

# Modulation Estimators and Confidence Sets

Rudolf Beran\* and Lutz Dümbgen†

University of California, Berkeley and Universität Heidelberg

January 1996, Revised August 1997

**Abstract.** An unknown signal plus white noise is observed at  $n$  discrete time points. Within a large convex class of linear estimators of  $\xi$ , we choose the estimator  $\hat{\xi}$  that minimizes estimated quadratic risk. By construction,  $\hat{\xi}$  is nonlinear. This estimation is done after orthogonal transformation of the data to a reasonable coordinate system. The procedure adaptively tapers the coefficients of the transformed data. If the class of candidate estimators satisfies a uniform entropy condition, then  $\hat{\xi}$  is asymptotically minimax in Pinsker's sense over certain ellipsoids in the parameter space and shares one such asymptotic minimax property with the James-Stein estimator. We describe computational algorithms for  $\hat{\xi}$  and construct confidence sets for the unknown signal. These confidence sets are centered at  $\hat{\xi}$ , have correct asymptotic coverage probability, and have relatively small risk as set-valued estimators of  $\xi$ .

*AMS 1991 subject classifications.* Primary 62H12; secondary 62M10.

*Key words and phrases.* Adaptivity, asymptotic minimax, bootstrap, bounded variation, coverage probability, isotonic regression, orthogonal transformation, signal recovery, Stein's unbiased estimator of risk, tapering.

---

\*Research supported in part by National Science Foundation Grant DMS95-30492 and in part by Sonderforschungsbereich 373 at Humboldt-Universität zu Berlin.

†Research supported in part by European Union Human Capital and Mobility Program ERB CHRX-CT 940693.

# 1 Introduction

The problem of recovering a signal from observation of the signal plus noise may be formulated as follows. Let  $X = X_n = [X(t)]_{t \in T}$  be a random function observed on the set  $T = T_n = \{1, 2, \dots, n\}$ . The components  $X(t)$  are independent with  $\mathbb{E} X(t) = \xi(t) = \xi_n(t)$  and  $\text{Var}[X(t)] = \sigma^2$  for every  $t \in T$ . Working with functions on  $T$  rather than vectors in  $\mathbf{R}^n$  is very convenient for the present purposes. As just indicated, we will usually drop the subscript  $n$  for notational simplicity. The signal  $\xi$  and the noise variance  $\sigma^2$  are both unknown. For simplicity we assume throughout that  $X$  is Gaussian. Portions of the argument that hold for non-Gaussian  $X$  are expressed by the lemmas in Section 6.2.

For any  $g \in \mathbf{R}^T$ , the space of real-valued functions defined on  $T$ , let

$$\text{ave}(g) := n^{-1} \sum_{t \in T} g(t).$$

The loss of any estimator  $\hat{\xi}$  for  $\xi$  is defined to be

$$(1.1) \quad L(\hat{\xi}, \xi) := \text{ave}[(\hat{\xi} - \xi)^2]$$

and the corresponding risk of  $\hat{\xi}$  is

$$\rho(\hat{\xi}, \xi, \sigma^2) := \mathbb{E} L(\hat{\xi}, \xi).$$

The first goal is to devise an estimator that is efficient in terms of this risk. If  $\xi$  and  $X$  are electrical voltages, then  $\text{ave}(\xi^2)$  and  $L(\hat{\xi}, \xi)$  are the time-averaged powers dissipated in passing the signal  $\xi$  and the error  $\hat{\xi} - \xi$  through a unit resistance.

Any estimator  $\hat{\xi}$  of  $\xi$  is governed by the asymptotic minimax bound

$$(1.2) \quad \liminf_{n \rightarrow \infty} \inf_{\hat{\xi}} \sup_{\text{ave}(\xi^2) \leq c} \rho(\hat{\xi}, \xi, \sigma^2) \geq \frac{\sigma^2 c}{\sigma^2 + c}$$

for every positive  $c$  and  $\sigma^2$ . Inequality (1.2) follows from a more general bound proved by Pinsker (1980) for signal recovery in Gaussian noise (see Nussbaum 1996 and Section 2). It may also be derived from ideas in Stein (1956) by considering best orthogonally equivariant estimators in the submodel where  $\text{ave}(\xi^2) = c$  (see Beran 1996b). Let  $\hat{\sigma}^2 = \hat{\sigma}_n^2$  be an estimator of  $\sigma^2$  that is consistent as in display (2.2) of Section 2. Then

$$\hat{\xi}_S := [1 - \hat{\sigma}^2 / \text{ave}(X^2)]^+ X$$

is essentially the James-Stein (1961) estimator, where  $[\cdot]^+$  denotes the positive part function. It achieves the Pinsker bound (1.2) because

$$(1.3) \quad \lim_{n \rightarrow \infty} \sup_{\text{ave}(\xi^2) \leq c} \rho(\hat{\xi}_S, \xi, \sigma^2) = \frac{\sigma^2 c}{\sigma^2 + c}$$

for every positive  $c$  and  $\sigma^2$ . The limit (1.3) follows from Corollary 2.3 or from asymptotics in Casella and Hwang(1982). For the maximum likelihood estimator  $\hat{\xi}_{\text{ML}} = X$ , the risk is always  $\sigma^2$ , which is strictly greater than the Pinsker bound.

Section 2 of this paper constructs estimators of  $\xi$  that are asymptotically minimax over a variety of ellipsoids in the parameter space while achieving, in particular, the asymptotic minimax bound (1.2) for every  $c > 0$ . These *modulation* estimators take the form  $\hat{f}X = [\hat{f}(t)X(t)]_{t \in T}$ . Here  $\hat{f} : T \rightarrow [0, 1]$  depends on  $X$  and is chosen to minimize the estimated risk of the linear estimator  $fX$  over all functions  $f$  in a class  $\mathcal{F} = \mathcal{F}_n \subset [0, 1]^T$ . Many well-known estimators are of this form with special classes  $\mathcal{F}$ . In the present paper we analyze such estimators under rather general assumptions on  $\mathcal{F}$ . How large this class may be is at the heart of the analysis. Taking  $\mathcal{F}$  to be the set of all functions from  $T$  to  $[0, 1]$  leads to a poor modulation estimator. Successful is to let  $\mathcal{F}$  be a closed convex set of functions with well-behaved uniform covering numbers. One example is the set of all functions in  $[0, 1]^T$  that are nonincreasing. The asymptotic theory of such modulation estimators, including links with the literature, is the subject of Section 2. Section 4 develops algorithms for computing  $\hat{f}X$  in the example of  $\mathcal{F}$  just cited.

Section 3 constructs confidence sets that are centered at a modulation estimator  $\hat{f}X$  and have asymptotic coverage probability  $\alpha$  for  $\xi$ . The risk of the modulation estimator at the center is shown to determine the risk of the confidence set, when that is viewed as a set-valued estimator for  $\xi$ . In this manner, efficiency of a modulation estimator determines the efficiency of the associated confidence set.

Before estimation of  $\xi$ , the data  $X$  may be transformed orthogonally without changing its Gaussian character. A modulation estimator computed in the new coordinate system can be transformed back into the original coordinate system to yield an estimator of  $\xi$ . Standard choices for such preliminary orthogonal transformation include Fourier transforms, wavelet transforms, or analysis-of-variance transforms. When applied in this

manner, modulation estimators perform data-driven tapering of empirical Fourier, wavelet or analysis-of-variance coefficients. Section 5 includes numerical examples of modulation estimators and confidence bounds after Fourier transformation.

## 2 Modulation estimators

After defining modulation estimators, this section obtains uniform asymptotic approximations to their risks. Let  $\mathcal{F} = \mathcal{F}_n$  be a given subset of  $[0, 1]^T$ . Each function  $f \in \mathcal{F}$  is called a *modulator* and defines a candidate linear estimator  $fX = [f(t)X(t)]_{t \in T}$  for  $\xi$ . The risk of this candidate estimator under quadratic loss (1.1) is

$$(2.1) \quad \rho(fX, \xi, \sigma^2) = \mathbb{E} L(fX, \xi) = \text{ave}[\sigma^2 f^2 + \xi^2 (1 - f)^2].$$

For brevity, we will write  $R(f, \xi, \sigma^2)$  in place of  $\rho(fX, \xi, \sigma^2)$ .

We will first construct a suitably consistent estimator  $\widehat{R}(f)$  of this risk. Suppose that  $\widehat{\sigma}^2 = \widehat{\sigma}_n^2$  is an estimator of  $\sigma^2$ , constructed (for instance) by one of the methods described later. Let  $X^*$  be a bootstrap random vector in  $\mathbf{R}^T$  such that  $\mathcal{L}(X^* | X, \widehat{\sigma}^2) = \mathcal{N}_T(X, \widehat{\sigma}^2 I)$ . The corresponding bootstrap risk estimator for  $R(f, \xi, \sigma^2)$  is

$$\mathbb{E} \left( L(fX^*, X) \mid X, \widehat{\sigma}^2 \right) = R(f, X, \widehat{\sigma}^2).$$

We call  $R(f, X, \widehat{\sigma}^2)$  the *naive risk estimator* because it is badly biased upwards, even asymptotically. The key point is

$$\mathbb{E} R(f, X, \sigma^2) = \text{ave}[f^2 \sigma^2 + (1 - f)^2 (\xi^2 + \sigma^2)] = R(f, \xi, \sigma^2) + \text{ave}[(1 - f)^2 \sigma^2].$$

Two possible corrections to the naive risk estimator are

$$\begin{aligned} \widehat{R}_C(f) &:= \text{ave}[f^2 \widehat{\sigma}^2 + (1 - f)^2 (X^2 - \widehat{\sigma}^2)] = R(f, X, \widehat{\sigma}^2) - \text{ave}[(1 - f)^2 \widehat{\sigma}^2], \\ \widehat{R}_B(f) &:= \max \left\{ \text{ave}(f^2 \widehat{\sigma}^2), \widehat{R}_C(f) \right\} = \text{ave}(f^2 \widehat{\sigma}^2) + \text{ave}[(1 - f)^2 (X^2 - \widehat{\sigma}^2)]^+. \end{aligned}$$

Risk estimator  $\widehat{R}_C$  is essentially Mallows' (1973)  $C_L$  criterion or Stein's (1981) unbiased estimator of risk, with estimation of  $\sigma^2$  incorporated. Risk estimator  $\widehat{R}_B$  corrects the possible negativity in  $\text{ave}[(1 - f)^2 (X^2 - \widehat{\sigma}^2)]$  as an estimator for  $\text{ave}[(1 - f)^2 \xi^2]$ . Let  $X^*$

be a random vector in  $R^T$  such that  $\mathcal{L}(X^* | X, \hat{\sigma}^2)$  is  $\mathcal{N}_T(\hat{\xi}, \hat{\sigma}^2 I)$ , where  $\hat{\xi} = \hat{\xi}(X, \hat{\sigma}^2)$  is a vector such that  $\text{ave}(\hat{\xi}^2) = \text{ave}[(1-f)^2(X^2 - \hat{\sigma}^2)]^+ / \text{ave}[(1-f)^2 X^2]$ . Then the bootstrap risk estimator  $\mathbb{E}[L(fX^*, \hat{\sigma}^2) | X, \hat{\sigma}^2]$  is precisely  $\hat{R}_B$ .

Let  $\hat{R}$  denote either  $\hat{R}_C$  or  $\hat{R}_B$ . We propose to estimate  $\xi$  by the *modulation* estimator  $\hat{f}X$ , where  $\hat{f}$  is any function in  $\mathcal{F}$  that minimizes  $\hat{R}(f)$ . Unless stated otherwise it is assumed throughout that

$\mathcal{F}$  is a closed convex subset of  $[0, 1]^T$  containing all constants  $c \in [0, 1]$ .

Because both  $\hat{R}_C(\cdot)$  and  $\hat{R}_B(\cdot)$  are convex functions on  $[0, 1]^T$ , the minimizer  $\hat{f}$  over  $\mathcal{F}$  exists in each case. These minimizers are unique with probability one because  $\hat{R}_C(f)$  is strictly convex in  $f$  whenever  $X(t) \neq 0$  for every  $t \in T$ . Similarly, the risk function  $R(f, \xi, \sigma^2)$  defined through (2.1) is strictly convex over  $[0, 1]^T$ , with unique minimizer  $\tilde{f}$ .

*REMARK A.* The modulation estimator  $\hat{f}X$  behaves poorly when the class  $\mathcal{F}$  is too large. For instance, let  $\mathcal{F}$  be the class of all functions in  $[0, 1]^T$ . The minimizer of  $R(\cdot, \xi, \sigma^2)$  over  $[0, 1]^T$  is the “oracle” modulator (cf. Donoho and Johnstone 1994)

$$\tilde{g} := \xi^2 / (\sigma^2 + \xi^2),$$

the division being componentwise, while the minimizer of  $\hat{R}(\cdot)$  over  $\mathcal{F}$  is now the *greedy* modulator  $\hat{g}^+$ , where

$$\hat{g} := (X^2 - \hat{\sigma}^2) / X^2.$$

To simplify the discussion, suppose that  $\sigma^2$  is known and  $\hat{\sigma}^2 \equiv \sigma^2$ . Then the estimator  $\hat{g}^+X$  is of the general form  $\hat{\xi} := \left( S[X(t)] \right)_{t \in T}$  for some measurable function  $S$  on the line. Since the maximum likelihood estimator  $X$  is componentwise admissible, the risk function  $\rho(\hat{\xi}, \cdot, \sigma^2)$  of  $\hat{\xi}$  is either identical to  $\rho(X, \cdot, \sigma^2) \equiv \sigma^2$  or there is a real number  $\theta$  such that  $\int (\theta - S)^2 d\mathcal{N}(\theta, \sigma^2) > \sigma^2$ . Then, if  $\xi(\cdot) \equiv \theta$ ,

$$\rho(\hat{\xi}, \xi, \sigma^2) > \sigma^2 = \rho(X, \xi, \sigma^2) > \sigma^2 \theta^2 / (\sigma^2 + \theta^2),$$

the latter being the asymptotic risk of the James-Stein estimator  $\hat{\xi}_S$ . Thus, the maximum risk of  $\hat{g}^+X$  is worse than that of estimators achieving Pinsker’s asymptotic minimax bound (1.2) and is even worse than that of the naive estimator  $X$ .

It should be mentioned that greedy modulation can be made successful in some sense if one overestimates the variance  $\sigma^2$  systematically. Donoho and Johnstone (1994) propose threshold estimators of the form  $\hat{\xi} = (1 - \lambda_n \sigma / |X|)^+ X$  or  $\hat{\xi} = 1\{|X| \geq \lambda_n \sigma\} X$ , and prove that they have surprising optimality properties if  $\lambda_n = (2 \log n)^{1/2} (1 + \delta_n)$  with a suitable sequence  $(\delta_n)_n$  tending to zero. These estimators are similar to  $\hat{g}^+ X$  if  $\hat{g}$  is computed with  $\hat{\sigma}_n^2 := \lambda_n^2 \sigma^2$ . While showing good performance in case of “sparse signals”, these estimators do not achieve the Pinsker bound (1.2) or the minimax bounds in Corollary 2.3 below. Also, the construction of confidence bounds for their loss seems to be intractable. Section 5 illustrates the possibly poor performance of hard thresholding for non-sparse signals.

*REMARK B.* Kneip’s (1994) ordered linear smoothers are equivalent to certain modulation estimators computed after suitable orthogonal transformation of  $X$ . The conditions that we impose on  $\mathcal{F}$  in this paper are substantially weaker than the ordering of  $\mathcal{F}$  required by Kneip. Consequently, our results also apply to the ridge regression, spline estimation, and kernel estimation examples discussed in Kneip’s paper. The earlier paper of Li (1987) treated non-diagonal linear estimators indexed by a parameter  $h$ . Li’s optimality result may be compared with Theorem 2.1 below. However, it does not seem easy to relate Li’s conditions on the range of  $h$  to our conditions on  $\mathcal{F}$ . The latter conditions give access to empirical process results that yield asymptotic distributions for the loss of  $\hat{f}X$  and hence confidence sets for  $\xi$  centered at modulation estimators.

*REMARK C.* Nussbaum (1996) surveyed constructions of adaptive estimators that achieve Pinsker-type asymptotic minimax bounds. For instance, Golubev and Nussbaum (1992) treated adaptive, asymptotically minimax estimation when  $\xi_i = g(x_i)$  and  $g$  lies in an ellipsoid of unknown radius within a Sobolev space of unknown order. Corollary 2.3 below is of related character. However, our results make no smoothness assumptions on  $\xi$ . For instance, sample paths up to time  $n$  of suitably scaled, discrete-time, independent white noise ultimately lie, as  $n \rightarrow \infty$ , within the ball  $\text{ave}(\xi^2) \leq c$ .

Useful classes of modulators  $\mathcal{F}$  can be characterized through their uniform covering numbers, which are defined as follows. For any probability measure  $Q$  on  $T$ , consider the

pseudo-distance  $d_Q(f, g)^2 := \int (f - g)^2 dQ$  on  $[0, 1]^T$ . For every positive  $u$ , let

$$N(u, \mathcal{F}, d_Q) := \min \left\{ \#\mathcal{F}_o : \mathcal{F}_o \subset \mathcal{F}, \inf_{f_o \in \mathcal{F}_o} d_Q(f_o, f) \leq u \forall f \in \mathcal{F} \right\}.$$

Define the uniform covering number

$$N(u, \mathcal{F}) := \sup_Q N(u, \mathcal{F}, d_Q),$$

where the supremum is taken over all probabilities on  $T$ . Let

$$J(\mathcal{F}) := \int_0^1 \sqrt{\log N(u, \mathcal{F})} du.$$

Throughout  $C$  denotes a generic universal real constant which does not depend on  $n$ ,  $\xi$ ,  $\sigma^2$  or  $\mathcal{F}$ , but whose value may be different in various places.

**THEOREM 2.1** *Let  $\mathcal{F}$  be any closed subset of  $[0, 1]^T$  containing 0, let  $\tilde{f}$  be a minimizer of  $R(f, \xi, \sigma^2)$  over  $f \in \mathcal{F}$ , and let  $\hat{f}$  minimize either  $\hat{R}_C(f)$  or  $\hat{R}_B(f)$  over  $f \in \mathcal{F}$ . Then*

$$\mathbb{E} \left| G - R(\tilde{f}, \xi, \sigma^2) \right| \leq C \left( J(\mathcal{F}) \frac{\sigma^2 + \sigma \sqrt{\text{ave}(\xi^2)}}{\sqrt{n}} + \mathbb{E} |\hat{\sigma}^2 - \sigma^2| \right),$$

where  $G$  is any one of the following quantities:

$$L(\hat{f}X, \xi), \quad \inf_{f \in \mathcal{F}} L(fX, \xi), \quad \hat{R}_C(\hat{f}), \quad \hat{R}_B(\hat{f}).$$

In particular,

$$\left| \rho(\hat{f}X, \xi, \sigma^2) - R(\tilde{f}, \xi, \sigma^2) \right| \leq C \left( J(\mathcal{F}) \frac{\sigma^2 + \sigma \sqrt{\text{ave}(\xi^2)}}{\sqrt{n}} + \mathbb{E} |\hat{\sigma}^2 - \sigma^2| \right).$$

This theorem is about convergence of losses and risks. The next result uses convexity of  $\mathcal{F}$  to establish that  $\hat{f}$  and  $\tilde{f}$ , as well as  $\hat{f}X$  and  $\tilde{f}X$ , converge to one another. Note that the second bound holds uniformly in  $\xi \in \mathbf{R}^T$ .

**THEOREM 2.2** *Let  $\hat{f}$  be the minimizer of  $\hat{R}_C$ . Then*

$$\begin{aligned} \mathbb{E} \text{ave} \left( (\sigma^2 + \xi^2) (\hat{f} - \tilde{f})^2 \right) &\leq C J(\mathcal{F}) \frac{\sigma^2 + \sigma \sqrt{\text{ave}(\xi^2)}}{\sqrt{n}} + \mathbb{E} |\hat{\sigma}^2 - \sigma^2|, \\ \mathbb{E} \text{ave} \left( (\hat{f}X - \tilde{f}X)^2 \right) &\leq C J(\mathcal{F}) \frac{\sigma^2}{\sqrt{n}} + \mathbb{E} |\hat{\sigma}^2 - \sigma^2|. \end{aligned}$$

Given consistency of  $\hat{\sigma}$  and boundedness of  $\sigma^2 + \text{ave}(\xi^2)$ , a key assumption on  $\mathcal{F}$  that ensures success of the modulation estimator  $\hat{f}X$  defined above is that  $J(\mathcal{F}) = o(n^{1/2})$ . Here are some examples of modulator classes  $\mathcal{F}$  to which Theorem 2.1 applies.

*EXAMPLE 1 (Stein shrinkage).* Suppose that  $\mathcal{F}$  consists of all constant functions in  $[0, 1]^T$ . The minimizer over  $\mathcal{F}$  of  $R(f, \xi, \sigma^2)$  is

$$\tilde{f}_S \equiv 1 - \sigma^2 / [\sigma^2 + \text{ave}(\xi^2)].$$

The minimizer of both  $\hat{R}_C$  and  $\hat{R}_B$  is

$$\hat{f}_S \equiv [1 - \hat{\sigma}^2 / \text{ave}(X^2)]^+.$$

The resulting modulation estimator  $\hat{f}_S X$  is the (modified) James-Stein (1981) estimator  $\hat{\xi}_S$ . Here one easily shows that  $N(u, \mathcal{F}) \leq 1 + (2u)^{-1}$ , whence  $J(\mathcal{F})$  is bounded by a universal constant.

*EXAMPLE 2 (Multiple Stein shrinkage).* Let  $\mathcal{B} = \mathcal{B}_n$  be a partition of  $T$  and define

$$\mathcal{F} := \left\{ \sum_{B \in \mathcal{B}} 1_B c(B) : c \in [0, 1]^{\mathcal{B}} \right\},$$

where  $1_B$  is the indicator function of  $B$ . The values of  $c(B)$  that define  $\tilde{f}$  and  $\hat{f}$ , respectively, are

$$\begin{aligned} \tilde{c}(B) &= \text{ave}(1_B \xi^2) / \text{ave}[1_B(\sigma^2 + \xi^2)], \\ \hat{c}(B) &= \text{ave}[1_B(X^2 - \hat{\sigma}^2)]^+ / \text{ave}(1_B X^2). \end{aligned}$$

The modulation estimator  $\hat{f}X$  now has the asymptotic form of the multiple shrinkage estimator in Stein (1966). Elementary calculations show that  $N(u, \mathcal{F}) \leq [1 + (2u)^{-1}]^{\#\mathcal{B}}$ . Thus  $J(\mathcal{F})$  is bounded by a universal constant times  $(\#\mathcal{B})^{1/2}$ , so that  $J(\mathcal{F}) = o(n^{1/2})$  follows from the intuitively appealing condition  $\#\mathcal{B} = o(n)$ .

*EXAMPLE 3 (Monotone shrinkage).* Let  $\mathcal{F}_{\text{mon}}$  be the set of all nonincreasing functions in  $[0, 1]^T$ . The class of candidate estimators  $\{fX : f \in \mathcal{F}_{\text{mon}}\}$  includes the nested model-selection estimators  $f_k X$ ,  $0 \leq k \leq n$ , defined by  $f_k(t) := 1\{t \leq k\}$ . In fact,  $\mathcal{F}_{\text{mon}}$  is the



convex hull of  $\mathcal{D}_{MS} := \{f_0, f_1, \dots, f_n\}$ . Elementary calculations show that

$$N(u, \mathcal{D}_{MS}) \leq 1 + u^{-2} \leq 2u^{-2}$$

for  $0 < u \leq 1$ . Together with Theorem 5.1 of Dudley (1987) it follows that

$$\log N(u, \mathcal{F}_{\text{mon}}) \leq Cu^{-1} \quad \text{for all } u \in ]0, 1].$$

*EXAMPLE 4 (Monotone shrinkage with respect to a quasi-order).* Let  $\preceq$  be a quasi-order relation on  $T$  (cf. Robertson *et al.* 1988, Chapter 1.3), and let  $\mathcal{F}_{\preceq}$  be the set of all functions in  $[0, 1]^T$  that are nonincreasing with respect to  $\preceq$ . That means, for all  $f \in \mathcal{F}_{\preceq}$  and  $s, t \in T$ ,

$$f(s) \geq f(t) \quad \text{if } s \preceq t.$$

Here one can easily deduce from the conclusion of Example 3 that

$$\log N(u, \mathcal{F}_{\preceq}) \leq CN_{\preceq} u^{-1}$$

for  $0 < u \leq 1$ , where  $N_{\preceq} = N_{\preceq, n}$  is the minimal cardinality of a partition of  $(T, \preceq)$  into totally ordered subsets. Thus  $J(\mathcal{F}_{\preceq})$  is of order  $O(N_{\preceq}^{1/2})$ . To give an example, suppose that  $X$  consists of  $n = 2^{k+1} - 1$  empirical Haar (or wavelet) coefficients, arranged as a binary tree. If this tree is equipped with its natural order  $\preceq$ , then the monotonicity constraint  $\hat{f} \in \mathcal{F}_{\preceq}$  means that  $\hat{f}X$  is a mixture of histogram estimators (cf. Engel 1994). Here  $N_{\preceq} = 2^k > n/2$ . Therefore, in order to apply our theory one has to replace the class  $\mathcal{F}_{\preceq}$  with suitable subclasses.

*EXAMPLE 5 (Shrinkage with bounded total variation).* Let  $\mathcal{F}_{(M)}$  be all functions  $f$  in  $[0, 1]^T$  with total variation not greater than  $M = M_n$ , i.e.

$$\sum_{t=2}^n |f(t) - f(t-1)| \leq M.$$

For instance, the class of functions  $f(t) := \max\{\min\{p(t), 1\}, 0\}$ , where  $p$  is a polynomial of degree less than or equal to  $M$ , belongs to  $\mathcal{F}_{(M)}$ . Any  $f \in \mathcal{F}_{(M)}$  can be written as  $(M+1)(f_1 - f_2)$  with  $f_1, f_2 \in \mathcal{F}_{\text{mon}}$ . Hence

$$\log N(u, \mathcal{F}_{(M)}) \leq 2 \log N\left([2(M+1)]^{-1}u, \mathcal{F}_{\text{mon}}\right) \leq C(M+1)u^{-1}$$

for  $0 < u \leq 1$ . In particular,  $J(\mathcal{F}_{(M)}) = O[(M + 1)^{1/2}]$ .

The minimizers  $\tilde{f}$  and  $\hat{f}$  in Examples 3-5 lack closed forms. Section 4 describes computational algorithms for  $\tilde{f}$  and  $\hat{f}$  in Examples 3-4. Example 5 differs from the remaining examples both theoretically as well as computationally and will be treated in detail elsewhere.

A particular consequence of Theorem 2.1 is that the modulation estimators are asymptotically minimax optimal for a large class of submodels for  $(\xi, \sigma^2)$ . Namely, for  $a \in [1, \infty]^T$  and  $c > 0$  define the linear minimax risk

$$\nu^2(a, c, \sigma^2) := \inf_{g \in [0, 1]^T} \sup_{\text{ave}(a\xi^2) \leq c} R(g, \xi, \sigma^2).$$

It is shown by Pinsker (1980) that the linear minimax risk approximates the unrestricted minimax risk in that

$$\inf_{\hat{\xi}} \sup_{\text{ave}(a\xi^2) \leq c} \rho(\hat{\xi}, \xi, \sigma^2) / \nu^2(a, c, \sigma^2) \rightarrow 1 \quad \text{as } n\nu^2(a, c, \sigma^2) \rightarrow \infty.$$

Moreover,

$$\nu^2(a, c, \sigma^2) = \sup_{\text{ave}(a\xi^2) \leq c} R(g_o, \xi, \sigma^2) = R(g_o, \xi_o, \sigma^2),$$

where  $g_o := [1 - (a/\mu_o)^{1/2}]^+$ ,  $\xi_o^2 := \sigma^2[(\mu_o/a)^{1/2} - 1]^+$ , and  $\mu_o > 0$  is the unique real number satisfying  $\text{ave}(a[(\mu_o/a)^{1/2} - 1]^+) = c/\sigma^2$ . The special case  $a \equiv 1$  yields (1.2).

If the minimax modulator  $g_o = g_o(\cdot | a, c/\sigma^2)$  happens to be in  $\mathcal{F}$ , which is certainly true for  $a \equiv 1$ , then

$$\sup_{\text{ave}(a\xi^2) \leq c} \rho(\hat{f}X, \xi, \sigma^2) \leq \sup_{\text{ave}(a\xi^2) \leq c} \left| \rho(\hat{f}X, \xi, \sigma^2) - R(\tilde{f}, \xi, \sigma^2) \right| + \nu^2(a, c, \sigma^2).$$

Thus Theorem 2.1 immediately implies the following minimax result, where the distribution of  $(X, \hat{\sigma}^2)$  is assumed to depend on  $(\xi, \sigma^2)$  only.

**COROLLARY 2.3** *Suppose that  $J(\mathcal{F}) = o(n^{1/2})$ , and that for every  $c, \sigma^2 > 0$ ,*

$$(2.2) \quad \epsilon_n(c, \sigma^2) := \sup_{\text{ave}(\xi^2) \leq c} \mathbb{E} |\hat{\sigma}^2 - \sigma^2| \rightarrow 0 \quad (n \rightarrow \infty).$$

*Then the modulation estimator  $\hat{f}X$  achieves the asymptotic minimax bound (1.2).*

More generally, let  $a = a_n \in [1, \infty]^T$  such that

$$(2.3) \quad [1 - (a/\mu)^{1/2}]^+ \in \mathcal{F} \quad \text{for all constants } \mu \geq 1.$$

Then for every  $c, \sigma^2 > 0$ ,

$$\sup_{\text{ave}(a\xi^2) \leq c} \rho(\hat{f}X, \xi, \sigma^2) \leq \nu^2(a, c, \sigma^2) + O[n^{-1/2}J(\mathcal{F}) + \epsilon_n(c, \sigma^2)]. \quad \square$$

Specifically, let  $a(t) = 1$  for  $t \in A \subset T$  and  $a(t) = \infty$  otherwise. Then  $\text{ave}(a\xi^2) \leq c$  is equivalent to  $\text{ave}(\xi^2) \leq c$  and  $\xi^2 = 0$  on  $T \setminus A$ . Here one can easily see that condition (2.3) is equivalent to  $1_A \in \mathcal{F}$ . The linear minimax risk equals

$$\nu^2(a, c, \sigma^2) = \frac{\sigma^2 \text{ave}(1_A)c}{\sigma^2 \text{ave}(1_A) + c},$$

which can be significantly smaller than the bound in (1.2).

In case of  $\mathcal{F} = \mathcal{F}_{\text{mon}}$  condition (2.3) is equivalent to  $a$  being nondecreasing on  $T$ .

We end this section with some examples for  $\hat{\sigma}$ . Internal estimators of  $\sigma^2$  depend only on  $X$  and require additional smoothness or dimensionality restrictions on the possible values of  $\xi$  to achieve the consistency property (2.2). One internal estimator of  $\sigma^2$ , analyzed by Rice (1984) and by Gasser et al. (1986) is

$$(2.4) \quad \hat{\sigma}_{(1)}^2 = [2(n-1)]^{-1} \sum_{t=2}^n [X(t) - X(t-1)]^2.$$

Here  $\mathbb{E}|\hat{\sigma}_n^2 - \sigma_n^2| \rightarrow 0$  as  $n \rightarrow \infty$  and

$$n^{-1} \sum_{t=2}^n [\xi_n(t) - \xi_n(t-1)]^2 \rightarrow 0.$$

External estimators of variance are available in linear models, where one observes an  $N$ -dimensional normal random vector  $Y$  with mean  $\mathbb{E}Y = D\xi$  and covariance matrix  $\text{Cov}(Y) = \sigma^2 I_N$  for some design matrix  $D \in \mathbf{R}^{N \times n}$ ,  $N = N_n > n$ . After suitable linear transformation of  $Y$  and  $\xi$  one may assume that  $\xi$  is the expectation of the vector  $X := (Y_1, Y_2, \dots, Y_n)$ . Then the standard estimator for  $\sigma^2$  is given by

$$\hat{\sigma}_{(2)}^2 := (N-n)^{-1} \sum_{i=n+1}^N Y_i^2,$$

which is independent from  $X$  with  $(N-n)\sigma^{-2}\hat{\sigma}_{(2)}^2 \sim \chi_{N-n}$ . This estimator also satisfies (2.2), provided that  $N-n \rightarrow \infty$ .

### 3 Confidence sets

Having replaced the maximum likelihood estimator  $X$  with  $\hat{f}X$ , a natural question is to what extent  $\hat{f}X$  is closer to the unknown signal  $\xi$  than  $X$ . More precisely we want to compare the distance  $L(X, \hat{f}X)^{1/2}$  with an upper confidence bound  $\hat{r} = \hat{r}(X, \hat{\sigma}^2)$  for  $L(\xi, \hat{f}X)^{1/2}$ . In geometrical terms, the confidence ball of primary interest is defined by

$$\hat{C} = \hat{C}_n := \{\gamma \in \mathbf{R}^T : L(\hat{f}X, \gamma) \leq \hat{r}^2\}.$$

The radius  $\hat{r}$  is chosen so that the coverage probability  $\mathbb{P}(\xi \in \hat{C})$  converges to  $\alpha \in ]0, 1[$  as  $n$  increases. The full definition of  $\hat{C}$  follows the theorem below. Underlying the construction is the confidence set idea sketched at the end of Stein (1981). The quality of  $\hat{C}$  as a set-valued estimator of  $\xi$  will be measured through the quadratic loss

$$(3.1) \quad L(\hat{C}, \xi) := \sup_{\gamma \in \hat{C}} L(\gamma, \xi) = [L(\hat{f}X, \xi)^{1/2} + \hat{r}]^2.$$

This is a natural extension of the quadratic loss defined in (1.1) and has an appealing projection-pursuit interpretation; see Beran (1996a).

One main assumption for this section is that

$$(3.2) \quad \begin{aligned} X_n \text{ and } \hat{\sigma}_n^2 \text{ are independent with } \mathcal{L}(\sigma_n^{-2}\hat{\sigma}_n^2) \text{ depending only on } n \\ \text{such that } \lim_{n \rightarrow \infty} m\left(\mathcal{L}[n^{1/2}(\sigma^{-2}\hat{\sigma}_n^2 - 1)], \mathcal{N}(0, \beta^2)\right) = 0. \end{aligned}$$

Here  $\beta^2 \geq 0$  is a given constant and  $m(\cdot, \cdot)$  metrizes weak convergence of distributions on the line. For instance, the estimator  $\hat{\sigma}_{(2)}^2$  of Section 2 satisfies Condition (3.2) with  $\beta := 2 \lim_{n \rightarrow \infty} n/(N_n - n)$ , provided that this limit exists. Condition (3.2) is made for the sake of simplicity. It could be replaced with weaker, but more technical conditions in order to include special internal estimators of variance such as  $\hat{\sigma}_{(1)}^2$ . A second key assumption is that

$$(3.3) \quad \int_0^1 \sqrt{\sup_n N(u, \mathcal{F}_n)} du < \infty.$$

Roughly speaking, this condition allows us to pretend that  $\hat{f}$  is equal to  $\tilde{f}$ . It is satisfied in all Examples 1-5, provided that  $\#\mathcal{B}_n = O(1)$  in Example 2,  $N_{\leq, n} = O(1)$  in Example 4, and  $M_n = O(1)$  in Example 5.

At first let us consider confidence balls centered at the naive estimator  $X$ . Since  $n\sigma^{-2} \text{ave}[(X - \xi)^2]$  has a chi-squared distribution with  $n$  degrees of freedom, we consider

$$\widehat{C}_N := \left\{ \gamma \in \mathbf{R}^T : \text{ave}[(X - \gamma)^2] \leq \widehat{\sigma}^2(1 + n^{-1/2}c) \right\}$$

for some fixed  $c$ . The inequality  $\text{ave}[(X - \xi)^2] \leq \widehat{\sigma}^2(1 + n^{-1/2}c)$  is equivalent to

$$n^{1/2} \left( \sigma^{-2} \text{ave}[(X - \xi)^2] - 1 \right) - n^{1/2}(\sigma^{-2}\widehat{\sigma}^2 - 1) \leq \sigma^{-2}\widehat{\sigma}^2 c = c + o_p(1).$$

Thus the Central Limit Theorem for the chi-squared distribution together with Condition 3.2 implies that  $c = (2 + \beta^2)^{1/2}\Phi^{-1}(\alpha)$  yields a confidence set  $\widehat{C}_N$  with

$$\lim_{n \rightarrow \infty} \sup_{\xi \in \mathbf{R}^T, \sigma^2 > 0} \left| \mathbb{P}\{\xi \in \widehat{C}_N\} - \alpha \right| = 0,$$

where  $\Phi^{-1}(\alpha)$  denotes the  $\alpha$ -th quantile of  $\mathcal{N}(0, 1)$ . Moreover,

$$\lim_{n \rightarrow \infty} \sup_{\xi \in \mathbf{R}^T} \mathbb{P}\left\{ |L(\widehat{C}_N, \xi) - 4\sigma^2| > \epsilon \right\} = 0 \quad \forall \epsilon > 0.$$

In what follows we shall see that confidence sets centered at a good modulation estimator  $\widehat{f}X$  dominate the naive confidence set  $\widehat{C}_N$  in terms of the loss  $L(\widehat{C}, \xi)$ .

To construct these confidence sets, we first determine the asymptotic distribution of

$$\widehat{d} = \widehat{d}_n := n^{1/2}[L(\widehat{f}X, \xi) - \widehat{R}_C(\widehat{f})].$$

This difference compares the loss of  $\widehat{f}X$  with an estimate for the expected loss of  $\widehat{f}X$ .

**THEOREM 3.1** *Under Conditions (3.2, 3.3),*

$$\lim_{n \rightarrow \infty} \sup_{\text{ave}(\xi^2) \leq c} m[\mathcal{L}(\widehat{d}), \mathcal{N}(0, \tau^2)] = 0,$$

for arbitrary  $c, \sigma^2 > 0$ , where

$$\tau^2 = \tau_n^2(\xi, \sigma^2) := 2\sigma^4 \text{ave}[(2\widetilde{f} - 1)^2] + \beta^2\sigma^4[\text{ave}(2\widetilde{f} - 1)]^2 + 4\sigma^2 \text{ave}[\xi^2(1 - \widetilde{f})^2].$$

A consistent estimator  $\widehat{\tau}^2 = \widehat{\tau}_n^2$  of  $\tau^2$  is obtained by substituting  $\widehat{\sigma}^2$  for  $\sigma^2$ ,  $\widehat{f}$  for  $\widetilde{f}$  and  $X^2 - \widehat{\sigma}^2$  for  $\xi^2$  in the expression for  $\tau^2$ . The implied estimator of the approximating normal

distribution  $\mathcal{N}(0, \tau^2)$  is  $\mathcal{N}(0, \hat{\tau}^2)$ . This leads to the following definition of a confidence ball for  $\xi$  that is centered at the modulation estimator  $\hat{f}X$ :

$$\hat{C} := \left\{ \gamma \in \mathbf{R}^T : L(\hat{f}X, \gamma) \leq \hat{R}_C(\hat{f}) + n^{-1/2} \hat{\tau} \Phi^{-1}(\alpha) \right\}.$$

The intended coverage probability of  $\hat{C}$  is  $\alpha$ . The next theorem establishes asymptotic properties of this confidence set construction. Beran (1994) treats in detail the example where  $\hat{f}X$  is the James-Stein estimator. That situation is much easier to analyze than the general case.

**THEOREM 3.2** *Under the conditions of Theorem 3.1, for arbitrary  $c, \sigma^2 > 0$ ,*

$$\begin{aligned} \lim_{n \rightarrow \infty, K \rightarrow \infty} \sup_{\text{ave}(\xi^2) \leq c} \mathbb{P} \left\{ |L(\hat{C}, \xi) - 4R(\tilde{f}, \xi, \sigma^2)| \geq Kn^{-1/2} \right\} &= 0 \\ \text{and} \quad \lim_{n \rightarrow \infty, K \rightarrow \infty} \sup_{\text{ave}(\xi^2) \leq c} \mathbb{P} \left\{ |\hat{\tau}^2 - R(\tilde{f}, \xi, \sigma^2)| \geq Kn^{-1/2} \right\} &= 0. \end{aligned}$$

Moreover,  $\hat{\tau}^2$  is consistent in that

$$\lim_{n \rightarrow \infty} \sup_{\text{ave}(\xi^2) \leq c} \mathbb{P} \left\{ |\hat{\tau}^2 - \tau^2| > \epsilon \right\} = 0 \quad \forall \epsilon > 0.$$

If

$$(3.4) \quad \liminf_{n \rightarrow \infty} \inf_{\text{ave}(\xi^2) \leq c} \tau_n^2(\xi, \sigma^2) > 0,$$

then

$$\lim_{n \rightarrow \infty} \sup_{\text{ave}(\xi^2) \leq c} \left| \mathbb{P} \{ \xi \in \hat{C} \} - \alpha \right| = 0.$$

A sufficient condition for (3.4) is the following: For every  $n$ ,  $\mathcal{F} = \mathcal{F}_n$  is such that

$$(3.5) \quad 1\{f \geq c\}f \in \mathcal{F} \quad \text{for all } f \in \mathcal{F} \text{ and } c \in [0, 1].$$

Condition (3.4) ensures that  $\mathcal{L}(\hat{d})$  does not approach a degenerate distribution. Note that Condition (3.5) is satisfied in Examples 1-4. When  $R(\tilde{f}, \xi, \sigma^2) = O(n^{-1/2})$  our confidence ball has loss  $L(\hat{C}, \xi) = O_p(n^{-1/2})$ . In fact, according to Theorem 2.1 of Li (1989) this is the smallest possible order of magnitude for a Euclidean confidence ball, unless one imposes further constraints on the signal  $\xi$ . The result (3.2) on asymptotic coverage of  $\hat{C}$  may be compared with the lower bound in Theorem 3.2 of Li (1989).

A key step in the proof of Theorem 3.1 is that in the definition of  $\hat{d}$  one may replace  $\tilde{f}$  with  $\hat{f}$ . Instead of the normal approximation underlying  $\hat{C}$  a bootstrap approximation of  $H = H_n := \mathcal{L}(\hat{d})$  that imitates the estimation of  $\tilde{f}$  seems to be more reliable in moderate dimension. Precisely, let  $\hat{H} = \hat{H}_n$  be the conditional distribution (function) of  $\hat{d}^*$  given  $(X, \hat{\sigma}^2)$ , where  $\hat{d}^*$  is computed as  $\hat{d}$  with the pair  $(X^*, \hat{\sigma}^{2*})$  in place of  $(X, \hat{\sigma}^2)$ . More precisely, let  $\hat{\xi} = \hat{\xi}(\cdot | X, \hat{\sigma}^2)$  be an estimator for  $\xi$ . Let  $S_n^2$  be a random variable with a specified distribution depending only on  $n$  such that

$$\lim_{n \rightarrow \infty} m\left(\mathcal{L}[n^{1/2}(S_n^2 - 1)], \mathcal{N}(0, \beta^2)\right) = 0,$$

where  $S_n$  and  $(X, \hat{\sigma}^2)$  are independent. Then

$$\mathcal{L}(X^*, \hat{\sigma}^{2*} | X, \hat{\sigma}^2) = \mathcal{N}(\hat{\xi}, \hat{\sigma}^2 I) \otimes \mathcal{L}(\hat{\sigma}^2 S_n^2 | X, \hat{\sigma}^2),$$

the product of the probability measures  $\mathcal{N}(\hat{\xi}, \hat{\sigma}^2 I)$  and  $\mathcal{L}(\hat{\sigma}^2 S_n^2 | X, \hat{\sigma}^2)$ . The resulting bootstrap confidence bound  $\hat{r}_b(\alpha)$  for  $L(\xi, \hat{f}X)$  is given by

$$\hat{r}_b^2(\alpha) = \hat{R}(\hat{f}) + n^{-1/2} \hat{H}^{-1}(\alpha).$$

The last theorem of this section states conditions, under which  $\hat{H}$  is a consistent estimator for  $H$ . An interesting fact is that neither  $\hat{\xi} = X$  nor  $\hat{\xi} = \hat{f}X$  satisfy these conditions.

**THEOREM 3.3** *Under the assumptions of Theorem 3.1,*

$$\lim_{n \rightarrow \infty} \sup_{\text{ave}(\xi^2) \leq c} \mathbb{P}\left\{|m(\hat{H}_n, H_n)| > \epsilon\right\} = 0 \quad \forall \epsilon > 0,$$

*provided that*

$$(3.6) \quad \hat{f} = \arg \min_{f \in \mathcal{F}} R(f, \hat{\xi}, \hat{\sigma}^2) \quad \text{almost surely,}$$

$$(3.7) \quad \limsup_{n \rightarrow \infty, K \rightarrow \infty} \sup_{\text{ave}(\xi^2) \leq c} \mathbb{P}\{\text{ave}(\hat{\xi}^2) > K\} = 0,$$

$$(3.8) \quad \lim_{n \rightarrow \infty} \sup_{\text{ave}(\xi^2) \leq c} \mathbb{P}\left\{\left|\text{ave}[\hat{\xi}^2(1 - \hat{f})^2] - \text{ave}[\xi^2(1 - \tilde{f})^2]\right| > \epsilon\right\} = 0 \quad \forall \epsilon > 0.$$

*In particular, suppose that each  $\mathcal{F}_n$  has the following property: For all  $X, \xi \in \mathbf{R}^T$  with  $X^2 > 0$  and any  $c \in [0, 1]$ ,*

$$c = \text{ave}[1\{\hat{f} = c\}(X^2 - \hat{\sigma}^2)]^+ / \text{ave}(1\{\hat{f} = c\}X^2) \quad \text{if } \{\hat{f} = c\} \neq \emptyset,$$

$$c = \text{ave}(1\{\tilde{f} = c\}\xi^2) / \text{ave}[1\{\tilde{f} = c\}(\sigma^2 + \xi^2)] \quad \text{if } \{\tilde{f} = c\} \neq \emptyset.$$

Then the function  $\hat{\xi} := (\hat{\sigma}^2 \hat{f} / (1 - \hat{f}))^{1/2}$  satisfies Conditions (3.6, 3.7, 3.8).

One can show that the last part of Theorem 3.3 applies to Examples 1-4. This follows from the representation of  $\hat{f}$  given in Section 4 and Robertson *et al.* (1988, Theorem 1.3.5). Here one can also show that  $\hat{\xi} := \hat{f}^{1/2} X$  satisfies these requirements, too. This yields a natural extension of the bootstrap method proposed by Beran (1994).

## 4 Computation of $\hat{f}$ and $\tilde{f}$

We restrict our attention to  $\hat{R} = \hat{R}_C$ . With the oracle modulator  $\tilde{g} := \xi^2 / (\sigma^2 + \xi^2)$  and its naive estimator  $\hat{g} := (X^2 - \hat{\sigma}^2) / X^2$  one can write

$$\begin{aligned} R(f, \xi, \sigma^2) &= \text{ave}[(f - \tilde{g})^2 (\sigma^2 + \xi^2)] + \text{ave}(\sigma^2 \tilde{g}), \\ \hat{R}_C(f) &= \text{ave}[(f - \hat{g})^2 X^2] + \text{ave}(\hat{\sigma}^2 \hat{g}). \end{aligned}$$

Hence both functions  $\tilde{f}$  and  $\hat{f}$  are metric projections of some function  $g$  onto  $\mathcal{F}$ .

Now we consider the family  $\mathcal{F}_{\preceq}$  of Example 4. Note that this case includes Examples 1-3 as special cases. The family  $\mathcal{F}_{\preceq}$  can be written as  $\mathcal{H}_{\preceq} \cap [0, 1]^T$ , where

$$\mathcal{H}_{\preceq} := \{\text{nonincreasing functions } h \in \mathbf{R}^T \text{ with respect to } \preceq\}.$$

Given arbitrary  $h \in \mathcal{H}_{\preceq}$  and  $a \in \mathbf{R}$ , the functions  $\max\{h, a\}$  and  $\min\{h, a\}$  also belong to  $\mathcal{H}_{\preceq}$ . Consequently, since  $0 \leq \tilde{g} \leq 1$  and  $\hat{g} \leq 1$ ,

$$\tilde{f} = \arg \min_{h \in \mathcal{H}_{\preceq}} \text{ave}[(h - \tilde{g})^2 (\sigma^2 + \xi^2)] \quad \text{and} \quad \hat{f} = \arg \min_{h \in \mathcal{H}_{\preceq}^+} \text{ave}[(h - \hat{g})^2 X^2],$$

where  $\mathcal{H}_{\preceq}^+ := \{h \in \mathcal{H} : h \geq 0\}$ . Hence  $\tilde{f}$  can be computed by any algorithm for projections onto  $\mathcal{H}_{\preceq}$ , while in case of  $\hat{f}$  one has to deal with a nonnegativity constraint. (Replacing  $\hat{g}$  with  $\hat{g}^+$  would yield an inconsistent estimator for  $\tilde{f}$  in general.) The latter problem is easy to solve. Let  $\hat{h}$  be the unique unrestricted projection

$$\hat{h} := \arg \min_{h \in \mathcal{H}_{\preceq}} \text{ave}[(h - \hat{g})^2 X^2]$$

of  $\hat{g}$ . Then

$$\hat{f} = \hat{h}^+.$$



For if  $\hat{h} = \hat{h}^+ - \hat{h}^-$ , then  $\hat{f} - \hat{h}^- \in \mathcal{H}_{\preceq}$  and  $\hat{h}^+ \in \mathcal{H}_{\succeq}^+$ . Hence

$$\begin{aligned}
\text{ave}[(\hat{h} - \hat{g})^2 X^2] &\leq \text{ave}[(\hat{f} - \hat{h}^- - \hat{g})^2 X^2] \\
&= \text{ave}[(\hat{f} - \hat{g})^2 X^2] + \text{ave}[(\hat{h}^- + \hat{g})^2 X^2] - \text{ave}[\hat{g}^2 X^2] \\
&\quad - 2 \text{ave}[\hat{f}\hat{h}^- X^2] \\
&\leq \text{ave}[(\hat{f} - \hat{g})^2 X^2] + \text{ave}[(\hat{h}^- + \hat{g})^2 X^2] - \text{ave}[\hat{g}^2 X^2] \\
&\leq \text{ave}[(\hat{h}^+ - \hat{g})^2 X^2] + \text{ave}[(\hat{h}^- + \hat{g})^2 X^2] - \text{ave}[\hat{g}^2 X^2] \\
&= \text{ave}[(\hat{h} - \hat{g})^2 X^2],
\end{aligned}$$

by definition of  $\hat{f}$  and  $\hat{h}$ . Thus  $\hat{f} - \hat{h}^-$  equals  $\hat{h}$ , whence  $\hat{f} = \hat{h}^+$ .

Explicit algorithms for computing  $\hat{h}$  are described in Robertson *et al.* (1988, Section 1). Because of the special form of  $\tilde{g}$  and  $\hat{g}$  one can even replace weighted least squares by ordinary least squares. Namely, let

$$H_{\xi} := \arg \min_{h \in \mathcal{H}_{\preceq}} \text{ave}((h - \xi^2)^2) \quad \text{and} \quad H_X := \arg \min_{h \in \mathcal{H}_{\preceq}} \text{ave}((h - X^2)^2).$$

Then

$$\tilde{f} = H_{\xi}/(\sigma^2 + H_{\xi}) \quad \text{and} \quad \hat{f} = (H_X - \hat{\sigma}^2)^+/H_X.$$

This follows from the min-max formula for antitonic regression (Robertson *et al.* 1988, Theorem 1.4.4). For let  $L$  and  $U$  be generic lower and upper sets, respectively. That means,

$$L = \{y \in T : y \preceq x \text{ for some } x \in L\} \quad \text{and} \quad U = \{y \in T : x \preceq y \text{ for some } x \in U\}.$$

Then

$$\begin{aligned}
\hat{h}(t) &= \max_{L:t \in L} \min_{U:t \in U} \text{ave}(1_{L \cap U} X^2 \hat{g}) / \text{ave}(1_{L \cap U} X^2) \\
&= \max_{L:t \in L} \min_{U:t \in U} \left(1 - \hat{\sigma}^2 \text{ave}(1_{L \cap U}) / \text{ave}(1_{L \cap U} X^2)\right) \\
&= 1 - \hat{\sigma}^2 / \max_{L:t \in L} \min_{U:t \in U} \left(\text{ave}(1_{L \cap U} X^2) / \text{ave}(1_{L \cap U})\right) \\
&= 1 - \hat{\sigma}^2 / H_X(t).
\end{aligned}$$

The formula for  $\tilde{f}$  is proved analogously.

## 5 Numerical examples

In this section we apply the proposed methods to empirical Fourier coefficients. We simulated data

$$Y(z) = \mu(z/n) + \epsilon(z), \quad z \in T,$$

where  $n = 1000$ ,  $\epsilon \sim \mathcal{N}_T(0, I)$ , and  $\mu$  is one of the following two functions on  $[0, 1]$ :

$$\begin{aligned} \text{Case 1: } \mu(u) &:= 2[6.75u^2(1-u)]^3, \\ \text{Case 2: } \mu(u) &:= \begin{cases} 1.5 & \text{if } 0.3 < u < 0.6, \\ 0.5 & \text{if } 0.6 < u \leq 0.8, \\ 2 & \text{if } 0.3 < u \leq 0.6, \\ 0 & \text{else,} \end{cases} \end{aligned}$$

With the orthonormal functions  $\phi_1 \equiv n^{1/2}$ ,  $\phi_{2k-1}(z) := (2n)^{1/2} \cos(2\pi kz/n)$  and  $\phi_{2k}(z) := (2n)^{1/2} \sin(2\pi kz/n)$  ( $1 \leq k < n/2$ ), and  $\phi_{n/2}(z) := n^{1/2}(-1)^z$  on  $T$  these data were transformed into

$$X := [\text{ave}(\phi_t Y)]_{t \in T} = [\text{ave}(\phi_t \mu(\cdot/n))]_{t \in T} + [\text{ave}(\phi_t \epsilon)]_{t \in T} =: \xi + E,$$

so that  $E \sim \mathcal{N}_T(0, I)$ . Then we computed the modulation estimator  $\hat{f}X$  using  $\hat{\sigma}^2 \equiv 1$  (for simplicity) and the modulator class

$$\mathcal{F} := \left\{ f \in \mathcal{F}_{\text{mon}} : f_{2k-1} = f_{2k} \text{ for } 1 \leq k < n/2 \right\}.$$

This yielded the estimator

$$\hat{\mu}(z/n) := \sum_{t \in T} \text{ave}(\phi_t \hat{f}X) \phi_t(z)$$

for  $\mu(z/n)$ . The additional requirement on  $f \in \mathcal{F}_{\text{mon}}$  takes into account the ambiguity of labeling sine and cosine functions of the same frequency. It also makes the resulting estimator  $\hat{\mu}$  equivariant under cyclical shifts of the data  $Y$ .

Figures 1a and 2a depict the data  $Y$  and the estimator  $\hat{\mu}(\cdot/n)$  in the two cases, respectively. Figures 1b and 2b show the estimator  $\hat{\mu}(\cdot/n)$  and the true function  $\mu(\cdot/n)$ . Figures 1a' and 2a' contain  $Y$  and  $\hat{\mu}(\cdot/n)$ , too, but this time the modulator  $\hat{f}$  was computed with  $\hat{\sigma} \equiv 0.9$  and  $\hat{\sigma} \equiv 1.1$ , respectively. A higher estimated variance leads to a

smoother estimate. These plots show clearly that estimating the variance is a crucial step. They also indicate the possibility to pick  $\hat{\sigma}$  visually.

Figures 1c and 2c show what is going on in the Fourier domain. The first plot shows the ideal greedy modulator  $\tilde{g}$  and the ideal monotone modulator  $\tilde{f} \in \mathcal{F}$ . The second plot gives the empirical counterparts  $\hat{g}^+$  and  $\hat{f}$ . Apparently the functions  $\tilde{g}$  and  $\hat{g}^+$  have little in common. On the other hand, the estimator  $\hat{f}$  and  $\tilde{f}$  are close to one another, as predicted by Theorem 2.2.

Note that a (hard) threshold estimator keeps all coefficients  $X(t)$  of  $X$  such that  $\hat{g}(t)$  is above a certain level while replacing the remaining coefficients with zero. In examples the authors looked at, this often led to peculiar estimators using only very few low frequencies or including some high frequencies. For cases 1 and 2, Figure 3 shows “oracle” (in a strong sense defined by the next display) threshold estimators  $\tilde{\mu}_{\text{th}}(\cdot/n) := \sum_{t \in T} \text{ave}(\phi_t \tilde{\xi}_{\text{th}}) \phi_t$ , where  $\tilde{\xi}_{\text{th}}(t) := 1\{|X(t)| \geq c_{\text{th}}\}X(t)$  and

$$c_{\text{th}} = c_{\text{th}}(X, \xi) := \arg \min_{c \geq 0} L(1\{|X| \geq c\}X, \xi).$$

In Case 1 the function  $\mu$  is very smooth, thus leading to a sparse signal  $\xi$ . In fact, the threshold fit is excellent. In Case 2 the threshold fit seems useless.

Finally we computed the bootstrap upper confidence bounds  $\hat{r}_b^2(0.9)$  for the actual loss  $L(\hat{f}X, \xi)$  as described in Section 3. The quantile  $\hat{H}^{-1}(0.9)$  was estimated in 3000 Monte-Carlo simulations. Table 1 contains the distance  $L(X, \hat{f}X)$ , the estimated risk  $\hat{R}(\hat{f})$ , the bootstrap bound  $\hat{r}_b^2(0.9)$ , the actual loss  $L(\hat{f}X, \xi)$  as well as the loss  $L(\tilde{\xi}_{\text{th}}, \xi)$ .

	$L(\hat{f}X, X)$	$\hat{R}(\hat{f})$	$\hat{r}_b^2(0.9)$	$L(\hat{f}X, \xi)$	$L(\tilde{\xi}_{\text{th}}, \xi)$
Case 1	0.9317	0.0108	0.0801	0.0109	0.0083
Case 2	0.9170	0.0755	0.1441	0.0546	0.1013

Table 1

## 6 Proofs

### 6.1 Auxiliary results

Our results utilize well-known techniques from empirical process theory. Theorem 6.1 below follows from standard symmetrization arguments and Pisier’s (1983) version of the

Chaining lemma (see also Pollard 1990, Sections 2 and 3). Theorem 6.2 is a simplified and modified version of Alexander's (1987) general results (see also Pollard 1990, Theorem 10.6).

Let  $S = \sum_{i=1}^n \phi_i$  with independent stochastic processes  $\phi_1, \phi_2, \dots, \phi_n$  on an index set  $\mathcal{T}$ . Examples for  $S$  are empirical processes and partial sum processes; see also Pollard (1990). For technical reasons we suppose that all  $\phi_i$  have continuous paths with respect to some metric  $d$  on  $\mathcal{T}$  such that  $(\mathcal{T}, d)$  is separable. Now define a random pseudodistance  $\hat{\rho}$  on  $\mathcal{T}$  via

$$\hat{\rho}(s, t)^2 := \sum_{i=1}^n [\phi_i(s) - \phi_i(t)]^2.$$

Further let

$$\rho(s, t) := \mathbb{E}(\hat{\rho}(s, t)^2)^{1/2}.$$

For any pseudo-metric  $\nu$  on  $\mathcal{T}$  define the covering numbers

$$N(u, \mathcal{T}, \nu) := \min \left\{ \#\mathcal{T}_o : \mathcal{T}_o \subset \mathcal{T}, \inf_{t_o \in \mathcal{T}_o} \nu(t_o, t) \leq u \forall t \in \mathcal{T} \right\}.$$

**THEOREM 6.1** *Suppose that  $S(t_o) \equiv 0$  for some  $t_o \in \mathcal{T}$ . Then*

$$\mathbb{E} \|S - \mathbb{E} S\|_{\mathcal{T}} \leq C \mathbb{E} \int_0^{\hat{D}} \sqrt{\log N(u, \mathcal{T}, \hat{\rho})} du,$$

where  $\hat{D} := \sup_{t \in \mathcal{T}} \hat{\rho}(t, t_o)$  and  $\|x\|_{\mathcal{T}} := \sup_{t \in \mathcal{T}} |x(t)|$ . □

**THEOREM 6.2** *Suppose that  $\mathcal{T} = \mathcal{T}_n$ ,  $\phi_i = \phi_{n,i}$  depend on  $n$  such that the following conditions are satisfied as  $n \rightarrow \infty$ :*

$$(6.1) \quad \mathbb{E} \sum_{i=1}^n \|\phi_i\|_{\mathcal{T}}^2 = O(1) \quad \text{and} \quad \mathbb{E} \sum_{i=1}^n 1\{\|\phi_i\|_{\mathcal{T}}^2 > u\} \|\phi_i\|_{\mathcal{T}}^2 = o(1) \text{ for all } u > 0;$$

$$(6.2) \quad \int_0^{\epsilon(n)} \sqrt{\log N(u, \mathcal{T}, \hat{\rho})} du \rightarrow_p 0 \quad \text{whenever } \epsilon(n) \downarrow 0.$$

Then

$$\begin{aligned} \mathbb{E} \|\hat{\rho}^2 - \rho^2\|_{\mathcal{T} \times \mathcal{T}} &= o(1) \quad \text{and} \quad N(u, \mathcal{T}, \rho) = O(1) \text{ for all } u > 0, \\ \sup_{s, t \in \mathcal{T} : \rho(s, t) \leq \alpha} \left| (S - \mathbb{E} S)(s) - (S - \mathbb{E} S)(t) \right| &\rightarrow_p 0 \quad \text{as } n \rightarrow \infty, \alpha \downarrow 0, \\ \|S - \mathbb{E} S\|_{\mathcal{T}} &= O_p(1). \end{aligned}$$

Moreover, let  $a_n : \mathcal{T} \rightarrow \mathbf{R}$  be arbitrary functions such that  $\sum_{t \in \mathcal{T}} |a_n(t)| = O(1)$ . Then

$$m\left(\mathcal{L}\left(\sum_{t \in \mathcal{T}} a_n(t)(S - \mathbb{E}S)(t)\right), \mathcal{N}\left(0, \text{Var}\left[\sum_{t \in \mathcal{T}} a_n(t)S(t)\right]\right)\right) \rightarrow 0. \quad \square$$

## 6.2 Proofs for Section 2

With the vector  $E := X - \xi$  of residuals we define random functions

$$W_1 := E^2 - \sigma^2 \quad \text{and} \quad W_2 := \xi E$$

on  $T$ . Then one can write

$$\begin{aligned} L(fX, \xi) - R(f, \xi, \sigma^2) &= \text{ave}\{[fE - (1-f)\xi]^2 - f^2\sigma^2 - (1-f)^2\xi^2\} \\ &= \text{ave}[f^2W_1 + 2(f^2 - f)W_2]. \end{aligned}$$

Moreover,  $X^2 = E^2 + 2\xi E + \xi^2 = W_1 + 2W_2 + \sigma^2 + \xi^2$ , and with

$$V := \hat{\sigma}^2 - \sigma^2$$

one obtains

$$\begin{aligned} \widehat{R}_C(f) - R(f, \xi, \sigma^2) &= \text{ave}[f^2\hat{\sigma}^2 + (1-f)^2(X^2 - \hat{\sigma}^2) - f^2\sigma^2 - (1-f)^2\xi^2] \\ &= \text{ave}[(f^2 - 2f + 1)(W_1 + 2W_2) + (2f - 1)V], \\ \widehat{R}_B(f) - R(f, \xi, \sigma^2) &= \text{ave}\left((1-f)^2\xi^2 + (1-f)^2(W_1 + 2W_2 - V)\right)^+ - \text{ave}[(1-f)^2\xi^2 + f^2V] \\ &= \lambda(f) \text{ave}\left((f^2 - 2f + 1)(W_1 + 2W_2) + [f^2 - \lambda(f)(1-f)^2]V\right), \end{aligned}$$

where  $0 \leq \lambda(f) \leq 1$ . Hence the following inequalities hold:

### LEMMA 6.3

$$\begin{aligned} \sup_{f \in \mathcal{F}} \left| \widehat{R}(f) - R(f, \xi, \sigma^2) \right| &\leq 4 \sup_{g \in \mathcal{G}} |\text{ave}(gW_1)| + 8 \sup_{g \in \mathcal{G}} |\text{ave}(gW_2)| + |V|, \\ \sup_{f \in \mathcal{F}} \left| L(fX, \xi) - R(f, \xi, \sigma^2) \right| &\leq \sup_{g \in \mathcal{G}} |\text{ave}(gW_1)| + 4 \sup_{g \in \mathcal{G}} |\text{ave}(gW_2)|, \end{aligned}$$

where  $\mathcal{G} := \{fg : f, g \in \mathcal{F}\}$ .

If  $\mathcal{F}$  is the convex hull of a family  $\mathcal{D}$  of (indicators of) subsets of  $T$ , which is closed under intersection, then  $\mathcal{G} = \mathcal{F}$ . In fact, in that case one may replace  $J(\mathcal{F})$  in Theorems 2.1, 2.2 and Corollary 2.3 with  $J(\mathcal{D})$ .

Theorem 2.1 follows easily from Lemma 6.3 and the following result, which is stated for potentially non-Gaussian error  $E$ .

**LEMMA 6.4** *Let  $X = \xi + E$ , where  $E$  has independent components with mean zero and variance  $\sigma^2$ . Then*

$$\begin{aligned}\mathbb{E} \sup_{g \in \mathcal{G}} |\text{ave}(gW_1)| &\leq C J(\mathcal{F}) \sqrt{\frac{\mathbb{E} \text{ave}(E^4)}{n}}, \\ \mathbb{E} \sup_{g \in \mathcal{G}} |\text{ave}(gW_2)| &\leq C J(\mathcal{F}) \sqrt{\frac{\sigma^2 \text{ave}(\xi^2)}{n}}.\end{aligned}$$

**Proof of Lemma 6.4.** Elementary calculations show that the uniform covering numbers of  $\mathcal{G}$  satisfy

$$N(u, \mathcal{G}) \leq N(u/2, \mathcal{F})^2 \quad \text{for all } u > 0.$$

Thus  $J(\mathcal{G}) \leq 4J(\mathcal{F})$ , and  $\mathcal{F}$  may be replaced with  $\mathcal{G}$ . Now the asserted inequalities for  $W_1, W_2$  are direct consequences of Theorem 6.1. For let  $\mathcal{T} := \mathcal{G}$  and

$$\phi_i(g) := \begin{cases} n^{-1} E(i)^2 g(i) & \text{for } j = 1, \\ n^{-1} \xi(i) E(i) g(i) & \text{for } j = 2. \end{cases}$$

Then  $\sup_{g \in \mathcal{G}} |\text{ave}(gW_j)| = \|S - \mathbb{E} S\|_{\mathcal{G}}$ .

For  $j = 1$ , and arbitrary  $g, h \in \mathcal{G}$ ,

$$\begin{aligned}\hat{\rho}(g, h)^2 &= n^{-1} \text{ave}(E^4(g - h)^2) \\ &= n^{-1} \text{ave}(E^4) d_{\hat{Q}}(g, h)^2 \\ &\leq n^{-1} \text{ave}(E^4)\end{aligned}$$

for some (random) probability measure  $\hat{Q}$  on  $T$ . Thus  $\hat{D} \leq n^{-1/2} \text{ave}(E^4)^{1/2}$  and

$$N(u, \mathcal{G}, \hat{\rho}) \leq N[n^{1/2} \text{ave}(E^4)^{-1/2} u, \mathcal{G}].$$

Hence

$$\begin{aligned}\int_0^{\hat{D}} \sqrt{\log N(u, \mathcal{G}, \hat{\rho})} du &\leq \int_0^{n^{-1/2} \text{ave}(E^4)^{1/2}} \sqrt{\log N[n^{1/2} \text{ave}(E^4)^{-1/2} u, \mathcal{G}]} du \\ &= n^{-1/2} \text{ave}(E^4)^{1/2} \int_0^1 \sqrt{\log N(u, \mathcal{G})} du,\end{aligned}$$

and  $\mathbb{E}[\text{ave}(E^4)^{1/2}] \leq [\mathbb{E} \text{ave}(E^4)]^{1/2}$ .

For  $j = 2$  the same arguments yield

$$\int_0^{\widehat{D}} \sqrt{\log N(u, \mathcal{G}, \widehat{\rho})} du \leq n^{-1/2} \text{ave}(\xi^2 E^2)^{1/2} \int_0^1 \sqrt{\log N(u, \mathcal{G})} du,$$

and  $\mathbb{E}[\text{ave}(\xi^2 E^2)^{1/2}] \leq \sigma \text{ave}(\xi^2)^{1/2}$ .  $\square$

**Proof of Theorem 2.2.** Let  $w_1 := X^2$ ,  $w_2 := \sigma^2 + \xi^2$ ,  $g_1 := \widehat{g}$  and  $g_2 := \widetilde{g}$ . Then  $\widehat{f} = f_1$  and  $\widetilde{f} = f_2$ , where generally

$$f_i := \arg \min_{f \in \mathcal{F}} \text{ave}(w_i(f - g_i)^2).$$

Let  $(i, j)$  be either  $(1, 2)$  or  $(2, 1)$ . By convexity of  $\mathcal{F}$ ,

$$(6.3) \quad \left. \frac{\partial}{\partial t} \right|_{t=0} \text{ave}(w_i[f_i + t(f_j - f_i) - g_i]^2) = 2 \text{ave}(w_i(f_i - g_i)(f_j - f_i)) \geq 0.$$

Hence,

$$\begin{aligned} & \left. \frac{\partial}{\partial t} \right|_{t=1} \text{ave}(w_i[f_i + t(f_j - f_i) - g_i]^2) \\ &= 2 \text{ave}(w_i(f_i - g_i)(f_j - f_i)) + 2 \text{ave}(w_i(f_j - f_i)^2) \\ &\geq 2 \text{ave}(w_i(f_j - f_i)^2). \end{aligned}$$

On the other hand,

$$\begin{aligned} \left. \frac{\partial}{\partial t} \right|_{t=1} \text{ave}(w_i[f_i + t(f_j - f_i) - g_i]^2) &= - \left. \frac{\partial}{\partial t} \right|_{t=0} \text{ave}(w_i[f_j + t(f_i - f_j) - g_i]^2) \\ &= -2 \text{ave}(w_i(f_j - g_i)(f_i - f_j)) \\ &\leq 2 \text{ave}([w_j(f_j - g_j) - w_i(f_j - g_i)](f_i - f_j)), \end{aligned}$$

where the latter inequality follows from (6.3) with  $i$  and  $j$  interchanged. Thus we end up with the inequality

$$\text{ave}(w_i(f_j - f_i)^2) \leq \text{ave}([w_j(f_j - g_j) - w_i(f_j - g_i)](f_i - f_j)).$$

Elementary algebra yields

$$w_j(f_j - g_j) - w_i(f_j - g_i) = \begin{cases} (\widehat{f} - 1)(W_1 + 2W_2) + V & \text{if } (i, j) = (1, 2), \\ (1 - \widetilde{f})(W_1 + 2W_2) - V & \text{if } (i, j) = (2, 1). \end{cases}$$

Consequently, Lemma 6.4 yields

$$\begin{aligned}
\mathbb{E} \operatorname{ave}\left(X^2(\widehat{f} - \widetilde{f})^2\right) &\leq 2 \mathbb{E} \sup_{f \in \mathcal{F}} \left| \operatorname{ave}\left((1 - \widetilde{f})(W_1 + 2W_2)f\right) \right| + \mathbb{E} |V| \\
&\leq CJ(\mathcal{F})n^{-1/2} \left( \sigma^2 + \sigma \operatorname{ave}[(1 - \widetilde{f})^2 \xi^2]^{1/2} \right) + \mathbb{E} |V| \\
&\leq CJ(\mathcal{F})n^{-1/2} \sigma^2 + \mathbb{E} |V|, \\
\mathbb{E} \operatorname{ave}\left((\sigma^2 + \xi^2)(\widehat{f} - \widetilde{f})^2\right) &\leq 4 \mathbb{E} \sup_{g \in \mathcal{G}} \left| \operatorname{ave}\left((W_1 + 2W_2)g\right) \right| + \mathbb{E} |V| \\
&\leq CJ(\mathcal{F})n^{-1/2} \left( \sigma^2 + \sigma \operatorname{ave}(\xi^2)^{1/2} \right) + \mathbb{E} |V|.
\end{aligned}$$

In the first case we applied Lemma 6.4 with  $(1 - \widetilde{f})\xi$  in place of  $\xi$  and utilized the fact that  $\operatorname{ave}\left((1 - \widetilde{f})^2 \xi^2\right) \leq R(\widetilde{f}, \xi, \sigma^2) \leq R(1, \xi, \sigma^2) = \sigma^2$ .  $\square$

### 6.3 Proofs for Section 3

Throughout this subsection asymptotic statements are meant as  $n \rightarrow \infty$  uniformly in  $\operatorname{ave}(\xi_n^2) \leq c$ , where  $c, \sigma^2 > 0$  are arbitrary and fixed.

**Proof of Theorem 3.1.** At first it is shown that

$$(6.4) \quad \widehat{d} = \widehat{d}(\widetilde{f}) + o_p(1),$$

where  $\widehat{d}(f) := n^{1/2}[L(fX, \xi) - \widehat{R}_C(f)]$ , i.e.  $\widehat{d} = \widehat{d}(\widehat{f})$ . The formulas of Section 6.2 for  $L(fX, \xi)$  and  $\widehat{R}_C(f)$  yield

$$\widehat{d}(f) = n^{1/2} \operatorname{ave}\left((2f - 1)W_1 + 2(f - 1)W_2 + (1 - 2f)V\right).$$

In particular,

$$\widehat{d}(f) - \widehat{d}(\widetilde{f}) = 2n^{1/2} \operatorname{ave}\left((f - \widetilde{f})(W_1 + W_2 - V)\right).$$

It follows from the proof of Theorem 2.2 that

$$(6.5) \quad \operatorname{ave}\left((\sigma^2 + \xi^2)(\widehat{f} - \widetilde{f})^2\right) = o_p(1),$$

Hence  $|n^{1/2}V \operatorname{ave}(\widehat{f} - \widetilde{f})| = o_p(1)$ , and it suffices to show that

$$(6.6) \quad n^{1/2} \operatorname{ave}[(\widehat{f} - \widetilde{f})W_j] = o_p(1) \quad \text{for } j = 1, 2.$$



For  $j = 1$  one can apply Theorem 6.2 with  $\mathcal{T} = \mathcal{F}$  and  $\phi_i(f) := n^{-1/2} f(i) E(i)^2$  similarly as in the proof of Theorem 2.1. The assumptions of Theorem 6.2 are satisfied, because

$$\begin{aligned} \mathbb{E} \sum_{i=1}^n \|\phi_i\|_{\mathcal{F}}^2 &= \mathbb{E} \text{ave}(E^4) = 3\sigma^4, \\ \mathbb{E} \sum_{i=1}^n 1\{\|\phi_i\|_{\mathcal{F}}^2 > u\} \|\phi_i\|_{\mathcal{F}}^2 &= \mathbb{E} \text{ave}(1\{E^4 > nu\} E^4) = o(1) \quad \forall u > 0, \\ \int_0^{\epsilon(n)} \sqrt{\log N(u, \mathcal{F}, \hat{\rho})} du &\leq \text{ave}(E^4)^{1/2} \int_0^{\epsilon(n)/\text{ave}(E^4)^{1/2}} \sqrt{N(u, \mathcal{F})} du \\ &\rightarrow_p 0 \quad \text{whenever } \epsilon(n) \downarrow 0. \end{aligned}$$

Hence assertion (6.6) for  $j = 1$  follows from the fact that

$$\rho(\hat{f}, \tilde{f})^2 = \mathbb{E} \text{ave}[(f - \tilde{f})^2 E^4] \Big|_{f=\hat{f}} = 3\sigma^4 \text{ave}[(\hat{f} - \tilde{f})^2] = o_p(1),$$

according to (6.5).

As for  $j = 2$ ,  $f \mapsto n^{1/2} \text{ave}(fW_2)$  is a centered Gaussian process with continuous paths with respect to  $\rho(f, g) = \sigma \text{ave}[\xi^2(f - g)^2]^{1/2} = \text{Var}\{n^{1/2} \text{ave}[(f - g)W_2]\}^{1/2}$ . Hence Pisier's (1983) maximal inequality yields

$$\mathbb{E} \sup_{f, g \in \mathcal{F}: \rho(f, g) \leq \epsilon(n)} \left| n^{1/2} \text{ave}((f - g)W_2) \right| \rightarrow_p 0 \quad \text{whenever } \epsilon(n) \downarrow 0,$$

and (6.6) follows from  $\rho(\hat{f}, \tilde{f}) = o_p(1)$ , again a consequence of (6.5).

Because of expansion (6.4) it suffices to show that the distribution of

$$\hat{d}(\tilde{f}) = n^{1/2} \text{ave}[(2\tilde{f} - 1)W_1 + 2(\tilde{f} - 1)W_2 + (1 - 2\tilde{f})V]$$

is asymptotically normal with mean zero and variance  $\tau^2$ . By assumption (3.2),  $(W_1, W_2)$  and  $V$  are independent with  $m(\mathcal{L}(n^{1/2}V), \mathcal{N}(0, \beta^2\sigma^4)) \rightarrow 0$ . Further, Lindeberg's Central Limit Theorem entails

$$m(\mathcal{L}(n^{1/2} \text{ave}[(2\tilde{f} - 1)W_1]), \mathcal{N}(0, 2\sigma^4 \text{ave}[(2\tilde{f} - 1)^2])) \rightarrow 0,$$

whereas  $\mathcal{L}(n^{1/2} \text{ave}[(\tilde{f} - 1)W_2]) = \mathcal{N}(0, \sigma^2 \text{ave}[\xi^2(1 - \tilde{f})^2])$ . Thus it remains to be shown that  $n^{1/2} \text{ave}[(2\tilde{f} - 1)W_1]$  and  $n^{1/2} \text{ave}[(\tilde{f} - 1)W_2]$  are asymptotically independent. For

that purpose we split  $T$  into the two subsets  $T^{(1)} := \{\xi^2 \leq n^{1/2}\}$  and  $T^{(2)} := T \setminus T^{(1)}$ . For any  $a \in \mathbf{R}^T$  let  $a^{(k)}(t) := 1\{t \in T^{(k)}\}a(t)$ . Note that  $(W_1^{(1)}, W_2^{(1)})$  and  $(W_1^{(2)}, W_2^{(2)})$  are independent. In addition,

$$n \mathbb{E} \text{ave}[(2\tilde{f} - 1)W_1^{(2)}]^2 \leq 2\sigma^4 \#T^{(2)}/n \leq 2\sigma^4 \text{ave}(\xi^2)/n^{1/2} = o(1),$$

while  $n^{1/2} \text{ave}[(2\tilde{f} - 1)W_1^{(1)}]$  and  $n^{1/2} \text{ave}[(\tilde{f} - 1)W_2^{(1)}]$  are asymptotically independent, according to the bivariate CLT.  $\square$

**Proof of Theorem 3.2.** It follows from Theorem 3.1 and the proof of Theorem 2.1 that

$$\begin{aligned} \hat{\tau}^2 &= \hat{R}_C(\hat{f}) + O_p(n^{-1/2}) = R(\tilde{f}, \xi, \sigma^2) + O_p(n^{-1/2}) \\ \text{and } L(\hat{f}X, \xi) &= R(\tilde{f}, \xi, \sigma^2) + O_p(n^{-1/2}). \end{aligned}$$

Consequently,

$$L(\hat{C}, \xi) = [L(\hat{f}X, \xi)^{1/2} + \hat{\tau}]^2 = 4R(\tilde{f}, \xi, \sigma^2) + O_p(n^{-1/2}).$$

As for the consistency of  $\hat{\tau}^2$ , we know already that  $\hat{\sigma}^2 = \sigma^2 + o_p(1)$  and  $\text{ave}[(\hat{f} - \tilde{f})^2] = o_p(1)$ . Thus the assertion follows from Theorem 2.1 via

$$\begin{aligned} \text{ave}[(X^2 - \hat{\sigma}^2)(1 - \hat{f})^2] &= \hat{R}_C(\hat{f}) - \hat{\sigma}^2 \text{ave}(\hat{f}^2) \\ &= R(\tilde{f}, \xi, \sigma^2) - \sigma^2 \text{ave}(\tilde{f}^2) + o_p(1) \\ &= \text{ave}[\xi^2(1 - \tilde{f})^2] + o_p(1). \end{aligned}$$

Given consistency of  $\hat{\tau}^2$ , the fact that (3.4) implies asymptotic level  $\alpha$  of  $\hat{C}$  follows from standard arguments. It remains to be shown that (3.5) implies (3.4). Suppose that the latter condition is violated. This is equivalent to

$$(6.7) \quad \liminf_{n \rightarrow \infty} \inf_{\text{ave}(\xi^2) \leq c} \left( \text{ave}[(\tilde{f} - 1/2)^2] + \text{ave}[\xi^2(1 - \tilde{f})^2] \right) = 0.$$

However,

$$R(\tilde{f}, \xi, \sigma^2) \geq \sigma^2 \text{ave}(\tilde{f}^2) \geq \sigma^2/4 - \text{ave}[(\tilde{f} - 1/2)^2]^{1/2}$$

while for  $f_o := 1\{\tilde{f} \geq 3/4\}\tilde{f} \in \mathcal{F}$ ,

$$\begin{aligned} R(\tilde{f}, \xi, \sigma^2) &\leq R(f_o, \xi, \sigma^2) \\ &= \text{ave}[\xi^2(1 - f_o)^2] + \sigma^2 \text{ave}(f_o^2) \\ &\leq 16 \text{ave}[\xi^2(1 - \tilde{f})^2] + 9\sigma^2 \text{ave}[(\tilde{f} - 1/2)^2]. \end{aligned}$$

These two inequalities are incompatible with (6.7).  $\square$

**Proof of Theorem 3.3.** It follows from the fact that  $\hat{\sigma}^2 = \sigma^2 + o_p(1)$  and Theorem 3.1, applied to  $\mathcal{L}(X^*, \hat{\sigma}^{2*} | X, \hat{\sigma}^2)$  in place of  $\mathcal{L}(X, \hat{\sigma}^2)$ , that  $m[\hat{H}, \mathcal{N}(0, \hat{\tau}^2)] = o_p(1)$ , where  $\hat{\tau}^2$  is defined as  $\tau^2$  with  $(\hat{\xi}, \hat{\sigma}^2, \hat{f})$  in place of  $(\xi, \sigma^2, \tilde{f})$ . One easily checks that  $\hat{\tau}^2 = \tau^2 + o_p(1) = O_p(1)$ , whence

$$m[\hat{H}, H] \leq m[\hat{H}, \mathcal{N}(0, \hat{\tau}^2)] + m[\mathcal{N}(0, \hat{\tau}^2), \mathcal{N}(0, \tau^2)] + m[H, \mathcal{N}(0, \tau^2)] = o_p(1).$$

Now suppose that  $\mathcal{F}$  has the special properties described in Theorem 3.3. If  $\hat{\sigma}^2 > 0$ , then  $\hat{f} < 1$ , so that  $\hat{\xi}$  is well-defined. Condition 3.6 follows from

$$\arg \min_{g \in [0,1]^T} R(g, \hat{\xi}, \hat{\sigma}^2) = \hat{\xi}^2 / (\hat{\sigma}^2 + \hat{\xi}^2) = \hat{f}.$$

The special construction of  $\hat{\xi}$  entails that

$$\hat{\xi}^2 = \sum_{c \in \hat{f}(T)} 1\{\hat{f} = c\} \text{ave}[1\{\hat{f} = c\}(X^2 - \hat{\sigma}^2)]^+ / \text{ave}(1\{\hat{f} = c\}).$$

Consequently the moment condition (3.7) follows from

$$\text{ave}(\hat{\xi}^2) = \sum_{c \in \hat{f}(T)} \text{ave}[1\{\hat{f} = c\}(X^2 - \hat{\sigma}^2)]^+ \leq \text{ave}(X^2) = O_p(1),$$

while (3.8) follows from

$$\begin{aligned} \text{ave}[\hat{\xi}^2(1 - \hat{f})^2] &= \text{ave}[\hat{\sigma}^2 \hat{f}(1 - \hat{f})] \\ &= \text{ave}[\sigma^2 \tilde{f}(1 - \tilde{f})] + o_p(1) \\ &= \text{ave}[\sigma^2 \tilde{f} / (1 - \tilde{f}) (1 - \tilde{f})^2] + o_p(1) \\ &= \sum_{c \in \tilde{f}(T)} \text{ave}[1\{\tilde{f} = c\} \text{ave}(1\{\tilde{f} = c\} \xi^2) / \text{ave}(1\{\tilde{f} = c\}) (1 - \tilde{f})^2] + o_p(1) \\ &= \sum_{c \in \tilde{f}(T)} \text{ave}[1\{\tilde{f} = c\} \xi^2 (1 - \tilde{f})^2] + o_p(1) \\ &= \text{ave}[\xi^2 (1 - \tilde{f})^2] + o_p(1). \quad \square \end{aligned}$$

## References

- ALEXANDER, K.S. (1987). Central limit theorems for stochastic processes under random entropy conditions. *Probab. Theory Rel. Fields* **75**, 351–378.
- BERAN, R. (1994). Stein confidence sets and the bootstrap. *Statistica Sinica* **5**, 109–127.
- BERAN, R. (1996a). Confidence sets centered at  $C_p$ -estimators. *Ann. Inst. Statist. Math.* **48**, 1–15.
- BERAN, R. (1996b). Stein estimation in high dimensions: a retrospective. *Festschrift in Honor of Madan L. Puri* (E. Brunner and M. Denker, eds.) 91–110. VSP, Zeist.
- CASELLA, G. AND J.T. HWANG (1982). Limit expressions for the risk of James-Stein estimators. *Canad. J. Statist.* **10**, 305–309.
- DONOHO, D.L. AND I.M. JOHNSTONE (1994). Ideal spatial adaptation by wavelet shrinkage. *Biometrika* **81**, 425–455.
- DUDLEY, R.M. (1987). Universal Donsker classes and metric entropy. *Ann. Prob.* **15**, 1306–1326.
- ENGEL, J. (1994). A simple wavelet approach to nonparametric regression from recursive partitioning schemes. *J. Multivariate Anal.* **49**, 242–254.
- GASSER, T., L. SROKA AND C. JENNEN-STEINMETZ (1986). Residual variance and residual pattern in nonlinear regression. *Biometrika* **73**, 625–633.
- GOLUBEV, G.K. AND NUSSBAUM, M. (1992). Adaptive spline estimates in a nonparametric regression model. *Theory Probab. Appl.* **37**, 521–529.
- JAMES, W. AND C. STEIN (1961). Estimation with quadratic loss. *Proc. Fourth Berkeley Symp. Math. Statist. Prob.* **1** (J. Neyman, ed.), 361–380. University of California Press, Berkeley.
- KNEIP, A. (1994). Ordered linear smoothers. *Ann. Statist.* **22**, 835–866.

- LI, K.-C. (1987). Asymptotic optimality for  $C_p$ ,  $C_L$ , cross-validation and generalized cross-validation: Discrete index set. *Ann. Statist.* **15**, 958–976.
- LI, K.-C. (1989). Honest confidence sets for nonparametric regression. *Ann. Statist.* **17**, 1001–1008.
- MALLOWS, C. (1973). Some comments on  $C_p$ . *Technometrics* **15**, 661–675.
- NUSSBAUM, M. (1996). The Pinsker bound: a review. Unpublished preprint.
- PINSKER, M.S. (1980). Optimal filtration of square-integrable signals in Gaussian noise. *Problems Inform. Transmission* **16**, 120–133.
- PISIER, G. (1983). Some applications of the metric entropy condition to harmonic analysis. *Banach Spaces, Harmonic Analysis, and Probability Theory* (R.C. Blei and S.J. Sidney, eds.). Lect. Notes Math. **995**, 123–154. Springer, Berlin.
- POLLARD, D. (1990). *Empirical Processes: Theory and Applications*. NSF-CBMS Regional Conf. Ser. Prob. Statist. **2**. IMS, Hayward.
- RICE, J. (1984). Bandwidth choice for nonparametric regression. *Ann. Statist.* **12**, 1215–1230.
- ROBERTSON, T., F.T. WRIGHT AND R.L. DYKSTRA (1988). *Order Restricted Statistical Inference*. Wiley, New York.
- STEIN, C. (1956). Inadmissibility of the usual estimator for the mean of a multivariate normal distribution. *Proc. Third Berkeley Symp. Math. Statist. Prob.* (J. Neyman, ed.), 197–206. University of California Press, Berkeley.
- STEIN, C. (1966). An approach to the recovery of inter-block information in balanced incomplete block designs. *Research Papers in Statistics. Festschrift for Jerzy Neyman* (F.N. David, ed.), 351–366. Wiley, London.
- STEIN, C. (1981). Estimation of the mean of a multivariate normal distribution. *Ann. Statist.* **9**, 1135–1151.

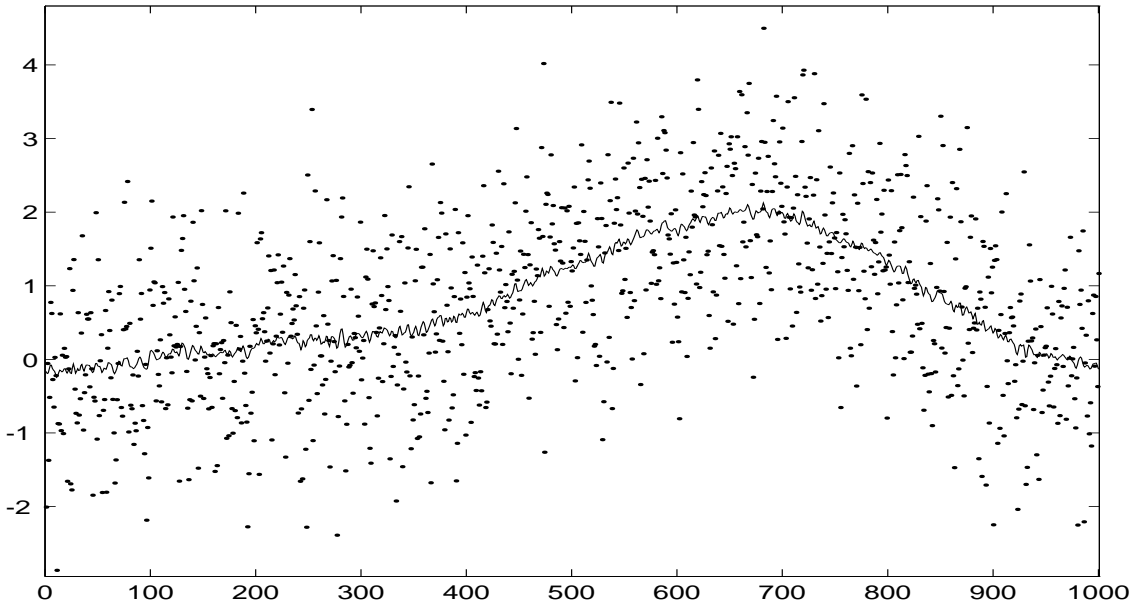


Figure 1a.  $Y$  and  $\hat{\mu}(\cdot/n)$ .

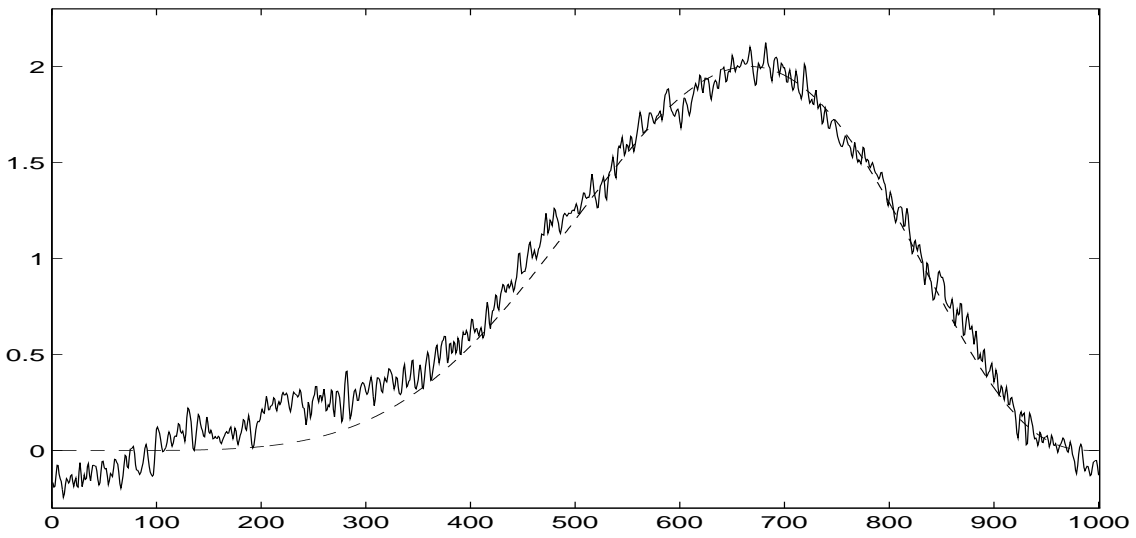
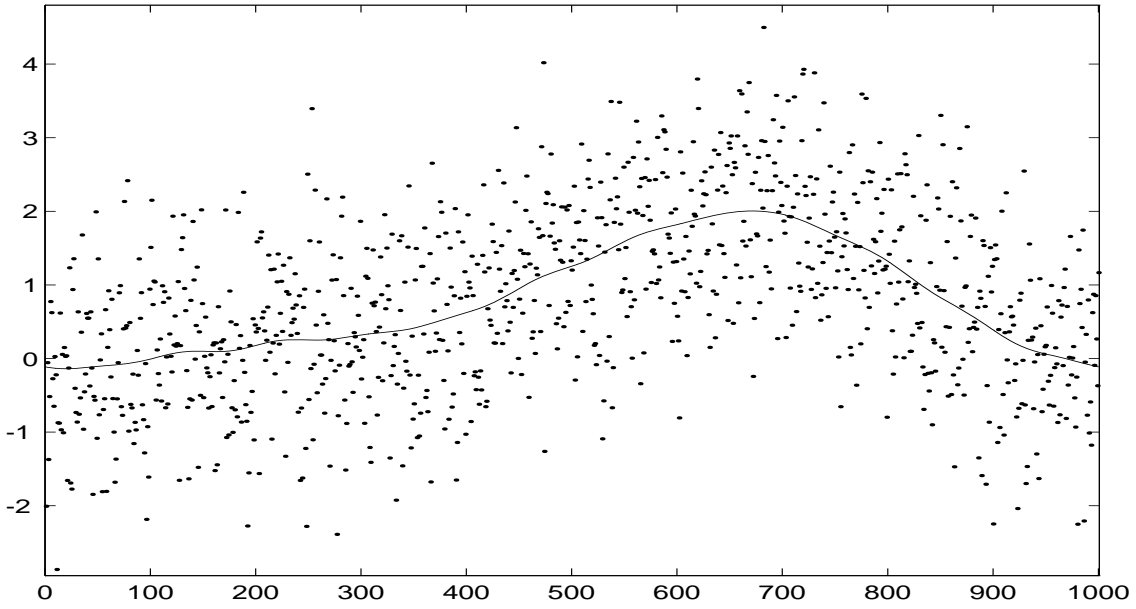
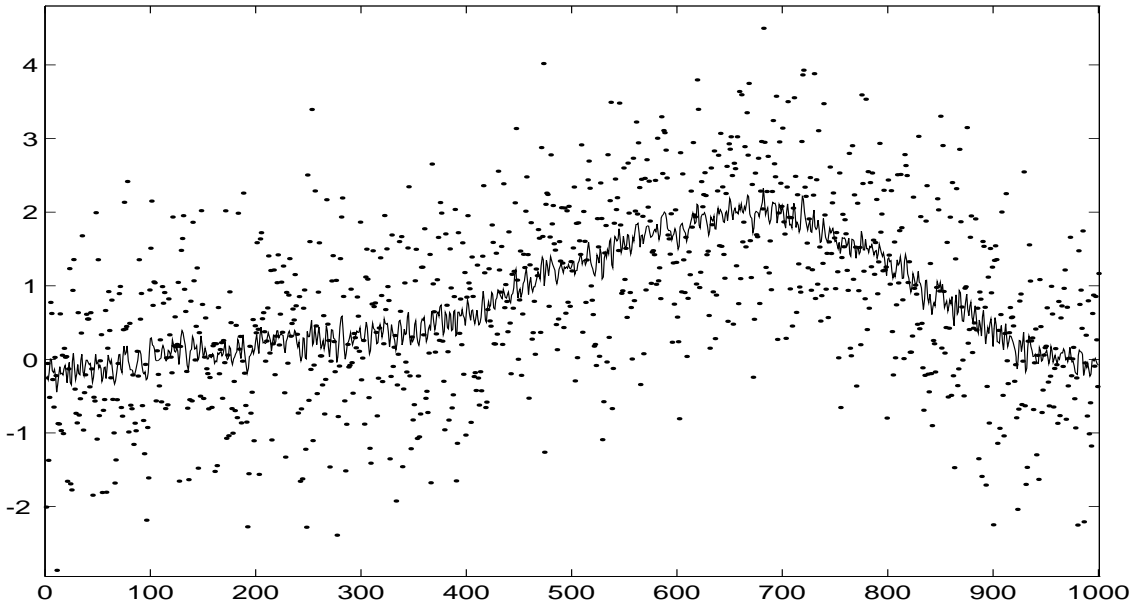
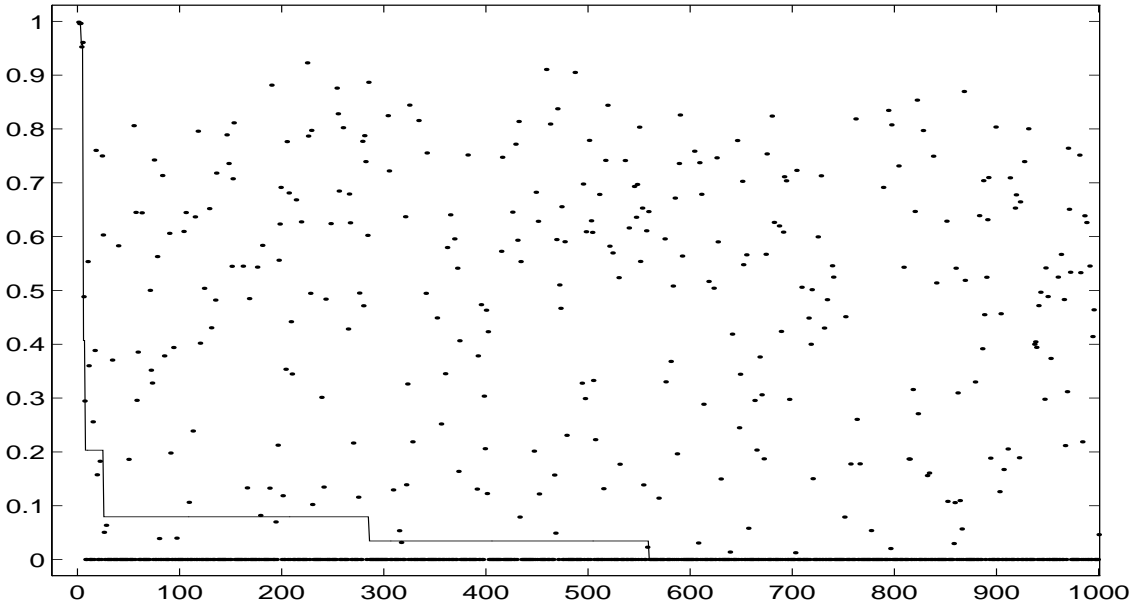
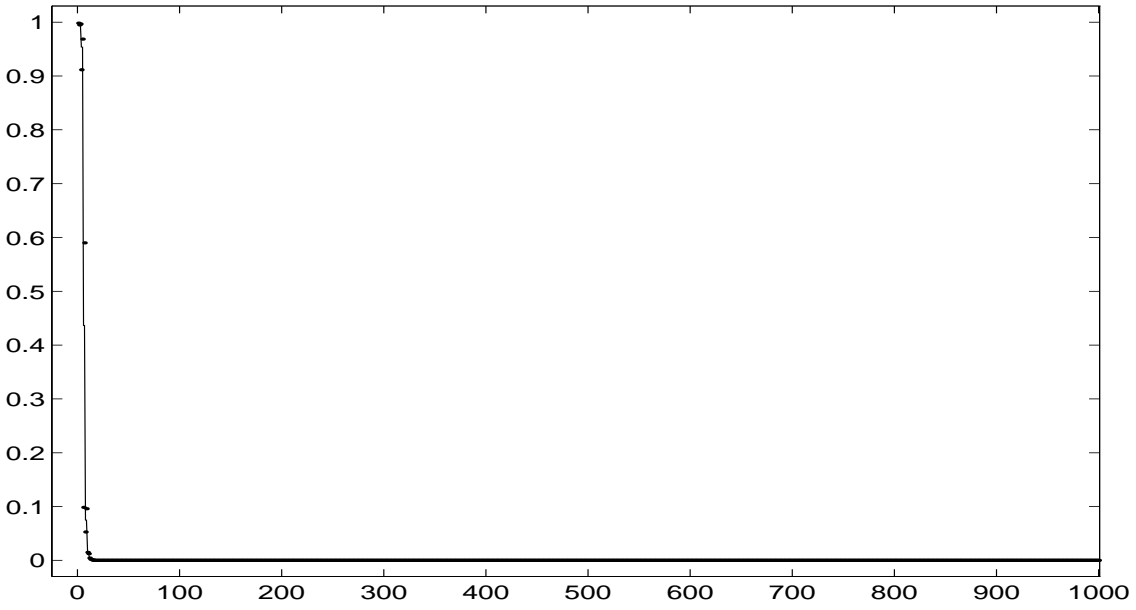


Figure 1b. Estimator  $\hat{\mu}(\cdot/n)$  (line) and true function  $\mu(\cdot/n)$  (broken line).



**Figure 1a'**.  $Y$  and  $\hat{\mu}(\cdot/n)$  with  $\hat{\sigma} = 0.9$  (1st plot) and  $\hat{\sigma} = 1.1$  (2nd plot).



**Figure 1c.** Greedy and monotone modulators:  $\tilde{g}$  and  $\tilde{f}$  (1st plot);  $\hat{g}^+$  and  $\hat{f}$  (2nd plot).



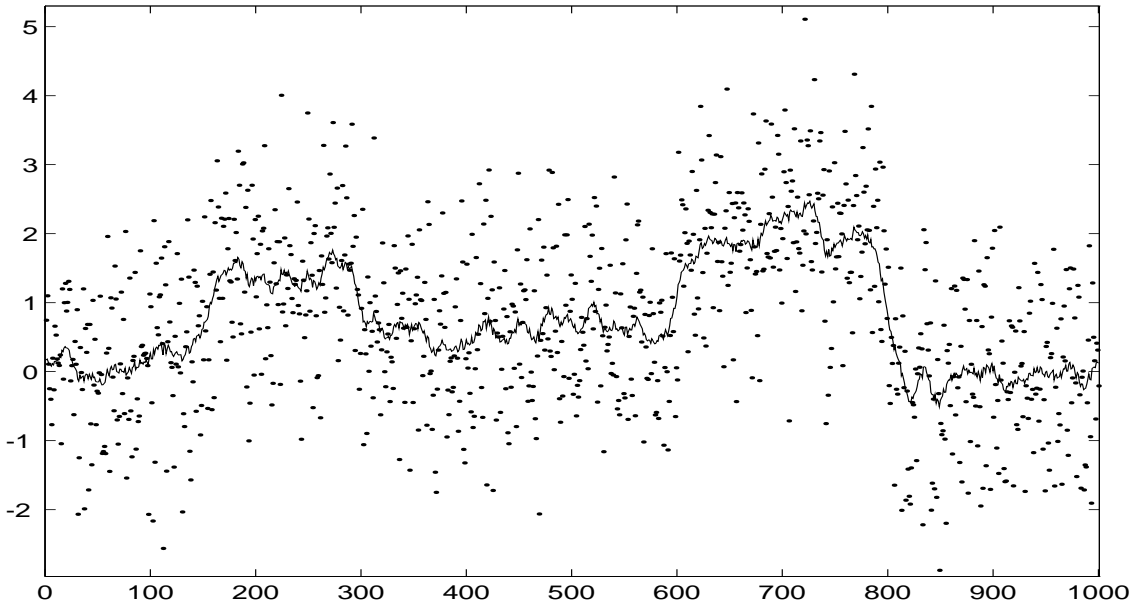


Figure 2a.  $Y$  and  $\hat{\mu}(\cdot/n)$ .

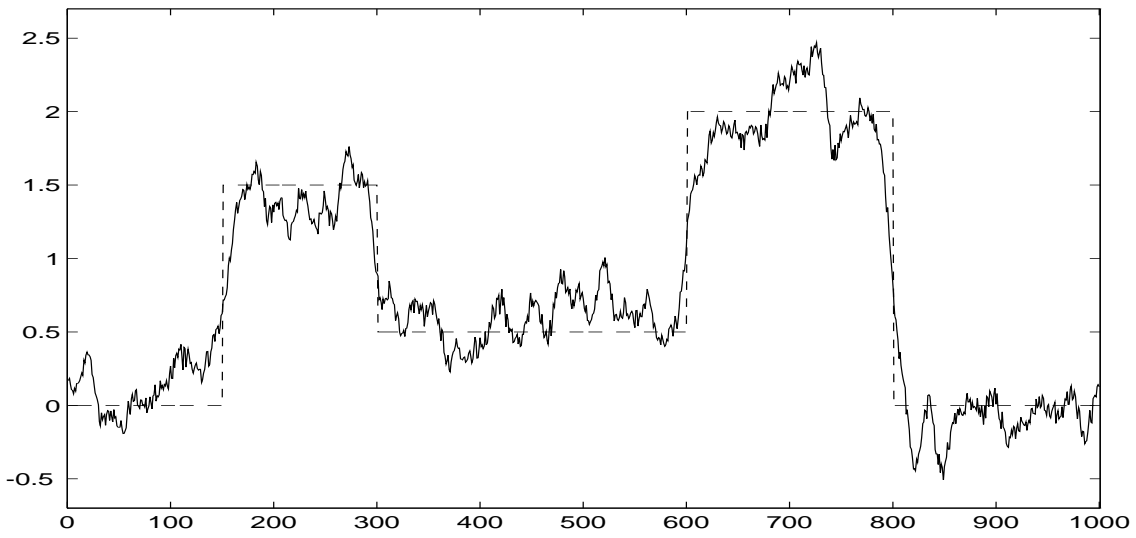
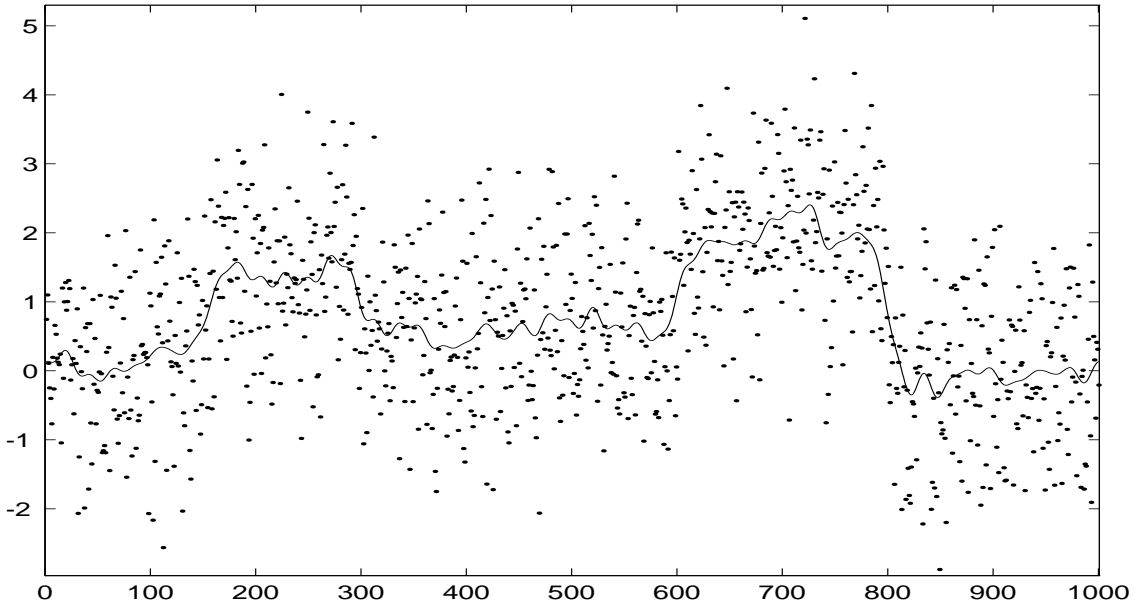
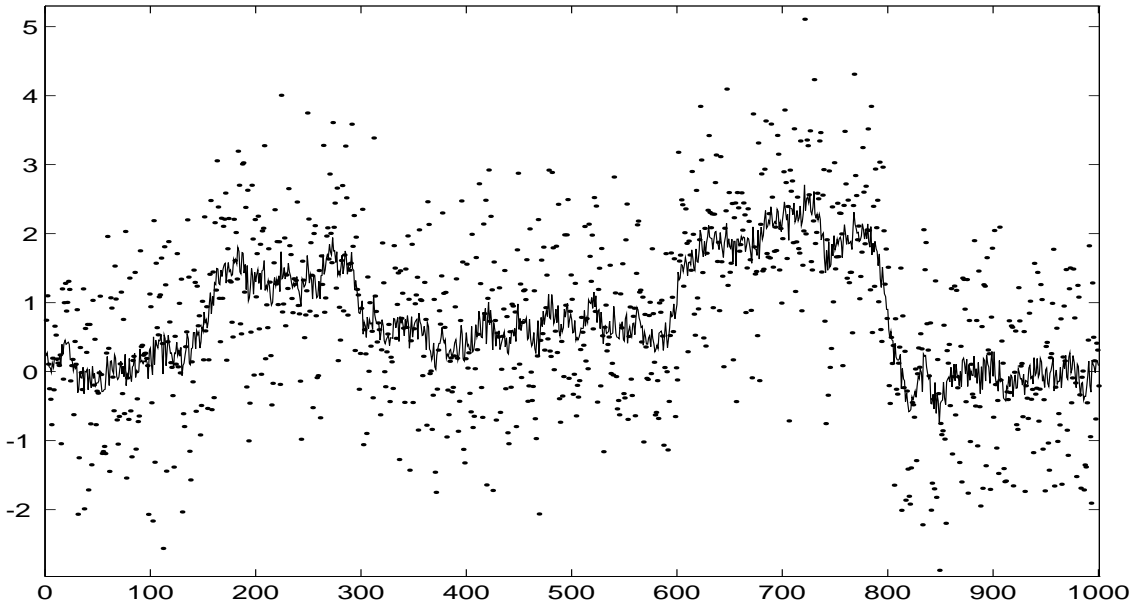
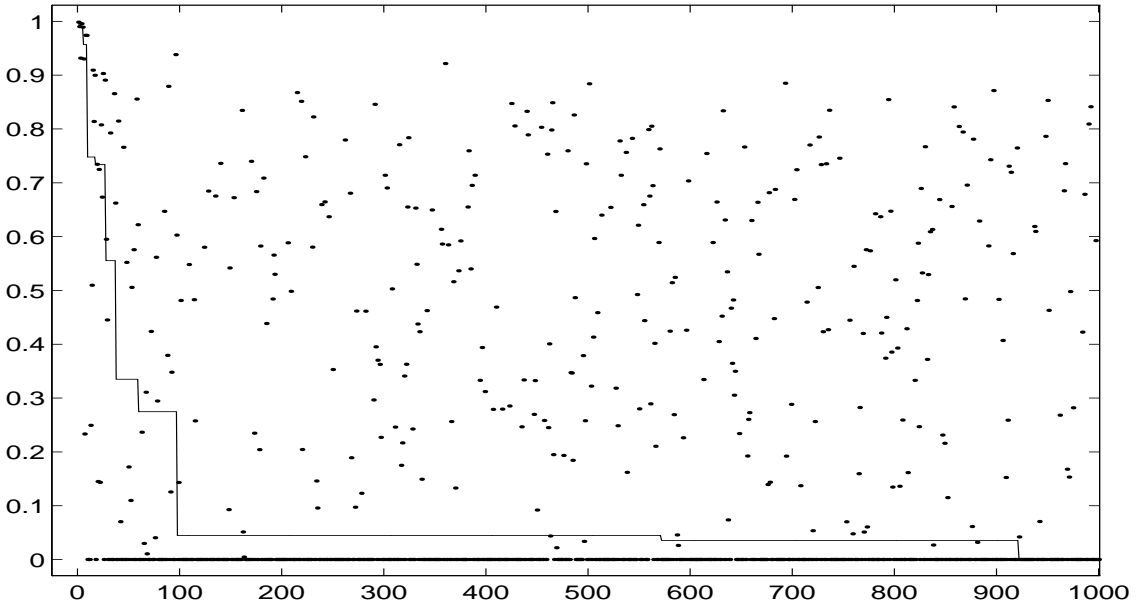
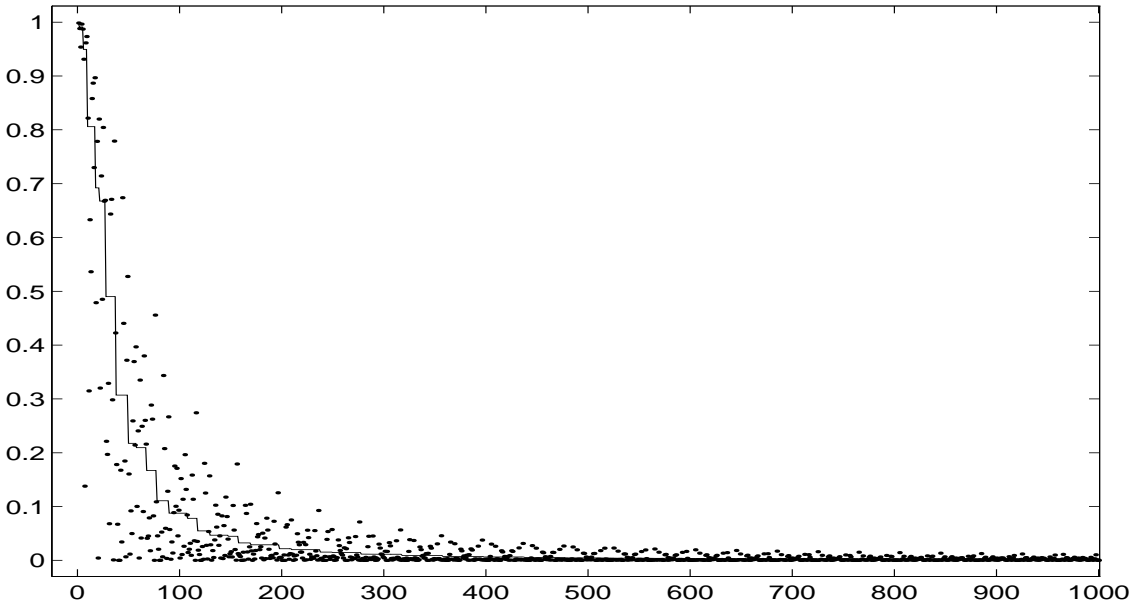


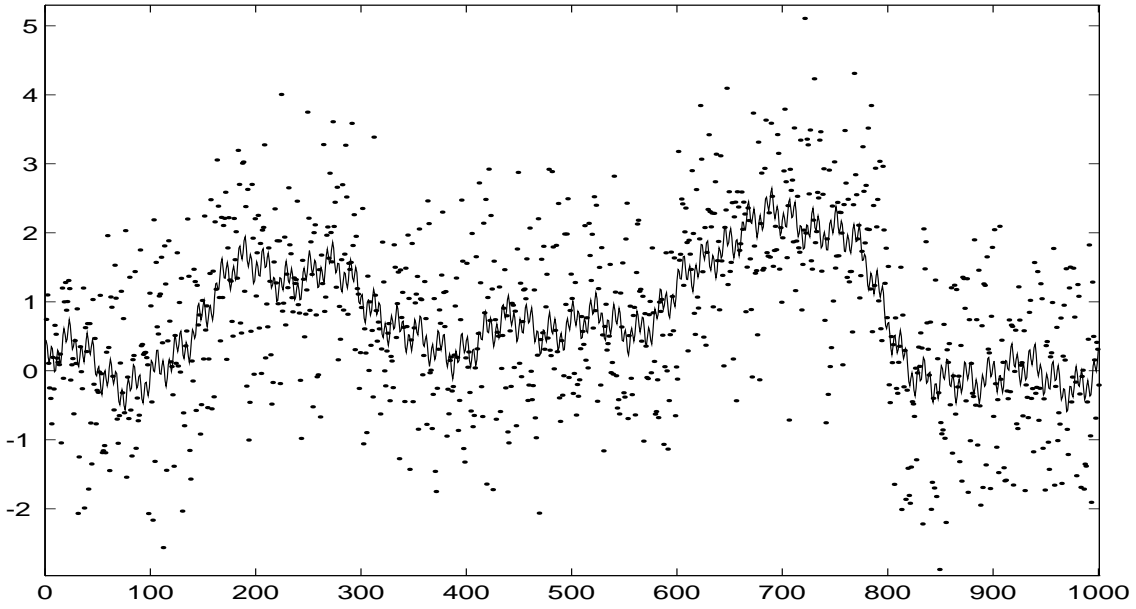
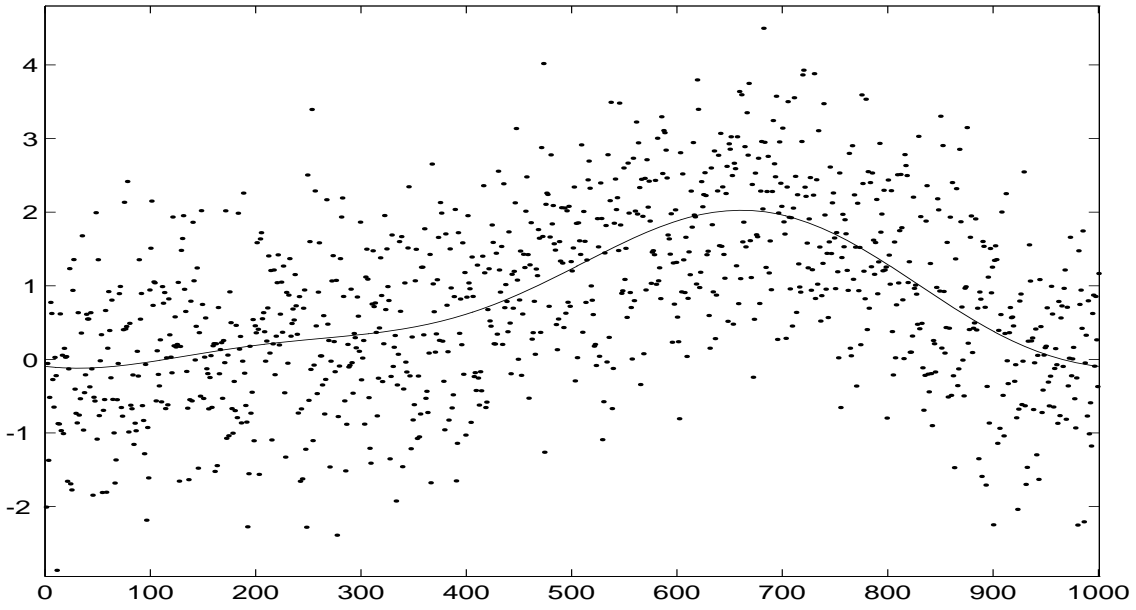
Figure 2b. Estimator  $\hat{\mu}(\cdot/n)$  (line) and true function  $\mu(\cdot/n)$  (broken line).



**Figure 2a'**.  $Y$  and  $\hat{\mu}(\cdot/n)$  with  $\hat{\sigma} = 0.9$  (1st plot) and  $\hat{\sigma} = 1.1$  (2nd plot).



**Figure 2c.** Greedy and monotone modulators:  $\tilde{g}$  and  $\tilde{f}$  (1st plot);  $\hat{g}^+$  and  $\hat{f}$  (2nd plot).



**Figure 3.** Oracle threshold estimator  $\tilde{\mu}_{\text{th}}(\cdot/n)$  for Case 1 (1st plot) and Case 2 (2nd plot).