

Adaptive Estimation of a Distribution Function and its Density in Sup-Norm Loss by Wavelet and Spline Projections

EVARIST GINÉ* AND RICHARD NICKL†

University of Connecticut and University of Cambridge

First Version: May 2008, This Version: August 2009

Abstract

Given an i.i.d. sample from a distribution F on \mathbb{R} with uniformly continuous density p_0 , purely-data driven estimators are constructed that efficiently estimate F in sup-norm loss, and simultaneously estimate p_0 at the best possible rate of convergence over Hölder balls, also in sup-norm loss. The estimators are obtained from applying a model selection procedure close to Lepski's method with random thresholds to projections of the empirical measure onto spaces spanned by wavelets or B -splines. The random thresholds are based on suprema of Rademacher processes indexed by wavelet or spline projection kernels. This requires Bernstein-analogues of the inequalities in Koltchinskii (2006) for the deviation of suprema of empirical processes from their Rademacher symmetrizations.

MSC 2000 subject classification: Primary: 62G07; Secondary: 60F05.

Key words and phrases: adaptive estimation, Rademacher processes, sup-norm loss, wavelet estimator, spline estimator, Lepski's method.

1 Introduction

If X_1, \dots, X_n are i.i.d. with unknown distribution function F on \mathbb{R} , then classical results of mathematical statistics establish optimality of the empirical distribution function F_n as an estimator of F . That is to say, if we assume no apriori knowledge whatsoever on F , and equip the set of all probability distribution functions with some natural loss function, such as sup-norm loss, then F_n is asymptotically sharp minimax for estimating F . (The same is true even if more is known about F , for instance if F is known to have a uniformly continuous density.) However, this does not preclude the existence of other estimators that are also asymptotically minimax for estimating F in sup-norm loss, but which improve upon F_n in other respects. What we have in mind is a purely-data driven estimator that is efficient for F but at the same time also estimates the density f of F at the best rate of convergence in some relevant loss-function over some prescribed classes of densities. More precisely, our goal in the present article is to construct estimators that satisfy the functional central limit theorem (CLT) for the distribution function *and* which adapt to the unknown smoothness of the density in *supnorm loss*. Whereas this article is concerned

*Department of Mathematics, University of Connecticut, Storrs, CT 06269-3009, USA. gine@math.uconn.edu

†Statistical Laboratory, Department of Pure Mathematics and Mathematical Statistics, University of Cambridge, Wilberforce Road, CB3 0WB Cambridge, U.K. nickl@statslab.cam.ac.uk

with the mathematical problem of the existence and construction of such estimators, it does not deal with the practical implementation of estimation procedures.

To achieve adaptation one can opt for several approaches, all of which are related. Among them we mention the penalization method of Barron, Birgé and Massart (1999), wavelet thresholding (Donoho, Johnstone, Kerkycharian and Picard (1996)), and Lepski's (1991) method. Our choice for the goal at hand consists of using Lepski's method, with random thresholds, applied to wavelet and spline projection estimators of a density.

The linear estimators underlying our procedure are projections of the empirical measure onto spaces spanned by wavelets, and wavelet theory is central to some of the derivations of this article. The wavelets most commonly used in statistics are those that are compactly supported (e.g., Daubechies' wavelets), and our results readily apply to these. However, for computational and other purposes, projections onto spline spaces are also interesting candidates for the estimators. Density estimators obtained from projecting the empirical measure onto Schoenberg spaces spanned by B -splines were studied by Huang and Studden (1993). As is well-known in wavelet theory, the Schoenberg spline spaces with equally spaced knots have an orthonormal basis consisting of the Battle-Lemarié wavelets, so that the spline projection estimator is in fact exactly equal to the wavelet estimator based on Battle-Lemarié wavelets. These wavelets do not have compact support but they are exponentially localized. Although we cannot handle in general exponentially decaying wavelets, we can still work with Battle-Lemarié wavelets because the B -spline expansion of the projections allows us to show that the relevant classes of functions are of Vapnik-Cervonenkis type, so that empirical process techniques can be applied. In particular, the adaptive estimators we devise in Theorem 3 may be based either on spline projections or on compactly supported wavelets. And in the process of proving the main theorem, we also provide new asymptotic results for spline projection density estimators similar to those for wavelet estimators in Giné and Nickl (2009b).

We need to use Talagrand's exponential inequality with sharp constants (Bousquet (2003), Klein and Rio (2005)) in the proofs, but to do this, we have to estimate the expectation of suprema of certain empirical processes that appear in the centering of Talagrand's inequality. The use of entropy-based moment inequalities for empirical processes typically results in too conservative constants (e.g., in Giné and Nickl (2009a)). In order to remedy this problem, we adapt recent ideas due to Koltchinskii (2001, 2006) and Bartlett, Boucheron and Lugosi (2002) to density estimation: the entropy based moment bounds are replaced by the sup norm of the associated Rademacher averages, which are, with high probability, better estimates of the expected value of the supremum of the empirical process. We derive a Bernstein-type analogue of an exponential inequality in Koltchinskii (2006) that shows how the supremum of an empirical process deviates from the supremum of the associated Rademacher processes. This Bernstein-type version allows to use partial knowledge on the variance of the empirical processes involved, which is crucial for applications in our context of adaptive density estimation. Moreover, we show that one can use, instead of the supremum of the Rademacher process, its conditional expectation given the data.

Adaptive estimation in sup-norm loss is a relatively recent subject. We should mention the results due to Tsybakov (1998), Golubev, Lepski and Levit (2001) – who only considered Sobolev-type smoothness conditions – and Goldenshluger and Lepski (2009). All these results were obtained in the Gaussian white noise model. If one is interested in adapting to a Hölder-continuous density in sup-norm loss in the i.i.d. density model on \mathbb{R} , this simplifying Gaussian structure is not available, and novel techniques are needed. In the i.i.d. density model on \mathbb{R} , a direct 'competitor' to the estimators constructed in this article is the hard thresholding wavelet density estimator introduced in Donoho et al. (1996): as proved in Giné and Nickl (2009b), its distribution function satisfies the functional CLT and it is adaptive in the sup-norm over Hölder balls; however, the proofs there require the additional assumption that dF integrates $|x|^\delta$ for some $\delta > 0$, and the

constants appearing in the threshold and the risk become quite large for δ small. The results in the present article hold under no moment condition whatsoever.

2 Wavelets expansions and estimators

We start with some basic notation. If (S, \mathcal{S}) is a measurable space, and for Borel-measurable functions $h : S \rightarrow \mathbb{R}$ and Borel measures μ on S , we set $\mu h := \int_S h d\mu$. We will denote by $L^p(Q) := L^p(S, Q)$, $1 \leq p \leq \infty$, the usual Lebesgue spaces on S w.r.t. a Borel measure Q , and if Q is Lebesgue measure on $S = \mathbb{R}$ we simply denote this space by $L^p(\mathbb{R})$, and its norm by $\|\cdot\|_p$ if $p < \infty$. We will use the symbol $\|h\|_\infty$ to denote $\sup_{x \in \mathbb{R}} |h(x)|$ for $h : \mathbb{R} \rightarrow \mathbb{R}$. For $s \in \mathbb{N}$, denote by $C^s(\mathbb{R})$ the spaces of functions $f : \mathbb{R} \rightarrow \mathbb{R}$ that are s -times differentiable with uniformly continuous $D^s f$, equipped with the norm $\|f\|_{s,\infty} = \sum_{0 \leq \alpha \leq s} \|D^\alpha f\|_\infty$, with the convention that $D^0 =: id$ and that then $C(\mathbb{R}) := C^0(\mathbb{R})$ is the space of bounded uniformly continuous functions. For noninteger $s > 0$ and $[s]$ the integer part of s , set

$$C^s(\mathbb{R}) = \left\{ f \in C^{[s]}(\mathbb{R}) : \|f\|_{s,\infty} := \sum_{0 \leq \alpha \leq [s]} \|D^\alpha f\|_\infty + \sup_{x \neq y} \frac{|D^{[s]}f(x) - D^{[s]}f(y)|}{|x - y|^{s-[s]}} < \infty \right\}.$$

2.1 Multiresolution analysis and wavelet bases

We recall here a few well-known facts about wavelet expansions, see, e.g., Sections 8 and 9 in Härdle, Kerkycharian, Picard and Tsybakov (HKPT, 1998). Let $\phi \in L^2(\mathbb{R})$ be a scaling function, that is, ϕ is such that $\{\phi(\cdot - k) : k \in \mathbb{Z}\}$ is an orthonormal system in $L^2(\mathbb{R})$, and moreover the linear spaces $V_0 = \{f(x) = \sum_k c_k \phi(x - k) : \{c_k\}_{k \in \mathbb{Z}} \in \ell^2\}$, $V_1 = \{h(x) = f(2x) : f \in V_0\}, \dots, V_j = \{h(x) = f(2^j x) : f \in V_0\}, \dots$, are nested ($V_{j-1} \subseteq V_j$ for $j \in \mathbb{N}$) and their union is dense in $L^2(\mathbb{R})$. In the case where ϕ is a bounded function that decays exponentially at infinity (i.e. $|\phi(x)| \leq C e^{-\gamma|x|}$ for some $C, \gamma > 0$) – which we assume for the rest of this subsection – the kernel of the projection onto the space V_j has certain properties: First, the series

$$K(y, x) := K(\phi, y, x) = \sum_{k \in \mathbb{Z}} \phi(y - k) \phi(x - k), \quad (1)$$

converges pointwise, and we set $K_j(y, x) := 2^j K(2^j y, 2^j x)$, $j \in \mathbb{N} \cup \{0\}$. Furthermore we have

$$|K(y, x)| \leq \Phi(|y - x|) \quad \text{and} \quad \sup_{x \in \mathbb{R}} \sum_k |\phi(x - k)| < \infty, \quad (2)$$

where $\Phi : \mathbb{R} \rightarrow \mathbb{R}^+$ is bounded and has exponential decay (cf. Lemma 8.6 in HKPT (1996)). For any j fixed, if $f \in L^p(\mathbb{R})$, $1 \leq p \leq \infty$, then the series

$$K_j(f)(y) := \int K_j(x, y) f(x) dx = \sum_{k \in \mathbb{Z}} 2^j \phi(2^j y - k) \int \phi(2^j x - k) f(x) dx, \quad y \in \mathbb{R}$$

converges pointwise, and for $f \in L^2(\mathbb{R})$, $K_j(f)$ coincides with the orthogonal projection $\pi_j : L^2(\mathbb{R}) \rightarrow V_j$ of f onto V_j . For $f \in L^1(\mathbb{R})$, which is the main case in this article, the convergence of the series in fact takes place in $L^p(\mathbb{R})$, $1 \leq p \leq \infty$. This still holds true if $f(x)dx$ is replaced by $d\mu(x)$, where μ is any finite signed measure. If now ϕ is a scaling function and ψ the associated

mother wavelet so that $\{\phi(\cdot - k), 2^{l/2}\psi(2^l(\cdot) - k) : k \in \mathbb{Z}, l \in \mathbb{N}\}$ is an orthonormal basis of $L^2(\mathbb{R})$, then any $f \in L^p(\mathbb{R})$ admits the formal expansion

$$f(y) = \sum_k \alpha_k(f) \phi(y - k) + \sum_{l=0}^{\infty} \sum_k \beta_{lk}(f) \psi_{lk}(y) \quad (3)$$

where $\psi_{lk}(y) = 2^{l/2}\psi(2^l y - k)$, $\alpha_k(f) = \int f(x) \phi(x - k) dx$, $\beta_{lk}(f) = \int f(x) \psi_{lk}(x) dx$. Since $(K_{l+1} - K_l)f = \sum_k \beta_{lk}(f) \psi_{lk}$, the partial sums of the series (3) are in fact given by

$$K_j(f)(y) = \sum_k \alpha_k(f) \phi(y - k) + \sum_{l=0}^{j-1} \sum_k \beta_{lk}(f) \psi_{lk}(y), \quad (4)$$

and, if ϕ, ψ are bounded and have exponential decay, then convergence of the series (4) holds pointwise, and it also holds in $L^p(\mathbb{R})$, $1 \leq p \leq \infty$, if $f \in L^1(\mathbb{R})$ or if f is replaced by a finite signed measure. Now, using these facts one can furthermore show that the wavelet series (3) converges in $L^p(\mathbb{R})$, $p < \infty$, for $f \in L^p(\mathbb{R})$, and we also note that if p_0 is a uniformly continuous density, then its wavelet series converges uniformly.

2.2 Density Estimation using wavelet and spline projection kernels

Let X_1, \dots, X_n be i.i.d. random variables with common law P and density p_0 on \mathbb{R} , and denote by $P_n = \frac{1}{n} \sum_{i=1}^n \delta_{X_i}$ the associated empirical measure. A natural first step is to estimate the projection $K_j(p_0)$ of p_0 onto V_j by

$$p_n(y) := p_n(y, j) = \frac{1}{n} \sum_{i=1}^n K_j(y, X_i) = \sum_k \hat{\alpha}_k \phi(y - k) + \sum_{l=0}^{j-1} \sum_k \hat{\beta}_{lk} \psi_{lk}(y) \quad y \in \mathbb{R}, \quad (5)$$

where K is as in (1), $j \in \mathbb{N}$, and where $\hat{\alpha}_k = \int \phi(x - k) dP_n(x)$, $\hat{\beta}_{lk} = \int \psi_{lk}(x) dP_n(x)$ are the empirical wavelet coefficients. We note that for ϕ, ψ compactly supported (e.g., Daubechies' wavelets), there are only finitely many k 's for which these coefficients are nonzero. This estimator was first studied by Kerkycharian and Picard (1992) for compactly supported wavelets.

If the wavelets ϕ and ψ do not have compact support, it may be impossible to compute the estimator exactly, since the sums over k consist of infinitely many summands. However, in the special case of the Battle-Lemarié family $\phi_r, r \geq 1$ (see, e.g., Section 6.1 in HKPT (1998)) – which is a class of non-compactly supported but exponentially decaying wavelets – the estimator has a simple form in terms of splines: the associated spaces $V_{j,r} = \{\sum_k c_k 2^{j/2} \phi_r(2^j(\cdot) - k) : \sum_k c_k^2 < \infty\}$ are in fact equal to the *Schoenberg* spaces generated by the Riesz-basis of B -splines of order r , so that the sum in (5) can be computed by

$$p_n(y, j) := \frac{1}{n} \sum_{i=1}^n \kappa_j(y, X_i) = \frac{2^j}{n} \sum_{i=1}^n \sum_k \sum_l b_{kl} N_{j,k,r}(X_i) N_{j,l,r}(y) \quad y \in \mathbb{R}, \quad (6)$$

where the $N_{j,k,r}$ (are suitably translated and dilated) B -splines of order r , the kernel κ is as in (29) below and the b_{kl} 's are the entries of the inverse of the matrix defined in (28) below. An exact derivation of this spline projection, their wavelet representation and detailed definitions are given in Section 3.2. It turns out that for every sample point X_i and for every y , each of the last two sums extends only over r terms. We should note that this 'spline projection' estimator was first studied (outside of the wavelet setting) by Huang and Studden (1993), who

derived pointwise rates of convergence. See also Huang (1999), where some comparison between Daubechies' and spline wavelets can be found.

In the course of proving the main theorem of this article, we will derive some basic results for the linear spline projection estimator (6), that we now state. For classical kernel estimators, results similar to those that follow were obtained in Deheuvels (2000) Giné and Guillou (2002) and Giné and Nickl (2009a), and for wavelet estimators based on compactly supported wavelets, this was done in Giné and Nickl (2009b).

Theorem 1 *Suppose that P has a bounded density p_0 . Assume $j_n \rightarrow \infty$, $n/(j_n 2^{j_n}) \rightarrow \infty$, $j_n/\log \log n \rightarrow \infty$ and $j_{2n} - j_n \leq \tau$ for some τ positive. Let $p_n(y) = p_n(y, j_n)$ be the estimator from (6) for some $r \geq 1$. Then*

$$\limsup_n \sqrt{\frac{n}{2^{j_n} j_n}} \sup_{y \in \mathbb{R}} |p_n(y) - E p_n(y)| = C \quad a.s.$$

and, for $1 \leq p < \infty$,

$$\sup_n \sqrt{\frac{n}{2^{j_n} j_n}} \left(E \sup_{y \in \mathbb{R}} |p_n(y) - E p_n(y)|^p \right)^{1/p} \leq C'$$

where C and C' depend only on $\|p_0\|_\infty$ and on r, p, τ .

For rates of convergence in probability the conditions on j_n can be weakened, cf. Proposition 3 below. The usual approximation error bounds for projections onto spaces spanned by the respective wavelet basis then also give, for $p_0 \in C^t(\mathbb{R})$ with $t \leq r$ and $2^{j_n} \simeq (n/\log n)^{1/(2t+1)}$, that

$$\sup_{y \in \mathbb{R}} |p_n(y) - p_0(y)| = O \left(\left(\frac{\log n}{n} \right)^{t/(2t+1)} \right) \quad \text{both a.s. and in } L^p(P).$$

For the following central limit theorem, we denote by $\rightsquigarrow_{\ell^\infty(\mathbb{R})}$ convergence in law for sample-bounded processes in the Banach space of bounded functions on \mathbb{R} , and by G_P the usual P -Brownian bridge (e.g., Chapter 3 in Dudley (1999)). We should emphasize that the optimal bandwidth choice $2^{-j_n} \simeq n^{-1/2t+1}$ (or, if sup-norm loss is considered, n replaced by $n/\log n$) is admissible for every $t > 0$ in the theorem below.

Theorem 2 *Assume that the density p_0 of P is a bounded function ($t = 0$) or that $p_0 \in C^t(\mathbb{R})$ for some t , $0 < t \leq r$. Let j_n satisfy $n/(2^{j_n} j_n) \rightarrow \infty$ and $\sqrt{n} 2^{-j_n(t+1)} \rightarrow 0$ as $n \rightarrow \infty$. If F is the distribution function of P and setting $F_n^S(s) := \int_{-\infty}^s p(y, j_n) dy$, then*

$$\sqrt{n}(F_n^S - F) \rightsquigarrow_{\ell^\infty(\mathbb{R})} G_P.$$

Proof. Given $\varepsilon > 0$, apply Proposition 4 below with $\lambda = \varepsilon$, so that $\|F_n^S - F_n\|_\infty = o_P(1/\sqrt{n})$ follows, and use the fact that $\sqrt{n}(F_n - F)$ converges in law in $\ell^\infty(\mathbb{R})$ to G_P . ■

3 The adaptive estimation procedures

In this section we construct data-driven choices of the resolution level j and state the main adaptation results. As mentioned in the introduction, we will use Rademacher symmetrization for this. Generate a Rademacher sequence ε_i , $i = 1, \dots, n$, independent of the sample, (that is, ε_i takes values $1, -1$ with probability $1/2$) and set, for $j < l$,

$$R(n, j) = 2 \left\| \frac{1}{n} \sum_{i=1}^n \varepsilon_i K_j(X_i, \cdot) \right\|_\infty \quad \text{and} \quad T(n, j, l) = 2 \left\| \frac{1}{n} \sum_{i=1}^n \varepsilon_i (K_j - K_l)(X_i, \cdot) \right\|_\infty, \quad (7)$$

where K_j is the kernel of the wavelet projection π_j onto V_j (both for Battle-Lemarié and compactly supported wavelets). In both cases, these are suprema of fixed random functions that depend only on known quantities that can be computed in a numerically effective way. For more details on Rademacher processes see Subsection 3.1.1 below.

To construct the estimators, we first need a grid indexing the spaces V_j onto which we project P_n . For $r \geq 1$, $n > 1$, choose integers $j_{\min} := j_{\min,n}$ and $j_{\max} := j_{\max,n}$ such that $0 < j_{\min} < j_{\max}$,

$$2^{j_{\min}} \simeq \left(\frac{n}{\log n} \right)^{1/(2r+1)} \quad \text{and} \quad 2^{j_{\max}} \simeq \frac{n}{(\log n)^2} \quad (8)$$

and set

$$\mathcal{J} := \mathcal{J}_n = [j_{\min}, j_{\max}] \cap \mathbb{N}.$$

Note that the number of elements in this grid is of order $\log n$. We will consider two preliminary estimators \bar{j}_n and \tilde{j}_n of the resolution level (of course, only one is needed, but we offer a choice amongst two as discussed below). Let $p_n(j)$ be as in (5) or (6). First, we set

$$\bar{j}_n = \min \left\{ j \in \mathcal{J} : \|p_n(j) - p_n(l)\|_{\infty} \leq T(n, j, l) + 7\|\Phi\|_2 \|p_n(j_{\max})\|_{\infty}^{1/2} \sqrt{\frac{2^l l}{n}} \quad \forall l > j, l \in \mathcal{J} \right\} \quad (9)$$

where the function Φ is as in (2), and we discuss an explicit way to construct Φ in Remark 2 below. If the minimum does not exist, we set \bar{j}_n equal to j_{\max} . An alternative estimator of the resolution level is

$$\tilde{j}_n = \min \left\{ j \in \mathcal{J} : \|p_n(j) - p_n(l)\|_{\infty} \leq (B(\phi) + 1)R(n, l) + 7\|\Phi\|_2 \|p_n(j_{\max})\|_{\infty}^{1/2} \sqrt{\frac{2^l l}{n}} \quad \forall l > j, l \in \mathcal{J} \right\} \quad (10)$$

where $B(\phi)$ is a bound, uniform in j , for the operator norm in $L^{\infty}(\mathbb{R})$ of the projection π_j , see Remark 3 below. Again, if the minimum does not exist, we set \tilde{j}_n equal to j_{\max} .

Before we state the main result, we briefly discuss these procedures: The data-driven resolution level \tilde{j}_n in (10) is based on tests that use Rademacher-analogues of the usual thresholds in Lepski's method: Starting with j_{\min} , the main contribution to $\|p_n(j) - p_n(l)\|_{\infty}$ is the bias $\|Ep_n(j) - p_0\|_{\infty}$. The procedure should stop when the 'variance term' $\|p_n(l) - Ep_n(l)\|_{\infty}$ starts to dominate. Since this is an unknown quantity, and since we know no good nonrandom upper bound for it, we estimate it by the supremum of the associated Rademacher process, i.e., by $R(n, l)$. The constant $B(\phi)$ is necessary in order to correct for the lack of monotonicity of the $R(n, l)$'s in the resolution level l .

The estimator \bar{j}_n in (9) is somewhat more refined: It tries to take advantage of the fact that in the 'small bias' domain, and using the results from Subsection 3.1.1,

$$\|p_n(j) - p_n(l)\|_{\infty} = \left\| \frac{1}{n} \sum_{i=1}^n (K_j - K_l)(X_i, \cdot) \right\|_{\infty}$$

should not exceed its Rademacher symmetrization

$$T(n, j, l) = 2 \left\| \frac{1}{n} \sum_{i=1}^n \varepsilon_i (K_j - K_l)(X_i, \cdot) \right\|_{\infty}.$$

We now state the main result, whose proof is deferred to the next section. As usual, we say that a wavelet basis is s -regular, $s \in \mathbb{N} \cup \{0\}$, if either the scaling function ϕ has s weak derivatives contained in $L^p(\mathbb{R})$ for some $p \geq 1$, or if the mother wavelet ψ satisfies $\int x^{\alpha} \psi(x) dx = 0$

for $\alpha = 0, \dots, s$. Note that any compactly supported element of $C^s(\mathbb{R})$, $s > 0$, is of bounded $(1/s)$ -variation, so that the p -variation condition in the following theorem is satisfied, e.g., for all Daubechies-wavelets. The estimators below achieve the optimal rate of convergence for estimating p_0 in sup-norm loss in the minimax sense (over Hölder balls), cf., e.g., Korostelev and Nussbaum (1999) for optimality of these rates.

Theorem 3 *Let X_1, \dots, X_n be i.i.d. on \mathbb{R} with common law P that possesses a uniformly continuous density p_0 . Let $p_n(j) := p_n(y, j)$ be as in (5), where ϕ is either compactly supported, of bounded p -variation ($p < \infty$) and $(r-1)$ -regular, or $\phi = \phi_r$ equals a Battle-Lemarié wavelet. Let the sequence $\{\hat{j}_n\}_{n \in \mathbb{N}}$ be either $\{\bar{j}_n\}_{n \in \mathbb{N}}$ or $\{\tilde{j}_n\}_{n \in \mathbb{N}}$, and let $F_n(\hat{j}_n)(t) = \int_{-\infty}^t p_n(y, \hat{j}_n) dy$. Then*

$$\sqrt{n} \left(F_n(\hat{j}_n) - F \right) \rightsquigarrow_{\ell^\infty(\mathbb{R})} G_P, \quad (11)$$

the convergence being uniform over the set of all probability measures P on \mathbb{R} with densities p_0 bounded by a fixed constant, in any distance that metrizes convergence in law. Furthermore, if C is any precompact subset of $C(\mathbb{R})$, then

$$\sup_{p_0 \in C} E \sup_{y \in \mathbb{R}} |p_n(y, \hat{j}_n) - p_0(y)| = o(1). \quad (12)$$

If, in addition, $p_0 \in C^t(\mathbb{R})$ for some $0 < t \leq r$ then also

$$\sup_{p_0: \|p_0\|_{t, \infty} \leq D} E \sup_{y \in \mathbb{R}} |p_n(y, \hat{j}_n) - p_0(y)| = O \left(\left(\frac{\log n}{n} \right)^{t/(2t+1)} \right). \quad (13)$$

Remark 1 *Relaxing the uniform continuity assumption.* The assumption of uniform continuity of the density of F can be relaxed by modifying the definition of \bar{j}_n (or \tilde{j}_n) along the lines of Giné and Nickl (2009a): The idea is to constrain all candidate estimators to lie in a ball of size $o(1/\sqrt{n})$ around the empirical distribution function F_n so that (11) holds automatically. Formally, this can be done by adding the requirement

$$\sup_{t \in \mathbb{R}} \left| \int_{-\infty}^t p_n(y, j) dy - F_n(t) \right| \leq \frac{1}{\sqrt{n} \log n}$$

in each test in (9) or (10). If this requirement does not even hold for j_{max} , it can be seen as evidence that F has no density, and one just uses F_n as the estimator, so as to obtain at least the functional CLT. If F has a bounded density, one can use the exponential bound in Proposition 4 in the proof to control rejection probabilities of these test in the 'small bias' domain $\hat{j}_n > j^*$, and Theorem 3 can then still be proved for this procedure, without *any* assumptions on F . See Theorem 2 in Giné and Nickl (2009a) for more details on this procedure and its proof.

Remark 2 *The constant $\|\Phi\|_2$.* Once the wavelet ϕ have been chosen, \hat{j}_n is purely data driven since the function Φ depends only on ϕ . For the Haar basis ($\phi = I_{[0,1)}$) we can take $\Phi = \phi$ because in this case $K(x, y) \leq I_{[0,1)}(|x - y|)$ so that $\|\Phi\|_2 = 1$. A general way to obtain majorizing kernels Φ is described in Section 8.6 of HKPT (1998). For Battle-Lemarié wavelets, the spline representation of the projection kernel is again useful for estimating $\|\Phi\|_2$. See Huang and Studden (1993) for explicit computations.

Remark 3 *The constant $B(\phi)$.* To construct \tilde{j}_n one requires knowledge of the constant $B(\phi)$ that bounds the operator norm $\|\pi_j\|'_\infty$ of π_j viewed as an operator $L^\infty(\mathbb{R})$. A simple way of obtaining a bound is as follows: for any $f \in L^\infty(\mathbb{R})$ we have, by (2),

$$|\pi_j(f)(x)| = \left| \int K_j(x, y) f(y) dy \right| \leq \|\Phi\|_1 \|f\|_\infty,$$

that is, $\|\pi_j\|'_\infty \leq \|\Phi\|_1$. In combination with the previous remark, one readily obtains possible values for $B(\phi)$. For instance, for the Haar wavelet, $B(\phi) \leq 1$. For spline wavelets, other methods are available. For example, for Battle-Lemarié wavelets arising from linear B -splines, $\|\pi_j\|'_\infty$ is bounded by 3, and Shadrin (2001, p.135) conjectures the bound $2r - 1$ for general order r . See DeVore and Lorentz (1996, Chapter 13.4), Shadrin (2001) and references therein for more information.

We also note that – as the results in Subsection 3.1.1, in particular Proposition 2 show – all our proofs go through if one replaces $R(n, j)$, $T(n, j, l)$ by their Rademacher expectations $E^\varepsilon R(n, j)$, $E^\varepsilon T(n, j, l)$ in the definition of \tilde{j}_n , \tilde{j}_n .

3.1 Estimating Suprema of Empirical Processes

Talagrand's (1996) exponential inequality for empirical processes (see also Ledoux (2001)), which is a uniform Prohorov type inequality, is not specific about constants. Constants in its Bernstein type version have been specified by several authors (Massart (2000), Bousquet (2003) and Klein and Rio (2005)). Let X_i be the coordinates of the product probability space $(S, \mathcal{S}, P)^\mathbb{N}$, where P is any probability measure on (S, \mathcal{S}) , and let \mathcal{F} be a countable class of measurable functions on S that take values in $[-1/2, 1/2]$, or, if \mathcal{F} is P -centered, in $[-1, 1]$. Let $\sigma \leq 1/2$ and V be any two numbers satisfying

$$\sigma^2 \geq \|Pf^2\|_{\mathcal{F}}, \quad V \geq n\sigma^2 + 2E \left\| \sum_{i=1}^n (f(X_i) - Pf) \right\|_{\mathcal{F}}, \quad (14)$$

in which case V is also an upper bound for $E \left\| \sum_{i=1}^n (f(X_i) - Pf)^2 \right\|_{\mathcal{F}}$ (Klein and Rio (2005)). Then, taking into account that $\sup_{f \in \mathcal{F} \cup (-\mathcal{F})} \sum_{i=1}^n f(X_i) = \sup_{\mathcal{F}} |\sum_{i=1}^n f(X_i)|$, Bousquet's (2003) version of Talagrand's inequality is as follows: For every $t > 0$,

$$\Pr \left\{ \left\| \sum_{i=1}^n (f(X_i) - Pf) \right\|_{\mathcal{F}} \geq E \left\| \sum_{i=1}^n (f(X_i) - Pf) \right\|_{\mathcal{F}} + t \right\} \leq \exp \left(-\frac{t^2}{2V + \frac{2}{3}t} \right). \quad (15)$$

In the other direction, the Klein and Rio (2005) result is: For every $t > 0$,

$$\Pr \left\{ \left\| \sum_{i=1}^n (f(X_i) - Pf) \right\|_{\mathcal{F}} \leq E \left\| \sum_{i=1}^n (f(X_i) - Pf) \right\|_{\mathcal{F}} - t \right\} \leq \exp \left(-\frac{t^2}{2V + 2t} \right). \quad (16)$$

These inequalities can be applied in conjunction with an estimate of the expected value obtained via empirical processes methods. Here we describe one such result for VC type classes, i.e., for \mathcal{F} satisfying the uniform metric entropy condition

$$\sup_Q N(\mathcal{F}, L^2(Q), \tau) \leq \left(\frac{A}{\tau} \right)^v, \quad 0 < \tau \leq 1, (A \geq e, v \geq 2). \quad (17)$$

with the supremum extending over all Borel probability measures on (S, \mathcal{S}) . We denote here by $N(\mathcal{G}, L^2(Q), \tau)$ the usual covering numbers of a class \mathcal{G} of functions by balls of radius less than or equal to τ in $L^2(Q)$ -distance. Then one has, for every n

$$E \left\| \sum_{i=1}^n (f(X_i) - Pf) \right\|_{\mathcal{F}} \leq 2 \left[15 \sqrt{2vn\sigma^2 \log \frac{5A}{\sigma}} + 1350v \log \frac{5A}{\sigma} \right], \quad (18)$$

see Proposition 3 in Giné and Nickl (2009a) with a change obtained by using V as in (14) instead of an earlier bound due to Talagrand for $E \left\| \sum (f(X_i) - Pf)^2 \right\|_{\mathcal{F}}$. This type of inequalities has also some history (Talagrand (1994), Einmahl and Mason (2000), Giné and Guillou (2001), Giné and Koltchinskii (2006), among others). The constants at the right hand side of (18) may be far from best possible, but we prefer them over unspecified 'universal' constants.

As is the case of Bernstein's inequality in \mathbb{R} , Talagrand's inequality is especially useful in the Gaussian tail range, and, combining (15) and (18), one can obtain such a 'Gaussian tail' bound for the supremum of the empirical process that depends only on σ (similar to a bound in Giné and Guillou (2001)).

Proposition 1 *Let \mathcal{F} be a countable class of measurable functions that satisfies (17), and is uniformly bounded (in absolute value) by $1/2$. Assume further that for some $\lambda > 0$,*

$$n\sigma^2 \geq \frac{\lambda^2 v}{2} \log \frac{5A}{\sigma}. \quad (19)$$

Set $c_1(\lambda) = 2[15 + 1350\lambda^{-1}]$ and let $c_2(\lambda) \geq 1 + 120\lambda^{-1} + 10800\lambda^{-2}$. Then, if

$$c_1(\lambda) \sqrt{2vn\sigma^2 \log \frac{5A}{\sigma}} \leq t \leq \frac{3}{2} c_2(\lambda) n\sigma^2, \quad (20)$$

we have

$$\Pr \left\{ \left\| \sum_{i=1}^n (f(X_i) - Pf) \right\|_{\mathcal{F}} \geq 2t \right\} \leq \exp \left(-\frac{t^2}{3c_2(\lambda)n\sigma^2} \right). \quad (21)$$

Proof. Under (19), inequality (18) gives $E \left\| \sum_{i=1}^n (f(X_i) - Pf) \right\|_{\mathcal{F}} \leq c_1(\lambda) \sqrt{2vn\sigma^2 \log \frac{5A}{\sigma}}$ and (14) implies that we can take $V = c_2(\lambda)n\sigma^2$. Now the result follows from (15), taking into account that in the range of t 's $E \left\| \sum_{i=1}^n (f(X_i) - Pf) \right\|_{\mathcal{F}} \leq t \leq 3V/2$, (15) becomes

$$\Pr \left\{ \left\| \sum_{i=1}^n (f(X_i) - Pf) \right\|_{\mathcal{F}} \geq 2t \right\} \leq \exp \left(-\frac{t^2}{3V} \right).$$

■

The constants here may be too large for some applications, but they are not so in situations where λ can be taken very large, in particular in asymptotic considerations. [Then $c_1(\lambda) \rightarrow 30$ and $c_2(\lambda) \rightarrow 1$ as $\lambda \rightarrow \infty$.]

3.1.1 Estimating the size of empirical processes by Rademacher averages

The constants one could obtain from Proposition 1 are not satisfactory for the applications to adaptive estimation we have in mind. We now propose a remedy for this problem, inspired by a nice idea of Koltchinskii (2001) and Bartlett, Boucheron and Lugosi (2002) consisting in replacing the expectation of the supremum of an empirical process by the supremum of the associated Rademacher process, that they used in other contexts, namely in risk minimization and model selection. An inequality of this type (see Koltchinskii (2006), page 2602), is

$$\Pr \left\{ \left\| \sum_{i=1}^n (f(X_i) - Pf) \right\|_{\mathcal{F}} \geq 2 \left\| \sum_{i=1}^n \varepsilon_i f(X_i) \right\|_{\mathcal{F}} + 3t \right\} \leq \exp \left(-\frac{2t^2}{3n} \right), \quad (22)$$

where ε_i , $i \in \mathbb{N}$, are i.i.d. Rademacher random variables, independent of the X_i 's, all defined as coordinates on a large product probability space. Note that this bound does not take the variance V in (15) into account, but in the applications to density estimation that we have in mind, V is much smaller than n (it is of the order $n2^{-j_n}$, $j_n \rightarrow \infty$). We need a similar inequality, with the quantity n in the bound replaced by V , valid over a large enough range of t 's.

It will be convenient to use the following well-known symmetrization inequality (e.g., Dudley (1999), p.343):

$$\frac{1}{2}E \left\| \sum_{i=1}^n \varepsilon_i f(X_i) \right\|_{\mathcal{F}} - \frac{\sqrt{n}}{2} \|Pf\|_{\mathcal{F}} \leq E \left\| \sum_{i=1}^n (f(X_i) - Pf) \right\|_{\mathcal{F}} \leq 2E \left\| \sum_{i=1}^n \varepsilon_i f(X_i) \right\|_{\mathcal{F}}, \quad (23)$$

The following exponential bound is the Bernstein-type analogue of (22). Denote by E^ε expectation w.r.t. the Rademacher variables only.

Proposition 2 *Let \mathcal{F} be a countable class of measurable functions, uniformly bounded (in absolute value) by $1/2$. Then, for every $t > 0$,*

$$\Pr \left\{ \left\| \sum_{i=1}^n (f(X_i) - Ef(X)) \right\|_{\mathcal{F}} \geq 2 \left\| \sum_{i=1}^n \varepsilon_i f(X_i) \right\|_{\mathcal{F}} + 3t \right\} \leq 2 \exp \left(-\frac{t^2}{2V' + 2t} \right), \quad (24)$$

as well as

$$\Pr \left\{ \left\| \sum_{i=1}^n (f(X_i) - Ef(X)) \right\|_{\mathcal{F}} \geq 2E^\varepsilon \left\| \sum_{i=1}^n \varepsilon_i f(X_i) \right\|_{\mathcal{F}} + 3t \right\} \leq 2 \exp \left(-\frac{t^2}{2V' + 2t} \right), \quad (25)$$

where $V' = n\sigma^2 + 4E \left\| \sum_{i=1}^n \varepsilon_i f(X_i) \right\|_{\mathcal{F}}$.

Proof. We have

$$\begin{aligned} & \Pr \left\{ \left\| \sum_{i=1}^n (f(X_i) - Pf) \right\|_{\mathcal{F}} \geq 2 \left\| \sum_{i=1}^n \varepsilon_i f(X_i) \right\|_{\mathcal{F}} + 3t \right\} \leq \\ & \Pr \left\{ \left\| \sum_{i=1}^n (f(X_i) - Pf) \right\|_{\mathcal{F}} \geq 2E \left\| \sum_{i=1}^n \varepsilon_i f(X_i) \right\|_{\mathcal{F}} + t \right\} + \Pr \left\{ \left\| \sum_{i=1}^n \varepsilon_i f(X_i) \right\|_{\mathcal{F}} \leq E \left\| \sum_{i=1}^n \varepsilon_i f(X_i) \right\|_{\mathcal{F}} - t \right\}. \end{aligned}$$

For the first term, combining (23) with (15) gives

$$\Pr \left\{ \left\| \sum_{i=1}^n (f(X_i) - Pf) \right\|_{\mathcal{F}} \geq 2E \left\| \sum_{i=1}^n \varepsilon_i f(X_i) \right\|_{\mathcal{F}} + t \right\} \leq \exp \left(-\frac{t^2}{2V' + (2/3)t} \right).$$

For the second term, note that (16) applies to the randomized sums $\sum_{i=1}^n \varepsilon_i f(X_i)$ as well by just taking the class of functions

$$\mathcal{G} = \{g(\tau, x) = \tau f(x) : f \in \mathcal{F}\},$$

$\tau \in \{-1, 1\}$, instead of \mathcal{F} and the probability measure $\bar{P} = 2^{-1}(\delta_{-1} + \delta_1) \times P$ instead of P . Hence

$$\Pr \left\{ \left\| \sum_{i=1}^n \varepsilon_i f(X_i) \right\|_{\mathcal{F}} \leq E \left\| \sum_{i=1}^n \varepsilon_i f(X_i) \right\|_{\mathcal{F}} - t \right\} \leq \exp \left(-\frac{t^2}{2V' + 2t} \right), \quad (26)$$

since $V' \geq n\sigma^2 + 2E \|\sum_{i=1}^n \varepsilon_i f(X_i)\|_{\mathcal{F}}$. Combining the bounds completes the proof of (24).

It remains to prove (25). Let \mathcal{G} , \bar{P} be as above, let $Y_i = (\varepsilon_i, X_i)$, and note that \bar{P} is the law of Y_i . By convexity,

$$Ee^{-tE^\varepsilon \|\sum_{i=1}^n \varepsilon_i f(X_i)\|_{\mathcal{F}}} \leq Ee^{-t \|\sum_{i=1}^n \varepsilon_i f(X_i)\|_{\mathcal{F}}} = Ee^{-t \|\sum_{i=1}^n g(Y_i)\|_{\mathcal{G}}}$$

for all t . The Klein and Rio (2005) version (16) of Talagrand's inequality is in fact established by estimating the Laplace transform $Ee^{-t \|\sum_{i=1}^n g(Y_i)\|_{\mathcal{G}}}$, and Theorem 1.2a in Klein and Rio (2005) implies

$$Ee^{-tE^\varepsilon \|\sum_{i=1}^n \varepsilon_i (f(X_i) - Pf)\|_{\mathcal{F}}} \leq -tE \left\| \sum_{i=1}^n g(Y_i) \right\|_{\mathcal{G}} + \frac{V}{9} (e^{3t} - 3t + 1),$$

for $V \geq n\sigma^2 + 2E \|\sum_{i=1}^n g(Y_i)\|_{\mathcal{G}}$, which, by their proof of the implication (a) \Rightarrow (c) in that theorem, gives

$$\Pr \left\{ E^\varepsilon \left\| \sum_{i=1}^n \varepsilon_i f(X_i) \right\|_{\mathcal{F}} \leq E \left\| \sum_{i=1}^n \varepsilon_i f(X_i) \right\|_{\mathcal{F}} - t \right\} \leq \exp \left(-\frac{t^2}{2V' + 2t} \right).$$

The proof of (25) now follows as in the previous case. ■

For \mathcal{F} of VC type, the moment bound (18) is usually proved as a consequence of a bound for the Rademacher process. In fact, the proof of Proposition 3 in Giné and Nickl (2009a) shows

$$E \left\| \sum_{i=1}^n \varepsilon_i f(X_i) \right\|_{\mathcal{F}} \leq 15 \sqrt{2vn\sigma^2 \log \frac{5A}{\sigma}} + 1350v \log \frac{5A}{\sigma}, \quad (27)$$

where σ is as in (14), which we use in the following corollary, together with the previous proposition. The constant $c_2(\lambda)$ in the exponent below is still potentially large, but tends to one if $\lambda \rightarrow \infty$.

Corollary 1 *Let \mathcal{F} be a countable class of measurable functions that satisfies (17), and assume it to be uniformly bounded (in absolute value) by $1/2$. Assume further (19) for some $\lambda > 0$. Then for $0 < t \leq \frac{1}{20}c_2(\lambda)n\sigma^2$ with $c_2(\lambda)$ as in Proposition 1, we have*

$$\Pr \left\{ \left\| \sum_{i=1}^n (f(X_i) - Ef(X)) \right\|_{\mathcal{F}} \geq 2 \left\| \sum_{i=1}^n \varepsilon_i f(X_i) \right\|_{\mathcal{F}} + 3t \right\} \leq 2 \exp \left(-\frac{t^2}{2.1c_2(\lambda)n\sigma^2} \right),$$

and the same inequality holds if $\|\sum_{i=1}^n \varepsilon_i f(X_i)\|_{\mathcal{F}}$ is replaced by its E^ε expectation.

Proof. By (19) and (27), we have $V' \leq c_2(\lambda)n\sigma^2$, and the condition on t together with (24) give the result. ■

3.2 Projections onto spline spaces and their wavelet representation

We briefly review in this section how the wavelet estimator (5) for Battle-Lemarié wavelets can be represented as a spline projection estimator (6). We shall need the spline representation in some proofs, while the wavelet representation will be useful in others.

Let $T := T_j = \{t_i(j)\}_{i=-\infty}^{\infty} = 2^{-j}\mathbb{Z}$, $j \in \mathbb{Z}$, be a bi-infinite sequence of equally spaced knots $t_i := t_i(j)$. A function S is a spline of order r , or of degree $m = r - 1$, if on each interval (t_i, t_{i+1}) , it is a polynomial of degree less than or equal to m (and exactly of degree m on at least one

interval), and, at each breakpoint t_i , S is at least $m - 1$ -times differentiable. The Schoenberg space $\mathcal{S}_r(T) := \mathcal{S}_r(T, \mathbb{R})$ is defined as the set of all splines of order (less than or equal to) r , and it coincides with the space $\mathcal{S}_r(T, 1, \mathbb{R})$ in DeVore and Lorentz (1993, p.135). The space $\mathcal{S}_r(T_j)$ has a Riesz-basis formed by B -splines $\{N_{j,k,r}\}_{k \in \mathbb{Z}}$ that we now describe. [See Section 4.4 in Schumaker (1993) and p.138f. in DeVore and Lorentz (1993) for more details.] Define

$$N_{0,r}(x) = 1_{[0,1)} * \dots * 1_{[0,1)}(x), \quad r\text{-times} \quad = \sum_{i=0}^r \frac{(-1)^i \binom{r}{i} (x-i)_+^{r-1}}{(r-1)!}.$$

For $r = 2$, this is the linear B -spline (the usual 'hat' function), for $r = 3$ it is the quadratic, and for $r = 4$ it is the cubic B -spline. Set $N_{k,r}(x) := N_{0,r}(x - k)$. Then the elements of the Riesz-basis are given by

$$N_{j,k,r}(x) := N_{k,r}(2^j x) = N_{0,r}(2^j x - k).$$

By the Curry-Schoenberg theorem, any $S \in \mathcal{S}_r(T_j)$ can be uniquely represented as $S(x) = \sum_{k \in \mathbb{Z}} c_k N_{j,k,r}(x)$. The orthogonal projection $\pi_j(f)$ of $f \in L^2(\mathbb{R})$ onto $\mathcal{S}_r(T_j) \cap L^2(\mathbb{R})$ is derived, e.g., in DeVore and Lorentz (1993, p.401f.), where it is shown that $\pi_j(f) = 2^{j/2} \sum_{k \in \mathbb{Z}} c_k N_{j,k,r}$ with the coefficients $c_k := c_k(f)$ satisfying $(Ac)_k = 2^{j/2} \int N_{j,k,r}(x) f(x) dx$ where the matrix A is given by

$$a_{kl} = \int 2^j N_{j,k,r}(x) N_{j,l,r}(x) dx = \int N_{k,r}(x) N_{l,r}(x) dx. \quad (28)$$

The inverse A^{-1} of the matrix A exists (see Corollary 4.2 on p.404 in DeVore and Lorentz (1993)), and if we denote its entries by b_{kl} so that $c_k = 2^{j/2} \int \sum_l b_{kl} N_{j,l,r}(x) f(x) dx$, we have

$$\pi_j(f)(y) = 2^j \int \sum_k \sum_l b_{kl} N_{j,l,r}(x) N_{j,k,r}(y) f(x) dx = \int \kappa_j(x, y) f(x) dx,$$

where $\kappa_j(x, y) = 2^j \kappa(2^j x, 2^j y)$ with

$$\kappa(x, y) = \sum_k \sum_l b_{kl} N_{l,r}(x) N_{k,r}(y), \quad (29)$$

is the spline projection kernel. Note that κ is symmetric in its arguments.

Diagonalisation of the kernel κ of the projection operator π_j led in fact to one of the first examples of wavelets, see, e.g., p.21f. and Section 2.3 in Meyer (1992), Section 5.4 in Daubechies (1992), or Section 6.1 in HKPT (1998). There it is shown that there exists a $r - 1$ -times differentiable scaling function ϕ_r with exponential decay, the Battle-Lemarié wavelet of order r , such that

$$\mathcal{S}_r(T_j) \cap L^2(\mathbb{R}) = V_{j,r} = \left\{ \sum_k c_k 2^{j/2} \phi_r(2^j(\cdot) - k) : \sum_k c_k^2 < \infty \right\}.$$

This necessarily implies that the kernels κ and $K = K(\phi_r)$ describe the same projections in $L^2(\mathbb{R})$, and the following simple lemma shows that these kernels are in fact pointwise the same.

Lemma 1 *Let $\{N_{k,r}\}_{k \in \mathbb{Z}}$ be the Riesz-basis of B -splines of order $r \geq 1$, and let ϕ_r be the associated Battle-Lemarié scaling function. If K is as in (1) and κ is as in (29), then, for all $x, y \in \mathbb{R}$, we have*

$$K(x, y) = \kappa(x, y).$$

Proof. If $r = 1$, then $N_{0,1} = \phi_1$ since this is just the Haar-basis. So consider $r > 1$. Since $\{\phi_r(\cdot - k) : k \in \mathbb{Z}\}$ is an orthonormal basis of $\mathcal{S}_r(\mathbb{Z}) \cap L^2(\mathbb{R})$ (cf., e.g., Theorem 1 on p. 26 in Meyer (1992)), it follows that K and κ are the kernels of the same L^2 -projection operator, and therefore, for all $f, g \in L^2(\mathbb{R})$

$$\int \int (K(x, y) - \kappa(x, y)) f(x) g(y) dx dy = 0.$$

By density in $L^2(\mathbb{R} \times \mathbb{R})$ of linear combinations of products of elements of $L^2(\mathbb{R})$, this implies that κ and K are almost everywhere equal in \mathbb{R}^2 . We complete the proof by showing that both functions are continuous in \mathbb{R}^2 . For K , this follows from the decomposition

$$|K(x, y) - K(x', y')| \leq \sum_k |\phi_r(x - k) - \phi_r(x' - k)| |\phi_r(y - k)| + \sum_k |\phi_r(y - k) - \phi_r(y' - k)| |\phi_r(x' - k)|,$$

the uniform continuity of ϕ_r ($r > 1$) and relation (2). For κ we use the relation (31) below,

$$|\kappa(x, y) - \kappa(x', y')| \leq \sum_i |N_{i,r}(x) - N_{i,r}(x')| |H(y - i)| + \sum_i |H(y - i) - H(y' - i)| |N_{i,r}(x')|,$$

which implies continuity of κ on \mathbb{R}^2 since $N_{0,r}$ and H are uniformly continuous (as $N_{0,r}$ is and $\sum_i |g(|i|)| < \infty$), and since $N_{0,r}$ has compact support. ■

3.3 An Exponential inequality for the uniform deviations of the linear estimator

To control the uniform deviations of the linear estimators from their means, one can use inequalities for the empirical process indexed by classes of functions \mathcal{F} contained in

$$\mathcal{K} = \{2^{-j} K_j(\cdot, y) : y \in \mathbb{R}, j \in \mathbb{N} \cup \{0\}\}, \quad (30)$$

together with suitable bounds on the 'weak' variance σ .

If ϕ has compact support (and is of finite p -variation), it is proved in Lemma 2 in Giné and Nickl (2009b) that the class \mathcal{K} also satisfies the bound (17). However, the proof there does not apply to Battle-Lemarié wavelets. A different proof, using the Toeplitz- and band-limited structure of the spline projection kernel, still enables us to prove that these classes of functions are of Vapnik-Cervonenkis type.

Lemma 2 *Let \mathcal{K} be as in (30), where ϕ_r is a Battle-Lemarié wavelet for some $r \geq 1$. Then there exist finite constants $A \geq 2$ and $v \geq 2$ such that*

$$\sup_Q N(\mathcal{K}, L^2(Q), \varepsilon) \leq \left(\frac{A}{\varepsilon}\right)^v$$

for $0 < \varepsilon < 1$ and where the supremum extends over all Borel probability measures on \mathbb{R} .

Proof. In the case $r = 1$, ϕ_1 is just the Haar wavelet, in which case the results follows from Lemma 2 in Giné and Nickl (2009b). Hence assume $r \geq 2$.

The matrix A is Toeplitz since, by change of variables in (28), $a_{kl} = a_{k+1, l+1}$ for all $k, l \in \mathbb{Z}$, and it is band-limited because $N_{0,r}$ has compact support. It follows that also A^{-1} is Toeplitz, and we denote its entries by $b_{kl} = g(|k - l|)$ for some function g . Furthermore it is known (e.g., Theorem 4.3 on p.404 in DeVore and Lorentz (1993)) that the entries of the inverse of any

positive definite band-limited matrix satisfy $|b_{kl}| \leq c\lambda^{|k-l|}$ for some $0 < \lambda < 1$ and c finite. Now, following Huang and Studden (1992), we write

$$\sum_k g(|l-k|)N_{k,r}(x) = \sum_k g(|l-k|)N_{k-l,r}(x-l) = \sum_k g(|k|)N_{k,r}(x-l),$$

so that

$$2^{-j}\kappa_j(\cdot, y) = \sum_{l \in \mathbb{Z}} N_{j,l,r}(y)H(2^j(\cdot) - l) \quad (31)$$

where $H(x) = \sum_{k \in \mathbb{Z}} g(|k|)N_{k,r}(x)$ is a function of bounded variation: To see the last claim, note that $N_{0,r}$ is of bounded variation, and hence $\|N_{k,r}\|_{TV} = \|N_{0,r}\|_{TV}$ (where $\|\cdot\|_{TV}$ denotes the usual total-variation norm) so that $\|H\|_{TV} \leq \|N_{0,r}\|_{TV} \sum_{k \in \mathbb{Z}} |g(|k|)| < \infty$ because $\sum_k |b_{l,l-k}| \leq \sum_k c\lambda^{|k|} < \infty$. The last fact implies that

$$\mathcal{H} = \{H(2^j(\cdot) - l) : l \in \mathbb{Z}, j \in \mathbb{N} \cup \{0\}\}$$

satisfies, for finite constants $B > 1$ and $w \geq 1$

$$\sup_Q N(\mathcal{H}, L^2(Q), \varepsilon) \leq \left(\frac{B\|H\|_\infty}{\varepsilon} \right)^w, \quad \text{for } 0 < \varepsilon < \|H\|_\infty$$

as proved in Nolan and Pollard (1987). Since $N_{j,0,r}$ is zero if y is not contained in $[0, 2^{-j}r]$, the sum in (31), for fixed y and j , extends only over the l 's such that $2^j y - r \leq l < 2^j y$, hence consists of at most r terms. This implies that \mathcal{K} is contained in the set \mathcal{H}_r of linear combinations of at most r functions from \mathcal{H} , with coefficients bounded in absolute value by $\|N_{j,l,r}\|_\infty = \|N_{0,r}\|_\infty < \infty$. Given ε , let $\varepsilon' = \varepsilon/(2r \max(\|H\|_\infty, \|N_{0,r}\|_\infty))$. Let $\alpha_1, \dots, \alpha_{n_1}$ be an ε' -dense subset of $[-\|N_{0,r}\|_\infty, \|N_{0,r}\|_\infty]$ which, for $\varepsilon' < \|N_{0,r}\|_\infty$, has cardinality $n_1 \leq 3\|N_{0,r}\|_\infty/\varepsilon'$. Furthermore, let h_1, \dots, h_{n_2} be a subset of \mathcal{H} of cardinality $n_2 = N(\mathcal{H}, L^2(Q), \varepsilon')$ which is ε' -dense in \mathcal{H} in the $L^2(Q)$ -metric. It follows that, for $\varepsilon' < \min(\|H\|_\infty, \|N_{0,r}\|_\infty)$, every $\sum_{l \in \mathbb{Z}} N_{j,l,r}(y)H(2^j(\cdot) - l)$ is at $L^2(Q)$ -distance at most ε from $\sum_{i=1}^{n_1} \alpha_{i(l)} h_{i'(l)}$ for some $1 \leq i(l) \leq n_1$ and $1 \leq i'(l) \leq n_2$. The total number of such linear combinations is dominated by $(n_1 n_2)^r \leq (B'/\varepsilon)^{(w+1)r}$. This shows that the lemma holds for $\varepsilon < 2r \min\{\|H\|_\infty, \|N_{0,r}\|_\infty\} \max\{\|H\|_\infty, \|N_{0,r}\|_\infty\} = 2r\|H\|_\infty\|N_{0,r}\|_\infty = U$, which completes the proof by taking $A = \max(B', U, e)$ (for $\varepsilon \in [U, A]$ one ball covers the whole set). ■

Proposition 3 *Let K be as in (1) and assume either that ϕ has compact support and is of bounded p -variation ($p < \infty$), or that ϕ is a Battle-Lemarié scaling function for some $r \geq 1$. Suppose P has a bounded density p_0 . Given $C, T > 0$, there exist finite positive constants $C_1 = C_1(C, K, \|p_0\|_\infty)$ and $C_2 = C_2(C, T, K, \|p_0\|_\infty)$ such that, if*

$$\frac{n}{2^j j} \geq C \quad \text{and} \quad C_1 \sqrt{\frac{2^j j}{n}} \leq t \leq T$$

then

$$\Pr \left\{ \sup_{y \in \mathbb{R}} |p_n(y, j) - Ep_n(y, j)| \geq t \right\} \leq \exp \left(-C_2 \frac{nt^2}{2^j} \right). \quad (32)$$

Proof. We first prove the Battle-Lemarié wavelet case. If $r > 1$, the function K is continuous (see the proof of Lemma 1), and therefore the supremum in (32) is over a countable set. That this is also true for $r = 1$ follows from Remark 1 in Giné and Nickl (2009b). We apply Proposition 1 and Lemma 2 to the supremum of the empirical process indexed by the classes of functions

$$\mathcal{K}_j := \{2^{-j}K_j(\cdot, y)/(2\|\Phi\|_\infty) : y \in \mathbb{R}\},$$

where Φ is a function majorizing K (as in (2)), so that \mathcal{K}_j is uniformly bounded by $1/2$. We next bound the second moments $E(2^{-2j}K_j^2(X, y))$. We have, using (2),

$$\begin{aligned} \int 2^{-2j}K_j^2(x, y)p_0(x)dx &\leq \int \Phi^2(|2^j(x - y)|)p_0(x)dx \\ &\leq 2^{-j} \int \Phi^2(|u|)p_0(y + 2^{-j}u)du \leq 2^{-j}\|p_0\|_\infty\|\Phi\|_2^2. \end{aligned} \quad (33)$$

Hence we may take $\sigma = \sqrt{2^{-j}\|\Phi\|_2^2\|p_0\|_\infty}/(2\|\Phi\|_\infty)$, and the result is then a direct consequence of Proposition 1, which applies by Lemma 2. For compactly supported wavelets, the same proof applies, using Lemma 2 (and Remark 1) in Giné and Nickl (2009b). ■

Proof. (Theorem 1) Using Lemma 2, the first two claims of the Theorem follow by the same proof as in Theorem 1 and Corollary 1, Giné and Nickl (2009b). For the bias term, we argue as in Theorem 8.1 in HKPT (1998) – using the fact that ϕ_r is $r - 1$ times differentiable – and obtain, for $p_0 \in C^t(\mathbb{R})$

$$|Ep_n(x) - p_0(x)| \leq 2^{-jt}\|p_0\|_{t,\infty}C \quad (34)$$

where $C := C(\Phi) = \int \Phi(|u|)|u|^t du$. ■

3.4 An exponential inequality for the distribution function of the linear estimator.

The quantity of interest in this subsection is the distribution function F_n^S of the linear projection estimator p_n from (6), more precisely, we will study the stochastic process

$$\sqrt{n}(F_n^S(s) - F(s)) = \sqrt{n} \int_{-\infty}^s (p_n(y, j) - p_0(y))dy, \quad s \in \mathbb{R}.$$

To prove a functional CLT for this process, it turns out that it is easier to compare F_n^S to F_n rather than to F . With $\mathcal{F} = \{1_{(-\infty, s]} : s \in \mathbb{R}\}$, the decomposition

$$(F_n^S - F_n)(s) = (P_n - P)(\pi_j(f) - f) + \int (\pi_j(p_0) - p_0)f, \quad f \in \mathcal{F}, \quad (35)$$

will be useful, since it splits the quantity of interest into a deterministic ‘bias’ term and an empirical process.

Lemma 3 *Assume that p_0 is a bounded function ($t = 0$), or that $p_0 \in C^t(\mathbb{R})$ for some $0 < t \leq r$. Let $\mathcal{F} = \{1_{(-\infty, s]} : s \in \mathbb{R}\}$. Then*

$$\left| \int_{\mathbb{R}} (\pi_j(p_0) - p_0)f \right| \leq C2^{-j(t+1)} \quad (36)$$

for some constant C depending only on r and $\|p_0\|_{t,\infty}$.

Proof. Let $\psi := \psi_r$ be the mother wavelet associated to ϕ_r . Since the wavelet series of $p_0 \in L^1(\mathbb{R})$ converges in $L^1(\mathbb{R})$ we have $\pi_j(p_0) - p_0 = -\sum_{l=j}^{\infty} \sum_k \beta_{lk}(p_0)\psi_{lk}$ in the $L^1(\mathbb{R})$ -sense, and then, since $f = 1_{(-\infty, s]} \in L^\infty(\mathbb{R})$,

$$-\int_{\mathbb{R}} (\pi_j(p_0) - p_0)f = \int_{\mathbb{R}} \left(\sum_{l=j}^{\infty} \sum_k \beta_{lk}(p_0)\psi_{lk}(x) \right) f(x)dx = \sum_{l=j}^{\infty} \sum_k \beta_{lk}(p_0)\beta_{lk}(f).$$

The lemma now follows from an estimate for the decay of the wavelet coefficients of p_0 and f , namely the bounds

$$\sup_{f \in \mathcal{F}} \sum_k |\beta_{lk}(f)| \leq c 2^{-l/2} \quad \text{and} \quad \sup_k |\beta_{lk}(p_0)| \leq c' 2^{-l(t+1/2)}, \quad (37)$$

The first bound is proved as in the proof of Lemma 3 in Giné and Nickl (2009b), noting that the identity before equation (37) in that proof also holds for spline wavelets by their exponential decay property. The second bound follows from

$$\begin{aligned} \sup_k |\beta_{lk}(p_0)| &\leq c'' 2^{-l/2} \|K_{l+1}(p_0) - K_l(p_0)\|_\infty \\ &\leq c'' 2^{-l/2} (\|K_l(p_0) - p_0\|_\infty + \|K_{l+1}(p_0) - p_0\|_\infty) \leq c' 2^{-l/2} 2^{-lt} \end{aligned}$$

where we used (9.35) in HKPT (1998) for the first inequality and (34) in the last. ■

To control the fluctuations of the stochastic term, one applies Talagrand's inequality to the empirical process indexed by the 'shrinking' classes of functions $\{\pi_j(f) - f : f \in \mathcal{F}\}$. These classes consist of differences of elements in \mathcal{F} and in

$$\mathcal{K}'_j := \left\{ \int_{-\infty}^t K_j(\cdot, y) dy : t \in \mathbb{R} \right\},$$

and we have to show that, for each j , this class satisfies the entropy condition (17). Again, for ϕ with compact support (and of finite p -variation), this result was proved in Lemma 2 in Giné and Nickl (2009b), and we extend it now to the Battle-Lemarié wavelets considered here.

Lemma 4 *Let \mathcal{K}'_j be as above where ϕ_r is a Battle-Lemarié wavelet for $r \geq 1$. Then there exist finite constants $A \geq e$ and $v \geq 2$ and independent of j such that*

$$\sup_Q N(\mathcal{K}'_j, L^2(Q), \varepsilon) \leq \left(\frac{A}{\varepsilon} \right)^v, \quad 0 < \varepsilon < 1,$$

where the supremum extends over all Borel probability measures on \mathbb{R} .

Proof. In analogy to the proof of Lemma 2, one can write

$$\int_{-\infty}^t K_j(\cdot, y) dy = \sum_{l \in \mathbb{Z}} \int_{-\infty}^t 2^j N_{j,l,r}(y) dy H(2^j(\cdot) - l),$$

since the series (31) converges absolutely (in view of

$$\sum_l |H(2^j x - l)| \leq \sum_k |g(|k|)| \sum_l N_{k,r}(2^j x - l) \leq r \|N_{0,r}\|_\infty \sum_k |g(|k|)| < \infty.)$$

Recall that $N_{j,l,r}$ is supported in the interval $[2^{-j}l, 2^{-j}(r+l)]$. Hence, if $l > 2^j t$, the last integral is zero. For $l \leq 2^j t - r$, the integral equals the constant $c = \int_{\mathbb{R}} N_{0,r}(y) dy$, and for $l \in [2^j t - r, 2^j t]$, the integral $c_{j,l,r}$ is bounded by c , so that this sum in fact equals

$$c \sum_{l \leq 2^j t - r} H(2^j(\cdot) - l) + \sum_{2^j t - r < l < 2^j t} c_{j,l,r} H(2^j(\cdot) - l).$$

The second sum is contained in the set \mathcal{H}_r from the proof of Lemma 2, which satisfies the required entropy bound independent of j . For the first sum, decompose H into its positive and negative

part, so that the two resulting collections of functions are linearly ordered (in t) by inclusion, and hence are VC-subgraph of index 1, see Theorems 4.2.6 and 4.8.1 in Dudley (1999). Moreover, we can take the envelope $r\|N_{0,r}\|_\infty \sum_k |g(|k|)|$ independent of j . Combining entropy bounds, this proves the lemma. ■

Combining these observations, one can prove the following inequality, which parallels Theorem 1 in Giné and Nickl (2009a) for the classical kernel density estimator, and Lemma 4 in Giné and Nickl (2009b) for the wavelet density estimator (with ϕ compactly supported).

Proposition 4 *Let $F_n(s) = \int_{-\infty}^s dP_n$ and $F_n^S(s) := F_n^S(s, j) = \int_{-\infty}^s p_n(y, j) dy$, where p_n is as in (6). Assume that the density p_0 of P is a bounded function ($t = 0$) or that $p_0 \in C^t(\mathbb{R})$ for some t , $0 < t \leq r$. Let $j \in \mathbb{Z}$ satisfy $2^{-j} \geq d(\log n/n)$ for some $0 < d < \infty$. Then there exist finite positive constants $L := L(\|p_0\|_\infty, K, d)$, $\Lambda_0 := \Lambda_0(\|p_0\|_{t,\infty}, K, d)$ such that for all $n \in \mathbb{N}$ and $\lambda \geq \Lambda_0 \max(\sqrt{j}2^{-j}, \sqrt{n}2^{-j(t+1)})$ we have*

$$\Pr(\sqrt{n}\|F_n^S - F_n\|_\infty > \lambda) \leq L \exp\left\{-\frac{\min(2^j \lambda^2, \sqrt{n} \lambda)}{L}\right\}.$$

Proof. Given the preceding lemmas, the proposition follows from Talagrand's inequality applied to the class $\{\pi_j(1_{(-\infty, x]} - 1_{(-\infty, x]})\}$ in the same way as in the proof of Lemma 4 in Giné and Nickl (2009b), so we omit it. ■

3.5 Proof of Theorem 3

We can now prove the main result, Theorem 3. We will prove it only for Battle-Lemarié wavelets. For compactly supported wavelets, the proof is exactly the same, replacing the results from steps I)-II) below and from Sections 3.3 and 3.4 for spline wavelets by the corresponding ones for compactly supported wavelets obtained in Giné and Nickl (2009b). Also, uniformity in p_0 – which is proved by controlling the respective constants – is left implicit in the derivations. We start with some preliminary observations.

I) Since, uniformly in $j \in \mathcal{J}$, we have $n/(2^j j) > c \log n$ for some $c > 0$ independent of n , we have from Proposition 3, integrating tail probabilities, that

$$E\|p_n(j) - Ep_n(j)\|_\infty^p \leq D^p \left(\frac{2^j j}{n}\right)^{p/2} := D^p \sigma^p(j, n) \quad (38)$$

for every $j \in \mathcal{J}$, $1 \leq p < \infty$ and some $0 < D < \infty$ depending only on $\|p_0\|_\infty$ and Φ .

For the bias, we recall from (34) that, for $0 < t \leq r$

$$|Ep_n(y, j) - p_0(y)| \leq 2^{-jt} \|p_0\|_{t,\infty} C(\Phi) := B(j, p_0). \quad (39)$$

If the density p_0 is only uniformly continuous, then one still has from (2) and integrability of Φ that, uniformly in $y \in \mathbb{R}$,

$$|Ep_n(y, j) - p_0(y)| \leq \left| \int |\Phi(|u|)| |p_0(y - 2^{-j}u) - p_0(y)| du \right| := B(j, p_0) = o(1). \quad (40)$$

II) Define $\tilde{M} := \tilde{M}_n = C\|p_n(j_{\max})\|_\infty$ and set $C = 49\|\Phi\|_2^2$. Define also $M = C\|p_0\|_\infty$ for the same C . We need to control the probability that $\tilde{M} > 1.01M$ or $\tilde{M} < 0.99M$ if p_0 is uniformly

continuous. For some $0 < L < \infty$ and n large enough we have

$$\begin{aligned}
\Pr\left(|\tilde{M} - M| > 0.01C\|p_0\|_\infty\right) &= \Pr(|\|p_n(j_{\max})\|_\infty - \|p_0\|_\infty| > 0.01\|p_0\|_\infty) \\
&\leq \Pr(\|p_n(j_{\max}) - p_0\|_\infty > 0.01\|p_0\|_\infty) \\
&\leq \Pr(\|p_n(j_{\max}) - Ep_n(j_{\max})\|_\infty > 0.01\|p_0\|_\infty - B(j_{\max}, p_0)) \\
&\leq \Pr(\|p_n(j_{\max}) - Ep_n(j_{\max})\|_\infty > 0.009\|p_0\|_\infty) \\
&\leq \exp\left\{-\frac{(\log n)^2}{L}\right\}
\end{aligned}$$

by Proposition 3 and Step I). Furthermore, there exists a constant L' such that $E\tilde{M} \leq L'$ for every n in view of

$$E\|p_n(j_{\max})\|_\infty \leq E\|p_n(j_{\max}) - Ep_n(j_{\max})\|_\infty + \|Ep_n(j_{\max})\|_\infty \leq c + \|\Phi\|_1\|p_0\|_\infty,$$

where we have used (2) and (38).

III) We need some observations on the Rademacher processes used in the definition of \hat{j}_n . First, for the symmetrized empirical measure $\tilde{P}_n = 2n^{-1} \sum_{i=1}^n \varepsilon_i \delta_{X_i}$, we have

$$R(n, j) = \|\pi_j(\tilde{P}_n)\|_\infty = \|\pi_j(\pi_l(\tilde{P}_n))\|_\infty \leq \|\pi_j\|'_\infty R(n, l) \leq B(\phi)R(n, l) \quad (41)$$

for every $l > j$: Here $\|\pi_j\|'_\infty$ is the operator norm in $L^\infty(\mathbb{R})$ of the projection π_j , which admits bounds $B(\phi)$ independent of j . (Clearly, π_j acts on finite signed measures μ by duality, taking values in $L^\infty(\mathbb{R})$ since $|\pi_j(\mu)| = |\int K_j(\cdot, y)d\mu(y)| \leq 2^j\|\Phi\|_\infty|\mu|(\mathbb{R})$.) See Remark 3 for details on how to obtain $B(\phi)$. Furthermore, for $j < l$,

$$T(n, j, l) \leq R(n, j) + R(n, l) \leq (1 + B(\phi))R(n, l), \quad (42)$$

and the same inequality holds for the Rademacher expectations of $T(n, j, l)$. We also record the following bound for the (full) expectation of $R(n, l)$, $l \in \mathcal{J}$: Using inequality (27) and the variance computation (33), we have that there exists a constant L depending only on $\|p_0\|_\infty$ and Φ such that, for every $l \in \mathcal{J}$, $ER(n, l) \leq L\sqrt{2^l}/n$.

Proof of (11). Let $\mathcal{F} = \{1_{(-\infty, s]} : s \in \mathbb{R}\}$, and let $f \in \mathcal{F}$. We have

$$\sqrt{n} \int (p_n(\hat{j}_n) - p_0)f = \sqrt{n} \int (p_n(j_{\max}) - p_0)f + \sqrt{n} \int (p_n(\hat{j}_n) - p_n(j_{\max}))f.$$

The first term satisfies the CLT from Theorem 2 for the linear estimator with $j_n = j_{\max}$. We now show that the second term converges to zero in probability. Observe first

$$p_n(\hat{j}_n)(y) - p_n(j_{\max})(y) = P_n(K_{\hat{j}_n}(\cdot, y) - K_{j_{\max}}(\cdot, y)) = - \sum_{l=\hat{j}_n}^{j_{\max}-1} \sum_k \hat{\beta}_{lk} \psi_{lk}(y),$$

with convergence in $L^1(\mathbb{R})$. Next, we have by (9.35) in HKPT (1998), for all $l \in [\hat{j}_n, j_{\max} - 1]$ and all k , by definition of \hat{j}_n , that for some $0 < D' < \infty$

$$\begin{aligned}
(1/D')2^{l/2}|\hat{\beta}_{lk}| &\leq \sup_{y \in \mathbb{R}} |P_n(K_{l+1}(\cdot, y)) - P_n(K_l(\cdot, y))| = \|p_n(l+1) - p_n(l)\|_\infty \\
&\leq \|p_n(l+1) - p_n(\hat{j}_n)\|_\infty + \|p_n(l) - p_n(\hat{j}_n)\|_\infty \\
&\leq (1 + B(\phi))(R(n, l+1) + R(n, l)) + 3\sqrt{\tilde{M}2^l/n},
\end{aligned}$$

in case $\hat{j}_n = \bar{j}_n$ using also the inequality $T(n, \bar{j}_n, l) \leq (1 + B(\phi))R(n, l)$ for $l \geq \bar{j}_n$ (see (42)). Consequently, uniformly in $f \in \mathcal{F}$,

$$\begin{aligned} E \left| \int (p_n(\hat{j}_n) - p_n(j_{\max}))f \right| &= E \left| \sum_{l=\hat{j}_n}^{j_{\max}-1} \sum_k \hat{\beta}_{lk} \int \psi_{lk}(y) f(y) dy \right| \\ &\leq E \sum_{l=j_{\min}}^{j_{\max}-1} D' 2^{-l/2} \left((B(\phi) + 1)(R(n, l+1) + R(n, l)) + 3\sqrt{\tilde{M} 2^l/n} \right) \sum_k |\beta_{lk}(f)| \\ &\leq \left(\frac{D''}{\sqrt{n}} \right) \sum_{l=j_{\min}}^{j_{\max}-1} 2^{-l/2} \sqrt{l} = o\left(\frac{1}{\sqrt{n}} \right) \end{aligned}$$

using the moment bounds in II), III), $\hat{j}_n \geq j_{\min} \rightarrow \infty$ as $n \rightarrow \infty$ (by definition of \mathcal{J}) and since $\sup_{f \in \mathcal{F}} \sum_k |\beta_{lk}(f)| \leq c 2^{-l/2}$ by (37) for some constant c .

Proof of (12) and (13): The proof of the case $t = 0$ follows from a simple modification of the arguments below as in Theorem 2 in Giné and Nickl (2009a), so we omit it. [In this case, one defines j^* as j_{\max} if $t = 0$ so that only the case $\hat{j}_n \leq j^*$ has to be considered.] For $t > 0$, define $j^* := j(p_0)$ by the balance equation

$$j^* = \min \left\{ j \in \mathcal{J} : B(j, p_0) \leq \sqrt{2 \log 2} \|p_0\|_{\infty}^{1/2} \|\Phi\|_2 \sigma(j, n) \right\}. \quad (43)$$

Using the results from I), it is easily verified that $2^{j^*} \simeq (n/\log n)^{\frac{1}{2t+1}}$ if $p_0 \in \mathcal{C}^t(\mathbb{R})$ for some $0 < t \leq r$, and that

$$\sigma(j^*, n) = O \left(\left(\frac{\log n}{n} \right)^{t/(2t+1)} \right)$$

is the rate of convergence required in (13).

We will consider the cases $\{\hat{j}_n \leq j^*\}$ and $\{\hat{j}_n > j^*\}$ separately. First, if \hat{j}_n is \bar{j}_n , then we have by the definition of \bar{j}_n , (42), the definitions of M and j^* , (38) and the moment bound in III)

$$\begin{aligned} &E \|p_n(\bar{j}_n) - p_0\|_{\infty} I_{\{\bar{j}_n \leq j^*\} \cap \{\tilde{M} \leq 1.01M\}} \\ &\leq E (\|p_n(\bar{j}_n) - p_n(j^*)\|_{\infty} + E \|p_n(j^*) - p_0\|_{\infty}) I_{\{\bar{j}_n \leq j^*\} \cap \{\tilde{M} \leq 1.01M\}} \\ &\leq (B(\phi) + 1)ER(n, j^*) + \sqrt{1.01M} \sigma(j^*, n) + \|p_n(j^*) - p_0\|_{\infty} \\ &\leq B' \sqrt{\frac{2^{j^*} j^*}{n}} + B'' \sigma(j^*, n) = O(\sigma(j^*, n)), . \end{aligned} \quad (44)$$

If \hat{j}_n is \tilde{j}_n then one has the same bound (without even using (42)).

Also, by the results in I), II),

$$\begin{aligned} &E \|p_n(\hat{j}_n) - p_0\|_{\infty} I_{\{\hat{j}_n \leq j^*\} \cap \{\tilde{M} > 1.01M\}} \\ &\leq \sum_{j \in \mathcal{J}: j \leq j^*} E \left([\|p_n(j) - Ep_n(j)\|_{\infty} + B(j, p_0)] I_{\{\hat{j}_n = j\}} I_{\{\tilde{M} > 1.01M\}} \right) \\ &\leq c \log n [D\sigma(j^*, n) + B(j_{\min}, p_0)] \cdot \sqrt{E 1_{\{\tilde{M} > 1.01M\}}} \\ &= o \left((\log n) \sqrt{\exp \left\{ -\frac{(\log n)^2}{L} \right\}} \right) = o(\sigma(j^*, n)). \end{aligned}$$

We now turn to $\{\hat{j}_n > j^*\}$. First,

$$\begin{aligned}
& E \left\| p_n(\hat{j}_n) - p_0 \right\|_{\infty} I_{\{\hat{j}_n > j^*\} \cap \{\tilde{M} < 0.99M\}} \\
& \leq \sum_{j \in \mathcal{J}: j > j^*} E \left(\left[\|p_n(j) - Ep_n(j)\|_{\infty} + B(j, p_0) \right] I_{\{\hat{j}_n = j\}} I_{\{\tilde{M} < 0.99M\}} \right) \\
& \leq c' \log n [D\sigma(j_{\max}, n) + B(j^*, p_0)] \cdot \sqrt{EI_{\{\tilde{M} < 0.99M\}}} \\
& = O \left(\sqrt{(\log n) \exp \left\{ -\frac{(\log n)^2}{L} \right\}} \right) = o(\sigma(j^*, n)),
\end{aligned}$$

again by the results in I), II), and second, for any $1 < p < \infty$, $1/p + 1/q = 1$, using (38) and the definition of j^*

$$\begin{aligned}
& E \left\| p_n(\hat{j}_n) - p_0 \right\|_{\infty} I_{\{\hat{j}_n > j^*\} \cap \{0.99M \leq \tilde{M}\}} \\
& \leq \sum_{j \in \mathcal{J}: j > j^*} (E \|p_n(j) - p_0\|_{\infty}^p)^{1/p} \left(EI_{\{\hat{j}_n = j\} \cap \{0.99M \leq \tilde{M}\}} \right)^{1/q} \\
& \leq \sum_{j \in \mathcal{J}: j > j^*} D'\sigma(j, n) \cdot \Pr \left(\{\hat{j}_n = j\} \cap \{0.99M \leq \tilde{M}\} \right)^{1/q}.
\end{aligned}$$

We show below that for n large enough, some constant c , some $\delta > 0$ and some $q > 1$,

$$\Pr(\{\hat{j}_n = j\} \cap \{0.99M \leq \tilde{M}\}) \leq c2^{-j(q/2+\delta)}, \quad (45)$$

which gives the bound

$$\sum_{j \in \mathcal{J}: j > j^*} D''\sigma(j, n) \cdot 2^{-j/2-j\delta/q} = O \left(\frac{1}{\sqrt{n}} \right) = o(\sigma(j^*, n)),$$

completing the proof, modulo verification of (45).

To verify (45), we split the proof into two cases. Pick any $j \in \mathcal{J}$ so that $j > j^*$ and denote by j^- the previous element in the grid (i.e. $j^- = j - 1$).

Case I, $\hat{j}_n = \bar{j}_n$: One has

$$\Pr(\{\bar{j}_n = j\} \cap \{0.99M \leq \tilde{M}\}) \leq \sum_{l \in \mathcal{J}: l \geq j} \Pr \left(\|p_n(j^-) - p_n(l)\|_{\infty} > T(n, j^-, l) + \sqrt{0.99M}\sigma(l, n) \right).$$

We first observe that

$$\|p_n(j^-) - p_n(l)\|_{\infty} \leq \|p_n(j^-) - p_n(l) - Ep_n(j^-) + Ep_n(l)\|_{\infty} + B(j^-, p_0) + B(l, p_0), \quad (46)$$

where, setting $\sqrt{2 \log 2} \|p_0\|_{\infty}^{1/2} \|\Phi\|_2 =: U(p_0, \Phi)$,

$$B(j^-, p_0) + B(l, p_0) \leq 2B(j^*, p_0) \leq 2U(p_0, \Phi)\sigma(j^*, n) \leq 2U(p_0, \Phi)\sigma(l, n)$$

by definition of j^* and since $l > j^- \geq j^*$. Consequently, the l -th probability in the last sum is bounded by

$$\Pr \left(\|p_n(j^-) - p_n(l) - Ep_n(j^-) + Ep_n(l)\|_{\infty} > T(n, j^-, l) + (\sqrt{0.99M} - 2U(p_0, \Phi))\sigma(l, n) \right), \quad (47)$$

and we now apply Corollary 1 to this bound. Define the class of functions

$$\mathcal{F} := \mathcal{F}_{j^-, l} = \{2^{-l}(K_{j^-}(\cdot, y) - K_l(\cdot, y))/(4\|\Phi\|_\infty)\},$$

which is uniformly bounded by $1/2$, and satisfies (17) for some A and v independent of l and j^- by Lemma 2 (and a computation on covering numbers). We compute σ , using (33) and $l > j^-$:

$$\begin{aligned} 2^{-2l}E(K_{j^-} - K_l)(X, y)^2 &\leq 2^{-2l+1} \left(EK_{j^-}^2(X, y) + EK_l^2(X, y) \right) \\ &\leq 2^{-2l+1} \|\Phi\|_2^2 \|p_0\|_\infty (2^{j^-} + 2^l) \leq 3 \cdot 2^{-l} \|\Phi\|_2^2 \|p_0\|_\infty \end{aligned}$$

so that we can take $\sigma^2 = 3 \cdot 2^{-l} \|\Phi\|_2^2 \|p_0\|_\infty / (16\|\Phi\|_\infty^2)$. Then the probability in (47) is equal to

$$\begin{aligned} \Pr \left(\frac{2^l 4 \|\Phi\|_\infty}{n} \left\| \sum_{i=1}^n f(X_i) - Pf \right\|_{\mathcal{F}} > \frac{2^l 4 \|\Phi\|_\infty}{n} 2 \left\| \sum_{i=1}^n \varepsilon_i f(X_i) \right\|_{\mathcal{F}} + (\sqrt{0.99M} - 2U(p_0, \Phi))\sigma(l, n) \right) \\ = \Pr \left(\left\| \sum_{i=1}^n f(X_i) - Pf \right\|_{\mathcal{F}} > 2 \left\| \sum_{i=1}^n \varepsilon_i f(X_i) \right\|_{\mathcal{F}} + 3 \frac{n(\sqrt{0.99M} - 2U(p_0, \Phi))\sigma(l, n)}{3 \cdot 2^l \cdot 4 \|\Phi\|_\infty} \right). \end{aligned}$$

Since $n\sigma^2/\log(1/\sigma) \simeq n/(2^l) \rightarrow \infty$ uniformly in $l \in \mathcal{J}$, there exists $\lambda_n \rightarrow \infty$ independent of l such that (19) is satisfied, and the choice

$$t = \frac{n(\sqrt{0.99M} - 2U(p_0, \Phi))\sigma(l, n)}{3 \cdot 2^l \cdot 4 \|\Phi\|_\infty}$$

is admissible in Corollary 1 for $c_2(\lambda_n) = 1 + 120\lambda_n^{-1} + 10800\lambda_n^{-2}$. Hence, using Corollary 1, the last probability is bounded by

$$\leq 2 \exp \left(- \frac{n^2(\sqrt{0.99M} - 2U(p_0, \Phi))^2 (2^l l / n) 16 \|\Phi\|_\infty^2}{9 \cdot 6.3 \cdot c_2(\lambda_n) 2^{2l} n 2^{-l} \|\Phi\|_2^2 \|p_0\|_\infty 16 \|\Phi\|_\infty^2} \right) \leq 2^{-l((q/2)+\delta)} \quad (48)$$

for some $\delta > 0$ and $q > 1$, by definition of M . Since $\sum_{l \in \mathcal{J}: l \geq j} 2^{-l((q/2)+\delta)} \leq c 2^{-j((q/2)+\delta)}$, we have proved (45).

Case II, $\hat{j}_n = \tilde{j}_n$: The proof reduces to the previous case since, by inequality (42), one has

$$\begin{aligned} &\Pr(\{\tilde{j}_n^\varepsilon = j\} \cap \{0.99M \leq \tilde{M}\}) \\ &\leq \sum_{l \in \mathcal{J}: l \geq j} \Pr \left(\|p_n(j^-) - p_n(l)\|_\infty > (B(\phi) + 1)R(n, l) + \sqrt{0.99M}\sigma(l, n) \right) \\ &\leq \sum_{l \in \mathcal{J}: l \geq j} \Pr \left(\|p_n(j^-) - p_n(l)\|_\infty > T(n, j^-, l) + \sqrt{0.99M}\sigma(l, n) \right). \end{aligned}$$

Acknowledgement We thank Patricia Reynaud-Bouret and Benedikt Pötscher for helpful comments. The idea of using Rademacher thresholds in Lepski's method arose from a conversation with Patricia Reynaud-Bouret.

References

- [1] BARRON, A., BIRGÉ, L., MASSART, P. (1999). Risk bounds for model selection via penalization. *Probab. Theory Related Fields* **113** 301-413.

- [2] BARTLETT, P., BOUCHERON, S. and LUGOSI, G. (2002). Model selection and error estimation. *Mach. Learn.* **48** 85-113.
- [3] BOUSQUET, O. (2003). Concentration inequalities for sub-additive functions using the entropy method. In: *Stochastic inequalities and applications.*, *Progr. Probab.* **56**, E. Giné, C. Houdré, D. Nualart, eds., Birkhäuser, Boston, 213-247.
- [4] DAUBECHIES, I. (1992). *Ten lectures on wavelets*. CBMS-NSF Reg. Conf. Ser. in Appl. Math. 61. Philadelphia, Society for Industrial and Applied Mathematics.
- [5] DEHEUVELS, P (2000). Uniform limit laws for kernel density estimators on possibly unbounded intervals. In: N. Limnios and M. Nikulin (Eds.), *Recent advances in reliability theory*, Birkhäuser, Boston, pp. 477-492.
- [6] DEVORE, R.A. and LORENTZ, G.G. (1993). *Constructive approximation*. Springer, Berlin.
- [7] DONOHO, D. L.; JOHNSTONE, I. M.; KERKYACHARIAN, G. and PICARD, D. (1996). Density estimation by wavelet thresholding. *Ann. Statist.* **24** 508-539.
- [8] DUDLEY, R.M. (1999). *Uniform central limit theorems*. Cambridge University Press; Cambridge, England.
- [9] EINMAHL, U. and MASON, D. M. (2000). An empirical process approach to the uniform consistency of kernel-type function estimators. *J. Theoret. Probab.* **13** 1-37.
- [10] GINÉ, E. and GUILLOU, A. (2001). On consistency of kernel density estimators for randomly censored data: rates holding uniformly over adaptive intervals. *Ann. Inst. H. Poincaré Probab. Statist.* **37** 503-522.
- [11] GINÉ, E. and GUILLOU, A. (2002). Rates of strong uniform consistency for multivariate kernel density estimators. *Ann. Inst. H. Poincaré Probab. Statist.* **38** 907-921.
- [12] GINÉ, E. and KOLTCHINSKII, V. (2006). Concentration inequalities and asymptotic results for ratio type empirical processes. *Ann. Probab.* **34** 1143-1216.
- [13] GINÉ, E. and NICKL, R. (2009a). An exponential inequality for the distribution function of the kernel density estimator, with applications to adaptive estimation. *Probab. Theory Related Fields* **143** 569-596.
- [14] GINÉ, E. and NICKL, R. (2009b). Uniform limit theorems for wavelet density estimators. *Ann. Probab.* **37** 1605-1646.
- [15] GOLDENSHLUGER, A. and LEPSKI, O. (2009). Structural adaptation via L_p -norm oracle inequalities. *Probab. Theory Related Fields* **143** 41-71.
- [16] GOLUBEV, Y., LEPSKI, O. and LEVIT, B. (2001). On adaptive estimation for the sup-norm losses. *Math. Methods Statist.* **10** 23-37.
- [17] HÄRDLE, W.; KERKYACHARIAN, G., PICARD, D. and TSYBAKOV, A. (1998). *Wavelets, approximation, and statistical applications*. Lecture Notes in Statistics **129**. Springer, New York.
- [18] HUANG, S.-Y. (1999). Density estimation by wavelet-based reproducing kernels. *Statist. Sinica* **9** 137-151.

- [19] HUANG, S.-Y. and STUDDEN, W. J. (1993). Density estimation using spline projection kernels. *Comm. Statist. Theory Methods* **22** 3263-3285.
- [20] KERKYACHARIAN, G., PICARD, D. (1992). Density estimation in Besov spaces. *Statist. Probab. Lett.* **13** 15-24.
- [21] KLEIN, T. and RIO, E. (2005). Concentration around the mean for maxima of empirical processes. *Ann. Probab.* **33** 1060-1077.
- [22] KOLTCHINSKII, V. (2001). Rademacher penalties and structural risk minimization. *IEEE Trans. Inform. Theory* **47** 1902-1914.
- [23] KOLTCHINSKII, V. (2006). Local Rademacher complexities and oracle inequalities in risk minimization. *Ann. Statist.* **34** 2593-2656.
- [24] KOROSTELEV, A. and NUSSBAUM, M. (1999). The asymptotic minimax constant for sup-norm loss in nonparametric density estimation. *Bernoulli* **5** 1099-1118.
- [25] LEDOUX, M. (2001) *The concentration of measure phenomenon*. Mathematical Surveys and Monographs, 89. American Mathematical Society.
- [26] LEPSKI, O.V. (1991). Asymptotically minimax adaptive estimation. I. Upper bounds. Optimally adaptive estimates. *Theory Probab. Appl.* **36** 682-697.
- [27] LEPSKI, O.V. (1992). On problems of adaptive estimation in white Gaussian noise. In: *Topics in nonparametric estimation* (R.Z. Khasminksi, ed.) 87-106. Amer. Math. Soc., Providence.
- [28] MASSART, P. (2000) About the constants in Talagrand's concentration inequalities for empirical processes. *Ann. Probab.* **28** 863-884.
- [29] MEYER, Y. (1992). *Wavelets and operators I*. Cambridge University Press; Cambridge, England.
- [30] NOLAN, D. and POLLARD, D. (1987). U -processes: rates of convergence. *Ann. Statist.* **15** 780-799.
- [31] SHADRIN, A. Y. (2001). The L_∞ -norm of the L_2 -spline projector is bounded independently of the knot sequence: a proof of de Boor's conjecture. *Acta Math.* **187** 59-137.
- [32] SCHUMAKER, L.L. (1993). *Spline functions: basic theory*. Correlated reprint of the 1981 original. Krieger, Malabar.
- [33] TALAGRAND, M. (1994). Sharper bounds for Gaussian and empirical processes. *Ann. Probab.* **22** 28-76.
- [34] TALAGRAND, M. (1996). New concentration inequalities in product spaces. *Invent. Math.* **126** 505-563.
- [35] TSYBAKOV, A.B. (1998). Pointwise and sup-norm sharp adaptive estimation of the functions on the Sobolev classes. *Ann. Statist.* **26** 2420-2469.