

# Weierstraß-Institut für Angewandte Analysis und Stochastik

im Forschungsverbund Berlin e.V.

Preprint

ISSN 0946 – 8633

## Varying coefficient regression modeling

Jörg Polzehl and Vladimir Spokoiny

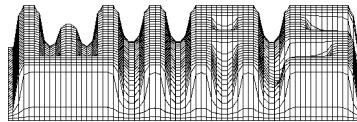
Weierstrass Institute for Applied Analysis and Stochastics

E-Mail: polzehl@wias-berlin.de, spokoiny@wias-berlin.de

submitted: 19th February 2003

No. 818

Berlin 2003



---

2000 *Mathematics Subject Classification.* 62G05.

*Key words and phrases.* adaptive weights; local structure; local polynomial regression.

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**

The adaptive weights smoothing (AWS) procedure was introduced in Polzehl and Spokoiny (2000) in the context of image denoising. The procedure has some remarkable properties like preservation of edges and contrast, and (in some sense) optimal reduction of noise. The procedure is also fully adaptive and dimension free. Simulations with artificial images show that AWS is superior to classical smoothing techniques especially when the underlying image function is discontinuous and can be well approximated by a piecewise constant function. However, the latter assumption can be rather restrictive for a number of potential applications. Here the AWS method is generalized to the case of an arbitrary local linear parametric structure. We also establish some important results about properties of the AWS procedure including the so called “propagation condition” and spatial adaptivity. The performance of the procedure is illustrated by examples for local polynomial regression in univariate and bivariate situations.

## 1 Introduction

Polzehl and Spokoiny (2000), referred to as PS2000 in what follows, offered a new method of nonparametric estimation, *Adaptive Weights Smoothing (AWS)*, in the context of image denoising. The main idea of the procedure is to describe in a data-driven and iterative way the largest local vicinity of every design point  $X_i$  in which the underlying model function can be well approximated by a constant. The procedure possesses some remarkable properties: it is fully adaptive in the sense that no prior information about the structure of the model is required, it is also design adaptive and has no boundary problem. A very important feature of the method is that it is dimension free and computationally straightforward. The results demonstrate that the new method is very efficient in situations when the underlying model allows a piecewise constant approximation with large homogeneous regions. This assumption seems to be reasonable e.g. in image analysis or for statistical inference in magnet resonance imaging, as shown in Polzehl and Spokoiny (2001), referred to as PS2001. Some other applications to density, volatility, tail index estimation can be found in Polzehl and Spokoiny (2002), referred to as PS2002. However, for many applications the assumption of a local constant structure can be too restrictive. A striking example is given by estimation of a smooth or piecewise smooth univariate regression function where a piecewise constant approximation is typically too rough. Local linear (polynomial) smoothing delivers much better results in this situation, see Fan and Gijbels (1996).

The aim of the present paper is to propose an extension of the AWS procedure to the case of varying coefficient regression models. Such models generalize classical nonparametric models and have got much attention within the last years, see e.g. Hastie and Tibshirani (1993), Fan and Zhang (1999), Carroll, Ruppert and Welsh (1998), Cai, Fan and Yao (2000) and references therein. The traditional approach is based on a local approximation of the varying coefficient model by a linear one in the varying parameter. The model is estimated for every localization point independently by the local least squares or local maximal likelihood. Accuracy of estimation is typically studied asymptotically as the localization parameter (bandwidth) tends to zero. Such an approach has serious drawbacks of being unable to incorporate special important cases like a global parametric model, a change-point model or more generally, models with inhomogeneous variability w.r.t. the varying parameter. In this paper we propose a completely different

approach based on the adaptive weights idea that allows to treat all the mentioned special cases in a unified way and to get a nearly optimal accuracy of estimation in every such situation.

The next section discusses the notions of global and local modeling. The basic idea of the generalized AWS and the description of the procedure are given in Section 3. The important special case of a local polynomial regression is discussed in Section 4. The performance of the method is studied for some simulated examples of univariate and bivariate regression in Section 5. Section 6 discusses theoretical properties of the procedure. We particularly prove an important property of the procedure called a “propagation condition” which means a free extension of every local model in a nearly homogeneous situation. Then we show that this condition automatically leads to a nearly optimal accuracy of estimation for a smooth regression function. Proofs and some technical results are provided in the Appendix.

## 2 Local modeling by weights

This section discusses our approach to local linear modeling. We start specifying the setup. Suppose that data  $Y_i$  are observed at design points  $X_i$  from the Euclidean space  $\mathbb{R}^d$ ,  $i = 1, \dots, n$ . In this paper we restrict ourselves to the regression setup with fixed design. The target of statistical analysis is the mean regression function  $f(x) = E(Y|X = x)$ . We will use a representation

$$Y_i = f(X_i) + \varepsilon_i \tag{2.1}$$

where  $\varepsilon_i$  can be naturally interpreted as additive random noise with zero mean. The distribution of the  $\varepsilon_i$ ’s is typically unknown but in many situations noise homogeneity can be assumed, that is, all the  $\varepsilon_i$ ’s are independent and satisfy  $E\varepsilon_i = 0$  and  $E\varepsilon_i^2 = \sigma^2$  for some  $\sigma > 0$ . For the case of exposition simplicity we restrict ourselves to this homoscedastic situation. Heteroscedastic noise can be considered as well, see PS2001 for some examples. We assume that an estimate  $\hat{\sigma}^2$  of  $\sigma^2$  from the data is available, see again PS2000 or PS2001 for specific examples.

A pure nonparametric estimate of the target function  $f(x)$  usually performs very poorly, especially in case of a multivariate design. The reason is that the underlying target function  $f(x)$  often is too complex to be estimated with a reasonable quality without further specifications of its structure.

The approach proposed in PS2000 and PS2001 can be called *structural adaptation*. One assumes that the underlying model has a relatively simple structure in some vicinity of every point  $X_i$  and the procedure attempts to recover this local structure using a pilot estimate of the model function and then to utilize this estimated local structural information for constructing a new improved estimate of the model function. These two steps are iterated several times extending at each iteration the degree of locality for every considered local model.

The original AWS method from PS2000 is based on the simplest local structural assumption: the function  $f$  is nearly constant within some neighborhood  $U(X_i)$  of the point  $X_i$ . In this paper the method is extended to the more general situation of a local linear structure.

## 2.1 Global linear modeling

Suppose we are given a parametric family  $\mathcal{F} = \{f_\theta, \theta \in \Theta\}$  where  $\Theta$  is a subset of a  $p$ -dimensional Euclidean space. A global parametric structure for the model (2.1) would mean that the underlying function  $f$  belongs to the family  $\mathcal{F}$ . The simplest example is a one-parameter family given by  $f_\theta(x) \equiv \theta$ . This family corresponds to a constant approximation of the underlying function and it was used in a local form in PS2000. In this paper we consider the more general situation of a linear parametric family. Let  $\psi_1(x), \dots, \psi_p(x)$  be some prescribed functions on  $\mathbb{R}^d$ . We define

$$\mathcal{F} = \{f_\theta(x) = \theta_1\psi_1(x) + \dots + \theta_p\psi_p(x), \quad \theta \in \mathbb{R}^p\}.$$

Under the global parametric assumption  $f \in \mathcal{F}$ , the corresponding parameter  $\theta$  can be easily estimated from the sample  $Y_1, \dots, Y_n$ . A natural estimate of  $\theta$  is given by the ordinary least squares method:

$$\hat{\theta} = \operatorname{arginf}_{\theta} \sum_{i=1}^n (Y_i - f_\theta(X_i))^2.$$

For an explicit representation of this estimate vector notation is useful. Define vectors  $\Psi_i$  in  $\mathbb{R}^p$  with entries  $\psi_m(X_i)$ ,  $m = 1, \dots, p$ , and the  $p \times n$ -matrix  $\Psi$  whose columns are  $\Psi_i$ . Let also  $Y$  stand for the vector of observations:  $Y = (Y_1, \dots, Y_n)^\top \in \mathbb{R}^n$ . Then

$$\hat{\theta} = \left( \sum_{i=1}^n \Psi_i \Psi_i^\top \right)^{-1} \sum_{i=1}^n \Psi_i Y_i = \left( \Psi \Psi^\top \right)^{-1} \Psi Y$$

provided that the  $p \times p$  matrix  $\Psi \Psi^\top$  is nondegenerated.

## 2.2 Local linear modeling

The global parametric assumption can be too restrictive in many situations and it does not allow to model complex statistical objects. A standard approach in nonparametric inference is to apply the parametric (linear) structural assumption locally.

Different possibilities to describe a local model centered at a given point are discussed in PS2002. The most general one is *localization by weights*. Let, for a fixed  $x$ , a nonnegative weight  $w_i \leq 1$  be assigned to the observation  $Y_i$  at  $X_i$ . The weights  $w_i = w_i(x)$  determine a local model corresponding to the point  $x$  in the sense that every observation  $Y_i$  is used with the weight  $w_i$  when estimating the local parameter  $\theta$  at  $x$ . This leads to the local (weighted) least squares estimate

$$\hat{\theta}(x) = \operatorname{arginf}_{\theta \in \Theta} \sum_{i=1}^n w_i (Y_i - f_{\theta}(X_i))^2 = \left( \Psi W \Psi^{\top} \right)^{-1} \Psi W Y \quad (2.2)$$

with  $W = \operatorname{diag}\{w_1, \dots, w_n\}$ .

We mention two examples of choosing the weights  $w_i$ . *Localization by a bandwidth* is defined by the weights of the form  $w_i(x) = K_l(\mathbf{l}_i)$  with  $\mathbf{l}_i = |\rho(x, X_i)/h|^2$  where  $h$  is a bandwidth,  $\rho(x, X_i)$  is the Euclidean distance between  $x$  and the design point  $X_i$  and  $K_l$  is a *location kernel*.

*Localization by a window* simply restricts the model to some subset (window)  $U$  of the design space, that is,  $w_i = \mathbf{1}(X_i \in U)$  and all data points  $Y_i$  with  $X_i$  outside the region  $U$  are not taken into account when estimating the value  $\theta(x)$ .

We do not assume any special structure for the weights  $w_i$ , that is, any configuration of the weights is allowed. In what follows we identify the diagonal weight matrix  $W = \operatorname{diag}\{w_1, \dots, w_n\}$  and the local model defined by these weights.

## 3 Adaptive weights smoothing

This section describes the new method of locally adaptive estimation, the *Adaptive Weights Smoothing*, starting from a preliminary discussion. The idea of the procedure is to describe adaptively for every point  $X_i$  the largest possible neighborhood in which the model function  $f(x)$  can be well approximated by a parametric function  $f_{\theta}$  from  $\mathcal{F}$ . The local model at  $X_i$  is described by weights  $w_{ij}$  assigned to every observation  $Y_j$ .

The procedure is iterative. At every iteration step, the procedure tries to extend the local model at each design point. We first illustrate the idea for the local constant

structural assumption as considered in PS2000. Here the estimate  $\hat{\theta}_i = \hat{f}(X_i)$  is defined as the mean of the observations  $Y_j$  with some weights  $w_{ij}$ :

$$\hat{f}(X_i) = \sum_{j=1}^n w_{ij} Y_j / \sum_{j=1}^n w_{ij}. \quad (3.1)$$

These weights  $w_{ij}$  are calculated iteratively. For the initial step, the estimate  $\hat{f}^{(0)}(X_i)$  is computed from a smallest local model defined by a bandwidth  $h^{(0)}$ , that is,

$$\hat{f}^{(0)}(X_i) = \hat{\theta}_i^{(0)} = \sum_{j=1}^n K_l(\mathbf{l}_{ij}^{(0)}) Y_j / \sum_{j=1}^n K_l(\mathbf{l}_{ij}^{(0)})$$

with  $\mathbf{l}_{ij}^{(0)} = |\rho(X_i, X_j)/h^{(0)}|^2$ . In other words, the algorithm starts with the usual kernel estimate with the bandwidth  $h^0$ , which is taken very small. If  $K_l = \mathbf{1}(u \leq 1)$  as in PS2000, then for every point  $X_i$  the weights  $w_{ij}$  vanish outside the ball  $U_i^{(0)}$  of radius  $h^{(0)}$  with the center at  $X_i$ , that is, the local model at  $X_i$  is concentrated on  $U_i^{(0)}$ . Next, at each iteration  $k$ , a ball  $U_i^{(k)}$  with a larger bandwidth  $h^{(k)}$  is considered and every point  $X_j$  from  $U_i^{(k)}$  gets a weight  $w_{ij}^{(k)}$  which is defined by comparing the estimates  $\hat{f}^{(k-1)}(X_i)$  and  $\hat{f}^{(k-1)}(X_j)$  obtained in the previous iteration. The weights are then used to compute new improved estimates  $\hat{f}^{(k)}(X_i)$  due to (3.1).

One possible interpretation of this procedure is that at each iteration step the location penalty  $\mathbf{l}_{ij}^{(k)}$  is relaxed by increasing  $h^{(k)}$  at cost of introducing a data-driven statistical penalty which comes from comparison of different local models.

Note that under the local constant assumption  $f(x) = \theta$ , the value  $\theta$  uniquely determines the model function and the comparison of the values  $\hat{f}^{(k-1)}(X_i)$  and  $\hat{f}^{(k-1)}(X_j)$  is equivalent to a comparison of two model functions. The extension of this approach to the more general local parametric assumption leads to a check of homogeneity for two local models  $W_i^{(k-1)} = \text{diag}\{w_{i1}^{(k-1)}, \dots, w_{in}^{(k-1)}\}$  and  $W_j^{(k-1)} = \text{diag}\{w_{j1}^{(k-1)}, \dots, w_{jn}^{(k-1)}\}$ , to specify the weight  $w_{ij}^{(k)}$ . Now we discuss how a statistical penalty (distance) for two local models can be computed.

### 3.1 Measuring the statistical difference between two local models

Consider two local models corresponding to points  $X_i$  and  $X_j$  and defined by diagonal weight matrices  $W_i$  and  $W_j$ . We suppose that the structural assumption is fulfilled for both, that is, the underlying regression function  $f$  can be well approximated by some  $f_\theta \in \mathcal{F}$  within every local model. However, the value of the parameter  $\theta$  determining the



approximating function  $f_\theta$  may be different for the two local models. We aim to develop a rule to judge from the data, whether the local model corresponding to the point  $X_j$  and described by  $W_j$  is not significantly different (in the value of the underlying parameter  $\theta$ ) from the model at  $X_i$  described by  $W_i$ . More precisely, we want to quantify the difference between these two local models in order to assign a weight  $w_{ij}$  with which the observation  $Y_j$  will enter into the local model at  $X_i$  in the next iteration of the algorithm.

A natural way is to consider the data from two local models as two different populations and to apply the two population likelihood ratio test for testing the hypothesis  $\theta_i = \theta_j$ . Suppose that the errors  $\varepsilon_i$  are normally distributed with parameters  $(0, \sigma^2)$ . The log-likelihood  $L(W_i, \theta, \theta')$  for the local regression model at  $X_i$  with the weights  $W_i$  is, for any pair  $\theta, \theta' \in \Theta$ , defined by

$$\begin{aligned} L(W_i, \theta, \theta') &= \frac{1}{2\sigma^2} \sum_{l=1}^n w_{il} \left[ (Y_l - \Psi_l^\top \theta')^2 - (Y_l - \Psi_l^\top \theta)^2 \right] \\ &= \frac{1}{2\sigma^2} \sum_{l=1}^n w_{il} \left[ 2(Y_l - \Psi_l^\top \theta') \Psi_l^\top (\theta - \theta') - (\theta - \theta')^\top \Psi_l \Psi_l^\top (\theta - \theta') \right] \end{aligned}$$

yielding

$$L(W_i, \hat{\theta}_i, \theta') = \frac{1}{2\sigma^2} (\hat{\theta} - \theta')^\top B_i (\hat{\theta} - \theta'),$$

with  $B_i = \Psi W_i \Psi^\top$ .

The classical likelihood-ratio test statistic is of the form

$$\begin{aligned} T_{ij}^\circ &= \max_{\theta} L(W_i, \theta, \theta') + \max_{\theta} L(W_j, \theta, \theta') - \max_{\theta} L(W_i + W_j, \theta, \theta') \\ &= L(W_i, \hat{\theta}_i, \theta') + L(W_j, \hat{\theta}_j, \theta') - L(W_i + W_j, \hat{\theta}_{ij}, \theta') \end{aligned} \quad (3.2)$$

where  $\hat{\theta}_i = \operatorname{argmax}_{\theta} L(W_i, \theta, \theta')$  is the maximum likelihood estimate (MLE) corresponding to the local model described by the weight matrix  $W_i$  and similarly for  $\hat{\theta}_j$ . Also  $\hat{\theta}_{ij} = \operatorname{argmax}_{\theta} L(W_i + W_j, \theta, \theta')$  is the local MLE corresponding to the combined model which is obtained by summing the weights from the both models. The value  $T_{ij}^\circ$  characterizes the difference between the two models in the statistical sense: if  $T_{ij}^\circ$  is larger than some prescribed value  $\lambda$ , then these two models are significantly different in the value of the underlying parameter  $\theta$ .

The criterion based on  $T_{ij}^\circ$  has a serious drawback of giving more weight to the “smaller” model. For instance, in the “unbalanced” situation when the model  $W_i$  is

much “larger” than  $W_i$  (that is,  $B_i \gg B_j$ ), the distribution of  $T_{ij}^\circ$  is mostly determined by the distribution of  $\hat{\theta}_j$ . This feature is not desirable when we define new weights  $w_{ij}$  for the model centered at  $X_i$ . To avoid this problem, we standardize the weights  $W_j = \{w_{jl}\}$  by multiplying with some factor  $\alpha$  and then optimize the test statistic w.r.t. this factor  $\alpha$ . The use of the factor  $\alpha$  leads to the test statistics

$$T_{ij}(\alpha) = L(W_i, \hat{\theta}_i, \theta') + L(\alpha W_j, \hat{\theta}_j, \theta') - L(W_i + \alpha W_j, \hat{\theta}_{ij}, \theta')$$

where

$$\hat{\theta}_{ij} = \operatorname{argmax}_{\theta} L(W_i + \alpha W_j, \theta, \theta') = \left( \Psi(W_i + \alpha W_j) \Psi^\top \right)^{-1} \Psi(W_i + \alpha W_j) Y.$$

The use of  $\theta' = \hat{\theta}_j$  yields

$$\begin{aligned} T_{ij}(\alpha) &= L(W_i, \hat{\theta}_i, \hat{\theta}_j) - L(W_i + \alpha W_j, \hat{\theta}_{ij}, \hat{\theta}_j) \\ &= (\hat{\theta}_i - \hat{\theta}_j)^\top \Psi W_i \Psi^\top (\hat{\theta}_i - \hat{\theta}_j) - (\hat{\theta}_{ij} - \hat{\theta}_j)^\top \Psi (W_i + \alpha W_j) \Psi^\top (\hat{\theta}_{ij} - \hat{\theta}_j) \\ &\leq (\hat{\theta}_i - \hat{\theta}_j)^\top \Psi W_i \Psi^\top (\hat{\theta}_i - \hat{\theta}_j). \end{aligned}$$

Moreover,

$$\lim_{\alpha \rightarrow +\infty} T_{ij}(\alpha) = T_{ij} = (\hat{\theta}_i - \hat{\theta}_j)^\top \Psi W_i \Psi^\top (\hat{\theta}_i - \hat{\theta}_j) = (\hat{\theta}_i - \hat{\theta}_j)^\top B_i (\hat{\theta}_i - \hat{\theta}_j). \quad (3.3)$$

Indeed, simple algebra provides

$$\begin{aligned} T_{ij}(\alpha) &= (2\sigma^2)^{-1} \left[ (\hat{\theta}_i - \hat{\theta}_j)^\top B_i (\hat{\theta}_i - \hat{\theta}_j) + (\hat{\theta}_{ij} - \hat{\theta}_j)^\top \alpha B_j (\hat{\theta}_{ij} - \hat{\theta}_j) \right] \\ &= (2\sigma^2)^{-1} (\hat{\theta}_i - \hat{\theta}_j)^\top B_i (B_i + \alpha B_j)^{-1} \alpha B_j (\hat{\theta}_i - \hat{\theta}_j) \rightarrow T_{ij}, \quad \alpha \rightarrow \infty. \end{aligned}$$

We consider the value  $T_{ij}$  as a ‘statistical penalty’, that is, when computing the new weight  $w_{ij}$  at the next iteration step we strongly penalize for a large value of  $T_{ij}$ .

For the case of a one-dimensional parameter  $\theta$ , that is, with  $p = 1$ , the expression for the statistical penalty can be simplified. Indeed,  $\Psi$  is a vector in  $\mathbb{R}^n$  and  $B_i = \Psi^\top W_i \Psi$  is a number yielding  $T_{ij} = B_i |\hat{\theta}_i - \hat{\theta}_j|^2 / (2\sigma^2)$ .

### 3.2 Penalization for extending a local model

An important feature of the original AWS method from PS2000 is its stability against iteration. It turns out that the generalization of the local constant procedure to the local linear case requires to introduce an additional penalty to prevent from leverage

problems. To clarify the idea, suppose for the moment that for every iteration step  $k$ , each local model  $W_i^{(k)}$  is restricted to the ball  $U_i^{(k-1)}$  of the radius  $h^{(k-1)}$  centered at  $X_i$ . Suppose also that the first  $k-1$  iterations of the algorithm have been carried over. As a result, we obtain for every point  $X_i$  a local model described by the weights  $w_{ij}^{(k-1)}$  for each  $X_j \in U_i^{(k-1)}$ . At the next iteration the procedure tries to extend every local model by increasing the bandwidth  $h^{(k)}$  and assigning the weights  $w_{ij} = w_{ij}^{(k)}$  for every point  $X_j$  from the larger neighborhood  $U_i^{(k)}$  of  $X_i$  with the radius  $h^{(k)}$ . If  $X_j \in U_i^{(k)} \setminus U_i^{(k-1)}$ , then giving  $X_j$  a significantly positive weight  $w_{ij}$  can be interpreted as including the point  $X_j$  into the local model centered at  $X_i$ . In some cases, including even one point  $X_j$  with a relatively large value  $\rho(X_i, X_j)$  into the local model at  $X_i$  may significantly change the estimate  $\hat{\theta}_i$ . Such leverage problem does not arise in the local constant modeling but it becomes crucial for local linear (polynomial) regression. To prevent from this danger, we introduce a special penalty for including an influence point.

To measure the influence of the observation  $Y_j$  at  $X_j$  in the local model described by the weight matrix  $W_i$ , one can consider the extended model obtained by adding a single observation at the point  $X_j$  and look at the relative difference between the original and the extended model. This leads to the value

$$\begin{aligned} \gamma_{ij} &= \text{tr} \left\{ \left( \Psi \overline{W}_i \Psi^\top \right)^{-1} \left( \Psi \overline{W}_i \Psi^\top + \Psi_j \Psi_j^\top \right) \right\} - p \\ &= \Psi_j^\top \left( \Psi \overline{W}_i \Psi^\top \right)^{-1} \Psi_j = (\text{tr} W_i) \Psi_j^\top \left( \Psi W_i \Psi^\top \right)^{-1} \Psi_j. \end{aligned}$$

Here  $\Psi_j \in \mathbb{R}^p$  is the  $j$ th column of  $\Psi$  and, for a diagonal matrix  $W$ , we denote  $\overline{W} = (\text{tr} W)^{-1} W$ . A large value of  $\gamma_{ij}$  means that  $X_j$  is a leverage point. To make the procedure more stable w.r.t. such influential points, we additionally penalize for including points with a large value  $\gamma_{ij}$ , i.e. assign small weights even when the difference  $\hat{\theta}_i - \hat{\theta}_j$  is statistically insignificant and the statistical penalty  $s_{ij}$  is small.

For adjusting the penalty term one can use the ‘propagation’ principle which means a free extension of the model in the homogeneous situation when the coefficients of the linear model do not vary with location. In that situation, neither the statistical penalty nor the penalty for extending the model would significantly affect the estimate leading after the first  $k-1$  iterations to the classical location weights  $w_{ij, \text{ho}}^{(k-1)} = K_l \left( l_{ij}^{(k-1)} \right) = K_l \left( |\rho(X_i, X_j)/h^{(k-1)}|^2 \right)$ . The influence of the point  $X_j$  within the local homogeneous

model described by  $W_{i,\text{ho}}^{(k-1)}$  is given by

$$\gamma_{ij,\text{ho}} = \gamma_j \left( W_{i,\text{ho}}^{(k-1)} \right) = \left( \text{tr} W_{i,\text{ho}}^{(k-1)} \right) \Psi_j^\top \left( \Psi W_{i,\text{ho}}^{(k-1)} \Psi^\top \right)^{-1} \Psi_j$$

where  $W_{i,\text{ho}}^{(k-1)} = \text{diag}\{w_{i1,\text{ho}}^{(k-1)}, \dots, w_{in,\text{ho}}^{(k-1)}\}$ . This value  $\gamma_{ij,\text{ho}}$  can be used for adjusting the penalty for extending the model. Namely, we assign to every observation  $Y_j$  at  $X_j$  the penalty

$$\mathbf{e}_{ij}^{(k)} = \tau^{-1} \left( \frac{\gamma_{ij}}{\gamma_{ij,\text{ho}}} - 1 \right)_+$$

where  $a_+$  means  $\max\{0, a\}$  and  $\tau$  is a numerical tuning parameter.

### 3.3 Defining weights

Using the previously described methods, we compute for every pair  $(i, j)$  the penalties  $\mathbf{l}_{ij}^{(k)}$ ,  $\mathbf{s}_{ij}^{(k)}$  and  $\mathbf{e}_{ij}^{(k)}$ . It is natural to require that the influence of every such factor is independent of the other factors. This suggests to define the new weight  $w_{ij}^{(k)}$  using the product

$$\tilde{w}_{ij}^{(k)} = K_l(\mathbf{l}_{ij}^{(k)}) K_s(\mathbf{s}_{ij}^{(k)}) K_e(\mathbf{e}_{ij}^{(k)}),$$

where  $K_l, K_s$  and  $K_e$  are three kernel functions on the positive semiaxis satisfying the condition  $K_l(0) = K_s(0) = K_e(0) = 1$ .

In the algorithm presented below in this section, we use one more (memory) parameter  $\eta \in (0, 1)$  which controls the rate of changing the weights for every local model within the iteration process. Namely, we define the new weight  $w_{ij}^{(k)}$  as a convex combination of the previous step weight  $w_{ij}^{(k-1)}$  and the just defined product  $\tilde{w}_{ij}^{(k)}$ :

$$w_{ij}^{(k)} = \eta w_{ij}^{(k-1)} + (1 - \eta) \tilde{w}_{ij}^{(k)}.$$

### 3.4 The procedure

Now we present a formal description. Important ingredients of the method are:

- kernels  $K_l, K_s$  and  $K_e$ ;
- parameters  $\lambda, \tau$  and  $\eta$ ;
- the initial bandwidth  $h^{(0)}$ , the factor  $a > 1$  and the maximal bandwidth  $h_{\max}$
- the estimated error variance  $\hat{\sigma}^2$ .

The choice of the parameters is discussed in Section 3.5. The generalized procedure reads as follows:

**1. Initialization:** For every  $i$  define the diagonal matrix  $W_i^{(0)}$  with the diagonal entries  $w_{ij}^{(0)} = K_l(\mathbf{l}_{ij}^{(0)})$  and  $\mathbf{l}_{ij}^{(0)} = |\rho(X_i, X_j)/h^{(0)}|^2$ , that is,  $W_i^{(0)} = \text{diag}\{w_{i1}^{(0)}, \dots, w_{in}^{(0)}\}$ . Compute

$$N_i^{(0)} = \text{tr} W_i^{(0)}, \quad B_i^{(0)} = \Psi W_i^{(0)} \Psi^\top, \quad Z_i^{(0)} = \Psi W_i^{(0)} Y \quad \text{and} \quad \widehat{\theta}_i^{(0)} = \left(B_i^{(0)}\right)^{-1} Z_i^{(0)}.$$

Set  $k = 1$ .

**2. Iteration:** for every  $i = 1, \dots, n$  define  $W_{i,\text{ho}}^{(k-1)} = \text{diag}\{K_l(\mathbf{l}_{i1}^{(k-1)}), \dots, K_l(\mathbf{l}_{in}^{(k-1)})\}$ ,

• **calculate the adaptive weights:** For every point  $X_j$  compute

$$\begin{aligned} \gamma_{ij}^{(k)} &= N_i^{(k-1)} \Psi_j^\top \left(B_i^{(k-1)}\right)^{-1} \Psi_j, \\ \gamma_{ij,\text{ho}}^{(k)} &= \text{tr}\left(W_{i,\text{ho}}^{(k-1)}\right) \Psi_j^\top \left(\Psi W_{i,\text{ho}}^{(k-1)} \Psi^\top\right)^{-1} \Psi_j \end{aligned}$$

where  $\Psi_j$  is  $j$ th column of  $\Psi$ .

Compute the penalties

$$\begin{aligned} \mathbf{l}_{ij}^{(k)} &= |\rho(X_i, X_j)/h^{(k)}|^2, \\ \mathbf{s}_{ij}^{(k)} &= \frac{1}{2\widehat{\sigma}^2\lambda} (\widehat{\theta}_i^{(k-1)} - \widehat{\theta}_j^{(k-1)})^\top B_i^{(k-1)} (\widehat{\theta}_i^{(k-1)} - \widehat{\theta}_j^{(k-1)}), \\ \mathbf{e}_{ij}^{(k)} &= \tau^{-1} (\gamma_{ij}^{(k)} / \gamma_{ij,\text{ho}}^{(k)} - 1)_+. \end{aligned} \tag{3.4}$$

Compute the value  $\widetilde{w}_{ij}^{(k)}$ :

$$\widetilde{w}_{ij}^{(k)} = K_l(\mathbf{l}_{ij}^{(k)}) K_s(\mathbf{s}_{ij}^{(k)}) K_e(\mathbf{e}_{ij}^{(k)}), \tag{3.5}$$

Denote by  $\widetilde{W}_i^{(k)}$  the diagonal matrix whose diagonal elements are  $\widetilde{w}_{ij}^{(k)}$ , that is,  $\widetilde{W}_i^{(k)} = \text{diag}\{w_{i1}^{(k)}, \dots, w_{in}^{(k)}\}$ .

• **Compute new estimate:** Compute

$$\begin{aligned} N_i^{(k)} &= \eta N_i^{(k-1)} + (1 - \eta) \text{tr} \widetilde{W}_i^{(k)}, \\ Z_i^{(k)} &= \eta Z_i^{(k-1)} + (1 - \eta) \Psi \widetilde{W}_i^{(k)} Y, \\ B_i^{(k)} &= \eta B_i^{(k-1)} + (1 - \eta) \Psi \widetilde{W}_i^{(k)} \Psi^\top, \end{aligned}$$

and the estimate  $\widehat{\theta}_i^{(k)}$  (resp.  $\widehat{f}_i^{(k)}$ ) of  $\theta_i$  (resp. of  $f_i = f(X_i)$ ) by

$$\widehat{\theta}_i^{(k)} = \left(B_i^{(k)}\right)^{-1} Z_i^{(k)}, \quad \widehat{f}_i^{(k)} = \Psi_i^\top \widehat{\theta}_i^{(k)}.$$

**3. Stopping:** Increase  $k$  by 1, set  $h^{(k)} = ah^{(k-1)}$ . If  $h^{(k)} \leq h_{\max}$  continue with step 2. Otherwise terminate.

By  $k^*$  we denote the total number of iterations. Also define the final estimates  $\widehat{f}_i = \widehat{f}_i^{(k^*)}$ .

### 3.5 Choice of parameters

The parameters of the procedure are selected similarly to PS2000. We briefly discuss each of the parameters.

**Kernels  $K_s$ ,  $K_l$  and  $K_e$ :** The kernels  $K_s$  and  $K_l$  must fulfill  $K_s(0) = K_l(0) = K_e(0) = 1$ , with  $K_s$ ,  $K_e$  decreasing and  $K_l$  non-increasing on the positive semiaxis. We recommend to take  $K_s(u) = e^{-u}I_{\{u \leq 6\}}$ . We also recommend to apply a compactly supported localization kernel  $K_l$  to reduce the computational effort of the method. PS2000 applied a uniform kernel, here we apply the triangle kernel  $K_l(u) = (1 - u)_+$ . We also set  $K_e = K_s$ .

**Initial bandwidth  $h^{(0)}$ , parameter  $a$  and maximal bandwidth  $h_{\max}$ :** The starting bandwidth  $h^{(0)}$  should be taken possibly small. In the most of example we select  $h^{(0)}$  such that every starting local neighborhood  $U_i^{(0)}$  contains sufficiently many design points to get an initial estimate of the parameter  $\theta_i$ .

The parameter  $a$  controls the growth rate of the local neighborhoods for every point  $X_i$ . It should be selected to provide that the mean number of points inside a ball  $U_i^{(k)}$  with radius  $h^{(k)}$  grows exponentially in  $k$  with the factor  $a_{\text{grow}}$ . If  $X_i$  are from the unit cube in the space  $\mathbb{R}^d$ , then the parameter  $a$  can be taken as  $a = a_{\text{grow}}^{1/d}$ . Our default choice is  $a_{\text{grow}} = 1.25$ . The exponential grow of the design points within every ball  $U_i^{(k)}$  ensures that the number of iterations  $k^*$  is at most logarithmic in the sample size.

The maximal bandwidth  $h_{\max}$  may be taken very large. However, this parameter can be used to bound the numerical complexity of the procedure, see Section 3.6. In some applications, the use of a very large final bandwidth  $h_{\max}$  leads to some oversmoothing of the underlying object. For such situations, a data-driven method of optimal stopping, based, for instance, on cross-validation can be applied.

**Parameter  $\lambda$ :** The most important parameter of the procedure is  $\lambda$  which scales the statistical penalty  $s_{ij}$ . Small values of  $\lambda$  lead to overpenalization which may result in unstable performance of the method in the homogeneous situation. Large values of  $\lambda$  may result in loss of adaptivity of the method (less sensitivity to structural changes). A reasonable way to define the parameter  $\lambda$  for a specific application is based on the condition of free extension, which we also call the “propagation condition”. This condition means that in a homogeneous situation, i.e. when the underlying parameters for

every two local models coincide, the impact of the statistical penalty in the computed weights  $w_{ij}$  is negligible. This would result in a free extension of every local model. If the value  $h_{\max}$  is sufficiently large, at the end of iteration process all the weights  $w_{ij}$  will then be close to one and every local model will essentially coincide with the global one. Therefore, one can adjust the parameter  $\lambda$  simply selecting the minimal value of  $\lambda$  still providing a prescribed probability of getting the global model at the end of iteration process for the homogeneous (parametric) model  $\theta(x) = \theta$  using Monte-Carlo simulations. The theoretical justification is given by Theorem 6.1 in Section 6.1, that claims that the choice  $\lambda = C \log n$  with a sufficiently large  $C$  yields the “propagation” condition whatever the parameter  $\theta$  is.

Our default value is  $\lambda = q_\alpha(\chi_p^2)$ , that is the  $\alpha$ -quantile of the  $\chi^2$  distribution with  $p$  degree of freedom, where  $\alpha$  depends on the specified linear parametric family. Defaults for the case of local polynomial regression are given in Section 5.

**Parameter  $\tau$ :** The optimal choice of  $\tau$  depends on the method of smoothing. For the local constant AWS considered in PS2000, there are no influential points (see Section 4.1). For local polynomial smoothing the choice of  $\tau$  is discussed in more details in Section 4.

**Parameter  $\eta$  and the control step:** A value  $\eta \in (0, 1)$  can be used to control the stability of the AWS procedure w.r.t. iterations. An increase of  $\eta$  results in a higher stability, however, it decreases the sensitivity to changes of the local structure. The use of the memory parameter also guarantees that the estimates  $\hat{\theta}_i^{(k)}$  are well defined, that is, all the matrices  $B_i^{(k)}$  are positive definite. Our default choice is  $\eta = 1/2$ .

The original AWS procedure from PS2000 did not involve the “memory” parameter  $\eta$  (it corresponds to  $\eta = 0$ ). Instead it contained one additional *control* step in which the new estimate  $\hat{\theta}_i^{(k)}$  is compared with all the previous estimates  $\hat{\theta}_i^{(k')}$  for  $k' < k$ . If the difference  $\hat{\theta}_i^{(k)} - \hat{\theta}_i^{(k')}$  became significant, the new estimate was not accepted and the previous step estimate was used. This control step is a very useful device for proving some theoretical properties of the procedure, because it ensures that the gained quality of estimation will not be lost in further iterations, see Section 6 for more details. In the local linear case this control step would accept the estimate  $\hat{\theta}_i^{(k)}$  only if

$$(2\hat{\sigma}^2)^{-1} (\hat{\theta}_i^{(k')} - \hat{\theta}_i^{(k)})^\top B_i^{(k')} (\hat{\theta}_i^{(k')} - \hat{\theta}_i^{(k)}) \leq \eta^*, \quad k' = 1, \dots, k-1, \quad (3.6)$$

that is, when the new estimate  $\hat{\theta}_i^{(k)}$  lies inside all confidence ellipsoids of previous estimates at the point  $X_i$ . However, our numerical results (not reported here) indicate that

the usefulness of the control step for practical purpose is questionable. The use of the “memory” parameter  $\eta$  can be regarded as a soft version of the control step.

### 3.6 Computational complexity of the algorithm

We start with the following two important remarks. First note, that every estimate is defined as  $\hat{\theta}_i^{(k)} = \left(B_i^{(k)}\right)^{-1} Z_i^{(k)}$  using the matrix  $B_i^{(k)}$  and the vector  $Z_i^{(k)}$ . Similarly, the new weights  $\tilde{w}_{ij}^{(k)}$  are computed on the basis of the same statistics  $B_i^{(k-1)}$ ,  $Z_i^{(k-1)}$  and  $N_i^{(k-1)}$  from the previous step of the procedure. Therefore, the whole structural information is contained in these three basis elements. During the adaptation step, we compute for every  $i$  the weights  $\tilde{w}_{ij}^{(k)}$  with different  $j$  only with the aim to compute the new elements  $B_i^{(k)}$ ,  $Z_i^{(k)}$  and  $N_i^{(k)}$ . This reduces the memory requirements for the algorithm to  $\mathcal{O}(np^2)$  or even to  $\mathcal{O}(np)$  for local polynomial modeling, see the next section, while keeping all the weights  $w_{ij}^{(k)}$  would lead to the memory requirement  $\mathcal{O}(n^2)$ .

Secondly we notice, that the localization kernel  $K_l$  usually has a compact support, say,  $[0, 1]$ . This immediately implies that for every local model at  $X_i$ , all the weights  $\tilde{w}_{ij}^{(k)}$  for the points  $X_j$  outside the ball  $U_i^{(k)} = \{x : \rho(X_i, x) \leq h^{(k)}\}$  vanish. Therefore, it suffices at each step to compute the weights  $\tilde{w}_{ij}^{(k)}$  for pairs  $X_i, X_j$  with  $\rho(X_i, X_j) \leq h^{(k)}$ . Denote by  $M_k$  the maximal number of design points  $X_j$  within a ball of radius  $h^{(k)}$  centered at a design point. At the  $k$  step there are at most  $M_k$  positive weights  $\tilde{w}_{ij}^{(k)}$  for any  $X_i$ .

Therefore, for carrying out the  $k$ th adaptation step of the algorithm, we have to compute the penalties  $l_{ij}^{(k)}$ ,  $s_{ij}^{(k)}$  and  $e_{ij}^{(k)}$  and the value  $\tilde{w}_{ij}^{(k)}$ , for every pair  $(i, j)$  with  $\rho(X_i, X_j) \leq h^{(k)}$  due to (3.5). This requires a finite number of operations depending on the number of parameters  $p$  only, and the whole  $k$ th adaptation step of the algorithm requires of order  $nM_k$  operations.

To obtain the estimate we need, for every point  $X_i$ , to compute the  $d \times d$ -matrix  $B_i^{(k)} = \eta B_i^{(k-1)} + (1 - \eta) \Psi \tilde{W}_i^{(k)} \Psi^\top$ , the vector  $Z_i^{(k)} = \eta Z_i^{(k-1)} + (1 - \eta) \Psi \tilde{W}_i^{(k)} Y$  and the value  $N_i^{(k)} = \eta N_i^{(k-1)} + (1 - \eta) \text{tr} \tilde{W}_i^{(k)}$ . It is clear that the complexity of computing all these values is of order  $M_k$ . Computing  $\hat{\theta}_i^{(k)} = \left(B_i^{(k)}\right)^{-1} Z_i^{(k)}$  requires a finite number operations depending on  $p$  only. Therefore, the complexity of the whole estimation step is again of order  $nM_k$ .

Since typically the numbers  $M_k$  grow exponentially, the complexity of the whole



algorithm is estimated as

$$n(M_1 + \dots + M_{k^*}) \asymp nM_{k^*}$$

where  $k^*$  is the number of iteration steps.

## 4 Local polynomial regression

In this section we specify the procedure for nonparametric estimation of a regression function with univariate and multivariate covariates. The underlying regression function is assumed to be smooth or piecewise smooth leading to a polynomial approximation of the function within each local model.

### 4.1 Local constant regression

First we briefly consider a special cases of the above procedure corresponding to the local constant AWS procedure from PS2000.

The local constant approximation corresponds to the simplest family of basis functions  $\{\psi_m\}$  consisting of one constant function  $\psi_0 \equiv 1$ . The major advantage of this method is that the dimensionality of the regressors plays absolutely no role. In this situation  $\Psi = (1, \dots, 1)$  and, for every diagonal matrix  $W = \text{diag}(w_1, \dots, w_n)$ , it holds  $\Psi W \Psi^\top = \text{tr} W$  and  $\Psi W Y = \sum_{l=1}^n w_l Y_l$ . Hence, for the local constant case, the  $B_i^{(k)}$ 's coincide with the  $N_i^{(k)}$ 's. The statistical penalty  $\mathbf{s}_{ij}^{(k)}$  can be written in the form  $\mathbf{s}_{ij}^{(k)} = (2\sigma^2)^{-1} N_i^{(k-1)} |\hat{\theta}_i^{(k-1)} - \hat{\theta}_j^{(k-1)}|^2$ . Also, for all  $i$  and  $k$ , it holds  $\gamma_{ij}^{(k)} = \text{tr} W_i^{(k-1)} / \text{tr} W_i^{(k-1)} \equiv 1$ , and similarly for  $\gamma_{ij, \text{ho}}^{(k)}$ . Therefore, the penalty  $e_{ij}$  is always zero and can be dropped.

The weights  $\tilde{w}_{ij}^{(k)}$  can be computed as  $\tilde{w}_{ij}^{(k)} = K_l(\mathbf{l}_{ij}^{(k)}) K_s(\mathbf{s}_{ij}^{(k)})$  that essentially coincides with the proposal from PS2000 if the uniform kernel  $K_l$  is applied. A small difference remains in the use of the memory parameter  $\eta$  and in a slightly different form of the statistical penalty.

### 4.2 Local polynomial univariate regression

Local linear (polynomial) smoothing is known to be much more accurate when estimating a smooth function, see e.g. Fan and Gijbels (1996). A generalization of the original AWS to the local linear (polynomial) regression therefore is of special importance. We describe the corresponding procedure in more details.

One can specify the basis functions as  $\psi_1(x) = 1$ ,  $\psi_2(x) = x$ ,  $\dots$ ,  $\psi_p(x) = x^{p-1}$ . However, to improve the numerical stability of the procedure it is, for every local model, useful to apply the basis functions centered at the reference point  $X_i$ , that is, the functions  $(X_i - x)^m$ . This requires to slightly modify the description of the procedure.

Denote by  $\Psi(X_i)$  the  $p \times n$  matrix with the entries  $(X_i - X_l)^m$  for  $m = 0, 1, \dots, p-1$  and  $l = 1, \dots, n$ .

First we describe the estimation step of the algorithm. This step is performed similarly to the case described in Section 3.4. The only difference is that the family of basis functions (or, equivalently, the matrix  $\Psi$ ) depends on the central point  $X_i$ . Suppose that at the  $k$ th step of the procedure, for a point  $X_i$ , the diagonal weights matrix  $\widetilde{W}_i^{(k)}$  has been computed. Next we compute the  $p$ -vector  $Z_i^{(k)} = \eta Z_i^{(k-1)} + (1-\eta)\Psi(X_i)\widetilde{W}_i^{(k)}Y$  with the entries  $Z_{i,m}^{(k)}$  of the form

$$Z_{i,m}^{(k)} = \eta Z_{i,m}^{(k-1)} + (1-\eta) \sum_{l=1}^n \widetilde{w}_{il}^{(k)} (X_i - X_l)^m Y_l \quad m = 0, \dots, p-1,$$

and the matrix  $B_i^{(k)} = \eta B_i^{(k-1)} + (1-\eta)\Psi(X_i)\widetilde{W}_i^{(k)}\Psi^\top(X_i)$  whose entries are of the form  $B_{i,mm'}^{(k)} = b_{i,m+m'}^{(k)}$  for  $m, m' = 1, \dots, p$  where

$$b_{i,m}^{(k)} = \eta b_{i,m}^{(k-1)} + (1-\eta) \sum_{l=1}^n \widetilde{w}_{il}^{(k)} (X_i - X_l)^m \quad m = 0, \dots, 2p-2,$$

The estimate  $\widehat{\theta}_i^{(k)}$  in to the local model at  $X_i$ , is of the form  $\widehat{\theta}_i^{(k)} = \left(B_i^{(k)}\right)^{-1} Z_i^{(k)}$ .

For carrying out the  $k$ th adaptation step, we have to compare two estimates corresponding to the local models  $W_i^{(k-1)}$  and  $W_j^{(k-1)}$ . Note however, that this comparison can be done only if the both estimates are computed for the same basis system. Thus, the comparison requires to recompute the estimate for the local model  $W_j^{(k-1)}$  w.r.t. the basis centered at the point  $X_i$ . Let  $\widehat{\theta}_j = (\widehat{\theta}_{j,0}, \dots, \widehat{\theta}_{j,p-1})^\top$  be the estimate for the local model at  $X_j$ . This estimate leads to a local approximation of the unknown regression function by the polynomial  $\widehat{f}_j(x) = \widehat{\theta}_{j,0} + \widehat{\theta}_{j,1}(x - X_j) + \dots + \widehat{\theta}_{j,p-1}(x - X_j)^{p-1}$ . Now we represent this polynomial as a linear combination of the basis functions  $(x - X_i)^m$ ,  $m = 0, \dots, p-1$ , that is, we have to find new coefficients  $\widehat{\theta}_{ij} = (\widehat{\theta}_{ij,0}, \dots, \widehat{\theta}_{ij,p-1})^\top$  such that

$$\widehat{f}_j(x) = \widehat{\theta}_{ij,0} + \widehat{\theta}_{ij,1}(x - X_i) + \dots + \widehat{\theta}_{ij,p-1}(x - X_i)^{p-1}.$$

The coefficients  $\widehat{\theta}_{ij,m}$  can be computed from the formula  $\widehat{\theta}_{ij,m} = (m!)^{-1} d^m \widehat{f}_j(X_i) / dx^m$ .

The  $k$ th adaptation step of the procedure can be performed as follows. Suppose that all the estimates  $\widehat{\theta}_i^{(k-1)} = (\widehat{\theta}_{i,0}^{(k-1)}, \dots, \widehat{\theta}_{j,p-1}^{(k-1)})^\top$  have been computed in the previous step. Next, for a fixed  $i$  and every  $j$ , we compute the estimates  $\widehat{\theta}_{ij}^{(k-1)}$  by

$$\widehat{\theta}_{ij,m}^{(k-1)} = \sum_{q=0}^{p-m-1} \binom{q+m}{q} \widehat{\theta}_{j,q+m}^{(k-1)} (X_i - X_j)^q. \quad m = 0, 1, \dots, p-1.$$

The estimate  $\widehat{\theta}_{ij}^{(k-1)}$  is used in place of  $\widehat{\theta}_j^{(k-1)}$  for computing the statistical penalty  $s_{ij}^{(k)}$  in (3.4). For computing the extension penalty, we apply  $\Psi(X_i)$  in place of  $\Psi$  and  $\Psi_j$  has to be replaced by  $\Psi_j(X_i)$  which is the  $j$ th column of  $\Psi(X_i)$ . The remaining steps of the procedure are performed similarly to the basic algorithm.

### 4.3 Local linear multiple regression

Let  $X_1, \dots, X_d$  be points in the  $d$ -dimensional Euclidean space  $\mathbb{R}^d$ . Classical linear regression leads to an approximation of the regression function  $f$  by a linear combination of the constant function  $\psi_0(x) = 1$  and  $d$  coordinate functions  $\psi_m(x) = x_m$ , so that the family  $\{\psi_m\}$  consists of  $p = d+1$  basis functions. Our procedure attempts to apply this approximation locally for adaptively selected local models. The global linear modeling arises as a special case if the underlying model is entirely linear.

Similarly to the univariate case, we adopt for every design point  $X_i$  a local linear model with centered basis functions  $\psi_m(x, X_i) = X_{im} - x_m$  for  $m = 1, \dots, d$ . The corresponding  $p \times n$  matrix  $\Psi(X_i)$  has columns  $\Psi_l(X_i) = (1, X_{i1} - X_{l1}, \dots, X_{id} - X_{ld})^\top$  for  $l = 1, \dots, n$ . At the estimation step one computes the estimates  $\widehat{\theta}_i^{(k)}$  of the parameter  $\theta \in \mathbb{R}^p$  for every local model, leading to a local linear approximation of the function  $f$  by the linear function  $\widehat{f}_j(x)$  with

$$\widehat{f}_j(x) = \widehat{\theta}_{j,0} + \sum_{m=1}^d \widehat{\theta}_{j,m} (x_m - X_{j,m}).$$

This linear function can be rewritten in the form

$$\widehat{f}_j(x) = \widehat{\theta}_{j,0} + \sum_{m=1}^d \widehat{\theta}_{j,m} (X_{i,m} - X_{j,m}) + \sum_{m=1}^d \widehat{\theta}_{j,m} (x_m - X_{i,m}).$$

Therefore, only the first coefficient of the vector  $\widehat{\theta}_j$  has to be corrected when the basis system  $\Psi(X_i)$  is used in place of  $\Psi(X_j)$ . This means that at the  $k$ th adaptation step, the vector  $\widehat{\theta}_j^{(k-1)}$  is replaced by  $\widehat{\theta}_{ij}^{(k-1)}$  where  $\widehat{\theta}_{ij,m}^{(k-1)} = \widehat{\theta}_{j,m}^{(k-1)}$  for  $m = 1, \dots, d$  and  $\widehat{\theta}_{ij,0}^{(k-1)} = \widehat{\theta}_{i,0}^{(k-1)} + \sum_{m=1}^d \widehat{\theta}_{j,m} (X_{i,m} - X_{j,m})$ . The rest of the procedure is carried through similarly to the univariate case.

#### 4.4 Local quadratic bivariate regression

Finally we shortly discuss the bivariate case with  $d = 2$  for local quadratic approximation. The case of a larger  $d$  can be handled similarly. The family  $\{\psi_m\}$  of basis functions contains one constant function equal to 1, two linear coordinate functions  $x_1$  and  $x_2$  and three quadratic functions  $x_1^2, x_2^2$  and  $x_1x_2$ . It is useful to utilize the notation  $m = (m_1, m_2)$ ,  $|m| = m_1 + m_2$  and  $x^m = x_1^{m_1}x_2^{m_2}$  for  $x = (x_1, x_2)^\top \in \mathbb{R}^2$  and integers  $m_1, m_2$ . The family of basis functions can now be written in the form  $\{\psi_m(x), |m| \leq 2\}$ . For numerical stability the centered functions  $\psi_m(X_i - x)$  should be used within each local model.

At the  $k$ th estimation step one computes the entries  $\hat{\theta}_{i,m}^{(k)}$  of the vector  $\hat{\theta}_i^{(k)}$ . At the  $k$ th adaptation step we additionally need, for every  $i$ , to recompute the vectors  $\hat{\theta}_j^{(k-1)}$  for the basis system  $\Psi(X_i)$ . Similarly to the univariate case, we get

$$\hat{\theta}_{ij,m}^{(k-1)} = \sum_{m': |m'| \leq 2-|m|} \binom{m+m'}{m} \hat{\theta}_{j,m+m'}^{(k-1)} (X_i - X_j)^{m'}, \quad |m| \leq 2.$$

Here  $\sum_{m': |m'| \leq 2-|m|}$  means the sum over the set of all pair  $m' = (l'_1, l'_2)$  with  $m'_1 + m'_2 \leq 2 - m_1 - m_2$  and  $\binom{m}{m'} = \binom{m_1}{m'_1} \binom{m_2}{m'_2}$ . Particularly,  $\hat{\theta}_{ij,m}^{(k-1)} = \hat{\theta}_{j,m}^{(k-1)}$  for all  $m$  with  $|m| = 2$ , and  $\hat{\theta}_{ij,0} = \hat{f}_j(X_i)$ .

The rest of the procedure remains as before.

### 5 Numerical results

We now demonstrate the performance of the method for artificial examples in univariate and bivariate regression. The aim of this study is to illustrate two important features of the procedure: adaptability to large homogeneous regions and sensitivity to sharp changes in the local structure of the model. We also try to give some hints about the choice of the degree of local polynomial approximation.

Our univariate simulations are conducted as follows:

- Data are generated as  $(X_i, Y_i)$  with  $Y_i = f(X_i) + \varepsilon_i$ . Sample size is  $n = 1000$ . The design is chosen as an equidistant grid on  $(0, 1)$ . Errors  $\varepsilon_i$  are i.i.d. Gaussian with an unknown standard deviation  $\sigma$ .
- Local linear ( $p = 1$ ), local quadratic ( $p = 2$ ) and local cubic ( $p = 3$ ) AWS estimates are computed for each of 1000 simulated data sets. The parameters applied are given in Table 1.

Table 1: Parameters for AWS procedure in univariate regression

p	$\lambda$	$\eta$	$\tau$	$h_{\max}$
1	$q_{\chi^2; .65, 1}$	.5	4.5	0.25
2	$q_{\chi^2; .92, 2}$	.5	13.5	0.25
3	$q_{\chi^2; .92, 3}$	.5	40	0.25

- For a comparison a penalized cubic smoothing spline is fitted using the R-library *pspline*. The smoothing parameter is determined by generalized cross validation. See Heckman and Ramsey (2000) or the documentation of the R-library *pspline* (<http://www.r-project.org/>) for details.

**Remark 5.1.** The choice of the penalized cubic smoothing spline as the competitor for our procedure is explained by the excellent numerical results delivered by this method for many situations. We also tried other more sophisticated procedures like wavelets, point-wise adaptive procedures, Markow Random Fields methods but the numerical results (not reported here) were always in favor of smoothing splines, see also PJ2000.

### 5.1 Univariate Example 1

Our first example is based on a piecewise smooth function given by

$$f(x) = \begin{cases} 8x & : x < .125 \\ 2 - 8x & : .125 \leq x < .25 \\ 44(x - .4)^2 & : .25 \leq x < .55 \\ .5 \cos(6\pi(x - .775)) + .5 & : .55 \leq x \end{cases}$$

see the top left of Figure 1 for a graph. The upper row of Figure 1 shows plots of the first data set for  $\sigma = .125, .25$  and  $.5$ , respectively, together with a graph of the regression function. The bottom row reports the results in form of box-plots of Mean Absolute Error (MAE) obtained for the four procedures in 1000 simulation runs.

Figure 2, providing pointwise estimates of the Mean Absolute Error for three of the procedures in case of  $\sigma = .125$ , illustrates the behaviour of the procedures in their dependence on the features of the regression function. Note that especially the local linear AWS is superior to the cubic smoothing spline both near the discontinuities and within smooth regions. Local quadratic AWS seems to be more flexible near the first singularity, e.g.  $x = .125$ , and it behaves excellent for the rest of the design. Advantages are due to the local adaptivity of the AWS procedures in contrast to the global nature of the smoothing spline.

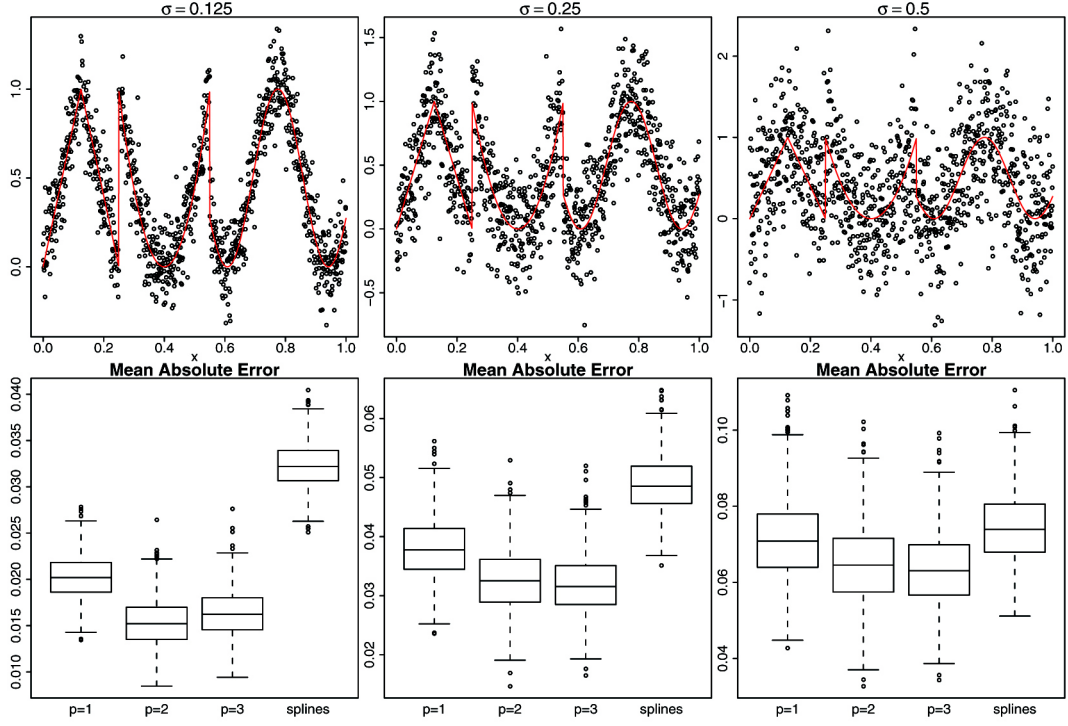


Figure 1: Univariate Example 1: Simulated data sets for  $\sigma = .125, .25$  and  $.5$  (upper row) and Box-Plots of MAE for local linear, local quadratic and local cubic AWS and penalized cubic smoothing splines, obtained from 1000 simulation runs (lower row).

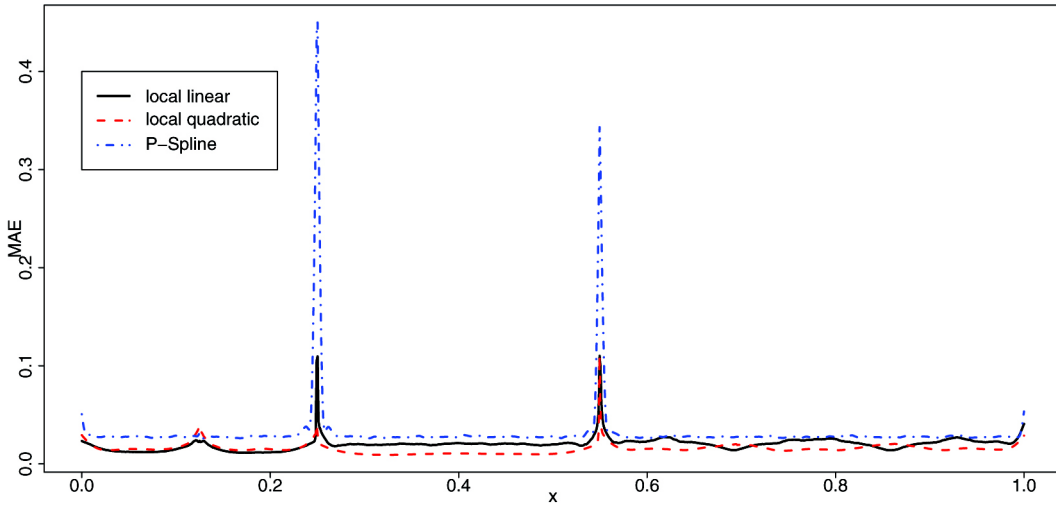


Figure 2: Example 1: Estimated pointwise Mean Absolute Error for local linear AWS, local quadratic AWS and penalized cubic smoothing splines in case of  $\sigma = 0.125$ .

## 5.2 Univariate example 2

The second univariate example uses a smooth regression function with varying second derivative:

$$f(x) = \sin(2.4\pi/(x + .2)) .$$

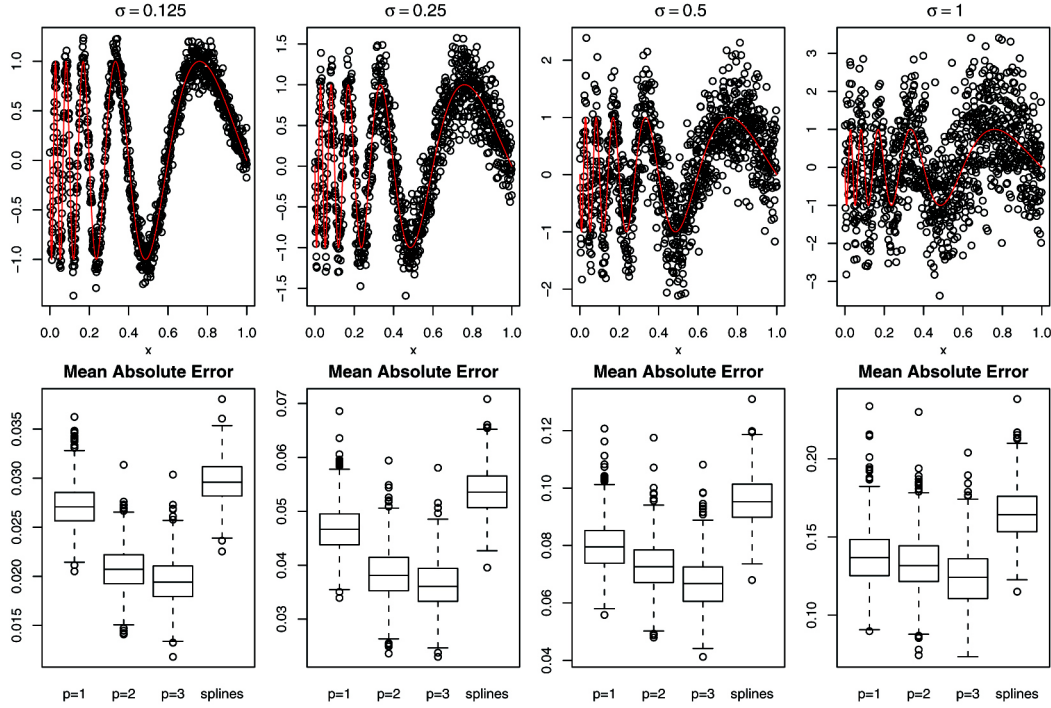


Figure 3: Univariate Example 2: Simulated data sets for  $\sigma = .125, .25, .5$  and  $1$  (upper row) and Box-Plots of MAE for local linear, local quadratic and local cubic AWS and penalized cubic smoothing splines, obtained from 1000 simulation runs (lower row).

The upper row of Figure 3 provides a view of a data set for  $\sigma = .125, .25, .5$  and  $1$ , respectively, and a graph of the regression function. The bottom row contains box-plots of Mean Absolute Error obtained for the four procedures in 1000 simulation runs.

Figure 4 again gives pointwise estimates of the Mean Absolute Error. Results are shown for local quadratic and cubic AWS and the penalized cubic smoothing splines in case of  $\sigma = .25$ . Note that the AWS procedures are superior in regions where the regression function is highly fluctuating or very smooth and loose compared to the smoothing spline only in case of medium fluctuation, i.e. for  $x \in (.05, .2)$ . For small values of  $x$  the spline suffers from high bias while for large values variability dominates. AWS delivers a good compromise, but with a little price for adaptation that can be seen in the region where the global smoothing parameter of the spline is nearly optimal.

### 5.3 Bivariate Example

We provide a bivariate example to demonstrate the behaviour of our procedure. Data are generated on a equidistant grid of  $100 \times 100$  points in  $[-1, 1]^2$  using the regression

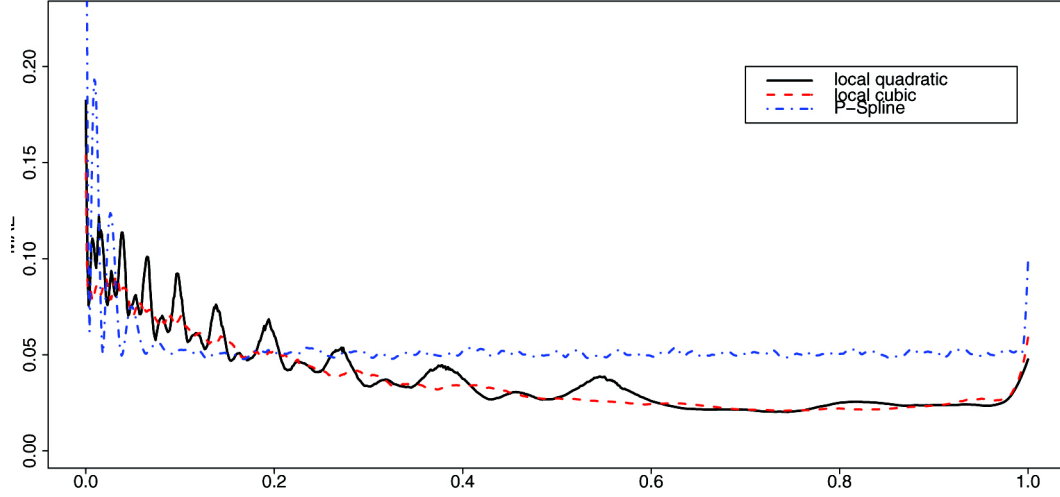


Figure 4: Example 2: Estimated pointwise Mean Absolute Error for local quadratic AWS, local cubic AWS and penalized cubic smoothing splines in case of  $\sigma = .25$ .

Table 2: Parameters for AWS procedure in bivariate regression

p	$\lambda$	$\eta$	$\tau$	$h_{\max}$
0	$q_{\chi^2; .966, 1}$	.5	—	0.2
1	$q_{\chi^2; .966, 3}$	.5	13.5	0.3
2	$q_{\chi^2; .966, 6}$	.5	100	0.4

function:

$$f(x, y) = (4x^2 + 8y^3)\text{sign}(4x^2 - 4xy - 6y^3).$$

The regression function is piecewise cubic with a discontinuity of varying strength along  $4x^2 - 4xy - 6y^3 = 0$ .

Table 2 provides the parameters used in the bivariate example. The upper left of Figure 5 shows a perspective plot of the data. The other three panels provide plots of the estimated surface obtained by local constant, local linear and local quadratic Adaptive Weights Smoothing. Local constant smoothing is not flexible enough to reasonably approximate the smooth part of the surface and introduces artificial segmentation. However the discontinuity is recovered rather well. The behavior of the local linear AWS is similar but the quality of approximation within the smooth part of the surface is drastically improved. The segmentation effect is only slightly observed. Local quadratic AWS delivers an almost perfect estimation quality both within smooth regions and near the edge.

For the global  $L_1$ -risk (Mean Absolute Error) of the considered estimates we got the



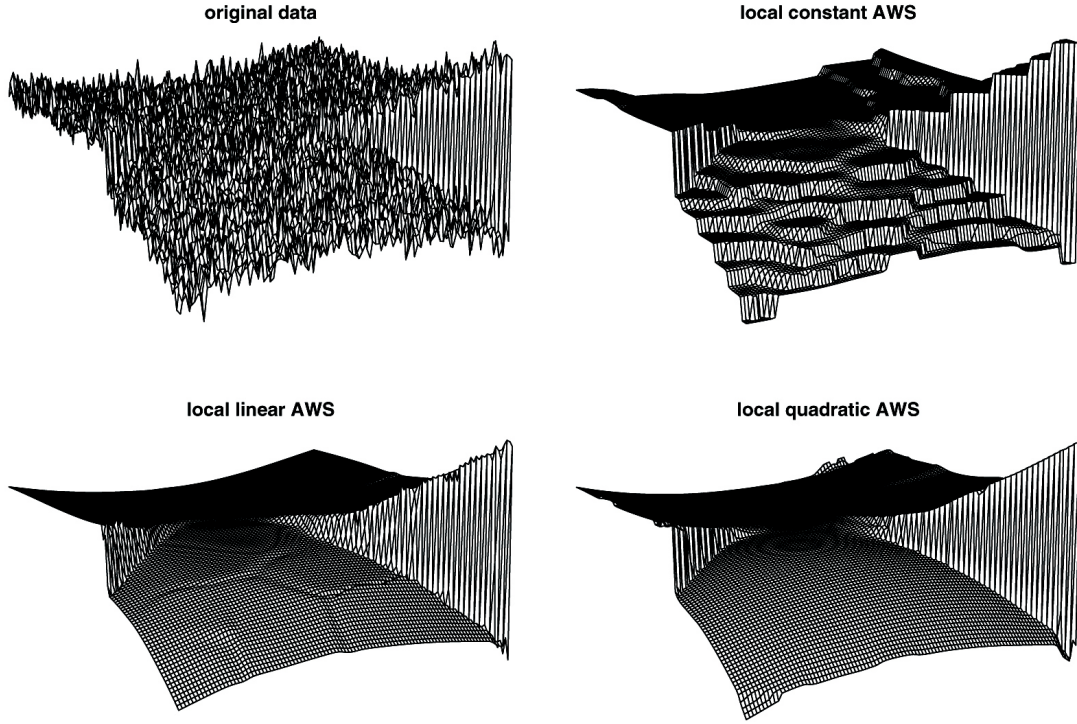


Figure 5: Bivariate Example: Perspective plots of data (upper left), local constant (upper right), local linear (lower left) and local quadratic (lower right) reconstruction.

following results: 0.090 for local constant AWS, 0.040 for local linear AWS, 0.030 for local quadratic AWS.

#### 5.4 Summary

The results of the simulation study can be summarized as follows. The performance of the method is completely in agreement with what was aimed: it is adaptive to variable smoothness properties of the underlying function and sensitive to discontinuities outperforming the classical smoothing methods. It demonstrates excellent results for a small or moderate noise and it is stable with respect to the large noise.

The local linear AWS seems to be a reasonable compromise for many situations combining a good approximating properties with very good quality of change-point or edge estimation. In situations with large homogeneous regions, local polynomial approximation of a higher order can be slightly preferable. The choice of the polynomial degree can be also done automatically using global cross-validation type criteria.

The procedure is rather stable w.r.t. to the choice of parameters as long as  $\lambda$  is not

chosen too small or too large.

## 6 Some important properties of AWS

This section discusses some properties of the AWS procedure. In particular we establish the “propagation condition” which means a free extension of every local model in a homogeneous situation, leading to a nearly parametric estimate at the end of iteration process. Further we discuss the rate of estimation for a smooth function  $\theta(x)$ . We start by listing some attractive features of the method which directly follow from the construction and are also justified by our numerical results.

**AWS applies in a unified way to a broad class of regression models:** This means that the procedure is able to adapt to the unknown and variable function structure without requiring any specific prior information like the degree of smoothness of the underlying regression function. Modeling by weights allows to proceed in a unified way with parametric and change point models, and also with models whose parameters vary smoothly.

**Weak model assumptions:** The procedure does not require any specific information about the noise distribution. The noise variance can be estimated from the data. If some prior information about the noise distribution or about the model is available, it can be easily incorporated in the method allowing a better estimation of the (local) noise variance. However, the present procedure is designed towards regression-like models with additive noise. Specific noise models like binary-response, exponential or Poisson etc. require a special treatment, see PS2002 for more details.

**AWS is design adaptive and has no boundary problem:** The method proceeds with the given “design”  $X_1, \dots, X_n$ , no assumptions or restrictions are imposed on it. The random design can be treated similarly to the case of a deterministic design. The local polynomial modeling applied in the algorithm does not suffer from nonregular design. This feature is important in connection with change point and edge estimation. The produced estimate does not indicate the usual Gibbs effect (high variability) near discontinuities like most of the other nonparametric methods.

**AWS applies for high dimensional models:** The procedure is described for the case of a design in a finite dimensional Euclidean space  $\mathbb{R}^d$ . This assumption however was only made for convenience of exposition. The dimensionality  $d$  of the regressors plays

absolutely no role for the procedure. Moreover, the procedure can be formulated for the design in an arbitrary metric space  $\mathcal{X}$ . Only the number  $p$  of parameters entering into the description of the underlying parametric model is essential. However, for local linear or local polynomial modeling, the number of parameters grows dramatically with the dimension  $d$ , and the procedure can face the so called “curse of dimensionality” problem: in high dimension, pure nonparametric modeling leads to strong oversmoothing. Specifically for the AWS method, if the number of parameters becomes too high (say, more than 6) then the procedure loses sensitivity to structural changes. For such situation, combining the procedure with some dimension reduction methods can be useful.

**AWS is computationally straightforward and the numerical complexity can be easily controlled:** Indeed, AWS requires of order  $nM_{k^*}$  operations with  $k^*$  being the number of iterations and  $M_k$  being the corresponding size of the typical neighborhood  $U_i^{(k)}$  at the step  $k$ . Therefore, the complexity of the method can be controlled simply by restricting  $k^*$ , or, equivalently the largest bandwidth  $h_{\max}$ , see Section 3.4.

Now we turn to the more involved properties of the method which require a theoretical justification.

### 6.1 Behaviour inside homogeneous regions. Propagation condition

The procedure is designed to provide a free extension of every local model within a large homogeneous region. An extreme case is given by a fully parametric homogeneous model. In that case, a desirable feature of the method is that the final estimate at every point coincides with high probability with the fully parametric global estimate. This property which we call the “propagation” condition is proved here under some simplifying assumptions.

The analysis of the properties of the iterative estimates  $\widehat{\theta}_i^{(k)}$  is very difficult. The main reason is that every estimate  $\widehat{\theta}_i^{(k)}$  solves the local likelihood problem for the local model defined by the weights  $w_{ij}^{(k)}$  which are random and depend on the same observations  $Y_1, \dots, Y_n$ . To tackle this problem we make the following assumption:

**(A0)** for every step  $k$  an independent sample  $Y_1, \dots, Y_n$  is available so that the weights  $w_{ij}^{(k)}$  are independent of the sample  $Y_1, \dots, Y_n$  for every  $k$ .

This assumption can be realized by splitting the original sample into  $k^*$  subsamples. Since the number of steps is only of logarithmic order this split can change of the quality

of estimation only by a logarithmic factor. Of course, such a split is only a theoretical device, the use of the same sample for all steps of the algorithm still requires a further justification.

In our study we restrict ourselves to the case of the varying coefficient model with homogeneous Gaussian noise:

**(A1)** The observations  $Y_1, \dots, Y_n$  follow the model  $Y_i = f(X_i) + \varepsilon_i$  where  $\varepsilon_1, \dots, \varepsilon_n$  are i.i.d.  $\mathcal{N}(0, \sigma^2)$ .

To simplify the presentation we also assume that

**(A2)** The statistical penalty  $s_{ij}^{(k)}$  is defined via the likelihood ratio test statistic  $T_{ij}^\circ$  from (3.2) in Section 3.1.

In our procedure the statistic  $T_{ij}$  from (3.3) is applied. However, an essential difference between  $T_{ij}$  and  $T_{ij}^\circ$  only occurs in the situations where the local models  $W_i$  and  $W_j$  are strongly unbalanced, which do not meet in the specific cases considered here.

**(A3)** The extension penalty  $e_{ij}^{(k)}$  is set to zero, that is,  $K_e(e_{ij}^{(k)}) = 1$ .

Again, this assumption is not restrictive because the extension penalty does not matter as long as the propagation condition is studied.

We first consider the homogeneous situation with  $\theta_i = \theta$  which corresponds to a global linear model  $f(x) = \theta_1 \psi_1(x) + \dots + \theta_p \psi_p(x)$ .

**Theorem 6.1.** *Let (A0), (A1), (A2) and (A3) be fulfilled. Suppose that  $\theta(X_i) \equiv \theta$ , i.e.  $f = \Psi\theta$ . If  $\lambda \geq C \log n$  with constant  $C$  depending on the kernel  $K_s$  only, then for every iteration  $k$*

$$\mathbf{P} \left( \min_{i,j=1,\dots,n} K_s(s_{ij}^{(k)}) > 1/2 \right) \geq 1 - 1/n.$$

*Proof.* Define  $b$  by the equation  $K_s(b) = 1/2$ . Theorem 7.2 from the Appendix yields for every iteration  $k$

$$\mathbf{P} \left( \min_{i,j=1,\dots,n} K_s(s_{ij}^{(k)}) > 1/2 \right) = \mathbf{P} \left( \max_{i,j=1,\dots,n} T_{ij}^{(k)} \leq b\lambda \right) \geq 1 - \sum_{i,j=1}^n q_p(b\lambda - p)$$

where  $q_p(u)$  is defined by  $\log q_p(u) = -u/2 + 0.5p \log(1 + u/p)$ . It is easy to see that  $q_p(u)$  fulfills  $\log q_p(u) \leq -2 \log n$  for  $u \geq C_p \log n$  with some constant  $C_p$  depending on  $p$  only. This yields the assertion as soon as  $b\lambda - p \geq C_p \log n$ , or, equivalently,  $\lambda \geq (p + C_p \log n)/b$ .  $\square$

This result means that the statistical penalty entering in the weights  $w_{ij}^{(k)}$  at every iteration  $k$  does not restrict a free extension of any local model.

**Corollary 6.1.** *Let the assumptions (A0), (A1), (A2) and (A3) be fulfilled and  $\theta(X_i) \equiv \theta$ . If  $\lambda \geq C \log n$  and if  $h_{\max}$  is sufficiently large then the last step estimate  $\hat{\theta}_i = \hat{\theta}_i^{(k^*)}$  fulfills for every  $z \geq 0$*

$$\mathbf{P} \left( (2\sigma^2)^{-1} (\hat{\theta}_i - \theta)^\top \Psi \Psi^\top (\hat{\theta}_i - \theta) > p + z \right) \leq q_p(z)$$

where  $\log q_p(u) = -u/2 + 0.5p \log(1 + u/p)$ .

*Proof.* If  $h_{\max}$  is sufficiently large then the location penalty  $K_l(\mathbf{l}_{ij}^{(k)})$  at the final iteration  $k = k^*$  fulfills  $K_l(\mathbf{l}_{ij}^{(k)}) \approx 1$  for every pair  $(i, j)$ . By Theorem 6.1 the statistical penalty  $K_s(\mathbf{s}_{ij}^{(k)}) \geq 1/2$ , hence  $w_{ij}^{(k)} \geq 1/2$  for all  $(i, j)$ . This yields  $\Psi W_i^{(k)} \Psi^\top \geq 0.5 \Psi \Psi^\top$  and the result follows from Theorem 7.1 from the Appendix.  $\square$

Due to this result the final estimate  $\hat{\theta}_i = \hat{\theta}_i^{(k^*)}$  delivers the same quality of estimation as the global LSE  $\hat{\theta} = (\Psi \Psi^\top)^{-1} \Psi Y$ . In fact, one can show an even stronger assertion: with a high probability it holds  $\hat{\theta}_i \approx \hat{\theta}$ . The explanation is as follows. Our way of computing the statistical penalty  $\mathbf{s}_{ij}^{(k)}$  does not take into account that two “local” models  $W_i$  and  $W_j$  have nonzero intersection. This means that there are some points  $X_l$  such that the weights  $w_{il}^{(k)}$  and  $w_{jl}^{(k)}$  are simultaneously positive and hence, the estimates  $\hat{\theta}_i^{(k)}$  and  $\hat{\theta}_j^{(k)}$  are dependent and positively correlated. In the homogeneous situation, for every two fixed points, this dependence grows with iteration, so that the estimates  $\hat{\theta}_i^{(k)}$  and  $\hat{\theta}_j^{(k)}$  become more and more positively correlated. In the extreme case they are almost identical at the end and the statistical penalty vanishes.

The propagation condition can be easily extended to the case of a large homogeneous region  $G$  in  $\mathcal{X}$ . Define for every  $x \in G$  the distance from  $x$  to the boundary of  $G$ , i.e.  $\rho_G(x) = \min\{\rho(x, X_j) : X_j \notin G\}$ . At every step  $k$  we consider only internal points  $X_i \in G$  which are separated from the boundary with the distance  $2h^{(k)}$ :

$$\mathcal{G}^{(k)} = \{X_i \in G : \rho_G(X_i) \geq 2h^{(k)}\}.$$

The next result claims the propagation condition (free extension) for all such points.

**Theorem 6.2.** *Let the assumptions (A0), (A1) and (A2) be fulfilled. Suppose that  $\theta(X_i) \equiv \theta$  for all  $X_i$  from some region  $G$  in  $\mathcal{X}$ . If  $\lambda \geq C \log n$  for some constant  $C$*

depending on the kernel  $K_s$  only, then for every iteration  $k$

$$\mathbf{P} \left( \min_{(i,j): X_i \in \mathcal{G}^{(k)}, \rho(X_i, X_j) \leq h^{(k)}} K_s(\mathbf{s}_{ij}^{(k)}) > 1/2 \right) \geq 1 - 1/n.$$

*Proof.* It suffices to note that if  $X_i \in \mathcal{G}^{(k)}$  then the local model  $W_i^{(k)}$  as well as all the models  $W_j^{(k)}$  for all  $X_j$  with  $\rho(X_i, X_j) \leq h^{(k)}$  are homogeneous. Hence, the result follows again by Theorem 7.2.  $\square$

## 6.2 Accuracy of estimation for a varying coefficient model

Here we consider the case of an arbitrary function  $f$  which allows a good linear approximation in a neighborhood of a point  $x \in \mathcal{X}$ . We first show that this condition ensures a free extension of all the local models within this neighborhood.

Let a design point  $x = X_i$  for some  $i$  be fixed, and let  $h$  be some bandwidth used in the iteration procedure. We define  $U_h(x) = \{x' : |x' - x| \leq h\}$ . We consider the following conditions which are specified for the fixed point  $x$  and the bandwidth  $h$ :

**(A4)** It holds  $|\Psi_j^\top [\theta(X_j) - \theta(x)]| \leq \delta$  for some  $\delta \geq 0$  and all  $X_j \in U_h(X_i)$ .

**(A5)** The kernel  $K_l$  is compactly supported on  $[0, 1]$ .

**(A6)** Define  $W_i^* = \text{diag}\{w_{i1}^*, \dots, w_{in}^*\}$  with  $w_{ij}^* = K_l(|\rho(X_i, X_j)/h|^2)$ ,  $N_i^* = \text{tr} W_i^*$  and  $B_i^* = \Psi W_i^* \Psi^\top$ . It holds

$$N_i^* \Psi_i \Psi_i^\top \leq C_B B_i^*.$$

Condition (A4) means that the function  $f(X_j)$  can be approximated by a linear function  $\Psi_j^\top \theta(x)$  with the precision  $\delta$  for every  $X_j \in U_h(x)$ . Condition (A6) guarantees a certain design regularity in a neighborhood of the reference point  $x$ . The next result claims the propagation condition (free extension) for the local models  $W_i^{(k)}$  as long as  $h^{(k)} \leq h$  provided that  $\delta$  is sufficiently small.

**Theorem 6.3.** *Let the assumptions (A0) through (A6) be fulfilled. Let  $\lambda \geq C \log n$  for some constant  $C$  depending on the kernel  $K_s$  only. If*

$$2\sigma^{-2} p \delta^2 (N_i^* + N_j^*) \leq b\lambda/6, \quad X_j \in U_h(x), \quad (6.1)$$

where  $b$  is defined by  $K(b) = 1/2$ , then for every iteration  $k$  with  $h^{(k)} \leq h$

$$\mathbf{P} \left( \min_{j: X_j \in U_h(X_i)} K_s(\mathbf{s}_{ij}^{(k)}) \geq 1/2 \right) \geq 1 - 1/n. \quad (6.2)$$

If  $h = h^{(k)}$ , then it holds with a probability of at least  $1 - 2/n$

$$\mathbf{P} \left( \left| \widehat{f}_i^{(k)} - f_i \right| > \sqrt{pC_B\delta} + \sigma\sqrt{2C_B\lambda/N_i^*} \right) \leq 2/n. \quad (6.3)$$

The proof is given in the Appendix. The result (6.3) indicates that the first  $k$  iterations of the procedure (for  $h^{(k)} \leq h$ ) lead to a reasonable quality of estimation of the function  $f(\cdot)$ . However, the procedure has to prevent from losing the obtained quality of estimation during further iterations. This is precisely what the *control step* of the original AWS procedure from PS2000 does, see the discussion at the end of Section 3.5. The procedure presented here applies this control step in a soft form, however we only show how the *hard* control step (3.6) can be used for proving the rate result.

**Theorem 6.4.** *Let the conditions of Theorem 6.3 be fulfilled and let the procedure involve the control step from (3.6) with  $\eta^* \geq \lambda$ . Then the last step estimate  $\widehat{f}_i = \Psi_i^\top \widehat{\theta}_i^{(k^*)}$  fulfills with a probability of at least  $1 - 2/n$*

$$\left| \widehat{f}_i - f_i \right| \leq \sqrt{pC_B\delta} + \sigma\sqrt{2C_B\lambda/N_i^*} + \sigma\sqrt{2C_B\eta^*/N_i^*}.$$

*Proof.* Let  $h = h^{(k)}$  for some  $k$ . The control step (3.6) ensures that

$$(2\sigma^2)^{-1} (\widehat{\theta}_i^{(k)} - \widehat{\theta}_i^{(k^*)})^\top B_i^{(k)} (\widehat{\theta}_i^{(k)} - \widehat{\theta}_i^{(k^*)}) \leq \eta^*.$$

This yields by (A6)

$$\begin{aligned} N_i^{(k)} |\widehat{f}_i^{(k)} - \widehat{f}_i|^2 &= N_i^{(k)} (\widehat{\theta}_i^{(k)} - \widehat{\theta}_i^{(k^*)})^\top \Psi_i \Psi_i^\top (\widehat{\theta}_i^{(k)} - \widehat{\theta}_i^{(k^*)}) \\ &\leq C_B (\widehat{\theta}_i^{(k)} - \widehat{\theta}_i^{(k^*)})^\top B_i^{(k)} (\widehat{\theta}_i^{(k)} - \widehat{\theta}_i^{(k^*)}) \leq 2\sigma^2 C_B \eta^*. \end{aligned}$$

By Theorem 6.3  $N_i^{(k)} \geq 0.5N_i^*$  with a high probability and the assertion follows directly from (6.3).  $\square$

### 6.3 Rate of estimation for a smooth function $f(\cdot)$ . Spatial adaptivity

Here we briefly discuss one important special case of the result of Theorem 6.4. Namely, we suppose that  $f(\cdot)$  is a smooth function in  $\mathbb{R}^d$  and consider the polynomial basis  $\{\psi_m\}$  of degree less than a given integer number  $s$ . In the univariate case  $d = 1$  there are exactly  $p = s$  basis functions, e.g.  $1, u - x, \dots, (u - x)^{s-1}$ . We also suppose that

**(A4s)** The function  $f(\cdot)$  is  $s$  times continuously differentiable and  $|f^{(s)}(u)| \leq Ls!$  for some constant  $L$  and all  $u \in U_h(x)$ .

**(A7)** For some positive constants  $C_{X1} \leq C_{X2}$  holds for all  $h \in [h^{(0)}, h_{\max}]$

$$C_{X1} \leq N_h^*/(nh^d) \leq C_{X2}.$$

where  $N_h^* = \sum_{j=1}^n K_l(|\rho(X_i, X_j)/h|^2)$ .

Note that condition (A4s) ensures (A4) with  $\delta = Lh^s$ . We now apply Theorem 6.4 to this situation with  $\eta^* = \lambda$ . The result is formulated as a separate statement.

**Theorem 6.5.** *Suppose that (A0), (A1), (A2), (A3), (A4s), (A5), (A7) are fulfilled and (A6) holds for all  $h \in [h^{(0)}, h_{\max}]$ . If  $\lambda \geq C \log n$  for some fixed  $C$ , then*

$$\mathbf{P} \left( |\hat{f}_i - f_i| > C_1 (\lambda \sigma^2 / n)^{s/(d+2s)} L^{d/(d+2s)} \right) \leq 2/n$$

where the constant  $C_1$  depends on  $C_{X1}, C_{X2}$  and  $C_B$  only.

*Proof.* The bound (6.3) and condition (A7) imply with a high probability

$$\left| \hat{f}_i - f_i \right| \leq \sqrt{pC_B} \delta + 2\sigma \sqrt{2C_B \lambda / N_i^*} \leq \sqrt{pC_B} Lh^s + 2\sigma \sqrt{2C_B \lambda / (C_{X1} n h^d)}.$$

Optimizing this expression w.r.t.  $h$  leads to the choice  $h = C_2 \{\lambda \sigma^2 / (n L^2)\}^{1/(d+2s)}$ . With this choice the condition (6.1) is fulfilled in view of (A7) provided that  $C_2$  is not too large. The use of such selected  $h$  results in the accuracy of order  $\{\lambda \sigma^2 / n\}^{s/(d+2s)} L^{d/(d+2s)}$  as required.  $\square$

The accuracy shown in Theorem 6.5 is optimal in rate for the problem of estimation of a smooth function  $f$  up to a logarithmic factor  $\lambda$ . Therefore, this result means that our procedure is pointwise adaptive in the sense that it automatically adapts to the unknown local smoothness degree measured by the exponent  $s$  and the Lipschitz constant  $L$ . As shown in Lepski, Mammen and Spokoiny (1997) this property automatically leads to rate optimality in the Sobolev and Besov function classes  $B_{p,q}^s$ .

## 7 Appendix

Here we present some general results on large deviation probabilities for local likelihood ratio test statistics in Gaussian regression.

We consider the varying coefficient regression model  $Y_i = f(X_i) + \varepsilon_i$  with Gaussian homogeneous errors  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ . The local model  $W$  is described by the weights  $w_1, \dots, w_n$ . Local linear modeling assumes the linear structure of the model function  $f$



within the local model  $W$ :  $f(x) = \theta_1\psi_1(x) + \dots + \theta_p\psi_p(x)$  for a given system  $\{\psi_m(x)\}$ . The corresponding local MLE  $\hat{\theta}$  can be represented in the form  $\hat{\theta} = (\Psi W \Psi^\top)^{-1} \Psi W Y$  with the notation from Section 2.2. The local likelihood ratio test statistic is defined for a given  $\theta$  by  $L(W, \hat{\theta}, \theta) = (\hat{\theta} - \theta)^\top B(\hat{\theta} - \theta) / (2\sigma^2)$  where  $B = \Psi W \Psi^\top$ .

### 7.1 Linear parametric case

Define  $\bar{\theta} = B^{-1} \Psi W f$ . Then  $\Psi \bar{\theta}$  is the best linear approximation of  $f$  within the local model  $W$ . In the homogeneous case  $f = \Psi^\top \theta$ , it obviously holds  $\bar{\theta} = \theta$ . The first result we present shows that  $\hat{\theta}$  is a good estimate of the vector  $\bar{\theta}$ . This particularly implies a nice properties of the estimate in a homogeneous situation when the local linear assumption is fulfilled and  $\theta$  is the true parameter.

**Theorem 7.1.** *For every  $z \geq 0$*

$$P\left(2L(W, \hat{\theta}, \bar{\theta}) > p + z\right) \leq q_p(z)$$

where

$$q_p(z) = \exp(-0.5z + 0.5p \log(1 + z/p)). \quad (7.1)$$

*Proof.* The model equation  $Y = f + \varepsilon$  immediately implies that  $\hat{\theta}_i = B_i^{-1} \Psi W_i Y = \bar{\theta}_i + B_i^{-1} \Psi W_i \varepsilon$ . Therefore,  $\hat{\theta}_i - \bar{\theta}_i = B_i^{-1} \Psi W_i \varepsilon$  does not depend on  $\theta$ , and we assume without loss of generality that  $\theta = 0$ , so that the observations  $Y_i$  coincide with the noise  $\varepsilon_i$ . This obviously implies  $E\hat{\theta} = 0$ . The covariance matrix  $V$  of the estimate  $\hat{\theta}$  can be represented as

$$V = E\hat{\theta}\hat{\theta}^\top = EB^{-1}\Psi\varepsilon\varepsilon^\top\Psi^\top B^{-1} = \sigma^2 B^{-1} D B^{-1}$$

where  $D = \Psi W^2 \Psi^\top$ . Therefore, the estimate  $\hat{\theta}$  can be represented as  $\hat{\theta} = V^{1/2} \zeta$  where  $\zeta$  is a standard Gaussian random vector in  $\mathbb{R}^p$ . This yields

$$L(W, \hat{\theta}, \theta) = (2\sigma^2)^{-1} \zeta^\top V^{1/2} B V^{1/2} \zeta = 0.5 \zeta^\top R \zeta$$

with  $R = B^{-1/2} D B^{-1/2}$ . Since  $w_i \leq 1$ , it holds  $D \leq B$  and  $\|R\| \leq 1$ , that is, the largest eigenvalue of  $R$  does not exceed one. Now the desired result follows from the general result for Gaussian quadratic forms in Lemma 7.1.  $\square$

**Lemma 7.1.** *Let a symmetric  $p \times p$ -matrix  $R$  fulfill  $\|R\| \leq 1$ . Then*

$$P\left(\zeta^\top R \zeta \geq p + z\right) \leq q_p(z).$$

*Proof.* Let  $r_1, \dots, r_p$  be the eigenvalues of  $R$  satisfying  $r_m \leq 1$  for all  $m$ . It holds for every  $\mu < 1$  by simple algebra

$$\log \mathbf{E} \exp(\mu \zeta^\top R \zeta / 2) = \log \prod_{m=1}^p \frac{1}{\sqrt{1 - \mu r_m}} = -\frac{1}{2} \sum_{m=1}^p \log(1 - \mu r_m) \leq -0.5p \log(1 - \mu).$$

Now the exponential Chebyshev inequality implies

$$\begin{aligned} \log \mathbf{P} \left( 0.5 \zeta^\top R \zeta \geq (p+z)/2 \right) &\leq -\mu(p+z)/2 + \log \mathbf{E} \left( 0.5 \mu \zeta^\top R \zeta \right) \\ &\leq -0.5\mu(p+z) - 0.5p \log(1 - \mu). \end{aligned}$$

This expression is maximized by  $\mu = z/(p+z)$  leading to

$$\log \mathbf{P} \left( \zeta^\top R \zeta \geq p+z \right) \leq -0.5z + 0.5p \log(1 + z/p)$$

as required.  $\square$

Next we consider the likelihood ratio test statistic  $T_{ij}^\circ$  defined in Section 3.6 for two local models  $W_i$  and  $W_j$ .

**Theorem 7.2.** *Let  $f = \Psi^\top \theta$ . Then for every  $z \geq 0$*

$$\mathbf{P} \left( T_{ij}^\circ > p+z \right) \leq q_p(z).$$

*Proof.* We use the representation  $2T_{ij}^\circ = \sigma^{-2}(\hat{\theta}_i - \hat{\theta}_j)^\top B_i(B_i + B_j)^{-1} B_j(\hat{\theta}_i - \hat{\theta}_j)$ . Note that

$$\text{Cov}(\hat{\theta}_i - \hat{\theta}_j) \leq 2 \text{Cov}(\hat{\theta}_i) + 2 \text{Cov}(\hat{\theta}_j) = 2V_i + 2V_j \leq 2\sigma^2(B_i^{-1} + B_j^{-1}).$$

Now the result follows from Lemma 7.1 similarly to the proof of Theorem 7.1.  $\square$

## 7.2 Sufficient conditions for free extension

Now we consider the general situation of a varying coefficient model. We show that if the difference between two local models defined in terms of the Kullback-Leibler distance, is sufficiently small, then  $T_{ij}^\circ$  is with a large probability smaller than  $b\lambda$  for some  $b \leq 1$ .

**Theorem 7.3.** *Let  $b \in (0, 1]$  be such that  $z = b\lambda/2 - p > 0$ . Then the condition*

$$\Delta := 0.5\sigma^{-2}(\bar{\theta}_i - \bar{\theta}_j)^\top B_i(B_i + B_j)^{-1} B_j(\bar{\theta}_i - \bar{\theta}_j) \leq b\lambda/6 \quad (7.2)$$

*with  $\bar{\theta}_i = B_i^{-1}\Psi W_i f$  and  $\bar{\theta}_j = B_j^{-1}\Psi W_j f$  implies*

$$\mathbf{P} \left( T_{ij}^\circ > b\lambda \right) \leq q_p(z) + e^{-b\lambda/12}.$$

*Proof.* We use the decomposition

$$\widehat{\theta}_i - \widehat{\theta}_j = \xi_i - \xi_j + \bar{\theta}_i - \bar{\theta}_j$$

where  $\xi_i = B_i^{-1}\Psi W_i \varepsilon$  and similarly for  $\xi_j$ . This implies with  $B_{ij} = B_i(B_i + B_j)^{-1}B_j$

$$2\sigma^2 T_{ij}^\circ = (\xi_i - \xi_j)^\top B_{ij}(\xi_i - \xi_j) + (\bar{\theta}_i - \bar{\theta}_j)^\top B_{ij}(\bar{\theta}_i - \bar{\theta}_j) + 2(\bar{\theta}_i - \bar{\theta}_j)^\top B_{ij}(\xi_i - \xi_j). \quad (7.3)$$

The result of Theorem 7.2 implies

$$\mathbf{P}\left(\sigma^{-2}(\xi_i - \xi_j)^\top B_{ij}(\xi_i - \xi_j) > p + z\right) \leq q_p(z).$$

Next,  $\zeta_{ij} = \sigma^{-2}(\bar{\theta}_i - \bar{\theta}_j)^\top B_{ij}(\xi_i - \xi_j)$  is a Gaussian random variable with zero mean satisfying

$$\begin{aligned} \mathbf{E}\zeta_{ij}^2 &= \sigma^{-4}(\bar{\theta}_i - \bar{\theta}_j)^\top B_{ij} \text{Cov}(\xi_i - \xi_j) B_{ij}(\bar{\theta}_i - \bar{\theta}_j) \\ &\leq 2\sigma^{-2}(\bar{\theta}_i - \bar{\theta}_j)^\top B_{ij}(\bar{\theta}_i - \bar{\theta}_j) \leq 4\Delta. \end{aligned} \quad (7.4)$$

Here we have used that  $\text{Cov}(\xi_i - \xi_j) \leq 2\sigma^2 B_{ij}$ , see the proof of Theorem 7.2. This and condition (7.2) imply

$$\mathbf{P}(\zeta_{ij} > b\lambda/3) \leq e^{-b\lambda/12}.$$

Since  $p + z = b\lambda$ , we finally obtain

$$\begin{aligned} \mathbf{P}(T_{ij}^\circ > b\lambda) &\leq \mathbf{P}\left(0.5\sigma^{-2}(\xi_i - \xi_j)^\top B_{ij}(\xi_i - \xi_j) \geq p + z\right) + \mathbf{P}(\zeta_{ij} > b\lambda/3) \\ &\leq q_p(z) + e^{-b\lambda/12} \end{aligned}$$

as required.  $\square$

The next assertion delivers some sufficient conditions ensuring (7.2). More precisely, we consider the situation when the function  $f$  can be well approximated by a linear function  $\Psi^\top \theta$  within both local models  $W_i$  and  $W_j$ . If  $|f(X_l) - \Psi_l^\top \theta| \leq \delta$  for some small positive  $\delta$  and all  $X_l$  entering with positive weight in the model  $W_i$ , then  $(f - \Psi^\top \theta)^\top W_i(f - \Psi^\top \theta) = \sum_l w_{il} |f(X_l) - \Psi_l^\top \theta|^2 \leq N_i \delta^2$  with  $N_i = \sum_l w_{il}$  and similarly for the model  $W_j$ .

**Lemma 7.2.** *The condition*

$$(f - \Psi^\top \theta)^\top W_i(f - \Psi^\top \theta) \leq \delta^2 N_i$$

implies

$$(\bar{\theta}_i - \theta)^\top B_i (\bar{\theta}_i - \theta) \leq p\delta^2 N_i.$$

If, in addition,  $(f - \Psi^\top \theta)^\top W_j (f - \Psi^\top \theta) \leq N_j \delta^2$ , then

$$(\bar{\theta}_i - \bar{\theta}_j)^\top B_{ij} (\bar{\theta}_i - \bar{\theta}_j) \leq 2p\delta^2 (N_i + N_j)$$

where  $B_{ij} = B_i(B_i + B_j)^{-1}B_j$ .

*Proof.* The use of  $B_i = \Psi W_i \Psi^\top$  and  $\bar{\theta}_i = B_i^{-1} \Psi W_i f$  gives

$$(\bar{\theta}_i - \theta)^\top B_i (\bar{\theta}_i - \theta) = (f - \Psi^\top \theta)^\top W_i \Psi^\top B_i^{-1} \Psi W_i (f - \Psi^\top \theta)$$

Define  $A = W_i^{1/2} \Psi^\top B_i^{-1} \Psi W_i^{1/2}$ . Then

$$\text{tr} A A^\top = \text{tr} W_i^{1/2} \Psi^\top B_i^{-1} \Psi W_i \Psi^\top B_i^{-1} \Psi W_i^{1/2} = \text{tr} B_i^{-1} \Psi W_i \Psi^\top = \text{tr} I_p = p.$$

Therefore, by the Cauchy-Schwarz inequality

$$|(\bar{\theta}_i - \theta)^\top B_i (\bar{\theta}_i - \theta)|^2 \leq \|W_i^{1/2} (f - \Psi^\top \theta)\|^2 \text{tr} A A^\top \leq N_i \delta^2 p.$$

and the first assertion follows.

Since  $B_{ij} \leq B_i$  and similarly  $B_{ij} \leq B_j$ , it holds

$$(\bar{\theta}_i - \bar{\theta}_j)^\top B_{ij} (\bar{\theta}_i - \bar{\theta}_j) \leq 2(\bar{\theta}_i - \theta)^\top B_i (\bar{\theta}_i - \theta) + 2(\bar{\theta}_j - \theta)^\top B_j (\bar{\theta}_j - \theta).$$

and the second assertion follows as well.  $\square$

### 7.3 Separability condition

Now we present some sufficient conditions for separability of two local models. Namely, we aim to establish conditions that ensure  $T_{ij}^\circ \geq A\lambda$  where  $A$  is the length of the support of the kernel  $K_s$  ( $K_s(u) = 0$  for  $u > A$ ). With this conditions, it holds  $K_s(T_{ij}/\lambda) = 0$  and hence the new weight  $w_{ij}$  will be equal to zero.

**Theorem 7.4.** *The condition*

$$\Delta := 0.5\sigma^{-2}(\bar{\theta}_i - \bar{\theta}_j)^\top B_i(B_i + B_j)^{-1}B_j(\bar{\theta}_i - \bar{\theta}_j) > A\lambda \quad (7.5)$$

implies with  $b = (\Delta - A\lambda)/\lambda$

$$P(T_{ij}^\circ < A\lambda) \leq e^{-\frac{b^2\lambda}{4(A+b)}}.$$

*Proof.* Similarly to the proof of Theorem 7.3, the decomposition (7.3) and the condition (7.4) imply

$$\mathbf{P}(T_{ij}^{\circ} < A\lambda) \leq \mathbf{P}(\Delta + \zeta_{ij} < A\lambda) \leq \mathbf{P}(-\zeta_{ij} > b\lambda) \leq e^{-b^2\lambda^2/(4\Delta)}.$$

□

### Proof of Theorem 6.3

The propagation condition (6.2) follows similarly to the proof of Theorem 6.2. The only difference is that in the local smooth case we apply Theorems 7.3 and 7.2 instead of Theorem 7.2. Let  $k$  be such that  $h^{(k)} \leq h$  and  $X_j \in U_h(X_i)$ . We apply Theorem 7.3 to the local models  $W_i^{(k)}$  and  $W_j^{(k)}$ . For this we have to check the condition (7.2) using Lemma 7.2. It holds with  $\theta = \theta(x)$  by the assumptions (A4) and (A6) that  $(f - \Psi^\top \theta)^\top W_j^{(k)}(f - \Psi^\top \theta) \leq N_j^{(k)} \delta^2 \leq N_j^* \delta^2$  for every  $X_j \in U_h(X_i)$ . Lemma 7.2 yields

$$(\bar{\theta}_i^{(k)} - \bar{\theta}_j^{(k)})^\top B_{ij}^{(k)} (\bar{\theta}_i^{(k)} - \bar{\theta}_j^{(k)}) \leq 2p\delta^2(N_i^* + N_j^*)$$

so that the condition (7.2) is fulfilled by (6.1).

Theorem 7.3 now applies yielding

$$\mathbf{P}\left(\min_{j=1,\dots,n} s_{ij}^{(k)} < 1/2\right) \leq n^{-1}$$

provided that  $\lambda = C \log n$  with a sufficiently large  $C$ .

The second assertion of the theorem follows from the next lemma.

**Lemma 7.3.** *Let the assumptions (A4), (A5) and (A6) hold true for some  $h$  and  $x = X_i$ . Let also the local model  $W_i$  be such that  $w_{ij} \geq 0.5\bar{w}_{ij} := K_l(l_{ij})$  for all  $j$ . If  $\lambda \geq C \log n$  for some fixed  $C$ , then*

$$\mathbf{P}\left(|\hat{f}_i - f_i| > \delta\sqrt{pC_B} + \sigma\sqrt{2C_B\lambda/N_i^*}\right) \leq 1/n.$$

*Proof.* Define  $W_i^* = \text{diag}\{w_{i1}^*, \dots, w_{in}^*\}$ ,  $B_i^* = \Psi W_i^* \Psi^\top$  and  $N_i^* = \text{tr} W_i^*$ . Then the conditions of the lemma yield  $N_i \geq 0.5N_i^*$  and  $B_i \geq 0.5B_i^*$ . Next, by Theorem 7.1

$$\mathbf{P}\left((\hat{\theta}_i - \bar{\theta}_i)^\top B_i(\hat{\theta}_i - \bar{\theta}_i) \geq \lambda\sigma^2\right) \leq 1/n$$

for  $\lambda \leq C \log n$  with a sufficiently large  $C$ . This implies by (A6) with a high probability

$$(\hat{\theta}_i - \bar{\theta}_i)^\top B_i^*(\hat{\theta}_i - \bar{\theta}_i) \leq 2\lambda\sigma^2$$

that, in its turn implies in view of (A6)

$$N_i^*(\widehat{\theta}_i - \bar{\theta}_i)^\top \Psi_i \Psi_i^\top (\widehat{\theta}_i - \bar{\theta}_i) \leq 2C_B \lambda \sigma^2$$

or equivalently

$$|\widehat{f}_i - \bar{f}_i| \leq \sigma \sqrt{2C_B \lambda / N_i^*}$$

where  $\bar{f}_i = \Psi_i^\top \bar{\theta}_i$ . Next, Lemma 7.2 and (A4) imply

$$(\bar{\theta}_i - \theta)^\top B_i^* (\bar{\theta}_i - \theta) \leq p \delta^2 N_i^*.$$

This and (A6) yield using the equality  $f_i = \Psi_i^\top \theta$

$$|\bar{f}_i - f_i|^2 = (\bar{\theta}_i - \theta)^\top \Psi_i \Psi_i^\top (\bar{\theta}_i - \theta) \leq C_B (\bar{\theta}_i - \theta)^\top B_i^* (\bar{\theta}_i - \theta) / N_i^* \leq p C_B \delta^2$$

and the assertion follows.  $\square$

## References

- [1] Cai, Z. Fan, J. and Li, R. (2000). Efficient estimation and inference for varying coefficients models. *J. Amer. Statist. Ass.* **95** 888–902.
- [2] Carroll, R.J., Ruppert, D. and Welsh, A.H. (1998). Nonparametric estimation via local estimating equation. *J. Amer. Statist. Ass.* **93** 214–227.
- [3] Fan, J. and Gijbels, I. (1996). *Local polynomial modelling and its applications*. Chapman & Hall, London.
- [4] Fan, J., Zhang, W. (1999). Statistical estimation in varying coefficient models. *Ann. Statist.* **27** 1491–1518.
- [5] Hastie, T.J. and Tibshirani, R.J. (1993). Varying-coefficient models (with discussion). *J. Royal Statist. Soc. Ser. B* **55** 757–796.
- [6] Heckman, N. and Ramsay, J.O. (2000). Penalized regression with model-based penalties. *Canadian Journal of Statistics* **28**, 241–258.
- [7] Lepski, O., Mammen, E. and Spokoiny, V. (1997). Ideal spatial adaptation to inhomogeneous smoothness: an approach based on kernel estimates with variable bandwidth selection. *Annals of Statistics*, **25**, no. 3, 929–947.

- [8] Polzehl, J. and Spokoiny, V. (2000). Adaptive weights smoothing with applications to image segmentation. *J. of Royal Stat. Soc.*, **62**, Series **B**, 335–354.
- [9] Polzehl, J. and Spokoiny, V. (2001). Functional and dynamic magnetic resonance imaging using vector adaptive weights smoothing. *Applied Statistics*, **50**, 485-501.
- [10] Polzehl, J. and Spokoiny, V. (2002). Local likelihood modeling by adaptive weights smoothing. Preprint 787. WIAS 2002. <http://www.wias-berlin.de/publications/preprints/787>.
- [11] Polzehl, J. and Spokoiny, V. (2003). Image denoising: pointwise adaptive approach. *Annals of Statistics*, **62**, in print.