

# Nonlinear Black-Box Models in System Identification : Mathematical Foundations

Anatoli Juditsky, Håkan Hjalmarsson,  
Albert Benveniste, Bernard Delyon,  
Lennart Ljung,  
Jonas Sjöberg, and Qinghua Zhang \*

March 24, 1995

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Basic mathematical problems . . . . .	2
1.2	Basic principles and limiting factors . . . . .	5
1.3	Structure of the paper . . . . .	6
<b>2</b>	<b>Approximation in function spaces</b>	<b>6</b>
2.1	Linear approximation schemes . . . . .	7
2.2	Besov spaces and classes of locally spiky and jumpy functions . . . . .	9
2.2.1	Spline approximations in Besov spaces . . . . .	10
2.2.2	Wavelet approximations in Besov spaces . . . . .	11
<b>3</b>	<b>Approximation in high-dimensional spaces</b>	<b>14</b>
3.1	The curse of dimensionality . . . . .	14
3.2	Function classes of lower effective dimension . . . . .	14
3.3	Neural networks . . . . .	15
3.4	Wavelets . . . . .	16
3.5	Breiman's Hinging Hyperplanes . . . . .	18
<b>4</b>	<b>Performance measures for non-parametric estimators</b>	<b>19</b>
4.1	Lower bounds for best achievable performance. . . . .	19
4.2	Some negative results. . . . .	20
4.3	Some positive results. . . . .	21
<b>5</b>	<b>Estimation in classes of uniformly smooth functions</b>	<b>22</b>
5.1	Kernel estimators for regression functions and densities . . . . .	22
5.2	Piecewise-polynomial estimators . . . . .	25
5.3	Projection estimates . . . . .	26
5.4	Selecting model complexity . . . . .	28

---

\*LL, HH, and JS are with Department of Electrical Engineering, Linköping University, S-581 83 Linköping, Sweden; AB, BD, AJ, QZ are with IRISA/INRIA, Campus de Beaulieu, 35042 RENNES cedex, FRANCE; e-mail: `name@isy.liu.se` and `name@irisa.fr`

<b>6</b>	<b>Estimation in classes of locally spiky and jumpy functions</b>	<b>29</b>
6.1	Spatial adaptivity . . . . .	29
6.2	Wavelet shrinkage algorithms . . . . .	30
<b>7</b>	<b>Application to non-parametric autoregression identification</b>	<b>34</b>
<b>8</b>	<b>Estimation of high-dimensional systems</b>	<b>35</b>
8.1	Wavelets . . . . .	35
<b>9</b>	<b>Conclusion : the gap between theory and everyday practice</b>	<b>36</b>

## Abstract

In this paper we discuss several aspects of the mathematical foundations of non-linear black-box identification problem. As we shall see that the quality of the identification procedure is always a result of a certain trade-off between the expressive power of the model we try to identify (the larger is the number of parameters used to describe the model, more flexible would be the approximation), and the stochastic error (which is proportional to the number of parameters). A consequence of this trade-off is a simple fact that good approximation technique can be a basis of good identification algorithm. From this point of view we consider different approximation methods, and pay special attention to *spatially adaptive* approximants. We introduce wavelet and “neuron” approximations and show that they are spatially adaptive. Then we apply the acquired approximation experience to estimation problems. Finally, we consider some implications of these theoretic developments for the practically implemented versions of the “spatially adaptive” algorithms.

**Keywords:** non-parametric identification, nonlinear systems, neural networks, wavelet estimators.

## 1 Introduction

The problem we are addressing in this paper is how to infer relationships between past input-output data and present/future outputs of a system when very little a priori knowledge is available. This is known as black-box modeling. There is a rich and well established theory for black-box modeling of linear systems, see *e.g.* (Ljung, 1987) and (Söderström and Stoica, 1989). It is not until the last few years that modeling and identification of non-linear systems have attracted a wide interest in the control community. Up to today, almost all attention has been concentrated on one single structure; neural networks. However, non-linear modeling have been studied for long in the statistics community where it is known under the label *non-parametric regression*. This area is quite rich and numerous methods exist. The purpose of this paper is to give an exposition of presently available techniques of non-linear modeling in a fairly unified and structured way. It is geared towards the mathematical foundations and exposes basic principles as well as presently available mathematical results. This exposition is, however, not exhaustive, neither with respect to presentation of existing methods nor with respect to mathematical results. In the companion paper (Sjöberg et al., 1995) the user aspects and the algorithmic aspects are extensively discussed.

### 1.1 Basic mathematical problems

The basic problem that we will address throughout this paper is now precisely stated. We first state the general problem and then we discuss some features of dynamic system modeling.

## The general problem

**Problem 1 (non-parametric regression)** Let  $(X, Y)$  be a pair of random variables with values in  $\mathcal{X} = \mathbf{R}^d$  and  $\mathcal{Y} = \mathbf{R}$  respectively. A function  $f : \mathcal{X} \mapsto \mathcal{Y}$  is said to be the **regression function of  $Y$  on  $X$**  if

$$\mathbf{E}(Y|X) = f(X) . \quad (1)$$

A typical case is  $Y = f(X) + e$ , where  $e$  is zero mean and independent of  $X$ . For  $N \geq 1$ ,  $\hat{f}_N$  shall denote an estimator of  $f$  based on the random sample  $\mathcal{O}_1^N = \{(X_1, Y_1), \dots, (X_N, Y_N)\}$  of size  $N$  from the distribution of  $(X, Y)$ , i.e., a map

$$\hat{f}_N : \mathcal{O}_1^N \mapsto \hat{f}_N(\mathcal{O}_1^N, \cdot) \quad (2)$$

where, for fixed  $\mathcal{O}_1^N$ ,  $x \mapsto \hat{f}_N(\mathcal{O}_1^N, x)$  is an estimate of the regression function  $f(x)$ . The family of estimators  $\hat{f}_N$ ,  $N \geq 1$  is said to be **parametric** if  $\hat{f}_N \in F$  for all  $N \geq 1$ , where  $F$  is some set of functions which are defined in terms of a fixed number of unknown parameters. Otherwise the family of estimators  $\hat{f}_N$ ,  $N \geq 1$  is said to be **non-parametric**.

For the sake of convenience, we shall often refer to  $X$  and  $Y$  as the *input* and *output* respectively (although they do not need to be such in actual applications).

Intuitively, the difference between the output and the regression function,  $Y - \mathbf{E}(Y|X)$ , is the part of the output that cannot be predicted from past data.

Two typical problems are considered in the statistical literature, namely the

- *non-parametric regression with random design (or sampling)*, where it is assumed that the variables  $X_i$  are random, independent, and identically distributed on  $[0, 1]^d$  with density  $g(x)$ , and the
- *non-parametric regression with deterministic design (or sampling)*, where it is assumed that the input variables  $X_i$  are nonrandom; the simplest case of deterministic design is the *regular design*, where the inputs  $X_i$  form a regular grid (for instance,  $f : \mathbf{R} \rightarrow \mathbf{R}$  and  $X_i = i/N$ ).

In the remainder of this section we consider the random design only.

## Non-parametric regression with dynamics.

Consider the following dynamical system:

$$Y_i = f(\Phi_i) + e_i, \quad i = 1, \dots, N ,$$

where  $Y_i \in \mathbf{R}$  and  $\Phi_i \in \mathbf{R}^d$  are observed, and  $e_i$  is a white noise as above. We assume that

$$\Phi_i = (Y_{i-1}, \dots, Y_{i-m}; U_i, \dots, U_{i-p}) , \quad (3)$$

where  $U_i \in \mathbf{R}$  denote the inputs ( $m + p = d$ ). For example, if  $\Phi_i = (Y_{i-1}, \dots, Y_{i-d})$ , then

$$Y_i = f(Y_{i-1}, \dots, Y_{i-d}) + e_i . \quad (4)$$

In analogy with the corresponding parametric model we call this system a *non-parametric autoregression* or a *functional autoregression* of dimension  $d$  (NAR( $d$ )). As an interesting application, we can consider a simple controlled NAR model for adaptive control :

$$Y_i = f(\Phi_i) + U_i + e_i , \quad (5)$$

where  $\Phi_i = (Y_{i-1}, \dots, Y_{i-m})$ , and  $U_i$  is the control. The following question can be considered : How to choose the control ( $U_i$ ) for the system (5) to track some reference trajectory  $y = (y_i)$ , or, at least, how to choose  $U_i$  in order to minimize  $\mathbf{E}Y_i^2$ , or, simply, to stabilize the system (5) ? If the function  $f(\Phi)$  was known, we could use the control

$$U_i = -f(\Phi_i) \quad (6)$$

to obtain  $Y_i = e_i$ . Clearly, this is a “minimum variance” control, since  $\mathbf{E}Y_i^2 \geq \sigma_e^2 = \mathbf{E}e_i^2$ . If  $f$  is unknown, a possible solution consists in performing non-parametric “certainty equivalence control”: compute an estimate  $\hat{f}_N$  of the regression function  $f$  based on the observations of the input/output pair  $(\Phi_i, Y_i - U_i)$ , and then take

$$U_i = -\hat{f}_i(\Phi_i) . \quad (6)$$

To analyze the certainty equivalence control (6), let us consider the control cost

$$Q_N = \frac{1}{N} \sum_{i=1}^N Y_i^2 = \frac{1}{N} \sum_{i=1}^N (f(\Phi_i) - \hat{f}_i(\Phi_i))^2 + \frac{1}{N} \sum_{i=1}^N e_i^2 .$$

It is easily checked that

$$\mathbf{E}(\hat{f}_i(\Phi_i) - f(\Phi_i))^2 \rightarrow 0 \text{ when } i \rightarrow \infty \quad (7)$$

implies  $\mathbf{E}Q_N \rightarrow \sigma_e^2$ , and  $\hat{f}_i(\Phi_i) - f(\Phi_i) \rightarrow 0$  a.s. implies  $Q_N \rightarrow \sigma_e^2$  a.s. Thus condition (7) is instrumental in analyzing this problem, and we shall informally discuss how it can be guaranteed.

Denote by  $\Phi_0^{i-1} = (\Phi_0, \dots, \Phi_{i-1})^T$  the vector of all available inputs up to time  $i-1$ , and by  $\varphi_0^{i-1} = (\varphi_0, \dots, \varphi_{i-1})^T$  the corresponding vector of integration variables. Let  $\mathbf{P}(\cdot)$  denote the distribution of the vector sequence  $(\Phi_i)$  when driven by the unknown “true” model (5)–(6), let for some  $1 \leq k \ll i$   $\mathbf{P}_{\Phi_0^{i-k}}(\cdot)$  be a distribution of  $\Phi_0^{i-k}$ , and let  $\mathbf{p}_{\Phi_i|\Phi_0^{i-k}}(\cdot)$  be a conditional density of the distribution of  $\Phi_i$  given  $\Phi_0^{i-k}$  (we assume that such a density exists). We have

$$\mathbf{E}|\hat{f}_i(\Phi_i) - f(\Phi_i)|^2 \sim \int |\hat{f}_i(x) - f(x)|^2 \mathbf{p}_{\Phi_i|\Phi_0^{i-k}}(x) dx \mathbf{P}_{\Phi_0^{i-k}}(d\Phi_0^{i-k}) .$$

Note that, if the closed-loop system (5)–(6) is stable, one would reasonably take equal weights for the observations  $\Phi_0, \dots, \Phi_i$  in the estimate  $\hat{f}_i$ . In such a case the estimate  $\hat{f}_i(\Phi)$  is asymptotically (as  $i \rightarrow \infty$ ) slowly varying, i.e.,  $\hat{f}_i \sim \hat{f}_{i-k}$ . On the other hand, the conditional density  $\mathbf{p}_{\Phi_i|\Phi_0^{i-k}}$  converges exponentially fast to the density  $\mathbf{p}_\Phi(x)$  of the invariant distribution of the Markov chain  $(\Phi_i)$  (again, we suppose that the correspondent quantities exist). Thus we can write

$$\mathbf{E}|\hat{f}_i(\Phi_i) - f(\Phi_i)|^2 \sim \int \mathbf{P}_{\Phi_0^{i-k}}(d\Phi_0^{i-k}) \int |\hat{f}_{i-k}(x) - f(x)|^2 \mathbf{p}_{\Phi_i|\Phi_0^{i-k}}(x) dx ,$$

and

$$\mathbf{E}|\hat{f}_i(\Phi_i) - f(\Phi_i)|^2 \sim \int \mathbf{E} \int |\hat{f}_i(x) - f(x)|^2 \mathbf{p}_\Phi(x) dx \equiv R_i(\hat{f}, f) .$$

Thus, as a conclusion, in any case, the crux in analyzing this adaptive minimum variance nonlinear control consists in getting bounds for the error in estimating the unknown function  $f$ . Note that the error measure we use in this case (the risk  $R_i(\hat{f}, f)$ ) is rather specific: *the error norm is weighted with the density of observation*. Hence, in addition to proving consistency for the estimates, getting such bounds is an important question.

REMARKS. The above discussion can be summarized as follows :

1. Non-parametric estimation of regression functions is instrumental in various problems such as adaptive identification and control.
2. Averaged  $L_p$ -norms of estimation error for various  $p$ 's are natural candidates as a figure of merit.
3. Having bounds for the estimation error is of paramount importance. This has been illustrated on the adaptive control example.

## 1.2 Basic principles and limiting factors

There are two factors that limit the accuracy with which the regression function  $f$  can be determined. Firstly, only a finite number of observation points<sup>1</sup>  $(x_k)_{k=1}^N$  are available. This means that  $f(x)$ , at other points  $x$  than those which are observed, must be obtained from the observed points by interpolation or extrapolation. Secondly, at the points of observation,  $x_k, k = 1, \dots, N$ ,  $f(x_k)$  is observed with an additive noise  $e_k = y_k - f(x_k)$ . Clearly, the observation noises  $e_i$  introduce a random component in the estimation error. A general approach to the problem is the following: we first choose an approximation method, i.e. substitute the function in question by its approximation; then we estimate the parameters involved in this approximation. This way we reduce the problem of function estimation to that of parametric estimation, though the number of parameters we have to estimate is not bounded a priori and can be large. To limit the number of parameters some *smoothness* or *regularity* assumptions have to be stated concerning  $f$ . Generally speaking, smoothness conditions require that the unknown function  $f$  belongs to a particular restricted functional class.

To see how the stochastic and deterministic approximation errors are combined in the estimation problem, consider, for simplicity, the following example. Suppose that  $N$  noisy observations of an unknown function  $f : \mathbf{R} \rightarrow \mathbf{R}$  are available:

$$Y_i = f(X_i) + e_i . \quad (8)$$

Suppose that  $f$  can be expressed in the form of some infinite expansion

$$f(x) = \sum_{j=0}^{\infty} \theta_j^* g_j(x) , \quad (9)$$

where  $(g_j(x))$ ,  $j = 0, \dots$  is a known family of basis functions. This is our approximated model. Our “smoothness” assumption about  $f$  is that the coefficients  $\theta_j^*$  are assumed to decrease in a certain way as  $j \rightarrow \infty$ . Thus the problem of estimating  $f$  reduces to the estimation of a suitable truncation of the vector of all parameters  $\Theta^* = (\theta_0^*, \dots)^T$ , using the observations  $(X_i, Y_i)$ . An estimate  $\hat{\Theta}_N$  of  $\Theta^*$  can be obtained by minimizing the following criterion

$$\sum_{k=1}^N \|Y_k - \hat{\Theta}_N^T \mathbf{g}(X_k)\|^2$$

where  $\|\cdot\|$  is some suitable norm (say Euclidean for instance), and the row vector  $\mathbf{g}(x)$  collects the components  $g_j(x)$  corresponding to the  $\theta_j$  selected in  $\hat{\Theta}_N$ . Suppose that  $(X_k)$  is a realization of a stationary stochastic process, then the following holds asymptotically in  $N \rightarrow \infty$ :

$$\mathbf{E}\|Y - \hat{\Theta}_N^T \mathbf{g}(X)\|^2 \approx \underbrace{\mathbf{E}\|e\|^2}_{\text{noise}} + \underbrace{\|(\Theta - \mathbf{E}\hat{\Theta}_N)^T \mathbf{g}(X)\|^2}_{\text{bias}} + \underbrace{\mathbf{E}\|(\mathbf{E}\hat{\Theta}_N) - \hat{\Theta}_N\|^2 \mathbf{g}(X)\|^2}_{\text{variance}}$$

---

<sup>1</sup>We distinguish between the random variables  $(Y_k, X_k)$  and the corresponding observations  $(y_k, x_k)$

As we shall see later, usually,  $\mathbf{E}\hat{\theta}_j = \theta_j^*$ , and the bias term in the right-hand side does not depend on the data record. Thus the “bias” is in fact the approximation error due to the truncation of the infinite vector  $\Theta^*$ . Let  $n$  be the dimension of  $\hat{\Theta}_N$ . Increasing  $n$  would reduce the bias down to zero. But the variance term is typically  $O(n/N)$ , thus increasing  $n$  increases the variance term as well. The optimum occurs when both bias and variance terms are balanced. In Section 5 — Section 8 we shall see several examples of how this compromise is met. For a further empiric discussion of the bias/variance trade-off and its implications for identification, the reader is referred to (Sjöberg et al., 1995).

The bottom line is that, to cope with the bias variance trade-off, it is important to use an efficient approximation technique, *i.e.* one that gives a small approximation error with few parameters. However, *which approximation method is effective depends on the particular function class to which the function is assumed to belong*. Since guessing appropriate function classes requires prior information which is hardly accessible to the engineer, *it is important to come up with approximation methods that are flexible and are as independent as possible from the particular function class*. This will be a recurring theme throughout the paper.

### 1.3 Structure of the paper

The paper is organized as follows. In the next section different smoothness classes are discussed together with associated approximation techniques from the point of view of their utility for estimation. “Spatially uniform” smoothness classes as well as classes of functions with sparse singularities are considered. For the latter classes it is shown that wavelets play an important role. There is a particular problem associated with the approximation of functions of a large number of input variables. How well observations data fill the input space decreases exponentially with the input dimension. Hence, it is necessary to further restrict the function class if approximations with reasonable accuracy is to be obtained for large input dimension and moderate sample sizes. This topic is given special attention in Section 3. Neural network, wavelets and other methods are discussed from this perspective.

The estimation problem is treated in Section 4 — Section 8. First, in Section 4, performance criteria are introduced; note that new tools are required since concepts such as Cramer-Rao bounds and Fisher Information matrix are not appropriate for non-parametric estimation. Section 5 discusses the estimation of uniformly smooth functions. Classical techniques such as kernel, piecewise-polynomial and projection estimates are reviewed. Techniques to estimate non-uniformly smooth functions are dealt with in Section 6 while estimation of highly multivariate functions is considered in Section 8. Finally, as a conclusion, the gap between theory and everyday practice is discussed in Section 9.

## 2 Approximation in function spaces

As we have seen in Section 1, due to the bias/variance trade-off, the number of parameters used in the expression of the regression function for estimation has to be kept as small as possible. Thus approximation methods performing good approximation with few parameters will be preferred. Not surprisingly, the approximation method should be selected according to the prior assumptions on the function class.

From an application point of view, these classes can be categorized as being either classes of uniformly smooth functions or classes of locally spiky and jumpy functions. Many real-life nonlinear systems are smooth with sparse singularities, e.g., mechanical systems, chemical systems, etc. Thus classes of locally spiky and jumpy functions are important in practice. A

basic problem is that one typically does not know which function class the unknown function belongs to. However, as we shall see, it is possible to come up with one single approximation scheme which has good approximation properties for a wide family of function classes covering both uniformly smooth classes and locally spiky and jumpy classes. In the sequel, we move progressively from simple to rather complex approximation problems.

## 2.1 Linear approximation schemes

We mean here approximation schemes that are linear in  $f$ , i.e., satisfy  $(f + g)_n = f_n + g_n$  if  $f_n$  denotes the  $n$ -th order approximant of  $f$ . Typical examples are of the following form. We have some function space  $\mathcal{F}$  and an increasing family of (closed) subspaces  $\mathcal{F}_n$  converging to  $\mathcal{F}$ . Then we consider some norm  $\|\cdot\|$  on  $\mathcal{F}$ . The  $n$ th-order *projection approximation* of  $f \in \mathcal{F}$  is the  $f_n \in \mathcal{F}_n$  minimizing the distance of  $f$  to the subspace  $\mathcal{F}_n$ . We call this type of approximations the projection approximants.

While such projections do not need to be associated with Hilbert space structures, the problem of determining an optimal approximation becomes particularly simple when the functional space is a Hilbert space. In that case, the best approximant is obtained by simple orthogonal projection of  $f$  onto some subspace. The following result can be obtained (cf, (Pinkus, 1985)).

**Proposition 1 (Optimal approximants in Hilbert spaces )** *Let  $\{g_k\}_{k=1}^\infty$  be an orthonormal basis for a Hilbert space  $\mathcal{H}$ . Let  $\{\mu_k\}_{k=1}^\infty$  be a sequence of non-increasing positive numbers. Consider the function class*

$$\mathcal{F} = \left\{ \sum_{k=1}^\infty c_k g_k : \sum_{k=1}^\infty \left| \frac{c_k}{\mu_k} \right|^2 \leq 1 \right\}. \quad (10)$$

*Then for any  $f \in \mathcal{H}$ , and  $f_n = \sum_{j=0}^n \mu_j g_j$ , we have  $\|f - f_n\|_{\mathcal{H}} \leq \mu_{n+1}$ .*

COMMENTS :

1. The convergence rate in Proposition 1 in some sense cannot be improved: there are “bad” functions in the space  $\mathcal{H}$  which lie at the distance  $\mu_{n+1}$  from any  $n$ -dimensional linear subspace.
2. If  $\mu_n \equiv 1, n = 1, 2, \dots$ , then  $\mathcal{F} = \mathcal{H}$ , and proposition 1 states that that the worst-case approximation error will not decrease when  $n$  is increased!! Indeed, in this case the corresponding “bad” function can be easily constructed: we can take simply  $f = g_{n+1}$  — such a situation occurs for instance if we take  $s = 0$  in the example below. This cannot happen however when the coefficients  $\mu_i$  vanish at a fixed rate when  $i \rightarrow \infty$ , i.e., when  $\mathcal{F}$  is compact in  $\mathcal{H}$ .
3. The optimal approximation  $f_{n+1}$  can be computed recursively from  $f_n$ .

As an application of Proposition 1, consider the following class of functions. Let  $\mathcal{W}_2^s(L)$  be the set of 1-periodic functions  $f(x)$ , which are defined by their Fourier series

$$f(x) = \sum_{j=1}^\infty c_j \varphi_j(x), \quad (11)$$

where  $\varphi_{2k}(x) = \sqrt{2}\sin(2\pi kx)$  and  $\varphi_{2k+1}(x) = \sqrt{2}\cos(2\pi kx)$ ,  $k = 1, \dots$  and where the Fourier coefficients satisfy

$$\sum_{j=1}^{\infty} |c_j|^2 (1 + |j|^{2s}) < L^2. \quad (12)$$

For this class of functions, the optimal approximant is given by  $f_n = \sum_{j=0}^n c_j \varphi_j$  and the rate of convergence is  $n^{-s}$ . Notice that  $\mathcal{W}_2^s(L)$  is a smoothness class. In fact (15) is one of the several equivalent definitions of the Sobolev class  $\mathcal{W}_2^s(L)$ , a particular case of the classes  $\mathcal{W}_p^s(L)$ ,  $p \geq 0$  which can be defined, for instance for  $s < 1$ , as the subset of the functions  $f \in L_1$ , such that

$$\frac{\|f(t+h) - f(t)\|_p}{|h|^s} \leq L. \quad (13)$$

As we have seen in Proposition 1, the best approximation of functions belonging to  $\mathcal{W}_2^s$ , when the error is measured in  $L_2$ -norm, is simply an orthogonal projection on some linear subspace. This is *not* true for other Sobolev classes. It can be shown that, for the class  $\mathcal{W}_p^s$  with  $p < 2$ , the optimal projection on a subspace of dimension  $n$  converges to  $f$  with the rate  $n^{-s'/d}$ , where  $s' = s - 1/p + 1/2$ . We shall see, however, that there are different approximations, which exhibit a better convergence rate equal to  $n^{-s/d}$ . The difference between these two rates becomes significant for small  $p$ .

From the mathematical point of view the problem can be explained as follows: the Sobolev classes  $\mathcal{W}_p^s(L)$  with  $s - 1/p + 1/2 > 0$  are compact subsets of  $L_2$ . For  $p \geq 2$  these subsets are convex. In contrast, when  $p < 2$ , these classes are not convex, and can hardly be approximated using linear (and thus convex) subspaces.

From the users point of view, the functions from the Sobolev classes  $\mathcal{W}_p^s$  with small  $p$  are essentially classes of functions with sparse singularities, or classes with spatially “non-uniform” smoothness. Hence, if linear approximation schemes are used, the approximation rate for locally spiky and jumpy functions will be slower than for uniformly smooth functions.

**Discussion.** Roughly speaking, linear approximations schemes use subspaces (or basis functions) for approximation, which are independent from the particular function to be approximated. Thus the question arises if one could not do better by selecting the basis functions *adaptively*, i.e., depending on the function to be approximated. To illustrate this point consider the function  $f(x) = 1_{\{0 \leq x < a\}}$  for some  $0 < a < 1$ . The Fourier coefficients of this function are

$$c_0 = a, \quad c_{2k} = \sqrt{2} \frac{\sin^2(\pi k a)}{\pi k}, \quad c_{2k+1} = \sqrt{2} \frac{\sin(\pi k a) \cos(\pi k a)}{\pi k},$$

From Proposition 1 we know that this function belongs to the subset  $\mathcal{F}$  of  $L^2[0, 1]$  which consists of the functions such that the coefficients  $\mu_k$  in decomposition (10) decrease slower than  $k^{-1/2}$ . This would provide a convergence rate not better than  $n^{-1/2}$  for the orthogonal projection using  $n$  basis functions. However, one would naturally expect being able to design a procedure which focuses on detecting the edges of  $f$ , thus exhibiting a much better convergence rate. We shall see that such methods cannot be linear schemes but must be *spatially adaptive*, i.e. the basis functions must be able to adapt to the function to be approximated. Spatial adaptation is our next important topic. But first we introduce some suitable functional classes.



## 2.2 Besov spaces and classes of locally spiky and jumpy functions

A suitable family of spaces to deal with functions that are locally spiky and jumpy is that of the Besov spaces. This is *a family* of functional spaces indexed with 3 parameters (in a way that the family of Sobolev spaces  $W_p^s$  is indexed with two parameters:  $p$  and  $s$ ). The interplay of these three indices gives to this family of spaces a great flexibility. For instance, for different combinations of the indices, we can obtain spaces of “uniformly” regular (or smooth) functions, as well as spaces of regular functions with sparse singularities. However, all Besov spaces possess the following important property (wavelets and wavelet expansions will be introduced later):

♡ norms in these spaces are easily evaluated using coefficients of the wavelet expansions of the functions of these spaces.

Let us move now on stating precise definitions. For the sake of clarity we consider only compactly supported<sup>2</sup> functions  $f$ :  $\text{supp } f \subseteq [0, 1]^d$ , though all the definitions below can be generalized for non-compact and multi-dimensional cases (we recommend (Triebel, 1983) and (Triebel, 1993) as extremely complete presentations of the current state of the art in the theory of functional spaces).

For  $f \in L_1$  and  $M \in \mathbf{N}$  we define the local oscillation of order  $M$  and radius  $t$  at the point  $x \in [0, 1]$  by

$$\text{osc}_M f(x, t) \triangleq \inf_P \frac{1}{t^d} \int_{|x-y|<t} |f(y) - P(y)| dy, \quad (14)$$

where the infimum is taken over all polynomials  $P$  of degree less than or equal to  $M$ . This quantity measures the quality of the local fit of  $f$  by polynomials on balls of radius  $t$ . Select  $p, q > 0$ ,  $s > d(p^{-1} - 1)$ , and take  $M = \lfloor s \rfloor$ . The following set of functions:

$$\mathcal{B}_{pq}^s = \left\{ f \in L_{1 \wedge p} : \|f\|_{\mathcal{B}_{pq}^s} = \|f\|_p + \left( \sum_{j=1}^{\infty} (2^{js} \|\text{osc}_M f(x, 2^{-j})\|_p)^q \right)^{1/q} < \infty \right\}, \quad (15)$$

(with the usual modification for  $p$  or  $q = \infty$ ) is identical to the *Besov spaces* of functions (Besov, 1959), and it is shown in (Triebel, 1983) that  $\|\cdot\|_{\mathcal{B}_{pq}^s}$  is equivalent to the classical Besov norm.

COMMENTS:

1. The triple parameterization using  $s, p$ , and  $q$  provides a very accurate characterization of smoothness properties. As usual for Hölder or Sobolev spaces, the index  $s$  indicates how many derivatives are smooth. Then, for larger  $p$ ,  $\|f\|_{\mathcal{B}_{pq}^s}$  is more sensitive to details. Finally, index  $q$  has no useful practical interpretation, but it is a convenient instrument that serves to compare Besov spaces with the more usual Sobolev spaces  $\mathcal{W}_p^s$ , as indicated next. It is interesting to notice that the indicator functions of intervals belong to the spaces  $\mathcal{B}_{s-1\infty}^s$  for all  $s > 0$ , this illustrates our claim in the title of this subsection.
2. It can be shown that (cf (Triebel, 1983)) for  $s \geq 0$ ,  $0 \leq p, q \leq \infty$ :
  - The family of Besov spaces includes some more classical spaces. for  $s$  non integer, Hölder classes  $\mathcal{C}^s = \mathcal{B}_{\infty\infty}^s$ , and Sobolev spaces  $\mathcal{W}_2^s = \mathcal{B}_{22}^s$ ;
  - $\mathcal{B}_{pq}^s \subset \mathcal{B}_{p'q'}^{s'}$  if  $p' \geq p$ ,  $q' \geq q$ ,  $s' \leq s - \frac{d}{p} + \frac{d}{p'}$  (strict inequality if  $p = \infty$ );

---

<sup>2</sup>  $\text{supp } f$  will be used to denote the support of  $f$ , *i.e.* the set where  $f$  is non-vanishing

- $\mathcal{B}_{pq}^0 \subseteq L_p \subseteq \mathcal{B}_{pq'}^0$ , where  $q = 2 \wedge p$  and  $q' = 2 \vee p$ ;
- $\mathcal{B}_{pp}^s \subset \mathcal{W}_p^s \subset \mathcal{B}_{p2}^s$  for  $p \leq 2$ ;
- $\mathcal{B}_{p2}^s \subset \mathcal{W}_p^s \subset \mathcal{B}_{pp}^s$  for  $p \geq 2$ .

In particular, if  $s > d/p$ , then  $\mathcal{B}_{pq}^s \subset \mathcal{C}$ .

### 2.2.1 Spline approximations in Besov spaces

We consider the  $d$ -dimensional case and  $\text{supp } f \subseteq [0, 1]^d$ . *Free knots spline approximations* have been analyzed in (Petrushev and Popov, 1987) (Theorems 7.3 and 7.4) using Besov spaces. Recall that a function  $f_n$  is called a spline function on  $[0, 1]$  of order  $k$  with  $n$  knots if  $f_n \in \mathcal{C}^{k-2}$  and there exist points (knots)  $0 = x_0 < x_1 \leq x_2 \leq \dots \leq x_{n-1} \leq x_n = 1$  such that  $f_n$  is an algebraic polynomial of degree  $k - 1$  in each interval  $(x_{i-1}, x_i)$ . Therefore, a spline is a smooth piecewise polynomial function. Free knot spline approximants are not linear schemes. The following result shows that any function from Besov space can be nicely approximated by splines. More surprisingly, any function which have a good spline approximant belongs to a certain Besov space.

We first state the so-called *Jackson inequality* for spline approximations. Consider  $f \in \mathcal{B}_{pq}^s$ ,  $p, q > 0$ . Then there exists a spline function with  $n$  free knots  $f_n$  such that the following bound holds:

$$\|f_n - f\|_u \leq C(s, p, q) n^{-s} \|f\|_{\mathcal{B}_{p\infty}^s}, \quad (16)$$

where  $u$  satisfies  $s - 1/p + 1/u > 0$ . The converse bound is provided by the *Bernstein inequality*: For any  $f \in L_u$ ,  $s - 1/p + 1/u = 0$ ,  $u < \infty$ ,

$$\|f\|_{\mathcal{B}_{pp}^s} \leq C(s, p, q) \left( 1 + n^s \inf_{f_n} \|f - f_n\|_u \right),$$

where the infimum ranges over the set of spline functions  $f_n$  of order  $k \geq s + 2$  with  $n$  free knots. A similar result holds in the multi-dimensional case and for  $n$ th-order *rational fraction approximations*, see Theorems 7.3 and 8.3 in (Petrushev and Popov, 1987).

In contrast, approximations using *fixed linear subspaces* perform poorly in Besov spaces. Consider some increasing family  $(\mathcal{L}_n)$  of  $n$ -dimensional linear subspaces of  $L_u$ ,  $u > p$ . Let  $f_n$  denote the linear projection of  $f \in \mathcal{B}_{pq}^s$  on  $\mathcal{L}_n$  using the  $L_u$ -norm. Then for any such family  $(\mathcal{L}_n)$ , there exists a least favorable  $f$  such that the following lower bound holds:

$$\|f - f_n\|_u \geq C n^{-s'} \|f\|_{\mathcal{B}_{uu}^{s'}}, \quad (17)$$

where  $s' = s - 1/p + 1/u$ . Note that the Sobolev injections stated at the end of section 2.2 imply that the same bounds hold for Sobolev classes  $\mathcal{W}_p^s$ .

Consider again the example of the indicator function  $f(x) = 1_{\{0 \leq x < a\}}$ . It can be easily verified that  $f \in \mathcal{B}_{s-1\infty}^s$  for any  $s > 0$ . On the one hand, it follows from (16) that  $f$  can be approximated using splines with  $n$  knots with asymptotic  $L_2$ -error of order  $o(n^{-l})$  for any  $l < \infty$  as  $n \rightarrow \infty$ . On the other hand, by (17), linear approximations of the same function have an  $L_2$ -error of order  $O(n^{-1/2})$ , where  $n$  is the dimension of the linear subspace. This remark would make rational approximations or splines with free knots very attractive for approximation in Besov spaces. Unfortunately, in the above result only the existence of such approximations is stated, and they are very hard to compute, for example, the optimal positioning of the knots of the spline approximation is very hard to find. It is amazing that *wavelet approximations are as good as spline ones, but are much more easily constructed*. We discuss this next.

### 2.2.2 Wavelet approximations in Besov spaces

The original objective of the theory of wavelets is to construct orthogonal bases of  $L_2(\mathbf{R})$  of the form  $(\psi(2^j x - k))_{j,k}$ , i.e. the bases which are constituted by translations and dilations of the same function  $\psi$ . It is preferable to take  $\psi$  localized and regular. We refer the reader to (Sjöberg et al., 1995) for the attractive computational features of orthonormal wavelet bases, and we concentrate here on the properties that are useful to understand how they perform for function approximation.

**The principle of wavelet construction** is the following: first construct a function  $\varphi \in L_2(\mathbf{R})$  such that

(S1) the functions  $\varphi(x - k)$  mutually orthogonal for  $k$  ranging over  $\mathbf{Z}$ .

(S2)  $\varphi$  is a *scale* function, i.e. there is a sequence  $h_k \in l_2$  such that

$$\varphi(x) = \sqrt{2} \sum_{k \in \mathbf{Z}} h_k \varphi(2x - k) . \quad (18)$$

This is an important step, several constructions of  $\varphi$  have been proposed (cf (Daubechies, 1992)). Next we define the *wavelet*

$$\psi(x) = \sqrt{2} \sum (-1)^k h_k \varphi(2x - k) . \quad (19)$$

It can be shown that the family  $\{\psi(2^j x - k), j \in \mathbf{N}, k \in \mathbf{Z}\}$  constitutes an orthonormal basis of  $L_2(\mathbf{R})$ . The crux in proving this property is to verify that, for any  $j_0$ , the family  $\{\varphi(2^{j_0} x - k), \psi(2^j x - k), k \in \mathbf{Z}, j \geq j_0\}$  also forms an orthogonal basis of  $L_2(\mathbf{R})$ ; this is achieved mainly by using the algebraic properties (S1), (S2), (19). If  $\varphi$  and  $\psi$  are compactly supported, they give us a local description, at different scales  $j$ , of the considered function:

$$f(x) = \sum_{k \in \mathbf{Z}} \langle f, \varphi_{j_0 k} \rangle \varphi_{j_0 k}(x) + \sum_{j \geq j_0, k \in \mathbf{Z}} \langle f, \psi_{jk} \rangle \psi_{jk}(x)$$

where  $\varphi_{jk}(x) = \varphi(2^j x - k)$ , and  $\langle \cdot, \cdot \rangle$  denotes inner product in  $L^2$ .

In what follows we assume that  $\varphi$  is a compactly supported piecewise continuous scale function satisfying the following condition:

$$\exists r > 0 : \varphi \in \mathcal{B}_{u\infty}^r , \quad (20)$$

and we move to the multidimensional case. Starting from  $\varphi$  one can construct the corresponding orthonormal basis of  $L_2(\mathbf{R}^d)$ , i.e. the functions  $\Phi, \Psi^{(i)}$ ,  $i = 1, \dots, 2^d - 1$  such that for any  $f \in L_2(\mathbf{R}^d)$  we have the formal expansion

$$\begin{aligned} f &= \sum_k \alpha_{0k} \Phi_{0k} + \sum_{j=0}^{\infty} \sum_{k \in \mathbf{Z}^d} \sum_{l=1}^{2^d-1} \beta_{jk}^{(l)} \Psi_{jk}^{(l)} \\ \alpha_{jk} &= \langle f, \Phi_{jk} \rangle, \quad \beta_{jk}^{(l)} = \langle f, \Psi_{jk}^{(l)} \rangle . \end{aligned} \quad (21)$$

Here

$$\left\{ \Phi_{0k}, \Psi_{jk}^{(l)} \right\}, \quad 0 \leq j < \infty, k \in \mathbf{Z}^d, 1 \leq l \leq 2^d - 1$$

is the corresponding orthonormal basis, formed by dilations and multi-dimensional translations of  $\Phi$  and  $\Psi^{(l)}$ . For details of definitions, the reader is again referred to (Sjöberg et al., 1995).

**Wavelet approximations.** We first state a result (Meyer, 1990) (Jaffard and Laurentçot, 1989) concerning functions that satisfy Hölder type conditions. Recall, that a function  $f$  is called *Hölder continuous with exponent  $s$  at point  $x_0$* , written  $f \in \mathcal{C}_{x_0}^s$ , if there is a polynomial  $P$  of degree at most  $\lfloor s \rfloor$  such that <sup>3</sup>

$$|f(x) - P(x - x_0)| \leq C|x - x_0|^s .$$

If  $f$  is Hölder continuous, with exponent  $s < r$  ( $r$  is the regularity of  $\varphi$  at  $x_0$ , see (20)), then there exists  $C < \infty$  such that, for any integer  $j > 0$ ,

$$\max_{\{k: x_0 \in \text{supp } \Psi_{jk}\}} \langle f, \Psi_{jk} \rangle \leq C 2^{-j(s+d/2)} . \quad (22)$$

Conversely, if (22) holds and  $f$  is known to be  $\mathcal{C}_{x_0}^\varepsilon$  for some  $\varepsilon > 0$ , then

$$|f(x) - P(x - x_0)| \leq C|x - x_0|^s \log \frac{2}{|x - x_0|} .$$

This result states that local smoothness of Hölder type can be characterized with the vanishing rate of the wavelet coefficients in the neighborhood of this point. This property is specific to the wavelet transform, and does not hold for other orthogonal bases. As we will see a similar property holds in Besov spaces. We have the following result (c.f. Theorem 4 in (Sickel, 1990)):

**Theorem 1 (Besov norms and wavelet decompositions)** *Let  $r > s > d(1/u - 1)$  and  $\varphi$  be a scale function satisfying conditions (20). For any  $f \in \mathcal{B}_{pq}^s$  define*

$$\|f\|_{spq} = \left( \sum_k |\alpha_k|^p \right)^{1/p} + \left( \sum_{j=0}^{\infty} \left[ 2^{j(s+d/2-d/p)} \|\beta_{j\cdot}\|_p \right]^q \right)^{1/q} \quad (23)$$

and  $\|\beta_{j\cdot}\|_p = (\sum_{l,k} |\beta_{jk}^{(l)}|^p)^{1/p}$ , see (21) for the definition of coefficients  $\alpha_k = \alpha_{0k}$  and  $\beta_{jk}^{(l)}$ . Then (23) is equivalent to the norm of Besov space  $\mathcal{B}_{pq}^s$ , i.e., there exist constants  $C_1$  and  $C_2$ , independent of  $f$ , such that

$$C_1 \|f\|_{\mathcal{B}_{pq}^s} \leq \|f\|_{spq} \leq C_2 \|f\|_{\mathcal{B}_{pq}^s} . \quad (24)$$

Theorem 1 states that norms in Besov spaces are suitably evaluated using orthonormal wavelet decompositions. This fact can be used to obtain very efficient approximations.

We now indicate how such a wavelet approximation of  $f$  can be constructed. Consider the full wavelet decomposition of  $f$ :

$$f(x) = \sum_{k \in \mathbf{Z}} \alpha_{0k} \Phi_{0k}(x) + \sum_{j=0}^{\infty} \sum_{k \in \mathbf{Z}^d} \sum_{l=1}^{2^d-1} \beta_{jk}^{(l)} \Psi_{jk}^{(l)}(x) . \quad (25)$$

1. Keep the projection of  $f$  on the subspace  $V_0$ , this corresponds to the left most sum in (25).

When  $f$  and  $\Phi$  are both compactly supported this requires computing only a fixed amount of coefficients, say  $m$ . And then

---

<sup>3</sup>recall that  $\lfloor s \rfloor$  denotes the largest integer  $\leq s$ .

2. Select in the second (triple) sum those coefficients  $\beta_\lambda$ ,  $\lambda = (i, j, k)$  with largest absolute value, denote by  $\Lambda$  the set of the  $n - m$  so selected wavelet coefficients. Finally
3. Add  $n - m$  detail terms  $\beta_\lambda \Psi_\lambda$  to the sum taken in step 1.

This procedure yields the approximation

$$w_n(x) = \underbrace{\sum_k \alpha_{0k} \Phi_{0k}(x)}_{\substack{m \text{ coeffs. } \neq 0 \\ (f, \Phi \text{ compact. supp.})}} + \underbrace{\sum_{j=0}^{\infty} \sum_{k \in \mathbf{Z}^d} \sum_{l=1}^{2^d-1} \beta_{jk}^{(l)} \Psi_{jk}^{(l)}(x)}_{\text{keep the largest } n-m \text{ coeffs.}} \quad (26)$$

and the following theorem provides corresponding approximation bounds.

**Theorem 2 ( DeVore, Jawerth, Popov, (DeVore et al., 1994) )** *Consider  $f \in \mathcal{B}_{pp}^s$ ,  $s, p > 0$  and  $s - d/p + d/u \geq 0$ . Let  $w_n$  denote the approximation (26) of  $f$ . If the scale function satisfies condition (20), then*

$$\|f - w_n\|_u \leq C(s, p) n^{-s/d} \|f\|_{\mathcal{B}_{pp}^s} \quad (27)$$

*holds. If, in addition,  $u$  satisfies  $s - d/p + d/u = 0$ ,  $u < \infty$ , and it is a priori known that  $f \in L_u$ , then the following converse bound holds.*

$$\|f\|_{\mathcal{B}_{pp}^s} \leq C(s, p, q) \left( 1 + n^{s/d} \|f - w_n\|_u \right) .$$

COMMENTS: This result is quite remarkable for the following reasons:

1. This approximation procedure gives the same rate of approximation for a wide variety of different Besov spaces (those satisfying  $s - d/p + d/u \geq 0$ ). Especially, it is not necessary to *a priori* know the extent of localized singularities, *i.e.*, index  $p$ .
2. When certain norms are used to measure the approximation error, and for functions with localized singularities, the approximation error tends to zero much faster if (26) is used than if linear approximation is performed. This follows by comparing (27) with the rate  $n^{-s+d/p-d/u}$  which is generic for linear approximations for the cases where  $0 < p \leq u$  and  $p \leq u$ .

Also, in the wavelet decomposition of a function with sparse singularities (e.g.,  $f \in \mathcal{B}_{pq}^s$ ,  $p < 2$ ), only a small number of basis functions are important, and the other ones can be neglected. In contrast to spaces of uniformly smooth functions, *it is not necessarily the first basis functions that should be used*. Let us go once more back to our example  $f(x) = 1_{\{0 \leq x < a\}}$ . Consider the wavelet decomposition of this function using a compactly supported wavelet  $\psi(x)$  such that  $\int \psi(x) dx = 0$ . It is evident that the coefficient  $\beta_{jk}$  vanishes for any wavelet  $\psi_{jk}(x)$  which does not cross the (local) singularities of  $f$ . Thus if we consider the projection of  $f$  on the subspace  $V_j$ , only  $O(j)$  coefficients of the decomposition significantly differ from zero (among  $2^j$  candidates).

### 3 Approximation in high-dimensional spaces

#### 3.1 The curse of dimensionality

The accuracy of an approximation depends on how densely observation points fill the input space. Thus having enough data points for good estimation would require the sample size to grow exponentially with the input dimension. This is referred to as the *curse of dimensionality*, a phrase coined by Bellman (Bellman, 1966). The curse of dimensionality is exhibited explicitly in the results (16) and (27) where the rate of convergence has order  $n^{-s/d}$ . For large input dimension  $d$ , this is exceedingly slow.

#### 3.2 Function classes of lower effective dimension

There are basically two ways to deal with the curse of dimensionality. Either to accept that a huge amount of data is necessary or to restrict the function class further. The latter means that, instead of having the dimension visible in the convergence rate, it will be hidden behind the function class. Kolmogorov (Kolmogorov, 1957) proved that every continuous function on  $[0, 1]^d$  can be represented as the additive superposition of continuous one-dimensional functions. Lorentz (Lorentz, 1976) gave an explicit scheme: Every continuous function  $f$  on  $[0, 1]^d$  can be written as

$$f(x_1, \dots, x_d) = \sum_{j=1}^{2d+1} g_j \left( \sum_{k=1}^d h_{jk}(x_k) \right) \quad (28)$$

for some continuous univariate functions  $(g_j)$ . Moreover, the functions  $(h_{jk})$  can be taken to be universal, *i.e.* they do not depend on  $f$ . Unfortunately, these results are not of great help for approximation, since the above functions  $g_j$  are usually extremely irregular even for a smooth  $f$  function<sup>4</sup>.

However, one way to bound the function class in its “effective dimension” is suggested by the generic decomposition (28). Introduce the variables  $z_j = \sum_{k=1}^d h_{jk}(x_k)$  which, since the  $h_{jk}$  are known, can be precomputed. The function  $f$  can then be written as

$$f(x_1, \dots, x_d) = \sum_{j=1}^{2d+1} g_j(z_j) \quad (29)$$

and the problem is to approximate  $2d + 1$  univariate functions  $g_j$ . Using  $m$  basis functions to approximate each  $g_j$  gives an approximation error of the order  $m^{-s}$  under usual smoothness assumptions on  $g$  (cf Section 2). Taking  $m = n/(2d + 1)$ , the total number of basis functions is  $n$  and the total approximation error will be of the order  $(2d + 1)(n/(2d + 1))^{-s}$  which, for large  $n$ , is of order  $n^{-s}$  and is much better than the above quoted  $n^{-s/d}$ .

Projection onto one-dimensional subspaces is the crux of the *Projection Pursuit Algorithm*, developed in (Friedman and Stuetzle, 1981) (a very good review of these results can be found in (Huber, 1985)), which consider estimates of  $f$  in the form

$$\hat{f}_N(x) = \sum_{j=1}^M \hat{g}_j(\alpha_j^T x),$$

where  $(\alpha_j)$  are unit vectors and each  $\alpha_j^T x$  may be thought of as a projection of  $x$ . The  $j$ -th term  $\hat{g}_j(\cdot)$  is constant along  $\alpha_j^T x = c$  and so is often called a *ridge function*: the estimate at a given

---

<sup>4</sup>This was already noted in the original paper by Kolmogorov.

point can be thought of as based on the averages over certain (in general adaptively chosen) strips  $\{x : |\alpha_j^T x - t_i| \leq \varepsilon\}$ . Other examples are *Recursive Partitioning* (Morgan and Sonquist, 1963), (Breiman et al., 1984) and related methods (c.f., for instance (Friedman, 1991) with discussion). These methods are derived from some mixture of statistic and heuristic arguments and give impressive results in simulations. Their drawback lies in the almost total absence of any theoretical results on their convergence rates. We refer the reader to the above references for additional information.

### 3.3 Neural networks

For a review of neural networks see the companion paper (Sjöberg et al., 1995). The following result was recently published in (Barron, 1993), it is the most accurate theoretical result on the neural network-based approximations today. Let  $\sigma(x)$  be a sigmoidal function (i.e. a bounded measurable function on the real line for which  $\sigma(x) \rightarrow 1$  as  $x \rightarrow \infty$  and  $\sigma(x) \rightarrow 0$  as  $x \rightarrow -\infty$ ). Consider a compactly supported function  $f$  with  $\text{supp}(f) \subseteq [0, 1]^d$ , and assume that

$$C_f = \int_{\mathbf{R}^d} |\omega| |\hat{f}(\omega)| d\omega < \infty, \quad (30)$$

where  $\hat{f}(\omega)$  denotes the Fourier transform of  $f$ . The main result of (Barron, 1993) can be roughly stated as follows: there exists an approximation  $f_n$  of the compactly supported function  $f$ , of the form

$$f_n(x) = \sum_{i=1}^n c_i \sigma(a_i^T x + t_i) + c_0 \quad (31)$$

(note that  $f_n$  is *not* compactly supported), such that

$$\|(f_n - f) 1_{[0,1]^d}\|_2 \leq 2\sqrt{d} C_f n^{-1/2}. \quad (32)$$

This result provides an upper bound of the minimum distance (in  $L_2$ -norm) between any  $f$  satisfying condition (30) and the class of all neural networks of size not larger than  $n$ . In the same article, the upper bound (32) is compared with the best achievable convergence rate for any linear approximation in class (30). It is shown that a lower rate for linear projections is  $n^{-1/d}$ , compare this with the much better rate  $n^{-1/2}$  for neural networks, especially for large dimension  $d$ .

**COMMENT** It is not easy to relate the function class defined by the condition (30) and more usual smoothness classes. For instance, one can show that if  $f \in W^{(d+1)/2+\epsilon}(L)$  for some  $\epsilon > 0$  and  $L < \infty$ , then (30) holds. Note that the same rate was obtained for “linear” projection estimators in section 2.1. On the other hand, it can be shown that even in the Sobolev class  $W^{(d+1)/2}$  there are “bad” functions, which does not verify (30).

This gives us an idea that the neural approximator outperform usual linear estimators on some special rather restricted functional classes. *What are these classes?* In the article (Barron, 1993) some classes are listed. In particular, if  $f(x) = g(x^T a)$  for some  $a \in \mathbf{R}^d$ ,  $|a| = 1$ , i.e.  $f(x)$  is a ridge function, then the Fourier transform  $\hat{f}(\omega)$  is concentrated in the direction of  $a$  and conditions (30) implies that  $\int |\omega| \hat{f}(\omega) d\omega < \infty$ .

It is much easier, however, to answer the question what are the classes on which neural nets behaves badly. Consider a class of function which are spherically symmetric, i.e.  $f(x) = g(|x|)$

for some  $g : \mathbf{R} \rightarrow \mathbf{R}$ . Using spheric coordinates  $\rho, \nu$  (here  $\rho = |x|$  and  $\nu$  is a vector on the unit sphere and  $\rho$  is the radius), we get from (30)

$$\int_{\mathbf{R}^d} |\omega| |\hat{f}(\omega)| d\omega = \int_0^\infty \rho |\hat{g}(\rho)| d\rho |S_\rho| \quad (33)$$

where  $|S_\rho|$  is the surface of the  $d$ -dimensional sphere of radius  $\rho$ , which is

$$S_\rho = \frac{\rho^{d-1} \pi^{d/2} d}{\Gamma(d/2 + 1)}$$

, where  $\Gamma(\cdot)$  is the standard  $\Gamma$ -function. We conclude from (30) and (33) that  $g(\cdot)$  should verify

$$\int \omega^d |\hat{g}(\omega)| d\omega < \infty.$$

This is a hard assumption, and it implies that the  $d$ -th derivative of  $g$  is bounded. We know (cf. section 2.2) that for such functions the rate  $n^{-1}$  can be attained by spline or wavelet approximations (and many other classical methods). The rate  $n^{-1/2}$  which is stated in (32) for neuron approximations is really not good in this case. These two examples illustrates the the following simple idea: the neural nets are not always good approximants. They behave badly on certain functional classes and outperform local estimators in some particular situations. This duality between local and “semi-local” methods has been discussed and developed in (Donoho and Johnstone, 1989).

When coming back to Barron’s result, it should be noted that no result is available which takes advantage of possible improved smoothness of the unknown function  $f$ . An iterative algorithm for the construction of the approximation (31) is also proposed. The true problem of system identification, i.e., that of neural network training based on noisy input/output data, is not addressed in this paper.

### 3.4 Wavelets

Note that in the orthonormal wavelet expansion

$$f(x) = \sum_k \alpha_{0k} \Phi_{0k}(x) + \sum_{ljk} \beta_{jk}^{(l)} \Psi_{jk}^{(l)}(x) ,$$

the dilation and translation parameters  $-2^{-dj}$  and  $k$  do not depend on the function to expand, only the linear weights  $\alpha_{jk}$  and  $\beta_{jk}^{(l)}$  depend on  $f$ . Suppose that we are able to construct an “adaptive wavelet basis”, i.e., with dilations and translations depending on the function  $f$ . The wavelet expansion of  $f$  using these basis functions is expected to use less wavelets, and thus we expect it to be more convenient for estimation purposes. To obtain such a basis we can *discretize the continuous wavelet transform* (36) which is given by the following theorem.

**Theorem 3** *Let  $\psi$  and  $\varphi$  be radial <sup>5</sup> functions satisfying*

$$\forall \omega \in \mathbf{R}^d : \int_0^\infty a^{-1} \hat{\varphi}(a\omega) \hat{\psi}(a\omega) da = 1 \quad (34)$$

---

<sup>5</sup>A function  $\varphi$  is radial if  $\varphi(x)$  depends only on  $|x|$ ; this implies that  $\hat{\varphi}(\omega)$  also depends only on  $|\omega|$ .



where we recall that  $\widehat{\varphi}(\omega)$  denotes the Fourier transform of function  $\varphi(x)$ . Then for any function  $f \in L_2(\mathbf{R}^d)$ , the following formulae define an isometry between  $L_2(\mathbf{R}^d)$  and a subspace of  $L_2(\mathbf{R}^d \times \mathbf{R}_+)$  (Daubechies, 1992) :

$$u(a, t) = a^{d-1/2} \int f(x) \varphi(a(x-t)) dx \quad (35)$$

$$f(x) = \int u(a, t) \psi(a(x-t)) a^{d-1/2} da dt . \quad (36)$$

Here  $a \in \mathbf{R}^+$  and  $t \in \mathbf{R}^d$  are respectively the dilation and translation factors.

We present the following algorithm proposed in (Delyon et al., 1994). Consider the continuous wavelet transform (36), which we rewrite as

$$\begin{aligned} f(x) &= \int u(a, t) \psi(a(x-t)) a^{d-1/2} da dt \\ &= \int \psi(a(x-t)) \text{sign}(u(a, t)) a^{(d-1)/2} |u(a, t)| da dt \\ &= \frac{1}{C} \int \psi(a(x-t)) \text{sign}(u(a, t)) w(a, t) da dt \end{aligned}$$

where we have renormalized  $u(a, t)$  by a constant factor  $C$  so that the function  $w(a, t) = C a^{(d-1)/2} |u(a, t)|$  can be considered as a probability density. Then we draw  $n$  independent random samples  $(a_i, t_i)_{i=1, \dots, n}$  from a distribution with density  $w(a, t)$ . Then we build

$$f_n(x) = \frac{1}{n} \sum_{i=1}^n a_i^{d/2} \psi(a_i(x-t_i)) \text{sign}(u(a_i, t_i)) , \quad (37)$$

which, thanks to the law of large numbers, converges to the true wavelet transform. Some faster implementations of this algorithm are given in (Delyon et al., 1994). Improving this estimate by some “bootstrapping” like technique, yields the following approximation result.

**Theorem 4 ((Delyon et al., 1994))**  *$\psi$  is any radial wavelet function such that there exists a related radial function  $\varphi$  which satisfies condition (34). Let  $p, \mu, l, s$  be real numbers satisfying*

$$1 < p < \left(1 - \frac{s-l}{d}\right)^{-1}, \quad \mu = \min\left(1 - \frac{1}{p}, \frac{1}{2}\right)$$

*and  $f$  be a function of the Sobolev space  $W_1^s(\mathbf{R}^d)$ ; then, for any  $n > 0$  there exists a function  $f_n$  of the form*

$$f_n(x) = \sum_{i=1}^n u_i \psi(a_i(x-t_i)) \quad (38)$$

*such that*

$$\|f_n - f\|_{W_p^l} \leq C n^{-\mu} \|f\|_{W_1^s} .$$

*In particular, if  $s > d/2$  then*

$$\|f_n - f\|_2 \leq n^{-1/2} C \|f\|_1^s . \quad (39)$$

COMMENT : Theorem 4 provides us with an upper bound for the rate of approximation when adaptive dilation/translation sampling is used to discretize the continuous wavelet transform. We should compare this rate with rates of convergence for approximations based on *fixed* dilation/translation sampling. For those approximations the best rate which can be attained for  $p = 2$  and  $s - d/2 = \varepsilon > 0$ , would be  $n^{-\varepsilon/d}$ , which is much slower than the rate (39).

Note also that in this case the rate of convergence for the “shrunk” wavelet approximants is also  $n^{-1/2-\varepsilon}$ . Thus these two types of nonlinear approximations have almost the same rate of convergence.

### 3.5 Breiman’s Hinging Hyperplanes

We now briefly discuss a recent technique due to Leo Breiman (Breiman, 1993), which practically combines some advantages of neural networks (in particular the ability to handle very large dimensional inputs) and of constructive wavelet based estimators (availability of fast training algorithms). Breiman’s technique is an elegant way of identifying piecewise linear models based on data collected from an unknown nonlinear system, see (Sontag, 1981) for the use of such models in control. Following (Breiman, 1993), we call *hinge function* a function  $y = h(x)$ ,  $x \in \mathbf{R}^d$  which consists of two hyperplanes continuously joined together, i.e., an open book. If the two hyperplanes are given as

$$y = \langle \beta^+, x \rangle + \beta_0^+, \quad y = \langle \beta^-, x \rangle + \beta_0^-,$$

where  $\langle \cdot, \cdot \rangle$  denotes scalar product in Euclidean spaces, then an explicit form for the hinge function is either

$$\begin{aligned} h(x) &= \max( \langle \beta^+, x \rangle + \beta_0^+, \langle \beta^-, x \rangle + \beta_0^- ) , \\ \text{or} \quad h(x) &= \min( \langle \beta^+, x \rangle + \beta_0^+, \langle \beta^-, x \rangle + \beta_0^- ) . \end{aligned}$$

It is proved in (Breiman, 1993), using the methods by Barron (Barron, 1993) that there is a constant  $C$  such that for any  $n$  there are hinge functions  $h_1, \dots, h_n$  such that

$$\|f - \sum_{i=1}^n h_i 1_{[0,1]^d}\|_2 \leq C n^{-1/2} \quad (40)$$

for any  $f$  such that

$$\int_{\mathbf{R}^d} |\omega|^2 |\hat{f}(\omega)| d\omega < \infty , \quad (41)$$

i.e., Breiman’s hinge model is as efficient as neural networks for the  $L_2$ -norm. Notice that, as in the neural network case, the condition (41) limits the function class by reducing the effective dimension, *c.f.* the comment after (32). However, as indicated at the beginning of this section, no convergence rate is given for models identified from noisy data (the bound (40) is not a convergence rate for identification, but only a rate of approximation of a given function by some finitely parameterized class of approximants). Further details are given in the companion paper (Sjöberg et al., 1995) concerning effective procedures for hinging hyperplanes estimation, we shall not discuss them any further.

## 4 Performance measures for non-parametric estimators

With the approximation results from the previous sections at hand, we are now ready to move on to analyze the behavior of estimation algorithms from noisy data. The performance analysis of non-parametric estimation algorithms and/or identification procedures is much more difficult than for parametric estimation. The following specific issues are important :

1. What plays the role of Cramer-Rao bound and Fisher Information Matrix in our case ? Recall that the Cramer-Rao bound reveals the best performance one can expect in identifying the unknown parameter  $\theta$  from sample data arising from some parameterized distribution  $\mathbf{p}_\theta, \theta \in \Theta$ , where  $\Theta$  is the domain over which the unknown parameter  $\theta$  ranges. In the non-parametric case, lower bounds for the best achievable performance are provided by *minimax risk functions*. We shall introduce these lower bounds and discuss associated notions of optimality.
2. For lower bounds, what is the class of systems on which best achievable performance is considered, is another important issue. For non-parametric representations of linear systems,  $L_2, L_\infty, H_2, H_\infty$ , with their associated norm are typical spaces to work with. For (even static) nonlinear systems, however, the choice is much wider. How wide should be the class  $\mathcal{F}$  of systems in consideration, what kind of smoothness should be required ? Are we interested in the behavior of the estimate at one particular point  $x$  of interest, or are we interested in the global behavior of the estimate ? Different distance measures should be used in these two different cases.

### 4.1 Lower bounds for best achievable performance.

In order to compare different non-parametric estimators it is necessary to introduce suitable figures of merit. It seems first reasonable to build on the mean square deviation (or mean absolute deviation) of some semi-norm<sup>6</sup> of the error, we denote it by  $\|\hat{f}_N - f\|$ . The following semi-norms are commonly used in non-parametric regression :  $\|f\| = (\int f^p(x)dx)^{1/p}$ ,  $0 < p < \infty$  ( $L_p$ -norm),  $\|f\| = \sup_x |f(x)|$  (uniform norm,  $C$ - or  $L_\infty$ -norm),  $\|f\| = |f(x_0)|$  (absolute value at a fixed point  $x_0$ ). Then we consider the *risk function*

$$R_{a_N}(\hat{f}_N, f) = \mathbf{E} \left[ a_N^{-1} \|\hat{f}_N - f\| \right]^2, \quad (42)$$

where  $a_N$  is a normalizing positive sequence. Letting  $a_N$  decrease as fast as possible so that the risk still remains bounded yields a notion of a convergence rate. Let  $\mathcal{F}$  be a set of functions which contains the “true” regression function  $f$ , then the maximal risk  $r_{a_N}(\hat{f}_N)$  of estimator  $\hat{f}_N$  on  $\mathcal{F}$  is defined as follows :

$$r_{a_N}(\hat{f}_N) = \sup_{f \in \mathcal{F}} R_{a_N}(\hat{f}_N, f) .$$

If the maximal risk is used as a figure of merit, the optimal estimator  $\hat{f}_N^*$  is the one for which the maximal risk is minimized, i.e., such that<sup>7</sup>

$$r_{a_N}(\hat{f}_N^*) = \min_{\hat{f}_N} \sup_{f \in \mathcal{F}} R_{a_N}(\hat{f}_N, f) .$$

---

<sup>6</sup>a semi-norm is a norm, except it does not satisfy the condition :  $\|f\| = 0$  implies  $f = 0$ .

<sup>7</sup>to properly understand the statement to follow, the reader should pay attention to definition (2) of an estimator.

We call  $\hat{f}_N^*$  the *minimax estimator* and the value

$$\min_{\hat{f}_N} \sup_{f \in \mathcal{F}} R_{a_N}(\hat{f}_N, f)$$

the *minimax risk* on  $\mathcal{F}$ . Notice that this concept is consistent with the mini-max concept used in the definition of  $n$ -widths in approximation theory in Section 2.

The construction of minimax non-parametric regression estimators for different sets  $\mathcal{F}$  is a hard problem. Presently, it has only solved asymptotically (for large samples) for some special cases (see, for instance, (Efroimovich and Pinsker, 1982), (Efroimovich and Pinsker, 1983), (Efroimovich and Pinsker, 1984)). However, letting  $a_N$  decrease as fast as possible so that the minimax risk still remains bounded yields a notion of a best achievable convergence rate, similar to that of parametric estimation. More precisely, we state the following definition:

**Definition 1 (lower rate and minimax rate of convergence)**

1. The positive sequence  $a_N$  is a **lower rate of convergence for the set  $\mathcal{F}$  in the semi-norm  $\|\cdot\|$**  if

$$\liminf_{N \rightarrow \infty} r_{a_N}(\hat{f}_N^*) = \liminf_{N \rightarrow \infty} \inf_{\hat{f}_N} \sup_{f \in \mathcal{F}} \mathbf{E} \left[ a_N^{-1} \|\hat{f}_N - f\| \right] \geq C_0 \quad (43)$$

for some positive  $C_0$ .

This notion can be refined as follows.

2. The positive sequence  $a_N$  is called **minimax rate of convergence for the set  $\mathcal{F}$  in semi-norm  $\|\cdot\|$** , if it is a lower rate of convergence, and if, in addition, there exists an estimator  $\hat{f}_N^*$  achieving this rate, i.e., such that

$$\limsup_{N \rightarrow \infty} r_{a_N}(\hat{f}_N^*) < \infty .$$

The inequality (43) is a kind of negative statement that says that no estimator of function  $f$  can converge to  $f$  faster than  $a_N$ . Thus, a coarser, but easier approach consists in assessing the estimators by their convergence rates. In this setting, by definition, optimal estimators reach the lower bound as defined in (43) (recall that the minimax rate is not unique: it is defined to within a constant).

## 4.2 Some negative results.

From the discussion in Section 1, it should be evident that it is the assumed smoothness class that dictates the minimax rate of convergence. Generally it holds that the larger the class of functions, the slower the convergence rate. Devroye and Györfi (Devroye and Györfi, 1985) (Devroye, 1982), have proven the following result<sup>8</sup>. Consider the following classes of functions on  $\mathbf{R}$ :

$\mathcal{F}^*$  : the class of all functions  $f$  such that  $f(x) = 0$  for  $x > 1$  or  $x < 0$ , and  $|f(x)| \leq C$  for  $x \in [0, 1]$ .

---

<sup>8</sup>Note, however, that convergence can sometimes be proved without any smoothness assumption (Devroye and Wagner, 1980).

$\mathcal{F}_0^*$  : the class of all continuous functions  $f \in \mathcal{F}^*$ .

$\mathcal{F}_\infty^*$  : the class of all functions  $f \in \mathcal{F}^*$  having all continuous derivatives on  $[0, 1)$  (notice that the interval is right open).

Let  $\hat{f}_N$  be an arbitrary estimate of  $f$ . Then for the classes  $\mathcal{F}^*$ ,  $\mathcal{F}_0^*$  and  $\mathcal{F}_\infty^*$  defined above (we denote them generically by  $\mathcal{F}$ ),

$$\sup_{\mathcal{F}} \limsup_{N \rightarrow \infty} \mathbf{E} \left[ a_N^{-1} \int_0^1 |\hat{f}_N(x) - f(x)| dx \right] = \infty$$

for any positive sequence  $a_N \rightarrow 0$ .

Thus, *no convergence rate does exist for any of the above classes  $\mathcal{F}^*$ ,  $\mathcal{F}_0^*$  and  $\mathcal{F}_\infty^*$* . In other words, the convergence can be arbitrary slow, depending on the unknown function or density  $f$  to be estimated! It is a natural consequence of the fact that the above classes  $\mathcal{F}^*$ ,  $\mathcal{F}_0^*$  and  $\mathcal{F}_\infty^*$  are too rich: they contain functions which are extremely difficult to approximate. *In other words, in order to obtain any interesting rate of convergence, smoothness conditions should be imposed.*

### 4.3 Some positive results.

Let us now concentrate on the case of deterministic uniform design, i.e., the input data  $X$  are uniformly sampled in the considered interval. The following result in the case of regular design can be acknowledged to (Ibragimov and Khasminskij, 1981) (for the random design case, see (Stone, 1982), (Korostelev and Tsybakov, 1981)).

**Theorem 5** *Let us consider the Hölder class  $\mathcal{C}^s(L)$  on  $[0, 1]^d$ , closely related to the Sobolev classes. is the The Hölder class  $\mathcal{C}^s(L)$  is the the family of functions  $f(x)$ ,  $x \in [0, 1]^d$  defined by*

$$\mathcal{C}^s(L) = \left\{ f : |f^{(k)}(x) - f^{(k)}(x')| \leq L|x - x'|^{s-k}, \text{ for any } x, x' \in [0, 1]^d \right\}, \quad k = \lfloor s \rfloor. \quad (44)$$

Consider

$$\|g\| = \left( \int |g(x)|^p dx \right)^{1/p}, \quad 0 < p < \infty \quad \text{or} \quad \|g\| = |g(x_0)|.$$

Then  $N^{-\frac{s}{2s+d}}$  is a lower rate of convergence for the class  $\mathcal{C}^s(L)$  in the semi-norm  $\|\cdot\|$ . Furthermore,  $\frac{N}{\ln N}^{-\frac{s}{2s+d}}$  is a lower rate of convergence for the class  $\mathcal{C}^s(L)$  in the norm  $\|g\| = \sup_{x \in [0, 1]} |g(x)|$ .

Note that to obtain the correct rate of convergence for the distance at a fixed point  $x_0$ , the corresponding Lipschitz property is required at  $x_0$  only. Similar results hold when the class  $\mathcal{C}^s(L)$  is replaced by the class  $\mathcal{W}_p^s(L)$ ,  $p \geq 2$  (see (13)). Then  $N^{-\frac{s}{2s+d}}$  is also a lower rate of convergence for this class in the  $L_p$ -norm of the error.

---

<sup>9</sup> $\lfloor s \rfloor$  denotes the maximal integer  $k < s$ .

## 5 Estimation in classes of uniformly smooth functions

Throughout this section, Problem 1 is considered. The discussion in Section 1.2 and Section 2 give the required background for understanding how to perform estimation of the unknown  $f$  function using a fixed basis function expansion. Let us take the Sobolev class  $\mathcal{W}_2^s(L)$  and our model (8), (9). In this case we take as  $(g_j)$  the Fourier basis, and to obtain the optimal approximation with  $n$  basis functions, we can simply take first  $n$  terms of the expansion (9) or (11). This gives a certain maximal bias error, *cf.* the first term in (10). In this particular example, this error is of order  $n^{-s/d}$ . The parameters in the function expansion (9) are estimated via empirical means based on  $N$  noisy observations. The mean square error of the estimate of each coefficient is  $O(1/N)$ . Thus the total mean square error of the estimate will be, as usual, the sum of the stochastic part and of the bias due to the approximation error: this yields  $O(n/N) + O(n^{-2s/d})$ . The optimal choice for  $n$  balances these two terms:  $n = N^{\frac{1}{2s+d}}$ . This choice for  $n$  yields a quadratic error of order  $N^{-\frac{2s}{2s+d}}$ . This is the typical scheme that is followed even in the cases where the basis function expansion is not as explicit as in this example.

The estimators we consider in this section are *linear*, i.e. they can be written

$$\hat{f}_N(x) = \sum_{i=1}^N Y_i W_{N,i}(x), \quad W_{N,i}(x) = W_{N,i}(x, X_1, \dots, X_N) \quad (45)$$

where we recall that  $\mathcal{O}_1^N = \{(X_1, Y_1), \dots, (X_N, Y_N)\}$  is the given random sample observation. It is only the weights  $W_{N,i}(x)$  that may differ from estimator to estimator. This means that such an estimator satisfy  $\widehat{f+g} = \hat{f} + \hat{g}$ ; notice, however, that functions  $f, g$ , and their estimates, are generally nonlinear as functions of their input  $x$ . Linear estimators build the folklore of non-parametric estimation: kernel estimators, projections on linear subspaces of functions, are typical instances we shall describe. We shall then discuss, both practically and theoretically, some severe practical limitations of linear estimators. Roughly speaking, linear estimators are suitable for systems with “uniform smoothness”; systems with sparse singularities (e.g., hard limiters, quantizers, some mechanical systems) are poorly handled. This follows from the discussion at the end of Section 2.1.

### 5.1 Kernel estimators for regression functions and densities

Kernel estimators were first proposed by Nadaraya and Watson in 1964 (Nadaraya, 1964) and (Watson, 1969). The Nadaraya-Watson kernel estimator is an interpolation procedure. It is given by

$$\hat{f}_N(x) = \frac{\sum_{i=1}^N Y_i K\left(\frac{X_i - x}{h_N}\right)}{\sum_{i=1}^N K\left(\frac{X_i - x}{h_N}\right)}, \quad (46)$$

where  $(h_N)$  is a sequence of positive numbers,  $h_N \rightarrow 0$  as  $N \rightarrow \infty$ , and  $K$  is a function on  $\mathbf{R}$  satisfying

$$\lim_{|u| \rightarrow \infty} |u| |K(u)| = 0, \quad \int_{-\infty}^{\infty} |K(u)| du < \infty, \quad \sup_{u \in \mathbf{R}} |K(u)| < \infty, \quad \int_{-\infty}^{\infty} K(u) du = 1. \quad (47)$$

The positive number  $h_N$  is called the *bandwidth* and the function  $K$  satisfying (47) is called a *kernel*; in fact,  $h_N$  is better interpreted as a scaling factor. Clearly, the Nadaraya-Watson estimator is linear, and has the form (45). Typical examples of kernels are  $K(u) = (1/2) 1_{\{|u| \leq 1\}}$

(rectangular window kernel), and  $K(u) = (1/\sqrt{2\pi}) \exp(-|u|^2/2)$  (Gaussian kernel), etc... Usually  $K$  is chosen to be an even function.

The idea of kernel estimation is simple, let us explain it for the case of the rectangular kernel in one dimension. In this case the estimator (46) is a simple moving average with equal weights: the estimate at point  $x$  is the average of observations  $Y_i$  corresponding to  $X_i$ 's belonging to the "window"  $[x - h_N, x + h_N]$ . If  $h_N \rightarrow \infty$  then the estimator tends to  $N^{-1} \sum_i Y_i$ , the average of all observations, and thus for functions  $f$  which are far from being constant, the bias becomes large. If  $h_N$  is very small (say, smaller than the pairwise distance between sample points  $X_i$ ) then the estimator reproduces the observations:  $\hat{f}_N = Y_i$ . In this extremal case the variance of the error becomes high. Thus increasing  $h_N$  tends to increase the bias of estimator, while reducing  $h_N$  leads to a larger variance. The optimal choice for  $h_N$  corresponds to an equal balance between bias and variance.

Also closely related to estimator (46) is the Parzen-Rosenblatt kernel estimator for densities. Let  $X_1, \dots, X_N$  be independent and identically distributed random variables with common density  $f(x), x \in \mathbf{R}^d$ . The Parzen-Rosenblatt estimator of density  $f(x)$  is a suitably smoothed histogram. It is defined as (Parzen, 1962), (Rosenblatt, 1956)

$$\hat{f}_N(x) = \frac{1}{Nh_N^d} \sum_{i=1}^N K\left(\frac{X_i - x}{h_N}\right), \quad (48)$$

where  $d$  is the state-space dimension of  $X$  and  $K$  is a kernel as in (47). Kernel estimate (46) can be easily derived from the Parzen-Rosenblatt one. Recall definition (1) of the regression function, take the Parzen-Rosenblatt estimator (48) for the joint density  $f(x, y)$  of  $(X, Y)$  and denote it by  $\hat{f}_N(x, y)$ . Then, replacing, in formula

$$f(x) = \frac{\int y f(x, y) dy}{\int f(x, y) dy},$$

$f(x)$  and  $f(x, y)$  by their corresponding Parzen-Rosenblatt estimates, yields Kernel estimate (46).

We now state a sample of results about the properties of kernel estimates for the  $d$ -dimensional case when it is known a priori that  $f$  belongs to the Hölder class  $\mathcal{C}^s(L)$  (see (44) for the definition of  $\mathcal{C}^s(L)$ ).

We acknowledge Rosenblatt (Rosenblatt, 1971) for the first two statements of it, though it probably belongs to the earlier folklore of non-parametric statistics.

**Theorem 6 ((Rosenblatt, 1971))** *Let  $\hat{f}_N$  be a kernel estimate with bandwidth  $h_N$  such that  $h_N \rightarrow 0$  and  $Nh_N \rightarrow \infty$ , with kernel  $K$  satisfying  $\int x^j K(x) dx = 0$  for  $j = 1, \dots, k$ . Here,  $x^j$  denotes any product of the form  $x_1^{j_1} x_2^{j_2} \dots x_d^{j_d}$  where  $j_1 + j_2 + \dots + j_d = j$ , and  $x_1, \dots, x_d$  are the coordinates of  $x$ . Assume that the observations  $X_i$  are independent and identically distributed on  $[0, 1]^d$  with density  $g(x) \geq c > 0$ ,  $g \in \mathcal{C}^s(L)$ , and that the noise satisfies  $\mathbf{E}e_i = 0$  and  $\mathbf{E}e_i^2 \leq \sigma_e^2 < \infty$ . Then*

1. *Uniformly over  $f \in \mathcal{C}^s(L)$  and  $x \in [0, 1]^d$ , we have the pointwise bound*

$$\mathbf{E}|\hat{f}_N(x) - f(x)|^2 \leq C \left( L^2 h_N^{2s} + \frac{\sigma_e^2}{Nh_N^d} \right). \quad (49)$$

The optimal value of  $h_N$  which minimizes the right-hand side of (49) is given by

$$h_N = \left( \frac{\sigma_e^2}{L^2 N} \right)^{1/(2s+d)} . \quad (50)$$

For this value of  $h_N$

$$\mathbf{E}|\hat{f}_N(x) - f(x)|^2 \leq CL^{2/(d+2s)} \left( \frac{\sigma_e^2}{N} \right)^{2s/(2s+d)} .$$

2. If we consider instead the global error measure  $\mathbf{E}\|\hat{f}_N - f\|_2^2$ , using again the same optimal value (50) for  $h_N$  yields the same bound, uniformly over  $f \in \mathcal{C}^s(L)$ .

COMMENTS :

1. As expected from the above informal discussion concerning the rectangular kernel, the bound for the estimation error variance given on the right hand side of (49) is decomposed into *bias* + *variance* terms. And, as expected, the optimal choice of  $h_N$  in (50) exactly balances these two terms.
2. Note that we have both pointwise and global bounds, which reflects the local nature of kernel estimates.
3. The properties of the Parzen-Rosenblatt algorithm of density estimation are identical when the unknown density  $f$  satisfies  $f \in \mathcal{C}^s(L)$ . Note that, since  $\text{supp } f \subseteq [0, 1]^d$ , the  $L_1$ -norm of the error (restricted to the  $[0, 1]^d$ ) is dominated by the  $L_2$ -norm. So we get from the second statement of the theorem

$$\mathbf{E}\|\hat{f}_N - f\|_1^2 \leq CL^{2/(d+2s)} \left( \frac{\sigma_e^2}{N} \right)^{2s/(2s+1)}$$

provided  $h_N$  is chosen as in (50).

4. Often the following recursive version of the kernel estimator is considered, (Duflo, 1993), (Oppenheim and Portier, 1990) :

$$\begin{aligned} \hat{f}_n(x) &= \Gamma_n^{-1}(x) \left( \sum_{i=0}^n Y_i h_i^{-d} K \left( \frac{X_i - x}{h_i} \right) \right) \quad \text{if } \Gamma_n(x) \neq 0 \text{ and } \hat{f}_n(x) = 0 \text{ if } \Gamma_n(x) = 0 , \\ \Gamma_n(x) &= \sum_{i=0}^n h_i^{-d} K \left( \frac{X_i - x}{h_i} \right) , \end{aligned}$$

or

$$\begin{aligned} \hat{f}_n(x) &= \hat{f}_{n-1}(x) + \Gamma_n^{-1}(x) \left( Y_n - h_n^{-d} K \left( \frac{X_n - x}{h_n} \right) \hat{f}_{n-1} \right) , \\ \Gamma_n(x) &= \Gamma_{n-1}(x) + h_n^{-d} K \left( \frac{X_n - x}{h_n} \right) . \end{aligned} \quad (51)$$

In this form the algorithm resembles very much the recursive Least Squares algorithm for estimating the parameters of linear models. When the bandwidth is such that  $h_i = h i^{-\alpha}$



for some  $0 < \alpha < 1$ , the properties of the algorithm (51) in the static regression problem are essentially the same as those of the “off-line version” (46). In (Oppenheim and Portier, 1990), (Portier, 1992) and (Duflo, 1993) this algorithm was used to identify stable non-parametric autoregression models of the form (4), and the convergence of this estimator was proved. Furthermore, the same algorithm was used to provide the estimates of  $\hat{f}_n$  in the closed loop system (5)–(6), and the stability of such an adaptive control scheme was proved — (Oppenheim and Portier, 1990) and (Portier, 1992) consider essentially one-dimensional case, and in (Duflo, 1993) the general multi-dimensional case is studied.

## 5.2 Piecewise-polynomial estimators

Another non-parametric regression estimator which is commonly used is the piecewise-polynomial one. The idea is the same as for the kernel estimator, though the averaging is made over *bins* (i.e., small cubes) of fixed size  $\delta_N$  rather than in  $h_N$ -neighborhood of the current point  $x$ . It is also closely related to the *radial-basis function (RBF) networks* with rigid location for the radial functions, see (Poggio and Girosi, 1990), (Wahba, 1991). The simplest example of this method is the piecewise-constant estimator or *regressogram*. The value of the estimate  $\hat{f}_N$  in each bin equals the average of observations  $Y_i$  such that corresponding  $X_i$  belong to the bin. For the sake of clarity we consider the one-dimensional case.

The piecewise-polynomial estimator can be formally defined in terms of the following optimization problem. Let  $\delta_N \rightarrow 0$  be a positive sequence, and we assume that  $\delta_N^{-1} = M$  is an integer. Define  $u_l = l\delta_N$ ,  $l = 0, \dots, M$ , and divide the interval  $[0, 1]$  into  $M$  cubes (bins) of the form  $U_1 = [0, u_1)$ ,  $U_2 = [u_1, u_2)$ ,  $\dots$ ,  $U_M = [u_{M-1}, 1]$ , so each bin has length  $\delta_N$ . Set  $F(x) = (1, x, \frac{x^2}{2}, \dots, \frac{x^k}{k!})^T$  and, for each bin  $U_l$ ,  $l = 1, \dots, M$ , solve for  $\theta \in \mathbf{R}^{k+1}$  in the least squares sense the system of equations

$$Y_i = \theta^T F\left(\frac{X_i - u_{l-1}}{\delta_N}\right), \quad X_i \in U_l \quad (52)$$

and denote by  $\hat{\theta}_{N,l}$  the corresponding solution. Then the piecewise-polynomial estimate  $\hat{f}_N$  of order  $k$  in each bin  $U_l$  is expressed as

$$\hat{f}_N(x) = \hat{\theta}_{N,l}^T F\left(\frac{x - u_{l-1}}{\delta_N}\right), \quad x \in U_l \quad (53)$$

The value  $\delta_N$  is called the *binwidth*. As for the bandwidth  $h_N$  of kernel estimate, the binwidth tunes the smoothness: larger  $\delta_N$  leads to a higher bias, and smaller  $\delta_N$  results in a higher variance. In order for the least-squares problem in (53) to be non-degenerate we require that the number of points  $X_i$  in each bin is larger than  $k + 1$ .

Stone, (Stone, 1982) has proved a result similar to Theorem 6 for this type of estimate (see (44) for the definition of the Hölder space  $C^s(L)$ ). We state this result in the general  $d$ -dimensional case. Assume that the observations  $X_i$  satisfy the assumptions of Theorem 6. Let  $\hat{f}_N$  be a piecewise polynomial estimate of order  $k = \lfloor s \rfloor$ , with binwidth  $\delta_N \rightarrow 0$  and  $N\delta_N \rightarrow \infty$  as  $N \rightarrow \infty$ . Then *statement 1 of Theorem 6 holds with binwidth  $\delta_N$  substituted for the bandwidth  $h_N$ .*

### COMMENTS :

1. Note that, unlike for Kernel estimates, piecewise polynomial estimates compute projections on the fixed set of functions  $F\left(\frac{x - u_{l-1}}{\delta_N}\right)$ ,  $x \in U_l$  (the  $l$ -th bin). The same remark holds for the projection estimate to follow.

2. As can be seen, piecewise polynomial and kernel estimates have the same asymptotic accuracy when  $N \rightarrow \infty$ .
3. If  $f$  is a smooth function (i.e.,  $s \geq 1$ ), the optimal number of bins is  $n_\delta \sim \delta_M^{-1}$  and is much less than the number of observations ( $n_\delta \sim N^{1/3}$  for  $s = 1$ ). This number is equivalent to the memory size required to implement the algorithm: to reconstruct the estimate,  $k = \lfloor s \rfloor$  coefficients are necessary. Thus, if  $N$  is large, this algorithm offers significant advantage, in terms of memory requirements, over kernel estimates in which all measurements should be kept to reconstruct  $f(x)$ . Also, computing (52)–(53) is of lower computational burden than computing (46). These two points make the piecewise polynomial estimate more attractive.
4. Unfortunately there is no reasonable recursive version of the estimate  $\hat{f}_n$ . Although one can use the recursive least squares algorithm to compute linear regression coefficients  $\hat{\theta}_{N,l}$  in (53), the derivations quickly become messy, because the number  $M$  of bins depends on  $N$ , and so does the number of equations in the algorithm.

### 5.3 Projection estimates

Another class of function estimates was introduced by Cencov (Cencov, 1982), who called them *projection estimates*. The idea consists in expanding the unknown function into its “empirical” Fourier series. Consider the Sobolev class  $\mathcal{W}_2^s(L)$  of functions  $f(x)$  (11) but now defined for  $x \in [0, 1]^d$ . In this case (11) becomes

$$f(x) = \sum_{|j|=1}^{\infty} c_j \varphi_j(x), \quad (54)$$

where  $j = (j_1, \dots, j_d)$  is a multi-index,  $x = (x^1, \dots, x^d)^T$ ,  $\varphi_j(x) = \varphi_{j_1}(x^1) \times \dots \times \varphi_{j_d}(x^d)$ ,  $\varphi_1 \equiv 1$ ,  $\varphi_{2k}(x) = \sqrt{2} \sin(2\pi kx)$  and  $\varphi_{2k+1}(x) = \sqrt{2} \cos(2\pi kx)$ ,  $k = 1, \dots$ . The condition (12) remains the same

$$\sum_{j=1}^{\infty} |c_j|^2 (1 + |j|^{2s}) < L^2. \quad (55)$$

We assume that the *input  $X$  is uniformly distributed*. We construct the estimate  $\hat{f}_N$  as follows:

$$\hat{f}_N(x) = \sum_{j=1}^m \hat{c}_j^N \varphi_j(x), \quad (56)$$

where  $m$  is the “model order”, and the empirical estimates  $\hat{c}_j^N$  of Fourier coefficients

$$\hat{c}_j^N = \frac{1}{N} \sum_{i=1}^N Y_i \varphi_j(X_i) \quad (57)$$

are substituted for the true ones  $c_j$ ,  $j = 1, \dots, m$ . Note that the assumption that  $X$  is uniformly distributed has been used. Note also that the estimate (56)–(57) is linear (cf. (45)) with weights given by

$$W_{N,i}(x) = \sum_{j=1}^m \frac{1}{N} \varphi_j(x) \varphi_j(X_i).$$

Cencov, (Cencov, 1982) has proved the following counterpart of statement 1 of Theorem 6 : Let  $\hat{f}_N$  be a projection estimate. Then, uniformly over  $f \in \mathcal{W}_2^s(L)$  and  $x \in [0, 1]^d$ ,

$$\mathbf{E} \|\hat{f}_N(x) - f(x)\|_2^2 \leq C \left( L^2 m^{-2s} + \frac{\sigma_e^2 m^d}{N} \right). \quad (58)$$

The optimal order  $m$  of the model is

$$m = \lfloor \left( \frac{L^2 N}{\sigma_e^2} \right)^{1/(2s+d)} \rfloor, \quad (59)$$

it balances bias and variance error estimates, and yields the bound

$$\mathbf{E} \|\hat{f}_N(x) - f(x)\|_2^2 \leq C L^{2/(d+2s)} \left( \frac{\sigma_e^2}{N} \right)^{2s/(2s+d)}. \quad (60)$$

The following result, due to Ibragimov and Khas'minskij (Ibragimov and Khasminskij, 1981), provides a global uniform bound. Take

$$m = \lfloor \left( \frac{N}{\ln N} \right)^{1/(2s+d)} \rfloor$$

for the model order (note that this is slightly different from (59)). Then, uniformly over  $f \in \mathcal{C}^s(L)$  (the class  $\mathcal{C}^s(L)$  is defined in (44)), it holds that

$$\mathbf{E} \|\hat{f}_N - f\|_\infty^2 \leq O \left( \frac{\ln N}{N} \right)^{2s/(2s+d)}. \quad (61)$$

COMMENTS :

1. Projection estimates have the same rate of convergence (up to a constant) as kernel or piecewise polynomial ones.
2. The bound (58) for the quadratic error of the algorithms appears rather naturally if we consider the following argument: when we approximate  $f \in \mathcal{W}_2^s$  using  $m$  terms of its Fourier decomposition, the approximation error is  $O(m^{-2s/d})$ . Furthermore, the stochastic error in each term is of order  $O(N^{-1})$ . This simple calculus can be repeated for any non-parametric estimate. Obviously, it is beyond our possibilities to reduce the stochastic component of the error. On the other hand, the bias part depends on the method we choose to approximate the function (piecewise polynomial, trigonometric series, etc.) with, and this choice of approximant is of primary importance.
3. From the computational point of view, projection estimates are more attractive than piecewise polynomial estimates, since it uses an orthonormal basis of functions (the Fourier basis), which dramatically simplifies the computation of the least-squares estimates  $\hat{c}_j$  of the Fourier coefficients  $c_j$ , cf. (57) with (52).

## 5.4 Selecting model complexity

As we have seen, the convergence of the estimates strongly depends on the choice of the bandwidth  $h_N$  for kernel estimator, the model order  $m$  for the projection estimator, and the binwidth  $\delta_N$  (or, equivalently, the “model order”  $M = \delta^{-1}$ ) for piecewise polynomial estimator. *These design parameters depend on the parameters of the smoothness class  $\mathcal{C}^s(L)$  or  $\mathcal{W}_2^s(L)$ , which are a priori unknown* — see definition (44) of this class and the use of parameters  $(s, L)$  in Theorem 6 and corresponding results for others estimators. Even if some information about the smoothness parameter  $s$  is available, the knowledge of the value  $L$  is of importance when the data sample is of bounded length. Let us illustrate this with the following example, where the input  $x$  is scalar. Consider the problem of estimating a function  $f(x)$  in additive white noise  $e$ , with  $\sigma_e^2 = 1$ . Assume that  $f$  has support  $[0, 1]$ , that all its derivatives are continuous, and that  $f(1/2) = 1$ ,  $f(0) = f(1) = 0$ . Note that in this case, typically,  $\sup_x |f^{(s)}(x)| \approx s^s$ , i.e., higher order derivatives become very large in uniform bound. In this case the bounds in Theorem 6 are of order  $a_N(s) = (s/N)^{2s/(2s+1)}$  when the parameter is selected for the smoothness  $s$ . Assume that the size of the observation sample is  $N = 10000$ , then  $a_N(2) = 0.0110$ ,  $a_N(3) = 0.0095$ , but we have already  $a_N(4) = 0.0122$  (the value of  $s$  which minimizes  $a_N$  is  $s \approx 3.4814$  with  $a_N(s) \approx 0.00946$ ). This illustrates the fact that the tightest bound is not obtained by taking the largest possible  $s$ , but rather by selecting the most favorable pair  $(s, L)$ , which obviously is much more difficult.

Given that we only have in practice samples of finite size  $N$ , we shall not try to estimate the most favorable pair  $(s, L)$ , but we shall proceed differently. The model order (or bandwidth, or binwidth, depending on the different estimates) shall be estimated from data using a procedure usually referred to as the *generalized cross validation* (GCV) test. GCV procedures were studied for kernel (see, for instance, (Rice, 1984), (Härdle and Marron, 1985)), spline ((Li, 1986), (Craven and Wahba, 1979)), and projection estimates (c.f. (Polyak and Tsybakov, 1990), (Li, 1987)). Let us consider, for instance, the procedure for the projection estimates<sup>10</sup>. To make the model order explicit in formula (56) we shall write  $\hat{f}_{m,N}$  instead of  $\hat{f}_N$ . Set  $S_{m,N}^2 = N^{-1} \sum_{i=1}^N \|Y_i - \hat{f}_{m,N}(X_i)\|^2$ . As for the prediction error variance estimate in parametric prediction error methods,  $S_{m,N}^2$  is a *biased* estimate of the error. Thus one cannot minimize  $S_{m,N}^2$  with respect to  $m$  directly: the result of such a brute force procedure would give a function  $\hat{f}_{m_N,N}(x)$  which perfectly fits the noisy data, this is known as “overfitting” in the neural network literature. The solution rather consists in introducing a penalty which is proportional to the model order  $m$ , i.e., we search for  $m_N$  such that

$$m_N = \arg \min_{m \leq N} \left( S_{m,N}^2 + \frac{2\sigma_e^2 m}{N} \right). \quad (62)$$

This technique is clearly equivalent to the celebrated Mallows-Akaike criterion (Mallows, 1973), (Akaike, 1970). The following result, due to Polyak and Tsybakov (Polyak and Tsybakov, 1990), shows the consistency of this procedure. Assume that the Fourier coefficients of  $f$  in expansion (11) satisfy  $|c_j| \leq \varepsilon_j$ ,  $\sum_{j=1}^{\infty} \varepsilon_j < \infty$ ,  $(j \varepsilon_j)$  is non-increasing, and  $\sigma_e^2$  is known. Set  $V_{m,N} = \|\hat{f}_{m,N} - f\|_2^2$ . Then for the estimate (56), (57), and (62), it holds that

$$\frac{V_{m_N,N}}{\min_m V_{m,N}} \rightarrow 1 \text{ a.e. as } N \rightarrow \infty.$$

Another “classical” adaptation approach is closely related to the problem of filtering of a Gaussian stationary process. The technique developed in (Efroimovich and Pinsker, 1982),

<sup>10</sup>In fact, a similar result holds for the spline or piecewise polynomial ones.

(Efroimovich and Pinsker, 1984) for the projection estimates is often referred to as Efroimovich-Pinsker filter. To understand the idea of the method let us consider the estimates  $\hat{c}_j$ ,  $j = 1, \dots, N$  of the Fourier coefficients  $c_j$  of  $f$  as in Section 5.3. If the observation noise ( $e_i$ ) is independent Gaussian, then the errors  $\hat{c}_j^N - c_j$  are Gaussian and un-correlated (and thus independent), and the Wiener filter can be applied to the sequence  $\hat{c}_j$  to obtain the estimates  $\tilde{c}_j$  of  $c_j$ :

$$\tilde{c}_j = \frac{c_j^2}{c_j^2 + \sigma_e^2/N} \hat{c}_j. \quad (63)$$

It can be shown that this choice of Fourier coefficients yields the least possible asymptotic  $L_2$  error among all “projection” estimates. Naturally, the exact values of the coefficients  $c_j$  are not available. To construct an adaptive filter Efroimovich and Pinsker proposed to use instead of the filter coefficients their estimates, which are obtained in the following way: we divide the set of indices  $j = 1, \dots, N$  into  $m$  groups:  $T_1 = \{j = 1, j = 2\}$ , ...,  $T_{k+1} = \{j = 2^k + 1, \dots, j = 2^{k+1}\}$ , ...,  $T_m = \{j = 2^{m-1} + 1, \dots, j = 2^m\}$  (we have supposed for the sake of clarity that  $N = 2^m$ ). We set

$$\Lambda_k = \frac{\Theta_k}{\Theta_k + \sigma_e^2/N},$$

where

$$\Theta_k = |T_k|^{-1} \sum_{j \in T_k} (\hat{c}_j^2 - \sigma_e^2/N)$$

(here  $|T_k|$  is a cardinality of  $T_k$ ). We put, finally,

$$\tilde{c}_j = \Lambda_k \hat{c}_j$$

for all  $j \in T_k$ . It was shown (cf. (Efroimovich and Pinsker, 1984)) that this adaptive filter is asymptotically equivalent to the optimal one (63).

**REMARK** Note that the adaptive algorithms, though we have developed them starting with linear methods, are not linear anymore. Quite naturally, when trying to infer some additional information from the data we loose the linearity of the estimates.

## 6 Estimation in classes of locally spiky and jumpy functions

### 6.1 Spatial adaptivity

The estimation of functions with sparse singularities should naturally be based on function approximation in the corresponding smoothness classes. This was discussed in Section 2.2 — Section 3.5. During the last fifteen years this topic has been a very active field in the statistical community and has been characterized by successful practical applications but, oddly enough, an almost complete lack of theoretical results. Spatially adaptive methods include all sorts of neural networks, projection pursuit (Friedman and Stuetzle, 1981), classification and regression trees (CART) (Breiman et al., 1984), Multivariate adaptive regression splines (MARS) (Friedman, 1991), Variable Bandwidth Kernel methods (Müller and Stadtmüller, 1987), and others.

These methods implicitly or explicitly attempt to adapt the fitting method to the form of the function being estimated, by ideas like recursive dyadic partitioning of the space on which the function is defined (CART and MARS) and adaptively estimating a local bandwidth function (Variable Kernel Methods).

We discuss now some issues related to the problem of spatial adaptivity. In fact, so-called “spatially adaptive” methods address two different problems:

## Estimating functions which have sparse singularities and otherwise are smooth.<sup>11</sup>

An interesting approach consists in *finding a parameterized family of functional classes which*

1. *fits our prior knowledge about the smoothness of the function to be estimated, (in particular, that  $f$  is smooth everywhere, except at a sparse set of points), and*
2. *has associated with it an estimation technique which is minimax within these classes.*

It was the merit of David Donoho and Iain Johnstone (Donoho and Johnstone, 1992a) to recognize that Besov spaces, which play a central role in Yves Meyer's mathematical theory of wavelets (Meyer, 1990), provide an adequate answer. They are perfectly suited to nonlinear systems which have sparse singularities and otherwise are smooth. The methods used can be qualified as local function expansions, they provide a combination of local averaging and short-range interpolation.

**Handling geometric issues in the multi-dimensional case.** If, for instance, the function of two variables  $f(x, y)$  is “regular” in  $x$  (in extreme case  $f(x, y) = g(y)$  does not depend on  $x$ ) and “irregular” in  $y$ , then an “intelligent” estimation algorithm would approximate thoroughly  $f$  only in  $y$ -direction. This way the problem can be reduced to that of function estimation in 1-dimensional, with the correspondent improvement in the rate of convergence. This effect is often called “the dimensionality reduction” in the statistical literature. Starting from early 1980's a variety of techniques have been proposed in the statistics literature, which exhibit this desirable feature of “dimensionality reduction”. The already mentioned *Projection Pursuit Algorithm* and neural network algorithms are examples of this idea. This contrasts with wavelet and other local averaging procedures, in which the smoothing is done over small balls  $\{x : |x - x_0| \leq \varepsilon'\}$ . It was shown in (Donoho and Johnstone, 1989) that, in a certain setting, projection based and local-averaging function estimates have complementary properties.

We will defer the discussion of the second issue to Section 8.

In the literature on non-parametric regression, the focus, so far, has been on locally adaptive bandwidths for kernel methods, see (Vieu, 1991) for an example. Adaptive local linear regression is treated in (Fan and Gijbels, 1992).

## 6.2 Wavelet shrinkage algorithms

Wavelet shrinkage algorithms offer the dual advantage of being 1/ spatially adaptive (and thus practically efficient) and of low comparatively computational cost, and 2/ supported by a complete mathematical analysis. This dual feature is rather unique, so we now concentrate on this technique.

**Non-parametric regression.** Assume that  $N$  samples of input/output observations of the following system are available:

$$Y_i = f(X_i) + w_i ,$$

---

<sup>11</sup>The CART and MARS algorithms, and the variable bandwidth kernel method were designed to handle this problem in multi-dimensional case. Surprisingly enough, the A.I. literature has proposed independently and at the same time different techniques with the same feature of “spatial adaptivity”. These include various forms of neural networks (Hunt et al., 1992).

where  $(X_i)$  and  $(w_i)$  are i.i.d. sequences of random variables,  $X_i$  is *uniformly distributed* on  $[0, 1]^d$  and  $Ew_i = 0$ ,  $Ew_i^2 \leq \sigma_w^2$ . These assumptions are introduced for the sake of simplicity. They can be weakened, in particular the (unusual) assumption that  $X$  is uniformly distributed can easily be relaxed, see (Delyon and Juditsky, 1995), this would introduce additional burden to our presentation, however.

For  $f \in L_2$ , recall the wavelet expansion

$$f(x) = \sum_{k \in \mathbf{Z}} \alpha_{0k} \Phi_{0k}(x) + \sum_{j=0}^{\infty} \sum_{k \in \mathbf{Z}^d} \sum_{l=1}^{2^d-1} \beta_{jk}^{(l)} \Psi_{jk}^{(l)}(x), \quad (64)$$

where

$$\alpha_{0k} = \int f(x) \Phi_{0k}(x) dx \quad \text{and} \quad \beta_{jk}^{(l)} = \int f(x) \Psi_{jk}^{(l)}(x) dx. \quad (65)$$

To construct an estimate of  $f$  a first idea consists in using the law of large numbers and replacing, in expansion (64), the coefficients  $\alpha_k$  and  $\beta_{jk}^{(l)}$  by their empirical estimates

$$\hat{\alpha}_{0k}(N) = \frac{1}{N} \sum_{i=1}^N Y_i \Phi_{0k}(X_i) \quad \text{and} \quad \hat{\beta}_{jk}^{(l)}(N) = \frac{1}{N} \sum_{i=1}^N Y_i \Psi_{jk}^{(l)}(X_i). \quad (66)$$

Note that the assumption that input  $X$  is uniformly distributed has been used at this point.

**Density estimation.** Assume independent observations  $X_1, \dots, X_N$  of some random variable  $X$  with unknown density  $f(x)$  are available. Again  $f$  can be expanded using (64) (65). But it turns out that

$$\alpha_{0k} = \int f(x) \Phi_{0k}(x) dx = \mathbf{E}_f \Phi_{0k}(X_i)$$

where  $\mathbf{E}_f$  denotes expectation with respect to density  $f$ , and the same holds for the  $\beta$ 's. Thus empirical estimates of the wavelet coefficients  $\alpha_k$  and  $\beta_{jk}$  are given by

$$\hat{\alpha}_{0k} = \frac{1}{N} \sum_{i=1}^N \Phi_{0k}(X_i) \quad \text{and} \quad \hat{\beta}_{jk}^{(l)} = \frac{1}{N} \sum_{i=1}^N \Psi_{jk}^{(l)}(X_i). \quad (67)$$

Thus both non-parametric regression and density estimation are faced with the same issue: in formulas (66) and (67), there may not even be  $X_i$ 's available within the support of many of the  $\Phi$ 's and  $\Psi$ 's! We shall now discuss this key point for the case of density estimation.

Obviously, in order to compute the empirical coefficient  $\hat{\beta}_{jk}^{(l)}$ , we need that at least several observations  $X_i$  hit the support of  $\Psi_{jk}^{(l)}(x)$ . Statistical laws of *loglog* type guarantee that this will generically hold for scales that are not too fine. More specifically, for  $j \leq j_{\max}$ , where

$$\frac{N}{\ln N} \leq 2^d j_{\max} \leq \frac{2N}{\ln N}$$

Thus we brute force set  $\hat{\beta}_{jk}^{(l)} = 0$  for  $j > j_{\max}$ . At this point we have built an estimator of the linear projection type, as in the case of Fourier series in Section 5. Since these estimators are linear, we cannot expect them to be efficient for Besov spaces (Kerkycharian and Picard, 1992).

**A first proposal.** Our first attempt to construct an “interesting estimate” is, following the intuition at the end of the previous section, to keep a properly chosen number of coefficients with largest absolute values, and set the others to zero. More precisely, let us consider the set  $\hat{\Lambda}_n$  of pairs  $\lambda = (j, k)$  corresponding to the  $n$  estimated wavelet coefficients  $\hat{\beta}_{jk}^{(l)}$  with largest absolute values. We construct the estimate  $\hat{f}_N$  as follows:

$$\hat{f}_N(x) = \underbrace{\sum_k \hat{\alpha}_{0k} \Phi_{0k}(x)}_{\substack{m \text{ coeffs. } \neq 0 \\ (f, \Phi \text{ compact. supp.})}} + \underbrace{\sum_{j=0}^{\infty} \sum_{k \in \mathbf{Z}^d} \sum_{l=1}^{2^d-1} \hat{\beta}_{jk}^{(l)} \Psi_{jk}^{(l)}(x)}_{\text{keep the largest } n-m \text{ coeffs.}} . \quad (68)$$

The following result can be proved about estimate (68) (see (15) for the definition of the Besov spaces) :

**Theorem 7** *Let  $f \in \mathcal{B}_{p\infty}^s$  with  $s \geq d/p$ ,  $\|f\|_\infty < \infty$ . If  $n = N^{1/(2s+d)}$  is selected in (68), then*

$$E\|\hat{f}_N - f\|_2^2 = O\left(\frac{\ln N}{N}\right)^{\frac{2s}{2s+d}} . \quad (69)$$

The idea of the proof of Theorem 7 is quite intuitive and typical for wavelet estimators. We follow the argument at the end of the previous section with the only following difference: since no information is available about the distribution of the error  $|\hat{\beta}_\lambda - \beta_\lambda|$  for  $\lambda \in \hat{\Lambda}_n$ , we take a cautious upper bound for it :

$$\mathbf{E} |\hat{\beta}_\lambda - \beta_\lambda|^2 1_{\{\hat{\beta} \neq 0\}} \leq \mathbf{E} \sup_{i,j,k} |\hat{\beta}_{jk}^{(l)} - \beta_{jk}^{(l)}|^2 = O\left(\frac{\ln N}{N}\right) ,$$

which explains the extra logarithmic factor in (69).

**The final solution.** Note that  $n$  in Theorem 7 depends on  $s$ , which is generally unknown. Hence, to complete the estimation algorithm, we need a method to estimate our model order  $n$ . Though Generalized Cross-Validation type techniques could be used, we prefer a somewhat different estimation approach developed by D. Donoho, I. Johnstone, G. Kerkycharian and D. Picard (see the references below). It uses simple thresholding rules <sup>12</sup> :

$$\tilde{\beta}_{jk}^{(l)} = \hat{\beta}_{jk}^{(l)} 1_{\{|\hat{\beta}| \geq \lambda_j\}} \quad (70)$$

where  $\lambda_j$  is a threshold parameter, so we set

$$\hat{f}(x) = \sum_k \hat{\alpha}_{0k} \Phi_{0k}(x) + \sum_{j=0}^{\infty} \sum_{k \in \mathbf{Z}^d} \sum_{l=1}^{2^d-1} \tilde{\beta}_{jk}^{(l)} 1_{\{|\hat{\beta}_{jk}^{(l)}| \geq \lambda_j\}} \Psi_{jk}^{(l)}(x) . \quad (71)$$

In other words, in expansion (64), we keep those empirical estimates of wavelet coefficients which exceed some properly selected threshold. How this threshold should be selected is provided by the following result :

---

<sup>12</sup>we consider here the so called “hard thresholding”, meanwhile, other rules can also be studied, for example the “soft thresholding” (Donoho and Johnstone, 1992b). See also the discussion in (Delyon and Juditsky, 1993).



**Theorem 8 ((Donoho et al., 1993a) and (Donoho et al., 1993b))** Let  $f \in \mathcal{B}_{p\infty}^s$  with  $s \geq d/p$ ,  $\|f\|_\infty < \infty$ . Select  $\lambda_j = \lambda = \sqrt{\frac{C \ln N}{N}}$ , with an appropriate  $C < \infty$ . Then

$$E\|\hat{f}_N - f\|_2^2 = O\left(\frac{\ln N}{N}\right)^{\frac{2s}{2s+d}}.$$

The constant  $C$  in the expression for the threshold parameter  $\lambda$  is a sort of an “hyperparameter” of the procedure, it can be easily estimated, see (Delyon and Juditsky, 1993) and (Donoho et al., 1993a) for related discussions. Note that the estimator  $\hat{f}_N$  is adaptive because *it does not require prior knowledge of the regularity parameter*.

#### DISCUSSION.

- Theorem 8 has the following intuitive explanation. As already mentioned, Besov classes  $\mathcal{B}_{pq}^s$  for  $p < 2$  have a special structure: a relatively small number of “important” wavelet coefficients are sufficient for obtaining a good function approximation. In the wavelet decomposition  $(\hat{\alpha}_k, \hat{\beta}_{jk}^{(l)})$  using noisy data, all coefficients are “contaminated” by noise. A Central Limit theorem argument suggests that this noise is approximately Gaussian with zero mean and variance  $O(1/N)$ . Thus a loglog law implies that the maximal error in the estimates has magnitude given by

$$\max_{j,k} |\hat{\beta}_{jk}^{(l)} - \beta_{jk}^{(l)}| \approx \sqrt{\frac{2 \ln N}{N}}.$$

Thus when small (according to threshold  $\lambda$  in Theorem 8) coefficients are shrunk to zero, noise is canceled with very high probability. On the other hand, coefficients exceeding this threshold are likely to be significantly different from zero. This property of thresholding explains another useful feature of the estimator: the estimate  $\hat{f}_N$  has the same regularity as the unknown function  $f$  to be estimated (cf. discussion in (Donoho et al., 1993b)).

- Let us now consider again our example of estimating the regression function or density  $f(x) = 1_{\{0 \leq x < a\}}$ . Theorem 8 states that the mean square rate of convergence of the wavelet estimator for any bounded function  $f \in \mathcal{B}_{s-1\infty}^s$  is very close to  $O(N^{-1})$ , which is nearly as good as the “parametric” rate of convergence, though the function we estimate is not even continuous. Let us compare the results above with the lower rate of convergence for this problem obtained in (Nemirovskij, 1985). Using the comments 2. of Section 2.2, the following lower bound is a direct corollary of the results of (Nemirovskij, 1985) which were originally formulated in terms of Sobolev spaces:

$$\inf_{\hat{f}_N} \sup_{f \in \mathcal{B}_{pq}^s} E\|\hat{f}_N - f\|_2 \geq CN^{-2s/(2s+d)} \quad (72)$$

for any estimator  $\hat{f}_N$ . As compared to (72), there is an extra logarithmic factor in the upper bound of Theorem 8. In the more subtle construction presented in (Donoho and Johnstone, 1992a), this logarithmic factor is eliminated (and even a precise minimax constant is obtained) in the case of Gaussian noises and deterministic design (observations are  $x_i = i/N$ ,  $i = 1, \dots, N$ ). In (Donoho and Johnstone, 1993a) a cross-validation procedure is proposed to adapt the optimal algorithm to unknown smoothness. Finally, in (Delyon and

Juditsky, 1993) the authors of this paper showed that properly selecting the threshold  $\lambda$  for shrinking provides the optimal rate of convergence (without a logarithmic factor). An adaptive version of this algorithm is developed in (Juditsky, 1994).

- Finally, a more practical version of this wavelet shrinking algorithm is presented and discussed in (Sjöberg et al., 1995). In particular, for the non-parametric regression problem, this version relaxes the unrealistic assumption that input  $X$  is uniformly distributed.

## 7 Application to non-parametric autoregression identification

Let us apply the above developed technique above to the problem of identification of a simple nonlinear dynamic system. Consider the system

$$y_i = f(y_{i-1}) + e_i, \quad y_0 \in \mathbf{R}, \quad i = 1, \dots, N. \quad (73)$$

where  $(e_i)$  is a sequence of independent and identically distributed random variables,  $Ee_1 = 0$ ,  $Ee_1^2 = \sigma_e^2$ . We want to estimate the function  $f(x)$ . It is clear that the accuracy of this estimate cannot be measured using usual  $L_p$  norms on  $\mathbf{R}$ . Indeed, if the system (73) is stable, then most of the observations are concentrated in a compact set, thus the value  $f(x)$  for large  $x$  cannot be reconstructed with reasonable precision. However, as we have already seen in section 1, a “common sense” measure of the accuracy (risk) of the estimate  $\hat{f}$  of  $f$  could be

$$R_N(\hat{f}, f) \equiv \mathbf{E}|\hat{f}_N(y_N) - f(y_N)|^2 \sim \int \mathbf{E} \int |\hat{f}_N(x) - f(x)|^2 \mathbf{p}_y(x) dx,$$

where  $\mathbf{p}_y$  is the invariant density of the Markov chain  $(y_i)$ . We can implement wavelet shrinkage algorithm for this system.

We are confronted here with two problems: the first one is that the observations  $y_i$  are dependent; the second one is that the density of observations  $y_i$  is not uniform. This first difficulty is handled by fine theoretical considerations without any impact on the estimation algorithm. To deal with it we require certain stability conditions of the Markov chain  $(y_i)$  (see (Delyon and Juditsky, 1995) for details). The second problem requires that the optimal threshold in the shrinkage algorithms now depends on the index  $k$  of the empiric coefficient  $\hat{\beta}_{jk}$ .<sup>13</sup> We can cope with this problem also: it appears that in this case the ideal threshold would be

$$\lambda_{jk} = \sigma_e \sqrt{\frac{2 \ln N}{N \mathbf{p}_y(2^{-j}k)}}, \quad (74)$$

where  $\mathbf{p}_y$  is the invariant density of the Markov chain  $(y_i)$ , i.e., higher threshold is used in low density regions. Though this density is not known *a priori*, it can be easily estimated using, for instance, a simple histogram estimator  $\hat{\mathbf{p}}$ .

It can be shown that under certain conditions on the underlying Markov chain this procedure supplies us with optimal (in the minimax sense) estimates of  $f$  (cf. (Delyon and Juditsky, 1995)). To illustrate the efficiency of this method we present here a simulation example: we estimate the function  $f(\cdot)$  using the fast wavelet estimator described in (Sjöberg et al., 1995) with adaptive thresholds  $\hat{\lambda}_{jk}$  estimated according to (74) with  $\mathbf{p}(\cdot)$  being substituted with  $\hat{\mathbf{p}}(\cdot)$ . We consider

---

<sup>13</sup>Recall, in the case of uniformly distributed inputs we used the uniform threshold  $\lambda_{jk} = \lambda$ .

the following three functions  $f$  in the model (73):

$$\begin{aligned} f(y) &= y - 2\text{sign}(y) && \text{“Sigmoidal” signal} \\ f(y) &= 0.9y, && \text{“Linear” signal} \\ f(y) &= y - 2(\text{sign}(y) - 0.9 \sin(y)) && \text{“Sine” signal.} \end{aligned}$$

$(e_i)$ ,  $i = 1, \dots, N$  is supposed to be a sequence of independent Gaussian random variables with  $\sigma_e^2 = 1$ . The estimate  $\hat{f}$  is computed using  $N = 1000$  observations  $y_i$ ,  $i = 1, \dots, N$ .

On figures 1–3 we present together the signal  $(y_i)$ , the resulting estimate, and the histogram estimate  $\hat{\mathbf{p}}$  of the density  $\mathbf{p}(y)$ .<sup>14</sup> The estimated values of the risk  $R_N(\hat{f}, f)$  are presented on the figures where it is called “error”. We consider that the proposed algorithm provides a good quality of visual reconstruction. Amazing quality of fit is obtained for  $f$  linear. Indeed, in this case the well known Cramer-Rao lower bound<sup>15</sup> gives

$$R_N(\hat{f}, f) \geq \frac{\sigma_e^2}{N} = 0.001.$$

One can see that the corresponding value for the wavelet estimator is of the same order of magnitude.

## 8 Estimation of high-dimensional systems

As we have seen in Section 3, the curse of dimensionality is encountered when estimating functions with large dimensional inputs, and one has to resort to methods that reduce the effective dimension of the input. All the methods reported in Section 3 can be, and have been (sometimes successfully) used in estimation. However, there are only a few theoretical results available today. In this section we will show how the approximation technique given in Section 3.4 can be adapted to the estimation problem.

### 8.1 Wavelets

We will construct now an estimation algorithm based on the approximation result of Theorem 4. It is essentially based on the idea of truncation proposed in (Donoho and Johnstone, 1992b).

**Estimation algorithm :**

1. **sample length  $N$ , set  $\beta = d - 1/2$  and**

$$n = \sqrt{N}, \quad a_{\min} = N^{-\frac{1}{2d}}, \quad a_{\max} = \left(\frac{N}{\ln N}\right)^{1/d}$$

2. **compute**

$$\begin{aligned} \hat{u}_N(a, t) &= \frac{b^d}{N} \sum_{i=1}^N y_i \varphi(a(x_i - t)) 1_{\{a_{\min} \leq a \leq a_{\max}\}} \\ \hat{v}_N(a, t) &= a^\beta \hat{u}_N(a, t) \end{aligned}$$

---

<sup>14</sup>We use the **TeachWave** package by D. Donoho and I. Johnstone (Donoho and Johnstone, 1993b).

<sup>15</sup>this is a parametric problem!

### 3. and truncate

$$\tilde{v}_N(a, t) = \begin{cases} \hat{v}_N(a, t) & \text{if } |\hat{v}_N(a, t)| \geq K \sqrt{\frac{\ln N}{N}}; \\ 0 & \text{if } |\hat{v}_N(a, t)| < K \sqrt{\frac{\ln N}{N}} \end{cases} \quad (75)$$

4. **compute**  $C_N = \|a^{(d-1)/2} \tilde{v}_N(a, t)\|_1$  **and draw**  $n$  **independent points**  $(a_1, t_1), \dots, (a_n, t_n)$  **according to the density**  $\tilde{w}_N(a, t) = |a^{(d-1)/2} \tilde{v}_N(a, t)|/C_N$ ;

5. **compute**

$$\hat{f}_N(x) = C_N n^{-1} \sum_{k=1}^n (a_k)^{d/2} \text{sign}(\hat{u}(a_k, t_k)) \psi(a_k(x - t_k)) \quad (76)$$

The key point of the algorithm is the truncation procedure of point **3** in which we shrink to 0 the “density”  $\hat{v}_N$  if its amplitude is less than  $K \sqrt{\ln N/N}$ . This step has a nice heuristic explanation (cf. Section 6.2). Note that  $\hat{v}_N(a, t)$  can be decomposed as  $\hat{v}_N(a, t) = v(a, t) + \varepsilon_N(a, t)$ , where  $v_N(\cdot)$  is the corresponding true density, and  $\varepsilon(\cdot)$  is a centered random noise with variance of order  $N^{-1}$ . On the other hand, the condition  $\|w\|_1 = \|a^{(d-1)/2} v\|_1 < \infty$  implies that  $v(a, t) = o(1/\sqrt{N})$  on a subset of  $\{a \leq (\frac{N}{\ln N})^{1/d}\}$  of the underlying measure. So, putting all these terms to the sum (76) would mean adding noise to our estimate and negligible useful information. Therefore, we can improve significantly our estimate when shrinking these terms to the zero.

The only parameter of the algorithm to be chosen is the truncation parameter  $K$  in (75). It depends on  $\|f\|_\infty$  and on  $Ee_1^2$ , and can be estimated “on line”. Furthermore, the algorithm is not too sensitive to this parameter, and it can be often chosen on the basis of the available a priori information.

The theoretic analysis reveals that if the function  $f(\cdot)$  is bounded and  $f \in \mathcal{W}_1^{d/2+\varepsilon} \cap \mathcal{W}_2^{d/4+\varepsilon}$  for some  $\varepsilon > 0$ , then the following bound holds for the estimate  $\hat{f}_N$ :

$$E\|\hat{f}_N - f\|_2^2 \leq C \left( \frac{\ln N}{N} \right)^{1/2}$$

for some  $C < \infty$ . This means that the proposed algorithm achieves “up to  $\varepsilon$ ” the rate of convergence of the estimators discussed in Section 6.2.

In (Delyon et al., 1994) another truncation algorithm is proposed which allows us to simplify significantly the very unpleasant part of the algorithm – the drawing independent samples from the density  $\tilde{w}_N$  (line 4 of the algorithm).

## 9 Conclusion : the gap between theory and everyday practice

In this paper, we have surveyed and discussed part of the mathematical foundations of nonlinear black-box modeling. What mathematics tells us can be summarized as follows :

- The bias/variance trade-off is a key factor, as it is always in system identification when model order is not known a priori. Error variance depends on how many parameters are used for fitting. Thus efforts concentrate on reducing the bias without increasing the number of parameters in the model. Since we are dealing with nonlinear systems, there is much more flexibility in handling this problem than in linear system identification. Thus we have paid a lot of attention to *function approximation* issues as a prior topic before considering estimation from noisy data.

- Function approximation very much depends on the function space in consideration. Thus we are faced with the new problem of having to specify the function space which our unknown system is supposed to belong to, this is prior information referred to as the *smoothness class*. Such kind of prior information is not easily at hand, however, thus it is of primary importance to design procedures which are equally efficient for various classes of functions.
- Approximation (and estimation) methods have been known for a long time, which perform equally well for various smoothness classes of systems, provided that they consist of *uniformly smooth* systems. Among such methods one finds classics of non-parametric statistics such as kernel or linear projection estimates.
- These methods perform poorly on systems that are smooth, but with some spikes and jumps. Unfortunately, such nonlinear systems are frequently encountered in practice. Approximation and estimation methods have been proposed, which are *spatially adaptive*, i.e., which are able to locally adapt the smoothness of the approximants or estimates to the function to be approximated. They proved efficient and successful in practice. They were more or less attractive, depending on their computational and memory cost (as well as psychological appeal). Within this large army of methods, neural networks reached the top in celebrity. Results are available which mathematically support the success of these methods for *approximation*, but few results are available to support their use for *estimation*.
- *Besov spaces* revealed to be an adequate parameterized family of function spaces to model spikyness and jumpiness in a tunable way. They support most of the mathematical results about approximation and convergence rates.
- *Wavelets* revealed to be nicely associated with Besov spaces since Besov norms are easily evaluated using wavelet decompositions. Thus wavelet based estimation algorithms are the only class of algorithms for which complete analysis is available today, both for approximation and estimation. These theoretical results show optimality of these algorithms.
- Getting *lower bounds* for convergence rates of estimation procedures is much more involved than for parametric estimation, in which key tools are Cramer-Rao bound and Fisher information. Lower rates and minimax optimality were introduced to this end. These tools are much more technical and difficult to use than Cramer-Rao bound and Fisher information.
- The *curse of dimensionality* refers to the fact that one barely has enough data for fitting when the input dimension is large. The notion of “effective dimension” which we discussed somewhat informally plays an important role. Function classes of low “effective dimension” can be found and used for analysis in the case of large dimensional inputs.
- Altogether, non-parametric estimation can also be considered as maximum likelihood estimation for unknown systems with unknown and unbounded model order. Consistency, convergence rates, and AIC/BIC/.../XIC criteria are the classics. Two difficulties occur in applying this point of view to black-box nonlinear system identification. First, the “good” parameter dimension is typically very large so that AIC/BIC/.../XIC are not very practical. Second, and most important from the theoretical point of view, asymptotic results involve the argument of the minimum of the likelihood. But the likelihood is very non-convex and no procedure is provably known to find the minimum. This is why we

did not pay a large attention to this point of view in this paper, where emphasis was on mathematical results only.

Obviously, the everyday practice in nonlinear black-box modeling is quite different from the direct implementation of mathematical advises. Now, there are several factors which may explain the gap between the mathematical foundations and the practical confidence in this or that method, in particular :

- Computational cost and memory requirements could be formally considered, but have not been discussed in this paper.
- How bad is a non convex functional for optimization, especially in large dimension, is hard to assess.

Since authors of this article and the companion one (Sjöberg et al., 1995) are the same, they would be quite schizophrenic if they would blindly and only trust the mathematics as a guideline for implementation. Nonetheless, mathematical results have the advantage of providing theoretically sound intuition and guidelines about how and why methods perform good or bad in various situations.

## References

- Akaike, H. (1970). Statistical predictor identification. *Ann. Inst. Math. Statist.*, 22:203–217.
- Barron, A. (1993). Universal approximation bounds for superpositions of a sigmoidal function. *IEEE Trans. on Information Theory*, 39(3).
- Bellman, R. (1966). *Adaptive Control Processes*. Princeton University Press, Princeton.
- Besov, O. (1959). On a family of functional spaces: Embedding theorems and applications. *Doklady Acad. Nauk SSSR*, 126, 1163-1165.
- Breiman, L. (1993). Hinging hyperplanes for regression, classification and function approximation. *IEEE Trans. on Information Theory*, 39(3):999–1013.
- Breiman, L., Friedman, J., Olshen, J., and Stone, C. (1984). *Classification and regression trees*. Wadsworth, Belmont, California.
- Cencov, N. (1982). Statistical decision rules and optimal inference. *Amer. Math. Soc. Transl.*, 53. Providence, R.I.
- Craven, P. and Wahba, G. (1979). Smoothing noisy data with spline functions. *Numer. Math.*, 31:337–403.
- Daubechies, I. (1992). *Ten lectures on wavelets*. CBMS-NSF regional series in applied mathematics, CBMS-NSF regional conference.
- Delyon, B. and Juditsky, A. (1993). Wavelet estimators, global error measures revisited. Technical report 782, irisa.
- Delyon, B. and Juditsky, A. (1995). Optimal estimators for functional autoregression. Tech. rep. irisa, in preparation.

- Delyon, B., Juditsky, A., and Benveniste, A. (1994). Accuracy analysis for wavelet networks. *IEEE Transactions on neural networks*. to appear.
- DeVore, R., Jawerth, B., and Popov, V. (1994). Compression of wavelet decompositions. *Amer. J. Math.*, To appear.
- Devroye, L. (1982). Any discrimination rule can have an arbitrary bad probability of error for final sample size. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, PAMI-4:154–157.
- Devroye, L. and Györfi, L. (1985). *Nonparametric Density Estimation  $L_1$  View*. J. Wiley, New-York.
- Devroye, L. and Wagner, T. (1980). Distribution free consistency result in nonparametric discrimination and regression function estimation. *The Annals of Statistics*, 8(2):231–239.
- Donoho, D. and Johnstone, I. (1989). Projection-based approximation and a duality with kernel methods. *Ann. Statist.*, 17:58–106.
- Donoho, D. and Johnstone, I. (1992a). Minimax estimation via wavelet shrinkage. Technical report, department of statistics, stanford university [ftp playfair.stanford.edu](ftp://playfair.stanford.edu).
- Donoho, D. and Johnstone, I. (1992b). Minimax risk over  $l_p$ -balls. Technical report, department of statistics, stanford university [ftp playfair.stanford.edu](ftp://playfair.stanford.edu).
- Donoho, D. and Johnstone, I. (1993a). Adapting to unknown smoothness via wavelet shrinkage. Technical report, department of statistics, stanford university, [ftp playfair.stanford.edu](ftp://playfair.stanford.edu).
- Donoho, D. and Johnstone, I. (1993b). Teachwave v 0.550. Department of statistics, stanford university [wavelab@playfair.stanford.edu](mailto:wavelab@playfair.stanford.edu).
- Donoho, D., Johnstone, I., Kerkycharian, G., and Picard, D. (1993a). Density estimation by wavelet thresholding. Technical report, department of statistics, stanford university [ftp playfair.stanford.edu](ftp://playfair.stanford.edu).
- Donoho, D., Johnstone, I., Kerkycharian, G., and Picard, D. (1993b). Wavelet shrinkage: Asymptopia. Manuscript on [ftp playfair.stanford.edu](ftp://playfair.stanford.edu).
- Duflo, M. (1993). *Recursive Stochastic Methods*. Springer-Verlag, Berlin.
- Efroimovich, S. and Pinsker, M. (1982). Estimation of square-integrable spectral density based on a sequence of observations. *Problems of Information Transmission (in Russian)*, pp. 182-196.
- Efroimovich, S. and Pinsker, M. (1983). Estimation of square-integrable probability density of a random variable. *Problems of Information Transmission (in Russian)*, pp. 175-189.
- Efroimovich, S. and Pinsker, M. (1984). A learning algorithm for nonparametric filtering. *Avtomatika i Telemekhanika (in Russian)*, 11, 58-65.
- Fan, J. and Gijbels, I. (1992). Variable bandwidth and local linear regression smoothers. *Annals of Statistics*, 20:2008–2036.

- Friedman, J. (1991). Multivariate adaptive regression splines (with discussion). *The Annals of Statistics*, 19:1–141.
- Friedman, J. and Stuetzle, W. (1981). Projection pursuit regression. *J. Amer. Stat. Assoc.*, 76:817–823.
- Härdle, W. and Marron, J. (1985). Optimal bandwidth selection in nonparametric regression function estimation. *The Annals of Statistics*, 13:1465–1481.
- Huber, P. (1985). Projection pursuit (with discussion). *The Annals of Statistics*, 13:435–475.
- Hunt, K., Sbarbaro, D., Zbikowski, R., and Gawthrop, P. (1992). Neural networks for control systems — a survey. *Automatica*, 28 n°6:1083–1112.
- Ibragimov, I. and Khasminskij, R. (1981). *Statistical Estimation Asymptotic Theory*. Springer-Verlag, Berlin.
- Jaffard, S. and Laurentçot, P. (1989). *Wavelets : A Tutorial*, chapter Wavelets and P.D.E.’s. Academic Press.
- Juditsky, A. (1994). Adaptive wavelet estimators. Technical report 815, irisa.
- Kerkycharian, G. and Picard, D. (1992). Density estimation in besov spaces. *Stat. and Prob. Letters*, 13, 15-24.
- Kolmogorov, A. (1957). On the representation of continuous functions of several variables by superposition of continuous functions of one variable and addition. *Dokl*, 114(1):679–681.
- Korostelev, A. and Tsybakov, A. (1981). *Minimax Theory of Image Reconstruction*. Springer-Verlag, Berlin.
- Li, K. (1986). Asymptotic optimality of  $c_L$  and generalized cross-validation in ridge regression and application to the spline smoothing. *The Annals of Statistics*, 14:1101–1112.
- Li, K. (1987). Asymptotic optimality of  $c_L$  and generalized cross-validation : discrete index set. *The Annals of Statistics*, 15:958–975.
- Ljung, L. (1987). *System Identification : Theory for the User*. Prentice Hall Information and System Sciences Series. Prentice Hall.
- Lorentz, G. (1976). The 13th problem of hilbert. In Browder, F., editor, *Mathematical Developments Arising from Hilbert Problems*. Am. MATH. Soc., Providence, R.L.
- Mallows, C. (1973). Statistical predictor identification. *Technometrics*, 15:661–675.
- Meyer, Y. (1990). *Ondelettes et Opérateurs*. Hermann.
- Morgan, J. and Sonquist, J. (1963). Problems in the analysis of survey data, and a proposal. *J. Amer. Stat. Assoc.*, 58:415–434.
- Müller, H.-G. and Stadtmüller, U. (1987). Variable bandwidth kernel estimators of regression curves. *The Annals of Statistics*, 15(1), 182–201.
- Nadaraya, E. (1964). On estimating regression. *Theory of Prob. and Appl.*, 9:141–142.



- Nemirovskij, A. (1985). Nonparametric estimation of smooth regression functions. *Izv. Acad. Nauk SSSR, Techn. Kibernet. (in Russian)*, 3, 50-60.
- Oppenheim, G. and Portier, B. (1990). Commande adaptative du processus de markov  $x_{t+1} = f_t + u_t + x_t$ ,  $t \in n$ . Technical Report 90-18, Université d'Orsay.
- Parzen, E. (1962). On estimation of probability density function and the mode. *Ann. of Math. Stat.*, 33:1065-1076.
- Petrushev, P. and Popov, V. (1987). *Rational Approximation Of Real Functions*. Cambridge University Press, Cambridge.
- Pinkus, A. (1985). *n-Widths in Approximation Theory*. Springer-Verlag, Berlin.
- Poggio, T. and Girosi, F. (1990). Networks for approximation and learning. *Proceedings of the IEEE*, 78(9):1481-1497.
- Polyak, B. and Tsybakov, A. (1990). Asymptotical optimality of  $c_p$  criterion for projection regression estimates. *Theory of Prob. and Appl.*, 35:305-317.
- Portier, B. (1992). *Estimation non paramétrique et commande adaptative de processus Markoviens non linéaires*. Ph. d. thesis, Université Paris Sud, Orsay.
- Rice, J. (1984). Bandwidth choice for nonparametric regression. *The Annals of Statistics*, 12:1215-1230.
- Rosenblatt, M. (1956). Remarks on some nonparametric estimates of density functions. *Ann. of Math. Stat.*, 27:832-835.
- Rosenblatt, M. (1971). Curve estimation. *Ann. of Math. Stat.*, 42(6):1815-1842.
- Sickel, W. (1990). Spline representations of functions in besov-triebel-lizorkin spaces on  $\mathbf{r}^n$ . *Forum Math.*, 2, 451-476.
- Sjöberg, J., Zhang, Q., Ljung, L., Benveniste, A., Deylon, B., Hjalmarsson, P.-Y. G. H., and Juditsky, A. (1995). Nonlinear black-box modeling in system identification: a unified overview. *Automatica*, page to appear.
- Söderström, T. and Stoica, P. (1989). *System Identification*. Prentice-Hall International, Hemel Hempstead, Hertfordshire.
- Sontag, E. (1981). Nonlinear regulation: the piecewise linear approach. *IEEE Trans. on Automatic Control*, 26:346-358.
- Stone, C. (1982). Optimal global rates of convergence for nonparametric regression. *The Annals of Statistics*, 10:1040-1053.
- Triebel, H. (1983). *Theory of Function Spaces*. Birkhäuser Verlag, Berlin.
- Triebel, H. (1993). *Theory of Function Spaces II*. Birkhäuser Verlag, Berlin.
- Vieu, P. (1991). Nonparametric regression: Optimal local bandwidth choice. *J.R. Statist. Soc. Ser. B*, 53:453-464.
- Wahba, G. (1991). *Spline functions for observational data*. SIAM, Philadelphia, PA.
- Watson, G. (1969). Smooth regression analysis. *Sankhya, Series, A*(26):359-372.

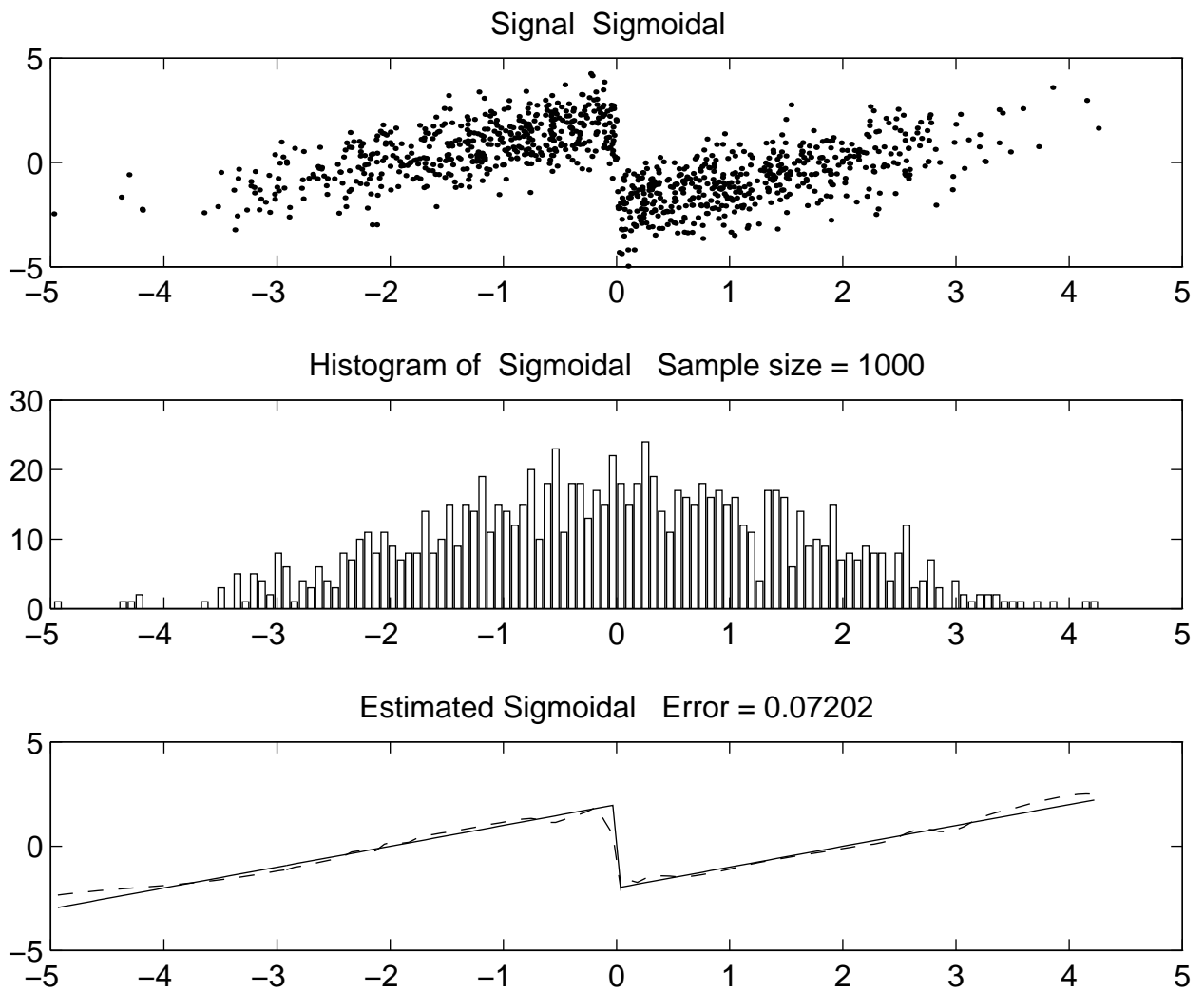


Figure 1:

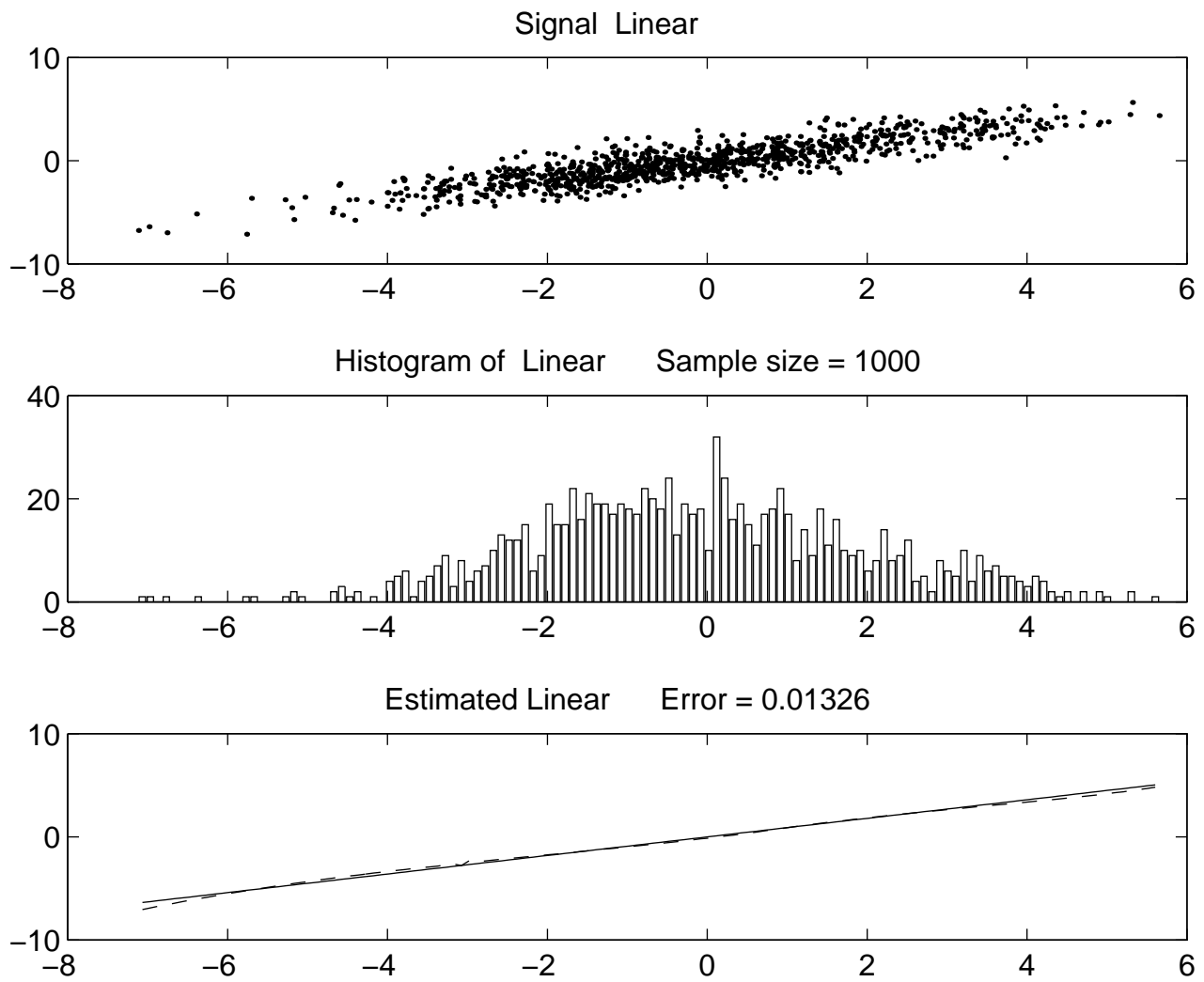


Figure 2:

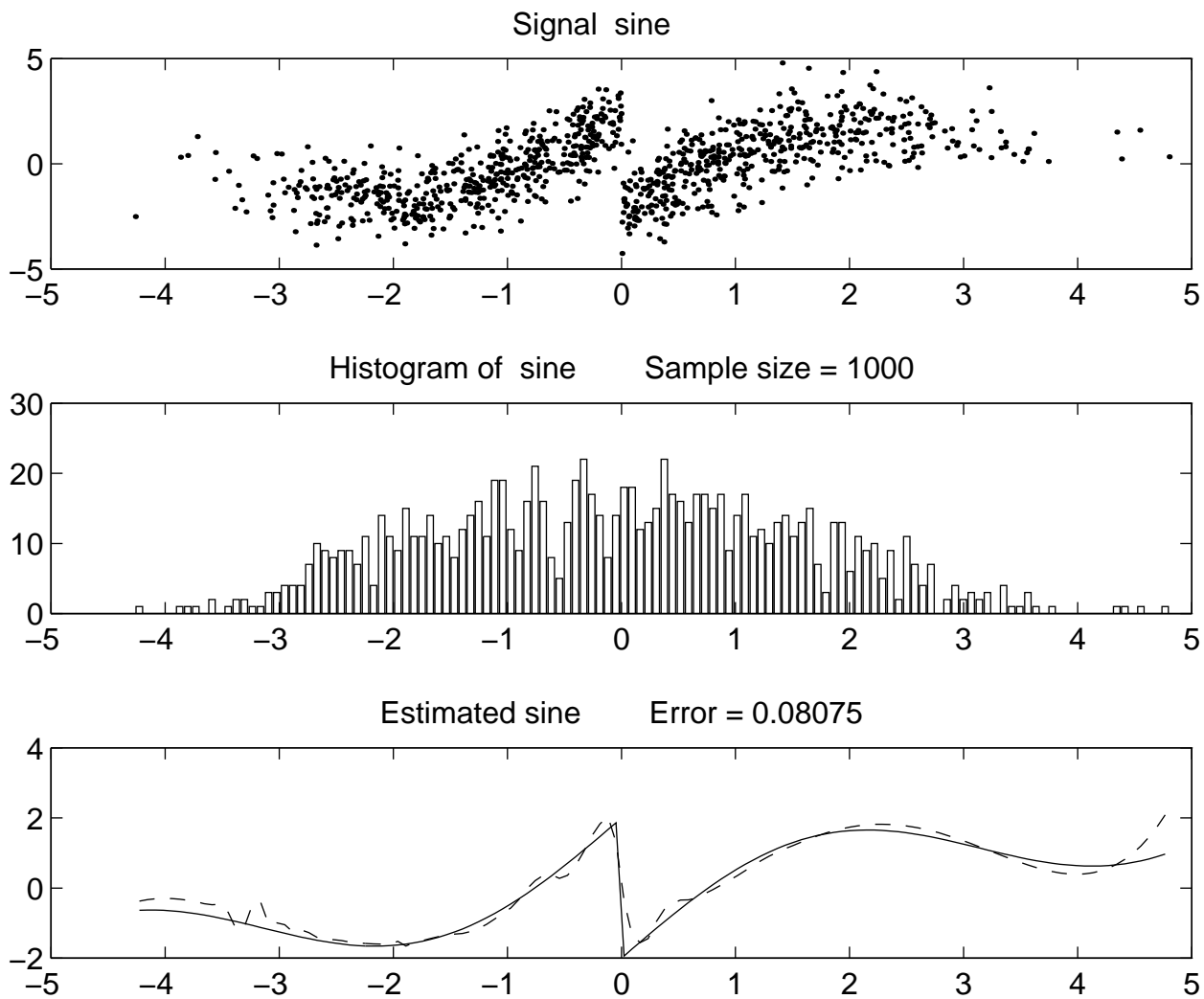


Figure 3: