}_{k(n)}$  such that*

$$\mathbf{P}(\mathcal{A}_k) \geq 1 - \frac{3k}{n}$$

and it holds on  $\mathcal{A}_k$

$$\|P_{\rho_{k+1}}^*(\widehat{\mathcal{B}}_k - \mathcal{B}^*)\|_2 \leq \delta_k, \quad k = 1, \dots, k(n) - 1.$$

**Proof.** See Lemma 6 in HJS98. ■

Now the result of Theorem 5.3 can be proved by one more application of Proposition 7.1 to the last step estimate  $\widehat{\mathcal{B}} = \widehat{\mathcal{B}}_{k(n)}$  with  $h = h_{k(n)} \geq 1$  and  $\rho = \rho_{k(n)} \approx n^{-1/3}$ , see again HJS98 for the detailed derivation.

## 7.4 Proof of Theorem 5.1

Let  $\widehat{\mathcal{B}}$  be the last step estimate of the matrix  $\mathcal{B}^*$ . We know from Theorem 5.3 that, with probability close to one,  $\widehat{\mathcal{B}}$  fulfills the conditions

$$\|P_{\rho}^*(\widehat{\mathcal{B}} - \mathcal{B}^*)\|_2 \leq \tau, \quad (7.9)$$

with  $\rho = \rho_{k(n)}$  and some small  $\tau$ . This implies by Lemma 7.7 from Appendix B

$$\|\widehat{\mathcal{B}} - \Pi^* \widehat{\mathcal{B}}\|_2 \leq \tau \quad (7.10)$$

where  $\Pi^*$  denotes the projector on the index space  $\mathcal{I}$ .

Recall that  $\widehat{\mathcal{B}}$  approximates the  $d \times L$ -matrix  $\mathcal{B}^*$  of rank  $m$ . However, it is typically of rank  $d$ . It is useful to introduce another  $d \times L$ -matrix  $\mathcal{B}_m$  of rank  $m$  which minimizes the expression  $\|\widehat{\mathcal{B}} - \mathcal{B}_m\|_2$  over all such matrices. The solution to this optimization problem can be described explicitly via the diagonal decomposition of the matrix  $\widehat{\mathcal{B}}^\top \widehat{\mathcal{B}} = \widehat{O} \widehat{\Lambda}_L \widehat{O}^\top$  with an orthogonal matrix  $\widehat{O}$  and a diagonal matrix  $\widehat{\Lambda}_L$  with non increasing eigenvalues, cf. Lemma 2.1. We use the notation  $I_m$  for the diagonal  $L \times L$ -matrix with the first  $m$  diagonal elements equal to 1 and the remaining ones equal to zero.

**Lemma 7.6 (cf. Harville (1997, Theorem 21.12.4))** *The  $d \times L$ -matrix  $\mathcal{B}_m = \widehat{\mathcal{B}} \widehat{O} I_m \widehat{O}^\top$  minimizes the norm  $\|B - \widehat{\mathcal{B}}\|_2$  over all  $d \times L$ -matrices  $B$  of rank  $m$ :*

$$\mathcal{B}_m = \widehat{\mathcal{B}} \widehat{O} I_m \widehat{O}^\top = \operatorname{arg\,inf}_{B \in \mathfrak{B}_m} \| \widehat{\mathcal{B}} - B \|_2 \quad (7.11)$$

where  $\mathfrak{B}_m$  denotes the set of  $d \times L$ -matrices of rank  $m$ .

**Proof.** Let  $\widehat{\mathcal{B}}^\top \widehat{\mathcal{B}} = \widehat{O} \widehat{\Lambda}_L \widehat{O}^\top$ . Then it holds for the  $d \times L$ -matrix  $\widetilde{\mathcal{B}} = \widehat{\mathcal{B}} \widehat{O}$

$$\widetilde{\mathcal{B}}^\top \widetilde{\mathcal{B}} = O^\top \widehat{\mathcal{B}}^\top \widehat{\mathcal{B}} \widehat{O} = \widehat{O}^\top \widehat{O} \widehat{\Lambda}_L \widehat{O}^\top \widehat{O} = \widehat{\Lambda}_L$$

that is, the columns of the matrix  $\widetilde{\mathcal{B}}$  are orthogonal and they are ranged in a way that their norms decrease. This clearly implies

$$\operatorname{arg\,inf}_{B \in \mathfrak{B}_m} \|\widetilde{\mathcal{B}} - B\|_2 = \widetilde{\mathcal{B}} I_m$$

and the assertion of the lemma follows by usual change-of-basis argument.  $\blacksquare$

Recall that we define the e.d.r. matrix  $\mathcal{R}_m$  by  $\mathcal{R}_m = (\widehat{\mathcal{B}} \widehat{O}_m)^\top$ , see (3.1). It follows from the last lemma that  $\mathcal{R}_m = (\mathcal{B}_m \widehat{O}_m)^\top$ . Also, (7.10) and the definition of  $\mathcal{B}_m$  (see (7.11)) imply

$$\|\widehat{\mathcal{B}} - \mathcal{B}_m\|_2 \leq \|\widehat{\mathcal{B}} - \Pi^* \widehat{\mathcal{B}}\|_2 \leq \tau,$$

and, since  $\|P_\rho^*\| \leq 1$ ,

$$\|P_\rho^*(\mathcal{B}_m - \mathcal{B}^*)\|_2 \leq \|\widehat{\mathcal{B}} - \mathcal{B}_m\|_2 + \|P_\rho^*(\widehat{\mathcal{B}} - \mathcal{B}^*)\|_2 \leq 2\tau. \quad (7.12)$$

This implies by Lemma 7.8 from Appendix B

$$\|P_{\rho,m}(\mathcal{B}_m - \mathcal{B}^*)\|_2 \leq 2\tau(1 - \varkappa)^{-1/2} \quad (7.13)$$

where  $P_{\rho,m} = (I + \rho^{-2} \mathcal{B}_m \mathcal{B}_m^\top)^{-1/2}$  and  $\varkappa = 4\tau/\rho + 4\tau^2/\rho^2$ . Now the result of Theorem 5.1 is a straightforward application of Theorem 5.3 and Lemma 7.9 from Appendix B.

## 7.5 Proof of Theorem 5.2

Let  $\widehat{\mathcal{B}}$  be the last step estimate of the matrix  $\mathcal{B}^*$ . We know from Theorem 5.3 that, with probability close to one,  $\widehat{\mathcal{B}}$  fulfills the condition (7.9) with  $\rho = \rho_{k(n)}$  and some small  $\tau$ . Next, let the matrices  $\widehat{\mathcal{B}}_m$ , and  $\mathcal{R}_m$  of rank  $m$  be defined as in the proof of Theorem 5.1 so that the condition (7.12) is fulfilled. The projectors  $\Pi^*$  and  $\mathcal{P}_m$  are defined as

$$\begin{aligned} \Pi^* &= (\mathcal{R}^*)^\top \left( \mathcal{R}^* (\mathcal{R}^*)^\top \right)^{-1} \mathcal{R}^*, \\ \mathcal{P}_m &= \mathcal{R}_m^\top \left( \mathcal{R}_m \mathcal{R}_m^\top \right)^{-1} \mathcal{R}_m. \end{aligned}$$

The use of Lemma 7.11 of Appendix B provides

$$\|\Pi^* - \mathcal{P}_m\|_2 \leq \sqrt{2} \lambda_m^{-1/2} 2\tau(1 - 4\tau/\rho - 4\tau^2/\rho^2)^{-1/2}$$

and we end up as in the proof of Theorem 5.1.

## Appendix B: Some matrix inequalities

Let  $B$  and  $B_1$  be two  $d \times L$ -matrices and  $\rho$  be some positive number. Define the  $d \times d$ -matrix  $P_\rho$  as

$$P_\rho = \left( I + \rho^{-2} B B^\top \right)^{-1/2}.$$

Here we collect some facts which can be obtained from the inequality

$$\|P_\rho(B_1 - B)\|_2 \leq \delta \tag{7.14}$$

with some small  $\delta \geq 0$ . Here and in what follows  $\|A\|_2$  denotes the  $L_2$ -norm of the matrix  $A$ , i.e.  $\|A\|_2^2 = \text{tr} A A^\top$ , and  $\|A\|$  is the sup-norm:  $\|A\| = \sup_{v \in \mathbb{R}^d} |Av|/|v|$ .

**Lemma 7.7** *The condition (7.14) implies*

$$\left\| P_\rho \left( B B^\top - B_1 B_1^\top \right) P_\rho \right\| \leq 2\rho\delta + \delta^2.$$

**Proof.** Since

$$\|P_\rho B\|^2 = \|P_\rho B B^\top P_\rho\| = \left\| \left( I + \rho^{-2} B B^\top \right)^{-1} B B^\top \right\| \leq \rho^2$$

(7.14) yields

$$\begin{aligned} & \left\| P_\rho \left( B_1 B_1^\top - B B^\top \right) P_\rho \right\| \\ & \leq 2 \left\| P_\rho (B_1 - B) B^\top P_\rho \right\| + \left\| P_\rho (B_1 - B) (B_1 - B)^\top P_\rho \right\| \\ & \leq 2 \left\| P_\rho (B_1 - B) \right\|_2 \|P_\rho B\| + \left\| P_\rho (B_1 - B) \right\|_2^2 \\ & \leq 2\delta\rho + \delta^2 \end{aligned}$$

as required. ■

Define also

$$P_{\rho,1} = \left( I + \rho^{-2} B_1 B_1^\top \right)^{-1/2}.$$

**Lemma 7.8** *Let  $B$  and  $B_1$  fulfill (7.14) for some  $\delta < \rho/4$ . Then*

$$\left\| P_{\rho,1} (B - B_1) \right\|_2 \leq \frac{\delta}{\sqrt{1 - 2\delta/\rho - \delta^2/\rho^2}}.$$

**Proof.** Let  $\alpha = 2\delta/\rho + \delta^2/\rho^2$ . By Lemma 7.7

$$\|P_\rho P_{\rho,1}^{-2} P_\rho - I\| = \rho^{-2} \|P_\rho (BB^\top - B_1 B_1^\top) P_\rho\| \leq \alpha$$

and hence,

$$\begin{aligned} \|P_{\rho,1}^{-1} P_\rho\|^2 &= \|P_\rho P_{\rho,1}^{-2} P_\rho\| \leq 1 + \alpha, \\ \|P_{\rho,1} P_\rho^{-1}\|^2 &= \|(P_\rho P_{\rho,1}^{-2} P_\rho)^{-1}\| \leq (1 - \alpha)^{-1}. \end{aligned}$$

Now

$$\begin{aligned} \|P_{\rho,1}(B - B_1)\|_2 &= \|P_{\rho,1} P_\rho^{-1} P_\rho (B - B_1)\|_2 \\ &\leq \|P_{\rho,1} P_\rho^{-1}\| \|P_\rho (B - B_1)\|_2 \leq \|P_{\rho,1} P_\rho^{-1}\| \delta \leq \delta (1 - \alpha)^{-1/2}. \end{aligned}$$

■

Next we consider the situation when both matrices  $B$  and  $B_1$  are of rank  $m$  with some  $m < d$ . By  $\Pi$  we denote the projector in  $\mathbb{R}^d$  onto the subspace  $\mathcal{L} = \text{Im } B$ . Similarly  $\Pi_1$  is the projector in  $\mathbb{R}^d$  onto the subspace  $\mathcal{L}_1 = \text{Im } B_1$ .

**Lemma 7.9** *Let  $d \times L$ -matrices  $B$  and  $B_1$  of rank  $m$  fulfill  $\|P_\rho(B - B_1)\|_2 \leq \delta$ . Then it holds*

$$\|(I - \Pi)B_1\|_2 \leq \delta.$$

**Proof.** Since  $P_\rho$  is the unity operator within the subspace  $\mathcal{L}^\perp = \text{Im}(I - \Pi)$ , it easily follows  $(I - \Pi)P_\rho = I - \Pi$  (this fact is obvious when  $BB^\top$  and hence  $P_\rho$  is a diagonal matrix, and the general case can be reduced to that one by an orthogonal transform). Since also  $(I - \Pi)B = 0$ , we derive

$$\begin{aligned} B_1 &= (\Pi + I - \Pi)B_1 \\ &= \Pi B_1 + (I - \Pi)(B_1 - B) \\ &= \Pi B_1 + (I - \Pi)P_\rho(B_1 - B) \end{aligned}$$

so that  $\|(I - \Pi)B_1\|_2 \leq \|P_\rho(B_1 - B)\|_2 \leq \delta$ .

■

**Lemma 7.10** *Let  $\Pi$  and  $\Pi_1$  be two projectors in  $\mathbb{R}^d$  of rank  $m < d$ . Then*

$$\|\Pi_1 - \Pi\|_2 = \sqrt{2} \|\Pi(I - \Pi_1)\|_2.$$



**Proof.** Note first that since  $\Pi$  and  $I - \Pi$  are orthogonal, it holds

$$\|\Pi_1 - \Pi\|_2^2 = \|\Pi_1(I - \Pi) - (I - \Pi_1)\Pi\|_2^2 = \|\Pi_1(I - \Pi)\|_2^2 + \|(I - \Pi_1)\Pi\|_2^2.$$

Now, since  $\|\Pi\|_2^2 = \|\Pi_1\|_2^2 = m$ , we derive

$$\begin{aligned} \|\Pi_1(I - \Pi)\|_2^2 &= \|\Pi_1\|_2^2 - \|\Pi_1\Pi\|_2^2 = m - \|\Pi_1\Pi\|_2^2, \\ \|(I - \Pi_1)\Pi\|_2^2 &= \|\Pi\|_2^2 - \|\Pi_1\Pi\|_2^2 = m - \|\Pi_1\Pi\|_2^2, \end{aligned}$$

so that  $\|\Pi_1(I - \Pi)\|_2 = \|(I - \Pi_1)\Pi\|_2$  and the assertion follows.  $\blacksquare$

Let now  $B^\top B = O\Lambda O^\top$  be the single value decomposition (SVD) of the matrix  $B$  where  $O$  is the unitary  $L \times L$ -matrix and  $\Lambda$  is the diagonal matrix with non-increasing eigenvalues. Let then  $m \times d$ -matrix  $R$  be constructed due to (2.6) on the base of  $B$ , that is,  $R = (BO_m)^\top$  where  $O_m$  is the block of the first  $m$  columns of  $O$ . Clearly it holds  $|Rv| = |v^\top B|$  for every  $v \in \mathbb{R}^d$ . Similarly we define  $R_1$  via the SVD of  $B_1$ .

The projector  $\Pi$  in  $\mathbb{R}^d$  onto the value space of  $B$ , can be represented in the form  $\Pi = R^\top (RR^\top)^{-1} R$ . Similarly  $\Pi_1 = R_1^\top (R_1 R_1^\top)^{-1} R_1$ . Let  $\lambda_m$  denotes the smallest eigenvalue of  $RR^\top$ .

**Lemma 7.11** *Let the matrices  $B, B_1$  of rank  $m$  fulfill (7.14) with some  $\delta < \rho/4$ . Then the associated projectors  $\Pi$  and  $\Pi_1$  fulfill*

$$\|\Pi - \Pi_1\|_2 \leq \sqrt{2}\lambda_m^{1/2}\delta_1$$

where  $\delta_1 = \delta(1 - 2\delta/\rho - \delta^2/\rho^2)^{-1/2}$ .

**Proof.** The condition (7.14) implies by Lemma 7.8  $\|P_{\rho,1}(B - B_1)\|_2 \leq \delta_1$  which yields by Lemma 7.9

$$\|R_1(I - \Pi)\|_2 = \|(I - \Pi)B_1\|_2 \leq \delta_1.$$

This and Lemma 7.10 provide

$$\begin{aligned} \|\Pi - \Pi_1\|_2 &= \sqrt{2}\|\Pi(I - \Pi_1)\|_2 \\ &= \sqrt{2}\|R^\top (RR^\top)^{-1} R(I - \Pi_1)\|_2 \\ &\leq \sqrt{2}\|R^\top (RR^\top)^{-1}\| \|R_1(I - \Pi)\|_2 \\ &\leq \delta_1\sqrt{2}\|R^\top (RR^\top)^{-1}\|. \end{aligned}$$

It remains to note that

$$\left\| R^\top (RR^\top)^{-1} \right\|^2 = \left\| (RR^\top)^{-1} RR^\top (RR^\top)^{-1} \right\| = \left\| (RR^\top)^{-1} \right\| = \lambda_m^{-1}$$

and the assertion follows. ■

## References

- [1] Brillinger, D.R. (1983). A generalized linear model with “Gaussian” regressor variables. In *A Festschrift for E.H.Lehmann* (Bickel, Doksum and Hodges eds.) 97–114. Wadsworth.
- [2] Cook, D. (1998). Principal Hessian Directions Revisited. *J. Amer. Statist. Ass.* **93**, no. 441, 84–93.
- [3] Doksum, K. and Samarov, A. (1995). Nonparametric estimation of global functionals and a measure of explanatory power of covariates in regression. *Ann. Statist.* **23** 1443–1473.
- [4] Donoho, D.L. and Nussbaum, M. (1990). Minimax quadratic estimation of a quadratic functional. *J. of Complexity*, **6** 290–323.
- [5] Fan, J. (1991). On the estimation of quadratic functionals. *Ann. Statist.* **19** 1273–1294.
- [6] Horn, R.A., Johnson, C.R., *Matrix Analysis*, Cambridge University Press, N.Y., 1985.
- [7] Harville, A.A., *Matrix Analysis from a Statistician Perspective*, Springer, N.Y., 1997.
- [8] Hristache, M., Juditsky, A. and Spokoiny, V. (1998). Direct estimation of the index coefficients in a single-index model. Preprint 409, Weierstrass-Institute, Berlin. [www.wias-berlin.de/publications/preprints](http://www.wias-berlin.de/publications/preprints)
- [9] Huang, L.-S. and Fan, J. (1998). Nonparametric estimation of quadratic regression functionals. *Benoulli* **5** 927–949.
- [10] Ibragimov, I. and Khasminski, R. (1987). Estimation of linear functionals in Gaussian noise. *Theory Probab. Appl.*, **32** 30–39.
- [11] Li, K.-C. and Duan, N. (1989). Regression analysis under link violation. *Ann. Statist.*, **17**, 1009–1052.
- [12] Li, K.-C. (1991). Sliced inverse regression for dimension reduction. (With discussion). *J. Amer. Statist. Ass.* **86**, no. 414, 316–342.
- [13] Li, K.-C. (1992). On principal Hessian directions for data visualization and dimension reduction: Another application of Stein’s lemma. *J. Amer. Statist. Ass.* **87**, no. 420, 1025–1039.
- [14] Samarov, A. (1993). Exploring regression structure using nonparametric functional estimation. *J. Amer. Statist. Ass.* **88**, no. 423, 836–847.