

## Estimation of Piecewise-smooth Functions by Amalgamated Bridge Regression Splines

Felix Abramovich

*Tel Aviv University, Israel*

Anestis Antoniadis

*University Joseph Fourier, France*

Marianna Pensky

*University of Central Florida, USA*

---

### Abstract

We consider nonparametric estimation of a one-dimensional piecewise-smooth function observed with white Gaussian noise on an interval. We propose a two-step estimation procedure, where one first detects jump points by a wavelet-based procedure and then estimates the function on each smooth segment separately by bridge regression splines. We prove the asymptotic optimality (in the minimax sense) of the resulting amalgamated bridge regression spline estimator and demonstrate its efficiency on simulated and real data examples.

*AMS (2000) subject classification.* Primary 62G05, 62G08; secondary 65D07.  
*Keywords and phrases.* Amalgamation, bridge regression, jumps detection, nonparametric regression, penalized regression splines, wavelets.

---

### 1 Introduction

In a variety of nonparametric regression applications, the underlying response function is piecewise-smooth with abrupt changes between smooth segments. Examples include i) seismology, where the density of the sedimentary layers of the earth's crust can be locally approximated by a step function, ii) image processing, where discontinuities are present at the edges and iii) econometric models, where structural changes due to governmental policies are not rare. "Direct" methods for estimating piecewise-smooth functions in nonparametric regression include wavelets that are known to efficiently tackle local singularities. However, in practice wavelets often produce pseudo-Gibbs phenomena and other local artifacts in reconstructing

smooth regions (e.g. Coifman and Donoho, 1995; Antoniadis and Gijbels, 2002). Alternatively, following a two-step segmentation approach, one first detects the locations of change points and then applies some smooth non-parametric techniques on each segment separately (e.g. Oudshoorn, 1998; Antoniadis and Gijbels, 2002; Lee, 2002; Fink and Wells, 2004; Gill and Baron, 2004). Somewhat similar ideas are considered in Chu et al. (1998).

In this paper, we consider the latter approach and combine methods that are most suitable for each step. In particular, we present a wavelet-based method for detecting discontinuities (jumps) of a function and then introduce amalgamated penalized regression splines for estimating the function at smooth regions. The multi-resolutional nature of wavelet analysis makes it an excellent tool for detecting local singularities (Mallat and Hwang, 1992; Wang, 1995), while penalized regression splines are popular statistical techniques for recovering smooth functions from noisy data due to their various optimal properties, good practical performance and computational simplicity (Eilers and Marx, 1996; Eubank, 1999, Section 6).

The developed wavelet-based jumps detection procedure is somewhat similar in spirit to that of Wang (1995). Its error adds a negligible contribution to the overall quadratic risk of the resulting amalgamated regression spline estimator and allows one to obtain the same optimal convergence rates for the latter as for the case with known jumps. In addition, for the fixed knots the traditional  $l_2$ -penalty leads to a linear shrinkage (essentially ridge regression) estimator. In this paper, we consider a more general  $l_\rho$ -type penalty for  $\rho > 0$ . Such an approach has a direct analogy with the bridge regression of Frank and Friedman (1993) and we will call the resulting splines *bridge regression splines*. In particular,  $\rho = 1$  corresponds to the LASSO estimator of Tibshirani (1996). Generally,  $l_\rho$ -penalties for  $0 < \rho \leq 1$  lead to (nonlinear) spline estimators with fewer knots.

The proposed two-step procedure has some similarities with that of Lee (2002) and the recent adaptive multi-order penalized splines (AMPS) hybrid procedure of Fink and Wells (2004). Lee (2002) suggests to choose the number and placement of discontinuity points by several model selection criteria. However, he provides no theoretical results on the optimality of the resulting spline estimator. Fink and Wells (2004) estimate the locations of jumps of a piecewise-smooth function on the basis of the first differences of the data and then fit regression splines using a quadratic penalty. In fact, in terms of wavelet analysis, such jump detection corresponds to an application of the Haar wavelets at the finest resolution level. As a result, the AMPS

procedure is not powerful enough, does not attain the optimal rates and can detect only sufficiently sharp jumps.

In what follows, we propose a two-step amalgamated bridge regression spline (ABS) estimation procedure for piecewise-smooth functions and show its optimality (in the minimax sense) over amalgam Sobolev balls. We also demonstrate the good performance of ABS on several simulated and a real data examples. In particular, it produces smooth curves between estimated jumps and, unlike “direct” wavelet denoising, does not suffer from a pseudo-Gibbs phenomena.

The rest of the paper is organized as follows. We present the two-step ABS estimation procedure in Section 2 and establish its optimality in Section 3. Section 4 illustrates the performance of the ABS on several simulated and a real data examples. Some concluding remarks are made in Section 5. All the proofs are given in the Appendix.

## 2 Amalgamated Bridge Regression Spline Estimation Procedure

*2.1. The model.* Consider the standard nonparametric regression model with equidistant design

$$Y_i = f(x_i) + \sigma Z_i, \quad (2.1)$$

where  $f$  is an unknown response function,  $x_i = i/n$  and  $Z_i$  are i.i.d. standard normal random variables. Assume also that the sample size  $n = 2^J$  for some integer  $J > 0$ .

Assume that  $f$  is a piecewise-smooth function and belongs to the amalgam Sobolev ball  $\mathcal{H}(m, R, \kappa, S)$  of radius  $R$  of functions satisfying the following conditions :

- M1.  $f \in L_\infty([0, 1])$ .
- M2.  $f$  has  $D$  discontinuity (jump) points at locations  $0 < \theta_1 < \dots < \theta_D < 1$ , where the integer  $D$  and the real  $\theta_l$ 's are unknown and  $\theta_{l+1} - \theta_l > \kappa$ ,  $l = 1, \dots, D - 1$  for some  $\kappa > 0$ . In particular,  $D = 0$  corresponds to a continuous  $f$ .
- M3. At each discontinuity point  $\theta_l$ , the left and right limits  $f(\theta_l-)$  and  $f(\theta_l+)$  exist and  $|f(\theta_l+) - f(\theta_l-)| \geq S$  for some  $S > 0$ .

M4.

$$\sum_{l=0}^D \int_{\theta_l}^{\theta_{l+1}} [f^{(m)}(x)]^2 dx \leq R,$$

where integer  $m \geq 1$ ,  $\theta_0 = 0$  and  $\theta_{D+1} = 1$ .

Note that the condition M2 implies that the number of change points  $D$  is finite and bounded from above by  $1/\kappa < \infty$ .

The statistical challenges in estimating a piecewise-smooth function from  $\mathcal{H}(m, R, \kappa, S)$  are:

1. estimating the number of jumps  $D$  and their locations  $\theta_l$ ,  $l = 1, \dots, D$ ;
2. recovering the function at smooth regions without degrading its discontinuities.

In this regard, we propose the following two-step procedure: first, to detect the jump points by wavelet-based procedure and then to apply amalgamated bridge regression splines for estimating  $f$  between them. We start with presenting the estimation of a piecewise-smooth function with the *known* jump points by an amalgamated bridge regression spline and then provide a wavelet-based procedure for adaptive estimation of jump points from the data.

*2.2. Amalgamated bridge regression splines.* Assume that  $f \in H(m, R, \kappa, S)$ . Consider a standard polynomial regression spline estimator  $\tilde{f}$  of order  $m$  with the fixed knots  $\xi_1 < \dots < \xi_K$ . It is a continuous piecewise polynomial of degree  $m - 1$  with  $m - 2$  continuous derivatives at the knots and can be represented as

$$\tilde{f}(x) = \sum_{k=0}^{m-1} \beta_k x^k + \sum_{j=1}^K \beta_{m-1+j} (x - \xi_j)_+^{m-1}, \quad (2.2)$$

where  $z_+ = \max(0, z)$ . The unknown coefficients  $\beta$  are estimated from the data.

Polynomial splines are useful for approximating smooth functions but evidently inappropriate for fitting functions with abrupt local changes. There have been proposed various knot selection algorithms to adapt to inhomogeneous smoothness of the unknown response function by placing more knots where it shows rapid changes (e.g., Zhou and Shen, 2001; He, Shen and Shen, 2001). However, the resulting spline estimator is still a polynomial spline of the same order and therefore cannot provide a satisfactory remedy for

fitting piecewise-smooth functions with jumps. To model such sharp local features of a function efficiently, one can consider a more general and flexible *multi-order* regression spline of the form

$$\tilde{f}(x) = \sum_{k=0}^{m-1} \beta_k x^k + \sum_{j=1}^K \beta_{m-1+j} (x - \xi_j)_+^{m_j}, \quad (2.3)$$

where the smoothness  $0 \leq m_j \leq m - 1$  at different knots  $\xi_j$  may vary (Koo, 1997; Fink and Wells, 2004).

Multi-order splines allow jumps in the  $m_j$ -th derivative at  $\xi_j$ . In particular, zero-order knots ( $m_j = 0$ ) model discontinuities of the function while first and second order knots allow one to represent sharp changes in local linear trend and local curvature, respectively. Standard polynomial splines (2.2) of order  $m$  correspond to the particular case when  $m_j = m - 1$  for all  $j = 1, \dots, K$ . A piecewise-smooth function with  $D$  jumps  $\theta_1, \dots, \theta_D$  can be approximated by a multi-order spline with  $D$  zero-order knots at jump points  $\theta_l$ ,  $l = 1, \dots, D$  and a set of  $m - 1$ -order knots at smooth segments (Fink and Wells, 2004). However, as it follows from (2.3), such a multi-order spline necessarily implies the conditions on one-sided derivatives at jump points, namely,  $\tilde{f}^{(j)}(\theta_l^-) = \tilde{f}^{(j)}(\theta_l^+)$ ,  $j = 1, \dots, m - 1$ . Additional flexibility can be achieved if one considers *amalgamated polynomial regression splines* of order  $m$  with zero-order knots  $\theta_1, \dots, \theta_D$  obtained by amalgamation of separate  $m$ -order splines at each segment. An amalgamated polynomial regression spline  $\tilde{f}(x)$  of order  $m$  with  $D$  zero-order knots  $\theta_1, \dots, \theta_D$  and  $q$  knots  $\xi_1, \dots, \xi_q$  of order  $m - 1$  can be represented then as

$$\tilde{f}(x) = \tilde{f}_0(x)I_{\{0 \leq x < \theta_1\}} + \tilde{f}_1(x)I_{\{\theta_1 \leq x < \theta_2\}} + \dots + \tilde{f}_D(x)I_{\{\theta_D \leq x \leq 1\}}, \quad (2.4)$$

where each  $\tilde{f}_l$ ,  $l = 0, \dots, D$  is a polynomial regression spline of order  $m$  with  $q_l$  knots located at  $\xi_{1,l}, \dots, \xi_{q_l,l}$  and  $\sum_{l=0}^D q_l = q$ .

Re-number the observations and the  $m - 1$ -order knots  $\xi_1, \dots, \xi_{q_n}$  using the double indices  $(x_{i,l}, Y_{i,l})$  and  $\xi_{\nu,l}$ ,  $i = 1, \dots, n_l$ ,  $\nu = 1, \dots, q_{n_l}$ ,  $l = 0, \dots, D$ , respectively, where  $\theta_l \leq x_{i,l} < \theta_{l+1}$  and  $\theta_l \leq \xi_{\nu,l} < \theta_{l+1}$ . Using (2.2) and (2.4),  $\tilde{f}$  can be represented by:

$$\tilde{f}(x) = \sum_{l=0}^D \left[ \sum_{k=0}^{m-1} \beta_{k,l} x^k + \sum_{\nu=1}^{q_{n_l}} \beta_{m-1+\nu,l} (x - \xi_{\nu,l})_+^{m-1} \right] I(\theta_l \leq x < \theta_{l+1}), \quad (2.5)$$

where  $\sum_{l=0}^D q_{n_l} = q_n$ .

By the definition of an amalgamated spline, on each interval  $\theta_l \leq x < \theta_{l+1}$ ,  $\tilde{f}(x)$  is a usual  $m$ -order polynomial spline with  $q_{n_l}$  knots located at  $\xi_{1,l}, \dots, \xi_{q_{n_l},l}$ . Unless some prior information is available, the  $m - 1$ -order knots  $\xi_{1,l}, \dots, \xi_{q_{n_l},l}$  are usually placed on the sufficiently dense equidistant grid. An excessive number of  $m - 1$ -order knots might imply too much variability in the resulting spline estimator, so one needs some regularization procedure to remove superfluous  $\xi_{\nu,l}$  within each segment. We present an example of such a procedure below. The jump points are assumed meanwhile to be known, and, hence, using the representation (2.5), one can estimate the vector of unknown coefficients  $\beta_l$  on each  $l$ -th segment separately.

Let  $\mathbf{X}^{(l)}$  be the  $n \times (m + q_{n_l})$  matrix with the rows  $(1, x_{i,l}, \dots, x_{i,l}^{m-1}, (x_{i,l} - \xi_{1,l})_+^{m-1}, \dots, (x_{i,l} - \xi_{q_{n_l},l})_+^{m-1})$  and  $\mathbf{Y}_j$  be the vector with components  $Y_{i,l}$ ,  $i = 1, \dots, n_l$ . Consider the penalized maximum likelihood estimator of  $\beta_l$  with  $l_\rho$ -penalty,  $\rho > 0$ , derived by minimizing

$$Q_l(\beta_l, \mathbf{Y}_l) = \|\mathbf{Y}_l - \mathbf{X}^{(l)}\beta_l\|^2 + n_l \lambda_{n_l} \sum_{k=m}^{m-1+q_{n_l}} |\beta_{k,l}|^\rho \quad (2.6)$$

with respect to  $\beta_l$ , where  $\lambda_{n_l} > 0$  is a smoothing parameter.

The idea of  $l_\rho$ -penalty in regression was introduced by Frank and Friedman (1993) and the corresponding technique is known as *bridge regression* estimation. The traditional  $l_2$ -penalty yields a ridge regression estimator which is based on linear shrinkage, while  $\rho = 1$  leads to the LASSO estimator of Tibshirani (1996). Any choice  $0 < \rho \leq 1$  implies a thresholding estimator of  $\beta_l$  and, therefore, results in a spline with fewer  $m - 1$ -order knots (see Antoniadis and Fan, 2001). Plugging the coefficients  $\tilde{\beta}_{k,l}$  into (2.5) leads to the amalgamated bridge regression spline estimator  $\tilde{f}(x)$  of  $f(x)$ .

A closed form solution of (2.6) is available for  $\rho = 2$ . For  $\rho = 1$  the minimizer of (2.6) is unique and can be found either by a LASSO-type algorithm (Tibshirani, 1996; Osborne et al., 2000), once the matrices  $\mathbf{X}^{(l)}$  are normalized to have columns of norm 1, or via surrogate functionals, a method recently introduced by Daubechies et al. (2004) in the context of wavelet shrinkage methods for deblurring. When  $\rho < 1$ , the objective function is no longer convex but one can still find a local minimizer using, for example, an approximate algorithm of Ruppert and Carroll (2000), a backfitting type algorithm of Fu and Kneight (2000) or a recently developed algorithm of Amato et al. (2006). We discuss these issues in more details in Section 4.1 below.

To conclude this section, note that we have used truncated power bases for a clearer exposition of a spline-based regression. However, the truncated power bases may lead to numerical instabilities, especially when a larger number of knots and a small penalty parameter are involved. Equivalent bases with more stable numerical properties are the B-spline bases, and it is easy to transform the matrices  $\mathbf{X}^{(l)}$  to a corresponding B-spline version. For this reason, we shall not further discuss numerical stability issues when we formulate the ABS estimator.

*2.3. Jumps detection.* In this section, we present a wavelet-based procedure for adaptive estimation of jump points. As it has been mentioned in the Introduction, due to their multiresolutional nature, wavelets have become an efficient and widely used tool for detecting local abrupt changes (jumps, in particular) of a function. The general idea behind wavelet-based detection is based on the characterization of function's local regularity at a point by the rate of decay of its wavelet coefficients across scales around this point. Local singularities can be then identified by the presence of large wavelet coefficients at high scales in their neighbourhood.

Let  $\psi$  be a mother wavelet, and  $d_{jk}$  and  $\hat{d}_{jk}$ ,  $j = 0, \dots, J-1$ ,  $k = 0, \dots, 2^j - 1$  be the corresponding sets of discrete wavelet transform (DWT) coefficients of the unknown response vector  $\mathbf{f} = (f_1, \dots, f_n)'$  and data  $\mathbf{Y} = (y_1, \dots, y_n)'$  respectively. For equispaced design and the sample size  $n = 2^J$ , fast algorithms of Mallat (1989) allow one to perform the DWT in  $O(n)$  operations.

It is well-known, that the DWT coefficients cannot be generally used to detect a local singularity point, since the discrete grid on scales for the DWT might be too "crude" in view of the possible presence of other singularities or strong oscillations around this point (Mallat and Hwang, 1992). In this respect, the DWT coefficients differ from the coefficients of the continuous wavelet transform or its discrete analog – the non-decimated wavelet transform (NWT) (Shenza, 1992) that generates an equal number of  $n$  coefficients at each of  $J$  levels. However, we will show that for *piecewise-smooth* functions from amalgam Sobolev balls, the DWT can still be used to detect jump points.

The detection algorithm described below analyses the DWT coefficients at an appropriately chosen scale and selects a threshold large enough to prevent the coefficients corresponding to smooth segments to penetrate by, but still small enough to allow coefficients corresponding to singularities to

pass it through. The locations of jumps are then estimated by the locations of coefficients which exceed the threshold.

We will assume hereafter that the variance  $\sigma^2$  of the noise is known. Otherwise, for regression functions from amalgam Sobolev balls, it can be estimated at a parametric rate in the wavelet domain by the median of the absolute deviation of the empirical wavelet coefficients of the data at the highest resolution level divided by 0.6745.

Assume that the mother wavelet  $\psi$  has a compact support  $[L; U]$ ,  $L < 0 < U$  and denote  $\psi_{j^*k}(x) = 2^{j^*/2}\psi(2^{j^*}x - k)$ ,  $j = 0, \dots, J-1$ ,  $k = 0, \dots, 2^j - 1$ . Fix an arbitrarily small  $\delta > 0$ . Define  $j^*$  such that  $2^{j^*} = (U - L)n/(\ln n)^{1+\delta}$  and a sequence of indices  $\tau(k) = -L + (U - L)k$ ,  $k = 0, \dots, 2^{j^*}/(U - L) - 1$ . Without loss of generality, we may assume that  $j^*$  and  $2^{j^*}/(U - L)$  are integers; otherwise, we take the corresponding integer parts. Note that  $\Omega_{j^*k} = \text{supp } \psi_{j^*\tau(k)} = [2^{-j^*}(U - L)k; 2^{-j^*}(U - L)(k + 1)]$  and, therefore, the intersection  $\Omega_{j^*k} \cap \Omega_{j^*(k+1)}$  is a zero-measure set containing a single boundary point  $2^{-j^*}(U - L)(k + 1)$ . Hence, the unit interval is divided into a grid of  $N = 2^{j^*}/(U - L)$  non-overlapping intervals of lengths  $2^{-j^*}(U - L) = (\ln n)^{1+\delta}/n$ . Due to M2, for sufficiently large  $n$ , each of these intervals can contain only a single jump point.

Let  $T_{j^*}$  be a set of indices  $\tau(k)$  such that the corresponding interval  $\Omega_{j^*k}$  does not contain a jump point. For an arbitrary  $0 < \alpha < \delta/2$ , define a threshold

$$t_n^* = \sigma \sqrt{n^{-1}(\log n)^{1+\delta-2\alpha}}. \quad (2.7)$$

**PROPOSITION 2.1.** *Consider the model (2.1), where the unknown  $f \in \mathcal{H}(m, R, \kappa, S)$  defined in Section 2.1. Let the wavelet  $\psi$  be differentiable with a compact support and  $\hat{d}_{jk}$ ,  $j = 0, \dots, J-1$ ,  $k = 0, \dots, 2^j - 1$  be the set of the DWT coefficients of the data  $\mathbf{Y} = (Y_1, \dots, Y_n)'$ . Then, for the threshold  $t_n^*$  defined in (2.7) one has uniformly in  $f \in \mathcal{H}(m, R, \kappa, S)$ :*

1.  $\mathbb{P}(\max_{\tau(k) \in T_{j^*}} |\hat{d}_{j^*\tau(k)}| > t_n^*) = o(n^{-\gamma})$
2.  $\mathbb{P}(\min_{\tau(k) \notin T_{j^*}} |\hat{d}_{j^*\tau(k)}| < t_n^*) = o(n^{-\gamma})$

as  $n \rightarrow \infty$  for an arbitrarily large  $\gamma > 0$ .

Proposition 2.1 shows that we can track down the jumps by the presence of large DWT coefficients  $d_{j^*\tau(k)}$  with very high accuracy.

REMARK 2.1. The proposed jump detection procedure is performed at such a high resolution level  $j^*$ , that there are essentially no differences between the DWT and NWT coefficients. This explains why for a piecewise-smooth function  $f$  satisfying the conditions of Section 2.1 the DWT can still be used.

Based on the Proposition 2.1, we suggest the following jumps estimation procedure:

1. Consider the DWT coefficients  $\hat{d}_{j^*\tau(k)}$  at the level  $j^*$  and find all  $\tau(k)$  such that  $|\hat{d}_{j^*\tau(k)}| > t_n^*$ . If the set  $\{|\hat{d}_{j^*\tau(k)}| > t_n^*\}$  is empty, set  $\hat{D} = 0$ . Otherwise,
2. Estimate the number of jump points  $D$  by  $\hat{D} = \#\{|\hat{d}_{j^*\tau(k)}| > t_n^*\}$  and the locations  $\theta_\ell$  of the jumps by the mid-points  $\hat{\theta}_\ell$  of the corresponding intervals  $\Omega_{j^*k}$ , i.e.,  $\hat{\theta}_\ell = 2^{-j^*}(U - L)(k + 1/2)$ ,  $l = 1, \dots, \hat{D}$ .

Proposition 2.1 immediately implies that  $\mathbb{P}\{\hat{D} \neq D\} = o(n^{-\gamma})$  for an arbitrarily large  $\gamma > 0$ . Note also that

$$\mathbb{E}(|\hat{\theta}_l - \theta_l|^2 \mathbf{I}_{\hat{D}=D}) \leq 2^{-2j^*} + \mathbb{P}\{|\hat{\theta}_l - \theta_l| > 2^{-j^*}\}, \quad (2.8)$$

where the first term in the right-hand side of (2.8) is  $O((\ln n)^{2+2\delta}/n^2)$ , while the second one is negligible due to the first statement of Proposition 2.1. Hence, the following Corollary holds.

COROLLARY 2.1. *Under the assumptions of Proposition 2.1, as  $n \rightarrow \infty$ ,*

1.  $\mathbb{P}\{\hat{D} \neq D\} = o(n^{-\gamma})$  for any  $\gamma > 0$ .
2. *Uniformly in  $f \in \mathcal{H}(m, R, \kappa, S)$ ,*

$$\mathbb{E}\left(|\hat{\theta}_l - \theta_l|^2 \mathbf{I}_{\hat{D}=D}\right) = O(2^{-2j^*}) = O\left(n^{-2} (\ln n)^{2(1+\delta)}\right), \quad l = 1, \dots, \hat{D}.$$

2.4. *The ABS procedure.* The resulting two-step ABS procedure naturally combines amalgamated bridge regression spline estimation with jumps detection and can be summarized as follows.

1. Estimate the number of jump points  $\hat{D}$  and their locations  $\hat{\theta}_1, \dots, \hat{\theta}_{\hat{D}}$  by the DWT-based procedure described in Subsection 2.3.
2. Plug in  $\hat{D}$  and  $\hat{\theta}_1, \dots, \hat{\theta}_{\hat{D}}$  into (2.5) and minimize the resulting expression (2.6) to obtain an amalgamated bridge regression spline estimator  $\hat{f}$ .

### 3 Optimality of the ABS procedure

In this section, we prove the optimality (in the minimax sense) of the proposed ABS procedure.

Consider the quadratic risk ( $L_2$ -loss)  $R(\hat{f}, f) = \mathbb{E}\{\|\hat{f} - f\|_2^2\}$  for an estimator  $\hat{f}$  of  $f \in H(m, R, \kappa, S)$ : The minimax quadratic risk over  $\mathcal{H}(m, R, \kappa, S)$  is then defined by

$$R(\mathcal{H}(m, R, \kappa, S)) = \inf_{\hat{f}} \sup_{f \in \mathcal{H}(m, R, \kappa, S)} R(\hat{f}, f),$$

where the infimum is taken over all estimators  $\hat{f}$ . Antoniadis and Gijbels (2002) derived the minimax rate over  $R(\mathcal{H}(m, R, \kappa, S))$  and showed that as  $n$  increases,

$$R(\mathcal{H}(m, R, \kappa, S)) = \mathcal{O}\left(n^{-\frac{2m}{2m+1}}\right) \quad (3.1)$$

Note that the optimal rate (3.1) for estimating piecewise-smooth functions from amalgam Sobolev classes satisfying M1–M4 is the same as for homogeneously smooth functions from the usual Sobolev spaces.

We now show that the proposed ABS estimator attains the minimax rate (3.1). We first prove that for the *fixed* zero-order knots the amalgamated bridge regression spline estimator from Section 2.2 achieves the optimal rate  $\mathcal{O}(n^{-2m/(2m+1)})$  and then demonstrate that the accuracy of the zero-knots estimation procedure of Section 2.3 is sufficiently high not to damage it.

Consider the amalgamated bridge regression spline estimator  $\tilde{f}$  from Section 2.2 with  $q_n$  equally spaced  $m - 1$ -order knots and *fixed* zero-order knots. Impose the following asymptotic assumptions on the design matrix  $\mathbf{X}$ , the number of  $m - 1$ -order knots  $q_n$  and the smoothing parameters  $\lambda_{n_l}$  in (2.6) as  $n \rightarrow \infty$  and  $q_n \rightarrow \infty$ :

- M5. There exists  $C_1 > 0$  and  $C_2 > 0$  such that  $0 < C_1 n < \lambda_{\min} < \lambda_{\max} < C_2 n$ , where  $\lambda_{\min}$  and  $\lambda_{\max}$  are the minimal and the maximal eigenvalues of the matrix  $\mathbf{X}^T \mathbf{X}$ .
- M6.  $q_n = Cn^{1/(2m+1)}$  for some  $C > 0$ .
- M7.  $\lambda_{n_l} n_l^{1-\rho/2} = O(1)$  as  $n_l \rightarrow \infty$ ,  $l = 1, \dots, D$ .

Due to the assumption M2, the number of observations  $n_l$  and the number of  $m - 1$ -order knots  $q_{n_l}$  on each  $l$ -th segment are of the order  $n$  and  $q_n$

respectively, the same as on the entire unit interval, and, therefore, the assumptions M5 and M6 holds on each segment as well. Thus, when the zero-order knots are fixed, the asymptotic properties of the amalgamated estimator  $\tilde{f}$  on the entire unit interval are the same as on each of its segments.

Consider then minimization of (2.6) under the assumptions M5–M7. The following proposition guarantees the existence of a local  $\sqrt{n_l/q_{n_l}}$ -consistent penalized maximum likelihood estimator  $\tilde{\beta}_l$  of (2.6) of  $\beta_l$ .

**PROPOSITION 3.1.** *Under the assumptions M5 and M7, there exists a local minimizer  $\tilde{\beta}_l$  of (2.6) such that  $\|\tilde{\beta}_l - \beta_l\| = O_p(\sqrt{q_{n_l}/n_l}) = O_p(\sqrt{q_n/n})$ .*

Proposition 3.1 only establishes the existence of a local  $\sqrt{q_n/n}$  consistent minimizer of (2.6) but does not provide any tools to obtain it. We have discussed briefly the computational issues in Section 2.2 and will give more technical details later in Section 4.1.

The resulting ABS estimator  $\tilde{f}$  is obtained by amalgamation of the corresponding estimators at each segment, that is,  $\tilde{f} = X\tilde{\beta}$ . The following proposition shows that  $\tilde{f}$  achieves the optimal rates (3.1).

**PROPOSITION 3.2.** *Let the assumptions M5–M7 hold and  $\tilde{f} = X\tilde{\beta}$ . Then, as  $n \rightarrow \infty$ ,*

$$\sup_{f \in \mathcal{H}(m, R, \kappa, S)} R(\tilde{f}, f) = O\left(n^{-2m/(2m+1)}\right). \quad (3.2)$$

So far we considered an idealized situation, where the jump locations were assumed to be known. The following proposition shows that when zero-order knots are *estimated* by the wavelet-based jumps detection procedure from Section 2.3, the resulting ABS estimator  $\hat{f}$  still attains the optimal rates (3.1). The high accuracy of estimating jump points makes the additional error contribution to be negligible in the overall estimation error.

**PROPOSITION 3.3.** *Let the assumptions M1–M7 hold and let  $\hat{f}$  be the ABS estimator with zero-order knots estimated by the wavelet-based procedure proposed in Section 2.3. Then, as  $n \rightarrow \infty$ ,*

$$\sup_{f \in \mathcal{H}(m, R, \kappa, S)} R(\hat{f}, f) = O\left(n^{-2m/(2m+1)}\right). \quad (3.3)$$

**REMARK 3.1.** As we have already mentioned, due to the high accuracy of estimating jumps locations, the resulting optimal rates for the ABS estimator

$\hat{f}$  are the same as that of  $\tilde{f}$  with known jumps. Hence, for a finite number  $D$  of jumps these rates are essentially the same as those optimal within Sobolev balls with the known smoothness index  $m$ . Furthermore, using the corresponding results of Agarwal and Studden (1980) under the assumptions M5–M7, one can show that the constant in the upper bound of the risk in (3.3) is, in fact, the Pinsker constant (Pinsker, 1980)

$$C_P = (R(2m + 1))^{1/(2m+1)} \left( \frac{m}{\pi(m + 1)} \right)^{2m/(2m+1)}.$$

For simplicity of exposition, we considered the case of equidistant design in the model (2.1). As long as M5 holds, this assumption does not really matter on the spline step. Nevertheless, some technical modifications are needed for the proposed wavelet-based jumps detection procedure to adapt the discrete wavelet transform to a non-equidistant design. For  $n = 2^J$  one can use an interpolating wavelet basis (Donoho, 1992) to interpolate the observed values and map then the data to a set of dyadic equidistant points. For a sufficiently regular design density, all the required DWT properties remain valid. When  $n$  is not a power of two, one can still use the above approach by adapting the interpolating wavelet basis to the sampling grid using an appropriate subdivision interpolation scheme (Cohen et al., 2003).

Finally, note that, in fact, one can use any jumps detection method on the first step that achieves the detection rates established in Corollary 2.1.

## 4 Numerical Analysis

In this section, we discuss the computational/practical implementation of the proposed ABS procedure and demonstrate its performance on several simulated examples and a real data set.

*4.1. Computational implementation.* Following the proposed two-step ABS algorithm, we first estimate the number and locations of jumps of an unknown function by the wavelet-based detection procedure. We then place zero order knots on estimated jump points and a relatively large number of quadratic knots ( $m = 3$ ) at locations fixed at “equally-spaced sample quantiles” similarly to standard penalized splines designed for estimating smooth functions (e.g. Ruppert and Carroll, 2000). As it has been noted in the previous sections, this step serves to refine the regression spline basis by allowing for additional smoothness between any zero-order features. During the second step, the ABS regression spline is fitted using either penalized

least squares with a quadratic penalty or with a more general  $l_\rho$ -type penalty, penalizing coefficients of the basis functions which are the least supported by the data.

Recall that the bridge regression estimators  $\tilde{\beta}_j$  of coefficients  $\beta_j$  within each segment are obtained by minimizing  $Q_l$  in (2.6) separately for each  $l = 0, \dots, D$  with respect to  $\beta_j$ :

$$Q_l(\beta_l, \mathbf{Y}_l) = \|\mathbf{Y}_l - \mathbf{X}^{(l)}\beta_l\|^2 + n_l\lambda_{n_l} \sum_{k=m}^{m-1+q_{n_l}} |\beta_{k,l}|^\rho$$

Suppressing index  $l$  in the equation above for simplicity, we consider minimization of

$$Q(\beta, \mathbf{Y}) = \|\mathbf{Y} - \mathbf{X}\beta\|^2 + \lambda \sum_{k=m}^{m-1+q_{n_l}} |\beta_k|^\rho.$$

As we have mentioned above, a closed solution is available for  $\rho = 2$ , while for  $\rho = 1$  there exist numerical algorithms (e.g. LASSO) for computing the estimator. For any  $0 < \rho \leq 1$ , a possible numerical solution is to minimize  $Q$  iteratively, one component of  $\beta$  at a time (backfitting). Assume for simplicity that  $\bar{Y} = 0$  (or replace  $Y_i$  by  $Y_i - \bar{Y}$ ). The algorithm we have used in our numerical implementation for  $\rho \leq 1$  can be described then as follows:

- (0). Center the columns of  $\mathbf{X}$  to have the mean 0 and scale them to have unit variance. Using centered columns, define an initial value  $\hat{\beta}$  by using the least squares algorithm. Set  $k = 1$ .
- (1). Define  $Q_k(\beta_k) = \sum_{i=1}^n (Y_i - \sum_{j \neq k} \hat{\beta}_j x_{ij} - \beta_k x_{ik})^2 + \lambda |\beta_k|^\rho$ .
- (2). Set  $\hat{\beta}_k = \arg \min Q_k$ . The minimization of  $Q_k$  with respect to  $\beta_k$  is solved by Newton-Raphson or fixed-point iteration.
- (3). Set  $k = k + 1$  if  $k < m - 1$ , and  $k = 1$  otherwise.
- (4). Repeat (1), (2) and (3) until convergence occurs.

The above algorithm always converges (see Nikolova and Ng, 2005) and works very well if the design is not “too collinear” (hence the interest in using B-splines). Otherwise, it might get stuck at a local minima. The problem is less severe when  $\rho$  is not too close to 0. For  $\rho = 1$ , it may be also computationally simpler than LASSO that involves linear programming techniques.

In the simulation and real data examples below, we considered  $\rho = 2$  (ABS2) and  $\rho = 1$  (ABS1), where in the latter, we used the backfitting algorithm described above. The quadratic  $l_2$ -penalty (ABS2) is equivalent to placing quadratic penalties on finite differences of adjacent B-splines coefficients and it results in shrinking all coefficients toward zero. On the contrary, the  $l_1$ -penalty on adjacent B-splines coefficients (ABS1) not only shrinks the coefficients but also thresholds them, removing, in the process, the corresponding “superfluous” second order knots.

In both approaches, the smoothing parameter  $\lambda_n$  was automatically chosen from the data by generalized cross-validation (GCV) as is usual in spline smoothing (see, e.g. Fan and Li, 2001). For all examples, we found the reasonable choice for the tuning parameters of the jumps detection procedure to be  $\delta = 1/2$  and  $\alpha = 1/4$ .

*4.2. Simulations.* In this subsection, we compare the estimators based on the ABS1 and ABS2 procedures with another related method, namely, the Spatially Adaptive Regression Splines (SARS) developed by Zhou and Shen (2001) which is particularly suited for functions that have jumps by themselves or in their derivatives. SARS is locally adaptive to variable smoothness and automatically places more knots in the regions where the function is not smooth. It has been proved as an effective tool for estimating such functions. For completeness, we also compare the above estimators with a standard wavelet denoising procedure based on universal thresholding of Donoho and Johnstone (1994), since wavelet based procedures are known to efficiently denoise inhomogeneous functions.

To investigate the performance of the developed ABS estimators we conducted a simulations study based on synthetic data. We used two of the standard test functions of Donoho and Johnstone (1994) that are examples of piecewise-smooth functions and commonly used for various wavelet procedures, namely, the `blocks` and `heavisine`. In addition, we considered two other test functions punctuated by jump discontinuities called `burt` and `cosine`, defined on  $[0, 1]$ , respectively, as

$$\begin{aligned} \text{cosine}(x) &= \cos(5.5\pi x) - 4 \text{sign}(0.23 - x) - 2 \text{sign}(0.3 - x) \\ &\quad - 1.75 \text{sign}(0.55 - x) + 3 \text{sign}(0.7 - x) \end{aligned}$$

and

$$\text{burt}(x) = 20x \cos(16x^{1.2}) - 20I(x < 0.5).$$

The functions are depicted in Figure 4.1. The sample size used in the simulations was  $n = 256$  and the design points were uniformly spaced within the

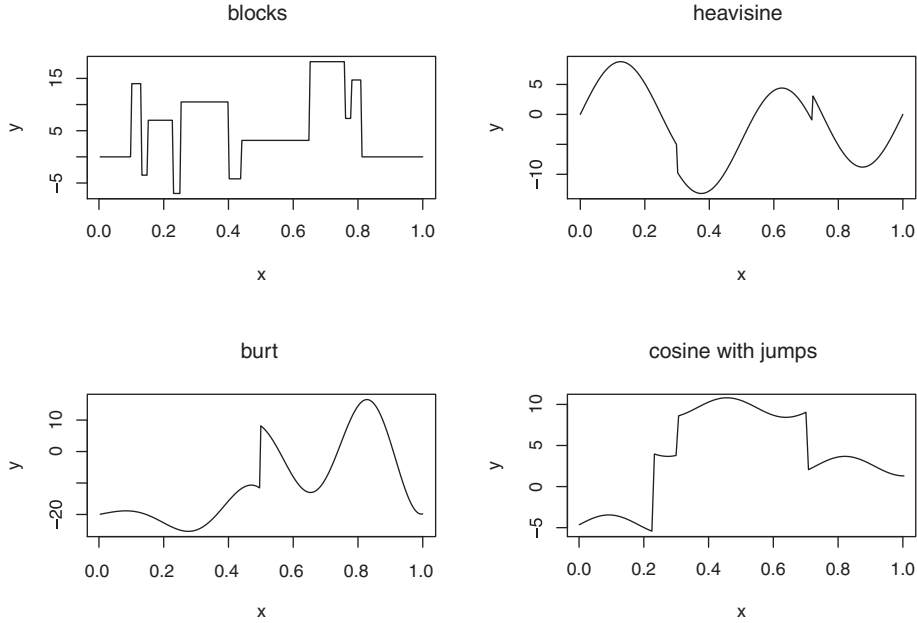


Figure 4.1: Four benchmark test functions punctuated by jump discontinuities used in the simulations.

unit interval. In each simulation, we added a normally distributed zero mean white noise with the standard deviation  $\sigma$  implying a chosen signal-to-noise-ratio (SNR). SNR are measured as  $sd(f)/\sigma$ , where  $sd(f)$  is the estimated standard deviation of the regression function over the grid. For each function and two values of SNR (4 and 6), we ran 200 simulations. In each realization, we placed equally-spaced 50 quadratic knots. The noise level  $\sigma$  was assumed unknown and estimated by the median of the absolute deviation of the empirical wavelet coefficients of the data at the highest resolution level divided by 0.6745. For each realization, we calculated the ABS1 and ABS2 estimators using the developed procedures, where the wavelet-based jumps detection was based on Symmlets of order 6, the SARS estimator and a wavelet denoising estimator (Wav) also with Symmlets of order 6. All simulations were performed using Matlab 7 (Mathworks 2001) and R.

The accuracy of an estimator  $\hat{f}$  of  $f$  was measured by the average mean square error (MSE) averaged over 200 simulation runs defined as

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (\hat{f}(x_i) - f(x_i))^2,$$

where  $\{x_i\}$  is the set of design points. The average MSEs for each method are reported in Table 1 with corresponding boxplots in Figure 4.2.

Figure 4.2 and Table 4.1 show that all the three spline estimators with adaptively placed knots outperform a standard wavelet denoising procedure and, unlike the latter, do not suffer from pseudo-Gibbs phenomena. ABS1 and ABS2 provide similar results and in most of the cases, they are somewhat better than SARS. The unusually large MSE for ABS2 for the blocks example, SNR=6 is explained by an extremely high MSE obtained in one simulation run (cf. Figure 4.2). As one can also see, application of the  $l_1$ -type penalty in ABS provides a relatively small gain. This can be explained by the fact that, for each of the functions, the regions between any two jumps are relatively smooth.

TABLE 4.1. AVERAGE MEAN SQUARED ERRORS BASED ON 200 SAMPLES OBTAINED USING SIGNAL TO NOISE RATIOS OF 6 AND 4 (IN PARENTHESIS) FOR FOUR DIFFERENT PROCEDURES: ABS1 ( $l_1$ -PENALTY); ABS2 ( $l_2$  PENALTY), SARS AND WAV. THE HYPERPARAMETERS ARE CHOSEN BY GENERALIZED CROSS-VALIDATION.

estimate	blocks	heavisine	burt	cosine
ABS1	0.47 (1.52)	0.21 (0.35)	0.43 (0.87)	0.25 (0.33)
ABS2	1.94 (1.52)	0.22 (0.37)	0.46 (0.99)	0.25 (0.33)
SARS	0.62 (1.45)	0.29 (0.48)	0.49 (1.11)	0.25 (0.34)
Wav	6.09 (8.91)	0.52 (0.61)	2.24 (3.42)	0.64 (0.86)

*4.3. A real example.* Since the ABS1 has shown the better performance in our simulations we applied the ABS1 procedure to the real life data provided by Mike Battaglia from the Department of Forestry at Virginia Tech (see Battaglia, 2000). For comparison, we have also applied the SARS procedure.

The data set contains relative light transmittance data recorded at equal time intervals throughout the daylight hours for numerous days (see Figure 4.3 for a plot of the relative light transmittance data for one station during one day). In this data set, sun light from various forest stations in plots with different cutting treatments is compared to the sun light in a nearby open plot. Cloud interference and overstory patterns (the shades produced by the trees) are the two most common phenomena that cause jump

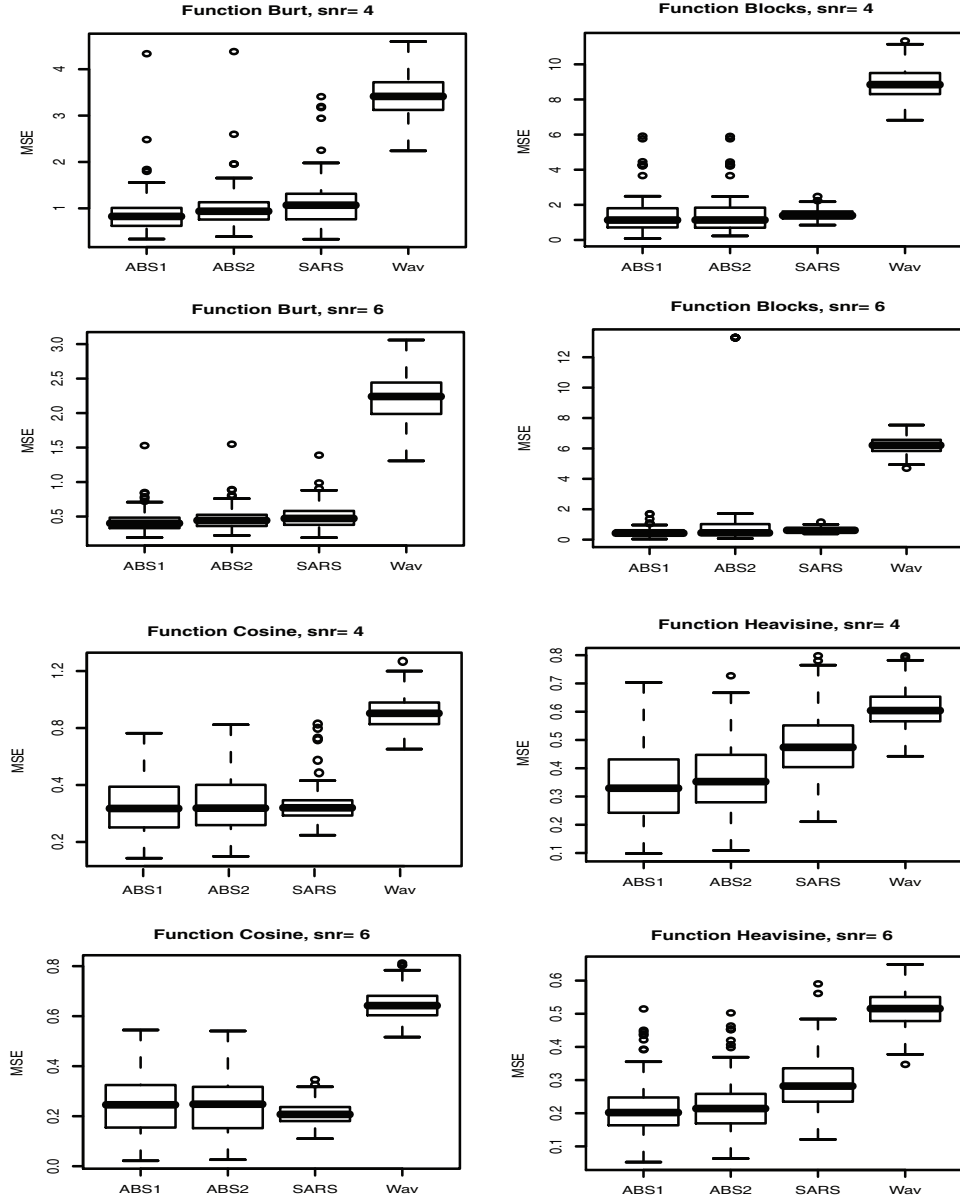


Figure 4.2: Boxplots of the mean squared errors on 100 samples obtained by four different procedures: ABS1 ( $l_1$ -penalty); ABS2 ( $l_2$  penalty); SARS, Spatially Adaptive Regression Splines; Wav, Wavelet denoising. The hyperparameters are chosen by generalized cross-validation.

points in the relative transmittance data. Jump points that remain consistent across days may be attributed to overstory pattern while jump points that do not remain consistent across days are probably caused by cloud interference. Since variation of the light availability increases as canopy gaps become larger, in order to predict the forest dynamics it is useful to estimate the relative transmittance taking into account such discontinuities.

We applied the ABS1 procedure to the data shown at Figure 4.3. Symmlets of order 6 were again used for jumps detection. The noise level  $\sigma$ , as usual, was estimated by the median of the absolute deviation of the empirical wavelet coefficients of the data at the highest resolution level divided by 0.6745. The procedure found three jump points located at 0.138, 0.494 and 0.873. On the spline step on each of the resulting four segments 12 quadratic knots were placed at equally spaced quantiles. The corresponding GCV-chosen smoothing parameters  $\lambda$  were 1.180e-06, 1.635e-05, 2.993e-06 and 4.816e-06. The ABS1 reduced the numbers of necessary quadratic knots to 3, 4, 3 and 3 respectively. The Figure 4.3 shows the resulting ABS1 and SARS estimates. ABS1 nicely fits the data, while SARS faces difficulties in detecting relatively small jumps and evidently misses a jump point around 0.14. The overall residual sums of squares were 0.003213 for ABS1 and 0.003826 for SARS. Summarizing, we can conclude that the developed ABS procedure successfully estimates trends of the daily relative light transmittance data and their jumps and can be therefore a powerful tool for the characterization of gap openings in forest stations.

## 5 Concluding Remarks

In this paper, we developed a two-step procedure for estimating piecewise-smooth functions by amalgamated bridge regression splines. It first detects the unknown jump points by a wavelet-based method and then estimates the regression function on each smooth segment separately by bridge regression splines. We showed that the resulting amalgamated estimator achieves min-max convergence rates over amalgam Sobolev balls. From a practical view, ABS is fairly accurate and computationally efficient: it requires a little more computation time than traditional penalized regression splines. The number of operations necessary to locate the zero-order knots grows linearly with  $n$ , inheriting the algorithmic complexity of the discrete wavelet transform. We demonstrated a good performance of ABS on several simulated and a real data examples. Summarizing, we believe that the proposed ABS estimators

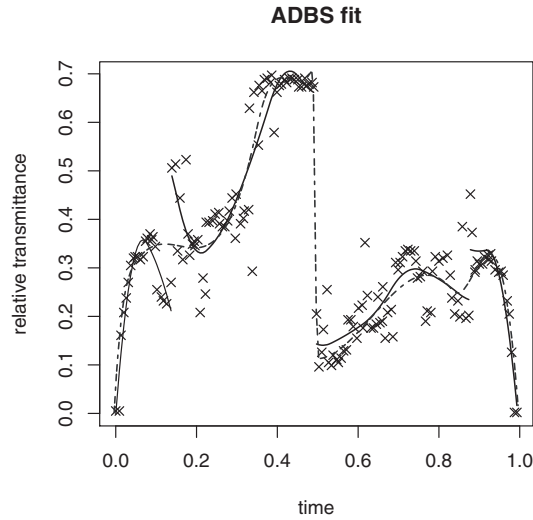


Figure 4.3: Display of the relative light transmittance data (data plotted as stars) for one station during one day from the Forestry Department of Virginia Tech. The x-axis is a daylight hours time interval rescaled to  $[0, 1]$ . The display also includes the ABS1 regression fit (solid line) to the transmittance data (with 3 jumps detected) and the fit obtained by the SARS procedure (dashed line).

are an attractive alternative to the existing estimators of piecewise-smooth functions.

## 6 Appendix

Throughout the proofs, we use the symbol  $C$  to denote a generic positive constant, not necessarily the same each time it is used, even within a single equation.

*A.1. Proof of Proposition 2.1.* Applying the DWT to the both sides of (2.1) we obtain  $\hat{d}_{jk} = d_{jk} + \epsilon_{jk}$ , where  $d_{jk}$  and  $\epsilon_{jk}$ ,  $j = 0, \dots, J-1$ ,  $k = 0, \dots, 2^j - 1$  are the DWT of the unknown  $f$  and the Gaussian noise respectively. Note that  $\epsilon_{jk}$  are independent Gaussian variables with  $\text{Var}(\epsilon_{jk}) = \sigma^2/n$ .

Consider the level  $j^*$  and the corresponding sequence of indices  $\tau(k)$  defined in Section 2.3. Consider first  $\tau(k) \in T_{j^*}$ . Since  $m \geq 1$ , the function  $f$  is at least of Lipschitz regularity one for all  $x \in \Omega_{j^*k}$ . Then, for sufficiently large  $n$ ,

$$\max_{k \in T_{j^*}} |d_{j^*\tau(k)}| < C2^{-\frac{3}{2}j^*} = O\left(\left(n^{-1} (\ln n)^{1+\delta}\right)^{3/2}\right) = o(t_n^*) \quad (\text{A.1})$$

(e.g. Daubechies, 1992, p.299). As  $n$  increases, for any  $0 < \alpha < \delta/2$  and for all  $\tau(k) \in T_{j^*}$

$$\begin{aligned} \mathbb{P}\{|\hat{d}_{j^*\tau(k)}| > t_n^*\} &\leq \mathbb{P}\{|\epsilon_{j^*\tau(k)}| > t_n^* - |d_{j^*\tau(k)}|\} \\ &\leq \mathbb{P}\{|\epsilon_{j^*\tau(k)}| > t_n^*/2\} \\ &\leq C(\log n)^{\alpha-(1+\delta)/2} \exp(-(\log n)^{1+\delta-2\alpha}/8) \\ &= o(n^{-\tilde{\gamma}}) \end{aligned} \quad (\text{A.2})$$

for any  $\tilde{\gamma} > 1$ . Since  $\text{card}\{T_{j^*}\} = O(2^{j^*}) = O(n/(\ln n)^{1+\delta})$ , we have

$$\mathbb{P}\{\max_{k \in T_{j^*}} |d_{j^*\tau(k)}| > t_n^*\} \leq \sum_{k \in T_{j^*}} \mathbb{P}\{|\hat{d}_{j^*k}| > t_n^*\} = o(n^{-\gamma}),$$

where  $\gamma = \tilde{\gamma} - 1 > 0$ .

Let now  $\tau(k) \notin T_{j^*}$ . The condition M2 guarantees that for sufficiently large  $n$  there is a single jump point  $\theta_l \in \Omega_{j^*k}$ . Similarly to the arguments of Wang (1995) and Antoniadis and Gijbels (2002) but applied for the DWT, under the conditions M1–M4 we derive that

$$\begin{aligned} |d_{j^*k}| &= 2^{-j^*/2} |f(\theta_l+) - f(\theta_l-)| \left\{ \left| \int \psi(\theta_l - u) \text{sign}(u) du \right| + O(1) \right\} \\ &\geq C2^{-j^*/2}. \end{aligned} \quad (\text{A.3})$$

For each  $\tau(k) \notin T_{j^*}$ , it then follows from (A.3) that

$$\begin{aligned} \mathbb{P}\{|\hat{d}_{j^*\tau(k)}| < t_n^*\} &\leq \mathbb{P}\{|\epsilon_{j^*\tau(k)}| > C2^{-j^*/2}\} \\ &< C(\log n)^{-(1+\delta)} \exp(-C^2(\log n)^{1+\delta}/2) \\ &= o(n^{-\gamma}) \end{aligned}$$

for any  $\gamma > 0$ . Hence,  $\mathbb{P}\{\min_{\tau(k) \notin T_{j^*}} |\hat{d}_{j^*\tau(k)}| < t_n^*\} \leq \sum_{k=1}^D \mathbb{P}\{|\hat{d}_{j^*\tau(k)}| < t_n^*\} = o(n^{-\gamma})$  whenever  $D$  is finite, that completes the proof.  $\square$

A.2. *Proof of Proposition 3.1.* As we have mentioned, the asymptotic properties of  $\tilde{\beta}_l$  are same for each  $l$  and for simplicity of exposition we omit the index  $l$  throughout the proof.

Let  $\beta$  be the true set of coefficients. Let  $M$  be a relatively large number and  $\alpha_n = \sqrt{q_n/n}$ . In order to prove Proposition 3.1 we show that

$$\lim_{n \rightarrow \infty} P \left\{ \sup_{\|\mathbf{u}\|=M} Q(\beta + \alpha_n \mathbf{u}, \mathbf{Y}) > Q(\beta, \mathbf{Y}) \right\} = 0. \quad (\text{A.4})$$

Equality (A.4) implies that with probability tending to one, there is a local minimum  $\tilde{\beta}$  of (2.6) in the ball with the center  $\beta$  and radius  $\alpha_n M$  such that  $\|\tilde{\beta} - \beta\| = O_p(\alpha_n)$ . Let  $\boldsymbol{\epsilon} = \mathbf{Y} - \mathbf{X}\beta$  and

$$\begin{aligned} D_n(\mathbf{u}) &= Q(\beta + \alpha_n \mathbf{u}, \mathbf{Y}) - Q(\beta, \mathbf{Y}) = \Delta_1 + \Delta_2 + \Delta_3 \\ &= \alpha_n^2 \|\mathbf{X}\mathbf{u}\|^2 - 2\alpha_n \mathbf{u}^T \mathbf{X}^T \boldsymbol{\epsilon} + n\lambda_n \sum_{k=1}^{q_n} [|\beta_{p+k} + \alpha_n u_k|^\rho - |\beta_{p+k}|^\rho] \end{aligned} \quad (\text{A.5})$$

Note that  $\Delta_1$  is a constant and, by the assumption M5,  $\Delta_1 \geq C_1 \alpha_n^2 n \|\mathbf{u}\|^2 = C_1 M^2 \alpha_n^2 n = C_1 M^2 q_n$ . The second term,  $\Delta_2$  is a random variable with  $E\Delta_2 = 0$ , so that, for any  $\delta > 0$  by Markov inequality, and M5 we have

$$P(|\Delta_2| > \delta) \leq 2\delta^{-1} \alpha_n \sqrt{E(\mathbf{u}^T \mathbf{X}^T \boldsymbol{\epsilon} \boldsymbol{\epsilon}^T \mathbf{X} \mathbf{u})} \leq 2\delta^{-1} M \sigma \alpha_n \sqrt{C_2 n}.$$

Therefore, setting  $\delta = 0.5 C_1 M^2 \alpha_n^2 n$ , we obtain  $P(|\Delta_2| > 0.5 C_1 M^2 \alpha_n^2 n) \leq 4\sigma(C_1 M)^{-1} (\alpha_n \sqrt{n})^{-1} \rightarrow 0$ . For an upper bound for  $\Delta_3$ , note that since  $|x + y|^\rho - |x|^\rho \leq |y|^\rho$  as  $0 < \rho \leq 1$ , by the condition M7 when  $M$  is large enough we have

$$|\Delta_3| \leq n\lambda_n \alpha_n^\rho \sum_{k=1}^{q_n} |u_k|^\rho \leq n\lambda_n \alpha_n^\rho M^{\rho/2} q_n^{1-\rho/2} < 0.5 C_1 M^2 q_n. \quad \square$$

A.3. *Proof of Proposition 3.2.* The proof Proposition 3.2 is based on the following lemma.

LEMMA 1.1. *Let  $w(z, \alpha) = \arg \min_w [w^2 - 2wz + \alpha|w|^\rho]$ ,  $0 < \rho \leq 1$ . Then,  $w(z, \alpha) = 0$  whenever  $|z| < a_\rho \alpha^{1/(2-\rho)}$ . If  $\rho \neq 1$ , then  $|w(z, \alpha)| \geq b_\rho \alpha^{1/(2-\rho)}$  whenever  $|z| \geq a_\rho \alpha^{1/(2-\rho)}$ . Here  $a_\rho = (2 - \rho)(\rho/2)^{1/(2-\rho)} (1 - \rho)^{(\rho-1)/(2-\rho)}$  and  $b_\rho = [\rho(1 - \rho)/2]^{1/(2-\rho)}$ .*

PROOF OF LEMMA 1.1. Since  $w(-z, \alpha) = -w(z, \alpha)$ , it is enough to consider  $z > 0$ . The derivative of the objective function is of the form  $h(w, z) = 2w - 2z + \alpha\rho|w|^{\rho-1}\text{sgn}(w)$ . Note that the function  $\phi(w) = 2w + \alpha\rho|w|^{\rho-1}\text{sgn}(w)$  is odd and for  $w > 0$  has a minimum at  $w_0 = b_\rho\alpha^{1/(2-\rho)}$  equal to  $a_\rho\alpha^{1/(2-\rho)}$ . Hence, whenever  $0 < z < a_\rho\alpha^{1/(2-\rho)}$  equation  $h(w, z) = 0$  has no solutions and  $h(w, z)$  is negative for any  $w < 0$  and positive for any  $w > 0$ . Thus, in this case  $w(z, \alpha) = 0$ . If  $z > a_\rho\alpha^{1/(2-\rho)}$ , equation  $h(w, z) = 0$  has two solutions  $0 < w_1 < w_2$  where solution  $w_1 < w_0$  corresponds to the local maximum of the objective function while  $w_2 > w_0$  corresponds to its absolute minimum.  $\square$

We now complete the proof of Proposition 3.2. Let  $\beta_j$  be the  $j$ -th component of  $\beta$  and  $\beta^{(-j)}$  be the vector  $\beta$  without component  $\beta_j$ . Similarly, let  $\mathbf{X}_j$  be the  $j$ -th column of matrix  $\mathbf{X}$  and  $\mathbf{X}^{(-j)}$  be the matrix  $\mathbf{X}$  without column  $j$ . Note that  $n^{-1}Q(\beta, \mathbf{Y})$  in (2.6) can be rewritten as  $n^{-1}Q(\beta_j, \beta^{(-j)}, \mathbf{Y}) = n^{-1}\|\mathbf{Y} - \mathbf{X}_j\beta_j - \mathbf{X}^{(-j)}\beta^{(-j)}\|^2 + \lambda_n \sum_{k=m}^{m-1+q_n} |\beta_k|^\rho$  and  $\tilde{\beta}_j = \arg\min_{\beta_j} Q(\beta_j, \beta^{(-j)}, \mathbf{Y})$ . If  $0 \leq j \leq m-1$ , equating derivative of  $n^{-1}Q(\beta, \mathbf{Y})$  over  $\beta_j$  to zero we obtain

$$n^{-1}\mathbf{X}_j^T\mathbf{X}\tilde{\beta} - n^{-1}\mathbf{X}_j^T\mathbf{Y} = 0. \quad (\text{A.6})$$

If  $m+1 \leq j \leq m+q_n$ , then application of Lemma 1.1 with  $\alpha = \lambda/(n^{-1}\mathbf{X}_j^T\mathbf{X}_j)$  and  $z = n^{-1}\mathbf{X}_j^T(\mathbf{Y} - \mathbf{X}^{(-j)}\tilde{\beta}^{(-j)})/(n^{-1}\mathbf{X}_j^T\mathbf{X}_j)$  yields that  $\tilde{\beta}_j = 0$  whenever

$$|n^{-1}\mathbf{X}_j^T(\mathbf{Y} - \mathbf{X}^{(-j)}\tilde{\beta}^{(-j)})| < a_\rho\lambda^{1/(2-\rho)}(n^{-1}\mathbf{X}_j^T\mathbf{X}_j)^{(1-\rho)/(2-\rho)}. \quad (\text{A.7})$$

In this case, taking into account the condition M5,  $\tilde{\beta}_j = 0$  and  $\mathbf{X}\tilde{\beta} = \mathbf{X}^{(-j)}\tilde{\beta}^{(-j)} + \mathbf{X}_j\tilde{\beta}_j$ , we obtain

$$|n^{-1}\mathbf{X}_j^T(\mathbf{Y} - \mathbf{X}\tilde{\beta})| = O\left(\lambda^{1/(2-\rho)}\right). \quad (\text{A.8})$$

If inequality (A.7) does not hold, then  $\hat{\beta}_j$  is the solution of the equation  $n^{-1}\mathbf{X}_j^T(\mathbf{X}\hat{\beta} - \mathbf{Y}) + \lambda\rho|\hat{\beta}_j|^{\rho-1}\text{sgn}(\hat{\beta}_j) = 0$ . For  $0 < \rho < 1$ , by Lemma 1.1 we obtain

$$|n^{-1}\mathbf{X}_j^T(\mathbf{Y} - \mathbf{X}\tilde{\beta})| = \lambda\rho|\tilde{\beta}_j|^{\rho-1} = O\left(\lambda^{1+(\rho-1)/(2-\rho)}\right) = O\left(\lambda^{1/(2-\rho)}\right). \quad (\text{A.9})$$

Observe that for  $\rho = 1$ , (A.9) continues to hold. Combining (A.6), (A.8) and (A.9), we derive that for any  $j$

$$n^{-1}|\mathbf{X}_j^T \mathbf{X} \tilde{\boldsymbol{\beta}} - n^{-1} \mathbf{X}_j^T \mathbf{Y}| = O\left(\lambda^{1/(2-\rho)}\right). \quad (\text{A.10})$$

Now, denote by  $\tilde{\boldsymbol{\beta}}_*$  the global minimizer of  $\|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2$ . It is easy to see that

$$n^{-1} \mathbf{X}_j^T \mathbf{X} \tilde{\boldsymbol{\beta}}_* - n^{-1} \mathbf{X}_j^T \mathbf{Y} = 0. \quad (\text{A.11})$$

Moreover, Agarwal and Studden (1980) showed that for  $\tilde{f}_* = \mathbf{X} \tilde{\boldsymbol{\beta}}_*$  under condition M7, one has

$$\sup_{f \in \mathcal{H}(m, R, \kappa, S)} R(\tilde{f}_*, f) = O\left(n^{-2m/(2m+1)}\right). \quad (\text{A.12})$$

Subtracting equation (A.11) from equation (A.10) one derives  $n^{-1} |(\mathbf{X}^T \mathbf{X}(\tilde{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}}_*))_j| = O(\lambda^{1/(2-\rho)})$  for all  $j = 0, \dots, m-1+q_n$ , so that (A.10) together with condition M7 imply that  $\|\mathbf{X} \tilde{\boldsymbol{\beta}} - \mathbf{X} \tilde{\boldsymbol{\beta}}_*\|^2 = O(q_n \lambda^{2/(2-\rho)})$ . To complete the proof combine the last equality with (A.12) and the conditions M6 and M7.  $\square$

*A.4. Proof of Proposition 3.3.* Let  $\hat{T} = \{\hat{\theta}_j, j = 1, \dots, \hat{D}\}$  be the set of estimated jump points of  $f$  and define the event  $F = \{\hat{D} = D\} \cap \{|\hat{\theta}_j - \theta_j| < n/(\log n)^{1+\delta}, j = 1, \dots, \hat{D}\}$ . We have

$$R(\hat{f}, f) = \mathbb{E}\{\|\hat{f} - f\|_2^2\} = \mathbb{E}\{\|\hat{f} - f\|_2^2 \mathbf{I}_F\} + \mathbb{E}\{\|\hat{f} - f\|_2^2 \mathbf{I}_{F^c}\}. \quad (\text{A.13})$$

Note that

$$\mathbb{E}\{\|\hat{f} - f\|_2^2 \mathbf{I}_{F^c}\} = \mathbb{E}\left(\mathbb{E}\{\|\hat{f} - f\|_2^2 \mathbf{I}_{F^c} | \hat{T}\}\right) = \mathbb{E}\left(\mathbf{I}_{F^c} \mathbb{E}\{\|\hat{f} - f\|_2^2 | \hat{T}\}\right).$$

Exploiting the results of Section 3 and the fact that both  $\hat{f}$  and  $f$  belong to amalgam Sobolev ball  $\mathcal{H}(m, R, \kappa, S)$ , one easily gets  $\mathbb{E}\{\|\hat{f} - f\|_2^2 | \hat{T}\} = O_P(1)$ . Furthermore, by Proposition 2.1,  $\mathbb{P}\{F^c\} = o(n^{-\gamma})$  for an arbitrarily large  $\gamma > 0$ . Hence the second term in the right-hand side of (A.13) is  $o(n^{-\gamma})$  and is negligible.

Consider now the first term in the right-hand side of (A.13). For simplicity of exposition, consider first the case where there is a single jump point  $\theta$  and  $D = 1$ . The unknown  $f$  can be decomposed as  $f(x) = \mathbf{I}_{[0, \theta]}(x) f_0(x) + \mathbf{I}_{(\theta, 1]}(x) f_1(x)$ , where  $f_0$  and  $f_1$  both belong to the Sobolev balls  $\mathcal{H}(m, R)$  of radius  $R$ . From the known properties of spline approximation, we can

approximate  $f$  by an amalgamated  $m$ -order polynomial spline  $s$  as  $s(x) = \mathbf{I}_{[0,\theta]}(x)s_0(x) + \mathbf{I}_{(\theta,1]}(x)s_1(x)$ , where  $s_0$  and  $s_1$  are spline approximations of  $f_0$  and  $f_1$  respectively, with the approximation error  $\|f - s\|_2^2 \leq O(n^{-2m/(2m+1)})$ . Let  $\hat{\theta}$  the estimate of  $\theta$ . Given the data, define amalgamated spline estimators  $\tilde{f}$  and  $\hat{f}$  with a zero-knot at true  $\theta$  and estimated  $\hat{\theta}$  respectively. Then,  $\tilde{f}(x) = \mathbf{I}_{[0,\theta]}(x)\tilde{s}_0(x) + \mathbf{I}_{(\theta,1]}(x)\tilde{s}_1(x)$  and  $\hat{f}(x) = \mathbf{I}_{[0,\hat{\theta}]}(x)\hat{s}_0(x) + \mathbf{I}_{(\hat{\theta},1]}(x)\hat{s}_1(x)$ , where  $\tilde{s}_j$ ,  $j = 0, 1$  and  $\hat{s}_j$ ,  $j = 0, 1$  are the corresponding spline estimates.

In order to simplify the notation, denote hereafter  $\mathbb{E}\{\cdot\mathbf{I}_F\}$  by  $\mathbb{E}_F\{\cdot\}$ . We have  $\mathbb{E}_F\{\|\hat{f} - f\|_2^2\} = \mathbb{E}_F\{\|f - s\|_2^2\} + O(n^{-2m/(2m+1)})$ .

Consider only the case  $\hat{\theta} \leq \theta$  since the opposite case can be treated in a similar way. Then,

$$\begin{aligned} \mathbb{E}_F\{\|\hat{f} - s\|_2^2\} &= \mathbb{E}_F\left\{\int_0^{\hat{\theta}} (s(x) - \hat{f}(x))^2 dx\right\} + \mathbb{E}_F\left\{\int_{\hat{\theta}}^{\theta} (s(x) - \hat{f}(x))^2 dx\right\} \\ &\quad + \mathbb{E}_F\left\{\int_{\theta}^1 (s(x) - \hat{f}(x))^2 dx\right\} \\ &= (A) + (B) + (C). \end{aligned} \tag{A.14}$$

For the first term (A) in (A.14) we have

$$\begin{aligned} &\mathbb{E}_F\left\{\int_0^{\hat{\theta}} (s(x) - \hat{f}(x))^2 dx\right\} \\ &= \mathbb{E}_F\left\{\int_0^{\hat{\theta}} (\hat{s}_0(x) - s_0(x))^2 dx\right\} \\ &\leq \mathbb{E}_F\left\{\int_0^{\hat{\theta}} (\hat{s}_0(x) - \tilde{s}_0(x))^2 dx\right\} + \mathbb{E}_F\left\{\int_0^{\hat{\theta}} (\tilde{s}_0(x) - s_0(x))^2 dx\right\} \\ &= (A_1) + (A_2). \end{aligned}$$

By Proposition 3.2, the term  $(A_2)$  is  $O(n^{-2m/(2m+1)})$ . For  $(A_1)$ , note that by construction and assumption M6, the sup-norm distance between the knots of  $\hat{s}$  and those of  $\tilde{s}$  is bounded above by  $O\left(q_n^{-1}(\theta - \hat{\theta})\right) \simeq o_P(n^{-2m/(2m+1)})$ . The knots of  $\hat{s}$  may be viewed as the set of knots of  $\tilde{s}$  perturbed by an amount  $o(n^{-2m/(2m+1)})$ . Using Theorem 6.2 of Lyche and Mørken (1999) it follows that  $\|\hat{s}_0(x) - \tilde{s}_0(x)\|_\infty^2 \leq O(n^{-2m/(2m+1)})$  and therefore  $(A_1)$  is also  $O(n^{-2m/(2m+1)})$ .

The third term (C) in the right-hand side of (A.14) can be handled exactly in the same way to verify that it is  $O(n^{-2m/(2m+1)})$ . Finally, it is

easy to see that the remaining term (B) is  $O\left(\mathbb{E}_F(|\theta - \hat{\theta}|)\right)$  which is also  $O(n^{-2m/(2m+1)})$  by Corollary 2.1.

Summarizing, so far we have proved the proposition for  $D = 1$ . For a general (but still finite!)  $D$ , one may partition the unit interval by  $\theta_l$  and  $\hat{\theta}_l$ ,  $l = 1, \dots, D$  in a similar way and use the established result for a single jump to obtain the rates stated in Proposition 3.3 similarly to the proof of Proposition 2 in Antoniadis and Gijbels (2002).  $\square$

*Acknowledgments.* This research was supported by the “Projet d’Actions de Recherches Concertées”, No. 93/98-164 of the Belgian Government and by NSF grants DMS 0505133 and DMS-0652524. Felix Abramovich and Marianna Pensky would like to thank Anestis Antoniadis and Department of Statistics of the Laboratory Jean Kuntzmann – University of Joseph Fourier for warm hospitality while visiting Grenoble to carry out part of this work. The authors would like to thank the anonymous referees for their helpful comments.

### References

- AGARWAL, G. and STUDDEN, W. (1980). Asymptotic integrated mean square error using least squares and bias minimizing splines. *Ann. Statist.*, **8**, 1307–1325.
- ANTONIADIS, A. and FAN, J. (2001). Regularization by wavelet approximations, *J. Amer. Statist. Assoc.*, **96**, 939–967.
- ANTONIADIS, A. and GIJBELS, I. (2002). Detecting abrupt changes by wavelet methods. *J. Nonparametr. Stat.*, **14**, 7–29.
- AMATO, U. , ANTONIADIS, A. and PENSKY, M. (2006). Wavelet kernel penalized estimation for non-equispaced design regression. *Stat. Comput.*, **16**, 37–56.
- BATTAGLIA, M. (2000). *The Influence of Overstory Structure on Understory Light Availability in a Longleaf Pine (Pinus Palustris Mill) Forest*. Master’s Thesis, Department of Forestry, Virginia Tech, USA.
- BERKNER, K. and WELLS, R.O. (1997). A fast approximation to the continuous wavelet transform with applications, *Conference Record of the Thirty-First Asilomar Conference on Signals, Systems & Computers, 1997*, **2**. IEEE, Pacific Grove, CA, 1235–1239.
- CHU, C.K., GLAD, I.K., GODTLIEBSEN, F. and MARRON, J.S. (1998). Edge-preserving smoothers for image processing (with discussion). *J. Amer. Statist. Assoc.*, **93**, 526–556.
- COHEN, A., DYN, N. and MATEI, B. (2003). Quasilinear subdivision schemes with applications to ENO interpolation. *Appl. Comput. Harmon. Anal.*, **15**, 89–116.
- COIFMAN, R.R. and DONOHO, D.L. (1995). Translation invariant denoising, in *Wavelets in Statistics*, Antoniadis, A. and Oppenheim, G. (eds.), Lecture Notes in Statistics **103**, Springer-Verlag, New York, 125–150.

- DAUBECHIES, I. (1992). *Ten Lectures on Wavelets*, CBMS-NSF Series in Applied Mathematics, SIAM, Philadelphia.
- DAUBECHIES, I., DEFRISE, M. and DE MOL C. (2004). An iterative thresholding algorithm for linear inverse problems with a sparsity constraint, Technical Report, Department of Mathematics, Princeton University.
- DONOHO, D.L. (1992). Interpolating wavelet transforms. Technical Report, Department of Statistics, Stanford University.
- DONOHO, D.L. and JOHNSTONE, I.M. (1994). Ideal spatial adaptation by wavelet shrinkage. *Biometrika*, **81**, 425–455.
- EILERS, P.H. and MARX, B.D. (1996). Flexible smoothing with B-splines and penalties. *Statist. Sci.*, **11**, 89–121.
- EUBANK, R.L. (1999). *Nonparametric Regression and Spline Smoothing* (Second Edition). Marcel Dekker, New York.
- FAN, J. and LI, R.Z. (2001). Variable selection via penalized likelihood. *J. Amer. Statist. Assoc.*, **96**, 1348–1360.
- FINK, D. and WELLS, M. (2004). Adaptive multiorder penalized regression splines. Technical report, Department of Statistical Science, Cornell University.
- FRANK, I.E. and FRIEDMAN, J.H. (1993). A statistical view on some chemometrics regression tools (with discussion). *Technometrics*, **35**, 109–148.
- GILL, R. and BARON, M. (2004). Consistent estimation in generalized broken-line regression. *J. Statist. Plann. Inference*, **126**, 441–460.
- HE, X., SHEN, L. and SHEN, Z. (2001). A data-adaptive knot selection scheme for fitting splines. *IEEE Signal Processing Letters*, **8**, 137–139.
- JAFFARD, S. (1989). Exposants de Hölder en des points donnés et coefficients d'ondelettes. *C.R. Acad. Sci. Paris*, **308**, Série I, 79–81.
- KOO, J.-Y. (1997). Spline estimation of discontinuous regression functions. *J. Comput. Graph. Statist.*, **6**, 266–284.
- LEE, T.C.M. (2002). Automatic smoothing for discontinuous regression functions. *Statist. Sinica*, **12**, 823–842.
- LYCHE, T. and MØRKEN, K. (1999). The sensitivity of a spline function to perturbations of the knots. *BIT*, **39**, 305–322.
- MALLAT, S. (1989). Theory for multiresolution signal decomposition: the wavelet representation. *IEEE Trans. PAMI*, **11**, 674–692.
- MALLAT, S. and HWANG, W.L. (1992). Singularity detection and processing with wavelets. *IEEE Trans. Inform. Theory*, **2**, 617–643.
- NIKOLOVA, M. and NG, M. (2005). Analysis of half-quadratic minimization methods for signal and image recovery. *SIAM J. Sci. Comput.*, **27**, 937–966.
- OSBORNE, M.R., PRESNELL, B. and TURLACH, B.A. (2000). On the LASSO and its dual. *J. Comput. Graph. Statist.*, **9**, 319–337.
- ODUSHOORN, C.G.M. (1998). Asymptotically minimax estimation of a function with jumps. *Bernoulli*, **4**, 15–33.
- PINSKER, M. (1980). Optimal filtration of square-integrable signals in Gaussian noise. *Problems Inform. Transmission*, **16**, 120–133.
- RUPPERT, D. and CARROLL, R.J. (2000). Spatially-adaptive penalties for spline fitting. *Aust. N.Z.J. Stat.*, **42**, 205–223.

SHENSA, M.J. (1992). The discrete wavelet transform: Wedding the À Trous and Mallat algorithms. *IEEE Trans. Inform. Theory*, **40**, 2464–2482.

TIBSHIRANI, R. (1996). Regression shrinkage and selection via lasso. *J. Roy. Statist. Soc. Ser. B*, **58**, 267–288.

WANG, Y. (1995). Jump and sharp cusp detection by wavelets. *Biometrika*, **82**, 385–397.

ZHOU, S. and SHEN, X. ( 2001 ). Spatially adaptive regression splines and accurate knot selection schemes. *J. Amer. Statist. Assoc.* **96**, 247–259.

FELIX ABRAMOVICH  
 DEPARTMENT OF STATISTICS  
 AND OPERATIONS RESEARCH  
 TEL AVIV UNIVERSITY  
 TEL AVIV 69978, ISRAEL  
 E-mail: felix@post.tau.ac.il

ANESTIS ANTONIADIS  
 LABORATOIRE LJK  
 UNIVERSITY JOSEPH FOURIER  
 BP 53, 38041 GRENOBLE CEDEX 9  
 FRANCE  
 E-mail: anestis.antoniadis@imag.fr

MARIANNA PENSKY  
 DEPARTMENT OF STATISTICS  
 UNIVERSITY OF CENTRAL FLORIDA  
 ORLANDO, FL 32816 -1364, USA  
 E-mail: mpensky@pegasus.cc.ucf.edu

Paper received January 2007; revised April 2007.