

# Neural Networks in System Identification

J. Sjöberg, H. Hjalmarsson, L. Ljung

*Department of Electrical Engineering, Linköping University, S-581 83 Linköping, Sweden  
E-mail: hakan@isy.liu.se, sjoberg@isy.liu.se, ljung@isy.liu.se*

**Abstract.** Neural Networks are non-linear black-box model structures, to be used with conventional parameter estimation methods. They have good general approximation capabilities for reasonable non-linear systems. When estimating the parameters in these structures, there is also good adaptability to concentrate on those parameters that have the most importance for the particular data set.

**Key Words.** Neural Networks, Parameter estimation, Model Structures, Non-Linear Systems.

## 1. EXECUTIVE SUMMARY

### 1.1. Purpose

The purpose of this tutorial is to explain how Artificial Neural Networks (NN) can be used to solve problems in System Identification, to focus on some key problems and algorithmic questions for this, as well as to point to the relationships with more traditional estimation techniques. We also try to remove some of the “mystique” that sometimes has accompanied the Neural Network approach.

### 1.2. What's the problem?

The identification problem is to infer relationships between past input-output data and future outputs. Collect a finite number of past inputs  $u(k)$  and outputs  $y(k)$  into the vector  $\varphi(t)$

$$\varphi(t) = [y(t-1) \dots y(t-n_a) u(t-1) \dots u(t-n_b)]^T \quad (1)$$

For simplicity we let  $y(t)$  be scalar. Let  $d = n_a + n_b$ . Then  $\varphi(t) \in \mathbb{R}^d$ . The problem then is to understand the relationship between the next output  $y(t)$  and  $\varphi(t)$ :

$$?y(t) \leftrightarrow \varphi(t)? \quad (2)$$

To obtain this understanding we have available a set of observed data (sometimes called the “training set”)

$$Z^N = \{[y(t), \varphi(t)] | t = 1, \dots, N\} \quad (3)$$

From these data we infer a relationship

$$\hat{y}(t) = \hat{g}_N(\varphi(t)) \quad (4)$$

We index the function  $g$  with a “hat” and  $N$  to emphasize that it has been inferred from (3). We also place a “hat” on  $y(t)$  to stress that (4) will in practice not be an exact relationship between  $\varphi(t)$  and the observed  $y(t)$ . Rather  $\hat{y}(t)$  is the “best guess” of  $y(t)$  given the information  $\varphi(t)$ .

### 1.3. Black boxes

How to infer the function  $\hat{g}_N$ ? Basically we search for it in a parameterized family of functions

$$\mathcal{G} = \{g(\varphi(t), \theta) | \theta \in D_{\mathcal{M}}\} \quad (5)$$

How to choose this parameterization? A good, but demanding, choice of parameterization is to base it on physical insight. Perhaps we know the relationship between  $y(t)$  and  $\varphi(t)$  on physical grounds, up to a handful of physical parameters (heat transfer coefficients, resistances, ...). Then parameterize (5) accordingly.

This tutorial only deals with the situation when physical insight is *not* used; i.e. when (5) is chosen as a flexible set of functions capable of describing almost any true relationship between  $y$  and  $\varphi$ . This is the *black-box approach*.

Typically, function expansions of the type

$$g(\varphi, \theta) = \sum_k \theta(k) g_k(\varphi) \quad (6)$$

are used, where

$$g_k(\varphi) : \mathbb{R}^d \rightarrow \mathbb{R}$$

and  $\theta(k)$  are the components of the vector  $\theta$ . For

example, let

$$g_k(\varphi) = \varphi_k \quad (k\text{th component of } \varphi) \quad k = 1, \dots, d.$$

Then, with (1)

$$y(t) = g(\varphi(t), \theta)$$

reads

$$\begin{aligned} y(t) + a_1 y(t-1) + \dots + a_{n_a} y(t-n_a) = \\ b_1 u(t-1) + \dots + b_{n_b} u(t-n_b) \end{aligned}$$

if

$$a_i = -\theta(i) \quad b_i = \theta(n_a + i)$$

so the familiar ARX-structure is a special case of (6), with a linear relationship between  $y$  and  $\varphi$ .

#### 1.4. Nonlinear black box models

The challenge now is the non-linear case: to describe general, non-linear, dynamics. How to select  $\{g_k(\varphi)\}$  in this general case? We should thus be prepared to describe a “true” relationship

$$\hat{y}(t) = g_0(\varphi(t))$$

for any reasonable function  $g_0 : \mathbb{R}^d \rightarrow \mathbb{R}$ . The first requirement should be that  $\{g_k(\varphi)\}$  is a basis for such functions, i.e. that we can write

$$[R1] : \quad g_0(\varphi) = \sum_{k=1}^{\infty} \theta(k) g_k(\varphi) \quad (7)$$

for any reasonable function  $g_0$  using suitable coefficients  $\theta(k)$ . There is of course an infinite number of choices of  $\{g_k\}$  that satisfy this requirement, the classical perhaps being the basis of polynomials. For  $d = 1$  we would then have

$$g_k(\varphi) = \varphi^k$$

and (7) becomes Taylor or Volterra expansion. In practice we cannot work with infinite expansions like (7). A second requirement on  $\{g_k\}$  is therefore to produce “good” approximations for finite sums: In loose notation:

$$[R2] : \quad \| g_0(\varphi) - \sum_{k=1}^n \theta(k) g_k(\varphi) \|$$

“decreases quickly as  $n$  increases” (8)

There is clearly no uniformly good choice of  $\{g_k\}$  from this respect: It will all depend on the class of functions  $g_0$  that are to be approximated.

#### 1.5. Estimating $\hat{g}_N$

Suppose now that a basis  $\{g_k\}$  has been chosen, and we try to approximate the true relationship by a finite number of the basis functions:

$$\hat{y}(t|\theta) = g(\varphi(t), \theta) = \sum_{k=1}^n \theta(k) g_k(\varphi(t)) \quad (9)$$

where we introduce the notation  $\hat{y}(t|\theta)$  to stress that  $g(\varphi(t), \theta)$  is a “guess” for  $y(t)$  given the information in  $\varphi(t)$  and given a particular parameter value  $\theta$ . The “best” value of  $\theta$  is then determined from the data set  $Z^N$  in (9) by

$$\hat{\theta}_N = \arg \min_{k=s}^N |y(t) - \hat{y}(t|\theta)|^2 \quad (10)$$

The model will be

$$\hat{y}(t) = \hat{y}(t|\hat{\theta}_N) = \hat{g}_N(\varphi(t)) = g(\varphi(t), \hat{\theta}_N) \quad (11)$$

#### 1.6. Properties of the estimated model

Suppose that the actual data have been generated by

$$y(t) = g_0(\varphi(t)) + e(t) \quad (12)$$

where  $\{e(t)\}$  is white noise with variance  $\lambda$ . The estimated model (11) (i.e. the estimated parameter vector  $\hat{\theta}_N$ ) will then be a random variable that depends on the realizations of both  $e(t), t = 1, \dots, N$  and  $\varphi(t), t = 1, \dots, N$ . Denote its expected value by

$$\mathbb{E} \hat{g}_N = g_n^* = \sum_{k=1}^n \theta^*(k) g_k \quad (13)$$

where we used subscript  $n$  to emphasize the number of terms used in the function approximation.

Then under quite general conditions

$$\mathbb{E} |\hat{g}_N(\varphi(t)) - g_n^*(\varphi(t))|^2 = \lambda \cdot \frac{m}{N} \quad (14)$$

where  $\mathbb{E}$  denotes expectation both with respect to  $\varphi(t)$  and  $\hat{\theta}_N$ . Moreover,  $m$  is the number of estimated parameters, i.e.,  $\dim \theta$ . The total error thus becomes

$$\begin{aligned} \mathbb{E} |\hat{g}_N(\varphi(t)) - g_0(\varphi(t))|^2 = \\ \| g_0(\varphi(t)) - g_n^*(\varphi(t)) \|^2 + \lambda \cdot \frac{m}{N} \end{aligned} \quad (15)$$

The first term here is an approximation error of the type (8). It follows from (15) that there is a trade-off in the choice of how many basis functions to use. Each included basis function increases the variance error by  $\lambda/N$ , while it decreases the bias error by an amount that could be less than so. A third requirement on the choice of  $\{g_k\}$  is thus to

[R3] Have a scheme that allows the exclusion of spurious basis functions from the expansion.

Such a scheme could be based on a priori knowledge as well as on information in  $Z^N$ .

### 1.7. Basis functions

Out of the many possible choice of basis functions, a large family of special ones have received most of the current interest. They are all based on just one fundamental function  $\sigma(\varphi)$ , which is scaled in various ways, and centered at different points, *i.e.*

$$g_k(\varphi) = \sigma(\beta_k^T(\varphi + \gamma_k)) = \sigma(\beta_k^T \varphi + \eta_k) = \sigma(\varphi, \eta_k) \quad (16)$$

where  $\gamma_k = \beta_k^T \tilde{\gamma}_k$  and  $\eta_k$  is the  $d+1$ -vector

$$\eta_k = [\beta_k, \gamma_k] \quad (17)$$

Such a choice is not at all strange. A very simplistic approach would be to take  $\sigma(\varphi)$  to be the indicator function (in the case  $d=1$ ) for the interval  $[0, 1]$ :

$$\sigma(\varphi) = \begin{cases} 1 & \varphi \in [0, 1] \\ 0 & \varphi \notin [0, 1] \end{cases}$$

For a countable collection of  $\eta_k$  (e.g. assuming all rational numbers) the functions  $g_k(\varphi)$  would then contain indicator functions for any interval, arbitrarily small and placed anywhere along the real axis. Not surprisingly, these  $\{g_k\}$  will be a basis for all continuous functions. Equivalently, it could be threshold function

$$\sigma(\varphi) = \begin{cases} 1 & \varphi > 0 \\ 0 & \varphi \leq 0 \end{cases} \quad (18)$$

since the basic indicator function is just the difference between two threshold functions.

### 1.8. What is the Neural Network Identification Approach?

The basic Neural Network (NN) used for System Identification (one hidden layer feedforward net) is indeed the choice (16) with a smooth approximation for (18), often

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

Include the parameter  $\eta$  in (16)-(17) among the parameters to be estimated,  $\theta$ , and insert into (9). This gives the Neural Network model structure

$$\hat{y}(t|\theta) = \sum_{k=1}^n \alpha_k \sigma(\beta_k \varphi + \gamma_k)$$

$$\theta = [\alpha_k, \beta_k, \gamma_k], \quad k = 1, \dots, n \quad (19)$$

The  $n \cdot (d+2)$ -dimensional parameter vector  $\theta$  is then estimated by (10).

### 1.9. Why has Neural Networks attached so much interest?

This tutorial points at two main facts.

1. The NN function expansion has good properties regarding requirement [R2] for nonlinear functions  $g_0$  that are “localized”; *i.e.* there is not much nonlinear effects going to the infinity. This is a reasonable property for most real life physical functions. More precisely, see (69) in Section 9.

2. There is a good way to handle requirement [R3] by *implicit or explicit regularization* (See Section 3)

### 1.10. Why has there been so much hesitation about Neural Networks in the statistics and system identification communities?

Basically, because NN is, formally, just one of many choices of basis functions. Algorithms for achieving the minimum in (10), and the statistical properties of the type (15) are all of general character and well known for the more traditional model structures used. They have typically been reinvented and rediscovered in the NN literature and been given different names there. This certainly has had an alienating effect on the “traditional estimation” communities.

### 1.11. Related approaches

Actually, the general family of basis functions (16), is behind both *Wavelet Transform Networks* and estimation of *Fuzzy Models*. A companion tutorial, Benveniste *et al.* (1994), explains these connections in an excellent manner.

### 1.12. Organization of the tutorial

In Section 2 we shall give some general background about function approximation. This overlaps and complements the corresponding discussion in Benveniste *et al.* (1994). Sections 3 and 4 deal with the fundamentals of estimation theory with relevance for Neural Networks. The basic Neural Network structures are introduced in Section 5. The question of how to fit the model structure to data is discussed in Section 6. In Sections 8 and 9 the perspective is widened

to discuss how Neural Networks relate to other black box non-linear structures. Also these sections deal with similar questions as Benveniste *et al.* (1994). Section 11 describes the typical structures that Neural Networks give rise to when applied to dynamical systems. The final Section 12 describes applications and research problems in the area.

## 2. THE PROBLEM

### 2.1. Inferring relationships from data

A wide class of problems in disciplines such as classification, pattern recognition and system identification can be fit into the following framework.

A set of observations (data)

$$Z^N = \{y(t), \Phi(t)\}_{t=1}^N$$

of two physical quantities  $y \in \mathbb{R}^p$  and  $\Phi \in \mathbb{R}^r$  is given. It may or may not be known which variables in  $\Phi$  influence  $y$ . There may also be other, non-measured, variables  $v$  that influence  $y$ . Based on the observations  $Z^N$ , infer how the variables in  $\Phi$  influence  $y$ .

Let  $\varphi$  be the variables in  $\Phi$  that influence  $y$ , then we could represent the relation between  $\varphi$ ,  $v$  and  $y$  by a function  $g_0$

$$y = g_0(\varphi, v) \quad (20)$$

The problem is thus two-fold:

1. Find which variables in  $\Phi$  that should be used in  $\varphi$ .
2. Determine  $g_0$ .

In identification of dynamical systems, finding the right  $\varphi$  is the model order selection problem. Then  $t$  represents the time index and  $\Phi(t)$  would be the collection of all past inputs and outputs.

There are two issues that have to be dealt with when determining  $g_0$ :

1. Only finite observations in the  $\varphi$ -space are available.
2. The observations are perturbed by the non-measurable variable  $\{v(t)\}$ .

1) represents the function approximation problem, *i.e.* how to do interpolation and extrapolation, which in itself is an interesting problem. Notice that there would be no problem at all if  $y$  was given for all values of  $\varphi$  (if we neglect the non-measurable input) since the function then in

fact would be defined by the data. 2) increases the difficulty further since then we cannot infer exactly how  $\varphi$  influences  $y$  even at the points of observations. Blended together, these two problems are very challenging. Below we will try to disclose the essential ingredients. For further insight in this problem see also Benveniste *et al.* (1994).

### 2.2. Prior assumptions

Notice that as stated, the problem is ill-posed. There will be far too many un-falsified models, *i.e.* models satisfying (20), if any function  $g$  and any non-measurable sequence  $\{v(t)\}$  is allowed. Thus, it is necessary to include some *a priori* information in order to limit the number of possible candidates. However, often it is difficult to provide *a priori* knowledge that is so precise that the problem becomes well-defined. To ease the burden it is common to resort to some general principles:

- 1) *Non-measurable inputs are additive.* This means that  $g_0$  is additive in its second argument, *i.e.*

$$g_0(\varphi, v) = g_0(\varphi) + v$$

This is, for example, a relevant assumption when  $\{v(t)\}$  is due mainly to measurement errors. Therefore  $v$  is often called disturbance or noise.

- 2) *Try simple things first (Occam's razor).* There is no reason to choose a complicated model unless needed. Thus, among all unfalsified models, select the simplest one. Typically, the simplest means the one that in some sense has the smoothest surface. An example is spline smoothing. Among the class  $C^2$  of all twice differentiable functions on an interval  $I$ , the solution to

$$\min_{g \in C^2} \sum (y(t) - g(\varphi(t)))^2 + \lambda \int_I (g''(\varphi))^2 d\varphi$$

is given by the cubic spline, Wahba (1990). Other ways to penalize the complexity of a function are information based criteria, such as AIC, BIC and MDL, regularization (or ridge penalty), cross-validation and shrinkage. We shall discuss these in Section 9. Part of this smoothness paradigm is that the roughness should be allowed to increase with the number of observations. If there is compelling evidence in the observations that the function is non-smooth, then the approximating function should be allowed to be more flexible. This also holds for which variables in  $\Phi(t)$  that should be included in  $\varphi(t)$ . Thus both the dimension and the entries in  $\varphi(t)$  could depend on the observations  $Z^N$ . In pure approximation theory all these smoothness

criteria are rather *ad hoc*. It is first when the non-measurable inputs are taken into account that they can be given meaningful interpretations. This will be the main topic in Section 4, see also Sections 8-9.

### 2.3. Function classes

Thus,  $g_0$  is assumed to belong to some quite general family  $\mathcal{G}$  of functions. The function estimate  $\hat{g}_N^n$  however, is restricted to belong to a possibly more limited class of functions,  $\mathcal{G}_n$  say. This family  $\mathcal{G}_n$ , where  $n$  represents the complexity of the class<sup>1</sup>, is a member of a sequence of families  $\{\mathcal{G}_n\}$  that satisfy  $\mathcal{G}_n \rightarrow \mathcal{G}$ . As explained above, the complexity of  $\hat{g}_N^n$  is allowed to depend on  $Z^N$ , i.e.  $n$  is a function of  $Z^N$ . We will indicate this by writing  $n(N)$ .

In this perspective, an identification method can be seen as a rule to choose the family  $\{\mathcal{G}_n\}$  together with a rule to choose  $n(N)$  and an estimator that given these provides an estimate  $\hat{g}_N^{n(N)}$ . Notice that both the selection of  $\{\mathcal{G}_n\}$  and  $n(N)$  can be driven by data. This possibility is, as we shall see in Section 9, very important.

Typical choices of  $\mathcal{G}$  are Hölder Balls which consist of Lipschitz continuous functions:

$$\Lambda^\alpha(C) = \{f : |f(x) - f(y)| \leq C \cdot |x - y|^\alpha\} \quad (21)$$

and  $L_p$  Sobolev Balls which have derivatives of a certain degree which belongs to  $L_p$ :

$$W_p^m(C) = \{f : \int |f^{(m)}(t)|^p dt \leq C^p\} \quad (22)$$

Recently, Besov classes and Triebel classes, Triebel (1983) have been employed in wavelet analysis. The advantage with these classes are that they allow for spatial inhomogeneity. Functions in these classes can be locally spiky and jumpy.

### 2.4. Noise assumptions

The non-measurable input  $\{v(t)\}$  is also restricted to some family  $\mathcal{V}$ . It is possible to classify these families into two categories:

1) *Deterministic*. Here  $\{v(t)\}$  is usually assumed to belong to a ball

$$|v(t)| \leq C_v \quad \forall t.$$

This is known as unknown-but-bounded disturbances and dates back to the work of Schweißpepe (1973). This assumption leads to set

---

<sup>1</sup>Typically  $n$  is the number of basis functions in the class

estimation methods, see Milanese and Vicino (1991).

2) *Stochastic*. Here  $\{v(t)\}$  is a stochastic process with certain properties. This type, which we shall focus on, is the most common one. However, for a connection between deterministic and stochastic disturbances see Hjalmarsson and Ljung (1994). The advantage with this type is that it fits with the smoothness principle. A stochastic disturbance is typically non-smooth as opposed by the function of interest  $g_0$ . This can be used to decrease the influence of the disturbance.

The challenge is to find identification methods that give good performance for as general families  $\mathcal{G}$  and  $\mathcal{V}$  as possible. For a chosen criterion (figure of merit) it is the interplay between the approximating properties of the method and the way that the disturbance corrupt the approximation that has to be considered. We shall delve into that issue in the sections that follow. Especially we shall examine what Artificial Neural Networks can offer in this respect.

### 2.5. Figures of merit

Since one is working in a function space it is natural to consider some norm of the error between the function estimate  $\hat{g}_N$  and the true function  $g_0$ . It is quite standard to use  $L_p$ -norms

$$\begin{aligned} J_p(\hat{g}_N, g_0) &= \|\hat{g}_N - g_0\|_{L_p}^p = \\ &\int |\hat{g}_N(\varphi) - g_0(\varphi)|^p dP_\varphi(\varphi) \end{aligned} \quad (23)$$

where  $P_\varphi$  is the probability distribution of  $\varphi$ . An estimator is almost surely *convergent* if

$$J_p(\hat{g}_N, g_0) \rightarrow 0 \text{ w.p. 1 as } N \rightarrow \infty \quad \forall g_0 \in \mathcal{G}.$$

In order to compare different estimators one can consider rates of convergence:  $\{\hat{g}_N\}$  converges to  $g$  with rate  $\{f_N\}$  if  $J_p(\hat{g}_N, g_0) \asymp f_N$  and  $f_N \rightarrow 0$ .<sup>2</sup>

Another figure of merit is the expected value

$$V_p^P(\hat{g}_N, g) = E[J_p(\hat{g}_N, g_0)] \quad (24)$$

where the expectation is taken over the probability space  $P$  of  $\{v(t)\}$ . With  $p = 2$  one gets the integrated mean square error (IMSE)

$$V_2^P(\hat{g}_N, g) = \int E[|\hat{g}_N(\varphi) - g_0(\varphi)|^2] dP_\varphi(\varphi). \quad (25)$$

---

<sup>2</sup> $a_N \asymp b_N$  means that  $-\infty < \liminf \frac{a_N}{b_N} < \limsup \frac{a_N}{b_N} < \infty$

This type of criteria is known as risk measures in statistics. Based on the risk, various optimality properties can be defined:

It is natural to try to minimize the risk for the worst-case: An estimator  $\hat{g}_N^*$  is said to be minimax if

$$\sup_{g_0 \in \mathcal{G}} V_p^P(\hat{g}_N^*, g_0) = \inf_{\hat{g}} \sup_{g_0 \in \mathcal{G}} V_p^P(\hat{g}_N, g_0).$$

Often it is too difficult to derive the minimax estimator and one has to resort to asymptotic theory: The estimator  $\hat{g}_N^*$  is asymptotically minimax if

$$\sup_{g_0 \in \mathcal{G}} V_p^P(\hat{g}_N^*, g_0) = \inf_{\hat{g}_N} \sup_{g_0 \in \mathcal{G}} V_p^P(\hat{g}_N, g_0)$$

as  $N \rightarrow \infty$ . An even weaker concept is the minimax rate. The estimator  $\hat{g}_N^*$  attains the minimax rate if

$$\sup_{g_0 \in \mathcal{G}} V_p^P(\hat{g}_N^*, g_0) \asymp \inf_{\hat{g}_N} \sup_{g_0 \in \mathcal{G}} V_p^P(\hat{g}_N, g_0). \quad (26)$$

Notice that the risk will depend on the assumed distribution. To safeguard against uncertainty about the distribution it is possible to consider a whole family of distributions  $\mathcal{P}$  and use

$$V_p^{\mathcal{P}}(\hat{g}_N, g_0) = \sup_{P \in \mathcal{P}} V_p^P(\hat{g}_N, g_0).$$

This is thus a minimax problem and is considered in robust statistics Huber (1981). Notice that when rates of convergence are considered, the shape of the distribution is less important. For the class of distributions where the support is unbounded and some mixing condition, the rate of convergence will be the same.

### 3. SOME GENERAL ESTIMATION RESULTS

The basic estimation set-up is what is called *non-linear regression* in statistics. The problem is as follows. We would like to estimate the relationship between a scalar  $y$  and  $\varphi \in \mathbb{R}^d$ . For a particular value  $\varphi(t)$  the corresponding  $y(t)$  is assumed to be

$$y(t) = g_0(\varphi(t)) + e(t) \quad (27)$$

where  $\{e(t)\}$  is supposed to be a sequence of independent random vectors, with zero mean values and variance

$$\mathbb{E} e(t)e^T(t) = \lambda \quad (28)$$

To find the function  $g_0$  in (27) we have the following information available:

#### 1. A parameterized family of functions

$$\mathcal{G}_m = \{g(\varphi(t), \theta) | \theta \in D_M \subset \mathbb{R}^m\} \quad (29)$$

2. A collection of observed  $y, \varphi$ -pairs:

$$Z^N = \{[y(t), \varphi(t)], t = 1, \dots, N\} \quad (30)$$

The typical way to estimate  $g_0$  is then to form the scalar valued function

$$V_N(\theta) = \frac{1}{N} \sum_{t=1}^N |y(t) - g(\varphi(t), \theta)|^2 \quad (31)$$

and determine the parameter estimate  $\hat{\theta}_N$  as its minimizing argument:

$$\hat{\theta}_N = \arg \min V_N(\theta) \quad (32)$$

The estimate of  $g_0$  will then be

$$\hat{g}_N(\varphi) = g(\varphi, \hat{\theta}_N) \quad (33)$$

Sometimes a general, non-quadratic, norm is used in (30)

$$V_N(\theta) = \frac{1}{N} \sum_{t=1}^N \ell(\varepsilon(t, \theta)) \quad (34)$$

$$\varepsilon(t, \theta) = y(t) - g(\varphi(t), \theta)$$

Another modification of (30) is to add a *regularization term*,

$$W_N(\theta) = V_N(\theta) + \delta|\theta - \theta^\#|^2 \quad (35)$$

(and minimize  $W$  rather than  $V$ ) either to reflect some prior knowledge that a good  $\theta$  is close to  $\theta^\#$  or just to improve numerical and statistical properties of the estimate  $\hat{\theta}_N$ . Again, the quadratic term in (35) could be replaced by a non-quadratic norm.

Now, what are the properties of the estimated relationship  $\hat{g}_N$ ? How close will it be to  $g_0$ ? Following some quite standard results, see e.g. Ljung (1987); Söderström and Stoica (1989), we have the following properties. We will not state the precise assumptions under which the results hold. Generally it is assumed that  $\{\varphi(t)\}$  is (quasi)-stationary and has some mixing property (i.e. that  $\varphi(t)$  and  $\varphi(t+s)$  become less and less dependent as  $s$  increases). The estimate  $\hat{\theta}_N$  is a random variable that depends on  $Z^N$ . Let  $\mathbb{E}$  denote expectation with respect to both  $e(t)$  and  $\varphi, t = 1, \dots, N$ . Let

$$\theta^* = \mathbb{E} \hat{\theta}_N$$

and

$$g^*(\varphi) = g(\varphi, \theta^*)$$

Then  $g^*(\varphi)$  will be as close as possible to  $g_0(\varphi)$  in the following sense:

$$\arg \min_{g \in \mathcal{G}_m} \mathbb{E} |g(\varphi) - g_0(\varphi)|^2 = g^*(\varphi) \quad (36)$$

where expectation  $E$  is over the distribution of  $\varphi$  that governed the observed sample  $Z^N$ . We shall call

$$g^*(\varphi) - g_0(\varphi)$$

the *bias error*. Moreover, if the bias error is small enough, the variance will be given approximately by

$$E |\hat{g}_N(\varphi) - g^*(\varphi)|^2 \approx \frac{m}{N} \lambda \quad (37)$$

Here  $m$  is the dimension of  $\theta$  (number of estimated parameters),  $N$  is the number of observed data pairs and  $\lambda$  is the noise variance. Moreover, expectation in both over  $\hat{\theta}_N$  and over  $\varphi$ , assuming, the same distribution for  $\varphi$  as in the sample  $Z^N$ . The total integrated mean square error (IMSE) will thus be

$$E |\hat{g}_N(\varphi) - g_0(\varphi)|^2 = ||g^*(\varphi) - g_0(\varphi)||^2 + \frac{m}{N} \lambda \quad (38)$$

Here the double bar norm denotes the functional norm, integrating over  $\varphi$  with respect to its distribution function when the data were collected. Now, what happens if we minimize the regularized criterion  $W_N$  in (35)?

1. The value  $g^*(\varphi)$  will change to the function that minimizes

$$E |g(\varphi, \theta) - g_0(\varphi)|^2 + \delta |\theta - \theta^\#|^2 \quad (39)$$

2. The variance (37) will change to

$$E |\hat{g}_N(\varphi) - g^*(\varphi)|^2 \approx \frac{r(m, \delta)}{N} \cdot \lambda \quad (40)$$

where

$$r(m, \delta) = \sum_{k=1}^m \frac{\sigma_k^2}{(\sigma_k + \delta)^2} \quad (41)$$

where  $\sigma_i$  are the eigenvalues (singular values) of  $E V_N''(\theta)$ , the second derivative matrix (the Hessian) of the criterion.

How to interpret (41)? A redundant parameter will lead to a zero eigenvalue of the Hessian. A small eigenvalue of  $V''$  can thus be interpreted as corresponding to a parameter (combination) that is not so essential: "A spurious parameter". The regularization parameter  $\delta$  is thus a threshold for spurious parameters. Since the eigenvalues  $\sigma_i$  often are widely spread we have

$$r(m, \delta) \simeq m^\# = \# \text{ of eigenvalues of } V'' \text{ that are larger than } \delta$$

We can think of  $m^\#$  as "the efficient number of parameters in the parameterization". Regularization thus decreases the variance, but typically increases the bias contribution to the total error.

#### 4. THE BIAS/VARIANCE TRADE-OFF

Consider now a sequence of parameterized function families

$$\mathcal{G}_n = \{g_n(\varphi(t), \theta) | \theta \in D_M \subset \mathbb{R}^m\}$$

$$n = 1, 2, 3 \dots \quad (42)$$

where  $n$  denotes the number of basis function (9).

In the previous section we saw that the integrated mean square error is typically split into two terms the *variance term* and the *bias term*

$$V_2(\hat{g}_N^n, g_0) = V_2(\hat{g}_N^n, g_n^*) + V_2(g_n^*, g_0) \quad (43)$$

where, according to (37),

$$V_2(\hat{g}_N^n, g_n^*) \sim \frac{m}{N}. \quad (44)$$

The bias term, which is entirely deterministic, decreases with  $n$ . Thus, for a given family  $\{\mathcal{G}_n\}$  there will be an optimal  $n = n^*(N)$  that balances the variance and bias terms.

Notice that (44) is a very general expression that holds almost regardless of how the sequence  $\{\mathcal{G}_n\}$  is chosen. Thus, it is in principle only possible to influence the bias error. In order to have a small integrated mean square error it is therefore of profound importance to choose  $\{\mathcal{G}_n\}$  such that the bias is minimized. An interesting possibility is to let the choice of  $\{\mathcal{G}_n\}$  be data driven. This may not seem like an easy task but here wavelets have proven to be useful, see Section 9.

When the bias and the variance can be exactly quantified, the integrated mean square error can be minimized w.r.t.  $n$ . This gives the optimal model complexity  $n^*(N)$  as a function of  $N$ . However, often it is only possible to give the rate with which the bias decreases as a function of  $n$  and the rate with which the variance increases with  $n$ . Then it is only possible to obtain the rate with which  $n^*(N)$  increases with  $N$ . Another problem is that if  $g_0$  in reality belongs not to  $\mathcal{G}$  but to some other class of functions, the rate will not be optimal. These considerations has lead to the development of methods where the choice of  $n$  is based on the observations  $Z^N$ . Basically,  $n$  is chosen so large that there is no evidence in the data that  $g_0$  is more complex than the estimated model, but not larger than that. Then, as is shown in Guo and Ljung (1994), the bias and the variance are matched. These adaptive methods are discussed in Section 9.

To get an idea of upper bounds for the optimal rate of convergence consider a simple linear regression problem:  $g_0(\varphi) = \varphi_0^T \theta$ . The bias is

then zero and the results in Section 3 give that the minimax rate (26) is

$$\inf_{\hat{g}_N} V_2^{\mathcal{P}_0}(\hat{g}_N, g_0) = \mathcal{O}\left(\frac{1}{N}\right)$$

for  $\mathcal{P}_0$ , the class of distributions with unbounded support. Furthermore, this rate is obtained for the least-squares estimate. Thus one cannot in general expect any better rate than this unless the distributions have bounded support, Akçay *et al.* (1994).

## 5. NEURAL NETS

What is meant by the term neural nets depends on the author. Lately neural net has become a word of fashion and today almost all kinds of models can be found by the names neural network somewhere in the literature. Old types of models, known for decades by other names, have been converted to, or reinvented as neural nets. This makes it impossible to cover all types of neural networks and only what is called feedforward and recurrent will be considered, which are the networks most commonly used in system identification. Information about other neural network models can be found in any introductory book in this field, *e.g.*, Kung (1993); Rumelhart and McClelland (1986); Hertz *et al.* (1991).

In Hunt *et al.* (1992); Narendra and Parthasarathy (1990); Sontag (1993) alternative overviews of neural networks in system identification and control can be found. Also the books White and Sofge (1992); W.T. Miller *et al.* (1992) contain many interesting articles on this topic.

### 5.1. Feedforward Neural Nets

The step from the general function expansion (9) to what is called neural nets is not big. With the choice  $g_k(\varphi) = \alpha_k \sigma(\beta_k \varphi + \gamma_k)$  where  $\beta_k$  is a parameter vector of size  $\dim \varphi$ , and  $\gamma_k$  and  $\alpha_k$  are scalar parameters we obtain

$$g(\varphi) = \sum_{k=1}^n \alpha_k \sigma(\beta_k \varphi + \gamma_k) + \alpha_0 \quad (45)$$

where a mean level parameter  $\alpha_0$  has been added. This model is referred to as a *feedforward network* with one hidden layer and one output unit in the NN literature. In Figure 1 it is displayed in the common NN way. The basis functions, called *hidden units*, *nodes*, or *neurons*, are univariate which makes the NN to an expansion in simple functions. The specific choice of  $\sigma(\cdot)$

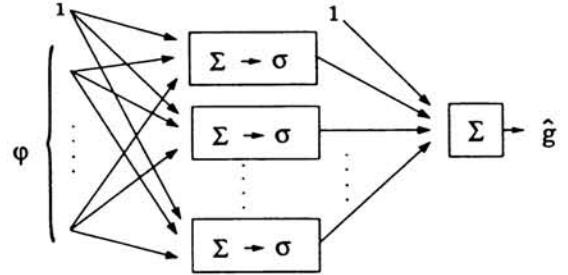


Fig. 1. Feedforward network with one hidden layer and one output unit. The arrows symbolize the parameters of the model.

is the *activation function* of the units which is usually chosen identically for all units.

The name *feedforward* is explained by the figure; there is a specific direction of flow in the computations when the output  $g$  is computed. First the weighted sums calculated at the input at each unit, then these sums pass the activation function and form the outputs of the hidden units. To form  $g$ , a weighted sum of the results from the hidden units is formed. This sum at the output is called the *output unit*. If  $g$  is vector function there are several output units forming an *output layer*. The input,  $\varphi$ , is sometimes called the *input layer*. The weights at the different sums are the parameters of the network.

In Cybenko (1989) it was shown that condition [R1], (7), holds if the activation function is chosen to be *sigmoidal* which is defined as

### Definition 5.1

Let  $\sigma(x)$  be continuous. Then  $\sigma(x)$  is called a *sigmoid function* if it has the following properties

$$\sigma(x) \rightarrow \begin{cases} a & \text{as } x \rightarrow +\infty \\ b & \text{as } x \rightarrow -\infty \end{cases} \quad (46)$$

where  $a, b, b < a$  are any real values.

The most common choice is

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad (47)$$

which gives a smooth, differentiable, model with the advantage that gradient based parameter estimate methods can be used, see Section 6. However, in Leshno *et al.* (1993) it is shown that (45) is a universal approximator, *i.e.*, [R1] holds for all non-polynomial  $\sigma(\cdot)$  which are continuous except at most in a set of measure zero.

The one hidden layer NN is related to the Projecting Pursuit (PP) model, see Section 9. In each unit a direction is estimated ( $\beta_k$ ) but, in

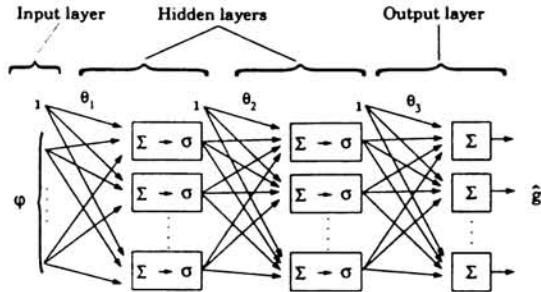


Fig. 2. Feedforward network with two hidden layers.

difference to PP, the function in this direction is fixed except for scaling and translation.

The one hidden layer network (45) can be generalized into a multi-layer network with several layers of hidden units. The outputs from the first hidden layer then feeds in to another hidden layer which feeds to the output layer - or another hidden layer. This is best shown with a picture; in Figure 2 such a net with two hidden layers and several outputs is shown. The formula for a NN with two hidden layer and one output becomes

$$g(\varphi) = \sum_i \theta_{1,i}^3 \sigma \left( \sum_j \theta_{i,j}^2 \sigma \left( \sum_m \theta_{r,m}^1 \varphi_m \right) \right) \quad (48)$$

The parameters have three indexes.  $\theta_{j,i}^M$ , is the parameter between the unit  $i$  in one layer and unit  $j$  in the following layer.  $M$  denotes which layer the parameter belongs to. The translation parameters corresponding to  $\gamma_k$  in (45) has not been written out.

At first, because of the general approximation ability of the NN with one hidden layer, there seems to be no reason to add more hidden layers. However, the rate of convergence might be very slow for some functions and it might be possible with a much faster convergence with two hidden layers (*i.e.*, condition [R2], (8) might favor two layers). Also, in Sontag (1990) it is shown that in certain control applications a two hidden layer NN can stabilize systems which cannot possibly be stabilized by NN with only one hidden layer.

### 5.2. Recurrent Neural Nets

If some of the inputs of a feedforward network consist of delayed outputs from the network, or some delayed internal state, then the network is called a *recurrent network*, or sometimes a dynamic network. In Figure 3 an example of a recurrent net with two past outputs fed back into the network.

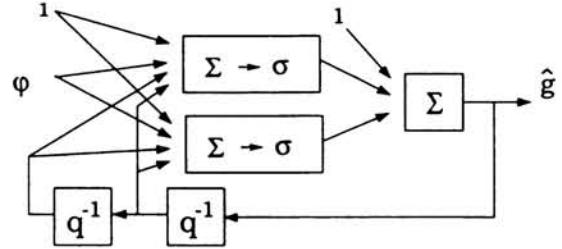


Fig. 3. Recurrent network.  $q^{-1}$  delays the signal one time sample.

The dynamic recurrent networks are especially interesting for identification, and in Section 11 two black-box models introduced based on recurrent networks.

Recurrent networks can also be used as a nonlinear state-space model. This is investigated in Matthews (1992).

### 6. ALGORITHMIC ASPECTS

In this section we shall discuss how to achieve the best fit between observed data and the model, *i.e.* how to carry out the minimization of (10).

$$V_N(\theta) = \frac{1}{2N} \sum_{t=1}^N |y(t) - g(\varphi(t), \theta)|^2 \quad (49)$$

No analytic solution to this problem is possible, so the minimization has to be done by some numerical search procedure. A classical treatment of the problem of how to minimize sum of squares is given in Dennis and Schnabel (1983). A survey of methods for the NN application is given in Kung (1993) and in van der Smagt (1994). Most efficient search routines are based on iterative local search in a “downhill” direction from the current point. We then have an iterative scheme of the following kind

$$\hat{\theta}^{(i+1)} = \hat{\theta}^{(i)} - \mu_i R_i^{-1} \nabla \hat{g}_i \quad (50)$$

Here  $\hat{\theta}^{(i)}$  is the parameter estimate after iteration number  $i$ . The search scheme is thus made up from the three entities

- $\mu_i$  step size
- $\nabla \hat{g}_i$  an estimate of the gradient  $V'_N(\hat{\theta}^{(i)})$
- $R_i$  a matrix that modifies the search direction

It is useful to distinguish between two different minimization situations

- (i) *Off-line or batch*: The update  $\mu_i R_i^{-1} \nabla \hat{g}_i$  is based on the whole available data record  $Z^N$ .
- (ii) *On-line or recursive*: The update is based only on data up to sample  $i$  ( $Z^i$ ), (typically done so that the gradient estimate  $\nabla \hat{g}_i$  is based only on data just before sample  $i$ .)

We shall discuss these two modes separately below. First some general aspects will be treated.

### 6.1. Search directions

The basis for the local search is the gradient

$$V'_N(\theta) = -\frac{1}{N} \sum_{t=1}^N (y(t) - g(\varphi(t), \theta)) \psi(\varphi(t), \theta) \quad (51)$$

where

$$\psi(\varphi(t), \theta) = \frac{\partial}{\partial \theta} g(\varphi(t), \theta) \quad (d \times 1 \text{-vector}) \quad (52)$$

It is well known that gradient search for the minimum is inefficient, especially close to the minimum. Then it is optimal to use the *Newton search direction*

$$R^{-1}(\theta) V'_N(\theta) \quad (53)$$

where

$$R(\theta) = V''_N(\theta) = \frac{1}{N} \sum_{t=1}^N \psi(\varphi(t), \theta) \psi^T(\varphi(t), \theta) + \frac{1}{N} \sum_{t=1}^N (y(t) - g(\varphi(t), \theta)) \frac{\partial^2}{\partial \theta^2} g(\varphi(t), \theta) \quad (54)$$

The true Newton direction will thus require that the second derivative

$$\frac{\partial^2}{\partial \theta^2} g(\varphi(t), \theta)$$

be computed. Also, far from the minimum,  $R(\theta)$  need not be positive semidefinite. Therefore alternative search directions are more common in practice:

- *Gradient direction*. Simply take

$$R_i = I \quad (55)$$

- *Gauss-Newton direction*. Use

$$R_i = H_i = \frac{1}{N} \sum_{t=1}^N \psi(\varphi(t), \hat{\theta}^{(i)}) \psi^T(\varphi(t), \hat{\theta}^{(i)}) \quad (56)$$

- *Levenberg-Maquard direction*. Use

$$R_i = H_i + \delta I \quad (57)$$

where  $H_i$  is defined by (56).

- *Conjugate gradient direction*. Construct the Newton direction from a sequence of gradient estimates. Loosely, think of  $V''_N$  as constructed by difference approximation of  $d$  gradients. The direction (53) is however constructed directly, without explicitly forming and inverting  $V''$ .

It is generally considered, Dennis and Schnabel (1983), that the Gauss-Newton search direction is to be preferred. For ill-conditioned problems the Levenberg-Maquard modification is recommended. However, good results with conjugate gradient methods have also been reported in NN applications (van der Smagt (1994)).

### 6.2. Back-Propagation: Calculation of the gradient

The only model-structure dependent quantity in the general scheme (50) is the gradient of the model structure (52). For a one-hidden-layer structure (45) this is entirely straightforward, since

$$\begin{aligned} \frac{d}{d\alpha} \alpha \sigma(\beta \varphi + \gamma) &= \sigma(\beta \varphi + \gamma) \\ \frac{d}{d\gamma} \alpha \sigma(\beta \varphi + \gamma) &= \alpha \sigma'(\beta \varphi + \gamma) \\ \frac{d}{d\beta} \alpha \sigma(\beta \varphi + \gamma) &= \alpha \sigma'(\beta \varphi + \gamma) \varphi \end{aligned}$$

For multi-layer NNs the gradient is calculated by the well known *Back-Propagation* (BP) method which can be described as the chain rule for differentiation applied to the expression (48). It also makes sure to re-use intermediate results which are needed at several places in the algorithm. Actually, the only complicated with the algorithm is to keep track of all indexes.

Backpropagation has been “rediscovered” several times, see e.g., Werbos (1974); Rumelhart *et al.* (1986).

Here the algorithm will be derived for the case where the network has two hidden layers and one output unit. For multi output models and with less- or more hidden layers only minor changes have to be done.

Consider the NN model (48). Denote by  $x_b^M$  and  $f_b^M$  the result at unit  $b$  in layer  $M$  before and after the activation function, respectively. That is

$$f_b^M = \sigma(x_b^M)$$

We can then write  $g(\varphi) = x_1^3 = \sum_i \theta_{1,i}^3 f_i^2$  and the derivative with respect to one of the parameters in the output layer becomes

$$\psi(\varphi)_{3,a,b} = \frac{\partial g}{\partial \theta_{1,b}^3} = f_b^2$$

In the same way  $x_a^2 = \sum_m \theta_{a,m}^2 f_m^1$  and the derivative of a parameter in the middle layer becomes

$$\psi(\varphi)_{2,a,b} = \frac{\partial g}{\partial \theta_{a,b}^2} = \delta_a^2 f_b^1$$

where

$$\delta_a^2 = \theta_{1,a}^3 \sigma'(x_a^2)$$

For the first layer we can write  $x_a^1 = \sum_m \theta_{a,m}^1 \varphi_m$  and the derivative of a parameter in this layer becomes

$$\psi(\varphi)_{1,a,b} = \frac{\partial g}{\partial \theta_{a,b}^1} = \delta_a^1 \varphi_b$$

where

$$\delta_a^1 = \sum_j \delta_j^2 \theta_{j,a}^2 \sigma'(x_a^1)$$

The nice feature is that  $\{f_b^M\}$  and  $\{x_b^M\}$  are obtained as intermediate results when  $g(\varphi)$  is calculated (forward propagation in the network). Calculating  $\{\delta_a^M\}$  can be viewed as propagating  $g(\varphi)$  backwards through the net, and this is the origin of the name of the algorithm.

The calculations are further simplified by the relation of the derivative of the sigmoid which follows from (47).

$$\sigma'(\cdot) = \sigma(\cdot)(1 - \sigma(\cdot)) \quad (58)$$

### 6.3. Implicit Regularization

Recall the discussion about regularization in Section 3. We pointed out that the parameter  $\delta$  in (35) acts like a knob that affects the “efficient number of parameters used”. It thus plays a similar role as the model size:

- Large  $\delta$ : small model structure, small variance, large bias
- Small  $\delta$ : large model structure, large variance, small bias

It is quite important for NN applications to realize that there is a direct link between the iterative process (50) and regularization in the sense that *aborting the iterations before the minimum has been found, has a quite similar effect as regularization*. This was noted in Wahba (1987) and pointed out as the cause of “overtraining” in Sjöberg and Ljung (1992). More precisely,

the link is as follows (when quadratic approximations are applicable)

$$(I - \mu R^{-1} V'')^i \sim \delta (\delta I + V'')^{-1}$$

so, as the iteration number increases, this corresponds to a regularization parameter that decreases to zero as

$$\log \delta \sim -i \quad (59)$$

How to know when to stop the iterations? As  $i \rightarrow \infty$  the value of the criterion  $V_N$  will of course continue to decrease, but as a certain point the corresponding regularization parameter becomes so small that increased variance starts to dominate over decreased bias. This should be visible when the model is tested on a fresh set – the *Validation data* (often called *generalization data* in the NN context). We thus evaluate the criterion function on this fresh data set, and plot the fit as a function of the iteration number. A typical such plot is shown in Figure 6. The point where the fit starts to be worse for the validation data is the iteration number (the degree of regularization or the effective model flexibility) where we are likely to strike the optimal balance between bias and variance error. Experience with NN applications has shown that this often is a very good way of limiting the actual model flexibility by effectively eliminating spurious parameters, *i.e.*, dealing with requirement [R3], mentioned in Section 1.

### 6.4. Off-line and on-line algorithms

The expressions (51) and (54) for the Gauss-Newton search clearly assume that the whole data set  $Z^N$  is available during the iterations. If the application is of an off-line character, *i.e.*, the model  $\hat{g}_N$  is not required during the data acquisition, this is also the most natural approach.

However, in the NN context there has been a considerable interest in on-line (or recursive) algorithms, where the data are processed as they are measured. Such algorithms are in NN contexts often also used in off-line situations. Then the measured data record is concatenated with itself several times to create a (very) long record that is fed into the on-line algorithm. We may refer to Ljung and Söderström (1983) as a general reference for recursive parameter estimation algorithm. In Solbrand *et al.* (1985) the use of such algorithms is the off-line case is discussed.

It is natural to consider the following algorithm as the basic one:

$$\hat{\theta}(t) = \hat{\theta}(t-1) + \mu_t R_t^{-1} \psi(\varphi(t), \hat{\theta}(t-1)) \varepsilon(t, \hat{\theta}(t-1)) \quad (60)$$

$$\varepsilon(t, \theta) = y(t) - g(\varphi(t), \theta) \quad (61)$$

$$R_t = R_{t-1} +$$

$$\mu_t [\psi(\varphi(t), \hat{\theta}(t-1))\psi^T(\varphi(t), \hat{\theta}(t-1)) - R_{t-1}] \quad (62)$$

The reason is that if  $g(\varphi(t), \theta)$  is linear in  $\theta$ , then (60) – (62), with  $\mu_t = 1/t$ , provides the analytical solution to the minimization problem (49). This also means that this is a natural algorithm close to the minimum, where a second order expansion of the criterion is a good approximation. In fact, it is shown in Ljung and Söderström (1983), that (60) – (62) in general gives an estimate  $\hat{\theta}(t)$  with the same (“optimal”) statistical, asymptotic properties as the true minimum to (49).

In the NN literature, often some averaged variants of (60) – (62) are discussed:

$$\hat{\theta}(t) = \hat{\theta}(t-1) + \mu_t R_t^{-1} \nabla \hat{g}_t \quad (63)$$

$$\nabla \hat{g}_t = \nabla \hat{g}_{t-1} +$$

$$\gamma_t [\psi(\varphi(t), \hat{\theta}(t-1))\varepsilon(t, \hat{\theta}(t-1)) - \nabla \hat{g}_{t-1}] \quad (64)$$

$$\hat{\theta}(t) = \hat{\theta}(t-1) + \rho_t [\hat{\theta}(t) - \hat{\theta}(t-1)] \quad (65)$$

The basic algorithm (60) – (62) then corresponds to  $\gamma_t = \rho_t = 1$ . Now, when do the different averages accomplish?

Let us first discuss (65) (and take  $\gamma_t \equiv 1$ ). This is what has been called “accelerated converges”. It was introduced by Polyak and Juditsky (1992) and has been extensively discussed by Kushner and others. The remarkable thing with this averaging is that we achieve the same asymptotic statistical properties of  $\hat{\theta}(t)$  by (63) – (65) with  $R_t = I$  (gradient search) as by (60) – (62) if

$$\gamma_1 = 1$$

$$\rho_t = 1/t$$

$$\mu_t >> \rho_t \quad \mu_t \rightarrow 0$$

It is thus an interesting alternative to (60) – (62), in particular if  $\dim \theta$  is large so  $R_t$  is a big matrix.

We now turn to the averaging in (64). For  $\gamma < 1$  this gives what is known as a *momentum* term. Despite its frequent use in NN applications, it is more debatable. An immediate argument for (64) is that the averaging makes the gradient  $\nabla \hat{g}_t$  more reliable (less noisy) so that we can take larger steps in (63). It is, however, immediate to verify that exactly the same averaging takes place in (63) if smaller steps are taken. A second argument is that (64) lends a “momentum” effect to the gradient estimate  $\nabla \hat{g}_t$ . That is, due to the low pass filter,  $\nabla \hat{g}_t$  will reflect not

only the gradient at  $\hat{\theta}(t-1)$ , but also at several previous values of  $\hat{\theta}(k)$ . This means that the update push in (63) will not stop immediately at value  $\theta$  where the gradient is zero. This could of course help to push away from a non-global, local minimum, which is claimed to be a useful feature. However, there seems to be no systematic investigation of whether this possible advantage is counter balanced by the fact that more iterations will be necessary for convergence.

### 6.5. Local Minima

A fundamental problem with minimization tasks like (49) is that  $V_N(\theta)$  may have several or many local (non-global) minima, where local search algorithms may get caught. There is no easy solution to this problem. It is usually well used effort to spend some time to come up with a good initial value  $\theta^{(0)}$  where to start the iterations. Other than that, only various global search strategies are left, such as random search, random restarts, simulated annealing, the genetic algorithm and whathaveyou.

## 7. INTERMISSION

We have so far described the “What” of Neural Networks: They are black-box, non-linear model structures, formed from certain basis functions. The estimation of the adjustable parameters in these structures follows traditional techniques, both in terms of algorithms and statistical properties. We have not yet discussed the far more difficult “Why” question: Are these structures to be preferred over other non-linear black box ones, and if so, why? To get some insight into this question we shall in the next two sections open up the perspective to general approximation issues, and then try to give a partial answer to the “Why” in Section 10.

We also still have to discuss how to apply the NN structures to System Identification problems. This will be dealt with in Section 11.

## 8. NON-ADAPTIVE METHODS

In this section we shall discuss the case when the sequence of approximation classes  $\{\mathcal{G}_n\}$  is chosen *a priori*. This leads to linear estimates, i.e. estimates that are linear in  $\{y(t)\}$ .

### 8.1. Basis function expansion

If it is assumed that  $g_0$  can be written as a function expansion

$$g_0(\varphi) = \sum_{k=1}^{\infty} \theta(k) b_k(\varphi),$$

where  $\{b_k\}$  forms a basis in  $\mathcal{G}$ , it is natural to use a truncated expansion

$$\hat{g}_N^n(\varphi, \hat{\theta}) = \sum_{k=1}^n \hat{\theta}_N(k) b_k(\varphi)$$

as estimate. Notice that this is a linear regression. This means that

$$\mathcal{G}_n = \text{span}\{b_k; k = 1, \dots, n\}$$

Usually, the parameters  $\{\hat{\theta}_N(k)\}$  are obtained by minimizing the quadratic norm (30). This gives the least-squares estimate

$$\hat{g}_N^n(\varphi) = \frac{1}{N} \sum_{t=1}^N b^T(\varphi) \left[ \frac{1}{N} \sum_{t=1}^N b(\varphi(t)) b^T(\varphi(t)) \right]^{-1} \times b(\varphi(t)) y(t) \quad (66)$$

where  $b(\varphi) = [b_1(\varphi), \dots, b_n(\varphi)]^T$ .

When the basis functions are orthonormal (w.r.t. the probability measure of  $\varphi$ ) the estimate is reduced to

$$\hat{g}_N^n(\varphi) = \sum_{k=1}^n \frac{1}{N} \sum_{t=1}^N b_k(\varphi(t)) b(\varphi) y(t)$$

For this type of models it is  $n$ , the number of basis functions, that controls the complexity of the model.

### 8.2. Kernel methods

Kernel methods, Härdle (1990), are based on the notion of smoothness. At a point  $\varphi$ , the value of  $g$  should be close to the value of the observations close to this point. Thus as estimate of  $g_0(\varphi)$  one could take a weighted average of the observations around  $\varphi$

$$\hat{g}(\varphi) = \frac{\frac{1}{N} \sum_{t=1}^N y(t) K(\varphi, \varphi(t))}{\frac{1}{N} \sum_{t=1}^N K(\varphi, \varphi(t))} \quad (67)$$

Here the kernel  $K$  is typically a function which has  $K(\varphi, \varphi)$  as its maximum and furthermore  $K(\varphi, \varphi')$  is decreasing with  $|\varphi - \varphi'|$ .

Consider now a family of kernels  $K_h$ ,  $h \in \mathbb{R}_+$ . The bandwidth  $h$  controls the size of the region around  $\varphi$  for which  $K$  is large. A smaller  $h$  corresponds to a more spiky kernel. Thus rougher

functions can be described with kernels using a smaller  $h$ . Hence, a class of kernels  $\{K_h\}$  with  $h$  decreasing corresponds to an increasing scale of spaces  $\{\mathcal{G}_{n(h)}\}$ . We denote the corresponding estimates by  $\hat{g}_N^h$ . A typical example of a family of kernels is

$$K_h(\varphi, \varphi') = \frac{1}{h} K \left( \frac{|\varphi - \varphi'|}{h} \right)$$

where  $K$  is a decreasing real-valued function. This type of kernel corresponds to estimator known as the Nadarya-Watson estimator, Nadarya (1964) and Watson (1964). So called *k-Nearest-Neighbor* methods, Yakowitz (1987), where the  $k$  nearest neighbors are used to compute the weighted average, also falls into this class. This procedure can be interpreted as a kernel estimate with a variable bandwidth.

Let us now compute the variance contribution to the mean-square error. Assume that  $\varphi(t)$  is uniformly distributed in the  $\varphi$ -space which we assume is a hypercube in  $\mathbb{R}^d$ . Let the kernel be normalized so that  $0 \leq h \leq 1$  represents the proportion of the  $\varphi$ -space for which the kernel is large. This means that  $\hat{g}_N^h(\varphi)$  is essentially the average over  $Nh$  data. Under suitable mixing conditions of the noise, the variance error will thus be of size  $\frac{C}{Nh}$  where  $C$  depends on the properties of the noise. Thus  $1/h$  plays the role of  $n$  for kernel methods.

Comparing (66) and (67) we see that a basis function expansion can be interpreted as a kernel method. Both are linear in the data.

### 8.3. Local linear regression

Closely related to kernel methods is local linear regression methods, Fan (1992). In local linear regression the least-squares problem

$$\frac{1}{N} \sum_{i=1}^N (y(t) - g^h(\varphi) - \theta^T(\varphi)(\varphi(t) - \varphi))^2 \times K\left(\frac{\varphi - \varphi(t)}{h}\right) \quad (68)$$

is solved w.r.t.  $g^h(\varphi)$  and  $\theta(\varphi)$  giving  $\hat{g}_N^h(\varphi)$  as estimate of  $g_0(\varphi)$ . The local character of this estimate is due to the local character of the kernel  $K$ . The method can be seen as a first order Taylor expansion of  $g_0$ .

### 8.4. Convergence results

For  $\mathcal{G}$  equal to the space of functions where the  $s-1$ th derivative belongs to a Hölder Ball with  $\alpha = 1$  (Lipschitz) then the bias for certain kernels can be shown to be proportional to

$h^{s/d}$  (recall that  $d$  is the dimension of  $\varphi$ ), see Härdle (1990). Using this and (44) to minimize the integrated mean square error gives the optimal bandwidth  $h(N) = N^{-\frac{d}{2s+d}}$  which gives an IMSE of order  $N^{-\frac{2s}{2s+d}}$ . This has been shown to be the minimax rate Stone (1982). Kernel methods applied to time-series data have been analyzed in Collomb and Härdle (1986), Vieu (1991b), Roussas (1990) and Truong and Stone (1992). An analysis of Nearest neighbor methods for time-series can be found in Yakowitz (1987).

For the Sobolev space  $\mathcal{W}_2^s(L)$ , Cencov (1982) has shown that a Fourier series expansion gives a IMSE of order  $N^{-\frac{2s}{2s+d}}$ . The optimal model order is of order  $N^{\frac{1}{2s+d}}$ . See Benveniste *et al.* (1994) for further details.

Notice that, for Sobolev Balls, a larger  $s$  corresponds to a smaller class of functions and that this corresponds to a faster convergence rate. This is of course natural but it illustrates a basic problem with these methods. Since one wants to make sure that the true system is contained within  $\mathcal{G}$ , this class  $\mathcal{G}$  will generally be chosen wider than necessary and as a consequence the convergence rate will be slower than necessary. However, observe the interesting fact that the Fourier expansion forms a basis for the whole scale of Sobolev spaces  $\mathcal{W}_2^s(L)$ ,  $s > 0$ . Thus, it is only the optimal number of basis functions  $n(N)$  that depends on which function class the true function belongs to. In the next section we shall discuss how  $n(N)$  can be selected using the observations  $Z^N$ .

### 8.5. Global vs local basis functions

One can classify the methods according to the support that the basis functions have:

1. *Global basis functions* whose support covers a large part of the space.
2. *Local basis functions* whose support only covers a small region in the space

The support of the basis functions determine the characteristics of the function expansion. Conventional expansions such as Fourier series are global. With such bases, properties of the estimated function at a point are inferred to distant points. This is advantageous if the model class does contain the true function  $g_0$ . The variance of the estimate in the region with sparse measurements is reduced since it “borrows” measurements from other regions. However, if the model class is inappropriate, extrapolation may give completely wrong estimates.

Kernel methods, on the other hand, give weighted local averages in the  $\varphi$  space. Thus, they provide a combination of smoothing and short-range interpolation. They can be viewed as local function expansions. Another example of local basis functions is so-called radial basis functions, see Chen and Billings (1992), which uses basis functions centered around the measurement points  $\{\varphi(t)\}$ . A local basis is suitable when the function is very spiky. Then a small bandwidth still gives a good local fit. However, a small bandwidth gives at the same time a high variance. This is especially a problem when the dimension of the  $\varphi$  vector is high.

It is also possible to combine global and local basis functions. Wavelet expansions is an example of this, see Benveniste *et al.* (1994).

## 9. ADAPTIVE METHODS

The use of data to select the basis functions characterize adaptive methods. The adaptation can be more or less sophisticated. In its simplest form, only the number of basis functions is selected. The merits and limitations of this procedure are explained in the first subsection while the second subsection deals with more advanced methods where also the basis functions themselves are adapted to the data.

### 9.1. Bandwidth and model order adaptation

For kernel methods (Section 8) it is the bandwidth  $h$  that should depend on  $\mathcal{G}$ . For a survey see Marron (1988). The basic idea is to estimate the IMSE (25) as a function of  $h$  and then select the minimizing  $h$ . The IMSE can be estimated in several ways. One method is cross-validation which means that the model is tested on a fresh data set and the empirical IMSE is computed. Instead of keeping a whole set of data for validation, the leave-one-out method, Stone (1974), can be used. Another possibility is so-called *Generalized Cross Validation* tests. These methods add a penalizing function to the quadratic prediction error criterion in the spirit of Akaike’s Information Criteria AIC, BIC.

The same ideas can be used for choosing the model order  $n$  in basis function expansion methods, Shibata (1980). See Benveniste *et al.* (1994) and Polyak and Tsybakov (1990) for further discussions of this.

There are however spaces where linear methods fail to attain the minimax rate. This was shown for certain Sobolev spaces in Nemirovskii *et al.* (1985) and Nemirovskii (1986) and for Besov and Triebel spaces, which allow for locally

non-smooth functions, in Donoho and Johnstone (1992). The basic problem for linear methods is that they use the same scale of resolution in the whole  $\varphi$ -space. If the function of interest is spatially inhomogeneous, these methods will either under-smooth the smooth part of the function or over-smooth the rough part. Thus, recently the interest has turned to non-linear methods which we shall discuss next.

### 9.2. Adaptive basis function expansion

Suppose that we have a set of basis functions  $\{b_k\}$  that span  $\mathcal{G}$ . Each set of  $n$  basis functions would generate a function class  $\mathcal{G}_n$  and a good idea would be to select these  $n$  basis function such that the approximation error is minimized among all possible choices of these sets of  $n$  basis functions. The problem of finding an  $n$ -dimensional subspace that minimizes the worst approximation error is known as Kolmogorov's  $n$ -width problem, Pinkus (1985). Depending on  $\mathcal{G}$ , the problem can be more or less complicated. For example, for the functions  $\sum_{k=0}^{\infty} a_k z^k$  that are analytic inside the disc of radius  $r \leq 1$  satisfying  $\sum_{k=0}^{\infty} |a_k|^2 r^{2k} < 1$ , the optimal subspace is given by  $\text{span}\{1, z, \dots, z^{n-1}\}$ .

*Wavelets.* For orthonormal basis functions, the basis functions that correspond to the largest coefficients in the expansion of  $g_0$  give the best approximation. Thus, an idea is to estimate a large number of coefficients and to select the  $n$  largest ones. It is interesting to note that with this procedure one get adaptation of the  $\{\mathcal{G}_n\}$  to the smoothness of  $\mathcal{G}$  for free; if the basis functions span several (or a scale of) spaces of functions, the approach will be optimal for all these spaces.

This approach has been exploited in the wavelet theory. Wavelet theory is based on orthonormal bases of  $L_2$  that also span a wide scale of function spaces with a varying degree of smoothness, Besov and Triebel spaces, Triebel (1983).

The basic problem with such a method is to determine which parameters are small and which are large, respectively. Donoho and Johnstone (1992) has shown that the use of *shrinkage* gives (near) minimax rates in these spaces. Shrinkage essentially means that a threshold is determined that depends on the number of data. Parameter estimates less than this threshold are set to zero. Often, for technical reasons, a soft threshold is used instead. In that case, every wavelet coefficient is "pulled" towards zero by a certain non-linear function. This is conceptually closely related to the *regularization* procedure outlined in Section 3. Then, parameters are attracted to-

wards the nominal value  $\theta^\#$ . However, so far explicit regularization does not seem to have been exploited in wavelet theory.

*Neural Networks.* Neural networks is an example of a structure where the basis functions appear more implicit. Consider the expression (45). This is an expansion with  $\{\sigma(\varphi, \eta_k)\}$  as basis functions. The fact that the  $\eta_k$ 's are estimated from data means that the basis functions are chosen adaptively. In other words, the basis functions are selected from data. Below we shall see that they have an important property when it comes to high-dimensional systems.

### 9.3. Adaptive kernel methods

As we illustrated in Section 8, basis function expansions correspond to kernel methods. This is still the case if the basis functions are selected in an adaptive fashion such as in wavelet theory. However, then the shape of the kernel and the bandwidth are locally chosen, Donoho and Johnstone (1992).

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

### 9.4. The "curse" of dimensionality

Almost all useful approximation theorems are asymptotic, i.e. they require the number of data to approach infinity,  $N \rightarrow \infty$ . In practical situations this cannot be done and it is of crucial importance how fast the convergence is. A general estimation of a function  $\mathbb{R}^d \rightarrow \mathbb{R}$  becomes slower in  $N$  when  $d$  is larger and in most practical situations it becomes impossible to do general estimation of functions for  $d$  larger than, say, 3 or 4. For higher dimensions the number of data required becomes so large that it is in most cases not realistic. This is the curse of dimensionality. This can be shown with the following example

#### Example 9.1

Approximate a function  $\mathbb{R}^d \rightarrow \mathbb{R}$  within the unit cube with the resolution 0.1. This requires that the distance between data is not larger than 0.1 in every direction, requiring  $N = 10^d$  data. This is hardly realistic for  $d > 4$ . When there are noisy measurements the demand of data increase further.  $\square$

### 9.5. Methods to avoid the “curse”

From the discussion in the preceding subsection it should be clear that general nonlinear estimation is not possible. Nevertheless, a number of methods have been developed to deal with the problems occurring for high-dimensional functions. The idea is to be able to estimate functions that in some sense have a low-dimensional character. *Projection pursuit regression*, Friedman and Stuetzel (1981), uses an approximation of the form

$$\hat{g}(\varphi) = \sum_{k=1}^n g_k (\varphi^T \theta(k))$$

where the  $g_k$ s are smooth univariate functions. The method thus expands the function in  $n$  different directions. These directions are selected to be the most important ones and, for each of these, the functions  $g_k$  are optimized. Thus, it is a joint optimization over the directions  $\{\theta(k)\}$  and the functions  $g_k$ . The claim is that for small  $n$  a wide class of functions can be well approximated by this expansion, Donoho and Johnstone (1989). The claim is supported by the fact that any smooth function in  $d$  variables can be written in this way, Diaconis and Shahshahani (1984). It is supposed to be useful for moderate dimensions,  $d < 20$ .

*Projection pursuit regression* is closely related to *neural networks* where the same function, any sigmoid function  $\sigma$  satisfying Definition 5.1, is used in all directions. The effectiveness of such methods has been illustrated in Barron (1993): Consider the class of functions  $\{g\}$  on  $\mathbb{R}^d$  for which there is a Fourier representation  $\tilde{g}$  which satisfies

$$C_f = \int |\omega| |\tilde{g}(\omega)| d\omega < \infty.$$

Then there is a linear combination of sigmoidal functions such that

$$\int_{B_r} |g(x) - g_n(x)|^2 dx \leq \frac{(2rC_f)^2}{n} \quad (69)$$

where  $B_r$  is a ball with radius  $r$ . The important thing to notice here is that the degree of approximation as a function of  $n$  does not depend on the dimension  $d$ .

This work originated with the result in Jones (1992) where sinusoidal functions were used to prove a similar result. The above result is not limited to sinusoidal and sigmoidal functions and the same idea has been applied to projection pursuit regression, Zhao (1992), *hinging hyperplanes*, Breiman (1993), and *radial basis functions*, Girosi and Anzellotti (1992).

Notice, however, that the result is only an approximation result and a stochastic counterpart still awaits its proof.

Barron (1993) also showed that  $1/n^{(2/d)}$  is a lower bound for the minimax rate for linear methods. For large  $d$ , this rate is exceedingly slow compared with  $1/n$ . Thus, this is a serious disadvantage for the methods described in the previous section. In higher dimensional spaces, the convergence rate of linear method is much slower compared with certain non-linear methods.

## 10. SPECIFIC PROPERTIES OF NN STRUCTURES

So, what are the special features of the Neural Net structure that motivate the strong interest? Based on the discussion so far, we may point to the following list of properties:

- The NN expansion is a basis, even for just one hidden layer, i.e. Requirement [R1] is satisfied.
- The NN structure does extrapolation in certain, adaptively chosen, directions and is localized across these directions. Like Projection Pursuit it can thus handle larger regression vectors, if the data pattern  $[y(t), \varphi(t)]$  cluster along subspaces.
- The NN structure uses adaptive bases functions, whose shape and location are adjusted by the observed data.
- The approximation capability (Requirement [R2]) is good as manifested in (69).
- Regularization, implicit (stopped iterations) or explicit (penalty for parameter deviations, usually from zero) is a useful tool to effectively include only those basis functions that are essential for the approximation, without increasing the variance. (Requirement [R3]).
- In addition, NNs have certain advantages in implementation, both in hardware and software, due to the repetitive structure. The basis functions are built up from only one core function,  $\sigma$ . This also means that the structure is resilient to failures, since any node can play any other node's role, by adjusting its weights.

## 11. MODELS OF DYNAMICAL SYSTEMS BASED ON NEURAL NETWORKS

We are now ready to take the step from general “curve fitting” to system identification. The choice of a model structure for dynamical systems contains two questions

- 1. What variables, constructed from observed past data, should be chosen as regressors, i.e., as components of  $\varphi(t)$ ?
- 2. What non-linear mapping should  $\varphi(t)$  be subjected to, i.e., How many hidden layers in (45) should be used, and how many nodes should each layer have?

The second question is related to more general NN considerations, as discussed in Section 5. The first one is more specific for identification applications. To get some guidance about the choice of regressors  $\varphi$ , let us first review the linear case.

### 11.1. A Review of Linear Black Box Models

The simplest dynamical model is the Finite Impulse Response model (FIR):

$$y(t) = B(q)u(t) + e(t) = b_1 u(t-1) + \dots + b_n u(t-n) + e(t) \quad (70)$$

Here we have used  $q$  to denote the shift operator, so  $B(q)$  is a polynomial in  $q^{-1}$ . The corresponding predictor is  $\hat{y}(t|\theta) = B(q)u(t)$  is thus based on a regression vector

$$\varphi(t) = [u(t-1), u(t-2), \dots, u(t-n)]$$

As  $n$  tends to infinity we may describe the dynamics of all (“nice”) linear systems. However, the character of the noise term  $e(t)$  will not be modeled in this way.

A variant of the FIR model is the Output Error model (OE):

$$y(t) = \frac{B(q)}{F(q)}u(t) + e(t) \quad (71)$$

where

$$F(q) = 1 + f_1 q^{-1} + \dots + f_{n_f} q^{-n_f}$$

The predictor is

$$\hat{y}(t|\theta) = \frac{B(q)}{F(q)}u(t) \quad (72)$$

Also this predictor is based on past inputs only. It can be rewritten

$$\hat{y}(t|\theta) = b_1 u(t-1) + \dots + b_{n_b} u(t-n_b) -$$

$$- f_1 \hat{y}(t-1|\theta) - \dots - f_{n_f} \hat{y}(t-n_f|\theta) \quad (73)$$

It is thus based on the regression variables

$$[u(t-1), \dots, u(t-n_b), \hat{y}(t-1|\theta), \dots, \hat{y}(t-n_f|\theta)] \quad (74)$$

Note that these regressors are partly constructed from the data, using a current model. As  $n_b$  and  $n_f$  tend to infinity, also this model is capable of describing all reasonable linear dynamic systems, but not the character of the additive noise  $e(t)$ . The advantage of (71) over (70) is that fewer regressors are normally required to get a good approximation. The disadvantage is that the minimization over  $\theta$  becomes more complicated. Also, the stability of the predictor (72) depends on  $F(q)$ , and thus has to be monitored during the minimization.

A very common variant is the ARX model

$$A(q)y(t) = B(q)u(t) + e(t) \quad (75)$$

with the predictor

$$\begin{aligned} \hat{y}(t|\theta) = & -a_1 y(t-1) - \dots - a_{n_a} y(t-n_a) \\ & + b_1 u(t-1) + \dots + b_{n_b} u(t-n_b) \end{aligned} \quad (76)$$

thus using the regressors

$$[y(t-1), \dots, y(t-n_a), u(t-1), \dots, u(t-n_b)] \quad (77)$$

As shown, e.g., in Ljung and Wahlberg (1992) this structure is capable of describing all (reasonable) linear systems, including their noise characteristics, as  $n_a$  and  $n_b$  tend to infinity. The ARX model is thus a “complete” linear model from the black box perspective. The only disadvantage is that  $n_a$  and  $n_b$  may have to be chosen larger than the dynamics require, in order to accommodate the noise description. Therefore, a number of variants of (75) have been suggested, where the noise model is given “parameters of its own”. The best known of these is probably the ARMAX model

$$A(q)y(t) = B(q)u(t) + C(q)e(t) \quad (78)$$

Its predictor is given by

$$\begin{aligned} \hat{y}(t|\theta) = & (c_1 - a_1)y(t-1) + \dots \\ & + (c_n - a_n)y(t-n) \\ & + b_1 u(t-1) + \dots + b_{n_b} u(t-n) \\ & + c_1 \hat{y}(t-1|\theta) + \dots + c_{n_f} \hat{y}(t-n_f|\theta) \end{aligned} \quad (79)$$

It thus uses the regression vector (77) complemented with past predictors, just as in (74) (although the predictors are calculated in a different way). A large family of black box linear models is treated, e.g. in Ljung and Söderström (1983). It has the form

$$A(q)y(t) = \frac{B(q)}{F(q)}u(t) + \frac{C(q)}{D(q)}e(t) \quad (80)$$

The special case  $A(q) = 1$  gives the well known Box-Jenkins (BJ) model. The regressors used for the corresponding predictor are given e.g. by equation (3.114) in Ljung and Söderström (1983). These regressors are based on  $y(t - k)$ ,  $u(t - k)$ , the predicted outputs  $\hat{y}(t - k|\theta)$  using the current model, as well as the simulated outputs  $\hat{y}_u(t - k|\theta)$ , which are predicted outputs based on an output error model (71).

Let us repeat that from a black box perspective, most variants of (80) are equivalent, in the sense that they can be transformed into each other, at the expense of changing the orders of the polynomials. The ARX model ( $C=D=F=1$ ) covers it all. The rationale for the other variants is that we may come closer to the true system using fewer regressors.

### 11.2. Choice of Regressors for Neural Network Models

The discussion on linear systems clearly points to the possible regressors:

- Past Inputs  $u(t - k)$
- Past Measured Outputs  $y(t - k)$
- Past Predicted Outputs, using current model,  $\hat{y}(t - k|\theta)$
- Past Simulated Outputs, using past inputs only and current model,  $\hat{y}_u(t - k|\theta)$

A rational question to ask would be: Given that I am prepared to use  $m$  regressors (the size of the input layer is  $m$ ), how should I distribute these over the four possible choices? There is no easy and quantitative answer to this question, but we may point to the following general aspects:

- Including  $u(t - k)$  only, requires that the whole dynamic response time is covered by past inputs. That is, if the maximum response time to any change in the input is  $\Upsilon$ , and the sampling time is  $T$ , then the number of regressors should be  $\Upsilon/T$ . This could be a large number. On the other hand, models based on a finite number of past inputs cannot be unstable in simulation, which often is an advantage.

A variant of this approach is to form other regressors from  $u^t$ , e.g. by Laguerre filtering, (e.g Wahlberg (1991)). This retains the advantages of the FIR-approach, at the same time as making it possible to use fewer regressors. It does not seem to have been discussed in the NN-context yet.

- Adding  $y(t - k)$  to the list of regressors makes it possible to cover slow responses with fewer regressors. A disadvantage is that past outputs bring in past disturbances into the model. The model is thus given an additional task to also sort out noise properties. A model based on past outputs may also be unstable in simulation from input only. This is caused by the fact that the past measured outputs are then replaced by past model outputs.

- Bringing in past predicted or simulated outputs  $\hat{y}(t - k|\theta)$  typically increases the model flexibility, but also leads to non-trivial difficulties. For neural networks, using past outputs at the input layer gives *recurrent* networks. See Section 5. Two problems must be handled:

- It may lead to instability of the network, and since it is a non-linear model, this problem is not easy to monitor.
- The simulated/predicted output depends on  $\theta$ . In order to do updates in (50) in the true gradient direction, this dependence must be taken into account, which is not straightforward. If the dependence is neglected, convergence to local minima of the criterion function cannot be guaranteed.

The balance of this discussion is probably that the regressors (77) should be the first ones to test.

### 11.3. Neural Network Dynamic Models

Following the nomenclature for linear models it is natural to coin similar names for Neural Network models. This is well in line with, e.g. Chen *et al.* (1990); Chen and Billings (1992). We could thus distinguish between

- *NNFIR*-models, which use only  $u(t - k)$  as regressors
- *NNARX*-models, which use  $u(t - k)$  and  $y(t - k)$  as regressors
- *NNOE*-models, which use  $u(t - k)$  and  $\hat{y}_u(t - k|\theta)$
- *NNARMAX*-models, which use  $u(t - k)$ ,  $y(t - k)$  and  $\hat{y}(t - k|\theta)$
- *NNBJ*-models, which use all the four regressor types.

In Narendra and Parthasarathy (1990) another notation is used for the same models. The NNARX model is called Series-Parallel model and the NNOE is called Parallel model.

From a structural point of view, these black-box models are just slightly more troublesome to handle than their linear counterparts. When the regressor has been decided upon, it only remains to decide how many hidden units which should be used. The linear ARX model is entirely specified by three structural parameters  $[n_a \ n_b \ n_k]$ . [ $n_k$  is here the delay, which we have taken as 1 so far. In general we would work with the regressors  $u(t - n_k), \dots, u(t - n_k - n_b + 1)$ .] The NNARX model has just one index more,  $[n_a \ n_b \ n_k \ n_h]$ , where  $n_h$  is the number of units in the hidden layer which in some way corresponds to "how non-linear" the system is. The notation for NNOE and NNARMAX models follow the same simple rule.

If more than one hidden layer is used there will be one additional structural parameter for each layer.

It follows from Section 5.2 that NNOE, NNBJ, and NNARMAX correspond to recurrent neural nets because parts of the input to the net (the regressor) consist of past outputs from the net. As pointed out before, it is in general harder to work with recurrent nets. Among other things it becomes difficult to check under what conditions the obtained model is stable, and it takes an extra effort to calculate the correct gradients for the iterative search.

#### 11.4. Some Other Structural Questions

The actual way that the regressors are combined clearly reflect structural assumptions about the system. Let us, for example, consider the assumption that the system disturbances are additive, but not necessarily white noise:

$$y(t) = g(u^t) + v(t) \quad (81)$$

Here  $u^t$  denotes all past inputs, and  $v(t)$  is a disturbance, for which we only need a spectral description. It can thus be described by

$$v(t) = H(q)e(t)$$

for some white sequence  $\{e(t)\}$ . The predictor for (81) then is

$$\hat{y}(t) = (1 - H^{-1}(q))y(t) + H^{-1}(q)g(u^t) \quad (82)$$

In the last term, The filter  $H^{-1}$  can equally well be subsumed in the general mapping  $g(u^t)$ . The structure (81) thus leads to a NNFIR or NNOE

structure, complemented by a *linear* term containing past  $y$ .

In Narendra and Parthasarathy (1990) a related Neural Network based model is suggested. It can be described by

$$\hat{y}(t) = f(\theta_1, \varphi_1(t)) + g(\theta_2, \varphi_2(t)) \quad (83)$$

where  $\varphi_1(t)$  consists of delayed outputs and  $\varphi_2(t)$  of delayed inputs. The parameterized functions  $f$  and  $g$  can be chosen to be linear or non-linear by a neural net. A further motivation for this model is that it becomes easier to develop controllers from (83) than from the models discussed earlier.

In McAvoy (1992), it is suggested first to build a linear model for the system. The residuals from this model will then contain all unmodelled non-linear effects. The Neural Net model could then be applied to the residuals (treating inputs and residuals as input and output), to pick up the non-linearities. This is attractive, since the first step to obtain a linear model is robust and often leads to reasonable models. By the second Neural Net step, we are then assured to obtain at least as good a model as the linear one.

The question of how many layers to use is not easy. Sontag (1993) contains many useful and interesting insights into the importance of second hidden layers in the NN structure. See also the comments on this in Section 5.1.

#### 11.5. The Identification Procedure

A main principle in identification is the rule *try simple things first*. The idea is to start with the simplest model which has a possibility to describe the system and only to continue to more complex ones if the simple model does not pass validation tests.

When a new more complex model is investigated the results with the simpler model give some guidelines how the structural parameters should be chosen in the new model. It is e.g. common to start with an ARX model. The delay and number of delayed inputs and outputs give a good initial guess how the structure parameters should be chosen for the more complex ARMAX model. In this way less combinations of structural parameters have to be tested and computer time is saved.

Many non-linear systems can be described fairly well by linear models and for such systems it is a good idea to use insights from the best linear model how to select the regressors for the NN model. To begin with, only the number of hidden units needs to be varied. Also,

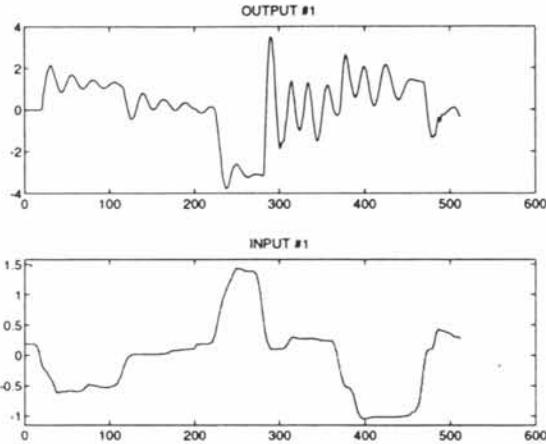


Fig. 4. Measured values of oil pressure (top) and valve position (bottom).

there might be more problems with local minima for the non-linear than for the linear models which makes it necessary to do several parameter estimates with different initial guesses. This further limits the number of candidate models which can be tested.

In the following example a hydraulic actuator is identified. First a linear model is proposed which does not capture all the fundamental dynamical behavior and then a NNARX model is tried. The same problem is considered in Benveniste *et al.* (1994) using wavelets as model structure.

#### Example 11.1 Modeling a Hydraulic Actuator.

The position of a robot arm is controlled by a hydraulic actuator. The oil pressure in the actuator is controlled by the size of the valve opening through which the oil flows into the actuator. The position of the robot arm is then a function of the oil pressure. In Gunnarsson and Krus (1990) a thorough description of this particular hydraulic system is given. Figure 4 shows measured values of the valve size and the oil pressure, which are input- and output signals, respectively. As seen in the oil pressure, we have a very oscillative settling period after a step change of the valve size. These oscillations are caused by mechanical resonances in the robot arm.

Following the principle "try simple things first" gives an ARX model with structural parameters  $[n_a \ n_b \ n_k] = [3 \ 2 \ 1]$ . In Figure 5 the result of a simulation with the obtained linear model on validation data is shown. The result is not very impressive.

Instead a NNARX model is considered with the same regressor as the linear model, i.e., with the

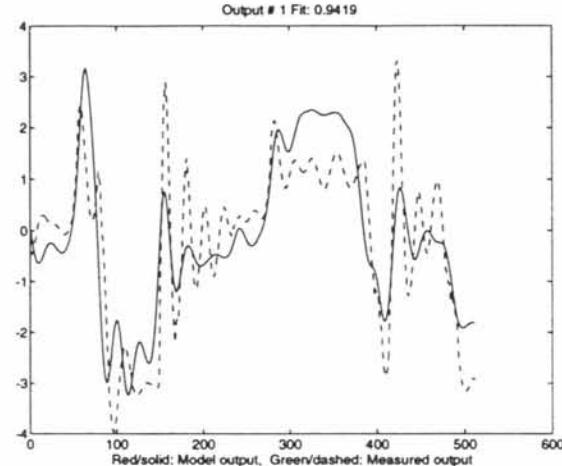


Fig. 5. Simulation of the linear model on validation data. Solid line: simulated signal. Dashed line: true oil pressure.

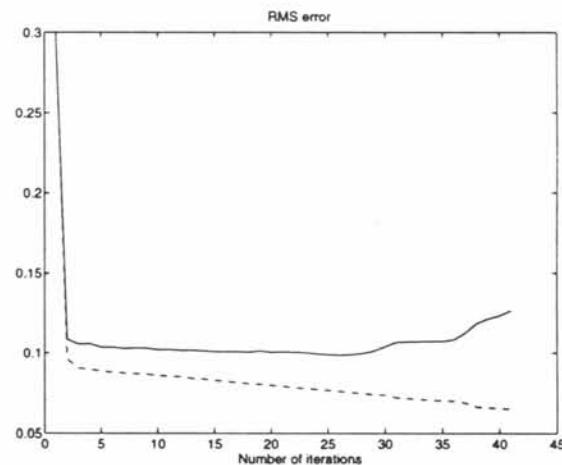


Fig. 6. Sum of squared error during the training of the NNARX model. Solid line: Validation data. Dashed line: estimation data.

same first three structural indexes, and with 10 hidden units,  $n_h = 10$ . In Figure 6 it is shown how the quadratic criterion develops during the estimation for estimation and validation data, respectively. For the validation data the criterion first decrease and then it starts to increase again. This is the overtraining which was described in Section 6.3. The best model is obtained at the minimum and this means that not all parameters in the non-linear model have converged and, hence, the "efficient number of parameters" is smaller than the dimension of  $\theta$ .

The parameters which give the minimum are then used in the non-linear model and in Figure 7 this NNARX model is used for simulation on the validation data.

This model performs much better than the linear model and it is compatible to the result ob-

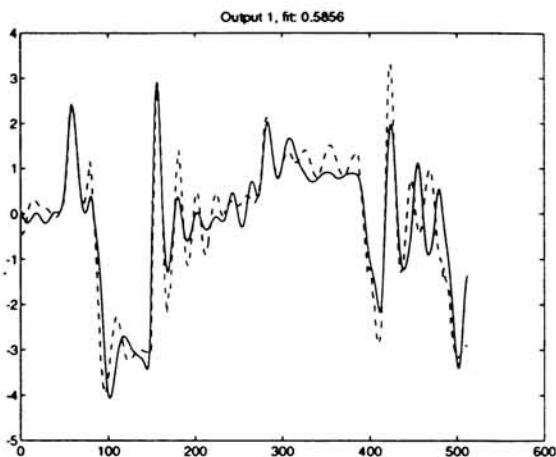


Fig. 7. Simulation of the non-linear model on validation data. Solid line: simulated signal. Dashed line: true oil pressure.

tained with a wavelet model in Benveniste *et al.* (1994).

□

## 12. USE OF NEURAL NET MODELS

In this section different uses of NN in system identification and some research topics will be discussed.

### 12.1. Modeling Controllers

If a good controller is available, which for some reason should replaced by a neural network, it can be used to produce training data. These are then used to estimate a NN model of the controller. It can, *e.g.*, be a human expert which one wants to duplicate or an optimal controller which one want to approximate with a feedback controller. See, *e.g.*, van Luenen (1993); McElvey (1992).

### 12.2. Modeling Inverse Systems

A popular approach in neurocontrol is to train the network to emulate the inverse of the system. The system is first used to produce training data for random choices of the input  $u(t)$ . The input- output signals are then interchanged and the network is trained to give the  $u(t)$  which produces the right  $y(t)$ . This gives a feedforward control scheme; after training the obtained neurocontroller is connected in cascade with the system. There are difficulties concerning the robustness with this approach and the inverse may not exist, see Sontag (1993).

In *e.g.*, Barto (1990), and the references given there, the idea of inverse control is further investigated.

### 12.3. Help for Linear Models

Neural nets can be used as a supporting system in situation where the plant is described by a linear model. It can be used to describe various non-linear relations, *e.g.*, in failure detection as in Kumamaru *et al.* (1994). In Chu and Shoureshi (1994) neural nets are used to identify the parameters of a linear model.

### 12.4. Model Choice

The choice of model set is important in linear identification and the issue becomes even more delicate in non-linear identification. Criteria for choosing the appropriate model set (polynomial, neural net, wavelet ...) would be of great help for the user. What makes one model a "good model" to a specific problem?

### 12.5. Initial Guess

For non-convex optimization there is always a problem with local minima of the criterion of fit. Several restarts from different initial guesses usually have to be done. Is there way to initialize the parameters in the network close to the optimum? This would speed up training and reduce the problem of local minima. To do this it might advantageous to make use of the fact that the data is generated of a dynamical system. Gotanda (1994) discusses this problem. One might also note that some of the parameters in the NN-structure ( $\alpha_k$  in (45)) enter linearly in the predictor. This should be utilized in the start-up procedure.

### 12.6. Estimation - or Training

Gauss-Newton type algorithms are typical the best way to minimize functions which are sum of squares. (compare the discussion in Section 6). Could one maybe make use of the sigmoid's very specific form to improve the algorithm? Very much has been written about different ways to speed up the training in different ways. See *e.g.*, Kung (1993) and the references there. Youlal and Kada (1994) also consider this topic.

### 12.7. Software

There are several free- and commercial software programs for neural networks available on the market. However, they are not primarily intended for system identification and the terminology used is usually from neural network community rather than from identification which make them less user-friendly. Surprisingly often the slow converging gradient-back-propagation method (see Section 6) is used instead of a faster Gauss-Newton type method. This has had the effect that many researchers write their own code.

### 13. REFERENCES

- Akçay, H., Hjalmarsson, H., and Ljung, L. (1994). On the choice of norms in system identification. In *10th IFAC Symposium on System Identification*.
- Barron, A. (1993). Universal approximation bounds for superpositions of a sigmoidal function. *IEEE Trans. Information Theory*, IT-39:930–945.
- Barto, A. (1990). Connectionist learning for control. In Miller, W., Sutton, R., and Werbos, P., editors, *Neural Networks for Control*, pages 5–58. MIT Press, Cambridge.
- Benveniste, A., Juditsky, A., Delyon, B., Zhang, Q., and Gorenne, P.-Y. (1994). Wavelets in identification. In *Preprint of the 10th IFAC Symposium on Identification*. (Copenhagen, 4-6 July).
- Breiman, L. (1993). Hinging hyperplanes for regression, classification and function approximation. *IEEE Trans. Information Theory*, IT-39:999–1013.
- Cencov, N. (1982). Statistical decision rules and optimal inference. *Amer. Math. Soc. Transl.*, 53.
- Chen, S. and Billings, S. (1992). Neural networks for nonlinear dynamic system modelling and identification. *Int. J. Control*, 56(2):319–346.
- Chen, S., Billings, S., and Grant, P. (1990). Non-linear system identification using neural networks. *Int. J. Control*, 51(6):1191–1214.
- Chu, S. and Shoureshi, R. (1994). Simultaneous parameter identification and states estimation using neural networks. In *Preprint, 10th IFAC Symposium on System Identification, Copenhagen*.
- Collomb, G. and Härdle, W. (1986). Strong uniform convergence rates in robust nonparametric time-series analysis and prediction: Kernel regression estimation from dependent observations. *Stochastic Processes and their Applications*, 23:77–89.
- Cybenko, G. (1989). Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals, and Systems*, 2:303–314.
- Dennis, J. and Schnabel, R. (1983). *Numerical Methods for Unconstrained Optimization and Nonlinear Equations*. Prentice-Hall, Englewood Cliffs, New Jersey.
- Diaconis, P. and Shahshahani, M. (1984). “On nonlinear functions of linear combinations”. *SIAM J. Sci. Statist. Comput.*, 5:175–191.
- 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. (1992). Minimax estimation via wavelet shrinkage. Technical report, Dept. of Statistics, Stanford University.
- Fan, J. (1992). Design-adaptive nonparametric regression. *Journal of the Amer. Stat. Assoc.*, 87:998–1004.
- Fan, J. and Gijbels, I. (1992). Variable bandwidth and local linear regression smoothers. *Annals of Statistics*, 20:2008–2036.
- Friedman, J. and Stuetzel, W. (1981). Projection pursuit regression. *J. Amer. Statist. Assoc.*, 76:817–823.
- Girosi, F. and Anzellotti, G. (1992). Convergence rates of approximation by translates. Technical report, Dept. Math. Art. Intell. Lab, M.I.T.
- Gotanda, H. (1994). Initialization of back propagation algorithms for multilayer neural networks. In *Preprint, 10th IFAC Symposium on System Identification, Copenhagen*.
- Gunnarsson, S. and Krus, P. (1990). Modelling of a flexible mechanical system containing hydraulic actuators. Technical report, Dep. of Electrical Engineering, Linköping University, S-581 83 Linköping, Sweden.
- Guo, L. and Ljung, L. (1994). The role of model validation for assessing the size of the unmodelled dynamics. Technical report, Report LiTH-ISY, Department of Electrical Engineering, Linköping University, Sweden.

- Härdle, W. (1990). *Applied Nonparametric Regression*. Cambridge University Press, Cambridge.
- Hertz, J., Krogh, A., and Palmer, R. (1991). *Introduction to the Theory of Neural Computation*. Addison-Wesley, 350 Bridge Parkway, Redwood City, CA 94065.
- Hjalmarsson, H. and Ljung, L. (1994). A unifying view of disturbances in identification. In *10th IFAC Symposium on System Identification*.
- Huber, P. J. (1981). *Robust Statistics*. Wiley.
- Hunt, K., Sbarbaro, D., Źbikowski, R., and Gawthrop, P. (1992). Neural networks for control systems - a survey. *Automatica*, 28(6):1083–1112.
- Jones, L. (1992). A simple lemma on greedy approximations in Hilbert space and convergence rates for projection pursuit regression and neural network training. *The Annals of Statistics*, 20:608–613.
- Kumamaru, K., Inoue, K., Nonaka, S., Ono, H., and Soderström, T. (1994). A neural network approach to failure decision of adaptively controlled systems. In *Preprint, 10th IFAC Symposium on System Identification, Copenhagen*.
- Kung, S. (1993). *Digital Neural Networks*. Prentice-Hall, Englewood Cliffs, New Jersey.
- Leshno, M., Lin, V., Pinkus, A., and Schocken, S. (1993). Multilayer feedforward networks with a nonpolynomial activation function can approximate any function. *Neural Networks*, 6:861–867.
- Ljung, L. (1987). *System Identification: Theory for the User*. Prentice-Hall, Englewood Cliffs, NJ.
- Ljung, L. and Söderström, T. (1983). *Theory and Practice of Recursive Identification*. MIT Press, Cambridge, Massachusetts.
- Ljung, L. and Wahlberg, B. (1992). Asymptotic properties of the least-squares method for estimating transferfunctions and disturbance spectra. *Adv. Appl. Prob.*, 24:412–440.
- Marron, J. (1988). Automatic smoothing parameter selection: a survey. *Empirical Economics*, 13:187–208.
- Matthews, M. B. (1992). *On the Uniform Approximation of Nonlinear Discrete-Time Fading-Memory Systems using Neural Network Models*. PhD thesis, ETH, Zürich, Switzerland.
- McAvoy, T. (1992). Personal communication.
- McKelvey, T. (1992). Neural networks applied to optimal control. In *Preprint IFAC/IFIP/IMAC Symposium on Artificial Intelligence in Real-Time Control*, pages 43–47.
- Milanese, M. and Vicino, A. (1991). Optimal estimation theory for dynamic systems with set membership uncertainty: An overview. *Automatica*, 27(6):997–1009.
- Nadarya, E. (1964). On estimating regression. *Theory Prob. Appl.*, 10:186–190.
- Narendra, K. and Parthasarathy, K. (1990). Identification and control of dynamical systems using neural networks. *IEEE Trans. Neural Networks*, 1:4–27.
- Nemirovskii, A. (1986). Nonparametric estimation of smooth regression functions. *J. Comput. Syst. Sci.*, 23:1–11.
- Nemirovskii, A., Polyak, B., and Tsybakov, A. (1985). Rate of convergence of nonparametric estimates of maximum-likelihood type. *Problems of Information Transmission*, 13:258–272.
- Pinkus, A. (1985). *n-Widths in Approximation Theory*. Springer-Verlag, Berlin.
- Polyak, B. and Juditsky, A. (1992). Acceleration of stochastic approximation by averaging. *SIAM J. Control*, 30(4):838–855.
- Polyak, B. and Tsybakov, A. (1990). Asymptotic optimality of  $c_p$  criterion for projection regression estimates. *Theory of Prob. and Appl.*, 35:305–317.
- Roussas, G. (1990). Nonparametric regression estimation under mixing conditions. *Stochastic Processes and their Applications*, 36:107–116.
- Rumelhart, D., Hinton, G., and Williams, R. (1986). Learning representations by back-propagating errors. *Nature*, 323(9):533–536.
- Rumelhart, D. and McClelland, J. (1986). *Parallel Distributed Processing: Explorations in the Microstructure of Cognition*, volume 1. MIT Press, Cambridge MA.
- Schweppen, F. (1973). *Uncertain Dynamic Systems*. Prentice-Hall, Englewood Cliffs, New Jersey.
- Shibata, R. (1980). Asymptotically efficient selection of the order of the model for estimating parameters of a linear process. *Ann. Statist.*, 8:147–164.

- Sjöberg, J. and Ljung, L. (1992). Overtraining, regularization, and searching for minimum in neural networks. In *Preprint IFAC Symposium on Adaptive Systems in Control and Signal Processing*, pages 669–674, Grenoble, France.
- Söderström, T. and Stoica, P. (1989). *System Identification*. Prentice-Hall International, Hemel Hempstead, Hertfordshire.
- Solbrand, G., Ahlén, A., and Ljung, L. (1985). Recursive methods for off-line identification. *Int. J. Control.*, 41:177–191.
- Sontag, E. (1990). Feedback stabilization using two-hidden-layer nets. Technical report, Report SYCON-90-11, Rutgers Center for Systems and Control, Dept. of Mathematics, Rutgers University, New Brunswick, NJ 08903.
- Sontag, E. (1993). Neural networks for control. In Trentelman, H. and Willems, J., editors, *Essays on Control: Perspectives in the Theory and its Applications*, volume 14 of *Progress in Systems and Control Theory*, pages 339–380. Birkhäuser.
- Stone, C. (1982). Optimal global rates of convergence for nonparametric regression. *Ann. Statist.*, 10:1040–1053.
- Stone, M. (1974). Cross-validatory choice and assessment of statistical predictions. *J.R. Statist. Soc. Ser. B*, 36:111–147.
- Triebel, H. (1983). *Theory of Function Spaces*. Birkhäuser Verlag, Basel.
- Truong, Y. and Stone, C. (1992). Nonparametric function estimation involving time series. *Ann. Statist.*, 20:77–97.
- van der Smagt, P. (1994). Minimisation methods for training feedforward neural networks. *Neural Networks*, 7(1):1–11.
- van Luenen, W. (1993). *Neural Networks for Control on Knowledge Representation and Learning*. PhD thesis, University of Twente, The Netherlands.
- Vieu, P. (1991a). Nonparametric regression: Optimal local bandwidth choice. *J.R. Statist. Soc. Ser. B*, 53:453–464.
- Vieu, P. (1991b). Quadratic errors for nonparametric estimates under dependence. *J. Multiv. Anal.*, 39:3227–347.
- Wahba, G. (1987). Three topics in ill-posed problems. In Engl, H. and Groetsch, C., editors, *Inverse and Ill-posed Problems*. Academic Press.
- Wahba, G. (1990). *Spline Models for Observational Data*. SIAM, 3600 University City Science Center, Philadelphia, PA 19104-2688.
- Wahlberg, B. (1991). System identification using laguerre models. *IEEE Trans. AC*, 36(5):551–562.
- Watson, G. (1964). Smooth regression analysis. *Sankya Series A*, 26:359–372.
- Werbos, P. (1974). *Beyond Regression: New Tools for Prediction and Analysis in the Behavioral Science*. PhD thesis, Harvard University.
- White, D. and Sofge, D., editors (1992). *Handbook of Intelligent Control, Neural, Fuzzy, and Adaptive Approaches*. Multiscience Press, Inc., Van Nostrand Reinhold, 115 Fifth Avenue, New York, NY 10003.
- W.T. Miller, I., Sutton, R., and Werbos, P., editors (1992). *Neural Networks for Control*. Neural Network Modeling and Connectionism. MIT Press. Series editor: J.L. Elman.
- Yakowitz, S. (1987). Nearest neighbor methods for time series analysis. *Journal of Time Series Analysis*, 18:1–13.
- Youla, H. and Kada, A. (1994). A class of recurrent neural networks for adaptive control of nonlinear systems. In *Preprint, 10th IFAC Symposium on System Identification, Copenhagen*.
- Zhao, Y. (1992). *On projection pursuit learning*. PhD thesis, Dept. Math. Art. Intell. Lab. M.I.T.