



# Exact oracle inequality for a sharp adaptive kernel density estimator

Clementine Dalelane

► **To cite this version:**

Clementine Dalelane. Exact oracle inequality for a sharp adaptive kernel density estimator. 2005. <hal-00004753>

**HAL Id: hal-00004753**

**<https://hal.archives-ouvertes.fr/hal-00004753>**

Submitted on 19 Apr 2005

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# EXACT ORACLE INEQUALITY FOR A SHARP ADAPTIVE KERNEL DENSITY ESTIMATOR

**Clementine Dalelane**

Laboratoire de Probabilités et Modèles Aléatoires  
Université Pierre et Marie Curie Paris VI  
dalelane@ccr.jussieu.fr

April 19, 2005

## Abstract

In one-dimensional density estimation on i.i.d. observations we suggest an adaptive cross-validation technique for the selection of a kernel estimator. This estimator is both asymptotic MISE-efficient with respect to the monotone oracle, and sharp minimax-adaptive over the whole scale of Sobolev spaces with smoothness index greater than  $1/2$ . The proof of the central concentration inequality avoids “chaining” and relies on an additive decomposition of the empirical processes involved.

**Keywords:** kernel density estimator, MISE-optimal kernel, monotone oracle, minimax-adaptivity

**Mathematical Subject Classification:** 62G07, 62G20

## 1 Introduction

For many years, adaptive estimation procedures have stimulated the statistical interest. Such estimates achieve minimax convergence rates relying on very little prior knowledge about the properties of the curves to be estimated. Oracle inequalities fit into this framework, but give much more precise information about the performance of an estimate. They compare the risk of an adaptive candidate not to the minimax, but to the best possible risk. Oracle inequalities have been proposed for a variety of problems and estimator, Kneip (1994) and Donoho, Johnstone (1994) presumably being the first papers to state some. More recent examples are Hall, Kerkycharian, Picard (1999), Cavalier, Tsybakov (2001), Cay (2003), Efromovich (2004).

In the following we will consider oracle inequalities in density estimation, and it was wavelet estimators that have received the main attention in this context. During the 90’s, authors like Donoho, Johnstone, Hall, Kerkycharian and Picard developed various estimation techniques that satisfied more and more refined oracle inequalities. Efromovich also examined Fourier series estimates. Of course, even the case of data controlled bandwidth selection investigated in the 80’s, can be regarded as kind of an oracle problem. But the only source for a more general oracle inequality for a kernel density estimator is Rigollet (2004). Remarkably, unlike the other oracle inequalities on density estimators, Rigollet’s is an exact one. Our contribution gives another exact oracle inequality for kernel density estimation, but the two do not cover one another in neither direction.

Rigollet’s application of Stein’s blockwise estimator to non-parametric density estimation is a sharp minimax-adaptive kernel selection rule. The procedure approximates the so-called

*monotone oracle* by the use of kernel functions with piecewise constant Fourier transform. The monotone oracle is a pseudo-estimator, which minimizes the quadratic risk (MISE) over the class of all kernel functions, whose Fourier transform is real, symmetric and decreases monotonously on  $\mathbb{R}^+$ .

When considering the concept of curve smoothing from the viewpoint of signal recognition, a monotone Fourier transform appears to be a natural assumption to a kernel. Given that the unknown density is square-integrable, it is equivalent with respect to MISE either to estimate the density itself or to reconstruct it from an estimate of its Fourier transform. On the other hand it is known that with increasing frequency, random influences overbalance the true value in the empirical Fourier transform. For this reason, the Fourier series projection estimator omits empirical Fourier coefficients beyond a critical frequency. In the famous Pinsker-filter, the rigid cut-off is weakened to a monotone shrinkage of the unreliable coefficients by the Pinsker-weights. The focus on kernels with monotonously decreasing, but otherwise arbitrary Fourier transform is just a further generalization of the this notion.

The objective of the present work is to propose a purely data dependent estimator that approximates the monotone oracle in an exact oracle inequality. In comparison to Rigollet (2004), we abandon the assumption of the kernel's piecewise constant Fourier transform, our kernels only being band-limited to  $[-n, n]$  and having a monotone Fourier transform. Asymptotic exact MISE-efficiency is shown to hold over the set of all bounded,  $L_2$ -integrable densities, which are not infinitely differentiable. Sharp asymptotic minimax-adaptivity on the whole scale of Sobolev spaces with smoothness index greater than  $1/2$  follows automatically.

There are essentially two quantities to determine the statement of an oracle inequality: the set of estimators disposable to the minimization of risk; and the set of true parameters, over which the oracle inequality is supposed to hold. Evidently, the larger these sets, the stronger the oracle inequality. In a non-parametric setting, regularity conditions are the natural way to specify the space of parameters. Classes of estimators considered have been quite diverse and cover many familiar non-parametric estimation methods. However, all the classes, for which oracle inequalities were proven so far, share an important property – whether fixed or growing with the number of observations, their dimension is finite.

This is natural, when dealing with wavelet or Fourier coefficients. In ordered linear smoothers and blockwise Stein's method, the assumptions "ordered" and "blockwise", respectively, assure the finite dimensionality. For penalized least squares estimators, the dimension is always explicitly determined.

Note that oracle inequalities rely on special concentration inequalities, because it is necessary to approximate the maximum of an empirical process indexed by a class of functions; functional limit theorems are imperative. With finite dimension we have access to the uniform entropy of the estimator class and chaining arguments provide us with a suitable bound for the process. Yet these approximations unavoidably contain the factor dimension in one way or another.

Estimators indexed by kernels with monotone Fourier transform is a class that has obviously not finite dimension. But it is known that the set of monotone functions also allows for an approximation of its uniform covering number. Unfortunately, the approximations do not carry over from the Fourier to the space domain. And so the chaining approach is obstructed to us.

Instead, we pursue an alternative way to approximate our empirical process, namely an additive decomposition. The process, indexed by the class of kernels, is decomposed into a linear combination of countably many basis processes. Separate arguments as regards the basis processes (exponential inequalities) and the size of the non-random coefficients are

combined. The resulting threshold is equivalent to those in finite dimensional model classes, except that the factor containing the dimension is replaced by  $\ln(n)$ .

For an outline of the exact procedure, see section 4, appendix A1 and A2. The theorem along with the hypothesis is formulated in section 2. Section 3 contains the proof of the theorem relying on the proposition that the empirical process can be bounded to an appropriate magnitude. Some practical considerations will be found in section 5.

## 2 Main results

Let  $(X_1, X_2, \dots, X_n) \in \mathbb{R}^n$  be an i.i.d. sample with common density function  $f$ . Let the density  $f$  be bounded,  $\|f\|_\infty < \infty$ , have finite  $L_2$ -norm,  $\|f\|_2 < \infty$ , and denote by  $\hat{f}(\omega) = \int f(x)e^{ix\omega} dx$  the characteristic function of  $f$ . Let  $\tilde{f}_K(x)$  be the standard kernel estimator with kernel  $K$

$$\tilde{f}_K(x) = \frac{1}{n} \sum_{i=1}^n K(X_i - x) \quad (1)$$

and consider the quadratic risk

$$MISE(K) = E \int \left( \tilde{f}_K(x) - f(x) \right)^2 dx. \quad (2)$$

The cross-validation criterion

$$CV(K) := \int \tilde{f}_K^2(x) dx - \frac{2}{n(n-1)} \sum_{i \neq j} K(X_i - X_j) \quad (3)$$

is an unbiased estimator for MISE up to the summand  $\|f\|_2^2$ . Let  $\mathcal{K}$  be the set of all  $L_2$ -integrable kernel functions with real, symmetric, non-negative and unimodal Fourier transform  $\hat{K}(\omega) := \int K(x)e^{i\omega x} dx$ . For technical reason let  $\|K\|_2$  be  $\leq \sqrt{n}$ . This does not represent a real constraint, since the MISE of a sequence of kernels with  $L_2$ -norm growing faster than  $n$  cannot approach 0.

Define  $K^*$  to be the MISE-optimal kernel function for  $f$  and  $n$  among the class  $\mathcal{K}$ , i.e. the monotone oracle, and let  $K_0$  be the CV-optimal kernel function among  $\mathcal{K}$  restricted to kernels, whose Fourier transform additionally has support in  $[-n, n]$ .

$$\begin{aligned} K^* &:= \arg \min \left\{ MISE(K) \mid K \in \mathcal{K} \right\} \\ K_0 &:= \arg \min \left\{ CV(K) \mid K \in \mathcal{K}, \text{supp } \hat{K} \subseteq [-n, n] \right\} \end{aligned} \quad (4)$$

**Theorem** *Under the aforementioned hypotheses, for all  $\delta > 0$  the following exact oracle inequality holds:*

$$|E[MISE(K_0)] - MISE(K^*)| = O(n^{-\delta})MISE(K^*) + O(n^{\delta-1} \ln^{5/2} n)$$

**Remark 1** Although the theorem is stated for a fixed density, we could of course let  $f$  vary in some appropriate set. Investigating the influence of  $f$  on the asserted oracle inequality, we find that both residuals  $O(n^{-\delta})$  and  $O(n^{\delta-1} \ln^{5/2} n)$  contain constants depending on  $f$ : namely  $\|f\|_2$  and  $\max f$ . Obviously, these are uniformly bounded within Sobolev classes  $\mathcal{S}_\beta(L)$  with smoothness index  $\beta > 1/2$  ( $\mathcal{S}_\beta(L) \iff f \in L_2$  and  $\frac{1}{2\pi} \int |\omega^\beta \hat{f}(\omega)|^2 d\omega \leq L$ ). We will explicitly indicate those steps in the proofs, where the dependence enters our approximations.

**Corollary**  $\tilde{f}_{K_0}$  is asymptotically sharp minimax-adaptive on the whole scale of Sobolev classes with smoothness index greater than  $1/2$ .

**Remark 2** In case the true density  $f$  is not infinitely smooth, the assertion of the theorem is equivalent to a general MISE-efficiency, analogously defined to Hall (1983) and Stone (1984):

$$\frac{E[ISE(K_0)]}{MISE(K^*)} \longrightarrow 1$$

### 3 Proofs

**Proof of the Theorem** First of all, it can be seen that for  $L_2$ -integrable  $f$  the difference between  $MISE(K^*)$  and the MISE of a truncated version of  $K^*$  is negligible in proportion to  $MISE(K^*)$ . So the minimization of MISE on  $\mathcal{K}$  is equivalent to that on

$$\mathcal{K}_n := \{K \in \mathcal{K} \mid \text{supp } K \subseteq [-n, n]\}$$

Next let us assume the following propositions, the validity of which will be shown in section 4 by wavelet decomposition of the empirical processes: For any  $\lambda < \infty$ , there exists a set  $A_n \subseteq \mathbb{R}^n$ , such that for an arbitrary observation  $X = (X_1, \dots, X_n) \in A_n$  and for  $\delta > 0$  it holds that:

$$\begin{aligned} \mathbf{A1} \quad & |ISE(K) - \widetilde{CV}(K)| = O(n^{-\delta})MISE(K) + O(n^{\delta-1} \ln^{5/2}n) & \forall K \in \mathcal{K}_n \\ \mathbf{A2} \quad & |ISE(K) - MISE(K)| = O(n^{-\delta})MISE(K) + O(n^{\delta-1} \ln^{5/2}n) & \forall K \in \mathcal{K}_n \\ \mathbf{A3} \quad & P(X \in A_n^c) = O(n^{-\lambda}) \end{aligned}$$

where  $O(n^{-\delta})$  and  $O(n^{\delta-1} \ln^{5/2}n)$  do not depend on  $K$ . In case  $f$  is a density function that can only be estimated at a rate  $n^{\varepsilon-1}$ ,  $n^{-\delta}MISE(K)$  will dominate  $n^{\delta-1}$  for small enough  $\delta > 0$  at the right-hand side of these equations. Otherwise, if either  $\delta$  is too big or if  $f$  can be estimated at a faster rate, the term  $n^{\delta-1} \ln^{5/2}n$  will be dominating.

$\widetilde{CV}$  is a criterion derived from  $CV$ , such that  $\widetilde{CV}(K) - \widetilde{CV}(K') = CV(K) - CV(K')$  for any  $K, K'$  in  $\mathcal{K}$ , and will be defined below. In addition, it holds that:  $ISE(K) \leq (\|K\|_2 + \|f\|_2)^2 \leq (n^{1/2} + \|f\|_2)^2$ . As a consequence, we can proceed in the following way:

$$\begin{aligned} & E[ISE(K_0)] - MISE(K^*) \\ &= E[ISE(K_0) - ISE(K^*)] \\ &\leq E_{A_n}[ISE(K_0) - ISE(K^*)] + P(A_n^c) \sup_{K \in \mathcal{K}} ISE(K) \\ &= E_{A_n}[ISE(K_0) - \widetilde{CV}(K_0) + \widetilde{CV}(K_0) - \widetilde{CV}(K^*) + \widetilde{CV}(K^*) - ISE(K^*)] \\ &\quad + O(n^{-\lambda}) (\sqrt{n} + \|f\|_2)^2 \tag{A3} \\ &\leq E_{A_n}[ISE(K_0) - \widetilde{CV}(K_0)] + 0 + E_{A_n}[\widetilde{CV}(K^*) - ISE(K^*)] + O(n^{-\lambda+1}) \\ &= O(n^{-\delta})E_{A_n}[MISE(K_0)] + O(n^{-\delta})MISE(K^*) + O(n^{\delta-1} \ln^{5/2}n) + O(n^{-\lambda+1}) \tag{A1} \\ &= O(n^{-\delta})E_{A_n}[ISE(K_0)] + O(n^{-\delta})MISE(K^*) + O(n^{\delta-1} \ln^{5/2}n) \tag{A2} \end{aligned}$$

for  $\lambda$  sufficiently large. In order to return to  $E[ISE(K_0)]$ , we exert again proposition **A3** so as to find  $|E_{A_n}[ISE(K_0)] - E[ISE(K_0)]| = O(n^{-\lambda+1})$ , and therewith

$$\begin{aligned} E[ISE(K_0)] - MISE(K^*) &= O(n^{-\delta})E[ISE(K_0)] + O(n^{-\delta})MISE(K^*) + O(n^{\delta-1} \ln^{5/2}n) \\ \implies E[ISE(K_0)] &= \left(1 + O(n^{-\delta})\right) MISE(K^*) + O(n^{\delta-1} \ln^{5/2}n) \end{aligned}$$

By **A2**, the opposite is also readily shown:

$$\begin{aligned} MISE(K^*) - E[ISE(K_0)] &= O(n^{-\delta})E[ISE(K_0)] + O(n^{\delta-1} \ln^{5/2}n) \\ \implies MISE(K^*) &= \left(1 + O(n^{-\delta})\right) E[ISE(K_0)] + O(n^{\delta-1} \ln^{5/2}n) \end{aligned}$$

which implies the desired result:

$$|E[ISE(K_0)] - MISE(K^*)| = O(n^{-\delta})MISE(K^*) + O(n^{\delta-1} \ln^{5/2}n) \quad \square$$

**Proof of the Corollary** The minimax risk of density estimation in Sobolev classes  $\mathcal{S}_\beta(L) = \{f \in L_2 \mid \|f^{(\beta)}\|_2^2 \leq L\}$ , where  $\beta \in \mathbb{N}^+$  and  $L < \infty$ , is known since Efroimovich, Pinsker (1983). It is also known that kernel estimators employing suitable kernels maintain the minimax risk. One of these so-called minimax kernels is  $K_\beta$  with

$$K_\beta(x) = \frac{\beta!}{\pi} \sum_{j=1}^{\beta} \frac{\sin^{(j)} x}{(\beta-j)! x^{j+1}} \quad \text{and} \quad \widehat{K}_\beta(\omega) = \left(1 - |\omega|^\beta\right)_+$$

Obviously, the Fourier transform of any  $K_\beta$ ,  $\beta \in \mathbb{N}^+$ , is unimodal, so it is contained in  $\mathcal{K}$ . That means, the MISE-optimal estimator (monotone oracle)  $\widetilde{f}_{K^*}$  cannot be worse than the minimax estimator  $\widetilde{f}_{K_\beta}$ . On the other hand, the CV-optimal estimator  $\widetilde{f}_{K_0}$  is asymptotically as good as  $\widetilde{f}_{K^*}$ , where the convergence is uniform on Sobolev classes with  $\beta > 1/2$ , as emphasized in Remark 1 in section 2. It follows that  $\widetilde{f}_{K_0}$  is asymptotically minimax simultaneously on the scale of Sobolev classes  $\mathcal{S}_\beta(L)$  with  $\beta \in \mathbb{N}^+$ .

The consideration of the minimax risk of density estimators can be extended to Sobolev type classes with non-integer smoothness index  $\beta \in \mathbb{R}^+$ , defined as

$$\mathcal{S}_\beta(L) := \left\{ f \in L_2 \mid \frac{1}{2\pi} \int |\omega^\beta \widehat{f}(\omega)|^2 d\omega \right\}$$

For  $\beta > 1/2$ , both the minimax risk and the minimax kernel  $K_\beta$  take forms analogous to those in ordinary Sobolev classes, although the proofs have to be adjusted (see Dalelane (2005)). The same idea as before leads to simultaneous asymptotic minimaxity of  $\widetilde{f}_{K_0}$  on the whole scale of Sobolev type classes  $\mathcal{S}_\beta(L)$  with  $\beta \in \mathbb{R}^+$ ,  $\beta > 1/2$ .  $\square$

## 4 The empirical process

As the proof of proposition **A2** is very much the same as the one for **A1**, we confine ourselves to a demonstration of how  $|ISE(K) - \widehat{CV}(K)|$  can be approximated by  $O(n^{-\delta})MISE(K) + O(n^{\delta-1} \ln^{5/2}n)$  simultaneously over  $\mathcal{K}_n$ . The first step towards this goal will be to split up the difference between ISE and CV into two empirical U-processes indexed by  $K_n$ , a degenerate U-process of order 2 and a U-process of order 1, i.e. a partial sum process. This splitting was already observed in Stone (1984), where the class of kernels consists but of one rescaled

kernel function:  $\mathcal{H}_K = \{K_h | h > 0\}$ . Obeying some assumptions on  $K$ , it is easy to bound the uniform covering number of  $\mathcal{H}_K$ , see Nolan, Pollard (1987). Chaining arguments apply to both the partial sum process and the empirical U-process. But for lack of an appropriate approximation on  $\mathcal{K}_n$ , a generalization of Nolan/Pollard's proof is not possible.

Instead, we define a wavelet inspired function basis for  $\mathcal{K}_n$ , such that every kernel  $K \in \mathcal{K}_n$  can be represented as a linear combination of the functions belonging to this basis. The linear decomposition is carried forward to the space of U-statistics made up by  $\mathcal{K}_n$ , such that each U-statistic in  $K$  is a weighted sum of all (countably many) U-statistics of the function basis. The values of the basic U-statistics can be controlled by means of exponential inequalities. On a set of "favorable events" with overwhelming probability (proposition **A3**), they do not exceed a comfortable threshold of  $n^{-1} \lambda \ln^{3/2} n$ . In turn, due to the unimodality of the kernels' Fourier transforms, we can bound the absolute sum of the (non-random) wavelet coefficients, assigning a linear combination of basic U-statistics to a given U-statistic in  $K$ , through  $\ln n \|K\|_2$ . Combining these arguments, we find that any U-statistic in  $K$  is an  $O(n^{-1} \ln^{5/2} n) \|K\|_2 = O(n^{-1/2} \ln^{5/2} n) \sqrt{MISE(K)}$ , the  $O$ 's neither depending on  $K$  nor on  $f$ .

To derive the desired bound of  $O(n^{-\delta})MISE(K^*) + O(n^{\delta-1} \ln^{5/2} n)$  therefrom, we have to differentiate several constellations between the true density  $f$  and the envisaged  $\delta$ . Recall that the monotone oracle-kernel  $K^*$  is not random and depends on nothing but  $f$  and  $n$ .

First consider  $f$  such that there exist constants  $0 < l_f, u_f < \infty$  and  $\varepsilon_f > 0$ , which satisfy  $l_f \cdot n^{\varepsilon_f - 1} \leq MISE(K^*) \leq u_f \cdot n^{\varepsilon_f - 1}$ . If  $\delta < \varepsilon_f/2$ , then we have immediately

$$O(n^{-1/2} \ln^{5/2} n) \sqrt{MISE(K^*)} < O(n^{-\delta})MISE(K^*).$$

If otherwise  $\delta \geq \varepsilon/2$  holds, it follows that

$$O(n^{-1/2} \ln^{5/2} n) \sqrt{MISE(K^*)} \leq O(n^{\delta-1} \ln^{5/2} n).$$

This second reasoning is also true, when the convergence rate of  $MISE(K^*)$  is inferior to  $n^{\varepsilon-1}$  for any  $\varepsilon > 0$ , i.e. if the density  $f$  has infinitely many derivatives.

By a similar procedure but employing a different function basis, we also approximate the partial sum process. But this is already proposition **A1**.

To be exact, let  $X_1, \dots, X_n$  be distributed as assumed in section 2. Let  $X$  and  $Y$  denote two further random variables with the same distribution, independent of  $X_1, \dots, X_n$  and of each other.

$$\begin{aligned} ISE(K) &:= \int \left( \tilde{f}_K(x) - f(x) \right)^2 dx = \int \tilde{f}_K^2(x) dx - \frac{2}{n} \sum_{i=1}^n E[K(X_i - X) | X_i] + E[f(X)] \\ CV(K) &= \int \tilde{f}_K^2(x) dx - \frac{2}{n(n-1)} \sum_{i \neq j} K(X_i - X_j) \end{aligned}$$

We obtain  $\widetilde{CV}$  from  $CV$  by adding a zero and a further term which does not depend on  $K$ . Define  $I_n(\omega) := I(|\omega| < n)$  and

$$h_f(x) := \frac{1}{2\pi} \int \hat{f}(\omega) \left( 1 - I_n(\omega) \right) e^{-i\omega x} d\omega \quad (5)$$

the high-frequency contribution of  $\widehat{f}$  to  $f$ .

$$\begin{aligned} \widetilde{CV}(K) &:= CV(K) + \left[ \frac{2}{n} \sum_{j=1}^n E[K(X - X_j)|X_j] - 2E[K(X - Y)] - \frac{2}{n} \sum_{j=1}^n E[K(X - X_j)|X_j] \right. \\ &\quad \left. + 2E[K(X - Y)] \right] + \frac{2}{n} \sum_{j=1}^n \left( f(X_j) - h_f(X_j) \right) - 2E[f(X_j) - h_f(X_j)] \end{aligned}$$

We can now split up the difference between the quadratic loss and the cross-validation criterion into two summands:

$$\begin{aligned} &ISE(K) - \widetilde{CV}(K) \\ &= -\frac{2}{n} \sum_{i=1}^n E[K(X_i - Y)|X_i] + E[f(X)] + \frac{2}{n(n-1)} \sum_{i \neq j} K(X_i - X_j) \\ &\quad - \frac{2}{n} \sum_{j=1}^n E[K(X - X_j)|X_j] + 2E[K(X - Y)] + \frac{2}{n} \sum_{j=1}^n E[K(X - X_j)|X_j] \\ &\quad - 2E[K(X - Y)] - \frac{2}{n} \sum_{j=1}^n \left( f(X_j) - h_f(X_j) \right) + 2E[f(X_j) - h_f(X_j)] \\ &= \frac{2}{n(n-1)} \sum_{i \neq j} \left( K(X_i - X_j) - E[K(X_i - X_j)|X_i] - E[K(X_i - X_j)|X_j] \right. \\ &\quad \left. + E[K(X_i - X_j)] \right) + \frac{2}{n} \sum_{j=1}^n \left( E[K(X - X_j)|X_j] - f(X_j) + h_f(X_j) \right. \\ &\quad \left. - E[K(X - X_j)] - E[f(X_j) - h_f(X_j)] \right) \\ &=: \frac{2}{n(n-1)} \sum_{i \neq j} U_K(X_i, X_j) + \frac{2}{n} \sum_{j=1}^n \left( b_K(X_j) + h_f(X_j) - E[b_K(X_j) + h_f(X_j)] \right) \quad (6) \end{aligned}$$

where  $b_K$  stands for the bias  $E\widetilde{f}_K - f$ . The first term corresponds to a degenerate U-statistic, since  $E[U_K(X, Y)|Y] = E[U_K(X, Y)|X] = E[U_K(X, Y)] = 0$  for all values of  $X$  and  $Y$ . In appendix A1, we will define a basis of father and mother wavelets for  $\mathcal{K}_n$ , which allows the following decomposition:

$$K(x) = \sum_t \alpha_t(K) \varphi_t(x) + \sum_{s,t} \beta_{st}(K) \psi_{st}(x)$$

This decomposition can also be assigned to the U-statistics, such that

$$\frac{1}{n(n-1)} \sum_{i \neq j} U_K(X_i, X_j) := \frac{1}{n(n-1)} \sum_{i \neq j} \left[ \sum_t \alpha_t(K) U_{\varphi_t}(X_i, X_j) + \sum_{s,t} \beta_{st}(K) U_{\psi_{st}}(X_i, X_j) \right]$$

A change of summation separates the stochastic processes from the deterministic coefficients.

$$\begin{aligned} \frac{1}{n(n-1)} \sum_{i \neq j} U_K(X_i, X_j) &= \sum_t \alpha_t(K) \left[ \frac{1}{n(n-1)} \sum_{i \neq j} U_{\varphi_t}(X_i, X_j) \right] \\ &\quad + \sum_{s,t} \beta_{st}(K) \left[ \frac{1}{n(n-1)} \sum_{i \neq j} U_{\psi_{st}}(X_i, X_j) \right] \end{aligned}$$



The basic U-statistics can be kept “small” on a set of “favorable events”  $A_{n1} \subseteq \mathbb{R}^n$  (see appendix A1), and in Lemma 1 we find bounds for the wavelet coefficients, so that on  $A_{n1}$  the following holds

$$\frac{1}{n(n-1)} \sum_{i \neq j} U_K(X_i, X_j) = O\left(\frac{\ln^{5/2} n}{n}\right) \left(\sqrt{\int K^2(x) dx} + 1\right) \quad (7)$$

and for sufficiently large  $\lambda < \infty$ , equation (10) in Lemma 2 shows that

$$P(A_{n1}^c) = O(n^{-\lambda^{2/3}+1})$$

On the other hand, we obtain in (6) a partial sum in  $b_K + h_f$ . Because of the bounded support of  $\widehat{K}$ ,  $b_K + h_f$  takes the form:

$$\begin{aligned} b_K(x) + h_f(x) &= f * K(x) - f(x) + h_f(x) \\ &= \frac{1}{2\pi} \int \widehat{f}(\omega) (\widehat{K}(\omega) - 1) e^{-i\omega x} d\omega + \frac{1}{2\pi} \int \widehat{f}(\omega) (1 - I_n(\omega)) e^{-i\omega x} d\omega \\ &= \frac{1}{2\pi} \int \widehat{f}(\omega) (\widehat{K}(\omega) - 1) I_n(\omega) e^{-i\omega x} d\omega \end{aligned}$$

It is the low-frequency component of the bias and exactly that part which really depends on the kernel. In appendix A2, the partial sum is bounded on another set of “favorable events”  $A_{n2} \subseteq \mathbb{R}^n$ .

$$\begin{aligned} &\frac{1}{n} \sum_{j=1}^n (b_K(X_j) + h_f(X_j)) - E[b_K(X_j) + h_f(X_j)] \\ &= O\left(\frac{\ln^2 n}{\sqrt{n}}\right) \left(\sqrt{\int b_K^2(x) dx} + \frac{\|f\|_2}{\sqrt{n}}\right) \end{aligned} \quad (8)$$

and in equation (12) Lemma 4 we see that

$$P(A_{n2}^c) = O(n^{-\lambda+1})$$

The intersection of these two sets of “favorable events”  $A_{n1} \cap A_{n2} =: A_n$  is the one used in section 3 to bound  $\widetilde{CV}$ -ISE (on the very same set  $A_n$ , ISE-MISE can be bounded to an identical size).

The threshold for the U-statistic is of order  $n^{-1/2} \ln^{5/2} n (\sqrt{MISE(K)} + n^{-1/2})$ . And the one for the bias is of order  $n^{-1/2} \ln^2 n (\sqrt{MISE(K)} + n^{-1/2})$ , but depends on  $\|f\|_2$ . When  $\|f\|_2$  is uniformly bounded, as in Sobolev classes with smoothness index greater than 1/2, also this approximation is uniform. Besides, MISE converges in any case not faster than  $n^{-1}$ . Hence:

$$\begin{aligned} &|ISE(K) - \widetilde{CV}(K)| \\ &\leq 2 \left| \frac{1}{n(n-1)} \sum_{i \neq j} U_K(X_i, X_j) \right| + 2 \left| \frac{1}{n} \sum_{j=1}^n (b_K(X_j) + h_f(X_j)) - E[b_K(X_j) + h_f(X_j)] \right| \\ &= O\left(\frac{\ln^{5/2} n}{n}\right) \left(\sqrt{\int K^2(x) dx} + 1\right) + O\left(\frac{\ln^2 n}{\sqrt{n}}\right) \left(\sqrt{\int b_K^2(x) dx} + \frac{\|f\|_2}{\sqrt{n}}\right) \end{aligned}$$

$$\begin{aligned}
&= O\left(\frac{\ln^{5/2}n}{\sqrt{n}}\right) \left(\sqrt{MISE(K)} + \frac{1 + \|f\|_2}{\sqrt{n}}\right) \\
&= O\left(\frac{\ln^{5/2}n}{\sqrt{n}}\right) \sqrt{MISE(K)}
\end{aligned}$$

which concludes the proof of proposition **A1** and

$$P(A_n^c) \leq P(A_{n1}^c) + P(A_{n2}^c) = O(n^{-\lambda'}) \quad \text{for an appropriate } \lambda' < \infty$$

which is proposition **A3**.

## 5 Practical computation

Once the statistical properties of the CV-optimal kernel function  $K_0$  have been examined, we would like to actually compute this kernel from a sample  $X_1, \dots, X_n$ .  $K_0$  is  $\arg \min CV(K)$  within the set  $\mathcal{K}_n := \{K \in \mathcal{K} | \text{supp } \widehat{K} \subseteq (-n, n)\}$ . Hence we face a minimization problem.

Note that the set  $\mathcal{K}$  is convex. With respect to the properties  $\widehat{K}$  real and non-negative and  $\|K\|_2^2 \leq n$ , convexity is obvious. Given that all  $\widehat{K}$  in  $\mathcal{K}$  are unimodal and symmetric around 0, their mode is 0. And a convex combination of any two  $\widehat{K}$  is again unimodal. Convexity is also preserved through the trimming of the support of  $\widehat{K}$ . On the other hand,  $CV(K)$  is a strictly convex function. Therefore  $\min CV(K)$  over  $\mathcal{K}_n$  is a convex optimization problem, where the argument is itself a non-increasing function,  $\widehat{K}: [0, n] \rightarrow [0, 1]$ .

Convex problems have a unique solution, so we are theoretically save. The question is of course to find the solution. Consider a discrete version of  $\mathcal{K}_n$ , say  $\mathcal{K}_n^t$ , which contains all real, symmetric and unimodal piecewise constant functions on  $[0, n]$ , with jumps at the points  $2^{-t}k$ ,  $k = 1 \dots 2^t n$ , and values  $\in [0, 1]$ . The minimization of  $CV(K)$  over  $\mathcal{K}_n^t$  is still a convex optimization problem, but this time with respect to a parameter of dimension  $2^t n$  (number of variables) and with  $2^t n + 2$  constraints (unimodality, positivity and  $L_2$ -norm).

The  $L_1$ -distance between a kernel function  $\widehat{K}$  in  $\mathcal{K}_n$  and its closest neighbor  $\widehat{K}^t$  in  $\mathcal{K}_n^t$  is not greater than  $2^{-t}$  (and thus the same applies for the supremum distance between  $K$  and  $K^t$ ). It follows with little effort that  $|CV(K) - CV(K^t)| \leq \frac{2}{\pi} \cdot 2^{-t}$  and therewith

$$CV(K_0^t) := \min_{\mathcal{K}_n^t} CV(K) \leq CV((K_0)^t) \leq CV(K_0) + \frac{2}{\pi} \cdot 2^{-t}$$

Since  $\mathcal{K}_n^t \subseteq \mathcal{K}_n^{t+1}$ , the sequence  $\{K_0^t\}_{t \in \mathbb{N}}$  converges towards  $K_0$ , the unique solution of the original problem.

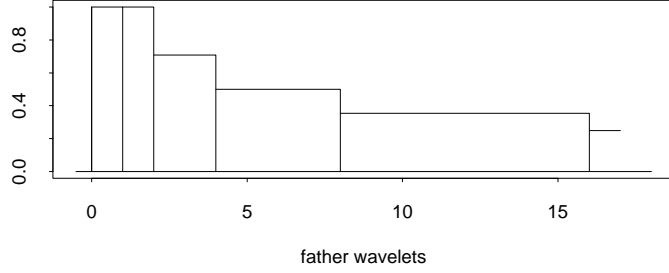
There is no doubt that a profound analysis in terms of optimization would yield a more sophisticated algorithm to solve to the problem, possibly avoiding discretization and giving convergence rates over classes of densities.

## A Appendix

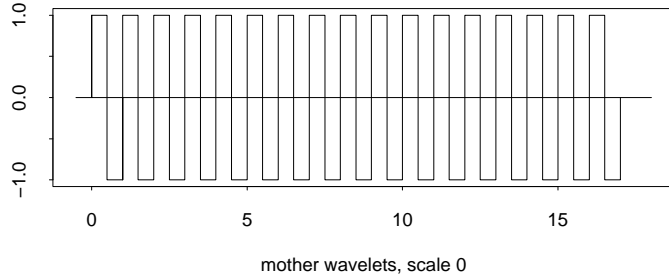
**A.1 Wavelet decomposition of the kernel** As the class  $\mathcal{K}_n$  itself, also the desired basis is constructed in the Fourier domain. We are searching for a way to compress most economically the information inherent to  $\widehat{K}$ . To this end, we utilize  $\widehat{K}$ 's assumed monotony on  $\mathbb{R}^+$ , which gives that for  $\|K\|_2$  fixed,  $\widehat{K}(\omega) \leq \|K\|_2 |2\omega|^{-1/2}$  must hold. Heuristically spoken, the further

out we reach on the line  $\mathbb{R}^+$ , the smaller will be the variation in  $\widehat{K}$ . But that means, we can allow for a rougher approximation without losing much of our approximating power.

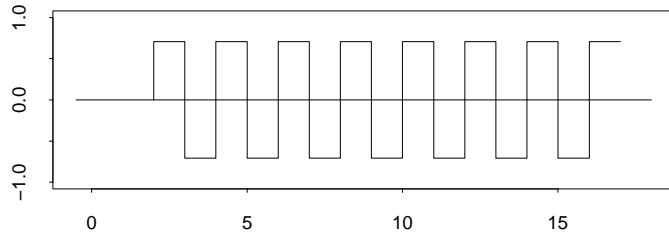
Technically we implement the idea as follows: Inspired by the well known Haar basis, symmetric father wavelets are defined on the interval  $[-n, n]$ :  $\widehat{\varphi}_{01}(\omega) := 2^{-1/2}I(|\omega| \in [0, 1))$ ,  $\widehat{\varphi}_{02}(\omega) := 2^{-1/2}I(|\omega| \in [1, 2))$ . After that, we let the supports of the wavelets grow: with negative scale index, we define  $\widehat{\varphi}_{-s,2}(\omega) := 2^{-(s+1)/2}I(|\omega| \in [2^s, 2^{s+1}))$ ,  $1 \leq s \leq d_n$ , where  $d_n \sim \ln n$ .



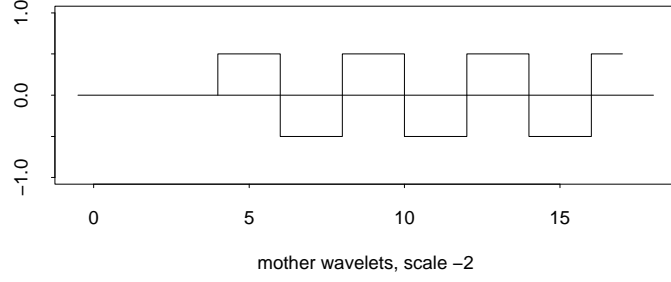
The sequence of father wavelets  $(\widehat{\varphi}_{01}, \widehat{\varphi}_{02}, \widehat{\varphi}_{-1,2}, \widehat{\varphi}_{-2,2}, \dots, \widehat{\varphi}_{-d_n,2})$  covers the whole interval  $[-n, n]$ , (the support of a function being defined as the closure of the set, where it is nonzero) and comprises  $d_n + 2$  elements. On the supporting interval of each father wavelet, the mother wavelets are defined on refining scales. With notation  $I_{ut}(\omega) := I(|\omega| \in [2^{-u}(t-1), 2^{-u}t))$ , the mother wavelets on  $(-2^{s+1}, -2^s] \cup [2^s, 2^{s+1})$  are  $\widehat{\psi}_{u,t}(\omega) := 2^{(u-1)/2}[I_{u+1,2t-1}(\omega) - I_{u+1,2t}(\omega)]$ ,  $u = -s, -s+1, \dots, 0, 1, 2, \dots$  and  $t = 2^{s+u} + 1, \dots, 2^{s+u+1}$ . When we combine all mother wavelets with the same scale index  $s$ , we arrive at a sequence of  $(\widehat{\psi}_{s,2}, \dots, \widehat{\psi}_{s,2^{s+1}})$  for  $s = -1, \dots, -d_n$ , and  $(\widehat{\psi}_{s,1}, \dots, \widehat{\psi}_{s,2^s})$ , for  $s \geq 0$ . We observe that for  $s < 0$ , the corresponding mother wavelets do not cover the whole interval  $[-n, n]$ , but only  $[-n, -2^s] \cup [2^s, n]$ .



mother wavelets, scale 0



mother wavelets, scale -1



Unifying the notation:

$$\begin{aligned} I_{st}(\omega) &:= I(|\omega| \in [2^{-s}(t-1), 2^{-s}t]) \\ \widehat{\varphi}_{st}(\omega) &:= 2^{(s-1)/2} I_{st}(\omega) \\ \widehat{\psi}_{st}(\omega) &:= 2^{(s-1)/2} [I_{s+1,2t-1}(\omega) - I_{s+1,2t}(\omega)], \end{aligned}$$

we have the following complete orthonormal function basis of  $L_2((-n, n))$ :  $\{\widehat{\varphi}_{01}\} \cup \{\widehat{\varphi}_{s2} | s = 0, \dots, -d_n\} \cup \{\widehat{\psi}_{st} | s = -1, \dots, -d_n \text{ and } t = 2, \dots, 2^s n\} \cup \{\widehat{\psi}_{st} | s \geq 0 \text{ and } t = 1, \dots, 2^s n\}$ . The decomposition of  $\widehat{K}$  results in:

$$\widehat{K}(\omega) = \alpha_{01}(K)\widehat{\varphi}_{01}(\omega) + \sum_{s=0}^{-d_n} \alpha_{s2}(K)\widehat{\varphi}_{s2}(\omega) + \sum_{s=-1}^{-d_n} \sum_{t=2}^{2^s n} \beta_{st}(K)\widehat{\psi}_{st}(\omega) + \sum_{s=0}^{\infty} \sum_{t=1}^{2^s n} \beta_{st}(K)\widehat{\psi}_{st}(\omega)$$

$$\alpha_{st}(K) := \int \widehat{\varphi}_{st}(\omega)\widehat{K}(\omega)d\omega \quad \text{and} \quad \beta_{st}(K) := \int \widehat{\psi}_{st}(\omega)\widehat{K}(\omega)d\omega$$

( $K$  and the wavelets are both symmetric, so conjugation can be dropped.) By an inverse Fourier transferred, the additive decomposition of  $\widehat{K}$  can be transformed to the space domain.

$$K(x) = \alpha_{01}(K)\varphi_{01}(x) + \sum_{s=0}^{-d_n} \alpha_{s2}(K)\varphi_{s2}(x) + \sum_{s=-1}^{-d_n} \sum_{t=2}^{2^s n} \beta_{st}(K)\psi_{st}(x) + \sum_{s=0}^{\infty} \sum_{t=1}^{2^s n} \beta_{st}(K)\psi_{st}(x)$$

Accordingly, the summands in the U-process decompose into:

$$\begin{aligned} U_K(X_i, X_j) &= \alpha_{01}(K)U_{\varphi_{01}}(X_i, X_j) + \sum_{s=0}^{-d_n} \alpha_{s2}(K)U_{\varphi_{s2}}(X_i, X_j) + \sum_{s=-1}^{-d_n} \sum_{t=2}^{2^s n} \beta_{st}(K)U_{\psi_{st}}(X_i, X_j) \\ &\quad + \sum_{s=0}^{\infty} \sum_{t=1}^{2^s n} \beta_{st}(K)U_{\psi_{st}}(X_i, X_j), \end{aligned}$$

where  $U_{\varphi_{st}}(X_i, X_j) := \varphi_{st}(X_i - X_j) - E[\varphi_{st}(X_i - X_j)|X_i] - E[\varphi_{st}(X_i - X_j)|X_j] + E[\varphi_{st}(X_i - X_j)]$ , and  $U_{\psi_{st}}$  equally defined for  $\psi_{st}$ . Interchanging the order of summation, we obtain that:

$$\begin{aligned} \frac{1}{n(n-1)} \sum_{i \neq j} U_K(X_i, X_j) &= \alpha_{01}(K) \left[ \frac{1}{n(n-1)} \sum_{i \neq j} U_{\varphi_{01}}(X_i, X_j) \right] \\ &\quad + \sum_{s=0}^{-d_n} \alpha_{s2}(K) \left[ \frac{1}{n(n-1)} \sum_{i \neq j} U_{\varphi_{s2}}(X_i, X_j) \right] \end{aligned}$$

$$\begin{aligned}
& + \sum_{s=-1}^{-d_n} \sum_{t=2}^{2^s n} \beta_{st}(K) \left[ \frac{1}{n(n-1)} \sum_{i \neq j} U_{\psi_{st}}(X_i, X_j) \right] \\
& + \sum_{s=0}^{\infty} \sum_{t=1}^{2^s n} \beta_{st}(K) \left[ \frac{1}{n(n-1)} \sum_{i \neq j} U_{\psi_{st}}(X_i, X_j) \right] \quad (9)
\end{aligned}$$

From this point onwards, the sums of wavelet coefficients and the U-statistics can be handled separately. The  $\alpha$ 's and  $\beta$ 's are deterministic and we show in Lemma 1:

$$\begin{aligned}
|\alpha_{01}(K)| + \sum_{s=0}^{-d_n} |\alpha_{s2}(K)| &\leq \sqrt{d_n + 2} \sqrt{2\pi \int K^2(x) dx} \\
\sum_{t=2}^{2^s n} |\beta_{st}(K)| &\leq \sqrt{2\pi \int K^2(x) dx} \quad \text{for } s < 0 \\
\sum_{t=1}^{2^s n} |\beta_{st}(K)| &\leq 2^{(-s+1)/2} \quad \text{for } s \geq 0
\end{aligned}$$

For a suitable constant  $\lambda < \infty$ , we choose our set of ‘‘favorable events’’ as:

$$\begin{aligned}
A_{n1} := \left\{ (X_1, \dots, X_n) : \right. & \left| \frac{1}{n(n-1)} \sum_{i \neq j} U_{\varphi_{st}}(X_i, X_j) \right| \leq \frac{\lambda \ln^{3/2} n}{n}, \quad (s, t) = (-d_n, 2), \dots, (0, 2), (0, 1); \\
& \left| \frac{1}{n(n-1)} \sum_{i \neq j} U_{\psi_{st}}(X_i, X_j) \right| \leq \frac{\lambda \ln^{3/2} n}{n}, \quad s = -d_n, \dots, -1, t = 2, \dots, 2^s n; \\
& \left. \left| \frac{1}{n(n-1)} \sum_{i \neq j} U_{\psi_{st}}(X_i, X_j) \right| \leq \frac{\lambda \ln n + s}{n}, \quad s \geq 0, t = 1, \dots, 2^s n \right\}
\end{aligned}$$

whereupon the U-statistics do not become excessively large. The fact that the complement of the set  $A_{n1}$  has probability tending to 0, as  $n \rightarrow \infty$ ,  $P(A_{n1}^c) = O(n^{-\lambda^{2/3} + 1})$  (uniformly for  $f \in \mathcal{S}_\beta(L)$  with  $\beta > 1/2$ ), will be shown in Lemma 2, equation (10). On  $A_{n1}$  it holds that (in connection with (9)):

$$\begin{aligned}
\left| \frac{1}{n(n-1)} \sum_{i \neq j} U_K(X_i, X_j) \right| &\leq \frac{\lambda \ln^{3/2} n}{n} \left[ |\alpha_{01}(K)| + \sum_{s=0}^{-d_n} |\alpha_{s2}(K)| + \sum_{s=-1}^{-d_n} \sum_{t=2}^{2^s n} |\beta_{st}(K)| \right] \\
& + \sum_{s=0}^{\infty} \frac{\lambda \ln n + s}{n} \sum_{t=1}^{2^s n} |\beta_{st}(K)| \\
&\leq \frac{\lambda \ln^{3/2} n}{n} \left[ \sqrt{d_n + 2} \sqrt{2\pi \int K^2(x) dx} + d_n \sqrt{2\pi \int K^2(x) dx} \right] \\
& + \sum_{s=0}^{\infty} \frac{\lambda \ln n + s}{n} 2^{(-s+1)/2} \\
& = O\left(\frac{\ln^{5/2} n}{n}\right) \left( \sqrt{\int K^2(x) dx} + 1 \right)
\end{aligned}$$

which completes (7). But two assertions are still left to be verified.

**Lemma 1** For the father and mother wavelet coefficients of  $K$  defined so far, it holds that

$$\begin{aligned} |\alpha_{01}(K)| + \sum_{s=0}^{d_n} |\alpha_{s2}(K)| &\leq \sqrt{d_n + 2} \sqrt{2\pi \int K^2(x) dx} \\ \sum_{t=2}^{2^s n} |\beta_{st}(K)| &\leq \sqrt{2\pi \int K^2(x) dx} \quad \text{for } s < 0 \\ \sum_{t=1}^{2^s n} |\beta_{st}(K)| &\leq 2^{(-s+1)/2} \quad \text{for } s \geq 0 \end{aligned}$$

**Proof** Since  $\widehat{\varphi}_{st}$  are orthonormal, we can deduce through Cauchy-Schwartz:

$$\begin{aligned} |\alpha_{01}(K)| + \sum_{s=0}^{d_n} |\alpha_{s2}(K)| &= \int \widehat{K}(\omega) \left[ \widehat{\varphi}_{01}(\omega) + \sum_{s=0}^{d_n} \widehat{\varphi}_{s2}(\omega) \right] d\omega \\ &\leq \sqrt{\int \widehat{K}^2(\omega) d\omega} \sqrt{d_n + 2} \end{aligned}$$

May  $\text{TV}(\widehat{K}|\text{supp } I_{st})$  denote the total variation of  $\widehat{K}$  on the support of  $I_{st}$ , and the like for max and min. It is known that for mother wavelet coefficients, it holds that:

$$\sum_t |\beta_{st}(K)| = 2^{-(s+1)/2} \text{TV} \left( \widehat{K} \Big| \bigcup_t \text{supp } I_{st} \right)$$

For  $s \geq 0$ , the supports of the mother wavelets cover the whole interval  $[-n, n]$ , and we obtain  $\text{TV}(\widehat{K}|\bigcup \text{supp } I_{st}) = \text{TV}(\widehat{K}) \leq 2$ , due to unimodality of  $\widehat{K}$  ( $\widehat{K}(0) = \int K(x) dx = 1$ ). For  $s < 0$ ,  $\bigcup \text{supp } I_{st} = [-n, -2^s] \cup [2^s, n]$ , and we use that

$$\int \widehat{K}^2(\omega) d\omega = \int_0^{\omega_0} \widehat{K}^2(\omega) d\omega + \int_{\omega_0}^n \widehat{K}^2(\omega) d\omega \leq \int_0^{\omega_0} \widehat{K}^2(\omega_0) d\omega = \omega_0 \widehat{K}^2(\omega_0)$$

This yields

$$\begin{aligned} \sum_t |\beta_{st}(K)| &\leq 2^{-(s+1)/2} \text{TV} \left( \widehat{K} \Big| \bigcup_t \text{supp } I_{st} \right) \\ &\leq 2^{-(s+1)/2} 2 \widehat{K}(2^{-s}) \\ &\leq 2^{-(s-1)/2} \frac{\sqrt{\int \widehat{K}^2(\omega) d\omega}}{\sqrt{2 \cdot 2^{-s}}} \\ &= \sqrt{\int \widehat{K}^2(\omega) d\omega} \quad \square \end{aligned}$$

**Lemma 2** For the father and mother wavelets defined above and arbitrary  $0 < \lambda < \infty$  it holds that

$$P \left( \frac{1}{n(n-1)} \left| \sum_{i \neq j} U_{\varphi_{st}}(X_i, X_j) \right| > \frac{\lambda \ln^{3/2} n}{n} \right) = O \left( n^{-\lambda^{2/3}} \right) \text{ for } (s, t) = (-d_n, 2), \dots, (0, 2)$$

$$\begin{aligned}
& \text{and } (0, 1) \\
P\left(\frac{1}{n(n-1)}\left|\sum_{i \neq j} U_{\psi_{st}}(X_i, X_j)\right| > \frac{\lambda \ln^{3/2} n}{n}\right) &= O\left(n^{-\lambda^{2/3}}\right) \text{ for } s = -1, \dots, -d_n \\
& \text{and } t = 2, \dots, 2^s n \\
P\left(\frac{1}{n(n-1)}\left|\sum_{i \neq j} U_{\psi_{st}}(X_i, X_j)\right| > \frac{\lambda \ln n + s}{n}\right) &= O\left(n^{-\lambda^{2/3}} e^{-s}\right) \text{ for } s = 0, 1, 2, \dots \\
& \text{and } t = 1, \dots, 2^s n
\end{aligned}$$

These bounds  $O(\cdot)$  are uniform in  $s$  and  $t$ .

$A_{n1}^c$  is the union of all complementary sets and the approximations of Lemma 2 give

$$\begin{aligned}
P(A_{n1}^c) &= O(n^{-\lambda^{2/3}}) \left[ (d_n + 2) + \sum_{s=-1}^{-d_n} \sum_{t=2}^{2^s n} 1 \right] + \sum_{s=0}^{\infty} \sum_{t=1}^{2^s n} O\left(n^{-\lambda^{2/3}} e^{-s}\right) \\
&= O(n^{-\lambda^{2/3}}) O(\ln n + n) + O(n^{-\lambda^{2/3}}) O(n)
\end{aligned} \tag{10}$$

**Remark** As we will see in the proof, the bounds of Lemma 2 are uniform in function sets with bounded  $\max f$ . This is the case for Sobolev classes  $\mathcal{S}_\beta(L)$  with  $\beta > 1/2$ . So over Sobolev classes, Lemma 2 holds uniformly.

**Proof** From the Bernstein type inequality for degenerate U-statistics, shown by Arcones, Giné (1993), it follows that for all  $\varphi_{st}$ , and analogously for all  $\psi_{st}$  with  $s < 0$ , there exist constants  $c_1$  and  $c_2$  independent from  $\varphi_{st}$  (and from  $\psi_{st}$  respectively), such that:

$$\begin{aligned}
& P\left(\frac{1}{n(n-1)}\left|\sum_{i \neq j} U_{\varphi_{st}}(X_i, X_j)\right| > \frac{\lambda \ln^{3/2} n}{n}\right) \\
& \leq c_1 \exp\left\{-\frac{c_2(n-1) \frac{\lambda \ln^{3/2} n}{n}}{\sqrt{E|U_{\varphi_{st}}|^2} + \left(\frac{n-1}{n} \|\varphi_{st}\|_\infty^2 \frac{\lambda \ln^{3/2} n}{n}\right)^{1/3}}\right\} \\
& \leq c_1 \exp\left\{-\frac{c_2(n-1) \frac{\lambda \ln^{3/2} n}{n}}{\frac{1}{\sqrt{2\pi}} \|f\|_\infty^{1/2} \|\widehat{\varphi}_{st}\|_2 + \left(\frac{n-1}{n} \frac{1}{(2\pi)^2} \|\widehat{\varphi}_{st}\|_1^2 \frac{\lambda \ln^{3/2} n}{n}\right)^{1/3}}\right\} \\
& = O\left(\exp\left\{-\frac{\lambda^{2/3} \ln n}{1 + \|f\|_\infty^{1/2} \ln^{-1/2} n}\right\}\right)
\end{aligned}$$

which is an  $O(n^{-\lambda^{2/3}})$ , not depending on  $s$  and  $t$ . By analogue calculations, we get for  $\psi_{st}$  with  $s \geq 0$ :

$$P\left(\frac{1}{n(n-1)}\left|\sum_{i \neq j} U_{\psi_{st}}(X_i, X_j)\right| > \frac{\lambda \ln n + s}{n}\right) = O\left(\exp\left\{-\frac{\lambda^{2/3} \ln n + \lambda^{-1/3} s}{\|f\|_\infty^{1/2}}\right\}\right)$$

an  $O(n^{-\lambda^{2/3}} e^{-s})$ , uniform in  $t$ . □

**A.2 Wavelet decomposition of the bias** We are now going to apply an additive decomposition to the bias term in the difference  $\widehat{CV}(K) - ISE(K)$ :

$$\frac{1}{n} \sum_{j=1}^n \left( b_K(X_j) + h_f(X_j) \right) - E \left[ b_K(X_j) + h_f(X_j) \right]$$

where  $b_K(x) = f * K(x) - f(x)$ ,  $h_f$  is the high-frequency component of  $f$  (definition (5)) and  $\widehat{b}_K + \widehat{h}_f = \widehat{b}_K \cdot I_n$ . In the bias, everything relates to the underlying density, so we construct basis functions depending on  $f$ . Let us define the integral of  $|\widehat{f}|^2$  over  $[-\omega, \omega]$  as a function  $F(\omega)$ .

$$F(\omega) := \int_{-\omega}^{\omega} |\widehat{f}(\tau)|^2 d\tau$$

This map transforms the  $\omega$ -halfaxis  $[0, \infty)$  by mapping  $\omega \mapsto F(\omega)$  to the interval  $[0, \|\widehat{f}\|_2^2)$ .

$$F(0) = 0, \quad F_n := F(n) = \int_{-n}^n |\widehat{f}(\tau)|^2 d\tau, \quad \lim_{\omega \rightarrow \infty} F(\omega) = \|\widehat{f}\|_2^2$$

The initial value of an interval, say  $[2^{-s}(t-1)F_n, 2^{-s}tF_n)$  with length  $2^{-s}$ , on this axis is the interval  $[F^{-1}(2^{-s}(t-1)F_n), F^{-1}(2^{-s}tF_n))$  on the original axis. The integral of  $|\widehat{f}|^2$  over the initial interval is obviously  $\frac{1}{2} 2^{-s} F_n$ .

$$2^{-s} F_n = \int_{-F^{-1}(2^{-s}tF_n)}^{F^{-1}(2^{-s}tF_n)} |\widehat{f}(\omega)|^2 d\omega - \int_{-F^{-1}(2^{-s}(t-1)F_n)}^{F^{-1}(2^{-s}(t-1)F_n)} |\widehat{f}(\omega)|^2 d\omega$$

Define the indicator functions:

$$I'_{st}(\omega) := I \left( |\omega| \in \left[ F^{-1}(2^{-s}(t-1)F_n), F^{-1}(2^{-s}tF_n) \right) \right)$$

satisfying  $\int |\widehat{f}(\omega)|^2 I'_{st}(\omega) d\omega = 2^{-s} F_n$ , and the orthonormal wavelet functions:

$$\begin{aligned} \widehat{\varphi}'_{st}(\omega) &:= 2^{s/2} F_n^{-1/2} \widehat{f}(\omega) I'_{st}(\omega), & \text{for } s = 1, \dots, s_n \text{ with } t = 2^s - 1 \text{ and} \\ & & s = s_n, t = 2^{s_n} \text{ where } s_n \sim \ln n \\ \widehat{\psi}'_{st}(\omega) &:= 2^{s/2} F_n^{-1/2} \widehat{f}(\omega) [I'_{s+1, 2t-1}(\omega) - I'_{s+1, 2t}(\omega)], & \text{for } s = 1, \dots, s_n - 1 \text{ with} \\ & & t = 1, \dots, 2^s - 1 \text{ and } s = s_n, s_n + 1, \dots \text{ with} \\ & & t = 1, \dots, 2^s \end{aligned}$$

$\{\widehat{\varphi}'_{st} | s = 1, \dots, s_n, t = 2^s - 1\} \cup \{\widehat{\varphi}'_{s_n, 2^{s_n}}\} \cup \{\widehat{\psi}'_{st} | s = 1, \dots, s_n - 1, t = 1, \dots, 2^s - 1\} \cup \{\widehat{\psi}'_{st} | s \geq s_n, t = 1, \dots, 2^s\}$  represent a complete orthonormal basis for the set of all functions  $\{\widehat{f} \cdot \widehat{g} \cdot I_n | \widehat{g} \in L_2\}$ , which the bias functions  $\widehat{b}_K \cdot I_n$  belong to for all  $K \in \mathcal{K}$ . After the inverse Fourier transform, we have

$$\begin{aligned} b_K(x) + h_f(x) &= \sum_{s=1}^{s_n} \alpha'_{s2^s-1}(b_K) \varphi'_{s2^s-1}(x) + \alpha'_{s_n 2^{s_n}}(b_K) \varphi'_{s_n 2^{s_n}}(x) + \sum_{s=1}^{s_n-1} \sum_{t=1}^{2^s-1} \beta'_{st}(b_K) \psi'_{st}(x) \\ &+ \sum_{s=s_n}^{\infty} \sum_{t=1}^{2^s} \beta'_{st}(b_K) \psi'_{st}(x) \end{aligned}$$



which gives in turn

$$\begin{aligned}
& \frac{1}{n} \sum_{j=1}^n \left( b_K(X_j) + h_f(X_j) \right) - E \left[ b_K(X_j) + h_f(X_j) \right] \\
&= \sum_{s=1}^{s_n} \alpha'_{s2^{s-1}}(b_K) \left[ \frac{1}{n} \sum_{j=1}^n \varphi'_{s2^{s-1}}(X_j) - E \varphi'_{s2^{s-1}}(X_j) \right] + \alpha'_{s_n 2^{s_n}}(b_K) \\
&\quad \times \left[ \frac{1}{n} \sum_{j=1}^n \varphi'_{s_n 2^{s_n}}(X_j) - E \varphi'_{s_n 2^{s_n}}(X_j) \right] \\
&\quad + \sum_{s=1}^{s_n-1} \sum_{t=1}^{2^s-1} \beta'_{st}(b_K) \left[ \frac{1}{n} \sum_{j=1}^n \psi'_{st}(X_j) - E \psi'_{st}(X_j) \right] + \sum_{s=s_n}^{\infty} \sum_{t=1}^{2^s} \beta'_{st}(b_K) \\
&\quad \times \left[ \frac{1}{n} \sum_{j=1}^n \psi'_{st}(X_j) - E \psi'_{st}(X_j) \right] \tag{11}
\end{aligned}$$

Again, we will proceed separately with the aim of finding bounds to the deterministic wavelet coefficients and the stochastic processes. Lemma 3 shows that

$$\begin{aligned}
\sum_{s=1}^{s_n} |\alpha'_{s2^{s-1}}(b_K)| + |\alpha'_{s_n 2^{s_n}}(b_K)| &\leq \sqrt{s_n + 1} \sqrt{2\pi \int b_K^2(x) dx} \\
\sum_{t=1}^{2^s-1} |\beta'_{st}(b_K)| &\leq 2 \sqrt{2\pi \int b_K^2(x) dx} \quad \text{for } s < s_n \\
\sum_{t=1}^{2^s} |\beta'_{st}(b_K)| &\leq 2 \cdot 2^{-s/2} \|f\|_2 \quad \text{for } s \geq s_n
\end{aligned}$$

Over a set of “favorable events”, whose complement has an asymptotically decreasing probability (Lemma 4, inequality (12)), the partial sum processes can be controlled. For  $\lambda < \infty$

$$\begin{aligned}
A_{n2} := \left\{ (X_1, \dots, X_n) : \frac{1}{n} \left| \sum_{j=1}^n \varphi'_{st}(X_j) - E \varphi'_{st}(X_j) \right| \leq \frac{\lambda \ln n}{\sqrt{n}}, \quad (s, t) = (1, 1), \dots, (s_n, 2^{s_n} - 1), (s_n, 2^{s_n}); \right. \\
\left. \frac{1}{n} \left| \sum_{j=1}^n \psi'_{st}(X_j) - E \psi'_{st}(X_j) \right| \leq \frac{\lambda \ln n}{\sqrt{n}}, \quad s = 1, \dots, s_n - 1, t = 1, \dots, 2^s - 1; \right. \\
\left. \frac{1}{n} \left| \sum_{j=1}^n \psi'_{st}(X_j) - E \psi'_{st}(X_j) \right| \leq \frac{\lambda \ln n + s}{\sqrt{n}}, \quad s \geq s_n, t = 1, \dots, 2^s \right\}
\end{aligned}$$

$$\text{and } P(A_{n2}^c) = O(n^{-\lambda+1})$$

Following (11) and taking into account that  $2^{s_n} \leq n^{-1}$ , it holds on  $A_{n2}$ :

$$\begin{aligned}
& \frac{1}{n} \left| \sum_{j=1}^n \left( b_K(X_j) + h_f(X_j) \right) - E \left[ b_K(X_j) + h_f(X_j) \right] \right| \\
&\leq \frac{\lambda \ln n}{\sqrt{n}} \left[ \sum_{s=1}^{s_n} |\alpha'_{s2^{s-1}}(b_K)| + |\alpha'_{s_n 2^{s_n}}(b_K)| + \sum_{s=1}^{s_n-1} \sum_{t=1}^{2^s-1} |\beta'_{st}(b_K)| \right] + \sum_{s=s_n}^{\infty} \frac{\lambda \ln n + s}{\sqrt{n}} \sum_{t=1}^{2^s} |\beta'_{st}(b_K)| \\
&= O\left(\frac{\ln^2 n}{\sqrt{n}}\right) \left( \sqrt{\int b_K^2(x) dx} + \frac{\|f\|_2}{\sqrt{n}} \right)
\end{aligned}$$

which completes (8). Now we proof the remaining assertions.

**Lemma 3** The coefficients of the bias defined through the  $f$ -depending function basis satisfy

$$\begin{aligned} \sum_{s=1}^{s_n} |\alpha'_{s2^{s-1}}(b_K)| + |\alpha'_{s_n 2^{s_n}}(b_K)| &\leq \sqrt{s_n + 1} \sqrt{2\pi \int b_K^2(x) dx} \\ \sum_{t=1}^{2^s-1} |\beta'_{st}(b_K)| &\leq 2\sqrt{2\pi \int b_K^2(x) dx} \quad \text{for } s < s_n \\ \sum_{t=1}^{2^s} |\beta'_{st}(b_K)| &\leq 2 \cdot 2^{-s/2} \|f\|_2 \quad \text{for } s \geq s_n \end{aligned}$$

**Proof** The father wavelet coefficients are bounded in the same way as in Lemma 1, such that

$$\sum_{s=1}^{s_n} |\alpha'_{s2^{s-1}}(b_K)| + |\alpha'_{s_n 2^{s_n}}(b_K)| = \sqrt{\int |\widehat{b}_K(\omega)|^2 d\omega} \sqrt{s_n + 1}$$

For every  $t$  in the summation range, choose an arbitrary  $\omega_{st} \in [F^{-1}(2^{-s}(t-1)F_n), F^{-1}(2^{-s}tF_n))$ . Again let  $\text{TV}(\widehat{K}|\text{supp } I'_{st})$  be the total variation of  $\widehat{K}$  over the support of  $I'_{st}$ .

$$\begin{aligned} \sum_t |\beta'_{st}(b_K)| &= \sum_t \left| \int \widehat{b}_K(\omega) \overline{\psi'_{st}} d\omega \right| \\ &= \sum_t 2^{s/2} F_n^{-1/2} \left| \int \widehat{f}(\omega) (1 - \widehat{K}(\omega)) \overline{\widehat{f}}(\omega) [I'_{s+1,2t-1}(\omega) - I'_{s+1,2t}(\omega)] d\omega \right| \\ &= \sum_t 2^{s/2} F_n^{-1/2} \left| \int |\widehat{f}(\omega)|^2 (1 - \widehat{K}(\omega_{st})) [I'_{s+1,2t-1}(\omega) - I'_{s+1,2t}(\omega)] d\omega \right. \\ &\quad \left. + \int |\widehat{f}(\omega)|^2 (\widehat{K}(\omega_{st}) - \widehat{K}(\omega)) [I'_{s+1,2t-1}(\omega) - I'_{s+1,2t}(\omega)] d\omega \right| \\ &\leq \sum_t 2^{s/2} F_n^{-1/2} \left[ (1 - \widehat{K}(\omega_{st})) \int |\widehat{f}(\omega)|^2 [I'_{s+1,2t-1}(\omega) - I'_{s+1,2t}(\omega)] d\omega \right. \\ &\quad \left. + \int |\widehat{f}(\omega)|^2 \text{TV}(\widehat{K}|\text{supp } I'_{s+1,2t-1}) I'_{s+1,2t-1}(\omega) d\omega \right. \\ &\quad \left. + \int |\widehat{f}(\omega)|^2 \text{TV}(\widehat{K}|\text{supp } I'_{s+1,2t}(|\omega|)) I'_{s+1,2t}(\omega) d\omega \right] \\ &= \sum_t 2^{s/2} F_n^{-1/2} \left[ 0 + \text{TV}(\widehat{K}|\text{supp } I'_{st}) \int |\widehat{f}(\omega)|^2 I'_{st}(\omega) d\omega \right] \\ &= \sum_t 2^{s/2} F_n^{-1/2} \text{TV}(\widehat{K}|\text{supp } I'_{st}) F_n 2^{-s} \\ &= 2^{-s/2} F_n^{1/2} \text{TV}\left(\widehat{K}|\text{supp } \bigcup_t I'_{st}\right) \end{aligned}$$

The mother wavelets on the scales  $s \geq s_n$  are defined over the whole interval  $[-n, n]$ , therefore  $\text{TV}(\widehat{K}|\text{supp } \bigcup_t I'_{st}) = \text{TV}(\widehat{K}) \leq 2$ . For  $s < s_n$ , the mother wavelets are supported on

$[-F^{-1}((1-2^s)F_n), F^{-1}((1-2^s)F_n)]$ . On this interval, the total variation amounts to at most  $2[1 - \widehat{K}(F^{-1}((1-2^s)F_n))]$ .

$$\begin{aligned}
\sum_t |\beta'_{st}(b_K)| &\leq 2^{-s/2} F_n^{1/2} 2 \left[ 1 - \widehat{K}(F^{-1}((1-2^s)F_n)) \right] \\
&= 2^{-s/2} F_n^{1/2} \left( \int |\widehat{\varphi}'_{s2^s}(\omega)|^2 d\omega \right) 2 \left[ 1 - \widehat{K}(F^{-1}((1-2^s)F_n)) \right] \\
&= 2^{s/2} F_n^{-1/2} \int |\widehat{f}(\omega)|^2 I'_{s2^s}(\omega) d\omega 2 \left[ 1 - \widehat{K}(F^{-1}(2^{-s}(2^s-1)F_n)) \right] \\
&\leq 2^{s/2} F_n^{-1/2} 2 \int \widehat{f}(\omega) \left[ 1 - \widehat{K}(\omega) \right] \overline{\widehat{f}}(\omega) I'_{s2^s}(\omega) d\omega \\
&= 2 \int \widehat{b}_K(\omega) \overline{\widehat{\varphi}'_{s2^s}(\omega)} d\omega \\
&\leq 2 \sqrt{\int |\widehat{b}_K(\omega)|^2 d\omega} \sqrt{\int |\widehat{\varphi}'_{s2^s}(\omega)|^2 d\omega} \\
&= 2 \sqrt{\int |\widehat{b}_K(\omega)|^2 d\omega} \quad \square
\end{aligned}$$

**Lemma 4** For any  $\lambda < \infty$ , the following inequalities hold uniformly for all indicated  $s$  and  $t$  and, exactly as in Lemma 2, as well uniformly for  $f \in \mathcal{S}_\beta(L)$ ,  $\beta > 1/2$ :

$$\begin{aligned}
P \left( \frac{1}{n} \left| \sum_{j=1}^n \varphi'_{st}(X_j) - E\varphi'_{st}(X_j) \right| > \frac{\lambda \ln n}{\sqrt{n}} \right) &= O(n^{-\lambda}), \quad (s, t) = (1, 1), \dots, (s_n, 2^{s_n} - 1), \\
&\text{and } (s_n, 2^{s_n}) \\
P \left( \frac{1}{n} \left| \sum_{j=1}^n \psi'_{st}(X_j) - E\psi'_{st}(X_j) \right| > \frac{\lambda \ln n}{\sqrt{n}} \right) &= O(n^{-\lambda}), \quad s = 1, \dots, s_{n-1} \text{ and} \\
&t = 1, \dots, 2^s - 1 \\
P \left( \frac{1}{n} \left| \sum_{j=1}^n \psi'_{st}(X_j) - E\psi'_{st}(X_j) \right| > \frac{\lambda \ln n + s}{\sqrt{n}} \right) &= O(n^{-\lambda} e^{-s}), \quad s = s_n, s_n + 1, \dots \text{ and} \\
&t = 1, \dots, 2^s
\end{aligned}$$

$A_{n2}^c$  is the union of all complementary sets and the approximations yield

$$\begin{aligned}
P(A_{n2}^c) &= O(n^{-\lambda}) \left[ \sum_{s=1}^{s_n} 1 + 1 + \sum_{s=1}^{s_n-1} \sum_{t=1}^{2^s-1} 1 \right] + \sum_{s=s_n}^{\infty} \sum_{t=1}^{2^s} O(n^{-\lambda} e^{-s}) \\
&= O(n^{-\lambda}) O(\ln n + n) + O(n^{-\lambda}) \tag{12}
\end{aligned}$$

so  $P(A_{n2}^c)$  is less than an  $O(n^{-\lambda+1})$ .

**Proof** According to Bernstein's inequality (e.g. Shorack, Wellner (1986), p. 855), for all  $\varphi'_{st}$  and analogously for all  $\psi'_{st}$  with  $s < s_n$  it holds that:

$$P \left( \frac{1}{n} \left| \sum_{j=1}^n \varphi'_{st}(X_j) - E\varphi'_{st}(X_j) \right| > \frac{\lambda \ln n}{\sqrt{n}} \right)$$

$$\begin{aligned}
&\leq 2 \exp \left\{ - \frac{\frac{n}{2} \left( \frac{\lambda \ln n}{\sqrt{n}} \right)^2}{E|\varphi'_{st}|^2 + \|\varphi'_{st}\|_\infty \frac{\lambda \ln n}{3\sqrt{n}}} \right\} \\
&\leq 2 \exp \left\{ - \frac{\lambda^2 \ln^2 n}{\frac{1}{\pi} \|f\|_\infty + \frac{1}{\pi} \sqrt{2n} \frac{\lambda \ln n}{3\sqrt{n}}} \right\} \\
&= O \left( \exp \left\{ - \frac{\lambda \ln n}{1 + \|f\|_\infty \lambda^{-1} \ln^{-1} n} \right\} \right)
\end{aligned}$$

which is a uniform  $O(n^{-\lambda})$ . For  $\psi'_{st}$  with  $s \geq s_n$ :

$$P \left( \frac{1}{n} \left| \sum_{j=1}^n \psi'_{st}(X_j) - E\psi'_{st}(X_j) \right| > \frac{\lambda \ln n + s}{\sqrt{n}} \right) = O \left( \exp \left\{ - \frac{\lambda \ln n + s}{1 + \|f\|_\infty \lambda^{-1} \ln^{-1} n} \right\} \right)$$

a uniform  $O(n^{-\lambda} e^{-s})$ . □

**Acknowledgement:** I kindly thank Prof. M. Neumann for initiating and supporting the present work. Further I thank Prof. A. Munk for some profound comments on curve smoothing.

## References

- Arcones, M.A. and Giné, E. (1993). Limit theorems for  $U$ -processes. *Ann. Probab.* **21** No. 3, 1494-1542.
- Cai, T. (2003) Rates of convergence and adaptation over Besov spaces under pointwise risk. *Stat. Sin.* **13** No.3, 881-902.
- Cavalier, L. and Tsybakov, A. (2001). Penalized blockwise Stein's method, monotone oracles and sharp adaptive estimation. *Math. Methods Stat.* **10** No.3, 247-282.
- Dalelane, C. (2005). Data driven kernel choice in non-parametric curve estimation. *Ph.D. dissertation*. TU Braunschweig. (available at <http://opus.tu-bs.de/opus/volltexte/2005/659/>)
- Donoho, D. and Johnstone, I. (1994) Ideal spatial adaptation by wavelet shrinkage. *Biometrika* **81** No.3, 425-455 (1994).
- Efremovich, S.Yu. and Pinsker, M.S. (1983). Estimation of square-integrable probability density of a random variable. *Probl. Inf. Transm.* **18**, 175-189.
- Efremovich, S. (2004). Oracle inequalities for Efremovich-Pinsker blockwise estimates. *Methodol. Comput. Appl. Probab.* **6** No.3, 303-322
- Hall, P. (1983). Large sample optimality of least squares cross-validation in density estimation. *Ann. Statist.* **11**, 1156-1174.
- Hall, P.; Kerkycharian, G. and Picard, D. (1999). Block threshold rules for curve estimation. *Ann. Statist* **43** No.4, 415-420.
- Kneip, A. (1994). Ordered linear smoothers. *Ann. Statist.* **22** No.2, 835-866.
- Nolan, D. and Pollard, D. (1987).  $U$ -processes: rates of convergence. *Ann. Statist.* **15**, 780-799.
- Rigollet, P. (2004). Adaptive density estimation using Stein's blockwise method. *Preprint PMA-913* (available at [www.proba.jussieu.fr](http://www.proba.jussieu.fr))
- Schipper, M. (1996). Optimal rates and constants in  $L_2$ -minimax estimation of probability density functions. *Math. Methods Stat.* **5** No.3, 253-274.

- Shorack, G.R. and Wellner, A.J. (1986). *Empirical Processes with applications to statistics*. John Wiley & Sons. New York.
- Stone, C.J. (1984). An asymptotically window selection rule for kernel density estimates. *Ann. Statist.* **12**, 1285-1297.