

<p>Provided for non-commercial research and educational use. Not for reproduction, distribution or commercial use.</p>
--

# Mathematica Balkanica

Mathematical Society of South-Eastern Europe  
A quarterly published by  
the Bulgarian Academy of Sciences – National Committee for Mathematics

---

The attached copy is furnished for non-commercial research and education use only. Authors are permitted to post this version of the article to their personal websites or institutional repositories and to share with other researchers in the form of electronic reprints.

Other uses, including reproduction and distribution, or selling or licensing copies, or posting to third party websites are prohibited.

For further information on Mathematica Balkanica visit the website of the journal  
<http://www.mathbalkanica.info>

or contact:

Mathematica Balkanica - Editorial Office;  
Acad. G. Bonchev str., Bl. 25A, 1113 Sofia, Bulgaria  
Phone: +359-2-979-6311, Fax: +359-2-870-7273,  
E-mail: [balmat@bas.bg](mailto:balmat@bas.bg)

## Penalized Wavelet Estimation with Besov Regularity Constraints <sup>1</sup>

*Lubomir T. Dechevsky \**, *James O. Ramsay \*\**,  
*Spiridon I. Penev \*\*\**

*Presented by P. Kenderov*

A wavelet counterpart to Wahba's spline smoothing technique is developed for the case of fitting less regular curves. Our approach is a realization of the smoothness penalty method of Tikhonov regularization of ill-posed inverse stochastic problems within a wavelet setting. The penalized cost functional to be minimized is Peetre's K-functional between Besov spaces. The regularity of the curve is discussed in terms of the size of its seminorm in homogeneous Besov spaces  $\dot{B}_{p,q}^s$  with a relatively small value of the smoothness index  $s > 0$ . Penalized  $L_2$ -estimation with Besov-type constraints, considered in the literature, is included as a partial case. The optimal solution of the penalization problem is in the form of a wavelet expansion whose coefficients are obtained by appropriate level- and/or space-dependent shrinking of the empirical wavelet coefficients. Thanks to the use of wavelets, both density and regression estimation can be treated in a somehow unified way. In the case of regression-function estimation, an enhanced version of cross validation (generalized full cross validation) is implemented for the choice of the smoothing parameter. Numerical examples illustrate the advantage of our procedures in comparison to more standard wavelet methods when the regularity is small and sample sizes are moderate. The approach is very flexible and allows for a diversity of extensions. Some first extensions can be found in Section 7, among which are the iterative individual shrinking estimator and the self-similar fractal estimator. The other extensions considered are grouped in Appendix B, which is in a sense the most advanced part of the paper, and can be considered as an outline of a plan for further research. For conciseness of presentation,

---

<sup>1</sup>Supported by the Natural Sciences and Engineering Research Council of Canada and the Australian Research Council



only the univariate case is considered in detail, but all statements and their proofs do admit multidimensional generalization.

*Mathematics Subject Classification:* Primary 42C15, 46E35, 62G05, 65D10; Secondary 26A30, 26D10, 35K10, 35K30, 39B12, 41A63, 42C20, 46N30, 47B10, 47D06, 49M45, 62J07, 65K10, 65R30, 65U05

*Key Words and Phrases:* penalized nonparametric regression and density estimation, consistency of estimation and cross validation, asymptotic-minimax rate of estimation, atomic decomposition of Besov spaces by orthonormal wavelets, K-functional Tikhonov regularization of ill-posed inverse stochastic problems, iterative individual shrinking of empirical wavelet coefficients, self-similar fractal estimator.

## 1. Introduction

This paper consists of three parts: main body, Appendix A and Appendix B. References to parts of the appendices are denoted in the text by  $A^2$  or  $A^7, B^3, B^8$ , and so on.

Our purpose is to suggest a wavelet-based approach that parallels spline smoothing techniques, as described by Wahba [117], and to apply it for estimating spatially inhomogeneous curves. Our method is thus an extension of the smoothness penalty approach to wavelets. The penalized cost functional to be minimized is Peetre's K-functional between Besov spaces. Its minimizer is in the form of a wavelet expansion with coefficients obtained by level- and/or space-dependent shrinking of the empirical wavelet coefficients. In the case of nonparametric regression, an enhanced version of cross validation (generalized full cross validation) is implemented. Measuring the estimation rates and the regularity of the curve in Besov spaces allows for optimizing the cross-validation criterion not only with respect to the penalizing smoothing parameter, but also with respect to the smoothness index of the Besov space. This brings about control over fractional-order derivatives<sup>A1</sup> of the curve and makes the new estimator superior to Wahba's smoothing natural spline estimator when estimating less regular curves. To quote Hall and Patil [69]: "we do not suggest that wavelet methods should supplant existing techniques for estimating very smooth curves,

but we do agree that they have ready application to problems where curves change relatively sharply. Thus, they complement existing methodology for nonparametric curve estimation, rather than replace it". In the specific context of the smoothness penalty approach there are two additional advantages of the utilization of wavelets. (a) While the optimization problem in Wahba's approach is intrinsically *non-parametric* (consisting of a large number of boundary-value problems for differential equations plus subsequent comparison of the boundary values), its analogue in the wavelet setting is *parametric*<sup>A2</sup>. (b) While the penalized smoothing spline technique is designed specifically for regression-function estimation, its wavelet analogue allows treating *both* non-parametric regression and density estimation in a somehow *unified* way. More precisely, we shall be considering the following regression and density models.

The regression model is:  $Y_i = f(x_i) + \epsilon_i$ ,  $i = 1, 2, \dots, n$ , where  $x_i = a_1 + (a_2 - a_1)i/n$ ,  $[a_1, a_2]$  being the definition interval. For the i.i.d. errors  $\epsilon_i$  we assume  $E\epsilon_1 = 0$ ,  $E\epsilon_1^2 = \delta^2$ .

In the case of density estimation, the sample consists of  $n$  i.i.d. observations  $X_i$ ,  $i = 1, 2, \dots, n$  with unknown density  $f(x)$ .

It is outside the scope of the paper to compare the merits of our solutions of these problems with the many alternative wavelet-based approaches (we can easily list at least ten alternatives) since this would require performing a very long simulation study and spacious discussion. This is why only one, well-known, method - the soft thresholding - was singled out for the sake of graphical comparison. The numerical simulations and graphical results presented in the sequel illustrate the advantages of the shrinking over the thresholding procedure when regularity is low and sample sizes are moderate. We should underscore though that our primary goal in the present paper is not demonstration of the superiority of the fits we get in comparison to soft thresholding and other standard wavelet methods. The main merits of our approach at this stage can rather be summarized as follows. (a) This approach offers an excellent trade-off between good estimation and *simplicity* of the estimator. (b) While satisfactory practical estimation with wavelet estimators (in particular, soft- and hard-thresholded)

which are asymptotic-minimax within the "asymptopia" paradigm (see Donoho, Johnstone, Kerkycharian and Picard [58]) can in general be only achieved for large sample sizes ( $n \approx 10^4$ ), the penalization approach considered here yields comparable quality of the fit already for moderate samples ( $n \approx 10^3$ ), under virtually no additional assumptions compared to soft thresholding provided that the Besov smoothness index is being *also optimized* in the cross-validation criterion. (c) The present approach is very flexible and can be extended simply and elegantly in diverse directions, thus allowing enriching and upgrading the features of the model, the expected ultimate reward being to get good fits for sample sizes as low as  $n \approx 10^2$  and, for some specific estimation problems, even less. Some of these extensions are briefly addressed in appropriate places (notably, Section 7, <sup>A3</sup>, Appendix B).

In the sequel only the univariate case is considered in detail for conciseness of presentation. However, all the main results can be generalized to the *multivariate* case. We shall make more detailed comments about this generalization in appropriate places later in the text.<sup>B21</sup>

Finally, a few more words on the organization of the paper. Sections 2 and 3 are introductory. Section 4 gives a brief exposition of our concept of *K-functional Tikhonov regularization* and explains why the *K-functional* approach is of fundamental importance for deterministic and nondeterministic smoothing and, in particular, for nonparametric penalized wavelet estimation. In Sections 5 and 6 a model case (the case of Hilbert spaces) is presented, and a model procedure for estimating the regularization parameter - cross validation - is studied in detail. In Section 7 numerical examples are discussed, together with some immediate extensions of the *K-functional* model. Appendix A contains the proofs, as well as "second reading" remarks to the text in the main body. Appendix B is essentially a list of further extensions and generalizations which can be considered as a program for further research on the topic of *K-functional* penalized estimation. In a sense, Appendix B is the most important and advanced part of this paper.

## 2. Thresholding

The study of statistical applications of wavelets received its initial boost in 1992-1994 with the development of an elegant general asymptotic-minimax theory in a series of papers by Donoho, Johnstone, Kerkycharian and Picard. In this theory thresholding of the empirical wavelet coefficients emerged as a key tool for smoothing and denoising. The loss functional with respect to which the estimation rates were measured was an  $L_p$ -norm. That essentially meant that only the function was being consistently estimated but not its derivatives. One reason for this was that in many cases the initially proposed "universal" threshold turned out to be locally too high or too low and did not allow for an optimal trade-off between bias and variance locally along the curve. Since then, Donoho and Johnstone have proposed several more refined threshold strategies for Gaussian white noise (notably the SURE shrink, WaveJS and oracle inequalities) that allow for better adaptivity (and have had important impact on the development of the currently most advanced adaptive thresholding strategies NeighBlock and NeighCoeff (see [28,29])). Further methodological advance in the asymptotic theory was achieved, in both density and regression estimation context, with the introduction of level-dependent thresholds by Donoho and Johnstone [57] for the case of nonparametric regression and by Delyon and Juditsky [53] for both the regression-function and density estimation problem. This refinement contributed to bringing down the additional log-factor appearing in the asymptotic rates when estimating sufficiently regular curves using wavelets. The loss functional in which these rates were estimated by Delyon and Juditsky was a Besov (quasi-)norm with possibly positive smoothness index. This meant that also derivatives of the function up to a certain order were being consistently estimated. Moreover, the estimation rates for the derivatives were also asymptotically optimal within the "asymptopia" paradigm of [58]. Despite these improvements, for practical purposes wavelet-based estimation still remained an essentially *large-sample* technique. This was part of the price to pay for making these estimators almost minimax in a broad range of function spaces simultaneously, as suggested in the "asymptopia"<sup>A4</sup>.

Recently, several thresholding approaches have been proposed to generate

wavelet estimators with improved performance on moderate samples. Variants of Bayesian approach have been suggested by several authors ([116,1,30]). Another approach was Nason's (see [92]). He considered the integrated square error of the thresholded estimator as a function of the threshold and tried to find its minimum. One problem there appeared to be that the cross-validation criterion depended non-smoothly on the threshold value and had numerous local extrema. This made it very difficult, if possible at all, to develop advanced rigorous theory for this estimator. A further refinement of the level-dependent approach of Delyon and Juditsky was suggested in the so-called block-thresholding procedure by Hall, Kerkycharian and Picard [68] in 1995 (followed by [70,89,25,28,29]), which helped to bring down the log-factor in the rates and to further improve the estimation of derivatives by combining wavelet coefficients into groups and using information about each one to assist in assessing its neighbours.

Yet, even after all these efforts, thresholded wavelet estimators still have their problems, especially when the sample sizes are small to moderate, the noise-to-signal ratio is relatively large, the white noise is not Gaussian, and the estimated function is not very smooth. The only realistic opportunity to systematically get good fits for small samples seems to be *to introduce in the model all available additional information about essential features of the curve*. Together with Bayesian approach, the *variational* approach seems to be the natural choice in such upgrading. However, optimization with respect to threshold values generally leads to non-linear, non-smooth, non-convex, multi-extremal problems with high computational complexity. An additional problem when estimating functions with low regularity and fractals is that thresholding methods tend to oversmooth the curve because they are well adapted for functions which gather their value on relatively few, relatively large wavelet coefficients only. Continuous fractals and functions with jumps do not generally fall into this class: they gather their value relatively uniformly from many wavelet coefficients on infinitely many levels and the vector of their wavelet coefficients does not necessarily look sparse. Due to these facts, for such types of functions a non-threshold shrinking procedure is appropriate. An alternative approach to thresholding, leading to smooth optimization criteria for determining the regu-

larization parameter, was proposed by Amato and Vuza [3] and, independently, by Antoniadis [9]. Their methods were implicitly based on a technique for explicit computation of K-functionals between Hilbert spaces, developed earlier by Dechevski [41] and Dechevsky [42,43]. Although of non-threshold type, the parameter settings in the methods of Amato and Vuza, and Antoniadis, made these methods appropriate for estimating very smooth and relatively spatially homogeneous functions only. On the contrary, in the context of spatially inhomogeneous functions with low regularity our penalized non-threshold shrinking wavelet estimator and its various extensions are at their best, and we recommend them as a valuable complement to existing thresholding techniques.

### 3. Splines, wavelets, function spaces

By the end of the 1960s the penalized smoothing-spline (or spline-smoothing) technique had already evolved into a very general penalization approach for solving *deterministic* smoothing problems (see [7]). Its *statistical* applications - to non-parametric regression - began approximately at that time. They were based on an intuitively appealing technique for selection of the penalizing (or smoothing) parameter: cross validation (notably *ordinary* and *generalized*) (see [117]). For the smoothing version of Schoenberg's problem ([117], formula (0.0.3)) there are examples where the performance of Wahba's natural-spline smoothing estimator is spectacular, even for sample sizes as small as 100 or so. However, getting such a good fit is only likely for sufficiently smooth functions, since the penalization term contains a homogeneous Sobolev semi-norm with integer smoothness index. Even if the curve is very smooth, it should not be too spatially inhomogeneous because the cross-validation criterion can only be optimized with respect to a single *global* smoothing parameter. It is for more spatially inhomogeneous and/or less regular curves that *wavelet* methods show their strength, in *both* regression and density estimation problems.

Let  $f(x)$ ,  $x \in \mathbf{R}^1$  be a locally integrable function. Its wavelet expansion

is

$$(1) \quad f(x) = \sum_{k \in Z} \alpha_{j_0 k} \varphi_{j_0 k}(x) + \sum_{j=j_0}^{\infty} \sum_{k \in Z} \beta_{jk} \psi_{jk}(x)$$

where  $j_0 \in Z$ ,  $\{\varphi_{j_0 k}$  and  $\psi_{jk}, j = j_0, j_0 + 1, j_0 + 2, \dots; k \in Z\}$  form an orthonormal basis of  $L_2(R)$ . Here  $\varphi_{jk}(x) = 2^{j/2} \varphi(2^j x - k)$ ,  $\psi_{jk}(x) = 2^{j/2} \psi(2^j x - k)$ ,  $j$  describes the resolution level. For better spatial and level localization,  $\varphi$  and  $\psi$  are assumed to be compactly supported. Convergence of the series in (1) is in the (quasi-)norm topology ([16], Section 3.10) of *Besov* spaces. The Besov scale  $B_{pq}^s(R^d)$  (resp.  $\dot{B}_{pq}^s(R^d)$ ) of *inhomogeneous* (resp. *homogeneous*) Besov spaces,  $0 < p \leq \infty$ ,  $0 < q \leq \infty$ ,  $s \in R$ , is the typical class of translation-invariant function spaces in which spatial inhomogeneity can be described. (For equivalent definitions of Besov spaces and their properties we refer to [16, 111–113].) In statistical context, Besov spaces have been introduced by Donoho for the purposes of asymptotic-minimax theory of wavelet estimators (see [58] for an overview). This approach is founded on a relatively recent discovery (Sickel [106]) of a remarkable family of equivalent (quasi-)norms in Besov spaces in terms of *wavelet coefficients*. In brief, the idea is as follows. In addition to the already made assumptions about  $\varphi$ ,  $\psi$ , suppose that for  $r > 0$  the linear span of  $\varphi(\cdot - k)$ ,  $k \in Z$ , contains all algebraic polynomials  $P$  with  $\deg P \leq [r]$  and  $\varphi \in B_{\infty\infty}^r$ . Then  $\varphi_{j_0 k}$ ,  $\psi_{jk}$ ,  $k \in Z$ ,  $j = j_0, j_0 + 1, \dots$ , form a Riesz basis simultaneously for all  $B_{pq}^s(R)$ ,  $0 < p \leq \infty$ ,  $0 < q \leq \infty$ ,  $\max(0, 1/p - 1) < s < r$ . Therefore, for the above range of parameters considered, if  $f \in B_{pq}^s(R)$ , then the series in (1) is always convergent in the (quasi-)norm topology of the space and

$$(2) \quad J_{spq}(\alpha, \beta) = [\|\alpha_{j_0 \cdot}\|_p^q + \sum_{j=j_0}^{\infty} (2^{j(s+(1/2)-(1/p))} \|\beta_{j \cdot}\|_p)^q]^{1/q}$$

is an equivalent norm in  $B_{pq}^s(R)$ . Here the notation

$$(3) \quad \|\alpha_{j_0 \cdot}\|_p = \left(\sum_{k \in Z} |\alpha_{j_0 k}|^p\right)^{1/p}; \quad \|\beta_{j \cdot}\|_p = \left(\sum_{k \in Z} |\beta_{jk}|^p\right)^{1/p}$$

has been used. Henceforward, we shall always make the aforementioned assumption that the analyzing wavelet has smoothness index  $r > s$ . For the purposes of Section 6 we also need to suppose that the wavelet basis is obtained via multiresolution analysis (see [36]) and in Subsection 6.2 an additional explicit assumption is required that the first  $[r]$  moments of  $\psi$  vanish. In Section 6 it is also implicitly assumed that  $\varphi$  and  $\psi$  have bounded variation  $\bigvee \varphi, \bigvee \psi$ , but this is true even for the Haar wavelet.

All notations related to wavelets and Besov spaces are the same as in [53].<sup>45</sup>

For simplicity, we shall always be working with wavelets obtained via multiresolution analysis on the whole  $\mathbb{R}$  (or  $\mathbb{R}^d$ ), and will be assuming that the support of the function in consideration is compact on its domain  $\Omega$ , because in this way we are working with spaces which are simultaneously closed subspaces of the function space over  $\Omega$ , and of the function space over  $\mathbb{R}^d$ . Because of their closedness, these subspaces enjoy all the essential properties of the entire function spaces. This trick is done only for simplicity of presentation; it is quite straightforward to extend our constructions to periodic functions in one and several dimensions by considering periodic wavelets and (when the sample size is a power of two) periodic orthogonal discrete wavelet transform (ODWT) (see, e.g., [57] and the references therein). It is also easy to further extend our considerations also for the finite interval (see [31,33]), as well as for a finite hyperrectangle in  $d$  dimensions. Moreover, it is even possible to extend our constructions, as described in more detail in <sup>B20,B21</sup>, to function classes on Lipschitz-graph domains  $\Omega$  in  $d$  dimensions, by utilizing wavelets generated via multiresolution analysis on  $\Omega$ , as proposed by Cohen, Dahmen and DeVore [32] in their wavelet-based approach to deriving Whitney-type extension theorems.

#### 4. The $K$ -functional

The  $K$ -functional has been introduced by Peetre [94] in 1963 as a tool for measuring smoothness in abstract spaces. Since then, it has had a rich variety of important applications, for example: interpolation of operators and interpo-



with moduli of smoothness ([24,80,39,41,42]); general characterization of best approximation ([97]), and many others. For the purposes of these applications it has also been generalized in several aspects. Here we consider one of these generalizations - the  $K_\rho$ -functional between two (quasi-semi-)normed spaces  $A$ ,  $B$ :

$$(4) \quad K_\rho(t, f; A, B) = \inf_{f=a+b} (\|a\|_A^\rho + t^\rho \cdot \|b\|_B^\rho)^{1/\rho}, 0 < t < \infty, 0 < \rho \leq \infty$$

(with max-modification for  $\rho = \infty$ ),  $f \in A + B$  (the sum of  $A$  and  $B$  as linear vector spaces)<sup>A6</sup>.

Taking in consideration the clear resemblance between Wahba's statistical spline-smoothing model and the deterministic model of the  $K$ -functional, it is amazing that, ever since the 1960s, two distinct and very *complementary* theories related to the two models have been developed, without almost any methodological interconnection. In the last several years a number of interesting papers on penalized  $L_2$ -estimation have appeared, some of them considering also Besov penalties (see [55,54,19,13,72,9,4]), but, to the best of our knowledge, so far the only one who has essentially exploited the analogy between Wahba's model and the  $K$ -functional has been Cox [35], Theorems 3.3,4. In a sequence of papers (see also [49,50]), we intend to pursue a study of the applications of the theory of the  $K$ -functional to statistical estimation.

To illustrate the conceptual power of the  $K$ -functional, and for further use, let us discuss the "smoothing paradigm" in a general  $K$ -functional setting. Our consideration will be an extension, generalization and further development of the approaches of Anselone and Laurent [7] and Cox [35]. Let  $A$ ,  $B$  be quasi-Banach spaces with  $B \hookrightarrow A$  (i.e.,  $B \subset A$  and there exists  $c < \infty$  such that  $\|b\|_A \leq c \cdot \|b\|_B$  for any  $b \in B$ ). If for any  $a \in A + B = A$  there exists  $b^* = b^*(a)$  so that  $K_\rho(t, a; A, B) = \min_{b \in B} (\|a-b\|_A^\rho + t^\rho \|b\|_B^\rho)^{1/\rho} = (\|a-b^*\|_A^\rho + t^\rho \|b^*\|_B^\rho)^{1/\rho}$ , then the  $K$ -functional measures the quality of the approximation of  $a$  by  $b \in B$  for a preassigned (through the  $t$ -value) trade-off between fidelity of the data and smoothness of the solution. The embedding  $B \hookrightarrow A$  ensures that the presence of the penalization term has smoothing effect. This general smoothing model can

be further enhanced by considering a *homogeneous version* of  $B$ . Assume that the quasi-norm in  $B$  is of the particular form  $\|b\|_B = (\|b\|_A^\rho + \|Tb\|_{B_1}^\rho)^{1/\rho}$ , where  $T$  is a densely defined on  $A$ , linear, possibly unbounded, closed operator from  $A$  to the quasi-Banach space  $B_1$ . Then, by the closedness of  $T$ ,  $B$  is quasi-Banach; clearly,  $B \hookrightarrow A$  and  $\|\cdot\|_B$  is the *graph* quasi-norm of  $T$ . The homogeneous version of  $B$  is the quasi-(semi-)normed complete space  $\dot{B}$  with  $\|\cdot\|_{\dot{B}} = \|T\cdot\|_{B_1}$ . Smoothing with respect to  $K_\rho(t, a; A, \dot{B})$  is generally not identical but always equivalent to smoothing via  $K_\rho(t, a; A, B)$  in the sense that

$$(5) \quad K_\rho(t, \cdot; A, B) \asymp K_\rho(t, \cdot; A, \dot{B}) + \min(1, t) \cdot \|\cdot\|_A, \quad 0 < t < \infty,$$

with equivalence constants *independent* of the smoothing parameter  $t$ .<sup>47</sup> In Wahba's model  $B$  and  $\dot{B}$  are an inhomogeneous and a homogeneous Sobolev space; in our wavelet model  $B = B_{pq}^s$ ,  $\dot{B} = \dot{B}_{pq}^s$ ,  $A = B_{\pi u}^\sigma$ , where the relation between  $(\pi, u, \sigma)$  and  $(p, q, s)$  must be such that: (a) for the selected  $r$  (see Section 3) (2,3) hold for both triples; (b)  $B \hookrightarrow A$ . Hence, by *embedding* theorems for Besov spaces (see [16,111]), the range of the two triples can be shown to *exactly coincide* with their respective range in [53], Theorem 1.<sup>48</sup> For a fixed preassigned value of  $t$ , the role of the parameters  $\pi, p, u, q, \sigma, s$  in the  $K$ -functional model is essentially the same as in [53]. In particular, the estimated function is assumed to be in  $B_{pq}^s$ . The loss functional is the quasi-norm in  $B_{\pi u}^\sigma$ . If, however, the problem is inverse, i.e.,  $t$  has to be *estimated* in order to recover  $f$ , then the meaning of  $\pi, p, u, q, \sigma, s$  is different and these parameters may eventually also be subject to estimation. To explain their meaning in this case, let us recall the general inhomogeneous  $K$ -functional model. According to it, the estimated  $f$  belongs to  $X$  - an *intermediate* quasi-Banach space between  $A$  and  $B$  (see [16]). In our case  $B \hookrightarrow A$ , therefore,  $B \hookrightarrow X \hookrightarrow A$ . For our purposes it suffices to assume that  $X = (A, B)_{\theta, \tau}$  or  $X = (A, B)_{[\theta]}$ ,  $0 < \theta < 1$ ,  $0 < \tau \leq \infty$ , i.e.,  $X$  is an interpolation space obtained by the real method of Lions and Peetre, or by the complex method of Calderon, Krein and Lions, respectively. We shall concentrate on the real method which is the

one directly related to the  $K$ -functional. Assume, additionally, that  $B$  is *dense* on  $A$ . (All the above assumptions, including the latter one, are fulfilled for Sobolev and Besov spaces.) Because of  $B \hookrightarrow A$ , by saturation (cf. <sup>A7</sup>) argument,  $\|f\|_X = \|f\|_{(A,B)_{\theta,\tau}} \asymp \|f\|_A + [\int_0^C (t^{-\theta} K_{\rho}(t, f; A, B))^{\tau} dt / t]^{1/\tau}$ , for any, henceforward fixed,  $C : 0 < C \leq \infty$ . Because  $B$  is dense in  $A$ , for any  $\epsilon > 0$  one can choose  $f_{\epsilon} \in B$  so that  $\|f - f_{\epsilon}\|_A < \epsilon$ . It can be seen that  $K(t, f_{\epsilon}; A, B) \asymp \min(1, t) \cdot \|f_{\epsilon}\|_B$ . Because of the latter, the *Generalized First Mean-value Theorem*, valid for improper integrals (see [64], Section 487) can be applied, which yields

$$(6) \quad \|f_{\epsilon}\|_X \asymp \|f_{\epsilon}\|_A + \frac{K_{\rho}(\tilde{t}, f_{\epsilon}; A, B)}{\min(1, \tilde{t})} \left[ \int_0^C (t^{-\theta} \cdot \min(1, t))^{\tau} \frac{dt}{t} \right]^{1/\tau}.$$

The integral is finite and its value does not exceed  $[(1 - \theta)\theta\tau]^{-1/\tau}$  for any  $C : 0 < C \leq \infty$ . (The case  $\tau = \infty$  is analogous, but simpler.) Here  $\tilde{t} = \tilde{t}_{\epsilon} : 0 \leq \tilde{t} \leq C \leq \infty$  is the *target value* of the smoothing parameter subject to estimation in the statistical smoothness-penalty models. In these models  $\epsilon = \epsilon(n) = o(1)$  as  $n \rightarrow \infty$ , where  $n$  is the sample size. The density of  $B$  on  $A$  ensures that the problem for penalized estimation of  $f$  is *correctly posed in the sense of Tikhonov*. For our particular model, however, we need further enhancement. Assuming  $X$  to be fixed, let  $A', B'$  be quasi-Banach spaces, such that  $B \hookrightarrow B' \hookrightarrow X \hookrightarrow A' \hookrightarrow A$ . For  $X = (A, B)_{\theta,\tau}$  one can choose  $A' = (A, B)_{\theta_0,\tau}$ ,  $B' = (A, B)_{\theta_1,\tau}$ , where  $0 < \theta_0 < \theta < \theta_1 < 1$ . By the Reiteration theorem (see [16]),  $X = (A', B')_{\theta',\tau}$ , where  $\theta' = \frac{\theta - \theta_0}{\theta_1 - \theta_0} \in (0, 1)$ . This means that (6) can be applied again, this time with  $A, B, \theta$  and  $\tilde{t}$  replaced by  $A', B', \theta'$  and  $\tilde{t}'$ , respectively. Therefore, estimation of  $f_{\epsilon}$ , hence, of  $f$ , can be *improved* by estimating  $t$  not only for  $A$  and  $B$  ( $\theta_0 = 0, \theta_1 = 1$ ), but also for  $A', B'$  corresponding to any  $\theta_0, \theta_1 : 0 < \theta_0 < \theta < \theta_1 < 1$ . By (5), the above considerations are also valid, *mutatis mutandis*, for the homogeneous version, too.

Returning to the interpretation of the Besov indices  $(\pi, u, \sigma)$  and  $(p, q, s)$ , we now see that penalized estimation via  $K_{\rho}(\tilde{t}, f; B_{\pi u}^{\sigma}, B_{pq}^s)$  is specialized for

$f \in B_{p',q'}^{s'}$ , where  $s' = (1 - \theta)\sigma + \theta s$ ,  $\frac{1}{p'} = \frac{1-\theta}{\pi} + \frac{\theta}{p}$ ,  $\frac{1}{q'} = \frac{1-\theta}{u} + \frac{\theta}{q}$ ,  $0 < \theta < 1$ ,  $q' = \tau \in (0, \infty]$ . In particular, if  $\pi = p$ , then  $\sigma$  and  $s$  are always a *lower* and an *upper bound* for the exact  $s'$ , respectively.<sup>A9,A10,A11</sup> The methods for penalized  $L_2$  estimation with Besov penalty (see, e.g., [9,4]) deal with the case  $\pi = u = 2$ .

### 5. Exact solution for the case $\pi = u = p = q = \rho = 2$ , $\sigma = 0$

In this and next section we shall be considering the Hilbert case  $\pi = u = p = q = 2$  where comparison can be made with the spline-smoothing model.  $\sigma$  is fixed to 0, that is, we shall be estimating the function itself, but not its derivatives. Besides  $t$ , only one additional continuous parameter -  $s$  - will be optimized most of the time. For the general case about  $\pi, p, u, q, \sigma, s$ , see <sup>B9</sup>. The estimator has the form

$$(1') \quad \tilde{f}(x) = \sum_{k \in Z} \tilde{\alpha}_{j_0 k} \varphi_{j_0 k}(x) + \sum_{j=j_0}^{j_1} \sum_{k \in Z} \tilde{\beta}_{jk} \psi_{jk}(x),$$

with a reasonable choice of the resolution levels  $j_0$  and  $j_1$ , variants of which will be discussed below.

The Hilbert case is of special importance, because then the estimator (1') is *linear* in the data for any fixed  $t$  in the  $K$ -functional. For the Hilbert case, a very general formula for the explicit computation of the  $K_2$ -functional between any semi-Hilbert spaces was obtained in [41-43], where some deterministic applications to the theory of real and complex interpolation spaces and to approximation theory were considered. For our purposes here it suffices to invoke only a partial case of these general results<sup>A13</sup>, as follows. The spaces  $A$  and  $B$  in the definition of the  $K_\rho$ -functional are any two Hilbert spaces such that  $A \cap B$  is dense in  $A$  and  $B$ , and the cardinality of  $B$  does not exceed that of  $A$ . Let  $L$  be any, henceforward fixed, isometric operator acting from  $A \cap B$  to  $A$ , with  $\langle La, La \rangle_A = \langle a, a \rangle_B$ ,  $a \in A \cap B$ , where  $\langle \cdot, \cdot \rangle_A$  and  $\langle \cdot, \cdot \rangle_B$  are the scalar products in  $A$ ,  $B$ , respectively. Then ([43], Theorem 2.1, (i,ii)),  $L$  is a densely defined, possibly unbounded, closed linear (DDPUCL) operator on  $A$ ;  $L^*$  (the Hilbert adjoint of  $L$ ) is also a DDPUCL operator with  $(L^*)^* = L$ ;  $L^*L$  is a

densely defined, possibly unbounded, positive definite self-adjoint linear operator on  $A$ . Moreover ([43], Corollary 2.1 and Remark 2.6),  $K_2(t, a; A, B) = S a$ , where  $S$  is the extension onto  $A + B$  of the densely defined continuous sublinear functional  $\bar{S} : A \cap B \rightarrow R_+$ , defined by

$$\bar{S}a := \left( \int_{Sp(L^*L)} \frac{t^2 \lambda}{1 + t^2 \lambda} d\langle \mathcal{E}_\lambda a, a \rangle_A \right)^{1/2}, \quad a \in A \cap B,$$

where  $\mathcal{E}_\lambda$  is a spectral function of  $L^*L$  (that is,  $\mathcal{E}_\lambda$  is a spectral resolution of the identity operator  $I$  on  $A \cap B$ , generated by  $L^*L$ ),  $Sp(L^*L) \subset R_+$  is the spectrum of  $L^*L$ .

An intermediate step to obtaining the above integral representation for  $K_2(t, a; A, B)$  is the computation of the optimal  $b^* \in B$  on which  $\min_{b \in B}$  in the definition of the  $K_2$ -functional is attained ([43], formula (5.2.13)):

$$b^* = (I + t^2 L^* L)^{-1} a = \int_{Sp(L^*L)} \frac{1}{1 + t^2 \lambda} d\mathcal{E}_\lambda,$$

where the integration is in the sense of the abstract Stieltjes integral with respect to the spectral measure  $d\mathcal{E}_\lambda$ .

In the case when  $\|b\|_B = (\|b\|_A^2 + \|Tb\|_A^2)^{1/2}$ , where  $T$  is defined in Section 4, a homogeneous semi-Hilbert version  $\dot{B}$  of the Hilbert space  $B$  can be considered, as discussed in Section 4. Then, the formulae about  $\bar{S}a$  and  $b^*$  continue to hold ([43], Theorems 3.1.1 and 3.1.1a) also for  $K_2(t, a; A, \dot{B})$ , with  $L = T$ . For the homogeneous case,  $L^*L$  is positive *semi-definite* ( $0 \in Sp(L^*L)$ ). The case when  $A = B_{22}^s$ ,  $B = B_{22}^s$ ,  $\dot{B} = \dot{B}_{22}^s$ ,  $0 \leq \sigma < s < r$ , with the equivalent wavelet-coefficient norms (2,3), starting from level  $j_0$ , is a very simple partial case of the above general consideration. For example, let us consider the homogeneous model. In the present context, the operator  $L = T$  is self-adjoint itself:  $L = L^*$ , therefore,  $L^*L = L^2$  and it suffices to study the spectral properties of  $L$ . The spectrum  $Sp(L)$  of  $L$  is discrete; every  $\lambda \in Sp(L)$  is an eigenvalue; the eigenvalues are:  $\lambda_0 = 0$ ,  $\lambda_j = 2^{js}$ ,  $j = j_0, j_0 + 1, \dots$ . The respective eigenspaces are the subspaces in the multiresolution analysis

generated by  $\varphi$  and  $\psi$  (with the usual notations  $V_{j+1} = V_j \oplus W_j$ ): the eigenspace of  $\lambda_0$  is  $V_{j_0}$ , that of  $\lambda_j$  is  $W_j$ ,  $j = j_0, j_0 + 1, \dots$ . The formula about  $b^*$  now implies the following

**Lemma 1.** <sup>A13, A14</sup> *Let  $s > 0$  be fixed,  $t > 0$  be fixed and the quality of fit be measured by the  $K$ -functional  $K_2(t, \hat{f}; B_{2,2}^0, \dot{B}_{2,2}^s)$ . Then, the  $K$ -functional is attained for the minimizing estimator (1') with coefficients*

$$(7) \quad \tilde{\alpha}_{j_0 k} = \hat{\alpha}_{j_0 k}, \quad k \in Z, \quad \tilde{\beta}_{jk} = \frac{\hat{\beta}_{jk}}{1 + t^2 \cdot 2^{2js}}, \quad j = j_0, \dots, j_1.$$

The empirical wavelet coefficients are:

$$\hat{\alpha}_{j_0 k} = \frac{1}{n} \sum_{i=1}^n \varphi_{j_0 k}(X_i), \quad \hat{\beta}_{jk} = \frac{1}{n} \sum_{i=1}^n \psi_{jk}(X_i)$$

in the density estimation case, and

$$\hat{\alpha}_{j_0 k} = \frac{a_2 - a_1}{n} \sum_{i=1}^n \varphi_{j_0 k}(x_i) y_i, \quad \hat{\beta}_{jk} = \frac{a_2 - a_1}{n} \sum_{i=1}^n \psi_{jk}(x_i) y_i$$

in the case of regression.

Formula (7) tells us that the resulting estimator (1') is of *level-dependent non-threshold shrinkage* type.

The parametrization of the nonparametric problem achieved by the atomic decomposition via orthogonal wavelets makes it very easy to give a direct proof of Lemma 1 (see Appendix A). This simple proof of (7) has been obtained, independently from the earlier results of Dechevski [41] and Dechevsky [42,43], by Amato and Vuza [3] and Antoniadis [9]. To the best of our knowledge, the non-threshold shrinking model (7) has not been explored in the case of density estimation. In the case of regression-function estimation, Amato and Vuza [3-6] and Antoniadis [9] seem to be the first to propose concrete statistical procedures for the estimation of the regularization parameter  $t^2$  for (7) (see also [2]). The methods of Antoniadis and of Amato and Vuza are both specialized for work

with sample sizes which are a power of two, when ODWT can be applied. In both methods the value of  $j_1$  is fixed on the maximal possible value in the ODWT model:  $j_1 = \log_2 n - 1$ . As we shall see, this turns out to be one major weakness of both of these methods when estimating continuous functions with fractal graphs and spatially inhomogeneous functions with relatively low regularity, in particular, piecewise smooth functions with discontinuities. We suggest that in all cases (i.e., when estimating fractal functions, functions with jumps, and even smooth functions, for both the regression and density estimation models)  $j_1$  should be a controlled parameter and that  $2^{j_1} = o(n)$  should hold. The choice  $2^{j_1} \asymp n$  can be acceptable for very smooth functions only. For such very smooth functions, however, thresholding techniques seem to offer quite a competitive alternative to (7) for consistent estimation of the function together with its derivatives. Thus, the methods of Antoniadis and of Amato and Vuza, which are not well adapted for work with small values of the smoothness index  $s$ , essentially miss the most important range of  $s$  for which non-threshold shrinking methods outperform threshold ones.

Considerations related to asymptotic minimax optimality theory (see [53]) suggest  $(\tilde{C}n)^{1/(2s+1)} \leq 2^{j_0} \leq 2(\tilde{C}n)^{1/(2s+1)}$ , where  $\tilde{C} = \tilde{C}(L)$  is defined via (16) and (21) in [53] for density and for regression-function estimation, respectively. (The reader is cautioned that there is a misprint in formula (16) in [53].) Here  $L$  is an *upper bound* for the  $B_{22}^s$ -norm of the estimated function  $f$ . Again by (16) and (21) in the same paper,  $n/\ln n \leq 2^{j_1} \leq 2n/\ln n$ .

**Remark 1.**<sup>A15</sup> The proofs of the theorems in Section 6 (see Appendix A) show that the cross-validation method for estimating  $t$  is consistent, for a broad range of  $s$ , including small values, under relaxed assumptions:  $2^{j_1} = o(n)$  (for both the regression and density case);  $2^{j_0} = O(1)$  or  $j_0 \rightarrow \infty$ ,  $j_0 \leq j_1$  for the regression case,  $2^{-j_0} = O(n^{-\frac{1}{2(1+s)}})$  in the density case. If no additional information about the function is available, we recommend the choice  $2^{j_1} \asymp \frac{n}{\ln n}$  because the estimator performs well for both smooth and piecewise smooth functions with isolated points of singularity. The type of 1-d singularity may

be, e.g., a "horn" or a "cusp" (discontinuity of the first derivative, which in the case of a cusp point is unbounded in a neighbourhood of this point), a "jump" (the left-hand and right-hand limits of the function are different at this point), or a "chirp" (the function does not have a left-hand and/or a right hand limit and its variation is unbounded in a neighbourhood of this point). If additional information is available that  $f$  is smooth, then choices of the type, say,  $2^{j_1} = \frac{n}{\ln n}$  are acceptable. If further additional information is available that  $f$  is not very spatially inhomogeneous and  $s$  is 'large' (at least  $s > 1$  but typically  $s \gg 1$ ), then  $2^{j_1} \asymp n$  is also an admissible choice. If, on the contrary, it is known additionally that  $s \leq 1/2$  and that the function is a 'monstrous' one (e.g., has a countable set of discontinuities with at least one density point), then choices  $2^{j_1} = o(\frac{n}{\ln n})$  should be considered. The option to control the performance of the estimator by varying  $j_1$  (that is, by applying a non-threshold shrinking approach via (7) in combination with a levelwise 'kill-or-keep' thresholding) is one realization of a simple partial case of the general idea about composite shrinking/thresholding estimators, considered in <sup>B8</sup>.

## 6. Cross validation in $L_2$

Besides asymptotic-minimax methods, there is also a variety of other methods for estimation of the smoothing parameter(s) in penalized estimation (e.g., the method of [9], SURE and James-Stein estimation for Gaussian white noise, etc.) but one of the oldest and most general methods, working also for non-Gaussian white noise, available also in the case of density estimation, is cross validation. Its theory poses serious challenges and asymptotic results hold sometimes in a weaker sense than in asymptotic minimax theory, but its performance on moderate samples is good most of the time. This method copes well even with correlated noise, at least when the correlation matrix is band-limited or its entries tend to vanish relatively fast away from the main diagonal. That is why, of all methods for determining the smoothness parameter in the general setting of (6), here we would like to consider in detail cross validation as a model example. In principle, cross-validation procedures can be formulated



for non-Hilbertian loss functions, too, but the resulting optimization problems are easier to analyze in the Hilbert setting. For this reason, it is desirable to study the Hilbert case in quite a detail, before looking at more general situations. This is going to be the direction of our research in this section, where we concentrate on work with the non-threshold shrinking estimator (7), adapted for the Hilbert case. For the generalization of (7) to the general quasi-Banach case, see Appendix B, formulae (B4-B6). Section 6 contains a more detailed and updated discussion of some earlier results (see [49]), as well as new results on consistency of the estimator and existence and asymptotic behaviour of the smoothing parameter. Subsection 6.1 includes also a comparison of our method for the regression case with the approach of Amato and Vuza [3,5,6] and the approach of [9].

Cross validation (CV) is performed in a different way for regression estimation and for density estimation. The reason for the difference is that in the latter case there are no output observations  $Y_i$ ,  $i = 1, 2, \dots, n$  that could be utilized to evaluate the predictive power of the model. Another important difference is that in the latter case  $\text{MISE}(v) = E \int (\tilde{f}_v(x) - f(x))^2 dx$  is being directly estimated, while in the former one the estimation is of a quadrature formula approximating  $\text{MISE}(v)$ . However, there is an important common feature between the two models: the dependence of the CV-criteria for both models is analytic in  $v = t^2$  on the compact  $[0, C]$  for any  $C : 0 < C < \infty$ . This is an essential difference between the shrinking and the thresholded wavelet estimator (cf. [79]). Because of it, existence and asymptotic behaviour of the  $v$ -minimizer of the CV-functional for the shrinking estimator can be studied in quite a detail. In general, consistency of the cross-validation procedure in  $L_2$  does not automatically imply consistency of the resulting estimator  $\tilde{f}_v$  in  $L_2$ ; for this to happen, the  $v$ -minimizer must have a specific asymptotic behaviour:  $v \rightarrow 0$  as  $n \rightarrow \infty$ , at least when  $j_0 = O(1)$ . If this is not fulfilled, the resulting estimator  $\tilde{f}_v$  cannot be consistent in  $L_2$ , because  $\beta_{jk}$  would not be consistently estimated. For an arbitrary wavelet-based estimator, consistent estimation of the coefficients in the wavelet expansion of  $f$  is only a very weak necessary condition for consistency of

the estimator itself: it is not sufficient for the estimator's consistency even in the weak topology on  $L_2$ . We show that the shrinking wavelet estimator  $\tilde{f}_v$  enjoys a remarkable property: when  $j_0 = O(1)$ , consistent estimation of the wavelet coefficients is for it not only necessary, but also *sufficient* for consistency of  $\tilde{f}_v$  in  $L_2$  and even in  $B_{22}^\sigma$ ,  $0 < \sigma < s$ . The range of  $\sigma$  depends on  $s$  and the choice of  $j_0$  and  $j_1$ .

The wavelet setting permits also explicit characterization of all functions for which there is no  $v$ -minimizer of the CV-criterion for arbitrarily large  $n$ . It turns out that any such function must be in  $V_j$  - the closure of the linear span of  $\varphi_{jk}$ ,  $k \in \mathbb{Z}$ , for some  $j \in \mathbb{Z} : j < j_0$ .

We shall be utilizing the homogeneous  $K$ -functional  $K_2(\sqrt{v}, f; L_2, \dot{B}_{22}^s)$  throughout, with  $f \in B_{22}^{s'}$ ,  $0 < s' < s$ . Under these conditions,  $B_{22}^{s'} \hookrightarrow \dot{B}_{22}^{s'} \hookrightarrow L_2 + \dot{B}_{22}^s$  (see also Section 4). For  $j_0$  we shall be considering two cases:  $j_0 = O(1)$  and  $j_0 \rightarrow \infty$ . In the proofs we shall be giving the details for the second case. In the first case the proof is analogous, but simpler. To obtain quantitative results about the second case we shall be considering only those  $f \in B_{22}^{s'}$  for which the index  $s'$  is sharp, that is,  $f \notin B_{2\infty}^s$  for any  $s > s'$ . Note that, for any compactly supported  $f \in B_{22}^{s_0} \setminus C^\infty$  for some  $s_0$ , there exists  $s' \geq s_0$  such that  $f \in B_{22}^{s'}$  and  $s'$  is sharp.<sup>A16</sup> The case  $f \in C^\infty$  is possible to incorporate, of course, but in  $K$ -functional context it is not very interesting.

For the purposes of estimation in regression problems, when  $f$  has low regularity ( $0 < s' \leq 1/2$ ), we discuss in subsection 6.1 the importance of the concept of  $A$ -spaces and *average moduli of smoothness* <sup>B12b</sup>.

Throughout the section we implicitly assume that  $f$  is compactly supported, as discussed in the end of Section 3.

### 6.1. Regression-function estimation

Cross validation is usually numerically intensive unless there are some updating formulae that allow to calculate the "leaving-out-one" predictions on the basis of the "full" predictions only. In this respect, very helpful is the "leaving-out-one" lemma (Lemma 4.2.1 of [117]) which allows deriving such updating

formulae. This lemma's conclusion can only be achieved under the so-called "compatibility condition" (see also [60]). Although the compatibility condition is easy to derive for projection-type estimators, it might be difficult to prove in other situations. In fact, it very often fails to hold. This is also the case with our estimator (1'). Remedy from this situation is to pretend that the compatibility condition "almost holds" and work further by bearing into consequence the imprecision that occurs due to the approximation error. This approach to cross-validation in wavelet regression has been adopted, in threshold setting, by Jansen, Malfait and Bultheel [79]. Similar is the situation with the cross-validation method of Nason [92] which is based on an ingenious exploitation of the symmetries existing only when  $n$  is a power of two. This method essentially relies on the brute force of the adopted "leave-out-half" strategy. At present, it is hard to say what price, in terms of, say, asymptotic rates, is being paid for leaving out half of the sample in this context, because this estimator has hardly any theory developed for it, partly due to its non-smooth optimization criterion. (We discuss Nason's estimator in more detail in  $B^{20}$ .) A similar approach to overcoming the failure of the compatibility condition is that of Amato and Vuza [5,6]. Since for the estimator (1') with (7) the compatibility condition is violated, they have chosen to work directly with a wavelet analogue of formula (4.3.1) in [117] for  $n = 2^N$  and, by doing so, have essentially adopted the same approach as of [79], but in the setting of (7). On its part, the method of [9], which is not of cross-validation type, completely avoids the compatibility problem. One big price to be paid for this is that the estimator depends explicitly on the noise variance. To deal with this problem, Antoniadis incorporates an estimator of the noise variance within the definition of his estimator, but this comes at the price of spoiling the performance of the estimator on less regular functions. This is one of several reasons why this estimator is only appropriate for estimating smooth functions.

Here we suggest to follow another possibility which leads to a stringent solution of the compatibility problem within the cross-validation setting, namely, to modify the cross-validation criterion itself by replacing it with the so-called

full cross-validation criterion (FCV) (see [60]).

Denote

$$\tilde{f}_v(x) = \sum_k \hat{\alpha}_{j_0 k} \varphi_{j_0 k}(x) + \sum_{j=j_0}^{j_1} \sum_k \hat{\beta}_{jk} \frac{1}{1 + v 2^{2js}} \psi_{jk}(x),$$

$$\tilde{f}_v^{(-i)}(x) = \sum_k \hat{\alpha}_{j_0 k(-i)} \varphi_{j_0 k}(x) + \sum_{j=j_0}^{j_1} \sum_k \hat{\beta}_{jk(-i)} \frac{1}{1 + v 2^{2js}} \psi_{jk}(x),$$

$$\hat{\alpha}_{jk(-i)} = \sum_{l=1, l \neq i}^n \frac{a_2 - a_1}{n - 1} \varphi_{jk}(x_l) Y_l, \quad \hat{\beta}_{jk(-i)} = \sum_{l=1, l \neq i}^n \frac{a_2 - a_1}{(n - 1)(1 + t^2 2^{2js})} \psi_{jk}(x_l) Y_l.$$

(For the sake of convenient comparison with our main reference source [117] here, the notation  $y_i := Y_i$  will be used henceforward.) Then, the typical *ordinary cross-validation* (OCV) functional to be minimized with respect to  $t^2 = v$  is given by ([117], p. 47):

$$(8) \quad V_0(v) = \frac{1}{n} \sum_{i=1}^n [y_i - \tilde{f}_v^{(-i)}(x_i)]^2$$

Computational reduction in minimizing (8) is achieved if the "compatibility condition"

$$(9) \quad \hat{y}_{(-i)}(v) = \tilde{f}_{n,v}(x_i; y_1, \dots, y_{i-1}, \hat{y}_{(-i)}(v), y_{i+1}, \dots, y_n)$$

holds, where

$$\hat{y}_i = \tilde{f}_v(x_i), \quad \hat{y}_{(-i)}(v) = \tilde{f}_v^{(-i)}(x_i) := \tilde{f}_{n-1,v}(x_i; y_1, \dots, y_{i-1}, y_{i+1}, \dots, y_n).$$

It is easy to show that, under the above condition, the OCV-functional can be expressed in terms of the ordinary residuals:

$$(10) \quad V_0(v) = \frac{1}{n} \sum_{i=1}^n [y_i - \tilde{f}_v(x_i)]^2 / (1 - h_{ii}(v))^2, \quad h_{ii}(v) = \frac{\partial \tilde{f}_v}{\partial y_i}(x_i; y_1, \dots, y_n)$$

To explain our notation, let us note that because of the linearity of our estimator, we can write  $\hat{\mathbf{y}}(v) = \mathbf{H}(v)\mathbf{y}$ ,  $\hat{\mathbf{y}}(v) = [\hat{y}_1(v), \hat{y}_2(v), \dots, \hat{y}_n(v)]'$ ,  $\mathbf{y} = [y_1, y_2, \dots, y_n]'$  with certain  $n \times n$  matrix  $\mathbf{H}(v)$  (called *influence matrix* (see [117])). In that case,  $\frac{\partial \tilde{f}_v}{\partial y_l}(x_i; y_1, \dots, y_n)$  would indeed be the  $(i, l)$ -th element of the matrix  $\mathbf{H}(v)$ .

Unfortunately, for a shrinking-type estimator like (1'), the compatibility condition is violated and a formula like (10) does not hold precisely. Simple observation shows that (10) still can be considered as approximately true, i.e.,  $h_{il}(v)$  can still be calculated as  $h_{il}(v) \approx \frac{\partial \tilde{f}_v}{\partial y_l}(x_i; y_1, \dots, y_n)$ , but the error of this approximation has yet to be evaluated. We opt for another alternative here by changing the cross-validation criterion itself. Note that standard cross-validation is defined entirely with respect to samples of size  $n - 1$ . One can adjust it for samples of size  $n$  as suggested in [23]. In the approach of [23], the value of  $y_i$  is replaced by  $\hat{y}_i(v)$  instead of leaving it out in defining the prediction of the  $i$ -th design point. The resulting functional is called the *FCV* (*full cross-validation*) functional:

$$(11) \quad FCV(v) = \frac{1}{n} \sum_{i=1}^n [y_i - \tilde{y}_i(v)]^2$$

with  $\tilde{y}_i(v) = \tilde{f}_{n,v}(x_i; y_1, \dots, y_{i-1}, \hat{y}_i(v), y_{i+1}, \dots, y_n)$ . Then it turns out that under the condition of *linearity only*, one gets in terms of the ordinary residuals:

$$(12) \quad FCV(v) = \frac{1}{n} \sum_{i=1}^n [y_i - \tilde{f}_v(x_i)]^2 \cdot (1 + h_{ii}(v))^2, \quad h_{il}(v) = \frac{\partial \tilde{f}_v}{\partial y_l}(x_i; y_1, \dots, y_n)$$

An easy calculation gives the following value for  $h_{il}(v)$ :

$$(13) \quad h_{il}(v) = \frac{a_2 - a_1}{n} \left[ \sum_k \varphi_{j_0 k}(x_i) \varphi_{j_0 k}(x_l) + \sum_{j=j_0}^{j_1} \sum_k \frac{1}{1 + v 2^{2js}} \psi_{jk}(x_i) \psi_{jk}(x_l) \right],$$

and the matrix  $\mathbf{H}(v)$  is symmetric positive semi-definite.

As it is often done in the cross-validation approach, OCV is replaced by *generalized cross validation* (GCV) (see [117], Section 4.3). The reason for this is that the generalized version has desirable invariance properties that do not generally hold for OCV. Numerically, going over from the ordinary to the generalized version means that one replaces the values of  $h_{ii}(v)$ ,  $i = 1, 2, \dots, n$  in (10) by one value only:  $\bar{h}(v) = \frac{1}{n} \sum_{i=1}^n h_{ii}(v) = \frac{1}{n} \text{tr } \mathbf{H}(v)$ . Also, proofs related to consistency of the cross-validation procedure are easier to derive for generalized, rather than ordinary cross-validation. Following the same idea, we introduce *generalized FCV* (GFCV) by replacing  $h_{ii}(v)$ ,  $i = 1, 2, \dots, n$  in (12) with  $\bar{h}(v)$ .

It is informative here to compare our approach to that of Amato and Vuza [5], who have tried to develop a straightforward analogue of the cross-validation procedure, as suggested by Wahba [117]. Because the compatibility condition fails for OCV when the shrinking is via (7), an analogue of our formula (12) and of formula (4.2.13) in [117] does not hold precisely. However, one can hope, in the spirit of [79], that for OCV such an analogous formula holds at least approximately. Under this assumption, following Wahba, it is possible to consider a generalized version of the approximate formula for OCV by replacing  $h_{ii}$  with  $\bar{h}$ . Following this line of reasoning for the case  $n = 2^N$ , and making use of the discrete wavelet transform, Amato and Vuza have proposed a method which they call generalized cross validation (GCV). In the sequel we shall refer to their method as quasi-GCV (or QGCV, for short) because it is not an exact generalization of OCV.

In the rest of this subsection we shall study the properties of GFCV and shall compare them to the respective properties of the existing non-threshold shrinking techniques QGCV of Amato and Vuza and WAVREG of [9].

We are in a position to show that GFCV is a consistent method in the traditional sense. To formulate precisely our results, let us introduce the mean square error  $ER(v) = E \frac{1}{n} \sum_{i=1}^n (\tilde{f}_v(x_i) - f(x_i))^2$  (the unknown quantity we are interested in controlling) and compare it to  $E(GFCV(v)) = E \frac{1}{n} \sum_{i=1}^n \{y_i -$

$\tilde{f}_v(x_i)\}^2(1 + \bar{h}(v))^2$ . Further, let us define  $\mu(v) = \frac{1}{n} \text{tr } \mathbf{H}'(v)\mathbf{H}(v) = \frac{1}{n} \text{tr } \mathbf{H}(v)^2$ . It can be seen that  $\mu(v) \geq \bar{h}(v)^2$  holds. ( $\bar{h}(v)$  and  $\mu(v)$  are also functions of  $n$  but this dependence has been suppressed in the notation.) We denote asymptotic equivalence by ' $\sim$ '.

**Theorem 1 .** (*Sufficient condition for existence of the  $v$ -minimizer.*) Let  $0 < s_1 \leq s' < s < r$ ,  $s_1 < 1$ . Assume that the penalized model is via  $K_2(\sqrt{v}, f; L_2, \dot{B}_{22}^{s'})$  and that  $f \in B_{22}^{s'} \cap B_{\infty\infty}^{s_1}$ . Assume also that  $j_0 \leq j_1$ ,  $j_1 \rightarrow \infty$  with  $2^{j_1} = O(\frac{n}{\ln n})$ . If  $j_0 = O(1)$  and  $f \notin V_j$  for any  $j < j_0$ , or if  $j_0 \rightarrow \infty$ ,  $f \notin V_j$  for any  $j \in \mathbb{Z}$ , the index  $s'$  is sharp and  $2^{j_0} = o(n^{\min(\frac{s_1}{(s-1/2)_+}, \frac{2}{1+2s})})$ , then, for sufficiently large  $n$  there exist  $\tilde{v} : E(\text{GFCV}(\tilde{v})) = \min_v E(\text{GFCV}(v))$ ,  $v^* : ER(v^*) = \min_v ER(v)$ , and  $0 \leq \tilde{v} < C$ ,  $0 \leq v^* < C$ , where  $C < \infty$  depends on  $f$  and  $\psi$  only.

The problem about existence of the minimizer has not been explored for QGCV and WAVREG and Theorem 1 has no analogue in the theory developed for these estimators.

Asymptotic optimality of the minimizer of the smoothing parameter is established by the following

**Theorem 2 .** (*Consistency of GFCV - cf. also [49].*) Assume that  $0 < s' \leq s < r$ , the penalized model is via  $K_2(\sqrt{v}, f; L_2, \dot{B}_{22}^{s'})$  and  $f \in B_{22}^{s'}$ . Assume also that  $j_0 \leq j_1$ ,  $j_1 \rightarrow \infty$  and  $2^{j_1} = o(n)$  holds. Then, for any  $v \geq 0$ ,  $E(\text{GFCV}(v)) \sim ER(v) + \delta^2$  holds,  $\delta^2$  being the error variance. If  $\tilde{v}$  and  $v^*$  exist for sufficiently large  $n$ , then  $\frac{ER(\tilde{v})}{ER(v^*)} \searrow 1$ , as  $n \rightarrow \infty$ .

Note that the asymptotic optimality of the minimizer  $\tilde{v}$  of GFCV is achieved without knowing the value of  $\delta$  or the type of the white noise. The respective result in [5], p.486, is also for unknown  $\delta$ , and is valid for  $s > 1/2$ . The respective results in [9] (Corollary 3.2 for the case of known noise variance  $\delta^2$  and Theorem 3.3 for the general case ( $\delta$  unknown)) are weaker: they describe the behaviour of WAVREG's  $v$ -minimizer only by measuring the relative error with respect to the asymptotic minimax rate  $n^{-\frac{2s}{1+2s}}$  instead of the relative error

with respect to the true value of the estimated quantity.

**Theorem 3 .** (*Consistency of  $\tilde{f}_v^*$  in  $B_{22}^\sigma$ ,  $0 \leq \sigma < s' \leq s < r$ .) Assume that the conditions of Theorem 1 hold, without the sharpness conditions on  $s'$ . Then, if  $j_0 = O(1)$  and if  $v = v_n$  is bounded, the condition  $v \rightarrow 0$  as  $n \rightarrow \infty$  is necessary for  $E\|\tilde{f}_v - f\|_{B_{22}^\sigma}^2 \rightarrow 0$  to hold; moreover, in both cases  $j_0 = O(1)$  and  $j_0 \rightarrow \infty$ , if  $2^{j_1} = o(n^{\min(\frac{s'}{\sigma}, \frac{1}{1+2\sigma})})$  holds, then  $v \rightarrow 0$  is also a sufficient condition for  $E\|\tilde{f}_v - f\|_{B_{22}^\sigma}^2 \rightarrow 0$ . In particular, if  $s'$  is sharp, then: the above necessity claim is true for  $v = v_n^*$ ; if  $2^{j_1} = o(n^{\min(\frac{s'}{\sigma}, \frac{1}{1+2\sigma})})$  holds, too, the above sufficiency claim is also true for  $v = v_n^*$ .*

Theorem 3 has no analogue for QGCV and WAVREG. One important reason for this is that in those two techniques  $j_1$  is fixed at  $2^{j_1} \asymp n$ . It should be noted, though, that if  $v \rightarrow 0$  and if the rate of  $v$  tending to zero is known, the restriction on  $j_1$  in Theorem 3 can be relaxed and, if the rate is high enough,  $2^{j_1} \asymp n$  becomes an admissible choice, too.

**Theorem 4 .** (*Asymptotic behaviour of  $v^*$  and  $\tilde{v}$ .) Under the premises of Theorem 1,  $v^* \rightarrow 0$  and  $\tilde{v} \rightarrow 0$  as  $n \rightarrow \infty$ .*

Theorem 4 has no analogue for WAVREG. For QGCV, Amato and Vuza study only the case about  $\tilde{v}$ , for  $j_0$  and  $j_1$  fixed at 0 and  $\log_2 n - 1$ , respectively, where  $n$  is a power of two. The restriction on  $s$  is  $s > 1/2$  (the range where QGCV has been shown to be consistent), and  $s'$  is set equal to  $s$ . Under these assumptions Amato and Vuza [6] find also the sharp quantitative rates for  $\tilde{v}$ . In our case, for GFCV, the sharp quantitative rates for  $v^* \rightarrow 0$  and  $\tilde{v} \rightarrow 0$  can be found in terms of  $s'$ ,  $s$ ,  $j_0$  and  $j_1$  for any  $j_0 \leq j_1$ ,  $j_1 \rightarrow \infty$ , and any  $s'$  and  $s$ , such that  $0 < s' \leq s < r$ . One of the interesting extensions of this topic is to analyze for which pairs  $(j_0, j_1)$  the expected rate for  $v^*$  and  $\tilde{v}$  is the asymptotic minimax one, i.e.,  $n^{-2s'/(1+2s')}$  (see also  $B^{14}$ ).

**Corollary 1 .** (*Consistency of  $\tilde{f}_{\tilde{v}}$  in  $B_{22}^\sigma$ ,  $0 \leq \sigma < s' \leq s < r$ .) Under the premises of Theorem 1,  $MISE(\tilde{v}) = E\|\tilde{f}_{\tilde{v}} - f\|_{L_2}^2 \rightarrow 0$  holds as  $n \rightarrow \infty$ .*



If, additionally,  $2^{j_1} = o(n^{\min(\frac{s_1}{\sigma}, \frac{1}{1+2\sigma})})$  is fulfilled, then also  $E\|\tilde{f}_{\tilde{v}} - f\|_{B_{22}^s}^2 \rightarrow 0$  holds.

The respective result for WAVREG ([9], Theorem 3.2) is valid in a weaker sense: convergence is in probability, for  $\sigma = 0$  only, under stronger restrictions on the white noise. In [6] an analogue of Corollary 1 is proved for QGCV when  $j_0 = 0$  and  $2^{j_1} = n - 1$ , by a detailed study of the convergence rate for  $\tilde{v}$  in their model. An analogous study for GFCV <sup>B14</sup> additionally allows to determine for which pairs  $(j_0, j_1)$  the rate of  $E\|\tilde{f}_v - f\|_{B_{22}^s}^2 \rightarrow 0$  (with  $v = E\tilde{v}$  or  $v = Ev^*$ ) is asymptotically minimax-optimal.

**Remark 2.** If  $f \in V_j$  for some  $j < j_0$ , then  $\{v_n^*\}$  and  $\{\tilde{v}_n\}$  can have more than one density point in  $[0, \infty)$ . It is also possible for the infima of  $\text{GFCV}(v)$  and  $\text{ER}(v)$  to get attained for  $v \rightarrow \infty$  for infinitely many  $n$ . In all these cases the estimator  $\tilde{f}_v$  with the cross-validated choice of  $v$  can still be consistent in  $B_{22}^s$ . It is important that for the *inhomogeneous* model (see Sections 4, 5) asymptotic behaviour  $v \rightarrow 0$  is guaranteed (for sharp  $s'$ ) in all cases except  $f \equiv 0$  a.e.

**Remark 3.** Let us summarize the results of the comparison of GFCV, QGCV and WAVREG, all of them based on non-threshold shrinking via (7). (a) In the theory of QGCV and WAVREG it is assumed that  $s' = s$ ; the considerations in our Section 4 show that  $s' = s$  is in fact only a limiting case, while the true relation between  $s'$  and  $s$  is  $0 \leq \sigma < s' \leq s$ . Therefore, in order to apply penalized shrinking, the true index  $s'$  of the estimated function does not have to be known exactly, but only approximately, which is, of course, a very important improvement of the penalization model. This also suggests that the formulations and proofs of all theorems in the theory of QGCV and WAVREG need to be modified so as to hold in the general case  $s' \leq s$ . Moreover, we believe that some of the key results in [9] can be improved considerably even in the case  $s' = s$ . (b) The three methods yield comparable results for smooth functions ( $s > 1$ ). The high order of vanishing moments of the coiflets used in WAVREG reduces the bias considerably when  $s \gg 1$ , but the variance term is generally

larger than that of GFCV (with  $2^{j_1} = o(n)$ ) and comparable to that of QGCV, which makes the use of coiflets efficient only for very small noise variance and very spatially homogeneous smooth curves. The performance of WAVREG on very smooth but spatially inhomogeneous curves can be visibly improved if  $j_1$  is reduced from its present value  $\log_2 n - 1$  to  $o(\log_2 n)$  because this would reduce the contribution of the variance term to the overall error of the method. In its present version proposed in [9], the distance between  $j_0$  and  $j_1$  is simply too large for a single smoothing parameter to handle for moderate sample size and a spatially inhomogeneous curve, even if the curve is very smooth. For the case of non-smooth functions ( $s \leq 1$ , especially  $s \leq 1/2$ ) GFCV is an easy winner over QGCV and both GFCV and QGCV are easy winners over WAVREG. The most important (though not the only) reason for the domination of GFCV in this range of the smoothness index is quite simple: for QGCV and WAVREG  $2^{j_1} \asymp n$  while for GFCV  $2^{j_1} = o(n)$  holds (see also Remark 4 below). If QGCV and WAVREG be modified so that  $2^{j_1} = o(n)$  then these estimators, especially QGCV, would become more competitive for small  $s$ . Let us support this claim by a more detailed comparison between GFCV and each of QGCV and WAVREG. In order to compare GFCV and QGCV, assume for a moment the (very unfavourable for GFCV) setting  $j_1 = \log_2 n - 1$  and  $j_0 = 0$ . Assume also that  $n = 2^N$  and ODWT is applied for both GFCV and QGCV. Under these very favourable for QGCV conditions, it is easy to compute

$$\bar{h}(v) = \frac{1}{n} \left( 1 + \sum_{j=0}^{j_1} \frac{2^j}{1 + v 2^{2js}} \right), \quad j_1 = \log_2 n - 1,$$

$$QGCV_s(v) = \frac{1}{(1 - \bar{h}(v)^2)^2} GFCV_s(v).$$

Taking in consideration that the asymptotic behaviour of the  $v$ -minimizer for both criteria is  $v \rightarrow 0$  as  $n \rightarrow 0$ , and that, clearly,  $\bar{h}(v) \rightarrow 1$  as  $v \rightarrow 0$ , it is seen that, with the increase of  $n$ ,  $QGCV_s$  becomes increasingly numerically unstable and the problem for optimization of  $v$  becomes increasingly ill-conditioned. It is also seen that when  $s$  decreases the instability problems with QGCV become

ever more severe, which means that the smaller  $s$ , the worse the performance of  $QGCV_s$  is, compared to that of  $GFCV_s$ . By replacing the sequence  $\{2^j\}$  in  $QGCV$  by another increasing sequence  $\{a_j\}$  (which is an additional justification for calling this method *quasi-GCV*), Amato and Vuza have reduced the instability problems, because for an appropriate choice of  $a_j$   $QGCV$  remains bounded as  $n \rightarrow \infty$ . However, even in this form,  $\bar{h}(v) > 0$  remains bounded away from zero as  $n \rightarrow \infty$  which means that the variability of  $QGCV$  is asymptotically bigger than that of  $GFCV$ . Returning to the present setting for  $j_1$  ( $2^{j_1} \asymp n$  for  $QGCV$ ,  $2^{j_1} = o(n)$  for  $GFCV$ ) leads to total domination of  $GFCV$ . If we now improve  $QGCV$  by resetting  $j_1$  in its definition to be such that  $2^{j_1} = o(n)$  holds, then  $\bar{h}(v) \rightarrow 0$  as  $n \rightarrow \infty$ , and the two methods become asymptotically equivalent. However, the variability of  $QGCV$  is still larger than that of  $GFCV$  (cf. also [60]), and, when  $s$  is small,  $GFCV$  performs better than the improved  $QGCV$  for moderate samples. If  $s$  is large, the performance of  $GFCV$  and the improved  $QGCV$  is very similar already for moderate sample sizes. These theoretical observations were confirmed by our numerical studies. With the increase of the noise variance the superiority of  $GFCV$  becomes more visible.  $GFCV$  is also more robust against outliers in the data. Summing up,  $GFCV$  can be considered as a stabilized version of the improved  $QGCV$ . Compared to the original version of  $QGCV$ ,  $GFCV$  is visibly superior not only for moderate, but also for relatively large sample sizes. In brief, we feel safer, both theoretically and computationally, when applying the  $FCV$  ( $GFCV$ ) approach rather than the  $OCV$  ( $GCV$ ) type approach. In fact, due to the above-mentioned reasons, we recommend it for Wahba's spline-smoothing model, too.

Comparing  $GFCV$  to  $WAVREG$ , we note that, in contrast to most wavelet estimation methods, *no estimator of the noise variance* is involved in the definition of our estimator. This is a very essential advantage: the definition of the estimator does not depend on the structure and parameters of the noise. On the contrary,  $WAVREG$  is a typical estimator involving noise variance. The explicit dependence of  $WAVREG$  (via the estimator  $\hat{\Lambda}_n$  ([9], p.323)) on the noise

variance is eliminated by passing on to  $\tilde{A}_n$  ([9], p.323), but only at the price of additional worsening of the performance of WAVREG for non-smooth functions. There is also an important error in Lemma 3.1 of [8,9]. We shall be discussing the consequences of this error in more detail in the end of this subsection.

In the proof of Theorem 2 the following lemma is of crucial importance.<sup>A17</sup>

**Lemma 2 .** (cf. [49].) Assume that for the estimator (1')  $j_0 \leq j_1$ ,  $j_0 = O(1)$  and  $2^{j_1} = o(n)$  hold as  $n \rightarrow \infty$ . Then, for any  $v \geq 0$ ,

$$(14) \quad \lim_{n \rightarrow \infty} \frac{\overline{h}(v)^2}{\mu(v)} = 0.$$

**Remark 4.** An analogue of (14) is crucial for the consistency of  $L_2$ -cross validation for any estimator - spline- or wavelet-based, thresholded or shrinking. Essentially, in (14) two invariants of  $\mathbf{H}(v)$  are compared. They are two distinct Schatten - von Neumann  $S_\gamma$ -norms of the linear compact finite-rank operators with matrices  $\mathbf{H}(v) = \mathbf{H}_n(v)$ ,  $n \in N$  (see [51]). It is the choice of the  $L_2$ -metric that determines the two  $S_\gamma$ -norms to be the Hilbert-Schmidt norm ( $\gamma = 2$ ) and the Fredholm (or nuclear) norm ( $\gamma = 1$ ) (for other metrics, see <sup>B1</sup>). Our proof is essentially different from the "classical" approach to proving (14) for the spline-smoothing model (see [117], pp. 57-58, for a short outline). The latter can only be applied under additional restrictive assumptions about the  $s$ -values of  $\mathbf{H}(v)$  ([114], p. 204, conditions (i) and (ii); [35], p. 701, Remark after Theorem 2.4; [117], pp. 57-59). The same is true for the "classical" approach to proving (14) adopted in wavelet context in [79] for thresholding via cross validation. The nature of the thresholding procedure requires that in this case - unlike ours -  $\mathbf{H}(v)$  also depends on the estimated function  $f$  via the sample  $\mathbf{y}$ . In this case the additional restrictions on the  $s$ -numbers of  $\mathbf{H}(v)$  result in an essential restriction on the regularity of the estimated  $f$ , formulated in Assumption 1 of [79]. This assumption virtually excludes from consideration continuous  $f$  with fractal graphs and also very spatially inhomogeneous  $f \in C^\infty$ . The wavelet version of the Riemann-Lebesgue Lemma also implies that Assumption 1 fails for

piecewise smooth  $f$  with discontinuities. The "classical" approach *ignores the key fact on which our proof relies*:  $2^{j_1} = o(n)$ . In the "classical" setting there is no parameter  $j_1$  to control, or it is ignored, which in our case is equivalent to the implicit assumption  $2^{j_1} = O(n)$ . Under the latter condition on  $j_1$ , we would also need restrictive assumptions on the  $s$ -numbers in order to prove (14). The above consideration shows that it is the *frequency localization* which is crucial here: it is the one that provides the controllable parameter  $j_1$ . As a consequence, the new approach to proving (14) can in principle be adapted to relax the restrictions on the  $s$ -numbers of  $\mathbf{H}(v)$  by a controlled truncation of the highest frequencies also for a rather wide subclass of spline-smoothing estimators, thereby improving their performance for less regular curves. Summing up: the non-threshold shrinking wavelet estimator with  $2^{j_1} = o(n)$  is expected to perform better than smoothing-spline and threshold-wavelet estimators, and also better than non-threshold shrinking estimators based on (7) for which  $2^{j_1} \asymp n$ , when the estimated  $f$  has low regularity. In this aspect, its main advantage to standard smoothing splines is the simultaneous spatial and frequency localization; its main advantage to thresholded-wavelet estimators is that its influence matrix  $\mathbf{H}$  does not depend on  $f$  or the sample  $\mathbf{y}$ . One of its main advantages to non-threshold estimators based on (7) with fixed  $j_1 : 2^{j_1} \asymp n$  is essentially the same as with the classical smoothing-spline technique, namely, with GFCV  $j_1$  is also a controlled parameter and is set at a lower value, with  $2^{j_1} = o(n)$  (see also Remark 3 (b)).

**Remark 5.** In Theorem 2 we considered the case of fixed  $s$ . This was done mainly for simplicity of the exposition. In practice, one would consider minimizing the FCV (or GFCV) functional as a function of both parameters  $s$  and  $v$ , together with the discrete parameter  $j_0$ . For both the homogeneous and inhomogeneous model, if  $f \in V_j$  for some  $j \in Z$ , then for  $n$  large enough the expected cross-validated value for  $j_0$  is  $j_0 = j$  (see also Remark 7).

Finally, let us note that the restriction  $s > 1/2$  in the analogues of our Theorem 2 for QGCV and WAVREG essentially eliminates discontinuous

functions  $f$  from consideration, because curves  $f$  which have at least one jump point do not belong to  $B_{22}^s$  for any  $s > 1/2$ . In this aspect, the performance of WAVREG is additionally worsened by the necessity to estimate the noise variance. It should be noted that there is a common error in the following results: [8], Lemma 3.1; [9], Lemma 3.1, Theorem 3.1, Theorem 3.2 and Corollary 3.2; [5], Theorem 3; [6], Theorem 1; in all these results it is assumed that  $s > 1/4$ , while the correct range for  $s$  in all of them is  $s > 1/2$ . The error consists in the assumption that the error of the quadrature formula admits the bound  $|\beta_{jk} - E\hat{\beta}_{jk}| = O(2^{-sj})$  for any  $f \in B_{22}^s$ , for any  $s > 0$ , thus including also the range  $0 < s \leq 1/2$ . In fact, the correct range of  $s$  for which the above bound holds is  $s > 1/2$  (see [56], section 2.3, pp. 11-12, section 2.4, p.13, section 6.1, pp. 29-30). A counterexample is given in [77]. For further discussion, see  $B^{12(b)}$ .

What about FCV and GFCV? Theorem 2 is clearly valid also for discontinuous functions. However, the other results in this subsection are valid under the constraint  $f \in B_{\infty\infty}^{s_1}$ ,  $s_1 > 0$ , which also eliminates discontinuous functions. That is why it is important to note here that in all these theorems the assumption  $f \in B_{22}^{s'} \cap B_{\infty\infty}^{s_1}$  can be relaxed by replacing it with  $f \in A_{22}^{s'}$ , where  $A_{pq}^{s'}$  is the A-space of V. A. Popov (see  $B^{12(b)}$ ). In this improvement of the results in this section, the condition

$$2^{j_1} = o(n^{\min(\frac{s_1}{\sigma}, \frac{1}{1+2\sigma})})$$

is replaced by

$$2^{j_1} = o(n^{\frac{2\min(s', 1/2)}{1+2\sigma}}), \quad 0 \leq \sigma < s' \leq s < r.$$

When  $s > 1/p$ , or  $s = 1/p$  and  $q \leq \min(1, p)$ ,  $A_{pq}^s$  and  $B_{pq}^s$  are isomorphic, with equivalence of the (quasi-)norms (Ivanov [77], Dechevski [39]), while if  $s < 1/p$  or  $s = 1/p$  and  $q > \min(1, p)$ , the essential embedding  $A_{pq}^s \hookrightarrow B_{pq}^s$  holds. For  $0 < p \leq \infty$ ,  $0 < q \leq \infty$ ,  $s > 0$ ,  $s_1 > 0$ ,  $B_{pq}^s \cap B_{\infty\infty}^{s_1} \subset A_{pq}^s \subset B_{pq}^s$  holds. All the elements of  $A_{pq}^s$  are bounded measurable functions defined pointwise everywhere on their domain, and for  $s < 1/p$  or  $s = 1/p$  and  $q = \infty$ ,  $A_{pq}^s$

contains also discontinuous functions. In fact, all interesting fractal curves are contained in some  $A_{pq}^s$ -space with  $s > 0$ ; for example, the fractal dimension of the function's graph is equal to  $2 - s$ , where  $s$  is the sharp smoothness index of the space  $A_{1,\infty}^s$  to which the function belongs. As for the set  $B_{pq}^s \setminus A_{pq}^s$ ,  $s < 1/p$ ,  $p < \infty$ , it contains unbounded and 'monstrous' bounded functions (see [77] for details).

Summarizing, our method GFCV is very well adapted for estimation of continuous functions with fractal graphs, as well as discontinuous functions.

## 6.2. Density estimation

The usual cross validation approach in density estimation is to replace  $V_0(v)$  in (8) by the  $ISE(v) = \int (\tilde{f}_v(x) - f(x))^2 dx = \int \tilde{f}_v^2(x) dx - 2 \int \tilde{f}_v(x) f(x) dx + \int f^2(x) dx$ . (For the sake of simpler and less spacious theoretical derivations, in this section the parameter  $s$  is considered to be fixed). In practice, one minimizes the empirical version of the ISE (up to a summand not depending on  $v$ ), i.e.,

$$M(\tilde{f}_v) = \int \tilde{f}_v^2(x) dx - \frac{2}{n} \sum_{i=1}^n \tilde{f}_v^{(-i)}(X_i) = \sum_k \hat{\alpha}_{j_0 k}^2 + \sum_{j=j_0}^{j_1} \sum_k \frac{\hat{\beta}_{jk}^2}{(1 + t^2 2^{2js})^2} - \\ - \frac{2}{n} \sum_{i=1}^n \left[ \sum_k \hat{\alpha}_{j_0 k}(-i) \varphi_{j_0 k}(X_i) + \sum_{j=j_0}^{j_1} \sum_k \hat{\beta}_{jk}(-i) \frac{\psi_{jk}(X_i)}{1 + t^2 2^{2js}} \right].$$

After a simple manipulation, we get (up to summands not depending on  $v$ ) the following expression to be minimized:

$$\tilde{M}(\tilde{f}_v) = \frac{1}{n^2} \sum_{i=1}^n \sum_{i'=1}^n \sum_{j=j_0}^{j_1} \sum_k \frac{\psi_{jk}(X_i) \psi_{jk}(X_{i'})}{(1 + t^2 2^{2js})^2} - \\ - \frac{2}{n(n-1)} \sum_{i=1}^n \sum_{i'=1, i' \neq i}^n \sum_{j=j_0}^{j_1} \sum_k \frac{\psi_{jk}(X_i) \psi_{jk}(X_{i'})}{1 + t^2 2^{2js}}.$$

**Theorem 5 .** (Sufficient conditions for existence of the  $v$ -minimizer.)

Let  $0 < s' < s < r$ . Assume that the penalized model is via  $K_2(\sqrt{v}, f; L_2, \dot{B}_{22}^s)$ ,

$f \in B_{22}^{s'} \cap L_\infty$  and that  $j_0 \leq j_1$ ,  $j_1 \rightarrow \infty$  and  $2^{j_1} = o(n)$ . If  $j_0 = O(1)$  and  $f \notin V_j$  for any  $j < j_0$ , or if  $j_0 \rightarrow \infty$ ,  $f \notin V_j$  for any  $j \in \mathbb{Z}$  and the index  $s'$  is sharp, then, for sufficiently large  $n$  there exists  $\tilde{v} : MISE(\tilde{v}) = \min_v MISE(v)$ , with  $0 < \tilde{v} < C$ . Under the same conditions, if in the case  $j_0 \rightarrow \infty$  additionally  $2^{j_0} = O(n^{\frac{1}{1+2s}})$ ,  $0 < s < r$ , holds, and, in particular, for  $0 < s \leq 1/2$   $2^{j_0} \leq 2^{j_1} \leq \text{const} \cdot 2^{j_0}$  is additionally fulfilled, then, with probability tending to 1 as  $n \rightarrow \infty$ , there exists  $v^* : M(\tilde{f}_{v^*}) = \min_v M(\tilde{f}_v)$ , and  $0 < v^* < C$  holds. Here  $C < \infty$  depends on  $f$  and  $\psi$  only.<sup>A18</sup>

Remark 2 is valid for  $\tilde{f}_{\tilde{v}}$  and  $\tilde{f}_{v^*}$  in the density case, too. However, it is not clear whether this remark is relevant when estimating compactly supported densities, because, although for any  $j \in \mathbb{Z}$  there exist functions in  $V_j$  which are non-trivially non-negative on  $R$  (see [118]), the functions in the examples available are not compactly supported.

Analogously to Theorem 2, we now show consistency of cross validation for the case of density, i.e., asymptotic equivalence of minimizing  $M(\tilde{f}_v)$  and  $MISE(v) : \frac{MISE(v^*)}{MISE(\tilde{v})} \xrightarrow{P} 1, n \rightarrow \infty$ . Corollary 1 in [110] implies that this follows from the next

**Theorem 6 .** (Consistency of CV - cf. also [49].) Assume that  $0 < s' \leq s < \infty$ , the penalized model is via  $K_2(\sqrt{v}, f; L_2, \dot{B}_{22}^s)$ ,  $f \in L_\infty \cap B_{22}^{s'}$ , and that the compactly supported  $\psi \in B_{\infty\infty}^r$ ,  $r > s$ , is orthogonal to algebraic polynomials  $P$  with  $\deg P \leq [r]$ . Assume also that  $j_0 \leq j_1$ ,  $j_1 \rightarrow \infty$ ,  $2^{j_1} = o(n)$ , and  $2^{-j_0} = O(n^{-\frac{1}{2(1+s')}})$ . Then, for any  $v \geq 0$ ,

$$\frac{M(\tilde{f}_v) - MISE(v) + T_n}{MISE(v)} \xrightarrow[n \rightarrow \infty]{P} 0,$$

where  $T_n = \int f(x)^2 dx + \frac{2(n+1)}{n} \sum_k \hat{\alpha}_{j_0 k} \alpha_{j_0 k} - \frac{2(n+1)}{n} \sum_k \alpha_{j_0 k}^2$  is a quantity that does not depend on  $v$  and  $T_n \xrightarrow[n \rightarrow \infty]{P} \int f(x)^2 dx$  holds.

**Remark 6.** The condition  $f \in L_\infty$  in Theorem 6 is essential for  $s \leq 1/2$ , that is, for less regular and possibly discontinuous curves. For  $s > 1/2$ ,



the *Sobolev embedding* (see, e.g., [16]) ensures that  $\|f\|_{L_\infty} \leq c_{s'} \|f\|_{B_{22}^{s'}}$ , where  $c_{s'}$  depends on  $s'$  only. Readers who utilize the reference source [110] should be cautioned that, for the same reasons, the additional constraint  $f \in L_\infty$  must be imposed in subsection 2.3 there to validate the proof of the bound for  $\text{Var}(T_{12})$  on p.50 when  $0 < s < d/2$ ,  $d$  being the dimension (in our case,  $d = 1$ ).

**Theorem 7 .** (*Consistency of  $\tilde{f}_{\tilde{v}}$  in  $B_{22}^\sigma$ ,  $0 \leq \sigma < s' \leq s < r$ .) Assume that the conditions of Theorem 5 hold, without the sharpness condition on  $s'$ . Then, if  $j_0 = O(1)$ , and if  $v = v_n$  is bounded, the condition  $v \rightarrow 0$  as  $n \rightarrow \infty$  is necessary for  $E\|\tilde{f}_v - f\|_{B_{22}^\sigma}^2 \rightarrow 0$  to hold; moreover, in both cases  $j_0 = O(1)$  and  $j_0 \rightarrow \infty$ , if  $2^{j_1} = o(n^{\frac{1}{1+2\sigma}})$  holds, then  $v \rightarrow 0$  is also a sufficient condition for  $E\|\tilde{f}_v - f\|_{B_{22}^\sigma}^2 \rightarrow 0$ . In particular, if  $s'$  is sharp, then: the above necessity claim is true for  $v = \tilde{v}_n$ ; if  $2^{j_1} = o(n^{\frac{1}{1+2\sigma}})$  is fulfilled, too, the above sufficiency claim is also true for  $v = \tilde{v}_n$ .*

**Theorem 8 .** (*Asymptotic behaviour of  $\tilde{v}$ .) Under the premises of Theorem 5,  $\tilde{v} \rightarrow 0$ .*

**Corollary 2 .** (*Consistency of  $\tilde{f}_{v^*}$  in  $B_{22}^0 = L_2$  in probability.) Under the premises of Theorem 5,  $\text{MISE}(v^*) \xrightarrow{P} 0$  as  $n \rightarrow \infty$  holds.*

The method of proof of Theorem 8 allows obtaining precise quantitative results, too: for any admissible choice of  $j_0$  and  $j_1$  the sharp asymptotic rates can be found for  $\tilde{v} \rightarrow 0$ . A complete analysis of the rates for  $\tilde{v}$  for any admissible  $j_0$  and  $j_1$ , the respective rates for  $E\|\tilde{f}_{\tilde{v}} - f\|_{B_{22}^\sigma}^2 \rightarrow 0$ , and comparison of these rates to the asymptotically minimax-optimal ones, will be carried out elsewhere.<sup>B14</sup> Here we only include the following simple model corollary.

**Corollary 3 .** (*Asymptotic-minimax rates in  $B_{22}^\sigma$ -norm,  $0 \leq \sigma < s' = s$ .) Under the premises of Theorem 5, assume that  $s = s'$  and  $c.n^{\frac{1}{1+2s}} \leq 2^{j_0} \leq 2^{j_1} \leq C.n^{\frac{1}{1+2s}}$ . Then,  $\tilde{f}_{\tilde{v}} = O(n^{-\frac{2s}{1+2s}})$  and  $E\|\tilde{f}_{\tilde{v}} - f\|_{B_{22}^\sigma}^2 = O(n^{-\frac{2(s-\sigma)}{1+2s}})$ , i.e., for this choice of  $s$ ,  $j_0$  and  $j_1$   $\tilde{f}_{\tilde{v}}$  achieves the asymptotic-minimax rate.*

Of course, Theorem 6 is not strong enough to imply, together with Corollary 3, that  $\tilde{f}_v^*$  also achieves the asymptotic-minimax rate, but the choice of  $s$ ,  $j_0$  and  $j_1$  in Corollary 3 is clearly of interest also for  $\tilde{f}_v^*$ .

**Remark 7.** The constraint on  $j_0$  in Theorem 6 is a lower bound depending on  $s'$ . It means that, if  $j_0$  is so selected that  $(\tilde{C}n)^{\frac{1}{2(s_0+1)}} \leq 2^{j_0} < 2(\tilde{C}n)^{\frac{1}{2(s_0+1)}}$ , for some  $s_0 : 0 < s_0 < r$ , then the theorem's statement holds uniformly in  $s' : s_0 \leq s' < r$ . Consideration of a bound  $s_0$  can be avoided by simultaneously optimizing the CV-functional in  $v$  and  $s$ , hence, also in  $j_0 = j_0(n, s)$ . In practice, for moderate samples, the relation  $j_0 = j_0(n, s)$  can be ignored and optimization can be in  $v$ ,  $s$ , and  $j_0$  as an independent discrete parameter. This is an analogue of our practical approach to the regression problem (see Remark 5), as well as of the practical cross-validation approach for estimating  $j_0$  by Tribouley [110].

Finally, comparing our results for the density and the regression case, we note that the condition  $f \in B_{22}^{s'} \cap L_\infty$  in the theory of density estimation is less restrictive than its counterpart  $f \in B_{22}^{s'} \cap B_{\infty\infty}^{s_1}$  (or the more general condition  $f \in A_{22}^{s'}$ ) in the case of regression. The reason is that in the density case  $\hat{\alpha}_{jk}$ ,  $\hat{\beta}_{jk}$  are unbiased estimators of  $\alpha_{jk}$ ,  $\beta_{jk}$ , respectively. For the same reasons, the same relaxation of the assumptions on  $f$  occurs for the Bowman-Rudemo type approach to cross validation for regression problems with random design (see  $B^1$ , [50,53]).

## 7. Examples and extensions

In this section we consider some model examples of estimation of a regression function or a density  $f$ , together with some first extensions of the penalized shrinking wavelet model (see also Appendix B). The comparison is between penalized shrinking and soft thresholding; numerical data will be added for hard thresholding in case that its performance turns out to be better than that of soft thresholding. To avoid introducing wavelets for a domain with boundary,

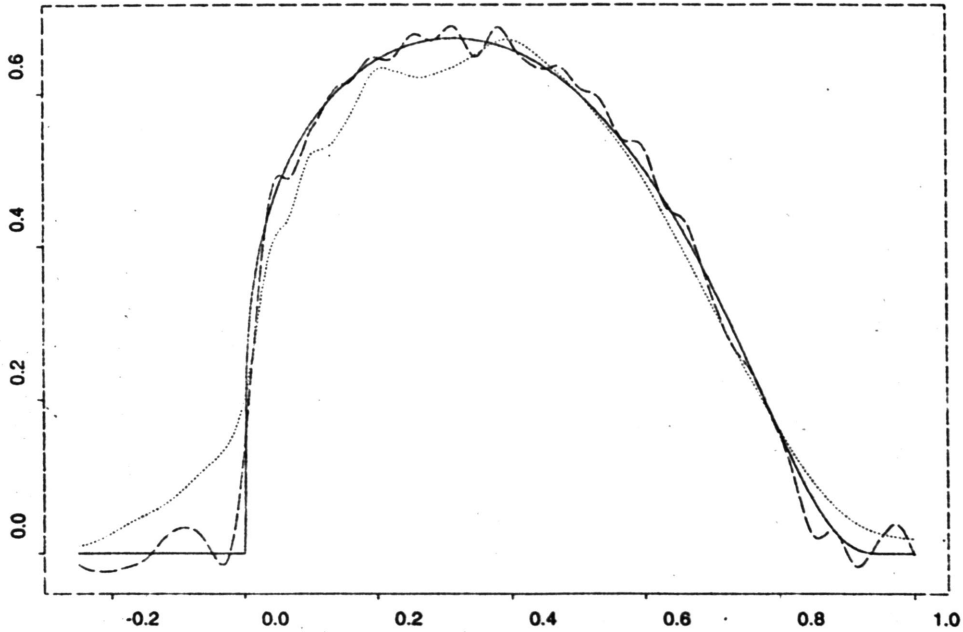


Figure 1: " $\lambda$ -tear":  $\lambda = 0.25$ ,  $n = 1024$ ,  $y_i = f[-\lambda + (1 + \lambda)i/n] + \epsilon_i$ ,  $\epsilon_i \sim N(0, 0.1^2)$ . True curve: solid line. Estimators: global penalized shrinking: dashed line; soft thresholding: dotted line.

$f$  is always assumed defined on  $R$ , with compact support. The combined requirement for orthogonality and compact support of the wavelet excludes the use of spline-wavelets; throughout Daubechies' extremal wavelets ([36], p. 195) of order  $M=6$  or  $8$  are utilized, calculated to considerable accuracy. The loss is quadratic and the values of  $\pi, p, u, q, \sigma, s$  are as in Section 5. The GFCV-functional is being minimized simultaneously in  $v, s$  and  $j_0$  (see Remarks 5, 7)<sup>A19</sup>;  $2^{j_1} \propto \frac{n}{\ln n}$  (which, by way of Corollary 1, means that  $\sigma = 0$ , i.e., only the function is being consistently estimated, but not its fractional derivatives of any positive order).

**Example 1.** *The "λ-tear" (Figure 1).*

$$f(x) = x_+^\lambda \exp\left(-\frac{x^2}{1-x^2}\right), \quad \lambda > 0,$$

where  $x_+ = \max(x, 0)$ ,  $\lambda = 0.25$ .  $f \in C^0(R) \cap C^\infty(R \setminus \{0\})$ ,  $\text{supp } f = [0, 1]$ . For  $1 \leq p \leq \infty$ ,  $0 < q \leq \infty$  and  $\lambda + 1/p > 0$ ,  $f \in B_{pq}^s$  if and only if  $s < \lambda + 1/p$  or  $s = \lambda + 1/p$ ,  $q = \infty$  ([22], Proposition 2.4.2).

Figure 1 displays clearly that around the critical point  $x = 0$  the shrinking estimator shows better adaptivity. Its graph can be made less "wiggly" by selecting a smaller value for  $j_1$ , as indicated in Corollary 1. (In this case,  $s_1 = \lambda = 1/4$ ,  $s' = \lambda + 1/2 = 3/4$ , and, by Corollary 1, a good selection for  $j_1$  is, e.g.,  $j_1 : 2^{j_1} \asymp n^{\min(\frac{1}{4\sigma}, \frac{1}{1+2\sigma})} / \ln \ln n$ , where  $\sigma : 0 < \sigma < s' = 3/4$  is small enough for  $2^{j_0} \leq 2^{j_1}$  to hold.) For further improvements, see  $B1, B6, B8$ .

**Example 2.** *The Weierstrass curve (Figure 2).*

$$f(x) = \sum_{k=0}^{\infty} 1.5^{-\tau k} \sin(1.5^k \times 5x), \quad \tau > 0, \quad x \in [0, 1],$$

where  $\tau = 0.5$ . As usual, we shall consider its restriction  $\bar{f}$  on  $[0, 1]$  (see also [119]). The graph of  $f$  is a typical *self-similar monofractal*: it has constant local Hölder index  $\tau$  and constant local fractal dimension  $2 - \tau$  (which is also its global fractal dimension on  $[0, 1]$ ). For such functions the universal thresholding and global penalized shrinking are expected to be at their best. For any compactly supported  $\chi \in C^\infty(R)$  such that  $[0, 1] \subset \text{supp } f$  and  $\chi \equiv 1$  on  $[0, 1]$ ,  $\chi \cdot f \in B_{pq}^s(R)$ ,  $1 \leq p \leq \infty$ ,  $0 < q \leq \infty$ , if and only if  $s < \tau$  or  $s = \tau$ ,  $q = \infty$  ([22], Proposition 2.4.1, the imaginary part of  $G_\tau$  with additional rescaling). By the Whitney-type trace theorem for Besov spaces (see, e.g., [96]), the same is true for  $\bar{f}$  with respect to  $B_{pq}^s([0, 1])$ .

The shrinking estimator on Figure 2 shows better adaptivity to abrupt changes in the curve's graph than the soft- (and hard-) thresholded. This is true

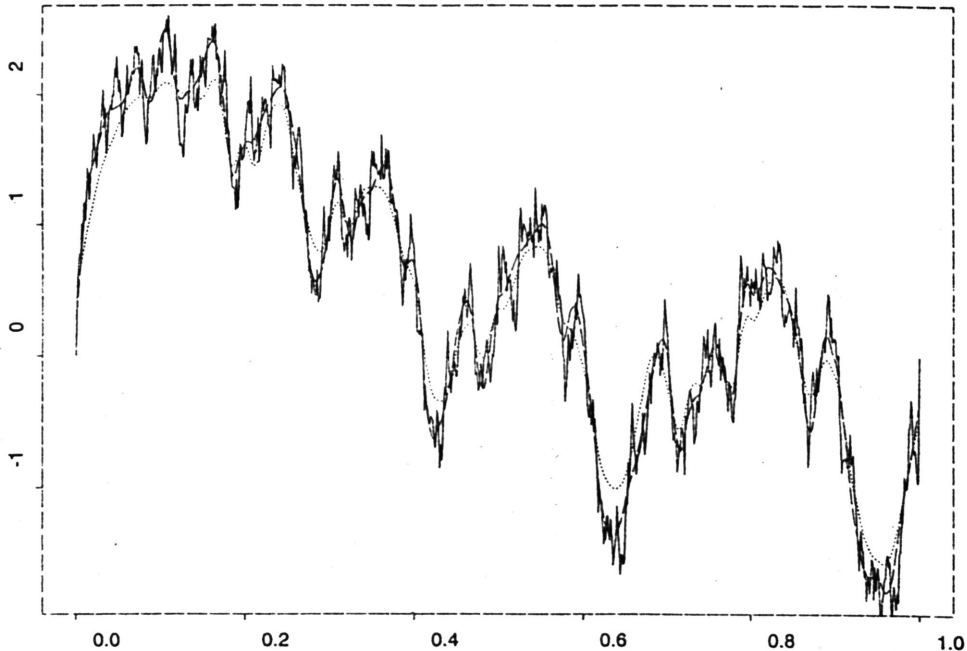


Figure 2: Weierstrass curve:  $\tau = 0.5$ ,  $n = 1024$ ,  $y_i = f(i/n) + \epsilon_i$ ,  $\epsilon_i \sim N(0, 0.2^2)$ . True curve: solid line. Estimators: global penalized shrinking: dashed line, estimated ISE  $\approx 3.442 \times 10^{-2}$ ; soft thresholding: dotted line, estimated ISE  $\approx 6.691 \times 10^{-2}$ ; hard thresholding: estimated ISE  $\approx 4.013 \times 10^{-2}$ . Data for  $f_1$  - discontinuous Weierstrass-type curve with vertical displacement at 0.5: estimated ISE  $\approx 3.450 \times 10^{-2}$  (global penalized shrinking);  $\approx 6.683 \times 10^{-2}$  (soft thresholding).

for discontinuities, too (numerical data are included for  $f_1(x) := f(x) + 0.4 \times (x - 0.5)_+^0$  in the caption to Figure 2).

**Extension: kernel-regularization of the shrinking wavelet estimator.** So far, we have been paying most attention to applications of the shrinking wavelet estimator in the case of non-smooth functions. What about very smooth functions? Recall that, for the choice  $j_1 : 2^{j_1} \asymp n/\ln n$ ,  $\tilde{f}_v$  in Section 6.1 and  $\tilde{f}_v^*$  in section 6.2 have been shown to be consistent in  $B_{22}^\sigma$  for  $\sigma = 0$  only, i.e., they estimate consistently the function, but not its derivatives of any positive order. Example 1 clearly illustrates this. Several ways can be

suggested to upgrade the shrinking wavelet estimator to make it estimate derivatives consistently. The simplest solution suggested by Theorem 3, Corollary 1 and Theorem 7 is to decrease  $j_1$ :  $2^{j_1} = o(n^{\frac{1}{1+2\sigma}})$  should hold, rather than the standard  $2^{j_1} = O(\frac{n}{\ln n})$ . Other possibilities will be discussed in  $B^1, B^6, B^8$ . Here we consider an alternative based on kernel-regularization. Let  $\Phi \in B_{\infty\infty}^r(R)$ ,  $\text{supp } \Phi = \{x : |x| \leq 1\}$ ,  $\int \Phi(x)dx = 1$  and, optionally,  $\Phi \geq 0$  on  $R$ . Denote  $\Phi_\varepsilon := \frac{1}{\varepsilon}\Phi(\frac{\cdot}{\varepsilon})$ ,  $\varepsilon > 0$ . Consider  $\tilde{f}_{v,\varepsilon}(x) := (\Phi_\varepsilon * \tilde{f}_v)(x)$ . Let  $v_{opt}$  be the  $v$ -minimizer with respect to any of the CV-functionals in Section 6. For the bandwidth  $\varepsilon$  we suggest the choice  $\varepsilon = O(v_{opt}^{1/2})$ . This topic is currently under investigation and results obtained so far indicate that  $\tilde{f}_{v,\varepsilon}$  with  $v = v_{opt}$  and  $\varepsilon = O(v_{opt}^{1/2})$  is a very competitive alternative to thresholded wavelet estimators when  $f$  is smooth and sample sizes are moderate. For a typical graphical result we refer to Example 4 below. For numerical integration of the convolution integral we recommend Romberg integration. This quadrature method is excellently suited to dyadic wavelets and can in fact be considered as an extension of Mallat's recursive algorithm "on sub-coefficient level".

The study of consistency of  $\tilde{f}_{v,\varepsilon}$  in Besov spaces is much facilitated by the equivalent definition of these spaces, considered, e.g., in [106], Definition 1 (ii), which is valid for a broader range of  $p$ ,  $q$  and  $s$  than the orthonormal-wavelet atomic decomposition of these spaces (see section 3).

**Extension: vector parameters of optimization.** There is an ongoing discussion in the literature about passing over from universal to *level-dependent*, *block* and even *individual* thresholding of the empirical wavelet coefficients. Level-dependent thresholding has been proposed by Donoho and Johnstone [57] and Delyon and Juditsky [53]. The block-thresholding approach was proposed and developed by Hall, Kerkycharian and Picard [68] and its numerical performance was studied in [70]. Individual shrinking was discussed, e.g., in [67]. Other relevant development and discussion of these topics can be found in [89,82,25,28,29]). The aims of these enhancements of the thresholding ap-

proach are twofold: (a) achieving better spatial and frequency adaptivity while preserving the smoothing and denoising properties of the wavelet estimator; (b) studying regression problems with correlated-noise error structure. This discussion can be readily extended to penalized shrinking, too.<sup>A20</sup> In the context of penalized shrinking, the new enhanced strategies are equivalent to replacing  $t$  (and, eventually,  $s$ ) in (7) by  $t_{jk}(s_{jk})$ , and then relaxing the global constraint  $t_{jk} = \text{const} = t$  ( $s_{jk} = \text{const} = s$ ),  $j = j_0, \dots, j_1$ ,  $k \in Z$ , by imposing a less stringent condition on  $t_{jk}(s_{jk})$ . The resulting extended version of the  $K$ -functional is an equivalent quasi-(semi-) norm in the quasi-Banach sum of *closed subspaces* of two (or more) Besov spaces.<sup>B12(a), B17</sup> The resulting cross-validation functionals for both the regression and the density case are now to be optimized in a vector parameter  $\mathbf{v}$  (in the optimization for the level-dependent case the parameter  $\mathbf{s}$  is eliminated by a change of variable, rescaling  $\mathbf{v}$ ).

**Level-dependent smoothing parameters.** Optimization of the cross-validation functional is in  $\mathbf{v} = (v_{j_0}, \dots, v_{j_1})'$ ,  $v_j = t_j^2$ . The expected improvement is better adaptivity in the *frequency* domain. Our first test was on the Weierstrass curve with white-noise error structure (Example 2), where the expected outcome for the optimal  $\mathbf{v}^*$  and  $v^*$  from Example 2 was  $v_j^* \approx \text{const} = v^*$ , because of the global self-similarity of the curve. Indeed, the new fit virtually coincided with the old fit from Example 2,  $\text{ISE}(\mathbf{v}^*)$  coinciding with  $\text{ISE}(v^*)$  up to the fourth decimal place. There was a slight change in favour of the level-dependent estimator when the white-noise model was replaced by fractional Brownian noise with correlation structure  $\rho(k) = 0.5[(k+1)^{2H} - 2k^{2H} + (k-1)^{2H}]$ ,  $k \geq 0$  and  $\rho(k) = \rho(-k)$  for  $k < 0$ , and with standard deviation of 0.2. This type of noise was simulated with  $H = 0.95$  by using the S-plus routine `simFGNO` in Chapter 12 of [15]. In one of the typical numerical results obtained,  $\text{ISE}(v)$  was  $4.633 \times 10^{-2}$  and  $\text{ISE}(\mathbf{v})$  was  $4.621 \times 10^{-2}$ .

**Block-penalized shrinking.** The constraints on  $t_{jk}$  are as follows: if  $\text{supp } \psi_{jk}$  and  $\text{supp } \psi_{j'k'}$  are contained in  $\text{supp } \psi_{j_0k_0}$  for the *same* value of  $k_0$ , then  $t_{jk} = t_{j'k'}$ . The expected improvement is better *spatial* adaptivity - see the

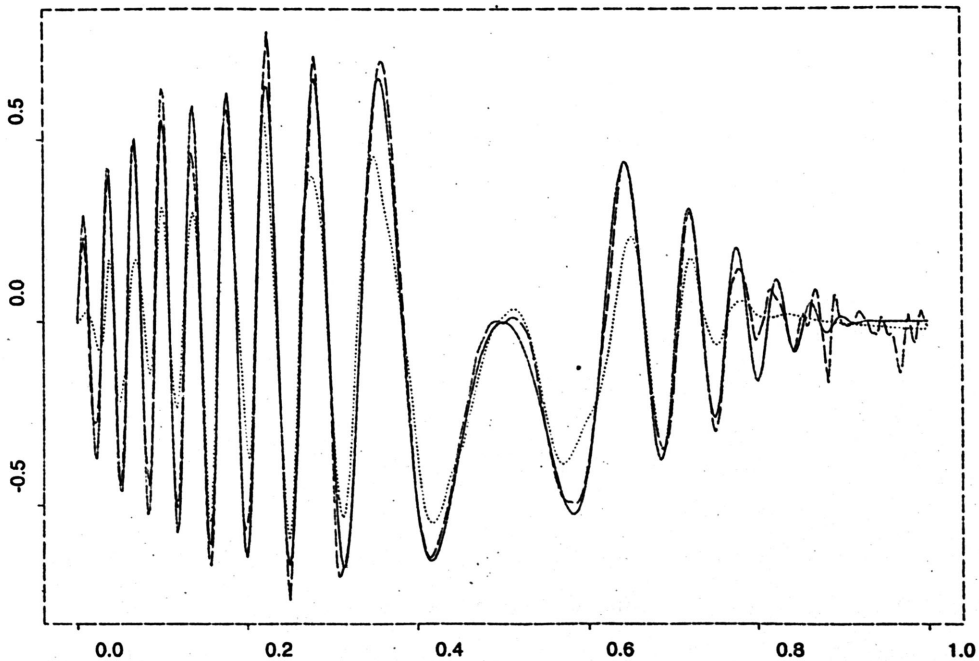


Figure 3: Double chirp:  $n = 1024$ ,  $y_i = f(i/n) + \epsilon_i$ ,  $\epsilon_i \sim N(0, 0.2^2)$ . True curve: solid line. Estimators: block penalized shrinking: dashed line; soft thresholding: dotted line.

following

Example 3. *Double chirp* (Figure 3).

$$f(x) = \sqrt{x} \exp\left(-\frac{x^2}{1-x^2}\right) \sin[64\pi x(1-x)],$$

$f \in C^0(R) \cap C^\infty(R \setminus \{0\})$ ,  $f^{(\nu)}(1) = 0$ ,  $\nu \in N \cup \{0\}$ ,  $f(0) = 0$ . The graph of  $f$  contains two single chirps (at 0 and 1) of very different nature. In a neighbourhood of 0  $f$  is with unbounded variation; in a neighbourhood of 1 (not containing 0)  $f$  is absolutely continuous. Its Besov regularity is bounded from above by that of its "profile": the " $\lambda$ -tear" for  $\lambda = 0.50$  (cf. Example 1).<sup>421</sup>



**Iterative individual penalized shrinking.** In this extreme case all  $t_{jk}$ ,  $j = j_0, \dots, j_1$ ,  $k \in Z$ , are independent parameters in which the cross-validation functional is being optimized. This extension is of interest for both the regression and density model. Here we shall consider density estimation. Then, optimization of  $\tilde{M}(\tilde{f}_v)$ ,  $v = \{v_{jk}, j = j_0, \dots, j_1, k \in Z\}$ , in  $v_{jk} = t_{jk}^2$  yields

$$\frac{[\sum_{i=1}^n \psi_{jk}(X_i)]^2}{n^2(1 + v_{jk}2^{2js})^3} = \frac{[(\sum_{i=1}^n \psi_{jk}(X_i))^2 - \sum_{i=1}^n \psi_{jk}(X_i)^2]}{n(n-1)(1 + v_{jk}2^{2js})^2}.$$

Hence, the following analogue of (7) is obtained:

$$\tilde{\beta}_{jk} = \frac{\hat{\beta}_{jk}}{1 + v_{jk}2^{2js}} = \hat{\beta}_{jk} \cdot [1 - \frac{1}{n-1} \frac{\overline{\psi_{jk}^2} - (\overline{\psi_{jk}})^2}{(\overline{\psi_{jk}})^2}],$$

where  $\overline{\psi_{jk}^2} = \frac{1}{n} \sum_{i=1}^n \psi_{jk}^2(X_i)$ ,  $\overline{\psi_{jk}} = \frac{1}{n} \sum_{i=1}^n \psi_{jk}(X_i)$ . Now the estimator (1') is very undersmoothed. To smooth it more, we reapply the individual penalization with  $\hat{\beta}_{jk}^{(1)} = \tilde{\beta}_{jk}$ . Repeating this procedure leads to ever smoother, highly adaptive estimators. Performing  $O(n)$  iterations has approximately the same effect for large  $n$  as the following *individually shrinking* rule:

$$(15) \quad \tilde{\beta}_{jk} = \hat{\beta}_{jk} C_1(s, v) \exp(-C_2(s, v) \frac{\overline{\psi_{jk}^2} - (\overline{\psi_{jk}})^2}{(\overline{\psi_{jk}})^2})$$

where the parameters  $C_1 : 0 < C_1 < \infty$  and  $C_2 : 0 < C_2 < \infty$  are explicit functions of  $s$  and  $v_{jk}$ ,  $j = j_0, \dots, j_1$ ,  $k \in Z$ . Substituting the quantities

$\frac{1}{1 + v_{jk}2^{2js}} = \frac{\tilde{\beta}_{jk}}{\hat{\beta}_{jk}} = C_1 \exp(-C_2 \frac{\overline{\psi_{jk}^2} - (\overline{\psi_{jk}})^2}{(\overline{\psi_{jk}})^2})$  back in the expression for  $\tilde{M}(\tilde{f}_v)$  yields a cross-validation functional of the form:

$$\begin{aligned} & \sum_{j=j_0}^{j_1} \sum_k \{ (\overline{\psi_{jk}})^2 C_1^2 \exp(-2C_2 \frac{\overline{\psi_{jk}^2} - (\overline{\psi_{jk}})^2}{(\overline{\psi_{jk}})^2}) - \\ & - \frac{2}{n-1} [n(\overline{\psi_{jk}})^2 - \overline{\psi_{jk}^2}] C_1 \exp(-C_2 \frac{\overline{\psi_{jk}^2} - (\overline{\psi_{jk}})^2}{(\overline{\psi_{jk}})^2}) \} \end{aligned}$$

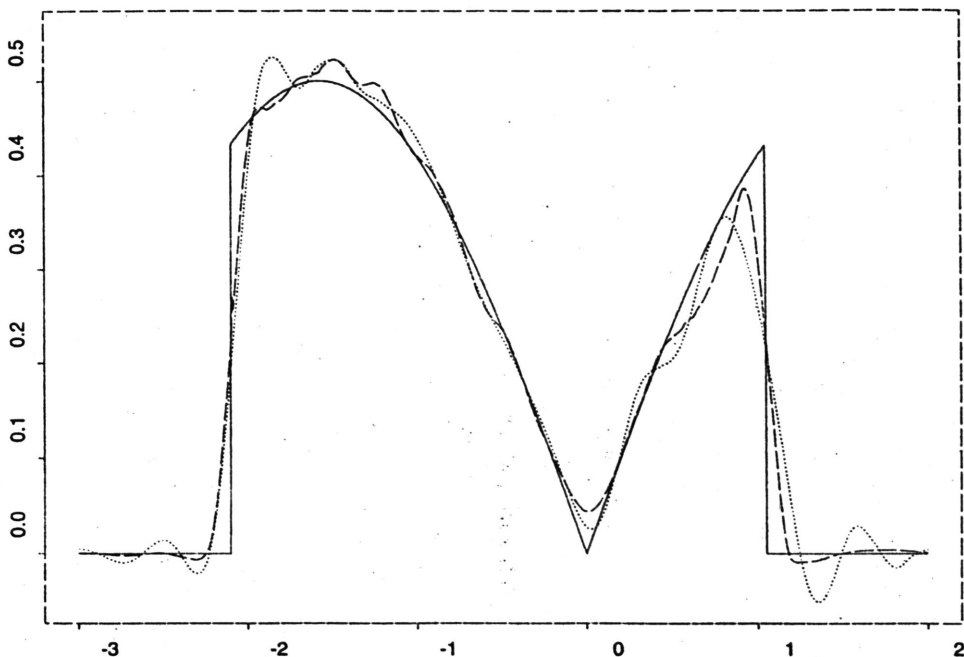


Figure 4: Sinusoidal density:  $n = 2000$ . True curve: solid line. Estimators: iterative individual penalized shrinking: dashed line; soft thresholding: dotted line. Kernel smoothing was by convolving with  $h_\epsilon(x) = \frac{1}{\epsilon}h(\frac{x}{\epsilon})$ ,  $h(x) = K \cdot \exp[-(\frac{x^2}{1-x^2})_+]$ , where  $K = 1/1.2069$  is a norming constant. The kernel bandwidth  $\epsilon$  was 0.2. The approximate computation of the convolution integral was via the trapezoidal rule.

which is to be minimized with respect to *two unknowns only*:  $C_1$  and  $C_2$ , with expected asymptotic behaviour  $C_1 \rightarrow 1$ ,  $C_2 \rightarrow 0$ . Hence, the strategy of choosing individual shrinkers is reduced to a minimization problem with respect to two variables only and can be tackled easily.

Formula (15) leads to an extremely adaptive estimator. Since it is generally still undersmoothed, additional smoothing is recommended, e.g., by convolving it with a smooth density kernel. The estimator is so adaptive that even the use of a simple constant-bandwidth kernel suffices to produce fits which are superior to the fits by standard thresholding techniques (for a typical example, see Figure 4). The iterative approach is currently under investigation, both

with respect to kernel smoothing<sup>B3,B4</sup> and with respect to composite shrinking/thresholding estimation<sup>B6,B8</sup>, and we hope to report some more detailed results soon. So far, we only note that the (constant) kernel bandwidth must be  $O(\sqrt{v^*})$ ,  $n \rightarrow \infty$ , where  $v^* = v_n^*$  is the minimizer from Theorem 5.

Another interesting opportunity here, in particular, in the case of Gaussian white noise, is to consider, levelwise and blockwise within each level, a composite shrinking/thresholding strategy by SURE (cf. <sup>B8</sup> (the second approach)).

**Example 4.** Sinusoidal density (Figure 4).

$$f(x) = \begin{cases} \frac{1}{2} |\sin x| & \text{for } x \in [-2\pi/3, \pi/3], \\ 0 & \text{elsewhere.} \end{cases}$$

Figure 4 shows that the new estimator adapts better to the spatial irregularities of  $f$  along the curve, including the discontinuity points  $x = -2\pi/3$  and  $x = \pi/3$  where the Gibbs effect is much smaller than with soft thresholding.

**Extension: self-similar fractal estimator.** If no additional information about self-similarity of  $f$  is available, between the  $\alpha$ - and  $\beta$ -coefficients only the relations from Mallat's recursive algorithm hold. For the wavelet estimators of  $f$ ,  $j_1 < \infty$  holds, which essentially means that  $f$  is being only perceived as a linear combination of wavelets and all information about the fractal structure of the graph of  $f$  comes from the fractal structure of the wavelet itself. Under these general assumptions, there are no forced relationships between the  $\beta$ -coefficients on neighbouring levels. Assume now that  $f$  is dyadically self-similar, i.e., it satisfies a dyadic functional equation

$$(16) \quad f(x) = \sum_{k=-L}^L c_k f(2x - k), \quad x \in R,$$

where  $L \in N$ . Then  $f$  is compactly supported at  $[-L, L]$  (see, e.g., [36,44] for details about the implied properties of  $f$ ) and there exists additional relationship between the betas on every two neighbouring levels. This relationship is given

by one and the same linear combination of the betas on the finer level, having for coefficients the  $c_k$ 's from (16). If  $f$  is being estimated in the non-parametric regression or density model, we see that for fixed  $v$  the problem about finding the optimal  $\tilde{\alpha}$ - and  $\tilde{\beta}$ -coefficients of  $\tilde{f}_v$  is now upgraded to an optimization problem with restrictions of type equality to which Lagrange-multiplier technique can be applied. The unknowns now are the  $\tilde{\alpha}$ 's and  $\tilde{\beta}$ 's and, additionally, the  $c_k$ 's from (16). The equality constraints are *bilinear* - they are linear with respect to two groups of unknowns in separate: the coefficients of  $\tilde{f}_v$  and the coefficients of the functional equation. This means that the optimization problem can be reduced to *convex*. It can be solved, e.g., via the *iterative quadratic penalty method* ([18], Section 2.1). This is a simple method, but the resulting problem about unconstrained multivariate optimization, solved in each iteration, may become very ill-conditioned. To avoid this, the *iterative Lagrange-multiplier method* ([18], Section 2.2) can be utilized. At every iteration it leads to an unconstrained optimization problem which is well-conditioned and, besides, its iterations converge superlinearly, much faster than with the quadratic-penalty method ([18], Subsection 2.2.5). More recently developed advanced superlinear iterative methods (of quadratic, of Newton type, etc.) can be found in literature from the last few years (see, e.g., [20]).

Since the optimization problem is (can be reduced to) convex, any vector of  $c_k$ 's is an admissible initial solution for starting the iterative process. To reduce the number of iterations, we suggest to obtain the initial solution in the following way: the initial  $\tilde{\alpha}$ - and  $\tilde{\beta}$ -coefficients are computed by (7), ignoring the information about self-similarity. The smoothing parameter is obtained as earlier, by cross validation. With the so-obtained values of the smoothing parameter and the respective  $\tilde{\alpha}$ 's and  $\tilde{\beta}$ 's, generating the initial values of the  $c_k$ 's is by *linear regression* (i.e., least-squares method) based on the equality constraints.

Although the class of integrable functions  $f$  that satisfy (16) is relatively small compared to, say, the whole of  $L_1$ , it contains functions which are interesting with respect to both applications and theory. For instance, such are

all Daubechies wavelets as well as orthogonal spline wavelets (for the latter,  $L = \infty$ ). Besides, our approach can be extended to all  $f$  which are solutions of *non-stationary dyadic refinement schemes* and form a much larger functional class. Here the  $c_k$ 's depend on the level  $j$  involved in the computation. A larger system of equations will result in this way, but the reward is that we can thus treat fractals with much more general structure. Having estimated the  $c_k$ 's means that we can construct a *fractal estimator* of  $f$  via the respective refinement scheme. In the stationary case ( $c_k$  independent of  $j$ ), this estimator is *globally self-similar*, like, say, the Weierstrass function. In the non-stationary case it is *self-similar locally in the frequency domain*. When the observations are noisy, the additional information about certain - global or local - self-similarity of  $f$  results in an estimator whose graph is a fractal as irregular as the noise itself, and yet it can be claimed that it is denoised to a considerable extent. Moreover, since self-similarity is a very spatially informative feature, it can be hoped that the fractal estimator will have good performance already for moderate samples, in particular, in the symmetric case, i.e., when it is known additionally that  $c_k = c_{-k}$  for all  $k$ .

For Gaussian and Poissonian white noise there are already publications discussing estimation of the local Hölder index and more detailed multifractal characteristics (see [12,76,14]). Under assumption of dyadic self-similarity and Gaussian white noise, the above suggested fractal estimator can now be upgraded to produce good fits already for small to moderate samples. The idea is to consider a penalized version of the NeighBlock and NeighCoeff estimators of [28] which takes into account dyadic self-similarity (with  $L_0 = 1$  in the notations of [28]). Assume first that the support of the dyadically self-similar function, that is, the value of  $L$  in (16) is known. (This is usually the case in related applications.) Then, the oracle inequality of NeighBlock is no longer needed; the length of each block, on every level, is  $d = 2L + 1$ ; the shrinking estimate is by the classical James-Stein estimator (JS)

$$\tilde{\beta}_{jk} = \left(1 - \frac{d-2}{\sum_{\nu=-L}^L \hat{\beta}_{j,k-\nu}^2}\right)_+ \hat{\beta}_{jk},$$

for any  $j = j_0, \dots, j_1$ , and any  $k : \hat{\beta}_{jk} \neq 0$ , assuming that the white noise is  $N(0, 1)$ .

We observe that JS can be obtained by a penalized model in the following way.

1) For any  $\theta \geq 0$ , minimize

$$\Phi(\tau_{-L}, \dots, \tau_0, \dots, \tau_L) = \sum_{l=-L}^L (\hat{\beta}_{j,k-l} - \tau_l)^2 + 2\theta \frac{\sum_{l=-L}^L |\hat{\beta}_{j,k-l}| |\tau_l|}{\sum_{\nu=-L}^L \hat{\beta}_{j,k-\nu}^2}$$

with respect to  $\tau_l$ ,  $l = -L, \dots, 0, \dots, L$ . The solution is

$$\tau_l = (1 - \frac{\theta}{\sum_{\nu=-L}^L \hat{\beta}_{j,k-\nu}^2})_+ \hat{\beta}_{j,k-l}.$$

2) Select optimal  $\theta$  by SURE, applied to the  $\tau_l$ 's as functions of the  $\hat{\beta}_{j,k-\nu}$ 's. The result is  $\theta = d - 2 = 2L - 1$ .

3) Set  $\tilde{\beta}_{jk} = \tau_0$ ; ignore the other  $\tau_l$ : continue in the same way for the next block, corresponding to  $\tilde{\beta}_{j,k+1}$ , and so on.

Now we shall demonstrate how the penalized model interpretation of JS allows to upgrade the estimator, so that the new version takes dyadic self-similarity into account.

Consider the problem of minimizing the function

$$\begin{aligned} & \sum_{j=j_0}^{j_1-1} \sum_k \{ (\hat{\beta}_{jk} - \tilde{\beta}_{jk})^2 + 2(d-2) \frac{|\hat{\beta}_{jk}| \cdot |\tau_{jk,0}|}{\sum_{\nu=-L}^L \hat{\beta}_{j,k-\nu}^2} + \\ & + \sum_{l=-L, l \neq 0}^L [(\hat{\beta}_{j,k-l} - \tau_{jk,l})^2 + 2(d-2) \frac{|\hat{\beta}_{j,k-l}| \cdot |\tau_{jk,l}|}{\sum_{\nu=-L}^L \hat{\beta}_{j,k-\nu}^2}] \}, \end{aligned}$$

with respect to the variables  $\tau_{jk,l}$ ,  $l = -L, \dots, -1, 1, \dots, L$ ,  $\tilde{\beta}_{jk}$ ,  $j = j_0, \dots, j_1 - 1$ ,  $k : \hat{\beta}_{jk} \neq 0$ ,  $c_\nu$ ,  $\nu = -L, \dots, 0, \dots, L$ , with the following constraints of type equality

$$\tilde{\beta}_{jk} = \frac{1}{\sqrt{2}} \sum_{\nu=-L}^L c_\nu \tilde{\beta}_{j+1,k+\nu}, \quad k : \hat{\beta}_{jk} \neq 0, \quad j = j_0, \dots, j_1 - 1,$$

$$\tilde{\beta}_{j_1 k} = (1 - \frac{d-2}{\sum_{l=-L}^L \hat{\beta}_{j_1, k-l}^2})_+ \hat{\beta}_{j_1 k},$$

i.e., the finest level is still estimated by the (unconstrained) JS.

If additional information about the relative importance of the different coefficients is available, a  $w_{jk}$ -weighted version of the cost functional can be considered.

The equality constraints come from (16), by argument similar to the one used in Mallat's reconstruction recursive algorithm. If these constraints are disabled, that is, if (16) is not fulfilled and, hence, there is no global dyadic self-similarity, it can be seen that the groups of variables  $\{\tilde{\beta}_{jk}, \tau_{jk,l}\}_{l=-L}^L$  (each one corresponding to the respective sliding block at position  $(j, k)$ ) become disconnected from each other, which leads to the classical James-Stein estimator for every  $(j, k)$ .

The problem can be solved numerically by a superlinear iterative method for nonsmooth convex optimization based on Clarke's subdifferential. The same is true if, for example, additional information about nonnegativity of  $f$  in (16) is available. In this case, additional constraints  $c_l \geq 0$ , of type inequality, are added to the optimization problem, which remains convex.

Further refinements of this model can be obtained by replacing the "hard" value  $d-2$  in the functional by a  $\theta_{jk}$  which varies with  $n$  in a small neighbourhood of  $d-2$ , and is reset by SURE-type optimization on every iteration of the optimization algorithm. However, developing a rigorous theory requires a much more detailed and elaborate consideration than is possible here.

If  $d$  is not known a priori, one good compromise between performance and simplicity of estimation of  $d$  is to apply a *global* oracle inequality, obtained by summing up in  $j$  the LHS and the RHS, respectively, of the levelwise oracle inequality (11) in [25], and thus select the same block size simultaneously for all levels  $j = j_0, \dots, j_1$ .

If the white noise is not Gaussian, this model can still be helpful for large samples, when the central limit theorem takes over. For moderate samples, however, the initially proposed variant via cross validation is expected to yield

better results, especially if the original noise distribution is far from Gaussian.

The requirement that the functional equation be dyadic is not essential. For example, one can consider triadically self-similar fractals, but then the wavelets  $\varphi$  and  $\psi$  must also be triadic.

### Appendix A: Proofs<sup>A0</sup> and details

**A0.** In all proofs, wherever this is relevant, we assume, with no loss of generality, that  $a_1 = 0$ ,  $a_2 = 1$ , in order to shorten the mathematical expressions.

**A1.** Fractional-order derivatives of  $f$  are (a) for the inhomogeneous model (see Sections 3, 4): in the sense of the *Bessel potential* of  $f$ , i.e., convolution of  $f$  with the Bessel-MacDonald kernel (see [103], Section 27); (b) for the homogeneous model (see Sections 3, 4): in the sense of the *Riesz potential* of  $f$ , i.e., convolution of  $f$  with the kernel in [103], Section 25, formula (25.25). Since we shall be estimating only functions  $f$  which are within the range of equivalence of the two models<sup>A14</sup>, we can consider the Bessel-MacDonald kernel only.

**A2.** The reason for this dramatic simplification is that atomic decomposition in Besov spaces (see [63,65]) leads to equivalent *sequence* norms (or semi-norms); in particular, when the atoms are orthonormal compactly supported wavelets (cf. [106]), the resulting sequences are comprised of the respective wavelet coefficients. In contrast to this, atomic decomposition in Triebel-Lizorkin spaces (including the Sobolev spaces associated with the smoothing spline technique) gives rise to norms that are not sequence norms except for those values of the indices for which the Besov and Triebel-Lizorkin space scales coincide.

**A3.** Here we single out only one of them: estimating functions whose graphs are self-similar fractals. Based on our approach, a new *self-similar fractal estimator* of noisy fractal curves can be constructed. Since it takes into account self-similarity, satisfactory estimation of simple fractal characteristics like, say, fractal dimension (see [119]) is expected to be achieved already for *moderate* samples. Moreover, considering a non-stationary, locally self-similar fractal estimator allows subsequent estimation of much more informative fractal



characteristics (e.g., Hölder spectrum, multifractal characteristics (see [11,91]), discontinuity points, etc.).

**A4.** Indeed, for *small and moderate samples*, the minimax estimators corresponding to two different metrics can be so far away from each other that any admissible estimator which is "near both of them" would necessarily have unsatisfactory performance with respect to *both* metrics!

**A5.** (1-3) can be generalized for  $B_{pq}^s(\mathbb{R}^d)$ , with preservation of the resolution-level expansion via a *tensor-product*  $d$ -variate wavelet basis (see [53] for a short overview; see also [36]). The range of the scale parameters is  $0 < p \leq \infty$ ,  $0 < q \leq \infty$ ,  $\max(0, d(1/p - 1)) < s < r$ , which is exactly the range for which Besov spaces contain *only regular distributions*, i.e., *only locally integrable functions*.

**A6.** (For two different  $\rho$ -values the respective  $K_\rho$ 's are equivalent but the extra parameter  $\rho$  provides useful additional flexibility.) For all  $t$  and all  $\rho$   $K_\rho(t, \cdot; A, B)$  is an equivalent (quasi-semi-)norm on  $A + B$ , the standard norm in this space being  $K_1(1, \cdot; A, B)$  (cf. [16]). Under very general assumptions the infimum in (4) is attained for a unique pair  $a^* \in A$ ,  $b^* \in B$ ,  $a^* = \alpha - b^*$ , for any  $\alpha \in A + B$ . (For example, this is true if  $A$  and  $B$  are *uniformly convex* Banach spaces.)

**A7.** Note that the second term on the right (sometimes known in approximation theory as a *saturation* term) is always subordinate to  $K_\rho(t, \cdot; A, \dot{B})$  which is typically  $c.t^\theta$  for some  $\theta \in [0, 1]$ .

**A8.** In fact, Delyon and Juditsky additionally require compact support of the estimated density or regression function. For such functions we can show that the range of admissible  $(\pi, u, \sigma)$  and  $(p, q, s)$  in the  $K$ -functional model is considerably *larger*.

**A9.** Then, the constraint on  $\tau$  via  $u$  and  $q$  can be removed. For  $\pi \neq p$  this is possible only if the model be refined by considering *four-indexed* Besov spaces.

**A10.** A comparison with the spline-smoothing model shows that the latter is essentially limited to the Hilbert range only ( $\pi = p = u = q = 2$ ). Out-

side this range the resulting Euler-Lagrange equations are *nonlinear differential* equations. The parametrization provided by the wavelet model makes it possible to compute the coefficients of the minimizing wavelet in an essentially *explicit* way for the *whole quasi-Banach range*.<sup>B9</sup> But even for  $\pi = p = u = q = 2$  the wavelet estimator is more adaptive: not only  $t$ , but also  $\sigma$ ,  $s$  and  $j_0$  (see (1)) can be optimized. (In principle, optimization in  $\sigma$  and  $s$  can be carried out for Wahba's spline model, too, but the problem is nonparametric - the *spectral theorem* has to be invoked. Then  $\sigma$  and  $s$  have the meaning of *fractional powers of a self-adjoint differential operator*. In practice, this means that a large number of eigenvalues and eigenfunctions have to be computed to a considerable precision, which can be a very ill-conditioned and numerically intensive problem.)

**A11.** The equivalence constants between the sequence norm (2) and any equivalent norm in the Besov space depend on  $j_0$ . This is important, because in some of our applications  $j_0 \rightarrow \infty$  as  $n \rightarrow \infty$ . This requires that resolution levels between, say,  $j = 0$ , and  $j = j_0$  be paid attention, too. However, the coefficients in these levels do not depend on the smoothing parameter, and thus these additional considerations only lead to additive rescaling of the relevant quantities and do not have essential impact on the statistical applications considered.

**A12.** These heuristic considerations are *pointwise*; this means that, in the context of Delyon and Juditsky's results, we must take  $p = 2$ ,  $\pi = \infty$ . This brings about the constraint  $s > 1/2$  (or  $s > d/2$  in the  $d$ -dimensional case).

**Proof of Lemma 1.** (See also [41-43].) Follows directly from the fact that in order to obtain the value of the  $K_2$ -functional, one needs to minimize the following expression:

$$\sum_k (\hat{\alpha}_{j_0 k} - \alpha_{j_0 k})^2 + \sum_{j=j_0}^{j_1} \sum_k (\hat{\beta}_{jk} - \beta_{jk})^2 + t^2 \sum_{j=j_0}^{j_1} \sum_k 2^{2js} \beta_{jk}^2 + (1+t^2) \sum_{j>j_1} \sum_k 2^{2js} \beta_{jk}^2$$

with respect to  $\alpha_{j_0 k}$ ,  $k \in Z$ ,  $\beta_{jk}$ ,  $j \geq j_0$ ,  $k \in Z$ . ■

**A13.** For the case when  $A \cap B$  may not be dense in  $A$  and/or  $B$ , see [43], Lemma 2.1. For the case when  $A$  and/or  $B$  are semi-Hilbert spaces, see

[43], Theorem 3.1.1. For the case when the cardinality of  $A$  is less than that of  $B$ , see [43], Theorem 2.1a (for Hilbert spaces), Theorem 3.1.1a (for semi-Hilbert spaces), and Remark 2.3. However, the generalizations with non-dense  $A \cap B$  are not needed for the purposes of Tikhonov regularization (see Section 4).

**A14.** By (5), the homogeneous and inhomogeneous models are equivalent on  $A \cap B = A \cap \dot{B} = B_{22}^s$ . For functions  $f$  with *heavy tails*, such that  $f \in \dot{B}_{22}^s \setminus B_{22}^s$ , the inhomogeneous model is more precise (see also [48], Remark 2.2.4.). In the sequel of this paper, compactly supported  $f$  will be considered only, so we shall be within the range of equivalence of the two models. We have chosen the homogeneous model because results about it are easier to directly compare to corresponding results for thresholding techniques, since in the homogeneous model, like with standard thresholding techniques, the empirical  $\alpha$ 's are left unchanged (see (7)).

**A15.** Asymptotic-minimax theory is beyond the scope of this paper,<sup>B2</sup> but we note that for  $\pi = p$  linear estimators do achieve the minimax rates (see [53], Theorems 3 and 4, case  $\epsilon > 0$ , as well as the sequence of two papers [47,48]), together with [85]. Moreover, it is possible that our estimators may achieve the asymptotic-minimax (within the "asymptopia") rates for functions in Besov spaces with  $p = q = 2$ ,  $2 \leq \pi \leq \infty$ ,  $0 < u \leq \infty$ , and also when  $\pi \neq p$  and  $\pi, p$  vary in a *neighbourhood of 2*. (A (very weak) necessary condition for this is that the estimator be *nonlinear*. In our case this condition is fulfilled: the minimizer of the smoothing parameter  $t$  with respect to cross validation (see Section 6) is nonlinear, hence the estimator (1'), with  $t$  estimated via cross validation, is also nonlinear.)

**A16.** If  $f \notin V_j$  for any  $j \in Z$  and  $f \in B_{22}^{s'}$ , but  $f \notin B_{2\infty}^s$  for any  $s > s'$ , then for any fixed  $s > s'$  there exists  $c_0 = c_0(f, \psi, s - s') > 0$  with  $\lim_{s \rightarrow s'+} c_0 = 0+$ , such that for any  $N$  there exists  $(j, k) \in Z^2$  with  $j \geq N$  so that  $\beta_{jk} \neq 0$  and  $|\beta_{jk}| \geq c_0 2^{js}$  hold. Here is an outline of the proof of this claim. Assume that its statement is not true. If for some  $N$  there is no  $(j, k)$  with  $j \geq N$  and  $\beta_{jk} \neq 0$ , this would contradict to  $f \notin V_j$  for any  $j \in Z$ . If for any  $c_0 > 0$

there exists  $N$  such that for any  $(j, k)$  with  $j \geq N$   $|\beta_{jk}| < c_0 2^{js}$  holds, then it can be shown that  $\|f\|_{B_{2\infty}^s} < \infty$ , which is in contradiction with  $f \notin B_{2\infty}^s$ . Finally, if there exists a decreasing sequence  $\{s_n\}$  such that  $c_0(f, \psi, s_n - s')$  remains bounded away from zero as  $n \rightarrow \infty$ , then it can be proved that there exists  $c_0(f, \psi, 0) > 0$  so that for any  $N$  there exists  $(j, k)$  with  $j \geq N$  and with  $|\beta_{jk}| \geq c_0(f, \psi, 0) 2^{-js'}$ . But then  $f \notin B_{22}^{s'}$ , which again is a contradiction with the claim's assumptions on  $f$ . This completes the proof of the claim.

**Proof of Theorem 1.** We give the proof for  $v^*$ . The statement about  $\tilde{v}$  can be obtained analogously, or can be derived from the asymptotic identity  $E(GFCV(v)) \sim ER(v) + \delta^2$  (see Theorem 2). Under the theorem's assumptions, it can be shown that for any  $v \in [0, \infty)$   $(\frac{d}{dv} ER)(v) \sim w'(v)$ , where

$$(A1) \quad w'(v) = 2 \sum_{j=j_0}^{j_1} \frac{2^{2js}}{(1 + v2^{2js})^3} \times \\ \times \sum_k [(v2^{2js} + 1 - \frac{1}{n} \sum_{i=1}^n \psi_{jk}(x_i)^2) \bar{\beta}_{jk}^2 - \frac{\delta^2}{n} (\frac{1}{n} \sum_{i=1}^n \psi_{jk}(x_i)^2)^2],$$

as  $n \rightarrow \infty$ . Here,

$$E\hat{\alpha}_{j_0k} = \bar{\alpha}_{j_0k} = \frac{1}{n} \sum_{i=1}^n f(x_i) \varphi_{j_0k}(x_i), \quad E\hat{\beta}_{jk} = \bar{\beta}_{jk} = \frac{1}{n} \sum_{i=1}^n f(x_i) \psi_{jk}(x_i), \\ E\hat{\alpha}_{j_0k}^2 = \bar{\alpha}_{j_0k}^2 + \frac{\delta^2}{n} \cdot \frac{1}{n} \sum_{i=1}^n \varphi_{j_0k}(x_i)^2, \quad E\hat{\beta}_{jk}^2 = \bar{\beta}_{jk}^2 + \frac{\delta^2}{n} \cdot \frac{1}{n} \sum_{i=1}^n \psi_{jk}(x_i)^2.$$

We outline the basic facts in this proof, as follows. (i) it holds

$$(A2) \quad |\bar{\beta}_{jk} - \beta_{jk}| \leq \|f\|_{B_{\infty\infty}^{s_1}} 2^{-j/2} n^{-s_1} (\|\psi\|_{L_1} + \frac{2^{j/2}}{n} \bigvee \psi) + \frac{2^{j/2}}{n} \|f\|_{L_\infty}$$

(for  $\hat{\alpha}_{j_0k}$  and  $\alpha_{j_0k}$   $\psi$  is replaced by  $\varphi$ ). An improvement of (A2) in terms of the average moduli of smoothness and A-spaces is considered in  $B^{12(b)}$ .

(ii) Orthogonality of the wavelet basis implies

$$(A3) \quad |\langle \psi_{jk}, \psi_{j'k'} \rangle - \frac{1}{n} \sum_{i=1}^n \psi_{jk}(x_i) \psi_{j'k'}(x_i)| \leq \frac{2^{(j+j')/2}}{n} (\bigvee \psi)^2$$

(for  $\varphi_{j_0 k} \cdot \varphi_{j_0 k'}$  and  $\varphi_{j_0 k} \cdot \psi_{j' k'}$   $(\bigvee \psi)^2$  is replaced by  $(\bigvee \varphi)^2$  and  $\bigvee \varphi \cdot \bigvee \psi$ , respectively).

(iii) Almost-diagonality (decorrelation) property of the wavelet transform with respect to the inner product  $\langle g, h \rangle_n := \frac{1}{n} \sum_{i=1}^n g(x_i) h(x_i)$ . Its essence is that the Gramian matrix of  $\{\varphi_{j_0 k}, \psi_{jk}\}$  with respect to  $\langle \cdot, \cdot \rangle_n$  is sparse because of the compactness of the wavelet support.

Let us compare the results obtained without and with the decorrelation property. The number  $N_j$  of non-zero  $\beta_{jk}$  on a fixed level does not exceed  $\sum_{-k \in \text{supp} \psi - 2^j \cdot \text{supp} f} 1 \leq c \cdot 2^j$ , where  $c$  depends on the lengths of the supports of  $f$  and  $\psi$  only. For fixed  $j$  and  $j'$  this estimate implies that the number  $N_{jj'}$  of non-zero products  $\beta_{jk} \cdot \beta_{j'k'} \cdot \langle \psi_{jk}, \psi_{j'k'} \rangle_n$  would be  $O(2^{j+j'})$ . On the other hand, by the decorrelation property,  $\langle \psi_{jk}, \psi_{j'k'} \rangle_n$  can be non-zero only if  $-k \in \text{supp} \psi - 2^{j-j'} \cdot \text{supp} \psi - 2^{j-j'} \cdot k'$ , i.e., only if  $|k - 2^{j-j'} \cdot k'| \leq c \cdot \max(2^{j-j'}, 1)$ ,  $c$  depending on the length of  $\text{supp} \psi$  only. This yields the more refined bound  $N_{jj'} = O(2^{\max(j, j')})$ . Estimating the quantity  $\sum_{j=j_0}^{j_1} \sum_{j'=j_0}^{j_1} N_{jj'}$  needed in the proof of (A1) via the former bound on  $N_{jj'}$  yields  $O(2^{2j_1})$ , while the latter bound implies  $\sum_{j=j_0}^{j_1} \sum_{j'=j_0}^{j_1} N_{jj'} \leq c[4(j_1 - j_0 - 1/2)2^{j_1} + 2^{j_0+1} + 1] = O(j_1 \cdot 2^{j_1})$ .

The proof of (A1) follows from (i-iii),  $2^{j_1} = O(\frac{n}{\ln n})$  and  $s_1 > 0$ .

Recalling that  $1 = \int \psi_{jk}(x)^2 dx$  and estimating the error of the quadrature formula for the latter integral yields (see (A3))  $|1 - \frac{1}{n} \sum_{i=1}^n \psi_{jk}(x_i)^2| \leq \frac{2^j}{n} \bigvee(\psi^2)$ . Fix  $v > 0$ . For  $n$  large enough,  $v2^{2j_s} + 1 - \frac{1}{n} \sum_{i=1}^n \psi_{jk}(x_i)^2 \geq \frac{1}{2} v2^{2j_s}$ , uniformly in  $j = j_0, \dots, j_1$ , where, with no loss of generality,  $j_0 \geq 0$ . By the lemma's assumptions about  $f$ , there exists  $(j', k')$  with  $j' \geq j_0$ , such that  $\beta_{j'k'} \neq 0$ . We shall be considering here the more difficult case  $j_0 \rightarrow \infty$ . It can be seen that, without loss of generality, one may assume that  $j' = j_0$ . By the sharpness of

$s'$  and  $s' < s$ , there exists<sup>A16</sup>  $c_0 : 0 < c_0 < \infty$ , depending on  $f$  and  $\psi$  only, such that  $|\beta_{j_0, k'}| \geq c_0 \cdot 2^{-j_0 s}$ . From here and (A2),  $|\bar{\beta}_{j_0, k'}| \geq |\beta_{j_0, k'}| - |\bar{\beta}_{j_0, k'} - \beta_{j_0, k'}| \geq c_0 \cdot 2^{-j_0 s} - O(2^{-j_0/2} \cdot n^{-s_1} + 2^{j_0/2} \cdot n^{-1}) = c_0 \cdot (1 - o(1)) \cdot 2^{j_0 s} > \frac{1}{2} \cdot c_0 \cdot 2^{-j_0 s}$  for sufficiently big  $n$ , by the theorem's assumptions on  $j_0$ . For  $n$  large enough  $w'(v)$  can now be bounded from below by

$$\begin{aligned} w'(v)/2 &> c_1 \cdot \frac{2^{4j_0 s} v}{(1 + v 2^{2j_0 s})^3} \cdot 2^{-2j_0 s} - \frac{\delta^2}{nv^3} \cdot \sum_{j=j_0}^{j_1} 2^{-4js} \cdot N_j \cdot (1 + O(\frac{2^j}{n}))^2 = \\ &= \frac{1}{v^2 2^{4j_0 s}} \left( \frac{c_1}{(1 + v^{-1} 2^{-2j_0 s})^3} - \frac{c_2}{nv} \sum_{j=j_0}^{j_1} 2^{j-4(j-j_0)s} \right) \geq \\ &\geq \frac{1}{v^2 2^{4j_0 s}} \left( \frac{c_1}{(1 + v^{-1} 2^{-2j_0 s})^3} - \frac{c_3}{v} \cdot \frac{2^{j_1}}{n} \right) = \frac{1}{v^2 2^{4j_0 s}} \left( \frac{c_1}{(1 - o(1))^3} - \frac{1}{v} \cdot o(1) \right) > 0 \end{aligned}$$

for sufficiently large  $n$ . Therefore, by continuity and asymptotic equivalence to  $(\frac{d}{dv} ER)(v)$  of  $w'(v)$ , there exists  $C = C(f, \psi) \in (0, \infty)$ , such that for  $n$  large enough  $\inf_{v \in R} w(v)$  and  $\inf_{v \in R} ER(v)$  are both achieved on  $[0, C]$ . By continuity of  $ER(v)$  and compactness of  $[0, C]$ ,  $v^*$  exists and, moreover,  $v^* \in [0, C]$  and either  $v^* = 0$  or  $(\frac{d}{dv} ER)(v^*) = 0$ . ■

**Proof of Theorem 2.** The proof is similar in spirit to the derivations in Section 4.4 of [117] (cf. [49]). First of all,

$$\begin{aligned} (A4) \quad E(GFCV(v)) &= (1 + \bar{h}(v))^2 \times \\ &\times [\delta^2 + \frac{1}{n} E \sum_{i=1}^n (f(x_i) - \tilde{f}_v(x_i))^2 - \frac{2}{n} E \sum_{i=1}^n (y_i - f(x_i))(\tilde{f}_v(x_i) - f(x_i))] = \\ &= (1 + \bar{h}(v))^2 [\delta^2 + ER(v) - \frac{2}{n} E \sum_{i=1}^n \{\epsilon_i(\tilde{f}_v(x_i) - f(x_i))\}]. \end{aligned}$$

Now we show that  $\bar{h}(v) \rightarrow 0$  as  $n \rightarrow \infty$ , uniformly in  $v \geq 0$ . Invoking again a familiar bound for the error of the quadrature formulae involved (cf.

the proof of Theorem 1), by the orthonormality of the wavelet basis and the compactness of the support of  $f$ ,

$$0 \leq \bar{h}(v) \leq c \left[ \frac{2^{j_0}}{n} \left( 1 + \frac{2^{j_0}}{n} \bigvee(\varphi^2) \right) + \frac{1}{n} \sum_{j=j_0}^{j_1} 2^j \left( 1 + \frac{2^j}{n} \bigvee(\psi^2) \right) \right] \leq c_1 \cdot \frac{2^{j_1}}{n} = o(1),$$

as  $n \rightarrow \infty$ , uniformly in  $v \geq 0$ . Here  $c_1 = c_1(\varphi, \psi, f)$ .

Fix  $v \geq 0$  and denote by  $\tilde{\epsilon}_i(v) = \tilde{f}_v(x_i) - f(x_i)$ . Since  $E\epsilon_i = 0$  and  $\epsilon_i$  is independent of any  $y_l, l \neq i$  we can evaluate:

(A5)

$$\begin{aligned} E[\epsilon_i \tilde{\epsilon}_i(v)] &= E\left\{ \epsilon_i \left[ \sum_k \sum_{l=1}^n \frac{\varphi_{j_0 k}(x_l) y_l}{n} \varphi_{j_0 k}(x_i) + \sum_k \sum_{j=j_0}^{j_1} \sum_{l=1}^n \frac{\psi_{j k}(x_l) \psi_{j k}(x_i) y_l}{n(1+v2^{2js})} \right] \right\} = \\ &= \sum_k E \frac{\varphi_{j_0 k}(x_i)^2 y_i \epsilon_i}{n} + \sum_k \sum_{j=j_0}^{j_1} E \frac{\psi_{j k}(x_i)^2 y_i \epsilon_i}{n(1+v2^{2js})} = \\ &= \frac{\delta^2}{n} \left[ \sum_k \varphi_{j_0}(x_i)^2 + \sum_{j=j_0}^{j_1} \sum_k \frac{\psi_{j k}(x_i)^2}{1+v2^{2js}} \right] = \delta^2 h_{ii}(v) \end{aligned}$$

Substitution from (A5) into (A4) yields  $E(GFCV(v)) = (1 + \bar{h}(v))^2 (\delta^2 + ER(v) - 2\delta^2 \bar{h}(v))$  from which also  $E(GFCV(v)) \sim ER(v) + \delta^2$  follows.

To show the required consistency, we first note that

$$(A6) \quad \frac{ER(v) - E(GFCV(v)) + \delta^2}{ER(v)} = -2\bar{h}(v) - \bar{h}(v)^2 + \frac{\bar{h}(v)^2 \delta^2}{ER(v)} (3 + 2\bar{h}(v))$$

holds. Denoting  $\mathbf{f} = [f(x_1), f(x_2), \dots, f(x_n)]'$ ,  $\mathbf{I}$  - the identity on  $R^n$ ,  $\|\cdot\|$  - the usual Hilbert norm on  $R^n$ ,  $\mathbf{f}_v = [\tilde{f}_v(x_1), \tilde{f}_v(x_2), \dots, \tilde{f}_v(x_n)]'$ ,  $\epsilon = [\epsilon_1, \epsilon_2, \dots, \epsilon_n]'$  we have

$$\begin{aligned} ER(v) &= E \frac{1}{n} \|\mathbf{f}_v - \mathbf{f}\|^2 = E \frac{1}{n} \|(\mathbf{I} - \mathbf{H}(v))\mathbf{f} - \mathbf{H}(v)\epsilon\|^2 = \\ &= \frac{1}{n} \|(\mathbf{I} - \mathbf{H}(v))\mathbf{f}\|^2 + \delta^2 \mu(v) \geq \delta^2 \mu(v). \end{aligned}$$

Therefore, (A6) yields:

$$\frac{|ER(v) - E(GFCV(v)) + \delta^2|}{ER(v)} \leq \frac{\bar{h}(v)^2}{\mu(v)}(3 + 2\bar{h}(v)) = \xi(v).$$

Now, by Lemma 2,  $\lim_{n \rightarrow \infty} \frac{h(v)^2}{\mu(v)} = 0$ , hence  $\xi(v) \rightarrow 0$ . Then it is easily seen that, by definition of  $\tilde{v}$  and  $v^*$ ,

$$(1 - \xi(\tilde{v}))ER(\tilde{v}) \leq E(GFCV(\tilde{v})) - \delta^2 \leq E(GFCV(v^*)) - \delta^2 \leq (1 + \xi(v^*))ER(v^*);$$

again by definition of  $\tilde{v}$  and  $v^*$ ,  $\frac{ER(\tilde{v})}{ER(v^*)} \geq 1$ , for any  $n \in N$ . Hence,  $1 \leq \frac{ER(\tilde{v})}{ER(v^*)} \leq \frac{1 + \xi(v^*)}{1 - \xi(\tilde{v})}$  holds. ■

**Proof of Lemma 2.** Fix  $v \geq 0$ . Denote  $m := \text{rank} \mathbf{H}(v)$ ;  $m \leq n$ . On the one hand,  $m = O(2^{j_1}) = o(n)$  holds. Indeed, while the domain of the linear operator  $\mathcal{L}_v(\mathbf{y}) := \mathbf{H}(v)\mathbf{y}$  is  $R^n$ , the range of  $\mathcal{L}_v$  is the linear span of the set of  $n$ -dimensional vectors

$$V_{j_1}(\text{supp} f) := \{(\varphi_{j_0 k}(x_1), \dots, (\varphi_{j_0 k}(x_n))', k : \text{supp} \varphi_{j_0 k} \cap \text{supp} f \neq \emptyset\} \cup$$

$$\cup \{(\psi_{j k}(x_1), \dots, (\psi_{j k}(x_n))', j = j_0, \dots, j_1, k : \text{supp} \psi_{j k} \cap \text{supp} f \neq \emptyset\},$$

which has dimension  $O(2^{j_1})$ .

On the other hand,

$$\frac{\bar{h}(v)^2}{\mu(v)} = \frac{(\text{tr} \mathbf{H}(v))^2}{n \cdot \text{tr}(\mathbf{H}'(v)\mathbf{H}(v))} = \frac{m}{n} \cdot \frac{(\frac{1}{m} \sum_{i=1}^m s_i)^2}{\frac{1}{m} \sum_{i=1}^m s_i^2},$$

where  $s_i$ ,  $i = 1, \dots, n$ , are the  $s$ -numbers of  $\mathbf{H}(v)$  (see, e.g., [51] and the references therein), and  $s_{m+1} = \dots = s_n = 0$ . Therefore, by the inequality between the arithmetic and the quadratic mean,  $\frac{\bar{h}(v)^2}{\mu(v)} \leq \frac{m}{n} = o(1)$  as  $n \rightarrow \infty$ , uniformly in  $v \geq 0$ . ■



**Proof of Theorem 3.** We prove necessity when  $j_0 = O(1)$  first. For any fixed  $(j, k)$ ,  $j \geq j_0$ , consistency of  $\tilde{f}_v$  in  $B_{22}^\sigma$ ,  $\sigma \geq 0$ , implies  $E(\tilde{\beta}_{jk} - \beta_{jk})^2 \rightarrow 0$ ,  $n \rightarrow \infty$ . Choose  $(j', k') : \beta_{j'k'} \neq 0$ . Simple computations show that

$$E(\tilde{\beta}_{jk} - \beta_{jk})^2 = \frac{1}{(1 + v2^{2js})^2} [(\bar{\beta}_{jk} - \beta_{jk})^2 + \frac{\delta^2}{n} \cdot \frac{1}{n} \sum_{i=1}^n \psi_{jk}(x_i)^2 - 2v2^{2js} \beta_{jk} (\bar{\beta}_{jk} - \beta_{jk}) + v^2 2^{4js} \beta_{jk}^2].$$

By Riemann integrability of  $f\psi_{jk}$ ,  $\bar{\beta}_{jk} - \beta_{jk} = o(1)$ ; by Riemann integrability of  $\psi_{jk}^2$  (or by the stronger (A3)),  $\frac{\delta^2}{n} \cdot \frac{1}{n} \sum_{i=1}^n \psi_{jk}(x_i)^2 = O(\frac{1}{n}) = o(1)$ ; by the boundedness of  $v$ ,  $2v2^{2js} \beta_{jk} (\bar{\beta}_{jk} - \beta_{jk}) = v \cdot o(1) = o(1)$  and  $\frac{1}{(1+v2^{2js})^2}$  is bounded away from zero. Therefore,

$$E(\tilde{\beta}_{jk} - \beta_{jk})^2 = \frac{1}{(1 + v2^{2js})^2} (o(1) + v^2 2^{4js} \beta_{jk}^2)$$

is only possible if  $v \rightarrow 0$ , which proves the necessity claim. Admissibility of the choice  $v = v_n^*$  here follows from the boundedness of  $v_n^*$ , proved in Theorem 1.

Now we prove sufficiency. After some computations, utilizing, as usual, the wavelet-coefficient equivalent norm of  $B_{22}^\sigma$ , and after a simple bound from above, by the inequality between the arithmetic and quadratic mean,

$$\begin{aligned} E \|\tilde{f}_v - f\|_{B_{22}^\sigma}^2 &= E \|\tilde{f}_v - E\tilde{f}_v\|_{B_{22}^\sigma}^2 + \|E\tilde{f}_v - f\|_{B_{22}^\sigma}^2 \leq \\ &\leq \frac{\delta^2}{n} \sum_k \frac{1}{n} \sum_{i=1}^n \varphi_{j_0 k}(x_i)^2 + \frac{\delta^2}{n} \sum_{j=j_0}^{j_1} \frac{2^{2j\sigma}}{(1 + v2^{2js})^2} \sum_k \frac{1}{n} \sum_{i=1}^n \psi_{jk}(x_i)^2 + \\ &+ \sum_k (\bar{\alpha}_{j_0 k} - \alpha_{j_0 k})^2 + 2 \sum_{j=j_0}^{j_1} \frac{2^{2j\sigma}}{(1 + v2^{2js})^2} \sum_k (\bar{\beta}_{jk} - \beta_{jk})^2 + \\ &+ 2 \sum_{j=j_0}^{j_1} 2^{2j\sigma} \frac{v^2 2^{4js}}{(1 + v2^{2js})^2} \sum_k \beta_{jk}^2 + \sum_{j=j_1+1}^{\infty} 2^{2j\sigma} \sum_k \beta_{jk}^2. \end{aligned}$$

Applying (A1,A2), the inequality between the arithmetic and quadratic mean,  $v \geq 0$ , utilizing  $\frac{s'-\sigma}{2s} \in (0, 1/2) \subset (0, 1)$ , which yields the inequality

$$\frac{v^2 2^{4js}}{(1 + v 2^{2js})^2} \leq \left( \frac{v^2 2^{4js}}{1} \right)^{\frac{s'-\sigma}{2s}} \cdot \left( \frac{v^2 2^{4js}}{v^2 2^{4js}} \right)^{1 - \frac{s'-\sigma}{2s}} = v^{\frac{s'-\sigma}{s}} \cdot 2^{2j(s'-\sigma)},$$

invoking also the bound

$$\sum_{j=j_1+1}^{\infty} 2^{2j\sigma} \sum_k \beta_{jk}^2 \leq 2^{2(\sigma-s')} 2^{2j_1(\sigma-s')} \sum_{j=j_1+1}^{\infty} 2^{2js'} \sum_k \beta_{jk}^2 \leq c 2^{2j_1(\sigma-s')} \|f\|_{B_{22}^{s'}}^2,$$

as well as the familiar upper bound on  $N_j$ , after obtaining

$$\sum_{j=j_0}^{j_1} 2^{2j\sigma} \frac{v^2 2^{4js}}{(1 + v 2^{2js})^2} \sum_k \beta_{jk}^2 \leq v^{\frac{s'-\sigma}{s}} \sum_{j=j_0}^{j_1} 2^{2js'} \sum_k \beta_{jk}^2 \leq v^{\frac{s'-\sigma}{s}} \|f\|_{B_{22}^{s'}}^2,$$

we arrive at

$$\begin{aligned} \|\tilde{f}_v - f\|_{B_{22}^{\sigma}} &\leq c(f, \varphi) \frac{\delta^2}{n} 2^{j_0} \left(1 + \frac{2^{j_0}}{n}\right) + c(f, \psi) \frac{\delta^2}{n} \sum_{j=j_0}^{j_1} 2^{j(1+2\sigma)} \left(1 + \frac{2^j}{n}\right) + \\ &+ c_1(f, \varphi) 2^{j_0} \left[ \frac{2^{-j_0}}{n^{2s_1}} \left(1 + \frac{2^{j_0}}{n^2}\right) + \frac{2^{j_0}}{n^2} \right] + c_1(f, \psi) \sum_{j=j_0}^{j_1} 2^{j(1+2\sigma)} \left[ \frac{2^{-j}}{n^{2s_1}} \left(1 + \frac{2^j}{n^2}\right) + \frac{2^j}{n^2} \right] + \\ &+ c_{s'}(f) v^{\frac{s'-\sigma}{s}} + c_{s',\sigma}(f) 2^{-2j_1(s'-\sigma)} \leq \\ &\leq c_{s',\sigma,\delta}(f, \varphi, \psi) \left( \frac{2^{j_0}}{n} + \frac{2^{j_1(1+2\sigma)}}{n} + \frac{2^{2j_1\sigma}}{n^{2s_1}} + v^{\frac{s'-\sigma}{s}} + 2^{-2j_1(s'-\sigma)} \right) = o(1), \end{aligned}$$

when  $2^{j_1} = o(n^{\min(\frac{21}{\sigma}, \frac{1}{1+2\sigma})})$ ,  $j_0 \leq j_1$ ,  $j_1 \rightarrow \infty$  and  $v \rightarrow 0$ .

Admissibility of the choice  $v = v_n^*$  follows from Theorem 4. ■

**Proof of Theorem 4.** We give the proof for  $v^*$ . The statement about  $\tilde{v}$  follows from the asymptotic equivalence  $E(GFCV(v)) - \delta^2 \sim ER(v)$ ,  $v \geq 0$  (see the proof of Theorem 2). Indeed, by the minimizing properties of  $\tilde{v}_n$

and  $v_n^*$  and the continuous dependence of  $E(GFCV(v))$  and  $ER(v)$  on  $v$ , this asymptotic equality implies (e.g., by proving contradiction with the opposite) that  $\tilde{v}_n$  and  $v_n^*$  have the same density points.

From the proof of Theorem 1 it follows that, for  $n$  large enough,  $v^* = v_n^*$  exists and  $v_n^* \in S_{n,1} := \{0\} \cup \{v = v_n \in [0, C] : (\frac{d}{dv}ER)(v) = 0\}$ , where  $C$  is defined in the same proof. Consider also  $S_{n,2} := \{0\} \cup \{v = v_n \in [0, C] : w'(v) = 0\}$ . By Bolzano-Weierstrass theorem, every sequence  $v_n \in S_{n,\nu}$  has at least one density point and every such point is in  $[0, C]$ ,  $\nu = 1, 2$ . Consider  $S_\nu = \bigcup_{n \in N} S_{n,\nu}$ ,  $\nu = 1, 2$ . By definition of  $S_{n,\nu}$ ,  $v = 0$  is a density point of  $S_{n,\nu}$ ,  $\nu = 1, 2$ . By (A1),  $S_1$  and  $S_2$  have the same density points. To prove  $v^* \rightarrow 0$ , it now suffices to show that the only density point of  $S_2$  is  $v = 0$ .

As in the proof of Theorem 1, we can assume that there exists  $(j_0, k_0)$  such that  $\beta_{j_0 k_0} \neq 0$ . Also, without loss of generality,  $C \geq 1$  and  $j_0 \geq 0$ . Assume that  $v^* = v_n^*$  remains bounded away from 0 for infinitely many values of  $n$ . Then, for these particular values of  $n$ , the same argument as in the proof of Theorem 1 yields

$$\begin{aligned} \frac{w'(v)}{2} &> c_1 \frac{v 2^{2j_0 s}}{(1 + v 2^{2j_0 s})^3} - \frac{c_2}{n v^3} \sum_{j=j_0}^{j_1} 2^{j(1-4s)} \geq \\ &\geq \frac{c_1}{8C^3} 2^{-4j_0 s} v - \frac{c_2}{n v^3} \sum_{j=j_0}^{j_1} 2^{j(1-4s)} = w_0(v). \end{aligned}$$

The equation  $w_0(v) = 0$  has a unique real positive root

$$v_0 : v_0^4 = \frac{c_3}{n} \sum_{j=j_0}^{j_1} 2^{j-4(j-j_0)s} \leq c_4 \cdot \frac{2^{j_1}}{n} = o(1)$$

as  $n \rightarrow \infty$ . Since  $w'(v) > 0$  for  $v \geq C$  and all sufficiently large  $n$ , for the same  $n$  the largest real zero of any lower bound of  $w'(v)$  is necessarily an upper bound for all real zeros of  $w'(v)$ . In particular,  $v_0$  is an upper bound for the zeros of

$w'(v)$  for any  $n$  considered. Therefore, contradiction has been achieved with the assumption that  $v^*$  remains bounded away from  $v = 0$  for infinitely many  $n$ . ■

**Proof of Corollary 1.** Follows from Theorem 4 and the sufficiency part of Theorem 3, applied to  $\tilde{v}$ . ■

**A17.** The proof of Lemma 2 can be simplified by relying on *selfadjointness* (cf. [49]). Indeed, in our case the influence matrix  $\mathbf{H}(v)$  is symmetric positive semi-definite, therefore, its  $s$ -numbers are just its eigenvalues (all nonnegative), arranged in order of decreasing magnitude. We have chosen a more general consideration because of envisaging an extension of our model for *biorthogonal wavelets* and *wavelet packets*.<sup>B18</sup> In these generalized cases the problem becomes non-selfadjoint. Our general proof shows that Lemma 2 is valid in this more general context, too.

**Proof of Theorem 5.** We outline briefly the proof for  $v = \tilde{v}$  and give a detailed proof of the more difficult case  $v = v^*$ . As earlier, we consider the more difficult case for  $j_0: j_0 \rightarrow \infty$ .

For any  $n = 1, 2, \dots$ ,

$$\left(\frac{d}{dv} MISE\right)(0) = -\frac{2}{n} \sum_{j=j_0}^{j_1} 2^{2js} \sum_k \left[ \int \psi_{jk}^2 f - \left( \int \psi_{jk} f \right)^2 \right] < 0,$$

by Cauchy-Schwartz inequality. Since  $s > s'$  and the index  $s'$  is sharp, one may assume that there exists (cf. the proof of Theorem 1 and Theorem 4)  $(j_0, k_0): \beta_{j_0 k_0} \neq 0$  and  $|\beta_{j_0 k_0}| \geq c_0 2^{-j_0 s}$  for some  $c_0 > 0$ . Hence, utilizing also  $\int \psi_{jk}^2 f \leq \|f\|_{L_\infty}$ , and invoking the bound  $N_j \leq c 2^j$  again,

$$\begin{aligned} \frac{1}{2} \left(\frac{d}{dv} MISE\right)(v) &\geq \frac{c_0^2 2^{2j_0 s} v}{(1 + v 2^{2j_0 s})^3} - \frac{\|f\|_{L_\infty}}{n} \sum_{j=j_0}^{j_1} \frac{2^{j(1+2s)}}{(1 + v 2^{2js})^3} \geq \\ &\geq \frac{1}{v^2 2^{4j_0 s}} \left( \frac{c_1}{(1 + v^{-1} 2^{-2j_0 s})^3} - \frac{c_2}{nv} \sum_{j=j_0}^{j_1} 2^{j-4(j-j_0)s} \right) \geq \end{aligned}$$

$$\geq \frac{1}{v^2 2^{4j_0 s}} \left( \frac{c_1}{(1 + v^{-1} 2^{-2j_0 s})^3} - \frac{c_3}{nv} 2^{j_1} \right) = \frac{1}{v^2 2^{4j_0 s}} \left( \frac{c_1}{(1 + v^{-1} 2^{-2j_0 s})^3} - \frac{o(1)}{v} \right) > 0$$

for  $v$  and  $n$  large enough. Now existence of  $C < \infty$  and of  $\tilde{v} : 0 < \tilde{v} < C$  follows in the same way as in Theorem 1.

The basic idea of the proof about  $v^*$  is the same. Denote  $w(v) := M(\tilde{f}_v)$ . For any  $n = 2, 3, \dots$ , by the inequality between the arithmetic and quadratic mean,

$$w'(0) = -\frac{2}{n-1} \sum_{j=j_0}^{j_1} 2^{2js} \sum_k \left[ \frac{1}{n} \sum_{i=1}^n \psi_{jk}(X_i)^2 - \left( \frac{1}{n} \sum_{i=1}^n \psi_{jk}(X_i) \right)^2 \right] \leq 0$$

and this inequality is strict almost surely. We proceed to study the behaviour of  $w'(v)$  for large  $v$ . By  $s' < s$ , there exists  $s_0 : 0 < s' < s_0 < \frac{s'+s}{2} < s$ . By sharpness of  $s'$ , one can find  $(j_0, k_0)$  with  $|\beta_{j_0 k_0}| \geq c_0 2^{-j_0 s_0}$ , for some  $c_0 = c_0(f, \psi, s_0 - s') > 0$ . Then, after computations, omitting positive summands, bounding positive terms from below and negative ones from above, one obtains

$$\begin{aligned} \frac{1}{2} E w'(v) &\geq \frac{v 2^{4j_0 s}}{(1 + v 2^{2j_0 s})^3} \beta_{j_0 k_0}^2 - \frac{c \|f\|_{L_\infty}}{n} \sum_{j=j_0}^{j_1} \frac{2^{j(1+2s)}}{(1 + v 2^{2js})^2} \geq \\ &\geq \frac{1}{v^2 2^{2j_0(s+s_0)}} \left( \frac{c_1}{(1 + v^{-1} 2^{-2j_0 s})^3} - \frac{c_2}{n} 2^{2j_0(s+s_0)} \sum_{j=j_0}^{j_1} 2^{j(1-2s)} \right). \end{aligned}$$

*Case*  $0 < s < 1/2$ . By conditions on  $j_0$  and  $j_1$ ,  $n^{-1} 2^{2j_0(s+s_0)} \sum_{j=j_0}^{j_1} 2^{j(1-2s)} = O(n^{-\frac{2(s-s_0)}{1+2s}}) = o(1)$ .

*Case*  $s = 1/2$ . Again by the conditions on  $j_0$  and  $j_1$ ,  $j_1 - j_0 + 1 = o(\ln n)$  and  $n^{-1} 2^{2j_0(s+s_0)} (j_1 - j_0 + 1) = O(n^{-\frac{1-2s_0}{1+2s}}) \cdot o(\ln n) = o(1)$ , since  $s_0 < s = 1/2$ .

*Case*  $s > 1/2$ . Now  $n^{-1} 2^{2j_0(s+s_0)} \sum_{j=j_0}^{j_1} 2^{j(1-2s)} = O(n^{-\frac{2(s-s_0)}{1+2s}}) = o(1)$ .

Therefore, in all three cases,

$$(A7) \quad E w'(v) \geq \frac{c_1}{v^2 2^{2j_0(s+s_0)}}$$

with certain positive constant  $c_1$  holds for sufficiently large  $v$  and  $n$  big enough.

On the other hand, by the triangle inequality, and after some computations,

$$(A8) \quad \frac{1}{2} E|w'(v) - Ew'(v)| \leq \sum_{j=j_0}^{j_1} \frac{2^{2js}}{(1 + v2^{2js})^2} \times \\ \times \sum_k [(E|\hat{\beta}_{jk}^2 - \beta_{jk}^2|)(\frac{v2^{2js}}{1 + v2^{2js}} + \frac{1}{n-1}) + (\int \psi_{jk}^2 f)(\frac{v2^{2js}}{n(1 + v2^{2js})} + \frac{2 + 1/n}{n-1})].$$

By Cauchy-Schwartz inequality, and by  $(|a| + |b|)^{1/2} \leq |a|^{1/2} + |b|^{1/2}$ ,

$$E|\hat{\beta}_{jk}^2 - \beta_{jk}^2| \leq (E(\hat{\beta}_{jk} - \beta_{jk})^2)^{1/2} (E(\hat{\beta}_{jk} + \beta_{jk})^2)^{1/2} \leq \\ \leq \frac{1}{\sqrt{n}} (\int \psi_{jk}^2 f)^{1/2} (2|\beta_{jk}| + \frac{1}{\sqrt{n}} (\int \psi_{jk}^2 f)^{1/2}).$$

Substituting this into (A8), after computations utilizing  $\int \psi_{jk}^2 f \leq \|f\|_{L_\infty}$ , yields

$$E|w'(v) - Ew'(v)| \leq c_4 \sum_{j=j_0}^{j_1} \frac{2^{2js}}{(1 + v2^{2js})^2} \sum_k (\frac{v2^{2js}}{1 + v2^{2js}} \frac{|\beta_{jk}|}{\sqrt{n}} + \frac{1}{n}) \leq \\ \leq \frac{c_4}{v^2} (\sum_{j=j_0}^{j_1} 2^{-2js} \sum_k \frac{|\beta_{jk}|}{\sqrt{n}} + \frac{c_5}{n} \sum_{j=j_0}^{j_1} 2^{j(1-2s)}) = \frac{c_4}{v^2} (I_1 + I_2),$$

where  $c_4 = c_4(f, \psi)$ . Therefore, by (A7) and Markov inequality,

$$P\{w'(v) < \frac{c_1}{2v2^{2j_0(s+s_0)}}\} \leq P\{|w'(v) - Ew'(v)| > \frac{c_1}{2v2^{2j_0(s+s_0)}}\} \leq \\ \leq \frac{2c_4}{c_1} 2^{2j_0(s+s_0)} . O(I_1 + I_2).$$

By Cauchy-Schwartz inequality, and by the bound on  $N_j$  (cf. the proof of Theorem 1),  $I_1 \leq \frac{c_4}{\sqrt{n}} (\sum_{j=j_0}^{j_1} 2^{j(1-2(2s+s'))})^{1/2} \|f\|_{B_{22}^{s'}}.$  Considering the three cases  $2s + s' < 1/2$ ,  $2s + s' = 1/2$ ,  $2s + s' > 1/2$ , and noting that in the first

two of these cases necessarily  $s < 1/2$  holds, it turns out that in all three cases  $2^{2j_0(s+s_0)}I_1 = O(n^{2^{2j_0-(s+s')/2}})$  holds, and this tends to zero as  $n \rightarrow \infty$  by the choice of  $s_0$ . As for the term  $2^{2j_0(s+s_0)}I_2$ , it was already shown while deriving the lower bound for  $Ew'(v)$  that this term tends to zero, too. Let  $C : 0 < C < \infty$ , be the constant obtained in the proof of the case about  $\bar{v}$ . Now we see that for  $v \geq C$

$$P\{w'(v) \leq 0\} \leq P\{w'(v) < \frac{c_1}{2v^{2^{2j_0(s+s_0)}}}\} \rightarrow 0$$

as  $n \rightarrow \infty$ . This implies the existence of  $v^*$  by already familiar continuity/compactness argument. ■

**A18.** In the case  $0 < s \leq 1/2$  (the range of  $s$  for which  $B_{22}^s$  contains discontinuous functions), the condition on  $j_1$  in Theorem 5 seems to be restrictive:  $2^{j_0} \leq 2^{j_1} < C_1 2^{j_0}$  indicates a rather narrow range for  $j_1$ . But first, Theorem 5 only gives sufficient conditions for the existence of the minimizer  $v^*$ . Second, a close inspection of the proof of this theorem shows that its statement in the case  $0 < s \leq 1/2$  remains true also for a broader range of  $j_1$ : it suffices to have  $2^{j_1} = o(n^{\frac{1-2s_0}{1+2s}} \cdot 2^{j_0})$  for any  $s_0$  such that  $s' < s_0 < \frac{s'+s}{2}$ . Third, the theorem may be sharpened by admitting a broader range for  $j_1$ , if one considers a more stringent notion of sharpness of  $s'$ , namely,  $f$  is such that  $f \in B_{22}^{s'}$  and  $f \notin B_{2q}^{s'}$  for any  $q < 2$ . The proof of this sharpened version is similar.

**Proof of Theorem 6.** In this proof,  $C$  denotes a positive constant which may vary along the lines. Fix  $v \geq 0$ . First, let us evaluate  $MISE(v)$  from below. Using Parseval's identity and neglecting some nonnegative summands, we get:

$$\begin{aligned} MISE(v) \geq & \frac{1}{n^2} \sum_k E\left\{\sum_{i=1}^n \varphi_{j_0 k}(X_i)\right\}^2 + \sum_{j=j_0}^{j_1} \sum_k \frac{1}{n^2} E\left\{\sum_{i=1}^n \psi_{jk}(X_i)\right\}^2 \frac{1}{(1+v^{2^{2js}})^2} \\ & - \sum_k \alpha_{j_0 k}^2 - \sum_{j=j_0}^{j_1} \sum_k \frac{\beta_{jk}^2 (1-v^{2^{2js}})}{1+v^{2^{2js}}} \end{aligned}$$

Since  $\frac{1-v2^{2js}}{1+v2^{2js}} < \frac{1}{(1+v2^{2js})^2}$ , we can continue the above inequality to obtain:

$$\begin{aligned} MISE(v) &\geq \\ \sum_k \frac{1}{n} \{ \int \varphi_{j_0 k}^2 f - (\int \varphi_{j_0 k} f)^2 \} &+ \sum_{j=j_0}^{j_1} \sum_k \frac{1}{n(1+v2^{2js})^2} \{ \int \psi_{jk}^2 f - (\int \psi_{jk} f)^2 \} \geq \\ &\geq \frac{1}{n} \{ \sum_k \int \varphi_{j_0 k}^2 f + \sum_{j=j_0}^{j_1} \sum_k \frac{1}{(1+v2^{2js})^2} \int \psi_{jk}^2 f - \|f\|_2^2 \}. \end{aligned}$$

Now we invoke Meyer's lemma (Lemma 1 in [110] - for the application of this lemma the periodic extension onto  $R$  of all relevant compactly supported functions is considered, with period equal to or greater than the support's diameter). By this lemma, and in view of  $\int f = 1$ ,

$$\begin{aligned} 2^{-j_0} \int f(x) \sum_k \varphi_{j_0 k}^2(x) dx &= \\ = \int f(x) \sum_k \varphi^2(2^{j_0} x - k) dx &\rightarrow_{j_0 \rightarrow \infty} \sum_k \int_0^1 \varphi^2(x - k) dx = C, \quad C = C(\varphi). \end{aligned}$$

Utilizing the compactness of  $\text{supp } f$  and  $j_0 \rightarrow \infty$  in a similar consideration of the expressions of type  $2^{-j} \int f(x) \sum_k \psi_{jk}^2(x) dx$ , we arrive at

$$(A9) \quad MISE(v) \geq \frac{2^{j_0} C}{n} \{ 1 + \sum_{j=j_0}^{j_1} \frac{2^{j-j_0}}{(1+v2^{2js})^2} \}, \quad C = C(\varphi, \psi),$$

for sufficiently large  $n$ , so that  $2^{-j_0} \|f\|_{L_2}^2$  is smaller than, say,  $1/2$ . Next, we consider  $\xi_v = M(\tilde{f}_v) - MISE(v) + T_n$  with  $T_n$  as defined in the premises of the theorem. Note that  $E\xi_v = 0$  and that  $T_n$  does not depend on  $v$ . After further computations, we obtain

$$\xi_v = -\frac{2}{n(n-1)} \sum_k \left[ \left\{ \sum_{i=1}^n \varphi_{j_0 k}(X_i) \right\}^2 - \sum_{i=1}^n \varphi_{j_0 k}^2(X_i) \right] -$$



$$\begin{aligned}
& -\frac{2}{n(n-1)} \sum_{j=j_0}^{j_1} \sum_k \frac{\{\sum_{i=1}^n \psi_{jk}(X_i)\}^2 - \sum_{i=1}^n \psi_{jk}^2(X_i)}{1+v2^{2js}} - \\
& -\frac{n+2}{n} \sum_k \alpha_{j_0k}^2 + \sum_{j=j_0}^{j_1} \sum_k \left\{ \frac{2\beta_{jk}^2}{1+v2^{2js}} - \frac{\beta_{jk}^2}{(1+v2^{2js})^2} \right\} + \\
& + \sum_k \hat{\alpha}_{j_0k}^2 + \sum_{j=j_0}^{j_1} \sum_k \frac{\hat{\beta}_{jk}^2}{(1+v2^{2js})^2} - \\
& - \sum_k \text{Var}(\hat{\alpha}_{j_0k}) - \sum_{j=j_0}^{j_1} \sum_k \frac{\text{Var}(\hat{\beta}_{jk})}{(1+v2^{2js})^2} + \frac{2(n+1)}{n} \sum_k \hat{\alpha}_{j_0k} \alpha_{j_0k}.
\end{aligned}$$

Substituting  $\text{Var}(\hat{\alpha}_{j_0k}) = \frac{1}{n} \{ \int \varphi_{j_0k}^2 f - (\int \varphi_{j_0k} f)^2 \}$ ,  $\text{Var}(\hat{\beta}_{jk}) = \frac{1}{n} \{ \int \psi_{jk}^2 f - (\int \psi_{jk} f)^2 \}$  leads to

$$\begin{aligned}
(A10) \quad \xi_v &= \left( -\frac{2}{n(n-1)} + \frac{1}{n^2} \right) \sum_k \sum_{i=1}^n \sum_{l=1, l \neq i}^n \varphi_{j_0k}(X_i) \varphi_{j_0k}(X_l) + \\
& + \sum_{j=j_0}^{j_1} \left\{ -\frac{2}{n(n-1)} + \frac{1}{n^2(1+v2^{2js})} \right\} \sum_k \sum_{i=1}^n \sum_{l=1, l \neq i}^n \frac{\psi_{jk}(X_i) \psi_{jk}(X_l)}{1+v2^{2js}} \\
& - \frac{n+1}{n} \sum_k (\int \varphi_{j_0k} f)^2 + \sum_{j=j_0}^{j_1} \sum_k (\int \psi_{jk} f)^2 \left\{ \frac{2}{1+v2^{2js}} - \frac{n-1}{n(1+v2^{2js})^2} \right\} + \\
& + \frac{2(n+1)}{n} \sum_k \hat{\alpha}_{j_0k} \alpha_{j_0k} + \frac{1}{n^2} \sum_{i=1}^n \sum_k \{ \varphi_{j_0k}^2(X_i) - \int \varphi_{j_0k}^2 f \} + \\
& + \frac{1}{n^2} \sum_{i=1}^n \sum_{j=j_0}^{j_1} \sum_k \frac{\psi_{jk}^2(X_i) - \int \psi_{jk}^2 f}{(1+v2^{2js})^2}.
\end{aligned}$$

Denoting  $A_1 = \frac{1}{n^2} \sum_{i=1}^n \sum_k \{ \varphi_{j_0k}^2(X_i) - \int \varphi_{j_0k}^2 f \}$ ,  $A_2 = \sum_{j=j_0}^{j_1} A_{2j}$ ,

$$A_{2j} = \frac{1}{n^2} \sum_{i=1}^n \sum_k \frac{\psi_{jk}^2(X_i) - \int \psi_{jk}^2 f}{(1+v2^{2js})^2},$$

we see that  $E(A_i) = 0$ ,  $i = 1, 2$  holds. With the above notation introduced, we can write for the RHS of (A10):

$$\begin{aligned} A_1 + A_2 - \frac{n+1}{n^2(n-1)} \sum_{i=1}^n \sum_{l=1, l \neq i}^n \sum_k \{ \varphi_{j_0 k}(X_i) - \int \varphi_{j_0 k} f \} \{ \varphi_{j_0 k}(X_l) - \int \varphi_{j_0 k} f \} + \\ + \sum_{j=j_0}^{j_1} \left[ -\frac{2}{n(n-1)(1+v2^{2js})} + \frac{1}{n^2(1+v2^{2js})^2} \right] \times \\ \times \sum_k \sum_{i=1}^n \sum_{l=1, l \neq i}^n (\psi_{jk}(X_i) - \int \psi_{jk} f)(\psi_{jk}(X_l) - \int \psi_{jk} f) \\ + \sum_{j=j_0}^{j_1} \left\{ -\frac{2}{n(n-1)(1+v2^{2js})} + \frac{1}{n^2(1+v2^{2js})^2} \right\} 2(n-1) \sum_k \int \psi_{jk} f \cdot \sum_{i=1}^n \psi_{jk}(X_i) \end{aligned}$$

Now we invoke the decorrelation property of the wavelet transform with respect to the inner product  $\langle g, h \rangle_f := \int g(x)h(x)f(x)dx$  (cf. the proof of Theorem 1). This property, together with the bound on  $N_j$ , implies that, in particular,  $\sum_k \sum_{k'} (\int \varphi_{j_0 k} \varphi_{j_0 k'} f)^2 \leq C2^{j_0}$  holds, rather than  $\leq C2^{2j_0}$ . Here, again,  $C = C(\varphi)$ . This observation is implicitly used in [110], p.49. Applying the above observation, let us evaluate the order of magnitude of

$$B_1 = -\frac{n+1}{n^2(n-1)} \sum_{i=1}^n \sum_{l=1, l \neq i}^n h(X_i, X_l)$$

with

$$h(x, y) = \sum_k \{ \varphi_{j_0 k}(x) - \int \varphi_{j_0 k} f \} \{ \varphi_{j_0 k}(y) - \int \varphi_{j_0 k} f \}.$$

Note that

$$(A11) \quad Eh(X_i, X_l) = E(h(X_i, X_l)|X_i) = E(h(X_i, X_l)|X_l) = 0$$

holds. Invoking (A11), the decorrelation property and the Cauchy-Schwartz inequality, we obtain like in [110], p.49 that  $Var(B_1) \leq C \frac{2^{j_0}}{n^2}$  and, hence,

$P\{B_1 MISE(v)^{-1} > \epsilon\} \leq \epsilon^{-2} C 2^{-j_0} \rightarrow 0$ . Let now  $B_2 = \sum_{j=j_0}^{j_1} B_{2j}$ ,

$$B_{2j} = \left[ -\frac{2}{n(n-1)(1+v2^{2js})} + \frac{1}{n^2(1+v2^{2js})^2} \right] \times \\ \times \sum_k \sum_{i=1}^n \sum_{l=1, l \neq i}^n (\psi_{jk}(X_i) - \int \psi_{jk} f)(\psi_{jk}(X_l) - \int \psi_{jk} f).$$

Applying the quasi-triangle inequality to the sum in  $j$  and continuing by the same type of argument as for  $B_1$ , one obtains

$$E(B_2^2) = Var B_2 \leq 2(j_1 - j_0 + 1) \sum_{j=j_0}^{j_1} E(B_{2j}^2) \leq C(f, \psi) j_1 \frac{2^{j_0}}{n^2} \left[ 1 + \sum_{j=j_0}^{j_1} \frac{2^{j-j_0}}{(1+v2^{2js})^2} \right],$$

and, on applying Chebyshev inequality,

$$P\{B_2 MISE(v)^{-1} > \epsilon\} \leq \epsilon^{-2} C j_1 2^{-j_0} = C(f, \psi, \epsilon) 2^{-j_0} o(\log_2 n) \rightarrow 0,$$

by the choice of  $j_0$ .

Next, we show that under the theorem's assumptions

$$Var\left\{ \sum_{j=j_0}^{j_1} \left[ -\frac{2}{n(n-1)(1+v2^{2js})} + \frac{1}{n^2(1+v2^{2js})^2} \right] 2(n-1) \times \right. \\ \left. \times \sum_k \int \psi_{jk} f \cdot \sum_{i=1}^n \psi_{jk}(X_i) \right\} = Var(D_1) \leq \frac{C 2^{-2j_0 s'}}{n} \cdot \gamma_{j_0}$$

where  $\sum_{j=j_0}^{\infty} \gamma_j^2 < \infty$ . To prove this, denote  $c_{j,n}(v) = \left[ -\frac{2}{n(n-1)(1+v2^{2js})} + \frac{1}{n^2(1+v2^{2js})^2} \right] 2(n-1)$ . It is easy to see that  $|c_{j,n}(v)| < 4/n$ , uniformly in  $v, s$  and  $j$ .

After computations,

$$Var(D_1) = n \sum_{j=j_0}^{j_1} \sum_{j'=j_0}^{j_1} c_{j,n}(v) c_{j',n}(v) \sum_k \sum_{k'} \beta_{jk} \beta_{j'k'} \int \psi_{jk} \psi_{j'k'} f -$$

$$-n \left( \sum_{j=j_0}^{j_1} c_{j,n}(v) \sum_k \beta_{jk}^2 \right)^2 \leq n \sum_{j=j_0}^{j_1} \sum_{j'=j_0}^{j_1} c_{j,n}(v) c_{j',n}(v) \sum_k \sum_{k'} \beta_{jk} \beta_{j'k'} \int \psi_{jk} \psi_{j'k'} f.$$

Now we apply consecutively a change of the order of summation and integration, Hölder's inequality taking into consideration that  $f \in L_\infty$ , Parseval's identity and the upper bound for  $|c_{j,n}(v)|$  to obtain

$$\begin{aligned} \text{Var}(D_1) &\leq n \int \left( \sum_{j=j_0}^{j_1} c_{j,n}(v) \sum_k \beta_{jk} \psi_{jk} \right)^2 f \leq n \|f\|_{L_\infty} \sum_{j=j_0}^{j_1} \sum_k c_{j,n}(v)^2 \beta_{jk}^2 \leq \\ &\leq \frac{16}{n} \|f\|_{L_\infty} \sum_{j=j_0}^{j_1} \sum_k \beta_{jk}^2 = \frac{16}{n} \|f\|_{L_\infty} \left\| \sum_{j=j_0}^{j_1} \sum_k \beta_{jk} \psi_{jk}(\cdot) \right\|_{L_2}^2. \end{aligned}$$

By using a theorem due to Kerkycharian and Picard (compare the citation in Theorem 1 in [110]) we can now claim that  $\text{Var}(D_1) \leq \frac{C 2^{-2j_0 s}}{n} \gamma_{j_0}$ , where  $C = C_s \|f\|_{L_\infty}$ ,  $C_s$  depending on  $s$  only, and where  $\{\gamma_{j_0}\}_{j_0=0}^\infty \in l_2$  and, therefore,  $\gamma_{j_0} \rightarrow 0$  as  $j_0 \rightarrow \infty$ , hence,  $P(D_1 \text{MISE}(v)^{-1} > \epsilon) \leq \epsilon^{-2} C 2^{-2j_0(1+s')} n \gamma_{j_0}$ . Having in mind that  $2^{-j_0} = O(n^{-\frac{1}{2+2s'}})$ , we see that  $\frac{D_1}{\text{MISE}(v)}$  also tends to zero in probability.

It remains to evaluate  $A_1$  and  $A_2$ .  $A_1$  is a sum of i.i.d. random variables and their variances can be bounded from above. Using again Meyer's lemma and applying the decorrelation argument (see also [110], p.50) yields

$$\text{Var}(A_1) = E(A_1^2) \leq \frac{C 2^{2j_0}}{n^3}.$$

The argument about  $A_2$  is analogous to the one with  $B_2$ , with additional invoking of the inequality between the  $l_2$  and  $l_1$ -norm:

$$\begin{aligned} \text{Var}(A_2) = E(A_2^2) &\leq 2(j_1 - j_0 + 1) \left[ \sum_{j=j_0}^{j_1} E(A_{2j}^2) \right] \leq \frac{C}{n^3} \cdot j_1 \cdot 2^{2j_0} \cdot \sum_{j=j_0}^{j_1} \frac{2^{2(j-j_0)}}{(1 + v 2^{2js})^4} \leq \\ &\leq \frac{C}{n^3} \cdot j_1 \cdot 2^{2j_0} \cdot \left[ \sum_{j=j_0}^{j_1} \frac{2^{j-j_0}}{(1 + v 2^{2js})^2} \right]^2 \end{aligned}$$

for a suitably chosen  $C = C(f, \psi) > 0$ . Hence, by virtue of (A9), using the Chebyshev inequality, we can claim that

$$P\{A_1 MISE(v)^{-1} > \epsilon\} \leq \epsilon^{-2} \frac{C}{n}, P\{A_2 MISE(v)^{-1} > \epsilon\} \leq \epsilon^{-2} \frac{C}{n} o(\log_2 n).$$

Finally, let us note that  $E(T_n) = \int f(x)^2 dx$  holds. Neglecting negative summands, interchanging the order of integration and summation and applying Hölder's inequality and Parseval's identity, we bound the variance of  $T_n$  from above:

$$\begin{aligned} Var(T_n) &= \frac{4(n+1)^2}{n^3} \left[ \sum_k \sum_{k'} \left( \int \varphi_{j_0 k} \varphi_{j_0 k'} f \right) \cdot \alpha_{j_0 k} \alpha_{j_0 k'} - \left( \sum_k \alpha_{j_0 k}^2 \right)^2 \right] \leq \\ &\leq \frac{4(n+1)^2}{n^3} \int \left( \sum_k \alpha_{j_0 k} \varphi_{j_0 k} \right)^2 f \leq \frac{4(n+1)^2}{n^3} \left( \sum_k \alpha_{j_0 k}^2 \right) \|f\|_{L_\infty} \leq \\ &\leq \frac{4(n+1)^2}{n^3} \cdot \|f\|_{L_2}^2 \|f\|_{L_\infty}. \end{aligned}$$

Thus, we arrive at  $Var(T_n) \leq \frac{C}{n}$  and  $T_n \xrightarrow{P} \int f(x)^2 dx$  as  $n \rightarrow \infty$ , by Chebyshev inequality. ■

**Proof of Theorem 7.** The idea of the proof of the necessity part is essentially the same as the necessity part of Theorem 3. The admissibility of the choice  $v = \tilde{v}_n$  in the necessity part of the theorem follows from the boundedness of  $\tilde{v}_n$  proved in Theorem 5. We outline the proof of the sufficiency part. Proceeding first analogously to the proof of the sufficiency part of Theorem 5, and then applying the Cauchy-Schwartz inequality and interchanging the order of summation and integration, we obtain

$$\begin{aligned} E \| \tilde{f}_v - f \|_{B_{22}^\sigma}^2 &= E \| \tilde{f}_v - E \tilde{f}_v \|_{B_{22}^\sigma}^2 + \| E \tilde{f}_v - f \|_{B_{22}^\sigma}^2 = \\ &= \frac{1}{n} \sum_k \left[ \left( \int \varphi_{j_0 k}^2 f \right) - \alpha_{j_0 k}^2 \right] + \frac{1}{n} \sum_{j=j_0}^{j_1} \frac{2^{2j\sigma}}{(1 + v 2^{2js})^2} \sum_k \left[ \left( \int \psi_{jk}^2 f \right) - \beta_{jk}^2 \right] + \end{aligned}$$

$$\begin{aligned}
 & + \sum_{j=j_0}^{j_1} 2^{2j\sigma} \frac{v^2 2^{4js}}{(1 + v^2 2^{2js})^2} \sum_k \beta_{jk}^2 + \sum_{j=j_1+1}^{\infty} 2^{2j\sigma} \sum_k \beta_{jk}^2 \leq \\
 & \leq \frac{1}{n} \int [(\sum_k \varphi_{j_0 k}^2) f] + \frac{1}{n} \sum_{j=j_0}^{j_1} 2^{2j\sigma} \frac{v^2 2^{4js}}{(1 + v^2 2^{2js})^2} \int [(\sum_k \psi_{jk}^2) f] + \\
 & + \sum_{j=j_0}^{j_1} 2^{2j\sigma} \frac{v^2 2^{4js}}{(1 + v^2 2^{2js})^2} \sum_k \beta_{jk}^2 + \sum_{j=j_1+1}^{\infty} 2^{2j\sigma} \sum_k \beta_{jk}^2.
 \end{aligned}$$

Bounding each of the last two terms from above in the same way as in the proof of Theorem 3, and applying Meyer's lemma to each of the first two terms, we obtain

$$\begin{aligned}
 & E \|\tilde{f}_v - f\|_{B_{22}^\sigma}^2 \leq \\
 & \leq c(\varphi) \frac{2^{j_0}}{n} + \frac{c(\psi)}{n} \sum_{j=j_0}^{j_1} \frac{2^{j(1+2\sigma)}}{(1 + v^2 2^{2js})^2} + (v^{\frac{s'-\sigma}{s}} + c_{s',\sigma} 2^{-2j_1(s'-\sigma)}) \|f\|_{B_{22}^{s'}}^2 \leq \\
 & \leq c_{s',\sigma}(f, \varphi, \psi) \left( \frac{2^{j_0}}{n} + \frac{2^{j_1(1+2\sigma)}}{n} + v^{\frac{s'-\sigma}{s}} + 2^{-2j_1(s'-\sigma)} \right) = o(1),
 \end{aligned}$$

when  $2^{j_1} = o(n^{\frac{1}{1+2\sigma}})$ ,  $j_0 \leq j_1$ ,  $j_1 \rightarrow \infty$  and  $v \rightarrow 0$ .

The admissibility of  $\tilde{v}_n$  in the sufficiency part follows from Theorem 8. ■

**Proof of Theorem 8. (Outline.)** We follow closely the idea of proof of Theorem 4. In this case, the resulting equation for the upper bound  $v_0$  is

$$(A12) \quad v = \frac{\|f\|_{L^\infty}}{c_0^2} \cdot \frac{1}{n} \cdot \sum_{j=j_0}^{j_1} 2^{2(j-j_0)s} \left( \frac{1 + v^2 2^{j_0 s}}{1 + v^2 2^{js}} \right)^3 \cdot 2^j,$$

where  $c_0$  is the sharpness constant. The proof can be completed by invoking the bound  $v < C$  again and proceeding as in the proof of Theorem 4. However, in this case, due to the unbiased estimation of the wavelet coefficients when the design is random, we can directly obtain the quantitative rates of  $\tilde{v} \rightarrow 0$ , for any  $s' > 0$ . In fact, for any  $j_0, j_1 : j_0 \leq j_1$  we can obtain the sharp rate in  $\tilde{v} = o(1)$

corresponding to these  $j_0$  and  $j_1$  in the following way: instead of utilizing the rough bound  $v \leq C$ , we replace the RHS of (A12) by

$$\frac{\|f\|_{L_\infty}}{c_0^2} \cdot \frac{1}{n} \cdot \sum_{j=j_0}^{j_1} 2^{2(j-j_0)s} \left( \frac{\max(1, v2^{j_0s})}{\max(1, v2^{j_1s})} \right)^3 \cdot 2^j.$$

The new RHS of (A12) is equivalent to the old one, with absolute equivalence constants. From the new expression, the sharp rate can be computed by considering the three cases: (i)  $v > 2^{-2j_0s}$ ; (ii)  $v \leq 2^{-2j_1s}$ ; (iii)  $2^{-2(j'+1)s} \leq v < 2^{-2j's}$ ,  $j' \in \mathbb{Z} : j_0 \leq j' < j_1$ . By solving (A12) for  $v$  for each of the cases (i-iii) and checking whether the solution  $v_0$  falls into the range of the respective case, we are able to eliminate two of the cases, and accept the remaining third case. If case (iii) is the accepted one, the substitution  $v = 2^{-j's}$ , where  $j' = (1 - \theta)j_0 + \theta j_1$ , allows locating  $j'$  within  $[j_0, j_1]$  by solving the resulting equation for  $\theta \in [0, 1]$ . (See also B14.) ■

**Proof of Corollary 2.** Theorem 8 and the sufficiency part of Theorem 7 are applied for  $v = \tilde{v}$  and  $\sigma = 0$ , and then the proof is completed via Theorem 6. ■

**Proof of Corollary 3.** The rate  $v_0 = O(\frac{2^{j_1}}{n})$  follows immediately from (A12) with  $j_0 = j_1$  and, therefore,  $cn^{\frac{1}{1+2s}} \leq 2^{j_1} \leq Cn^{\frac{1}{1+2s}}$  and  $\tilde{v} \leq v_0$  together imply  $\tilde{v} = O(n^{-\frac{2s}{1+2s}})$ . Now it follows from the proof of Theorem 7 that  $E \|\tilde{f}_{\tilde{v}} - f\|_{B_{22}^s}^2 = O(n^{-\frac{2(s-\sigma)}{1+2s}})$ . ■

**A19.** (See also [49].) For the purpose of optimizing the GFCV- (FCV-) functional, the excellent subroutine `amoeba` from Numerical Recipes [100] has been utilized. It is based on the downhill simplex method of Nelder and Mead. Since the minimization is performed only once, without iterations, so the computational burden is small, the above procedure has been selected because of its simplicity and robustness. It is also derivative-free and we intend to use this for some of the extensions in Appendix B.

**A20.** There follow some model examples where the universal-threshold and global-penalization shrinking strategies are *not* expected to yield equally satisfactory estimation of every part of the curve: Marron and Wand's "smooth comb" (see [96]), Examples 3 and 4, and, on fractal level, any *multifractal* function (see [11]).

**A21.** Notice the bad performance of our non-threshold estimator near  $x = 1$ , where the true signal is negligible compared to the noise level and the thresholded estimator is at its best. One optional modification of the shrinking estimator in the case of regression is to threshold it (after shrinking) with hard threshold equal to the noise variance of a purely noisy coefficient. This threshold level (typically,  $\frac{2\delta}{\sqrt{n}}$ ), is very low compared to levels used in threshold methods, and is aimed at removing pure noise for moderate samples. This is a most primitive example of a composite shrinking/thresholding estimator which will be discussed in more detail in  $B^8$ .

## Appendix B: Further extensions and generalizations

In Section 7 we discussed some extensions of the  $K$ -functional technique of penalized wavelet estimation. Here we give a (non-exhaustive) list of further developments and applications, together with short descriptions of some underlying ideas. Again, we assume that  $a_1 = 0$ ,  $a_2 = 1$ .

**B1. Cross-validation in Besov and potential spaces.** There is a common opinion expressed in literature that cross-validation in  $L_2$  may in some cases lead to a slightly underestimated value of the smoothing parameter leading to a slight tendency to overfit the curve. Analytically, this can be explained with the fact that in general  $L_2$ -cross validation leads to consistent estimation of  $f$  itself, but not of its fractional derivatives of any positive order. This can be overcome by considering cross-validation with respect to the  $B_{22}^\sigma$ -norm such that available results for  $L_2$ -cross validation correspond to the partial case  $\sigma = 0$ . The key fact for this extension is to notice that  $L_2$ -cross validation, as performed in the density case (see subsection 6.2) is of Bowman–Rudemo type and is equivalent to applying the cross-validation rule  $\tilde{\alpha}_{j_0 k} \alpha_{j_0 k} \mapsto \frac{1}{n} \sum_{i=1}^n \tilde{\alpha}_{j_0 k(-i)} \varphi_{j_0 k}(X_i)$ ,  $\tilde{\beta}_{jk} \beta_{jk} \mapsto$



$\frac{1}{n} \sum_{i=1}^n \tilde{\beta}_{jk(-i)} \psi_{jk}(X_i)$ , for every  $(j, k) : j_0 \leq j \leq j_1$ ,  $\text{supp} \varphi_{j_0 k} \cap \text{supp} f \neq \emptyset$ ,  $\text{supp} \psi_{jk} \cap \text{supp} f \neq \emptyset$ , where the notations are the same as in Subsection 6.2 (recall that the  $\tilde{\beta}$ 's depend on  $v$ ). This leads to optimization in  $v$  of the CV-criterion

$$CV_\sigma(v) = CV_{\sigma,s,j_0}(v) = \sum_k \tilde{\alpha}_{j_0 k}(v)^2 + \sum_{j=j_0}^{j_1} 2^{2j\sigma} \sum_k \tilde{\beta}_{jk}(v)^2 - \\ - \frac{2}{n} \sum_{i=1}^n \left( \sum_k \tilde{\alpha}_{j_0 k}(-i) \varphi_{j_0 k}(X_i) + \sum_{j=j_0}^{j_1} 2^{2j\sigma} \sum_k \tilde{\beta}_{jk}(-i) \psi_{jk}(X_i) \right),$$

with  $CV_0(v) = M(\tilde{f}_v)$ . For  $\sigma > 0$   $CV_\sigma(v)$  is the cross-validated version of  $\| \tilde{f}_v - f \|_{B_{22}}^2$  up to summands independent of  $v$ . For the regression case we suggest the same Bowman–Rudemo approach, based on the formulae  $\tilde{\alpha}_{j_0 k} \alpha_{j_0 k} \mapsto \frac{1}{n} \sum_{i=1}^n \tilde{\alpha}_{j_0 k}(-i) y_i \varphi_{j_0 k}(x_i)$ ,  $\tilde{\beta}_{jk} \beta_{jk} \mapsto \frac{1}{n} \sum_{i=1}^n \tilde{\beta}_{jk}(-i) y_i \psi_{jk}(x_i)$ . This yields

$$\tilde{CV}_\sigma(v) = CV_{\sigma,s,j_0}(v) = \sum_k \tilde{\alpha}_{j_0 k}(v)^2 + \sum_{j=j_0}^{j_1} 2^{2j\sigma} \sum_k \tilde{\beta}_{jk}(v)^2 - \\ - \frac{2}{n} \sum_{i=1}^n \left( \sum_k \tilde{\alpha}_{j_0 k}(-i) y_i \varphi_{j_0 k}(x_i) + \sum_{j=j_0}^{j_1} 2^{2j\sigma} \sum_k \tilde{\beta}_{jk}(-i) y_i \psi_{jk}(x_i) \right),$$

where the  $\tilde{\alpha}$ 's and  $\tilde{\beta}$ 's have meaning as of Subsection 6.1. For  $\sigma = 0$  (see [49]) this is an important and, as it seems, quite new alternative to Wahba's model, with "pointwise" estimation (i.e., estimation of the linear interpolation functionals  $L_i(f) = f(x_i)$ ) being replaced by "weak" estimation (i.e., estimation of the linear interpolation functionals  $f \mapsto \alpha_{j_0 k}$  and  $f \mapsto \beta_{jk}$ ). Besides admitting an extension for  $\sigma > 0$  in a natural way, the Bowman–Rudemo type approach to the regression case has some additional advantages compared to the model in Subsection 6.1: (a) it is available also for the model where  $x_i = X_i$ ,  $i = 1, \dots, n$ , are i.i.d. random variables, uniformly distributed in a neighbourhood

of  $\text{supp} f$ ; (b) this model admits *analogous theory of iterative individual penalized shrinking*, as considered in Section 7 for the density case.

It should be noted that in the case of nonparametric regression theoretical derivations for the Bowman–Rudemo approach are easier to carry through if *periodized* (see [90,37]) wavelets  $\varphi_{j_0k}^{per}$  and  $\psi_{jk}^{per}$  are considered. (Computing the  $\alpha_{j_0k}$ ’s and the  $\beta_{jk}$ ’s in (2,3) with respect to  $\varphi_{j_0k}^{per}$  and  $\psi_{jk}^{per}$  can be shown to yield an equivalent quasi-norm of  $J_{spq}(\alpha, \beta)$ .)

As for the model in Section 6.1, it can be extended for  $\sigma > 0$  by utilization of the fact that  $B_{22}^\sigma$  is isomorphic to the potential (Sobolev) space  $H_2^\sigma$  (see [16], 6.4.4). The CV-functional is now the expectation of a quadrature formula for the integral in  $\|\tilde{f}_v - f\|_{H_2^\sigma}^2 = \int [G_\sigma * (\tilde{f}_v - f)(x)]^2 dx$ , where  $G_\sigma$  is the Bessel–McDonald kernel<sup>41</sup>. The order of the quadrature formula can be higher than one. (In comparison, the optimization criteria of GCV type like QGCV and GFCV are essentially quadrature formulae of first order.) The higher order of the quadrature formula is helpful in reducing the bias for more smooth functions, and plays essentially the same role as the order of vanishing moments of the coiflet used in the method WAVREG of [9]. This approach is more technically involved because it includes regularization of the singular kernel  $G_\sigma$  depending on the sample size  $n$ . We conjecture that the analogues of  $\bar{h}(v)^2$  and  $\mu(v)$  (see (14)) are in the new case the norms in *weighted* (both with the same weight) Schatten–von Neumann  $S_\gamma$ -norms,  $\gamma = 2, 1$ , respectively, the weight depending on  $\sigma$  and tending to 1 as  $\sigma \rightarrow 0+$ .

The results in Section 6 show that consistency of cross validation for  $\sigma > 0$  is sufficient for consistency of the resulting estimator  $\tilde{f}_v$  in  $B_{22}^\sigma$ , i.e., with better estimation of the derivatives and less overfitting of the curve. One way to determine a best selection of  $\sigma$  for the concrete  $f$ , if simultaneous optimization in  $v$  and  $s$  is performed, is to increase  $\sigma$  (e.g., by using the dyadic-bisection trial-and-error method) until for the optimal value  $s^* = s^*(\sigma)$  the equality  $s^*(\sigma) \approx \sigma$  holds with sufficient precision. For  $f \notin C^\infty$  this is possible to achieve because  $s^*(\sigma)$  remains in a neighbourhood of the true smoothness index  $s'$  of the curve

for any  $\sigma > 0$ . Another approach here, for both the regression and density case, is to consider a *relative cross-validation criterion*:

$$CW_{\sigma}(v) = CW_{\sigma, s, j_0}(v) = (b_{22}^{\sigma}(\hat{f}))^{-1} CV_{\sigma}(v),$$

where  $b_{22}^{\sigma}(\hat{f}) = b_{22}^{\sigma}(\hat{f})_{j_0, j_1} = \|\hat{\alpha}_{j_0}\|_2^2 + \sum_{j=j_0}^{j_1} (2^{j\sigma} \|\hat{\beta}_j\|_2)^2$  (cf. (2)). The optimal value of  $\sigma$  is obtained by  $CW_{\sigma^*} = \min_{0 \leq \sigma \leq s^*} CW_{\sigma, s^*, j_0^*}(v^{opt})$ , where  $s^*$ ,  $j_0^*$  and  $v^{opt}$  are obtained as discussed in Remarks 5 and 7, for the regression and density case, respectively. An interesting modification is to minimize  $CW_{\sigma, s, j_0}(v)$  simultaneously in  $\sigma, s, j_0$  and  $v$ .

In the case of regression with deterministic design, the method of [9] can also be extended, similarly to the Bowman–Rudemo approach, so as to determine the optimal value of the smoothing parameter on the basis of an estimate of the upper bound for  $\|\tilde{f}_v - f\|_{B_{22}^{\sigma}}$  for  $0 \leq \sigma \leq s$ ,  $j_1 : 2^{j_1} = o(n)$ , which bound is the RIIS of

$$\|\tilde{f}_v - f\|_{B_{22}^{\sigma}} \leq \|(I - P_{j_1})f\|_{B_{22}^{\sigma}}^2 + 2\|(P_{j_1}f - \Pi_{j_1})f\|_{B_{22}^{\sigma}}^2 + 2\|(\Pi_{j_1}f - \tilde{f}_v)\|_{B_{22}^{\sigma}}^2 = RII S,$$

where  $P_{j_1}$  is the orthogonal projection of  $f$  onto  $V_{j_1}$ , and  $\Pi_{j_1}$  is defined in [9], p.318. An unbiased estimate  $CV_{\sigma, s, j_0}(v)$  of  $\|\Pi_{j_1}f - \tilde{f}_v\|_{B_{22}^{\sigma}}^2$  can be obtained for  $\sigma > 0$  in the way it has been proposed in [9] for the case  $\sigma = 0$ . The use of coiflets makes this approach quite promising, especially when the estimated function  $f \in B_{22}^{s'}$ ,  $1/2 < \sigma + 1/2 < s' \leq s < r$ , is very smooth (i.e.,  $s' \gg 1$ ) and  $\sigma \approx s' - 1/2$ . However, the support of the coiflet is larger than that of the respective Daubechies' minimal-support wavelet (see [36]), which reduces the adaptivity of the method for spatially inhomogeneous curves, especially, when the true smoothness index  $s'$  does not exceed 1. In order to improve the spatial adaptivity of the method, we suggest to replace the coiflets with 'multi-coiflets', that is, multi-scaling functions  $B^{18(c)}$  which are orthogonal to the same monomials as the coiflets and have approximately the same regularity, but are better spatially localized. The so proposed generalization of WAVREG of [9] (with the additional improvement that  $j_1$  is not fixed, but controlled, together with  $j_0$  and with a

coiflet or 'multi-coiflet') will be referred to as  $WAVREG_\sigma$  later in this text. As far as determining an optimal value of  $\sigma^*$  (based on  $WAVREG_\sigma$ ) is concerned, it does not seem possible to develop a relative optimization criterion of the type of  $CW_\sigma(v)$  defined for the Bowman-Rudemo approach. However, the first method considered above (i.e., finding a fixed point of the mapping  $\sigma \mapsto s^*(\sigma)$ ) is still available.

Finally, let us note that in the models of cross validation in  $B_{22}^\sigma$  with subsequent optimal choice of  $\sigma$  it is a good idea to optimize the criterion not only in  $j_0$ , but also in  $j_1$ . (For the usual  $L_2$ -cross validation optimization in  $j_1$  is not advisable, because the optimal  $j_1$  would tend to be too large most of the time, thus leading to overfitting the curve.)

**B2. Asymptotic-minimax theory in Besov spaces.**<sup>B9</sup> As with thresholded estimators, it can be expected that in the theoretical study of the asymptotic-minimax optimal choice of the smoothing parameter exhaustive quantitative results can be obtained. We expect that such exhaustive results can be obtained for the level-dependent shrinking estimator  $\tilde{f}_v$ ,  $v = (v_{j_0}, \dots, v_{j_1})'$ , for any  $(\pi, u, \sigma)$  and  $(p, q, s)$  such that  $B_{\pi u}^\sigma \hookrightarrow B_{pq}^s$ . Although not identical, the proof in the shrinking case can be based on the same basic ideas as the one for the threshold case. We do not expect major technical challenges in obtaining asymptotic-minimax results analogous to the results of [53], i.e., including optimal rates of estimation of the derivatives of  $f$ , too. However, our aim will be to remove some of the logarithmic factors in the rates and to obtain *more precise bounds for the constant factors of the rates*, depending on the concrete choice of  $(\pi, u, \sigma)$  and  $(p, q, s)$ .

**B3. Kernel regularization via quasi-linearization of  $K$ -functionals.** Following Peetre [94], we call the pair of quasi-normed spaces  $A, B$  quasi-linearizable if there exists a family of (sub)linear (see [16], 3.10, 3.11, and [39], Section 5) operators  $V(t) : A + B \rightarrow A + B$  such that  $V(t)(A + B) \subset B$ ,  $(I_{A+B} - V(t))(A + B) \subset A$  and

$$c[\| (I_{A+B} - V(t))f \|_A^p + t^p \| V(t)f \|_B^p]^{1/p} \leq K_\rho(t, f; A, B) \leq$$

$$\leq C[\| (I_{A+B} - V(t))f \|_A^\rho + t^\rho \| V(t)f \|_B^\rho]^{1/\rho}$$

for any  $f \in A + B$  and any  $t \in (0, \infty)$ . Here  $0 < c \leq C < \infty$  and  $c, C$  may depend on  $A, B$  and  $\rho : 0 < \rho \leq \infty$ , but are independent of  $f$  and  $t$ .

Let  $K(x)$  be the kernel associated with any  $r$ -th order kernel estimator, that is,  $\int K(x)dx = 1$  and  $\int K(x)x^l dx = 0, l = 1, \dots, [r]$  (with obvious modifications in the multivariate case, when  $l$  becomes a multiindex). Denote  $K_t = K(\frac{\cdot}{t})$ . Then, utilizing the real and complex interpolation methods (see Section 4), it can be shown that  $(B_{pu}^\sigma, B_{pq}^s)$  and  $(B_{pu}^\sigma, \dot{B}_{pq}^s)$  are quasi-linearizable pairs with  $V(t) : f \mapsto K_{t^{\frac{1}{s-\sigma}}} * f$ , uniformly in  $0 < p \leq \infty, 0 < u \leq \infty, 0 < q \leq \infty, \max\{0, d(\frac{1}{p} - 1)\} < \sigma < s < r$ , where  $d$  is the dimension. Here  $c$  and  $C$  depend on the concrete choice of  $K, p, q, u, \sigma, s$  and  $\rho$ . The same is true if Besov spaces are replaced by potential (Sobolev) spaces for the same values of the parameters  $\sigma, s$  and  $p : 1 \leq p \leq \infty$ . A similar result holds for wavelet estimators with  $\psi$  having vanishing first  $[r]$  moments. In fact, for the case  $p = u = q = 2$  and  $\rho = 2$  the  $\beta$ -coefficients of  $V(t)f$  are computed from the coefficients of  $f$  via (7) (for other values, see  $B^9$ ). These general quasi-linearization results show that the role of the parameter  $v = t^\rho$  is that of a smoothing parameter in a penalization model induced by the  $K$ -functional and, simultaneously,  $v^{\frac{1}{d(s-\sigma)}}$  has the role of a kernel estimator's bandwidth. From this point of view, the difference between kernel and wavelet estimators is mainly in the choice of the type of equivalent quasi-norms in the Besov spaces involved in the  $K$ -functional. In fact, rigorous results can be obtained, linking the optimal smoothing parameter  $\tilde{v}$  for the wavelet model with the bandwidth  $v^*$  of a kernel estimator. The relation is very simple:  $c_1 \tilde{v} \leq v^* \leq C_1 \tilde{v}$ , where  $c_1$  and  $C_1$  depend only on the constants of equivalence between the  $K$ -functional, involving the coefficient Besov quasi-norms, and its quasi-linearized version involving the kernel  $K$ . Taking into account that under very general assumptions  $\tilde{v} = \tilde{v}_n \rightarrow 0$  as the sample size  $n \rightarrow \infty$ , we see that a "good bandwidth"  $v^*$  can be found by only searching within  $[c_1 \tilde{v}_n, C_1 \tilde{v}_n]$ , which interval becomes narrower with the increase of  $n$ . Certainly this technique can be applied also for variable-kernel estimators,

provided that one starts from the block-shrinking wavelet estimator (see also Examples 3 and 4). Of course, rigorous formulations depend on the specific statistical model (density estimation, regression-function estimation (Wahba's model or the Bowman-Rudemo approach), etc.) but in all cases the underlying general idea is as outlined above. In particular, one or both spaces in the K-functional may depend on its step  $t$ :  $K(t, f) = K(t, f; A_t, B_t)$  (see, e.g. <sup>B10</sup> and [39]). In the most practically important statistical applications one faces the situation  $K(t, f) = K(t_n, \hat{f}_n; A_n, B_n)$ , where  $n$  and  $\hat{f}_n$  are, as usual, the sample size and the empirical estimator, respectively. In these more general situations  $t$  may still be equivalent to the bandwidth of a kernel estimator, but the simple rescaling  $t \mapsto t^{\frac{1}{s-\sigma}}$  considered above may be replaced now by the more general  $t \mapsto w(t)$ , where  $w(\cdot)$  is a warping function which depends on the type of dependence of  $A$  and/or  $B$  on the parameter. For some cases of dependence of  $A$  and  $B$  on the parameter, it may happen that  $w(\cdot)$  is not right- or left-invertible, which may be an indication that the dependence of the spaces  $A$  and/or  $B$  on the parameter takes control over the smoothing process, the penalization parameter  $t$  playing then the less significant role. The typical case when  $A$  and  $B$  depend on  $n$  is  $j_0 \rightarrow \infty$ ,  $j_1 = o(\log_2 n)$ .

**B4. Kernel estimation via solutions of evolutionary differential equations.** Here we outline a new, very general approach to kernel estimation in statistics. This approach falls into the general quasi-linearization scheme considered in <sup>B3</sup>, but the method of generating the kernel estimator is very different and seems to be complementary to traditional kernel-estimation techniques.

We return to the considerations in Section 4 leading to formula (5). In the notation there, assume that  $B_1 = A$  where  $A$  is a Banach space. To conform with commonly accepted notation, we denote by  $B = D(T)$  the domain of  $T$  (see [16], Section 6.7), where  $T$  has the properties assumed in Section 4. Consider the following Cauchy problem for the evolutionary differential equation in the Banach space  $A$

$$(B1) \quad \frac{d}{dt}g(t) = Tg(t), \quad g(0) = a \in A, \quad t > 0.$$

To ensure smoothing effect,  $\|g(t_2)\|_A \leq \|g(t_1)\|_A$ ,  $t_2 > t_1$ , should be fulfilled. A necessary assumption to ensure consistency is  $\lim_{t \rightarrow 0+} \|g(t) - a\|_A = 0$ . Both these assumptions are fulfilled if  $T$  is such that  $g(t) = G(t)a$ , where  $G(t)$  is an *equibounded, strongly continuous operator semigroup* (ESCOSG) on  $A$  (see [16], Section 6.7, (i-iii)), which, additionally, is *contractive*, that is, in the notations of [16],  $M \leq 1$ . In this case  $T$  is the infinitesimal generator (IG) of the semi-group  $\{G(t), t > 0\}$ . For necessary and sufficient conditions on  $T$  under which  $T$  is the IG of a contractive ESCOSG on  $A$ , we refer to [61,123,86,73,75]. A simple necessary (but not sufficient) condition is that the spectrum of  $T$  be contained in  $\{z \in \mathbb{C} : \operatorname{Re} z \leq 0\}$ . The key result is (see [94,24]): if  $\{G(t), t > 0\}$  is an ESCOSG on  $A$  with IG  $T$  and  $\nu \in N$  is fixed, then for any  $\rho \in (0, \infty]$  there exist  $c_\nu = c_\nu(\rho)$ ,  $C_\nu = C_\nu(\rho) : 0 < c_\nu \leq C_\nu < \infty$ , such that

(B2)

$$c_\nu K_\rho(t, a; A, D(T^\nu)) \leq \omega_\nu(a; t)_A + \min(1, t) \|a\|_A \leq C_\nu K_\rho(t, a; A, D(T^\nu)),$$

for any  $t > 0$  and any  $a \in A$ , where  $\omega_\nu(a; t)_A = \sup_{0 < s \leq t} \|(G(s) - I_A)^\nu a\|_A$ . We see from (B2) that the parameter  $t$  plays a twofold role again: on the one side, it is the smoothing parameter in a penalization model; on the other, it denotes the length of the "time interval" in which  $a$  has been subject to smoothing via the process described by (B1).

Let us consider one simple example. Assume that  $A \in L_2$ ,  $T_1 = k^2 \frac{d^2}{dx^2}$ ,  $k \in \mathbb{R}$ ,  $x \in \mathbb{R}$ ,  $\nu = 1$ ,  $\sigma = 0$ ,  $s = 2 = \rho$ ,  $a = f \in L_2(\mathbb{R})$ ,  $T$  is the closed self-adjoint densely defined extension of  $T_1$  on  $L_2(\mathbb{R})$ . In other words, (B1) is the Cauchy problem for the unidimensional heat and diffusion equation in homogeneous media with conductivity coefficient  $k^2$ .  $T$  is the IG of a contractive ESCOSG on  $L_2(\mathbb{R})$ . Then  $T$  has only continuous spectrum, and it is contained in  $\{z : \operatorname{Re} z \leq 0\}$  (if  $\operatorname{supp} f \subset [0, l]$ ,  $l < \infty$ , then we can consider  $L_2([0, l])$  instead of  $L_2(\mathbb{R})$  and the spectrum of  $T$  is discrete; it consists of the simple eigenvalues  $\lambda_\mu = -k^2 \pi^2 \mu^2 / l^2$ ). Then, by (5),  $K_2(t, a; A, D(T^\nu))$  is equivalent to  $K_2(\sqrt{v}, f; B_{22}^0, \dot{B}_{22}^s) + \min(1, \sqrt{v}) \|f\|_{B_{22}^0}$ , with equivalence constants in-

dependent of  $t = \sqrt{v}$  and  $a = f$ . Besides,  $[G(t)f](x) = K_{t^{1/2}} * f(x)$ ,  $x \in R$ , where  $K(\xi) = \frac{|k|}{2\sqrt{\pi t}} e^{-\frac{\xi^2}{4k^2}}$ , and we observe that this is indeed a partial case of the kernel-estimation model considered in B3. This last observation is in fact correct under quite general assumptions on  $T$  in (B1) (see, e.g., the example considered in [16], Section 6.7, pp. 157,158, and Theorem 6.7.4., as well as [108,109]).<sup>B21</sup>

This new approach to kernel estimation is of considerable practical importance. Indeed, in the last several decades an extensive theory and rich variety of software products has been developed for the solution of evolutionary differential equations via *finite-element* and *boundary-element* methods (FEM and BEM). The new approach described here allows applying the full computational power of FEM and BEM to statistical kernel-estimation, which is of particular benefit in the multivariate case.<sup>B15,B21</sup>

**B5. Estimation with linear constraints.** The procedure of shrinking the wavelet coefficients described by (7) can be generalized in a natural way if the  $K$ -functional penalization model considered in Section 5 be upgraded to take into account additional information about constraints of: (a) *interpolatory* type:  $f^{(\nu)}(\xi_\mu) = c_{\mu\nu}$ ,  $\nu = 0, \dots, \nu_\mu$ ,  $\mu = 0, \dots, m$ , where  $\xi_{\mu_1} \neq \xi_{\mu_2}$  for  $\mu_1 \neq \mu_2$ ; (b) *isoperimetric* type: e.g.,  $\int f^{(\nu)}(x) g_\nu(x) dx = c_\nu$ ,  $\nu = 0, \dots, \nu_0$ , where  $g_\nu(x)$  are known functions. If the constraints are to be obeyed *strictly*, then *Lagrange-multiplier* technique can be applied. If only an *approximate* tendency to obey the constraints is being assumed (which is often the case in statistical applications), then *quadratic penalization* is recommended. In the "mixed" case there will be a Lagrange multiplier for each strict constraint and a quadratic penalization term (multiplied by a respective large parameter) for every non-strict constraint. Since all constraints, in both the interpolatory and isoperimetric case, are linear, so will be the system of equations for  $\tilde{\alpha}_{j_0k}$ ,  $\tilde{\beta}_{jk}$  and the eventual Lagrange multipliers. The unique solution of this system provides the analogue of (7) for the so-upgraded model. The value of  $v$  utilized in the computations can be chosen as the optimal value obtained by cross validation or by asymptotic-



minimax considerations for the unconstrained model (7).

**B6. Iterative application of the  $K$ -functional criterion.** One instance of such iterative smoothing was considered in Section 7, related to the individual shrinking model in density estimation. In fact, this approach is of interest for all penalization models discussed in Section 6 and 7. For example, consider the GFCV model in Subsection 6.1. Since only linearity of the estimator is needed for GFCV, after determining  $\tilde{f}_{\tilde{v}}$  from  $\hat{f}$ , on a second iteration, one may replace  $\hat{f}$  by  $\tilde{f}_{\tilde{v}}$  and, denoting  $\tilde{v}_1 = \tilde{v}$ , apply GFCV again to obtain  $(\tilde{f}_{\tilde{v}_1})_{\tilde{v}_2}$ , and so on. In general, the latter estimator will be smoother than  $\tilde{f}_{\tilde{v}_1}$ . For the case of a single global smoothing parameter  $v$ , it can be expected that the iterative process will be quickly convergent (i.e., the iteration number  $l$  increasing,  $\tilde{v}_l$  will tend quickly to zero.) For the level-dependent and block-shrinking models the iterative process is expected to converge more gradually, leading to a more visible improvement in the adaptivity of the final iterated estimator. For individual shrinking the iterative process may not converge, but an optimal number of iterations can be determined by cross validation (cf. Section 7).

**B7. Estimation with smoothness constraints.** Consider  $\tilde{f}_v$ , with wavelet coefficients defined via (7), where  $s$  is assumed to be fixed. Suppose that for a certain  $p, q, s' \leq s$  it is additionally known for the true function  $f$  that  $\|f\|_{B_{pq}^{s'}} = c_f$  holds, where  $c_f$  is a known positive constant (there is a modification for the homogeneous model). Then, it is possible to determine a unique non-negative value  $v_n$  of the smoothing parameter, such that  $\|\tilde{f}_{v_n}\|_{B_{pq}^{s'}} = c_f$  for any  $n \in N$ . (If  $\|\hat{f}\|_{B_{pq}^{s'}} \leq c_f$ , then  $v_n = 0$ .) Thus, additional information about such a type of smoothness constraint on  $f$  proves to be decisive for the choice of the smoothing parameter, and we no longer need additional (cross-validational, asymptotic-minimax optimal, etc.) statistical procedure for determining  $v$ . (For the level-dependent model in Section 7, under the additional assumption  $q = \infty$ , one may determine in a unique way the smallest  $v_j$  on every level (i.e., the least regular  $\tilde{f}_v$ ), so that  $\|\tilde{f}_v\|_{B_{pq}^{s'}} = c_f$ .) A lower (upper) bound on  $c_f$  leads to an

upper (lower) bound on the smoothing parameter, respectively. In <sup>B9(c)</sup> utilizing a smoothness constraint will be important. In all these cases, the statistical aspect of the model can be transferred to estimating  $c_f$  itself, and obtaining  $v_n$  from the estimate  $\hat{c}_f = \hat{c}_{f,n}$  is then achieved by the above deterministic procedure. In principle, this procedure of smoothing is available for very general type of noise (even correlated).

**B8.** *Composition and comparison with other consistent estimators.* <sup>B6</sup> provided an example of a composite strategy of smoothing and denoising involving several estimators. Here we discuss two other composite strategies.

(a) *Shrinking by a combined cross-validation/asymptotic-minimax approach.* For  $\tilde{f}_v$  defined by (1',7), let  $v_{1,n}^*$  be the optimal value obtained for sample size  $n$  via cross validation; let  $v_{2,n}^*$  be the respective value recommended by asymptotic-minimax theory. Consider  $v_n^* = (1 - \tau_n)v_{1,n}^* + \tau_n v_{2,n}^*$ , where  $\tau_n \in [0, 1]$ . For small to moderate samples, take  $\tau_n = 0$ ; for very large  $n$ , choose  $\tau_n = 1$ ; for intermediate  $n$  let  $\tau_n$  have intermediate values. The reason for this type of choice is that for very large  $n$  the asymptotic-minimax choice is both computationally cheaper and more reliable than the cross-validated one, most of the time. If, additionally, there is a biased perception of  $f$  being a smooth function, then the choice  $v^* = \max(v_1^*, v_2^*)$  is recommended; if, on the contrary,  $f$  is perceived to have an irregular graph (e.g., fractal), then  $v^* = \min(v_1^*, v_2^*)$  can be suggested (cf. also <sup>B11</sup>).

(b) *Combined shrinking/thresholding approach.* This is a modification of the iterative approach <sup>B6</sup> where on some of the iterations thresholding is performed rather than shrinking. The most important case is that of two iterations - one shrinking, one thresholding, or vice versa (in fact, these iterations commute modulo rescaling the shrinking and/or thresholding parameter). Because the strategy is composite, estimation of the smoothness parameter for general types of noise is an open problem. We propose two approaches to solve it. The first one (for arbitrary white or weakly correlated noise, and for density estimation) is by estimation with smoothness constraints <sup>B7</sup>: for  $\tilde{f}_v$ , as defined via

(7), we apply the CV-techniques of Section 6 to obtain an optimal value  $v_0^*$  and optimal  $s^*$ . Thus, it is assumed that  $f \in B_{22}^{s^*}$ , and  $\|f\|_{B_{22}^{s^*}}$  is estimated by the coefficient norm  $\|\tilde{f}_{v_0^*}\|_{B_{22}^{s^*}}$ . For a fixed  $v \in [0, v_0^*)$  the  $\beta$ -coefficients of  $\tilde{f}_v$  are being thresholded with a threshold value such that the  $B_{22}^{s^*}$ -norm of the resulting estimator is approximately equal to  $\|\tilde{f}_{v_0^*}\|_{B_{22}^{s^*}}$ . This is repeated for a mesh of values of  $v$  in  $[0, v_0^*]$ . An optimal value  $v^* \in [0, v_0^*]$  is finally selected by a customized criterion.

It is possible to start from the other side, too. Begin by determining an optimal threshold value for the "pure" thresholded estimator, compute from there an estimate for  $\|f\|_{B_{22}^s}$  for some  $s$ , then decrease the threshold and start trials with shrinking via (7) where  $\hat{\alpha}_{j_0k}$  and  $\hat{\beta}_{jk}$  are the coefficients of  $\hat{f}$  after thresholding at the current threshold level. Select the optimal proportion between shrinking and thresholding by a customized criterion.

The second approach, in the regression case, for Gaussian white noise, is to consider, in the level-dependent case or blockwise, the composite penalized "shrink/threshold" estimator with a SURE strategy of determining both the smoothing parameter  $v_j$  and the threshold level  $\lambda_j$ . Denote by  $SURE_{v,\lambda}(v, \lambda, \mathbf{x})$  and  $SURE_{\lambda,v}(v, \lambda, \mathbf{x})$  the SURE functionals corresponding to the "first shrink, then threshold" and the "first threshold, then shrink" estimators, respectively. In the notations and under the assumptions of [57], subsection 2.3, with  $\tau = v_j 2^{2j^*}$ ,  $j$  - fixed,  $t = \lambda_j$ ,  $d = c2^j$  (the number of non-zero empirical coefficients  $x_i$ ,  $i = 1, \dots, d$  at level  $j$ ),

$$\begin{aligned} SURE_{v,\lambda}(\tau, t, \mathbf{x}) &= \\ &= \frac{1-\tau}{1+\tau}d - \frac{2}{1+\tau} \cdot \#\{i : |x_i| \leq t(1+\tau)\} + \sum_{i=1}^d \left[ \min\left(t + \frac{\tau}{1+\tau}|x_i|, |x_i|\right) \right]^2, \\ SURE_{\lambda,v}(\tau, t, \mathbf{x}) &= SURE_{v,\lambda}\left(\tau, \frac{t}{1+\tau}, \mathbf{x}\right). \end{aligned}$$

The limiting case  $\tau = 0$  corresponds to Donoho and Johnstone's SURE( $t, \mathbf{x}$ ); the

other limiting case  $t = 0$  - to linear estimators (cf. Li [87]). Minimizing with respect to both  $\tau$  and  $t$  clearly gives more flexibility than in the limiting cases.

**B9.** *Shrinking in the general quasi-Banach case.* Formula (7) refers to the case  $p = q = \pi = u = 2$ ,  $\sigma = 0$  only. Now we consider the general case, as discussed in the end of Section 4.

(a) The penalization model is via

$$L(v, \hat{f}; B_{\pi\pi}^\sigma, \dot{B}_{pp}^s) := K_1(v, \hat{f}; (B_{\pi\pi}^\sigma)^\pi, (\dot{B}_{pp}^s)^p),$$

where  $p, \pi, s, \sigma$  are such that  $p < \infty, \pi < \infty$  and  $B_{\pi\pi}^\sigma \hookrightarrow B_{pp}^s$ . (Here  $A^\tau$  is the space  $A$  endowed with the quasi-norm  $\|\cdot\|_{A^\tau} := \|\cdot\|_A^\tau$ .) From the proof of Theorem 5.5.1 in [16] it follows that the formula about  $\tilde{\beta}_{jk}$  in (7) is now generalized to  $\tilde{\beta}_{jk} = \mu_{jk} \cdot \hat{\beta}_{jk}$ , where  $\mu_{jk}(v)$ :

$$(B3) \quad N(|1 - \mu_{jk}|^\pi, v 2^{j\varepsilon/2} |\hat{\beta}_{jk}|^{p-\pi} |\mu_{jk}|^p) = \min_{\mu \in \mathbb{R}} N(|1 - \mu|^\pi, v 2^{j\varepsilon/2} |\hat{\beta}_{jk}|^{p-\pi} |\mu|^p),$$

$N(x_1, x_2)$  being any fixed norm (or, more generally, any fixed quasi-norm equivalent to any norm) in  $\mathbb{R}^2$ , and  $\varepsilon = 2ps - 2\pi\sigma + p - \pi$  being the *critical-regularity* index, well-known in asymptotic-minimax theory (see [53], Theorem 1, for  $d = 1$ ). Then, there exist  $c, C : 0 < c \leq 1 \leq C < \infty$ , depending on  $N(\cdot, \cdot)$  only, such that

$$c.L(v, \hat{f}; B_{\pi\pi}^\sigma, \dot{B}_{pp}^s) \leq \|\hat{f} - \tilde{f}_v\|_{B_{\pi\pi}^\sigma}^\pi + v \|\tilde{f}_v\|_{\dot{B}_{pp}^s}^p \leq C.L(v, \hat{f}; B_{\pi\pi}^\sigma, \dot{B}_{pp}^s),$$

for any  $v > 0$  and any  $\hat{f} \in B_{\pi\pi}^\sigma + \dot{B}_{pp}^s$ . Moreover,  $c = C = 1$  if and only if  $N(x_1, x_2) = |x_1| + |x_2|$ .

If, in particular,

$$N(x_1, x_2) = N_\eta(x_1, x_2) = (|x_1|^\eta + |x_2|^\eta)^{1/\eta}, \quad \max\{1/p, 1/\pi\} < \eta \leq \infty,$$

then  $\mu_{jk}$  exists, is unique, and is the only root in  $[0, 1]$  of the equation for  $\mu$

$$(B4) \quad \pi^{\frac{1}{\eta}} (1 - \mu)^{\pi - \frac{1}{\eta}} |\hat{\beta}_{jk}|^{\pi - p} = v 2^{\frac{j\varepsilon}{2}} p^{\frac{1}{\eta}} \mu^{p - \frac{1}{\eta}}.$$

(Note that  $B_{\pi\pi}^\sigma \hookrightarrow B_{pp}^s$  is only possible if  $\pi \geq p$ .) In general, (B4) cannot be solved explicitly. However, for any  $(j, k)$  it can be solved by very quickly convergent iterations of the dyadic or Fibonacci bisection method. Moreover, in many important partial cases (B4) can be solved explicitly. The most important such case is  $\pi = p < \infty$ ,  $s > \sigma$ . Then,

$$(B5) \quad \mu_{jk} = \mu_j = \frac{1}{1 + v^{\frac{1}{p-\eta}} \cdot 2^{j(s-\sigma)\frac{p}{p-\eta}}}.$$

It can be seen that (7) corresponds to the partial case  $p = 2$ ,  $\eta = 1$ . The case when  $\pi \leq 1/\eta$  and/or  $p \leq 1/\eta$  can be handled via the same argument, but this case requires special considerations and will not be discussed here.

(b)  $\pi = p \leq \infty$ . In this case the penalization is via  $K_p(t, f; B_{pp}^\sigma, \dot{B}_{pp}^s)$ , and the embedding  $B_{pu}^\sigma \hookrightarrow B_{pq}^s$  holds,  $\sigma < s$ ,  $0 < u \leq \infty$ ,  $0 < q \leq \infty$ . Assuming that in the computation of the shrinking factor  $\mu_{jk}$  the quasi-norm  $N_\eta(x_1, x_2)$  is being utilized (see case (a)),  $\eta : 1/p < \eta \leq \infty$ , we see that  $\mu_{jk}$  is given by (B5) with  $v = t^p$ , and the case  $p = \infty$  is also included:

$$(B6) \quad \mu_{jk} = \mu_j = \frac{1}{1 + (t2^{j(s-\sigma)})^{\frac{p}{p-\eta}}},$$

which yields  $\mu_{jk} = \mu_j = \frac{1}{1+t2^{j(s-\sigma)}}$  for  $p = \infty$ . Denote by  $\tilde{f}_{1,t}$  the estimator  $\tilde{f}$  defined by (1'), with (7) replaced by (B6). Then,

$$(B7) \quad c.K_p(t, f; B_{pp}^\sigma, \dot{B}_{pp}^s) \leq (\|\hat{f} - \tilde{f}_{1,t}\|_{B_{pp}^\sigma}^p + t^p \|\tilde{f}_{1,t}\|_{\dot{B}_{pp}^s}^p)^{1/p} \leq C.K_p(t, f; B_{pp}^\sigma, \dot{B}_{pp}^s)$$

(with max-modification for  $p = \infty$ ), where  $c$  and  $C$  depend on  $\eta$  and  $c = C = 1$  if and only if  $\eta = 1$  or  $p = \infty$ .

**Remark .** It is seen from (B4-B6) that, the closer  $p$  is to  $1/\eta$ , the more sensitive the model is to variations of  $s$ ,  $\sigma$ ,  $t$  or  $v$ . The case  $p = 1/\eta$  will be dealt with in  $B^{10}$ .

**Remark .** The inequalities

$$\begin{aligned} \min(1, 2^{\frac{1}{\eta_2} - \frac{1}{\eta_1}})(|a|^{\eta_1} + |b|^{\eta_1})^{1/\eta_1} &\leq (|a|^{\eta_2} + |b|^{\eta_2})^{1/\eta_2} \leq \\ &\leq \max(1, 2^{\frac{1}{\eta_2} - \frac{1}{\eta_1}})(|a|^{\eta_1} + |b|^{\eta_1})^{1/\eta_1} \end{aligned}$$

show that in (B7)  $c = \min(1, 2^{\frac{1}{p}(\frac{1}{\eta}-1)})$ ,  $C = \max(1, 2^{\frac{1}{p}(\frac{1}{\eta}-1)})$  holds, and, unless  $p = \infty$ ,  $C$  becomes much larger than  $c$  as  $\eta$  approaches 0. In statistical context this means that the smaller  $\eta : 0 < \eta < 1$ , the larger the sample sizes  $n$  for which  $\tilde{f}_{1,t}$  is expected to perform well. On the other hand, for  $\eta : 1 \leq \eta \leq \infty$  the upper constant  $C$  is 1 and the smallest possible value for  $c$  is  $1/2$ , so  $\tilde{f}_{1,t}$  is expected to perform well already for moderate sample sizes. The properties of the estimator for  $p = \infty$  deserve special attention.

(c) *Determining the smoothness parameter.* In order to estimate  $v$  or  $t$  when utilizing the general shrinking formulae in (a) and (b), one may apply asymptotic-minimax considerations. In particular, all comments in (B2) are valid here, for  $1/\eta < p = \pi \leq \infty$  (see case (b)). This approach is expected to be effective for large samples. To improve the performance of  $\tilde{f}$ , defined by (1') and (B4), (B5) or (B6), on moderate samples, we suggest estimation with smoothness constraints<sup>B7</sup>. To this end, cross validation is performed as proposed in <sup>B8(b)</sup>. This yields the estimate  $\|\tilde{f}_{v_0^*}\|_{B_{22}^*}$  for  $\|f\|_{B_{22}^*}$ , as discussed there. The smoothing parameter  $t$  of  $\tilde{f}_{1,t}$  (see case (b)) is determined by requiring that  $\|\tilde{f}_{1,t}\|_{B_{22}^*} \approx \|\tilde{f}_{v_0^*}\|_{B_{22}^*}$ . For Gaussian white noise, the method of <sup>B8(b)</sup> can be applied. An extension of <sup>B6</sup> for this case is also of interest.

**B10. Thresholding - partial case of the penalized regularization approach.** *Relationship between the smoothing parameter and the threshold value.* Amato and Vuza [4] show that soft and hard thresholding, as well as other familiar thresholding methods can be obtained as partial cases of penalized  $L_2$  estimation. In our more general parameter setting in <sup>B9</sup>, their considerations about soft thresholding correspond to studying the behaviour of the solution

of (B5) when  $\pi = 2$ ,  $\eta = 1$  and  $p \searrow 1+$ . For hard thresholding they consider a penalized  $L_2$  model which is not of K-functional type, hence, the underlying functional space cannot be precisely identified. Here we indicate how hard thresholding can be obtained within the K-functional setting. Looking at (4), we see that the abundance of admissible couples  $(a, b)$ :  $a \in A$ ,  $b \in B$ ,  $a + b = \alpha$ , among which the minimum (infimum) is being sought, depends on how spacious  $A \cap B$  is. The more spacious this intersection is, the larger the class of admissible  $(a, b)$ . On the other extreme, if  $A \cap B = \{0\}$ , then the  $K$ -functional is an equivalent quasi-norm in the *direct sum* of  $A$  and  $B$ , and any  $\alpha \in A + B = A \oplus B$  is a unique sum of the unique projections of  $\alpha$  onto  $A$  and  $B$ , respectively. Thus, there is only one admissible couple  $(a, b)$  in this case and, of course, it is the optimal one, or, in other words, the variational nature of the model is *eliminated*. It turns out that classical hard thresholding strategies correspond to exactly such  $K$ -functionals with  $A \cap B = \{0\}$ , where the choice of  $B$  depends on the concrete thresholding strategy. Typically,  $A \oplus B = C$  is a fixed finite-dimensional space with dimension  $O(2^{j_1})$  in which  $\{\varphi_{j_0 k} : \text{supp } \varphi_{j_0 k} \cap \text{supp } \hat{f} \neq \emptyset\} \cap \{\psi_{j_1 k} : j = j_0, \dots, j_1, \text{supp } \psi_{j_1 k} \cap \text{supp } \hat{f} \neq \emptyset\}$  is a basis;  $\|\cdot\|_A$  and  $\|\cdot\|_B$  are any two fixed (quasi-)norms in finite dimensional space with dimension equal to that of  $C$ ;  $B$  is the (closed) linear span of those  $\psi_{j_1 k}$  in  $C$  for which  $|\langle \hat{f}, \psi_{j_1 k} \rangle|$  is above the threshold value. The threshold value is a monotone function of the regularization parameter  $t$  of the K-functional, and thus the space  $B$  and its complement  $A$  to  $C$  depend on  $t$ , too. This dependence is easiest to express for the  $K_\infty$ -functional. Denoting by  $\lambda = \lambda(t)$  the threshold level, and by  $\tilde{f}_\lambda$  the thresholded estimator, it can be seen that  $K_\infty(t, \hat{f}; A_t, B_t) = \inf_{j=j_1+j_2} \max(\|\hat{f}_1\|_{A_t}, t\|\hat{f}_2\|_{B_t}) = \max(\|\hat{f} - \tilde{f}_{\lambda(t)}\|_{A_t}, t\|\tilde{f}_{\lambda(t)}\|_{B_t})$ , where the threshold  $\lambda(t)$  is selected so that  $|\|\hat{f} - \tilde{f}_{\lambda(t)}\|_{A_t} - t\|\tilde{f}_{\lambda(t)}\|_{B_t}|$  is minimal. Here:  $\lambda(0) = 0$ ;  $A_0 = \{0\}$ ;  $B_0 = C$  (endowed with  $\|\cdot\|_B$ ); the first term  $\|\hat{f} - \tilde{f}_{\lambda(t)}\|_{A_t}$  is zero at  $t = 0$ , increases with the increase of  $t \geq 0$  and is  $\|\hat{f}\|_A$  for sufficiently large  $t$ , while the second term  $\|\tilde{f}_{\lambda(t)}\|_{B_t}$  is  $\|\hat{f}\|_B$  at  $t = 0$ , decreases with the

increase of  $t$  and is zero for  $t$  large enough.

In fact, the  $K$ -functional setting proposes also a new, unexplored thresholding method, different from the classical ones, with remarkable properties deserving much attention. We shall briefly describe it in item (b) below. Before this, in item (a) we shall complete our analysis of the case  $p = \pi$ , begun in  $B^9$ .

(a) *The case  $p = \pi = 1/\eta < \infty$  in  $(B_4)$ .* In this case formulae (B5) and (B6) degenerate, but  $\mu_{jk}$  can still be computed explicitly from (B3) and it turns out that

$$(B5') \quad \mu_{jk} = \mu_j = (2^{-pj(s-\sigma)} - v)_+^0.$$

We see that this is a classical *thresholding method localizing the spectral window and having band-limiting effect*. This shows that classical filtering methods can also be described in terms of appropriate choice of equivalent quasi-norms in the spaces involved in the  $K$ -functional. This clearly demonstrates once again that the  $K$ -functional setting provides a very general and *unified* methodological basis for the study, description and classification of smoothing, denoising and filtering techniques from variational point of view.

Completing our discussion of the case  $p = \pi$  begun in  $B^9(b)$ , we note that, together with other applications of theirs, the shrinking formulae obtained have potential applications in Bayesian statistics. For example, the values  $\mu_j$  in (B6) can be used to generate the parameters  $\pi_j$  as defined in [1], formula (10) (note that for  $\mu_j$  in (B6)  $\mu_j \leq \min(1, (2^{j\sigma}/t)^{p/(p-1/\eta)}, 2^{-jsp/(p-1/\eta)}) \leq 2\mu_j$  holds). For instance, it would be interesting to test the performance of the Bayesian model proposed in [1] with  $\pi_j = \min(1, \tilde{v}^{-1} \cdot 2^{-2j\tilde{s}})$ , corresponding to  $p = 2$ ,  $\eta = 1$ ,  $\sigma = 0$ , with  $\tilde{v}$  and  $\tilde{s}$  obtained by optimization of the GFCV criterion (see Subsection 6.1).

(b) In  $B^9(b)$  we saw that in the case of the embedding  $B_{pp}^\sigma \hookrightarrow B_{pp}^s$ ,  $s > \sigma$ , there is an explicit shrinking formula for the shrinking factor  $\mu_{jk}$  (formula (B6)). Another major embedding within the Besov-space scale is the *Sobolev-type* embedding:  $B_{\pi\pi}^\sigma \hookrightarrow B_{pp}^s$  if  $\sigma - \frac{1}{\pi} = s - \frac{1}{p}$  and  $0 < p \leq \pi \leq \infty$ . From (B4) it



can be seen that there is no explicit formula of the type of (B6) for this setting of  $\sigma$ ,  $s$ ,  $\pi$  and  $p$ . This makes even more interesting the fact that the  $K$ -functional approach indicates that there is a remarkable new thresholding method narrowly specialized for:  $\sigma$ ,  $s$ ,  $\pi$ ,  $p$ :  $\sigma - \frac{1}{\pi} = s - \frac{1}{p} =: \tau \in R$ ,  $0 < p \leq \pi \leq \infty$ . The essential additional fact here is that  $\tau$  is the same for both spaces in the  $K$ -functional. Because of this invariance of  $\tau$  we can invoke the following formula (B8)

$$K_1(v, f; L_\pi(U, d\mu), L_p(U, d\mu)) \asymp \left( \int_{v^{-\lambda}}^{\infty} f^*(\xi)^\pi d\xi \right)^{1/\pi} + v \left( \int_0^{v^{-\lambda}} f^*(\xi)^p d\xi \right)^{1/p},$$

where  $(U, d\mu)$  is an arbitrary measure space ( $d\mu$  can be also an atomic measure),  $f^*(\xi)$  is the decreasing rearrangement of  $f$  with respect to the measure  $d\mu$  (see [16], Section 1.3), and  $\lambda: \frac{1}{\lambda} = \frac{1}{p} - \frac{1}{\pi}$ . Formula (B8) cannot be found explicitly in our reference sources but it can be obtained from the Peetre-Kr  e formula ([16], Theorem 3.6.1, together with 3.14.5,6 and Theorem 5.2.1 (2) for  $q = p$  in their notation).

Suppose, as usual, that  $f$  and  $\psi$  are compactly supported. For simplicity of presentation, assume that  $j_0 = 0$  for any sample size  $n$ . Formula (B8) implies the following strategy for thresholding the  $\beta$ -wavelet coefficients of  $\hat{f} = \tilde{f}_0$  (see (1', 7) for  $v = 0$ ):

(b1) Consider all  $(j, k): \text{supp} \psi_{jk} \cap \text{supp} f \neq \emptyset$ . Denote the set of all such  $(j, k)$  by  $I(f, \psi)$ . We already know (see the proof of Theorem 1 in Appendix A) that the number  $M$  of elements of  $I(f, \psi)$  does not exceed  $c.2^{j_1}$ , where  $c = c(f, \psi)$ .  
 (b2) Consider the decreasing rearrangement  $\{b_\nu, \nu = 1, \dots, M\}$  of the finite set  $\{2^{j(\tau+1/2)}|\hat{\beta}_{jk}| : (j, k) \in I(f, \psi)\}$ . By definition,  $b_\nu \geq b_{\nu+1}$ ,  $\nu = 1, \dots, M-1$ ; for any  $\nu = 1, \dots, M-1$ , there exists a unique  $(j_\nu, k_\nu) \in I(f, \psi)$ , such that  $2^{j(\tau+1/2)}|\hat{\beta}_{jk}| = b_\nu$ ;  $(j_{\nu_1}, k_{\nu_1}) \neq (j_{\nu_2}, k_{\nu_2})$  if and only if  $\nu_1 \neq \nu_2$ . The new, hard-threshold wavelet estimator  $f_v^*$  is defined by

$$(B9) \quad f_v^*(x) = \sum_k \hat{\alpha}_{0k} \varphi_{0k}(x) + \sum_{\nu=1}^{\lfloor v^{-\frac{\pi p}{\pi-p}} \rfloor} \hat{\beta}_{j_\nu k_\nu} \psi_{j_\nu k_\nu}(x), \quad x \in R.$$

(b3) If  $b_{[v-\frac{\pi p}{\pi-p}]_{-1}} = b_{[v-\frac{\pi p}{\pi-p}]}$  up to some resolution threshold (tending to zero as  $n \rightarrow \infty$  and much smaller than the threshold levels in standard thresholding methods), then the terms corresponding to all the last  $\nu$ 's included in the RHS of (B9) which fall into this resolution threshold (that is, are formally equal to  $b_{[v-\frac{\pi p}{\pi-p}]}$ ) are being *removed*. Step (b3) ensures the uniqueness of  $f_v^*$ . If the sequence  $\{b_\nu\}$  is strictly decreasing (relative to the resolution threshold chosen), step (b3) is not needed.

Let us discuss (B9) in brief. The significance of every  $|\hat{\beta}_{jk}|$  is being assessed with respect to its level  $j$ . As  $v \rightarrow 0$  with  $n \rightarrow \infty$ , the most significant features appear first, and less significant details emerge only for sufficiently small  $v$ . The criterion for significance of the coefficients is the regularity assumption expressed in the value of  $\tau$ . Notice the qualitative differences between the strategies for "more regular functions" ( $\tau > -1/2$ ) and for "less regular functions" ( $\tau < -1/2$ ). The boundary value  $\tau = -1/2$  exactly corresponds to the critical regularity  $\varepsilon = -2(\pi - p)(\tau + 1/2) = 0$  (cf. [53], Theorem 1,  $d = 1$ ). We believe that the new thresholding approach suggested by (B9) can prove to be quite useful in applications involving composite estimators.<sup>B8</sup> It is also very appropriate for the case when  $n$  is a power of two and the orthogonal discrete wavelet transform (ODWT) is being applied.

(c) *Determining the smoothing parameter.* We suggest the same approaches as in <sup>B9(c)</sup>, although the algorithmic realization of the same ideas would obviously be quite different, notably involving sorting procedures.

**B11. Estimation in intersections of Besov spaces.** There is a theorem about explicit computation of K-functionals between intersections of quasi-normed lattices, due to Karadzhov [84], Theorem 1.2.51. Later results of Feichtinger and Gröchenig [63], Frazier and Jawerth [65] and Sickel [106] on atomic decompositions of Besov and Triebel-Lizorkin spaces, in combination with Karadzhov's result, made it possible to compute the K-functional between intersections of such spaces (with norm  $\|\cdot\|_{B_0 \cap B_1} = \max(\|\cdot\|_{B_0}, \|\cdot\|_{B_1})$ ). As a con-

sequence one obtains the following statement. Consider  $K(t, \hat{f}; B_{\pi u}^\sigma, B_{p_0 q_0}^{s_0})$ ,  $K(t, \hat{f}; B_{\pi u}^\sigma, B_{p_1 q_1}^{s_1})$ ,  $t > 0$ , where the Besov quasi-norms are given by (2),  $\hat{f}$  is the noisy wavelet expansion (from  $j_0$  to  $j_1$ ) and let  $\mu_{jk}^{(\nu)}(t) \in [0, 1]$  be such that  $\tilde{\beta}_{jk}^{(\nu)}(t) = \mu_{jk}^{(\nu)}(t) \cdot \hat{\beta}_{jk}$ , where  $\tilde{\beta}_{jk}^{(\nu)}(t)$  is the respective coefficient of the minimizer of the  $\nu$ -th K-functional,  $\nu = 0, 1$ . If  $\tilde{\beta}_{jk}(t)$  is the respective coefficient of the minimizer of  $K(t, \hat{f}; B_{\pi u}^\sigma, B_{p_0 q_0}^{s_0} \cap B_{p_1 q_1}^{s_1})$ , and if  $\mu_{jk}(t)$  is such that  $\tilde{\beta}_{jk}(t) = \mu_{jk}(t) \cdot \hat{\beta}_{jk}$ , then  $\mu_{jk}(t) = \min_\nu \mu_{jk}^{(\nu)}(t)$ .

**B12. More general function spaces with coefficient norms.** Recent applications of wavelet theory to the study of Hölder, cusp, singularity and oscillation spectra and to the development of multifractal formalism for functions require the consideration of more general types of function spaces. Here we consider two model types of such spaces.

(a) *The Nikol'skii-Besov spaces  $B_{pq}^{\omega(\cdot)}(R^d)$ , with norm*

$$\|f\| = \|f\|_{L_p(R^d)} + \left\{ \int_0^1 \left[ \frac{\omega_r(f; t)_{L_p}}{\omega(t)} \right]^q \frac{d\omega(t)}{\omega(t)} \right\}^{1/q} < \infty,$$

$r \in N$ ,  $1 \leq p \leq \infty$ ,  $1 \leq q \leq \infty$ , where  $\omega : [0, \infty) \rightarrow [0, \infty)$  is a nondecreasing continuous function with  $\omega(0) = 0$  and  $\left\{ \int_0^1 [t^r / \omega(t)]^q \frac{dt}{t} \right\}^{1/q} < \infty$ ,  $\omega_r(f; t)_{L_p}$  being the  $L_p$ -modulus of smoothness of  $f$  of order  $r$ , with step  $t$ . The classical theory (embeddings, traces, etc.) of these spaces has been developed in the 1970s and 1980s, notably by M. L. Gol'dman. To our best knowledge, after the series of papers on atomic decompositions and Riesz bases of orthogonal wavelets in 1989-1990, it seems to have remained unnoticed that the Nikol'skii-Besov spaces (which are translation invariant) do have the same Riesz bases of sufficiently regular orthogonal wavelets as usual Besov spaces and, hence, also admit equivalent coefficient norms. The only difference with usual Besov spaces (which may be very useful in the context of penalized statistical estimation, in particular, level and block shrinking strategies) is that now the weights on each level  $j$  depend on the more general function  $\omega$ .

Other generalizations of the already discussed scales of function spaces which are of certain interest to statistics-oriented applications are the *anisotropic* Besov and Triebel-Lizorkin spaces on domains and on manifolds in  $R^d$ , which can be of several different types - see, e.g., [111], Sections 10.1 and 10.2. The classical anisotropic Besov and Sobolev spaces (see, e.g., [93], p.153) are defined as the intersection of usual Besov or Sobolev spaces, of the same or of different metrics, with respect to partial or more general directional finite differences and derivatives. Another type of anisotropic Besov and Triebel-Lizorkin spaces is defined by the mixed quasi-norms given in [111], Section 10.2, formulae (1-3). This second type of anisotropic function spaces have been studied in detail in [104]. The first of these two types of anisotropic spaces is expected to play essential role in the development of the theory of multivariate regression-function and density estimation and of multivariate *functional data analysis* (see [101]). The second type of anisotropic spaces is expected to play major role in the *non-parametric statistical estimation of operators*, in particular, of *integral operators*. For spaces on the whole  $R^d$ , or on *unbounded* domains or manifolds in  $R^d$ , the *homogeneous* versions of the various types of anisotropic spaces should be used to take into account the role of the tail weights (see also Remark 2.2.4 in [48] and  $A^{14}$ ).

(b) *The A-spaces of Vasil A. Popov* have initially been defined by V. A. Popov as analogues of Besov spaces, where the integral modulus of smoothness (or  $\omega$ -modulus)  $\omega_r(f; t)_{L_p}$  has been replaced by the so-called average modulus of smoothness (or  $\tau$ -modulus)  $\tau_r(f; t)_{L_p}$  introduced by Bl. Sendov (see, e.g., [105,97,99,77,38-40]). More precisely, the inhomogeneous A-space  $A_{pq}^s$  has the quasi-norm

$$\|f\| = \|f\|_{L_p} + \left\{ \int_0^1 [t^{-s} \tau_r(f; t)_{L_p}]^q \frac{dt}{t} \right\}^{1/q} < \infty,$$

$r \in N$ ,  $s > 0$ ,  $0 < p \leq \infty$ ,  $0 < q \leq \infty$ , while the quasi-seminorm in the

homogeneous A-space  $\dot{A}_{pq}^s$  is

$$\|f\| = \left\{ \int_0^\infty [t^{-s} \tau_r(f; t)_{L_p}]^q \frac{dt}{t} \right\}^{1/q} < \infty.$$

Although not directly discussed in the recent survey [62] (which addresses only issues related to harmonic analysis) A-spaces and the  $A_{p,h}$  spaces mentioned below are perhaps the most practically important examples of Wiener amalgam spaces, due to their intrinsic relevance to pointwise approximation processes. In the 1980s these spaces proved to be very useful in the study of rates of convergence in approximation theory and numerical analysis, as well as for the characterization of best *one-sided* approximation. It is to be noted that (see [39])  $A_{pq}^s = B_{pq}^s$ ,  $\dot{A}_{pq}^s = \dot{B}_{pq}^s$  (with equivalence of norms) when  $0 < p \leq \infty$ ,  $0 < q \leq \infty$ ,  $s > d/p$ . Recently, Jaffard [78] has rediscovered these spaces (the spaces  $V^{s,p}$  (in his notations) are isomorphic to  $\dot{A}_{p\infty}^{s+d/p}$ ) from quite a different viewpoint - the study of multifractal properties of functions. He has been able to find a new insight into the essence of A-spaces by obtaining a wavelet characterization for  $\tau_r(f; t)_{L_p}$ . These results of Jaffard show that clearly the A-spaces will be an important space scale in future study of deterministic fractals and trajectories of random processes. In fact,  $\tau$ -moduli and A-spaces are of immediate use in nonparametric regression estimation with deterministic design, because the bias in estimating the wavelet coefficients, which is nothing else but the error of a quadrature formula, cannot be expressed in terms of the usual  $\omega$ -moduli, but in terms of the  $\tau$ -moduli (see, e.g., [105,38,40]). Hence, the bias cannot be measured in Besov spaces but only in terms of A-spaces. Therefore, estimation rates in the case of regression with deterministic design can be studied in Besov spaces only because these spaces are isomorphic to A-spaces when  $s' > 1/p$  (for dimension 1). The importance of the bound  $s' > 1/p$  has been discussed in a relevant context by Donoho [56], sections 2.3 and 2.4. In section 6.1 of [56] Donoho discusses also the critical case  $s' = 1/p$ . For this case Donoho proposes that interpolation and sampling (hence, also quadrature formulae) be studied in the class  $V_p$  of Bergh and Peetre [17], which includes the usual class with bounded Jordan

variation ( $p = 1$ ), and Wiener's quadratic variation ( $p = 2$ ) (see [122]). We agree with Donoho's analysis, and make the following additional comments. (i) The definition of  $V_p$  is a partial case of the more general concept of the  $\Phi$ -variation of Young [124], and other, more refined, choices of  $\Phi(t)$  (notably,  $\Phi(t) = t^p \Phi_1(t)$ , where  $\Phi_1$  has logarithmic size) are also of interest in the relevant context of error bounds for quadrature formulae. (ii) Considerations involving  $V_p$  are a partial case of the more general approach of  $\tau$ -moduli, since  $\tau_r(f; h)_{L_p} \leq c_r h^{1/p} (V^p f)^{1/p}$  holds,  $1 \leq p < \infty$ , where  $V^p f$  is the  $p$ -variation of  $f \in V_p$  (see Hristov [74], Dechevski [40]). (iii) The isomorphism  $\dot{A}_{p, \min(1, p)}^{1/p} = \dot{B}_{p, \min(1, p)}^{1/p}$  has been proved in [77] for the case  $p \in [1, \infty)$ , and can be proved also for  $p \in (0, 1)$  by a refinement of the technique in [39]. This indicates that  $\dot{B}_{p, \min(1, p)}^{1/p}$  is for fixed  $p$  the largest space in the Besov scale for which the error of a quadrature formula can be bounded uniformly in all elements of the Besov space. (iv) For the critical index  $s' = 1/p$ , in the case  $1 \leq p < \infty$ , the theory of A-spaces indicates that the natural space in which pointwise procedures (interpolation, sampling, etc.) should be studied is  $\dot{A}_{p, \infty}^{1/p}$  (or  $\dot{A}_{p, \infty}^{d/p}$  in the  $d$ -dimensional case). Note that  $V_p \hookrightarrow \dot{A}_{p, \infty}^{1/p}$  holds. (v) The multivariate theory of  $\tau$ -moduli and A-spaces has not been studied in such a detail as in the univariate case, since the great variety of interpolation problems in  $d$  dimensions,  $d > 1$  (see [21]) can be studied in terms of several different extensions of the concept of  $\tau$ -modulus to the multivariate case. Our opinion is that the most important of these extensions is based on the K-functional involving  $A_{p, h}$ -space (see [39] for  $d = 1$ ), because the space  $A_{p, h}$  is easy to define on general domains in  $R^d$ , and the resulting theory is analogous, though not identical, to the results of [80] for the  $\omega$ -modulus. The above-mentioned isomorphism  $A_{pq}^s = B_{pq}^s$ ,  $\dot{A}_{pq}^s = \dot{B}_{pq}^s$ ,  $s > d/p$ , refers to this type of  $\tau$ -moduli in the definition of A-spaces in  $d$ -dimensions. (vi) Summing up (i-v), the natural three-index space scale, in which approximation processes involving pointwise functionals can be considered for any  $s > 0$ , is  $\dot{A}_{pq}^s$ . This,

in particular, refers to the evaluation of the expected rates of nonparametric regression estimation with deterministic (uniform or non-uniform) design.

To illustrate the use of  $\tau$ -moduli and  $\Lambda$ -spaces, let us analyze the derivation of formula (A2) in detail.

Step 1. The error of the first-order quadrature formula in the LHS of (A2) is bounded in terms of the  $\tau$ -modulus (see [105]):

$$|E\hat{\beta}_{jk} - \beta_{jk}| \leq \tau_r(f\psi_{jk}; \frac{1}{n})_{L_p}, \quad r = 1, p = 1.$$

Step 2. The Leibniz formula for  $\tau$ -moduli ([40], Lemma 3.8.1) is applied:

$$\tau_r(g_1 g_2; h)_{L_p} \leq \sum_{\nu=0}^r \binom{r}{\nu} \tau_{r-\nu}(g_1; \frac{2rh}{r_\nu})_{L_{p_\nu}} \cdot \tau_\nu(g_2; 2h)_{L_{p'_\nu}},$$

where  $h > 0$ ;  $r_0 = r_1$ ,  $r_\nu = \nu$ ,  $\nu = 1, \dots, r$ ;  $0 < p \leq \infty$ ,  $\frac{1}{p_\nu} + \frac{1}{p'_\nu} = \frac{1}{p}$ ,  $\nu = 0, \dots, r$ ;  $\tau_0(g; h)_{L_p} := \|g\|_{A_{p,h}}$ ; the definition of the  $A_{p,h}$ -spaces is given in [38-40]. In the case of the first-order quadrature formula considered here,

$$\tau_1(f\psi_{jk}; \frac{1}{n})_{L_1} \leq \|f\|_{A_{p_0,1/n}} \tau_1(\psi_{jk}; \frac{1}{n})_{L_{p'_0}} + \tau_1(f; \frac{1}{n})_{L_{p_1}} \|f\|_{A_{p'_1,1/n}}.$$

Step 3. The bound

$$\|g\|_{A_{p,h}} \leq \|g\|_{L_p} + \tau_1(g; h)_{L_p}$$

(see [38], Lemma 1 and [40]), is applied to the result in Step 2, for  $g = f$ ,  $p = p_0$ , and for  $g = \psi_{jk}$ ,  $p = p'_1$ ;  $h = 1/n$ .

Step 4. The properties of the  $\tau$ -modulus (see the general references given above) are applied to obtain a bound of the error in terms of  $\Lambda$ -spaces, bounded  $p$ -variation and other relevant quantities.

Formula (A2) has been obtained via Steps 1-4, for  $p_0 = p_1 = \infty$ , in order to remain entirely within the Besov scale, since  $A_{\infty q}^s = B_{\infty q}^s$  for any  $s > 0$  and  $0 < q \leq \infty$ . The more advantageous selection (which, however, brings us out

of the Besov scale) is  $p_0 = p_1 = 2$ . For this choice, as discussed in the end of subsection 6.1, one obtains a relaxation of the assumption  $f \in B_{22}^{s'} \cap B_{\infty\infty}^{s_1}$  by replacing it with  $f \in A_{22}^{s'}$ . Formula (A3) has been obtained following Steps 1-4, too.

It is also informative to apply the above Steps 1-4 to outline the proof of a correct version of Lemma 3.1 in [8,9]. In our notations, the original statement of the lemma claims that if the scaling function  $\varphi$  is a coiflet of accuracy order  $r$  (that is, if  $\varphi$  is orthogonal to the monomials  $x^\nu$  of degree  $\nu = 1, \dots, r-1$ ) then  $|E\hat{\alpha}_{j_0k} - \alpha_{jk}| \leq c_s(f)2^{-2j_0s}$ , where  $c_s(f) < \infty$  if  $\|f\|_{B_{22}^s} < \infty$ . Our observation is that for  $s \leq 1/2$  a counterexample can be constructed (see [77]) which shows that this claim is not true. The correct version follows via Steps 1-4, with  $r$  in Step 1 being equal to the accuracy order of the coiflet, and with  $p_\nu = 2$ ,  $\nu = 0, \dots, r$ , in Step 2. The resulting sharp condition for the above bound on  $|E\hat{\alpha}_{j_0k} - \alpha_{jk}|$  to hold for any  $0 < s \leq r$  is  $f \in A_{2\infty}^s$ ,  $0 < s < r$ , and  $f \in A_{22}^s$ ,  $s = r$ . Hence, since  $A_{22}^s \hookrightarrow A_{2\infty}^s$  and  $B_{22}^s = A_{22}^s$  for  $s > 1/2$ , Lemma 3.1 in [8,9] is true if  $1/2 < s \leq r$ .

**B13. *K-functional estimation with constraints.*** In some applications it is desirable to preserve certain informative spatial features (e.g., translation-invariance ([34]), positivity and integral equal to one in density estimation ([98,96]) and others). In  $B^5$  and in Section 7 we already encountered examples of optimization with constraints in the K-functional setting. Further examples, with convex constraints, can be found in [52] and the references therein. Utilizing the general results of [88] about K-functionals with convex restrictions, it is now possible to solve nonparametric estimation problems with quite general convex constraints.

On the other hand, recall that wavelet atomic decomposition reduces nonparametric problems to essentially parametric ones. It is, therefore, interesting to find out to what extent the theory of parametric models (see [102]) can be helpful for solving constrained optimization problems in a wavelet setting.

The constrained approximation approach can help improve the performance of nonparametric density and regression-function penalized estimators



in the cases when additional information is a priori available about the shape of the curve (e.g., positivity, monotonicity, convexity,  $k$ -monotonicity, etc.). In the case of density estimation, imposing a non-negativity constraint on the solution of the K-functional penalized optimization problem, together with an isoperimetric type of constraint of type equality to 1 of its integral, would yield a solution which is a density. It is remarkable that this is achieved entirely within the framework of convex optimization, which means that finding the global extremum is guaranteed when using *any* numerical method for search of *local* extrema (compare also with the approaches in [98,96,45,46]).

**B14. Rates for the smoothing parameter in cross validation.** Continuing the discussion begun in the proof of Theorem 8 (see (A12)), we can obtain, for the density case, by bounding  $\frac{d}{dv} MISE(v)$  from below (under a sharpness assumption on  $s'$  if  $j_0 \rightarrow \infty$ ), the sharp upper rates for the smoothing parameter  $\tilde{v} = \operatorname{argmin} MISE(v)$ . On the other hand, by bounding  $\frac{d}{dv} MISE(v)$  from above, we can obtain (by using  $f \in B_{22}^{s'}$ , without an assumption about sharpness of  $s'$ ) the sharp lower rates for  $\tilde{v}$ . (For each fixed choice of  $j_0$  and  $j_1$  there would be one upper and one lower rate for  $\tilde{v}$  which are expected to coincide (up to a constant factor)). These rates can be then utilized to improve the performance of cross-validation (also with respect to choice of  $j_0$  and  $j_1$ , in the spirit of the ideas of Johnstone and Hall [81]). This approach can be transferred to Bowman-Rudemo type of cross validation in the case of regression (see [50]). It can also be extended, for both the density and regression estimation, to cross validation in Besov spaces with  $\sigma > 0$ , as considered in  $B^1$ . In principle, this method works also in the case of GFCV, but the difficulty in this case is that the rate of the optimal  $v$  gets "contaminated" by the error of the quadrature formula involved. In other words, varying  $s'$  within the interval  $[0, 1/2]$  is no longer sufficient to influence the rate for  $v$ . If, however, it is assumed that  $f \in A_{22}^{s'}$  rather than  $f \in B_{22}^{s'}$ , then  $s'$  controls the rate for  $v$  also when  $s' \in [0, 1/2]$  (see  $B^{12(b)}$ ). The same restriction appears, for the same reasons, in the theory of QGCV. In this case, the condition  $s > 1/2$  is needed also to ensure consistency

of QGCV (see [5]). For  $s > 1/2$  and  $s' = s$ , Amato and Vuza [6] found the sharp rates of  $v^* \rightarrow 0$  for the QGCV criterion by a technique which refers to the case  $j_0 = 0$ ,  $j_1 = \log_2 n - 1$ . We can generalize these results to include also the case  $j_0 \rightarrow \infty$  (under assumption about sharpness of  $s'$ ), as well as the case  $0 < s' \leq 1/2$  (under the assumption that  $f \in \Lambda_{22}^{s'}$ ). From the rates for  $v \rightarrow 0$  it is, of course, easy to obtain the rates for  $E\|\tilde{f}_v - f\|_{L_2}^2 \rightarrow 0$ . This approach can be extended to the generalizations considered in  $B^1$ , too. In [6] it was shown that minimizing the QGCV criterion for  $\sigma = 0$  leads to a consistent estimate of  $E\|\Pi_{j_1} f - \tilde{f}_v\|_{B_{22}^{\sigma}}^2 \rightarrow 0$ , for  $0 < \sigma < s' = s$ , where  $\Pi_{j_1}$  is as in  $B^1$ ,  $j_0 = 0$ ,  $j_1 = \log_2 n - 1$ . On the basis of this observation Amato and Vuza claim that QGCV is good enough to estimate the derivatives of the function (despite the fact that they proved that the minimizer  $v = v_\sigma$  of  $E\|\Pi_{j_1} f - \tilde{f}_v\|_{B_{22}^{\sigma}}^2$  (see  $B^1$ ) has rates which asymptotically split apart from the rates of the minimizer  $v^*$  of QGCV). We disagree with this opinion of Amato and Vuza, because it is based only on an asymptotic result about consistency of the estimator (cf. our Theorem 3 and Corollary 1). The more important result which has to be established here is, in our opinion, to derive an analogue of our Theorem 2 about consistency of the choice of the smoothing parameter when  $0 < \sigma < s' \leq s$  (see also  $B^1$ ). It is for this type of result that the advantage of  $v_\sigma$  over  $v^*$  will show up. It is also clear that, if  $s \geq s' > \sigma + 1/2 \gg 1$ , the difference between the estimator  $WAVREG_\sigma$  (see  $B^1$ ) and WAVREG of [9], corresponding to  $\sigma = 0$  will be quite visible for moderate samples. That is why it is important to study in detail the estimator  $WAVREG_\sigma$  proposed in  $B^1$ , with controlled selection of  $j_0$  and  $j_1$ .

The key technical result in the evaluation of the rates is

$$\sum_{j=j_0}^{j_1} \frac{2^{j\lambda}}{(1 + v2^{j\mu})^\nu} \asymp v^{-\lambda/\mu} \int_{v^{1/\mu}2^{j_0}}^{v^{1/\mu}2^{j_1+1}} \min(\xi^\lambda, \xi^{\lambda-\mu\nu}) \frac{d\xi}{\xi},$$

$\lambda > 0$ ,  $\mu > 0$ ,  $\nu > 0$ ,  $v \geq 0$ . It follows by change of variable in the integral from

the bounds

$$2^{-\lambda} \frac{2^{\eta\lambda}}{(1 + v2^{\eta\mu})^\nu} \leq \frac{2^{j\lambda}}{(1 + v2^{j\mu})^\nu} \leq \max(1, 2^{\mu\nu-\lambda}) \frac{2^{\eta\lambda}}{(1 + v2^{\eta\mu})^\nu}$$

and  $1 + v2^{\eta\mu} \asymp \max(1, v2^{\eta\mu})$ , where  $j \leq \eta \leq j+1$ . The first of these two bounds follows from the strict monotonicity of  $g_1(\zeta) = 2^\zeta$  and  $g_2(\zeta) = \frac{\zeta}{1+\zeta}$  for  $\zeta \geq 0$  (cf. [6], Lemma 3.1, for another bound used in their computations).

The above technical result is sufficient for finding the sharp rates in the density case and the regression case with random design under the assumptions that  $f \in B_{22}^{s'}$ ,  $s' > 0$ . If  $j_0 \rightarrow \infty$ , then also sharpness of  $s'$  is required. When the design is deterministic, the additional condition  $s' > 1/2$  has to be imposed. The case  $0 < s' \leq 1/2$  can be included by assuming that  $f \in A_{22}^{s'}$ ,  $s' > 0$ , but an additional step in the proof is needed, involving the equivalence of the  $\tau$ -modulus to a K-functional (see [38-40]).

As for WAVREG of [9], the sharp rates for  $Ev^* \rightarrow 0$  ( $v^*$  being here the optimizer of  $\tilde{A}_n$ ) and the respective estimation rates have not been studied at all. They, together with the respective rates for  $WAVREG_\sigma$ , can be obtained in a similar way.

**B15. Redesigning the experiment by K-functional techniques.** Recently, the development of discrete wavelet transform techniques requiring dyadic uniform design aroused renewed interest in methods for redesigning nonparametric regression models with random or irregular deterministic design (see [71,26,27]). Another important application of redesigning is in medical imaging, for the purposes of uniformization of data sets and registration of images (e.g., positron emission tomography (PET) images). In this case images are usually very smooth and not too spatially inhomogeneous, and the convolution method of [66] seems to be appropriate. For this particular type of applications we suggest a modification of the convolution approach based on the K-functional. In <sup>B4</sup> we suggested a kernel estimation technique via solutions of evolutionary differential equations. The numerical solution proposed was by FEM. Now we propose to use a difference scheme over a uniform mesh. The resulting numerical solution

provides the redesigned data set corresponding to the selected uniform mesh. The bias term in the estimation of the values on the mesh has been studied in detail (see [95,120,121,108,109,22,99,38,40]).

**B16. Penalized regularization of non-K-functional type.** The K-functional is a very general functional characteristics. Most generally, the spaces  $A$  and  $B$  in (4) can be any quasi-normed abelian groups. Such groups are metrizable, and so the K-functional can be extended to any metric spaces  $A$  and  $B$ , such that the sum  $a + b$ ,  $a \in A$ ,  $b \in B$ , is meaningful. It is possible to go beyond that. Amato and Vuza [4] give an example of penalized regularization where the penalty is of entropy type. This is one of their examples illustrating their general setting of solving  $\min_{x \in K} (f(x - y) + g(x))$  (Formula (1) in [4]). Here  $K \subset X$ , where  $X$  is a Banach lattice with order continuous norm;  $K$  is weakly closed;  $y \in D \subset X$ ;  $f : X \rightarrow R$  is continuous for the norm topology and lower semi-continuous for the weak topology;  $g : K \rightarrow R \cup \{\infty\}$  is lower semi-continuous for the weak topology. Under these general (and some additional) assumptions they study in Theorems 2.1 and 2.2 the existence of a weak solution of the optimization problem when  $y_n \rightarrow y$  in the norm topology, for some  $y \in D$ , as  $n \rightarrow \infty$ . Our experience with the diverse aspects of statistical estimation considered in the present paper tells us that this general optimization model needs to be extended in several directions: (a) one should consider a sequence of penalized optimization problems depending on  $n$ :  $\min_{x \in K} (f_n(x - y_n) + g_n(x))$  as  $y_n \rightarrow y$  (this corresponds to the case  $j_0 \rightarrow \infty$  (see our Section 6)); (b) one should be able to characterize simultaneously the family  $\{(f, g)\}$  such that  $\min_{x \in K} (f_1(x - y) + g_1(x)) \asymp \min_{x \in K} (f_2(x - y) + g_2(x))$  for every admissible  $y$  and for every  $(f_1, g_1)$  and every  $(f_2, g_2)$  from this family (this is important, because some of the  $(f, g)$  in the family may have better properties (e.g., smoothness, strict convexity) than others, and for the worse ones there may not exist a solution to the respective optimization problem); (c)  $X$  should be allowed to be quasi-Banach (note that it was possible to derive formulae (B3-B6) only because quasi-Banach spaces were considered in the K-functional; the method would not work in the narrower Banach setting); (d) it should be possible to characterize

the order of approximation to zero of  $\min_{x \in K} (f_n(x - y) + g_n(x))$ ,  $n \rightarrow \infty$ , when the *intermediate* (cf. our Section 4) set  $D : K \cap D \subset D \subset X$  varies from  $X$  to  $K \cap D$  (this is necessary for studying the bias term when deriving asymptotic rates of statistical estimation).

In the K-functional setting, objectives (a-d) are achieved in a natural way.

Finally, it should be noted that if some of the general assumptions of Amato and Vuza on  $f$  and  $g$  be made more stringent, then others may be relaxed, or stronger results about existence and/or uniqueness of the solution of the optimization problem can be achieved (cf., e.g., the classical result in Theorem 28.1 of [115]).

**B17. *K-functional with several penalization parameters.*** The models of level-dependent and block-penalized shrinking are examples of generalization of the K-functional for the case of vector penalizing parameter. The respective properties of the K-functional and the associated real-interpolation theory have been studied (see, e.g., [107]) and can offer further insight into level-dependent and block-shrinking strategies.

**B18. *Biorthogonal wavelets, wavelet packets and multiwavelets.***

(a) It is possible to consider *biorthogonal* atomic decompositions of Besov and more general spaces.<sup>B12(a)</sup> There are some advantages of biorthogonal compared to orthogonal wavelet approximation. For example, in signal and image processing it is often observed that biorthogonally compressed images have somehow better quality than orthogonally compressed ones. Another instance is that in biorthogonal wavelet analysis one can utilize polynomial spline-wavelets with compact support, which is impossible in the orthogonal case.

(b) The utilization of *wavelet packets* becomes essential when studying local self-similarity and local periodicity in processes which are globally non-selfsimilar and non-periodic. (One typical example is *turbulence*.) The penalized model can be useful here, too, in particular, if the penalty is of *entropy* type (cf. [4]). In the context of statistical nonparametric estimation, this allows combining adaptive smoothing and denoising with the Coifman-Wickerhauser

recursive best-base selection algorithm. The advantages for fractal estimation can be considerable when *locally* self-similar fractals are studied because, as already noted in Section 7, the choice of the underlying wavelet is of crucial importance.

(c) *Multiwavelets* have more degrees of freedom than usual (mono)wavelets, thanks to which they enjoy a number of desirable properties that monowavelets do not have. This refers to both wavelets and scaling functions (mother and father wavelets). In particular, multiwavelets can combine sufficient regularity with short support, symmetry and orthogonality to polynomials up to a high degree, which is theoretically impossible for monowavelets. Besides their general computational advantages, also for statistical estimation (see, e.g., [59] and the references therein), they have the potential to become a very effective tool in fractal-function estimation. Recalling the algorithm for generating the penalized self-similar fractal estimator (see Section 7), one can study in an analogous way fractals which are not self-similar but are comprised of a finite number of interlaced or superposed self-similar fractals. This may prove to be useful, e.g., in the study of growth of crystalline structures (cf. [10]).

**B19. Pattern recognition.** In our opinion, the most adaptive block-shrinking estimators can be obtained by introducing elements of pattern recognition into the block selection strategy. The pioneering work of Hall, Kerkyacharian and Picard [68] and subsequent work in [25] and [28] (combining the block strategy of Hall, Kerkyacharian and Picard in the partial case of Gaussian white noise with oracle inequalities and the James-Stein estimator proposed by Donoho and Johnstone for level-dependent shrinking when the white noise is Gaussian) can be considered as prototypes of results in this direction.

**B20. The orthogonal discrete wavelet transform (ODWT), fractal estimation, and parallel computing architectures (PCA).** In this paper we have given preference to empirical wavelet coefficient estimation for sample sizes  $n$  which are not necessarily powers of 2 and, therefore, ODWT (called also *fast wavelet transform*) cannot be applied. The cross validation considered is also not adapted for ODWT inversion. One reason for our choice is that, when uti-

lizing ODWT, one works with approximations of  $\varphi_{j_0k}$  and  $\psi_{jk}$  which depend on  $n$ . When  $n$  takes small or moderate values, and when the estimated function is not smooth, but a fractal, it is important to work with the "exact"  $\varphi_{j_0k}$  and  $\psi_{jk}$  (that is, with their values computed to sufficient precision independent of  $n$ , by the respective functional equations of type (16) by which  $\varphi$  and  $\psi$  are defined - see [36]). For example, this is the case when working with the self-similar fractal estimator considered in Section 7.

We also note that: (a) most of our considerations can be modified for sample sizes  $n = 2^N$  and adapted for the use of ODWT, the resulting algorithms for  $n = 2^N$  being fast, i.e.,  $O(n)$ -methods (see the sequel of this paragraph); (b) while on sequential computers ODWT has essential advantages, being an  $O(n)$ -method, when utilizing PCA with  $O(n)$  processors (this is realistic for moderate and small samples), both the ODWT and the method without use of ODWT are  $O(\log_2 n)$ -methods.

In order to justify our claim in (a), for sequential computing architecture, when  $n$  is a power of two, let us compare GFCV, as proposed in subsection 6.1, with the 'leave-half-out' or 'twofold' cross validation of Nason [92], QGCV of Amato and Vuza [5,6] and WAVREG of Antoniadis [9]. Nason's algorithm exploits in a nice way the symmetry on the real line to produce an  $O(n)$  algorithm when  $n = 2^n$ , or, roughly, an  $O(n^2)$ -algorithm for general  $n \in N$  when the 'leave-one-out' CV algorithm has to be applied (see [92], p.570). Like Nason's twofold method, QGCV and WAVREG are specially designed for  $n = 2^N$  and the use of ODWT, and it can be shown that they are also  $O(n)$ -methods. What about FCV and GFCV? For arbitrary  $n \in N$ , without making use of the fulfillment of the compatibility condition, FCV, as defined in (11), is an  $O(n^3)$ -method. Because of the validity of the compatibility condition, however, FCV can be computed by (12,13), instead of (11), by  $O(n^2)$  operations. For general  $n$  GFCV is also an  $O(n^2)$ -method. Now assume that  $n = 2^N$  and that FCV and GFCV are applied in combination with ODWT. In this setting,  $x$  ranges over the discrete mesh  $\{x_i\}_{i=1}^n$ , the vectors  $\{\varphi_{j_0k}(x_i)\}_{i=1}^n$ ,  $\{\psi_{jk_j}(x_i)\}_{i=1}^n$ ,  $k_j = 1, \dots, 2^j$ ,  $j = j_0, \dots, j_1$ ,

are replaced by the respective first  $2^{j_1+1}$  orthonormal basis vector-lines of the matrix  $\mathbf{W}$  of the corresponding direct ODWT; the vector of noisy wavelet coefficients  $\mathbf{w}$  is computed from the data vector  $\mathbf{y}$  by  $\mathbf{w} = \mathbf{W}\mathbf{y}$ . (For simplicity of presentation, here we assume the periodic setting, the case of boundary-adjusted wavelets and corresponding wavelet transform being similar.) Then, it can be shown that FCV becomes an  $O(n)$ -method. GFCV becomes a 'faster'  $O(n)$ -method than FCV, the constant factor in the  $O(n)$ -rate being smaller because of the simple formula for  $\bar{h}(v)$  in this case (see Remark 3). Thus, the numerical performance of FCV and GFCV is comparable with that of Nason's method in both cases of  $n$  being, or not being, a power of two. The advantages of full cross validation clearly show up in the multidimensional case, when the design is not rectangular and not uniform. The FCV criterion can be defined also when the design set  $\{x_i\}_{i=1}^n$  is a (uniformly or non-uniformly) scattered set in  $R^d$ , for any  $n \in N$ . The construction is based on Voronoi simplectification (VS, for short) and is discussed in more detail in <sup>B21</sup>. Being able to handle the case of scattered design in the multivariate case is a *very important advantage* of FCV which makes the method well adapted to deal with the challenging computational geometric aspects of the rapidly developing wavelet theory for  $d$ -dimensional domains with general types of boundary. This theory has already a considerable achievement for Lipschitz-graph domains (see [32]), and it is also imminent that in the near future there will be also a complete wavelet-based analogue of the convolution/local Taylor expansion technique of [83] of deriving Whitney-type extension theorems for Besov spaces on generalized Cantor sets with non-integer Hausdorff dimension in  $R^d$ . VS is usually a computationally intensive procedure, but for fixed design VS has to be applied only once. Then, for any random vector  $\{Y_i\}$  corresponding to the same fixed design set  $\{x_i\}$ , FCV (in its simplest version (11)) is an  $O(n^3)$ -method, for any dimension  $d$ . Moreover, the fulfillment of the compatibility condition for FCV makes it possible to obtain a computational reduction from  $O(n^3)$  to  $O(n^2)$  for any dimension  $d \in N$ . Proving this reduction is based on the multivariate Hermite expansion for scattered data and the density of  $C^\infty$  in Besov spaces. A sufficiently detailed



consideration of this topic (as well as of the possibility to define an analogue of GFCV in this case) would be too spacious and too technical to present here.

Finally, comparing the complexity of the optimization algorithm, GFCV, FCV, QGCV and WAVREG are easy winners over Nason's twofold and 'leave-out-one' algorithms, because in the former group of algorithms the dependence on the regularization parameter is analytic. It is possible, however, to modify the definition of Nason's CV estimator by replacing thresholding with non-threshold shrinking via (7) or, more generally, (B5) or (B6). For this modification, the dependence on the smoothness parameter in Nason's twofold and 'leave-out-one' algorithms becomes analytic, and full-scale theoretical analysis of these algorithms can be carried out, in analogy to the results discussed in subsection 6.1 and  $B^{14}$ , with controlled  $j_0$  and  $j_1$ :  $j_0 = O(1)$  or  $j_1 \rightarrow \infty$ ,  $j_1 = o(\log_2 n)$  or  $j_1 \asymp \log_2 n$ ,  $j_0 \leq j_1$ . After this, a theoretical comparison with FCV and GFCV can be made. In the multidimensional case clearly FCV is the more flexible method, available for much more general data sets.

**B21. The multivariate case.** All the main results for the basic penalized regression-function and density estimation models in Sections 2-6, as well as all extensions of these models considered in Section 7 and Appendix B, are also available in the  $d$ -dimensional case,  $d > 1$ . Here we give some additional details.

*Besov spaces of multivariate functions and multivariate wavelets.* There are several ways to define multivariate orthonormal wavelet bases. We consider here the tensor-product bases.<sup>A5</sup> With respect to these bases Besov spaces admit a coefficient quasi-norm<sup>A5</sup> which is a straightforward generalization of (2,3). In  $d$  dimensions the index of the dyadic level weights is  $s + d(1/2 - 1/p)$ . The Sobolev embedding is now true when the following inequality between the indices of two Besov spaces is fulfilled:  $B_{pq}^s \hookrightarrow B_{p_1q_1}^{s_1}$ , if  $s - d/p \geq s_1 - d/p_1$ ,  $0 < p \leq p_1 \leq \infty$ ,  $0 < q \leq q_1 \leq \infty$ . (In this paper we have often been dealing with the case  $d = 1$ ,  $p_1 = q_1 = \infty$ ,  $s_1 = 0$ ,  $B_{p_1q_1}^{s_1} = L_\infty$ ,  $p = q = 2$ .) The restriction on  $s$  when utilizing  $B_{\infty\infty}^r$ -regular wavelets with compact support is<sup>A5</sup>  $\max[0, d(1/p - 1)] < s < r$ .

*Cross validation.* The method for cross validation in the density esti-

mation setting is extended to the multidimensional case without any essential changes<sup>A5</sup>. The same refers to the Bowman–Rudemo approach<sup>B1</sup> in the case of regression–function estimation with random design. This is also true for Bowman–Rudemo’s and Wahba’s approach (in the latter case we have in mind FCV and GFCV) when the design of the nonparametric regression estimation problem is deterministic, but in this case an explanation is needed how to generate such design in  $d$  dimensions. If the sample size is of the form  $n = \prod_{\nu=1}^d n_{\nu}$ , where  $n_{\nu} \rightarrow \infty$ ,  $\nu = 1, \dots, d$ , as  $n \rightarrow \infty$ , a rectangular mesh  $\{\{x_{i\nu}\}_{i=1}^{n_{\nu}}\}_{\nu=1}^d$  can be considered (see <sup>B20</sup>). However, this can only be of interest for very specific problems. In general, for any  $n \in N$  consider the mesh  $\{x_i\}_{i=1}^n$ ,  $x_{i_1} \neq x_{i_2}$ ,  $i_1 \neq i_2$ ,  $x_i \in R^d$ , and construct the *Voronoi simplectification of the convex hull* of  $\{x_i\}_{i=1}^n$ . Then, the empirical wavelet coefficients in the  $d$ -dimensional noisy wavelet expansion are of the form

$$\hat{\alpha}_{j_0 k} = \sum_{i=1}^n w_i y_i \varphi_{j_0 k}(x_i), \quad \hat{\beta}_{j k}^{[l]} = \sum_{i=1}^n w_i y_i \psi_{j k}^{[l]}(x_i), \quad l = 1, \dots, 2^d - 1,$$

where the weight  $w_i$  is defined by  $w_i = \frac{1}{(d+1)V} \sum_{\nu=1}^{m_i} V_{i\nu}$ ,  $V_{i\nu}$  being the  $d$ -dimensional volumes of those simplices which have  $x_i$  as a vertex,  $m_i$  being the number of all such simplices,  $V$  being the  $d$ -dimensional volume of the convex hull of  $\{x_i\}_{i=1}^n$ . The wavelets  $\varphi_{j_0 k}$ ,  $\psi_{j k}^{[l]}$  are now the elements of the tensor-product orthonormal wavelet basis on  $L_2(R^d)$ , discussed in <sup>A5</sup>. The same weights appear in the quadrature formula of the GFCV (FCV) criterion. It can be seen that  $\sum_{i=1}^n w_i = 1$ . If the mesh is uniform, all simplices are uniform and congruent to each other, hence,  $w_i = \text{const}$  for all  $i$  such that  $x_i$  is in the interior of the convex hull. In the uniform case one can also construct a modification where all  $w_i$  are equal to  $1/n$ .

*Applications in Section 7* can be extended to the  $d$ -variate case. In particular, this is true for the fractal estimator and its version for Gaussian white noise.

*Kernel regularization.* The kernel is usually taken radial (e.g.,  $\Phi(\mathbf{x}) =$

$K_d \exp[-(\frac{\|\mathbf{x}\|^2}{1-\|\mathbf{x}\|^2})_+]$ , where  $\|\mathbf{x}\|$  is the usual Hilbert norm of  $\mathbf{x} \in R^d$ ), or as a normalized tensor product of univariate kernels. In combination with the tensor-product wavelet basis<sup>A5</sup> the second alternative is computationally more advantageous.

*The extensions in Appendix B* can all be generalized to the multivariate case. Some remarks are due, as follows.

Multivariate periodized wavelets are defined as tensor products of univariate periodized wavelets, and yield equivalent quasi-norms in Besov spaces, with equivalence constants depending on  $d$ .

In the  $d$ -dimensional case, the most important class of operators  $T$  considered in (B1) are partial differential operators which are *parabolic in the sense of Petrovskii* (see [108,109]). The respective partial finite difference schemes<sup>B15</sup> are *parabolic in the sense of John* (see [108,109]).

Formulae (B3,B4) remain the same, but *the critical regularity index*  $\epsilon$  involved in (B3,B4) now generalizes to  $\epsilon = 2ps - 2\pi\sigma + d(p - \pi)$  (cf. [53]). Formulae (B5,B5',B6) remain unchanged.

Formula (B9) also remains unchanged, as well as the definition of *decreasing rearrangement* in  $B^{10b,(b2)}$ , but the relevant value  $\tau$  now generalizes to  $\tau := \sigma - d/\pi = s - d/p$  (thus, the formula about  $\epsilon$  generalizing to  $\epsilon = -2(\pi - p)(\tau + d/2) = 0$  for  $\tau = -d/2$ ).

The remarks about the critical value  $s' = 1/p$  in  $B^{12,B14}$  are now relevant when  $s' = d/p$ .

### Acknowledgements

This paper has been written by L. T. Dechevsky, in collaboration with J. O. Ramsay, and is a result of L. T. Dechevsky's 15-year work on the study of  $K$ -functional penalization.

S. I. Penev is included as a joint author in acknowledgement, and as an appreciation, of his contribution to the earlier unpublished paper [49].

The numerical and graphical results related to Figures 1-4 of the present paper have been taken from the data files of [49].

We wish to thank Brenda MacGibbon (UQAM and GERAD) for her support and encouragement.

### References

- [1] F. A b r a m o v i c h, T. S a p a t i n a s, B. S i l v e r m a n. Wavelet thresholding via a Bayesian approach, *J. Royal Statist. Soc., Ser. B*, **60**, No 4, 1998, 725-749.
- [2] U. A m a t o, M. R. O c c o r s i o, D. T. V u z a. Fully adaptive wavelet regularization for smoothing data. In: *Proc. 3-rd Intern. Conf. Functional Analysis and Approximation Theory (Acquafredda di Maratea, Potenza' 1996)*, vol. 1 = *Suppl. Rendic. Circ. Mat. Palermo, Ser. 2* **52**, 1998, 207-222.
- [3] U. A m a t o, D. T. V u z a. Wavelet regularization for smoothing data, *Techn. Report CNR 108/1994*, Istituto per Applicazioni della Matematica, 1994.
- [4] U. A m a t o, D. T. V u z a. Besov regularization, thresholding and wavelets for smoothing data, *Numer. Funct. Anal. and Optimiz.* **18**, No 5-6, 1997, 461-493.
- [5] U. A m a t o, D. T. V u z a. Wavelet approximation of a function from samples affected by noise, *Rev. Roumaine Math. Pures Appl.* **42**, No 7-8, 1997, 481-493.
- [6] U. A m a t o, D. T. V u z a. Wavelet simultaneous approximation from samples affected by noise, *Comput. Math. Appl.* **36**, No 5, 1998, 101-111.

- [7] P. M. A n s e l o n e, P. J. L a u r e n t. A general method for the construction of interpolating or smoothing spline-functions, *Numerische Math.* **12**, 1968, 66-82.
- [8] A. A n t o n i a d i s. Smoothing noisy data with coiflets, *Statistica Sinica* **4**, 1994, 651-678.
- [9] A. A n t o n i a d i s. Smoothing noisy data with tapered coiflet series. *Scand. J. Statist.* **23**, 1996, 313-330.
- [10] A. A r n e o d o, F. A r g o u l, E. B a c r y, J. E l e z g a r a y, J. F. M u z y. *Ondelettes, multifractales et turbulences: de l'ADN aux croissances cristallines*, Diderot, Paris, 1995.
- [11] A. A r n e o d o, E. B a c r y, S. J a f f a r d, J. F. M u z y. Oscillating singularities on Cantor sets: a grand-canonical multifractal formalism, *J. Statist. Phys.* **87**, No 1/2, 1997, 179-209.
- [12] A. A r n e o d o, E. B a c r y, J. F. M u z y. Random cascades on wavelet dyadic trees, *J. Math. Phys.* **39**, No 8, 1998, 4142-4164.
- [13] A. B a r r o n, L. B i r g é, P. M a s s a r t. Risk bounds for model selection via penalization, *Preprint*, 1995.
- [14] A. B e n a s s i, S. C o h e n, J. I s t a s. Identifying the multifractal function of a Gaussian process, *Statist. Probab. Lett.* **39**, No 4, 1998, 337-345.
- [15] J. B e r a n. *Statistics for Long Memory Processes*. Monogr. on Statist. and Appl. Probab., **61**, Chapman and Hall, New York, 1994.
- [16] J. B e r g h, J. L ö f s t r ö m. *Interpolation Spaces. An Introduction*, Grundle. der math. Wiss., **223**, Springer, Berlin-Heidelberg-New York, 1976.
- [17] J. B e r g h, J. P e e t r e. On the spaces  $V_p$ , ( $0 < p \leq \infty$ ), *Bol. Unione Mat. Ital.* (4), **10**, 1974, 632-648.
- [18] D. P. B e r t s e k a s. *Constrained Optimization and Lagrange Multiplier Methods*, Acad. Press, New York-London-Paris, 1982.
- [19] L. B i r g é, P. M a s s a r t. From model selection to adaptive estimation, *Festschr. for Lucien Le Cam.*, Springer, New York, 1997, 55-87.
- [20] P. T. B o g g s, J. W. T o l l e. Sequential quadratic programming, *Acta Numerica*, 1995, 1-51.

- [21] B. D. B o j a n o v, H. A. H a k o p i a n, A. A. S a h a k i a n. *Spline Functions and Multivariate Interpolations*, Math. and Its Appl., **248**, Kluwer, Dordrecht, 1993.
- [22] P. B r e n n e r, V. T h o m é e, L. W a h l b i n. *Besov spaces and Applications to Difference Methods for Initial Value Problems*, Lecture Notes in Mathematics, **434**, Springer, Berlin-Heidelberg-New York, 1975.
- [23] O. B u n k e, B. D r o g e, J. P o l z e h l. Model selection and variable transformations in nonlinear regression, *CORE Discussion paper*, 1993.
- [24] P. L. B u t z e r, H. B e r e n s. *Semi-groups of Operators and Approximation*, Grundle. der Math. Wiss., **145**, Springer, Berlin-Heidelberg-New York, 1967.
- [25] T. T. C a i. Adaptive wavelet estimation: a block thresholding and oracle inequality approach, *Technical Report*, 1998.
- [26] T. T. C a i, L. D. B r o w n. Wavelet shrinkage for nonequispaced samples, *The Annals of Statist.* **26**, 1998, 1783-1799.
- [27] T. T. C a i, L. D. B r o w n. Wavelet estimation for samples with random uniform design, *Statist. Probab. Lett.* **42**, No 3, 1999, 313-321.
- [28] T. T. C a i, B. W. S i l v e r m a n. Incorporating information on neighboring coefficients into wavelet estimation, *Technical Report*, 1998.
- [29] T. T. C a i, B. W. S i l v e r m a n. Finite-sample performance of NeighBlock and NeighCoeffestimators, *Preprint*, 1998.
- [30] H. C h i p m a n, E. K o l a c z y k, R. M c C u l l o c h. Adaptive Bayesian wavelet shrinkage, *J. Amer. Statist. Assoc.* **92**, No 440, 1997, 1413-1421.
- [31] C. K. C h u i, E. Q u a k. Wavelets on a bounded interval, In: D. Braess, L. L. Schumaker (Eds.) *Numerical Methods of Approximation Theory*, vol. **9**, Intern. Ser. Numer. Math. **105**, Birkhäuser, Basel, 1992, 53-75.
- [32] A. C o h e n, W. D a h m e n, R. D e V o r e. Multiscale decompositions on bounded domains, *Trans. Amer. Math. Soc.*, to appear.
- [33] A. C o h e n, I. D a u b e c h i e s, P. V i a l. Wavelets and fast wavelet transforms on the interval, *Appl. Comput. Harmon. Anal.* **1**, 1994, 54-81.
- [34] R. R. C o i f m a n, D. L. D o n o h o. Translation invariant denoising.

- In: A. Antoniadis, G. Oppenheim (Eds.) *Wavelets and Statistics*. Lecture Notes in Statist. **103**, Springer, New York, 1995, 125-150.
- [35] D. D. C o x. Approximation of method of regularization estimators, *The Annals of Statistics* **16**, 1988, 694-713.
- [36] I. D a u b e c h i e s. *Ten Lectures on Wavelets*, SIAM, Philadelphia, 1992.
- [37] W. D a h m e n. Wavelet and multiscale methods for operator equations, *Acta Numerica*, 1997, 55-228.
- [38] L. T. D e c h e v s k i. Network-norm error estimates of the numerical solution of evolutionary equations. *Serdica* **12**, 1986, 53-64.
- [39] L. T. D e c h e v s k i.  $\tau$ -moduli and interpolation, In: M. Cwikel, J. Peetre, Y. Sagher, H. Wallin (Eds.) *Theory of Function Spaces and Applications*, Lecture Notes in Math. **1302**, Springer, Berlin-Heidelberg-New York, 1988, 177-190.
- [40] L. T. D e c h e v s k i. *Some Applications of the Theory of Function Spaces to Numerical Analysis*, Ph. D. Dissertation, Sofia University, 1988 (In Bulgarian).
- [41] L. T. D e c h e v s k i. On the constants of equivalence between some functional moduli and  $K$ -functionals, *Compt. Rend. Acad. Bulg. Sci.* **42**, No 2, 1989, 21-24.
- [42] L. T. D e c h e v s k y. On the sharp constants of equivalence between integral moduli of smoothness and  $K$ -functionals, *Research Report 288*, Center for Approximation Theory, Dept. of Math., Texas A&M Univ., 1993.
- [43] L. T. D e c h e v s k y. Explicit computation of the  $K$ -functional between Hilbert spaces. *Preprint 1579*, Fachbereich Mathematik, TH Darmstadt, 1993.
- [44] L. T. D e c h e v s k y, S. D u b u c. Multidimensional dyadic iterative interpolation and Fourier multipliers on Lebesgue spaces, *Research Report CRM-2549*, Centre de Recherches Mathématiques, Université de Montréal, Montreal, 1998.
- [45] L. T. D e c h e v s k y, B. M a c G i b b o n. Asymptotically minimax

- non-parametric function estimation with positivity constraints, I, *GÉRARD Research Report G-99-24*, École des Hautes Études Commerciales, École Polytechnique, McGill University, Université du Québec à Montréal, Montréal, 1999.
- [46] L. T. D e c h e v s k y, B. M a c G i b b o n, S. I. P e n e v. Numerical methods for asymptotically minimax non-parametric function estimation with positivity constraints, I, *Preprint*, 1999.
  - [47] L. T. D e c h e v s k y, S. I. P e n e v. On shape-preserving probabilistic wavelet approximators, *Stochast. Anal. and Appl.* **15**, No 2, 1997, 187-215.
  - [48] L. T. D e c h e v s k y, S. I. P e n e v. On shape-preserving wavelet estimators of cumulative distribution functions and densities, *Stochast. Anal. and Appl.* **16**, No 3, 1998, 428-469.
  - [49] L. T. D e c h e v s k y, S. I. P e n e v. On penalized wavelet estimation, *Report S98-15*, Dept. of Statist., School of Math., Univ. of New South Wales, Sydney, 1998.
  - [50] L. T. D e c h e v s k y, S. I. P e n e v. Weak penalized least squares wavelet regression estimation, *Report S99-1*, Dept. of Statist., School of Math., Univ. of New South Wales, Sydney, 1999.
  - [51] L. T. D e c h e v s k y, L.-E. P e r s s o n. On sharpness, applications and generalizations of some Carleman-type inequalities, *Tôhoku Math. J.* **48**, 1996, 1-22.
  - [52] M. D e l e c r o i x, M. S i m i o n i, C. T h o m a s-A g n a n. Functional estimation under shape constraints, *Nonparam. Statist.* **6**, 1996, 69-89.
  - [53] B. D e l y o n, A. J u d i t s k y. On minimax wavelet estimators, *Appl. Comput. Harmon. Anal.* **3**, 1996, 215-228.
  - [54] R. A. D e V o r e, G. K y r i a z i s, D. L e v i a t a n. V. M. T i k h o m i r o v, Compression and nonlinear  $n$ -widths, *J. Adv. Comp. Math.* **1**, 1993, 197-214.
  - [55] R. A. D e V o r e, B. J. L u c i e r. Fast wavelet techniques for near-optimal image processing, In: *Proc. IEEE Military Communications Conf. IEEE Communications Society*, New York, 1992.



- [56] D. L. D o n o h o. Interpolating wavelet transforms, *Technical Report*, Dept. of Statist., Stanford Univ., 1992.
- [57] D. D o n o h o, I. J o h n s t o n e. Adapting to unknown smoothness via wavelet shrinkage, *J. American Statist. Assoc.* **90**, No 432, 1995, 1200-1224.
- [58] D. L. D o n o h o, I. J o h n s t o n e, G. K e r k y a c h a r i a n, D. P i c a r d. Wavelet shrinkage: asymptopia? (With discussion), *J. Royal Statist. Soc., Ser. B* **57**, No 2, 1995, 301-369.
- [59] T. R. D o w n i e, B. W. S i l v e r m a n. The discrete multiple wavelet transform and thresholding methods, *IEEE Transact. on Signal Processing* **46**, 1998, 2558-2561.
- [60] B. D r o g e. Some comments on cross validation, Discussion paper 7, *Sonderforschungsbereich* **373**, Humboldt University, Berlin, 1994.
- [61] N. D u n f o r d, J. T. S c h w a r t z. *Linear Operators*, vol. **1-3**, Wiley, New York, 1967.
- [62] H. G. F e i c h t i n g e r. Amalgam spaces and generalized harmonic analysis, In: Proc. of the Norbert Wiener Centenary Congress (East Lansing, MI, 1994), *Proc. Sympos. Appl. Math.* **52**, AMS, Providence R. I., 1997, 141-150.
- [63] H. G. F e i c h t i n g e r, K.-H. G r ö c h e n i g. Banach spaces related to integrable group representations and their atomic decompositions, II, *Monatshefte Math.* **108**, 1989, 129-148.
- [64] G. M. F i k h t e n g o l ' t s. *A Course in Differential and Integral Calculus*, vol. **2** (7th ed.), Nauka, Moscow, 1969 (In Russian).
- [65] M. F r a z i e r, B. J a w e r t h. A discrete transform and decomposition of distribution spaces, *J. Functional Analysis* **93**, 1990, 34-170.
- [66] T. G a s s e r, H.-J. M ü l l e r. Kernel estimation of regression functions, In: T. Gasser, M. Rosenblatt (Eds.), *Smoothing Techniques for Curve Estimation*, Lecture Notes in Math. **757**, Springer, Heidelberg, 1979, 23-68.
- [67] W. H ä r d l e, G. K e r k y a c h a r i a n, D. P i c a r d, A. T s y b a k o v. *Wavelets, Approximation, and Statistical Applications*, Lecture Notes in

- Statist. **129**, Springer, New York, 1998.
- [68] P. H a l l, G. K e r k y a c h a r i a n, D. P i c a r d. Block threshold rules for curve estimation using kernel and wavelet methods, *The Annals of Statist.* **26**, 1998, 922-942.
  - [69] P. H a l l, P. P a t i l. On the choice of smoothing parameter, threshold and truncation in nonparametric regression by non-linear wavelet methods, *J. Royal Statist. Soc., Ser. B* **58**, No 2, 1996, 361-377.
  - [70] P. H a l l, S. P e n e v, G. K e r k y a c h a r i a n, D. P i c a r d. Numerical performance of block thresholded wavelet estimators, *Statistics and Computing* **7**, No 2, 1997, 115-124.
  - [71] P. H a l l, B. A. T u r l a c h. Interpolation methods for nonlinear wavelet regression with irregularly spaced design, *The Annals of Statist.* **25**, 1997, 1912-1925.
  - [72] N. H e c k m a n, J. O. R a m s a y. Some general theory for spline smoothing, *Preprint*, 1996; See also: J. O. Ramsay, N. Heckman. Some theory for  $L$ -spline smoothing, In: S. Dubuc, G. Deslauriers (Eds.) *Spline Functions and the Theory of Wavelets*. CRM Proc. Lecture Notes **18**, Centre de recherches mathématiques, Université de Montréal, AMS, Providence R.I., 1999, 371-380.
  - [73] D. H e n r y. *Geometric Theory of Semilinear Parabolic Equations*, Springer, Berlin-Heidelberg-New York, 1981.
  - [74] V. H. H r i s t o v. On the coefficients of Fourier-Lagrange, *Preprint* **5-81-318**, JINR, Dubna, 1981 (In Russian).
  - [75] E. H o p p e n s t a d t. Asymptotic series solutions of some nonlinear parabolic equations with a small parameter, *Arch. Rat. Mech. Anal.* **35**, 1969, 284-298.
  - [76] J. I s t a s, G. L a n g. Quadratic variations and estimation of the local Hölder index of a Gaussian process, *Ann. Inst. Henry Poincaré Probab. Statist.* **33**, No 4, 1997, 407-436.
  - [77] K. G. I v a n o v. On the behaviour of two moduli of functions, *Compt. Rend. Acad. Bulg. Sci.* **38**, No 5, 1985, 539-542.

- [78] S. J a f f a r d. Oscillation spaces: properties and applications to fractal and multifractal functions, *J. Math. Phys.* **39**, No 8, 1998, 4129-4141.
- [79] M. J a n s e n, M. M a l f a i t, A. B u l t h e e l. Generalized cross validation for wavelet thresholding, *Signal Processing* **56**, 1997, 33-44.
- [80] H. J o h n e n, K. S c h e r e r. On the equivalence of the  $K$ -functional and moduli of continuity and some applications, In: W. Schempp, K. Zeller (Eds.), *Constructive Theory of Functions of Several Variables*, Lecture Notes in Math. **571**, Springer, Berlin-Heidelberg-New York, 1977, 119-140.
- [81] I. M. J o h n s t o n e, P. G. H a l l. Empirical functionals and efficient smoothing parameter selection (With discussion), *J. Royal Statist. Soc., Ser. B* **54**, No 2, 1992, 475-530.
- [82] I. J o h n s t o n e, B. S i l v e r m a n. Wavelet threshold estimators for data with correlated noise, *J. Royal Statist. Soc., Ser. B* **59**, No 2, 1997, 319-351.
- [83] A. J o n s s o n, H. W a l l i n. *Function Spaces on Subsets of  $R^N$* , Math. Reports, vol. **2**, Part 1, Harwood, London, 1984.
- [84] G. E. K a r a d z h o v. *The "Means" Interpolation Method for Quasi-normed Spaces and Its Applications*, Ph.D. Dissertation, St. Petersburg Univ., 1973 (In Russian).
- [85] G. K e r k y a c h a r i a n, D. P i c a r d. Density estimation in Besov spaces, *Statist. Probab. Lett.* **18**, 1993, 327-336.
- [86] S. G. K r e i n. *Linear Differential Equations in Banach Space*, Nauka, Moscow, 1967 (In Russian).
- [87] K. C. L i. From Stein's unbiased risk estimates to the method of generalized cross-validation, *The Annals of Statist.* **13**, 1985, 1352-1377.
- [88] J. L ö f s t r ö m. Real interpolation with constraints, *J. Approx. Theory* **82**, No 1, 1995, 30-53.
- [89] B. M a c G i b b o n, R. v o n S a c h s. Nonparametric curve estimation by wavelet thresholding with locally stationary errors, *GÉRAD Research Report G-97-58*, École des Hautes Études Commerciales, École Polytechnique,

McGill University, Université du Québec à Montréal, Montreal, 1997.

- [90] Y. M e y e r. *Wavelets and Operators*, Cambridge Univ. Press, Cambridge, 1991.
- [91] Y. M e y e r. *Wavelets, Vibrations and Scalings*, CRM Monogr. Ser. **9**, Centre de recherches mathématiques, Université de Montréal, AMS, Providence R.I., 1998.
- [92] G. N a s o n. Wavelet shrinkage using cross validation, *J. Royal Statist. Soc., Ser. B* **58**, No 2, 1996, 463-479.
- [93] S. M. N i k o l ' s k i i. *Approximation of Functions of Several Variables and Imbedding Theorems*, Grundle. der math. Wiss. **205**, Springer, New York-Heidelberg-Berlin, 1975.
- [94] J. P e e t r e. *A Theory of Interpolation of Normed spaces*, Lecture Notes, Brasilia, *Notas de matematica* **39**, 1968, 1-86.
- [95] J. P e e t r e, V. T h o m é e. On the rate of convergence for discrete initial-value problems, *Math. Scand.* **21**, 1967, 159-176.
- [96] S. P e n e v, L. D e c h e v s k y. On non-negative wavelet-based density estimators, *Nonparam. Statist.* **7**, 1997, 365-394.
- [97] P. P. P e t r u s h e v, V. A. P o p o v. *Rational Approximation of Real Functions*, Encycl. of Math. and its Appl. **28**, Cambridge Univ. Press, Cambridge, 1987.
- [98] A. P i n h e i r o, B. V i d a k o v i c. Estimating the square root of a density via compactly supported wavelets, *Comput. Statist. Data Anal.* **25**, No 4, 1997, 399-415.
- [99] V. A. P o p o v, L. T. D e c h e v s k i. On the error of numerical solution of the parabolic equation in network norms, *Compt. Rend. Acad. Bulg. Sci.* **36**, No 4, 1985, 429-432.
- [100] W. P r e s s, S. T e u k o l s k y, W. V e t t e r i n g, B. F l a n n e r y. *Numerical Recipes in Fortran: the Art of Scientific Computing* (2nd ed.), Cambridge Univ. Press, Cambridge, 1992.
- [101] J. O. R a m s a y, B. W. S i l v e r m a n. *Functional Data Analysis*, Springer Ser. in Statistics, Springer, Berlin-Heidelberg-New York, 1997.

- [102] T. Robertson, F. T. Wright, R. L. Dykstra. *Order Restricted Statistical Inference*, Wiley, Chichester, 1988.
- [103] S. G. Samko, A. A. Kilbas, O. I. Marichev. *Integrals and Derivatives of Fractional Order and some of Their Applications*, Nauka i Tehnika, Minsk, 1987 (In Russian); English transl.: *Fractional Integrals and Derivatives: Theory and Applications*, Gordon & Breach, New York-London-Paris, 1993.
- [104] H.-J. Schmeisser, H. Triebel. *Topics in Fourier Analysis and Function Spaces*, Wiley, Chichester, 1985.
- [105] Bl. Sendov, V. A. Popov. *The Averaged Moduli of Smoothness. Applications in Numerical Methods and Approximation*, Wiley, Chichester, 1988.
- [106] W. Sickel. Spline representations of functions in Besov-Triebel-Lizorkin spaces on  $R^n$ , *Forum Math.* **2**, 1990, 451-475.
- [107] G. Sparrr. Interpolation of several Banach spaces, *Ann. Mat. Pura Appl.* **99**, 1974, 247-316.
- [108] V. Thomée. Stability theory for partial difference operators, *SIAM Review* **11**, No 2, 1969, 152-195.
- [109] V. Thomée. Convergence estimates in discrete initial value problems, *Actes Congrès Intern. Math., Nice'70*, 1971, 321-329.
- [110] K. Triboley. Practical estimation of multivariate densities using wavelet methods, *Statistica Neerlandica* **49**, 1995, 41-62.
- [111] H. Triebel. *Theory of Function Spaces*, Monogr. in Math. **78**, Birkhäuser, Basel-Boston-Stuttgart, 1983.
- [112] H. Triebel. *Theory of Function Spaces, II*, Monogr. in Math. **84**, Birkhäuser, Basel-Boston-Stuttgart, 1992.
- [113] H. Triebel. *Fractals and Spectra Related to Fourier Analysis and Function Spaces*, Monogr. in Math. **91**, Birkhäuser, Basel-Boston-Berlin, 1997.
- [114] F. D. Utreras. Cross-validation techniques for smoothing spline-functions in one or two dimensions, In: T. Gasser, M. Rosenblatt (Eds.),

- Smoothing Techniques for Curve Estimation*, Lecture Notes in Math. **757**, Springer, Heidelberg, 1979, 196-231.
- [115] M. M. V a i n b e r g. *Functional Analysis*, Prosveshtenie, Moscow, 1979 (In Russian).
- [116] B. V i d a k o v i c. Nonlinear wavelet shrinkage with Bayes rules and Bayes factors, *Discussion Paper 94-24*, ISDS, Duke Univ., 1994.
- [117] G. W a h b a. *Spline Models for Observational Data*, CBMS-NSF Regional Conference Ser. in Appl. Math. **59**, SIAM, Philadelphia, 1990.
- [118] G. G. W a l t e r, X. S h e n. Continuous non-negative wavelets and their use in density estimation, *Comm. Statist. Theory Methods* **28**, No 1, 1999, 1-17.
- [119] Y. W a n g. Fractal function estimation via wavelet shrinkage, *J. Royal Statist. Soc., Ser. B* **59**, No 3, 1997, 603-613.
- [120] O. B. W i d l u n d. On the rate of convergence for parabolic difference schemes, I, In: *Numerical Solution of Field Problems in Continuum Physics* (Proc. AMS Sympos. Appl. Math., Durham, N.C., 1968), SIAM-AMS Proc., vol. **II**, 1970, 60-73.
- [121] O. B. W i d l u n d. On the rate of convergence for parabolic difference schemes, II, *Comm. Pure Appl. Math.* **23**, 1970, 79-96.
- [122] N. W i e n e r. The quadratic variation of a function and its Fourier coefficients, *Mass. J. Math. Phys.* **3**, 1924, 72-94.
- [123] K. Y o s i d a. *Functional Analysis*, Grundle. der math. Wiss. **123**, Springer, New York-Heidelberg, 1974.
- [124] L. C. Y o u n g. Sur une généralisation de la notion de puissance  $p$ -ième bornée au sens de Wiener, et sur la convergence de series de Fourier, *Compt. Rend. Acad. Sci.* **204**, No 7, 1937, 470-472.

\* *Département de mathématiques et de statistique*

*Received 13.08.1998*

*Université de Montréal*

*C. P. 6128, Succursale A*

*Revised: 22.03.1999*

*Montréal, Québec, CANADA H3C 3J7*

*e-mail: dechevsk@dms.umontreal.ca*

*Second Revised: 19.11.1999*

**\*\* Department of Psychology**  
*McGill University*  
*1205 Dr. Penfield Ave.*  
*Montreal, Quebec, CANADA H4A 1B1*  
*e-mail: ramsay@psych.mcgill.ca*

**\*\*\* Department of Statistics**  
*School of Mathematics*  
*The University of New South Wales*  
*Sydney NSW 2052 AUSTRALIA*  
*e-mail: s.penev@unsw.edu.au*