

**Figure 1.** Results from test flights of the new Swedish aircraft JAS-Gripen, developed by SAAB Military Aircraft AB, Sweden. From above (a) Pitch rate, (b) Elevator angle, (c) Canard angle, (d) Leading edge flap.

## SYSTEM IDENTIFICATION

The process of going from observed data to a mathematical model is fundamental in science and engineering. In the control area, this process has been termed *system identification*, and the objective is then to find dynamical models (difference or differential equations) from observed input and output signals. Its basic features are, however, common with general model building processes in statistics and other sciences.

System identification covers the problem of building models of systems when insignificant prior information is available and when the system's properties are known, up to a few parameters (physical constants). Accordingly, one talks about *black-box* and *gray-box* models. Among black-box models, there are familiar linear models such as ARX and ARMAX, and among nonlinear black-box models we have, for example, artificial neural networks (ANN).

### THE PROBLEM

The area of system identification begins and ends with real data. Data are required to build and to validate models. The result of the modeling process can be no better than what corresponds to the information contents in the data.

Look at two data sets:

#### **Example 1 An Unstable Aircraft**

Figure 1 shows some results from test flights of the new Swedish aircraft, JAS-Gripen, developed by SAAB Military Aircraft AB, Sweden. The problem is to use the information in these data to determine the dynamical properties of the aircraft for fine-tuning regulators, for simulations, and so on. Of particular interest are the aerodynamical derivatives.

#### **Example 2 Vessel Dynamics**

Figure 2 shows data from a pulp factory. They are collected from one of the buffer vessels. The problem is to determine the residence time in the vessel. The pulp spends about 48 h total in the process, and knowing the residence time in the different vessels is important in order to associate various portions of the pulp with the different chemical actions that have taken place in the vessel at different times. (The  $\kappa$ -num-

ber is a quality property that in this context can be seen as a marker allowing us to trace the pulp.)

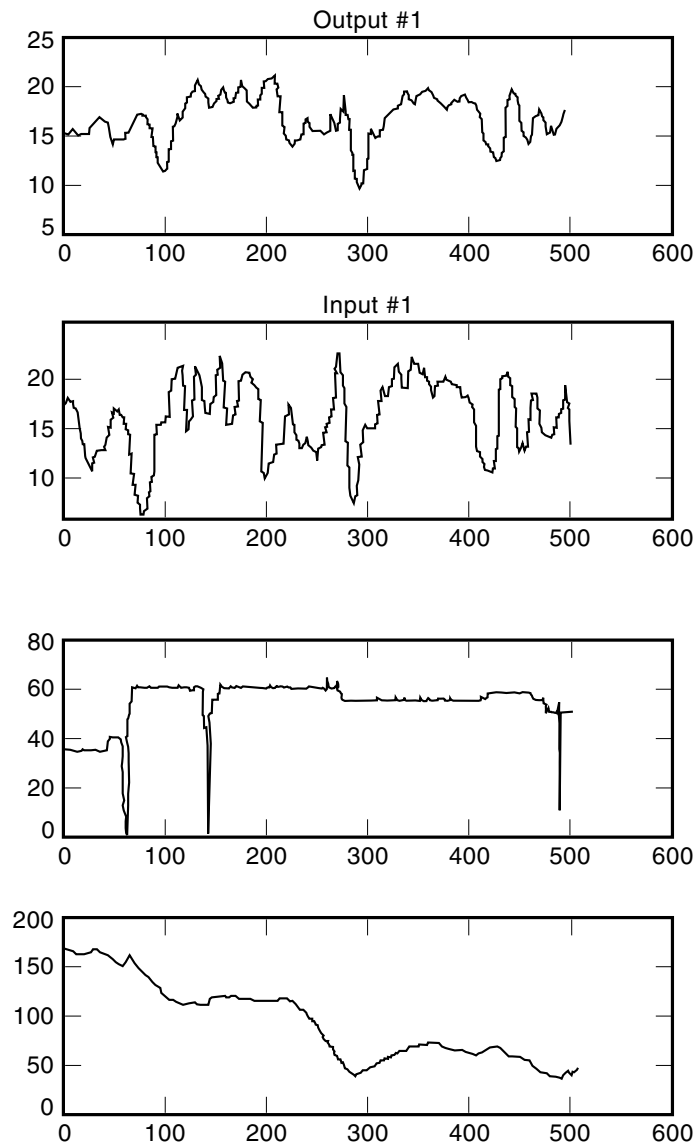
So, the bottom line of these examples is that we have collected input–output data from a process or a plant, and we need to extract information from these to find out (something about) the process' dynamical properties.

## BACKGROUND AND LITERATURE

System identification has its roots in standard statistical techniques, and many of the basic routines have direct interpretations as well-known statistical methods such as least squares and maximum likelihood. The control community took an active part in the development and application of these basic techniques to dynamic systems right after the birth of modern control theory in the early 1960s. Maximum likelihood estimation was applied to different equations (ARMAX models) by Ref. 1, and thereafter, a wide range of estimation techniques and model parameterizations flourished. By now, the area is well matured with established and well understood techniques. Industrial use and application of the techniques has become standard. See Ref. 2 for a common software package.

The literature on system identification is extensive. For a practical user oriented introduction, we may mention (3). Texts that go deeper into the theory and algorithms include Refs. 4 and 5. A classical treatment is Ref. 6.

These books all deal with the “mainstream” approach to system identification, as described in this article. In addition, there is substantial literature on other approaches, such as “set membership” (compute all those models that reproduce the observed data within a certain given error bound), estimation of models from given frequency response measurement (7), on-line model estimation (8), non-parametric frequency domain methods (9), etc. To follow the development in the



**Figure 2.** From the pulp factory at Skutskär, Sweden. The pulp flows continuously through the plant via several buffer tanks. From above: (a) the  $\kappa$ -number of the pulp flowing into a buffer vessel, (b) the  $\kappa$ -number of the pulp coming out from the buffer vessel, (c) flow out from the buffer vessel, and (d) level in the buffer vessel.

field, the IFAC series of Symposia on System Identification [Budapest, Hungary (1991), Copenhagen, Denmark (1994), Fukuoka, Japan (1997)] is also a good source.

## OUTLINE

The system identification procedure is characterized by four basic ingredients:

1. The observed data
2. A set of candidate models
3. A criterion of fit
4. Validation

The problem can be expressed as finding that model in the candidate set, that best describes the data, according to the criterion, and then evaluate and validate that model's properties. To do this we need to penetrate a number of things:

1. First, we give a preview of the whole process, as applied to the simplest set of candidate models.
2. Then, at some length, we display and discuss the most common sets of candidate models used in system identification. In general terms, a model will be a predictor of the next output  $y(t)$  from the process, given past observations  $Z^{t-1}$ , and parameterized in terms of a finite-dimensional parameter vector  $\theta$ :

$$\hat{y}(t|\theta) = g(\theta, Z^{t-1}) \quad (1)$$

3. We then discuss the criterion of fit for general model sets. This will have the character

$$V_N(\theta) = \sum \|y(t) - \hat{y}(t|\theta)\|^2 \quad (2)$$

We also discuss how to find the best model (minimize the criterion) and how to assess its properties.

4. We shall describe special methods for linear black-box models. This includes frequency analysis, spectral analysis and so called subspace methods for linear state-space models.
5. We then turn to the practical issues of system identification, to assure good quality of the data by proper experiment design, how to decide upon a good model structure, and how to deal with the data.

## DISPLAYING THE BASIC IDEAS: ARX MODELS AND THE LINEAR LEAST SQUARES METHOD

### The Model

We shall generally denote the system's input and output at time  $t$  by  $u(t)$  and  $y(t)$ , respectively. Perhaps the most basic relationship between the input and output is the *linear difference equation*

$$\begin{aligned} y(t) + a_1 y(t-1) + \dots + a_n y(t-n) \\ = b_1 u(t-1) + \dots + b_m u(t-m) \end{aligned} \quad (3)$$

We have chosen to represent the system in *discrete time*, primarily since observed data are always collected by sampling. It is thus more straightforward to relate observed data to discrete time models. Nothing prevents us, however, from working with continuous time models: we shall return to that later.

In Eq. (3), we assume the *sampling interval* to be one time unit. This is not essential but makes notation easier.

A pragmatic and useful way to see Eq. (3) is to view it as a way of determining the next output value given previous observations

$$\begin{aligned} y(t) = -a_1 y(t-1) - \dots - a_n y(t-n) \\ + b_1 u(t-1) + \dots + b_m u(t-m) \end{aligned} \quad (4)$$

For more compact notation, we introduce the vectors

$$\theta = [a_1, \dots, a_n, b_1, \dots, b_m]^T \quad (5)$$

$$\varphi(t) = [-y(t-1) \dots -y(t-n)u(t-1) \dots u(t-m)]^T \quad (6)$$

With these, Eq. (4) can be rewritten as

$$y(t) = \varphi^T(t)\theta$$

To emphasize that the calculation of  $y(t)$  from past data [Eq. (4)] indeed depends on the parameters in  $\theta$ , we shall rather call this calculated value  $\hat{y}(t|\theta)$  and write

$$\hat{y}(t|\theta) = \varphi^T(t)\theta \quad (7)$$

**The Least Squares Method**

Now suppose for a given system that we do not know the values of the parameters in  $\theta$ , but that we have recorded inputs and outputs over a time interval  $1 \leq t \leq N$ :

$$Z^N = \{u(1), y(1), \dots, u(N), y(N)\} \quad (8)$$

An obvious approach is then to select  $\theta$  in Eqs. (3) through (7) so as to fit the calculated values  $\hat{y}(t|\theta)$  as well as possible to the measured outputs by the least squares method:

$$\min_{\theta} V_N(\theta, Z^N) \quad (9)$$

where

$$\begin{aligned} V(\theta, Z^N) &= \frac{1}{N} \sum_{t=1}^N (y(t) - \hat{y}(t|\theta))^2 \\ &= \frac{1}{N} \sum_{t=1}^N (y(t) - \varphi^T(t)\theta)^2 \end{aligned} \quad (10)$$

We shall denote the value of  $\theta$  that minimizes Eq. (9) by  $\hat{\theta}_N$ :

$$\hat{\theta}_N = \arg \min_{\theta} V_N(\theta, Z^N) \quad (11)$$

(*arg min* means the minimizing argument, i.e., that value of  $\theta$  which minimizes  $V_N$ .)

Since  $V_N$  is quadratic in  $\theta$ , we can find the minimum value easily by setting the derivative to zero:

$$0 = \frac{d}{d\theta} V_N(\theta, Z^N) = -\frac{2}{N} \sum_{t=1}^N \varphi(t)(y(t) - \varphi^T(t)\theta)$$

which gives

$$\sum_{t=1}^N \varphi(t)y(t) = \sum_{t=1}^N \varphi(t)\varphi^T(t)\theta \quad (12)$$

or

$$\hat{\theta}_N = \left[ \sum_{t=1}^N \varphi(t)\varphi^T(t) \right]^{-1} \sum_{t=1}^N \varphi(t)y(t) \quad (13)$$

Once the vectors  $\varphi(t)$  are defined, the solution can easily be found by modern numerical software, such as MATLAB, see Ref. 2.

**Example 3 First-Order Difference Equation**

Consider the simple model

$$y(t) + \alpha y(t-1) = bu(t-1)$$

This gives us the estimate according to Eqs. (5), (6) and (13)

$$\begin{bmatrix} \hat{a}_N \\ \hat{b}_N \end{bmatrix} = \begin{bmatrix} \sum y^2(t-1) - \sum y(t-1)u(t-1) \\ -\sum y(t-1)u(t-1) \sum u^2(t-1) \end{bmatrix}^{-1} \begin{bmatrix} -\sum y(t)y(t-1) \\ \sum y(t)u(t-1) \end{bmatrix}$$

All sums are from  $t = 1$  to  $t = N$ . A typical convention is to take values outside the measured range to be zero. In this case, we would thus take  $y(0) = 0$ .

The simple model in Eq. (3) and the well-known least squares method in Eq. (13) form the archetype of system identification. Not only that, they also give the most commonly used parametric identification method and are much more versatile than perhaps perceived at first sight. In particular, one should realize that Eq. (3) can directly be extended to several different inputs [this just calls for a redefinition of  $\varphi(t)$  in Eq. (6)] and that the inputs and outputs do not have to be the raw measurements. On the contrary, it is often most important to think over the physics of the application and come up with suitable inputs and outputs for Eq. (3), formed from the actual measurements.

**Example 4 An Immersion Heater**

Consider a process consisting of an immersion heater immersed in a cooling liquid. We measure

- $v(t)$ : the voltage applied to the heater
- $r(t)$ : the temperature of the liquid
- $y(t)$ : the temperature of the heater coil surface

Suppose we need a model for how  $y(t)$  depends on  $r(t)$  and  $v(t)$ . Some simple considerations based on common sense and high school physics (“semiphysical modeling”) reveal the following:

- The change in temperature of the heater coil over one sample is proportional to the electrical power in it (the inflow power) minus the heat loss to the liquid.
- The electrical power is proportional to  $v^2(t)$ .
- The heat loss is proportional to  $y(t) - r(t)$ .

This suggests the model

$$y(t) = y(t-1) + \alpha v^2(t-1) - \beta(y(t-1) - r(t-1))$$

which fits into the form

$$y(t) + \theta_1 y(t-1) = \theta_2 v^2(t-1) + \theta_3 r(t-1)$$

This is a two input ( $v^2$  and  $r$ ) and one output model and corresponds to choosing

$$\varphi(t) = [-y(t-1) \ v^2(t-1) \ r(t-1)]^T$$

in Eq. (7).

### Some Statistical Remarks

Model structures, such as Eq. (7) that are linear in  $\theta$ , are known in statistics as *linear regression*, and the vector  $\varphi(t)$  is called the *regression vector* (its components are the *regressors*). “Regress” here alludes to the fact that we try to calculate (or describe)  $y(t)$  by “going back” to  $\varphi(t)$ . Models such as Eq. (3), where the regression vector  $\varphi(t)$  contains old values of the variable to be explained  $y(t)$  are then partly *auto-regressions*. For that reason, the model structure in Eq. (3) has the standard name ARX-model (auto-regression with extra inputs).

There is rich statistical literature on the properties of the estimate  $\hat{\theta}_N$  under varying assumptions (10). So far, we have just viewed Eqs. (9) and (10) as “curve-fitting.” Later, we shall deal with a more comprehensive statistical discussion, which includes the ARX model as a special case. Some direct calculations will be done in the following subsection.

### Model Quality and Experiment Design

Let us consider the simplest special case, that of a finite impulse response (FIR) model. That is obtained from Eq. (3) by taking  $n = 0$ :

$$y(t) = b_1 u(t-1) + \dots + b_m u(t-m) \quad (14)$$

Suppose that the observed data really have been generated by a similar mechanism

$$y(t) = b_1^0 u(t-1) + \dots + b_m^0 u(t-m) + e(t) \quad (15)$$

where  $e(t)$  is a white noise sequence with variance  $\lambda$ , but otherwise unknown. (That is,  $e(t)$  can be described as a sequence of independent random variables with zero mean values and variances  $\lambda$ .) Analogous to Eq. (7), we can write this as

$$y(t) = \varphi^T(t)\theta_0 + e(t) \quad (16)$$

We can now replace  $y(t)$  in Eq. (13) by the above expression and obtain

$$\begin{aligned} \hat{\theta}_N &= \left[ \sum_{t=1}^N \varphi(t)\varphi^T(t) \right]^{-1} \sum_{t=1}^N \varphi(t)y(t) \\ &= \left[ \sum_{t=1}^N \varphi(t)\varphi^T(t) \right]^{-1} \left[ \sum_{t=1}^N \varphi(t)\varphi^T(t)\theta_0 + \sum_{t=1}^N \varphi(t)e(t) \right] \end{aligned}$$

or

$$\tilde{\theta}_N = \hat{\theta}_N - \theta_0 = \left[ \sum_{t=1}^N \varphi(t)\varphi^T(t) \right]^{-1} \sum_{t=1}^N \varphi(t)e(t) \quad (17)$$

Suppose that the input  $u$  is independent of the noise  $e$ . Then,  $\varphi$  and  $e$  are independent in this expression, so it is easy to see

that  $E\tilde{\theta}_N = 0$ , since  $e$  has zero mean. The estimate is consequently *unbiased*. Here,  $E$  denotes *mathematical expectation*.

We can also form the expectation of  $\tilde{\theta}_N\tilde{\theta}_N^T$ , that is the covariance matrix of the parameter error. Denote the matrix within brackets by  $R_N$ . Take expectation with respect to the white noise  $e$ . Then,  $R_N$  is a deterministic matrix, and we have

$$\frac{\partial^2}{\partial \theta^2} \hat{y}(t|\theta)$$

since the double sum collapses to  $\lambda R_N$ .

We have, thus, computed the covariance matrix of the estimate  $\hat{\theta}_N$ . It is determined entirely by the input properties and the noise level. Moreover, define

$$\bar{R} = \lim_{N \rightarrow \infty} \frac{1}{N} R_N \quad (19)$$

This will be the *covariance matrix* of the input, that is, the  $i - j$ -element of  $\bar{R}$  is

$$R_{uu}(i-j) = E u(t-i)u(t-j)$$

If the matrix  $\bar{R}$  is non-singular, we find that the covariance matrix of the parameter estimate is approximately (and the approximation improves as  $N \rightarrow \infty$ )

$$P_N = \frac{\lambda}{N} \bar{R}^{-1} \quad (20)$$

A number of things follow from this. All of them are typical of the general properties to be described later

- The covariance decays like  $1/N$ , so the parameters approach the limiting value at the rate  $1/\sqrt{N}$ .
- The covariance is proportional to the noise-to-signal ratio. That is, it is proportional to the noise variance and inversely proportional to the input power.
- The covariance does not depend on the input's or noise's signal shapes, only on their variance/covariance properties.
- Experiment design, that is, the selection of the input  $u$ , aims at making the matrix  $\bar{R}^{-1}$  “as small as possible.” Note that the same  $\bar{R}$  can be obtained for many different signals  $u$ .

## MODEL STRUCTURES I: LINEAR MODELS

### Output Error Models

Starting from Eq. (3), there is actually another, quite different, way to approach the calculation of good values of  $a_i$  and  $b_i$  from observed data in Eq. (8).

Equation (3) describes a linear, discrete-time system with transfer function

$$G(z) = \frac{b_1 z^{n-1} + b_2 z^{n-2} + \dots + b_m z^{n-m}}{z^n + a_1 z^{n-1} + \dots + a_n} \quad (21)$$

(assuming  $n \geq m$ )

Here,  $z$  is the  $z$ -transform variable, and one may simply think of the transfer function  $G$  as a shorthand notation for the difference in Eq. (3).

We shall here use the shift operator  $q$  as an alternative for the variable  $z$  in Eq. (21). The shift operator  $q$  has the properties

$$qu(t) = u(t+1) \quad (22)$$

(just as multiplying a  $z$ -transform by  $z$  corresponds to a time shift).

Given only an input sequence

$$\{u(t), t = 1, \dots, N\}$$

we could then calculate the output for system in Eq. (21) by running  $u$  as input to this system

$$\hat{y}(t|\theta) = G(q)u(t) \quad (23)$$

### Example 5 A First-Order System

Consider the system

$$y(t+1) + ay(t) = bu(t)$$

The output according to Eq. (23) is then obtained as

$$\hat{y}(t|\theta) = \frac{b}{q+a}u(t) = b \sum_{k=1}^{\infty} (-a)^{k-1} u(t-k)$$

or

$$\hat{y}(t+1|\theta) + a\hat{y}(t|\theta) = bu(t) \quad (24)$$

Notice the essential difference between Eqs. (23) and (7). In Eq. (7), we calculated  $\hat{y}(t|\theta)$  using *both* past measured inputs and also *past measured outputs*  $y(t-k)$ . In Eq. (23),  $\hat{y}(t|\theta)$  is calculated from *past inputs only*. As soon as we use data from a real system [that does not *exactly obey* Eq. (13)], there will always be a difference between these two ways of obtaining the computed output.

Now, we could of course still say that a reasonable estimate of  $\theta$  is obtained by minimizing the quadratic fit,

$$\hat{\theta}_N = \arg \min_{\theta} \frac{1}{N} \sum_{t=1}^N [y(t) - \hat{y}(t|\theta)]^2 \quad (25)$$

even when  $\hat{y}(t|\theta)$  is computed according to Eq. (23). Such an estimate is often called an *output-error* estimate since we have formed the fit between a purely simulated output and the measured output. Note that  $\hat{y}(t|\theta)$  according to Eq. (23) is not linear in  $\theta$ , so the function to be minimized in Eq. (25) is not quadratic in  $\theta$ . Hence, some numerical search schemes have to be applied in order to find  $\hat{\theta}_N$  in Eq. (25). Most often in practice, a Gauss-Newton iterative minimization procedure is used.

It follows from the discussion that the estimate obtained by Eq. (25) will in general differ from the one from Eq. (9). What is the essential difference? To answer that question, we will have to discuss various ways of perceiving and describing the disturbances that act on the system.

### Noise Models and Prediction Filters

A linear, finite-dimensional dynamical system can be described by the equation

$$y(t) = \frac{B(q)}{A(q)}u(t) \quad (26)$$

See Eqs. (21)–(22). Based on Eq. (26) we can predict the next output from previous measurements either as in Eq. (23)

$$\hat{y}(t|\theta) = \frac{B(q)}{A(q)}u(t) \quad (27)$$

or as in Eqs. (4) and (7):

$$\hat{y}(t|\theta) = (1 - A(q))y(t) + B(q)u(t) \quad (28)$$

Which one shall we choose? We can make the discussion more general by writing for Eq. (26)

$$y(t) = G(q, \theta)u(t) \quad (29)$$

to indicate that the transfer function depends on the (numerator and denominator) parameters  $\theta$  [as in Eq. (5)]. We can multiply both sides of Eq. (29) by an arbitrary stable filter  $W(q, \theta)$  giving

$$W(q, \theta)y(t) = W(q, \theta)G(q, \theta)u(t) \quad (30)$$

Then, we can add  $y(t)$  to both sides of the equation and rearrange to obtain

$$y(t) = (1 - W(q, \theta))y(t) + W(q, \theta)G(q, \theta)u(t) \quad (31)$$

We assume that the filter  $W$  starts with a 1:

$$W(q, \theta) = 1 + w_1q^{-1} + w_2q^{-2} + \dots$$

so that  $1 - W(q, \theta)$  actually contains a delay. We thus obtain the predictor

$$\hat{y}(t|\theta) = (1 - W(q, \theta))y(t) + W(q, \theta)G(q, \theta)u(t) \quad (32)$$

Note that this formulation is now similar to that of Eq. (28).

We see that the method used in Eq. (27) corresponds to the choice  $W(q, \theta) \equiv 1$ , while the procedure in Eq. (28) is obtained for  $W(q, \theta) = A(q)$ .

Now, does the predictor in Eq. (32) depend on the filter  $W(q, \theta)$ ? Well, if the input-output data are exactly described by Eq. (29) and we know all relevant initial conditions, the predictor in Eq. (32) *produces identical predictions*  $\hat{y}(t|\theta)$ , regardless of the choice of stable filters  $W(q, \theta)$ .

To bring out the relevant differences, we must accept the fact that there will always be disturbances and noise that affect the system, so instead of Eq. (29), we have a true system that relates the inputs and outputs by

$$y(t) = G_0(q)u(t) + v(t) \quad (33)$$

for some disturbance sequence  $\{v(t)\}$ . So Eq. (32) becomes

$$\begin{aligned} \hat{y}(t|\theta) = & \{(1 - W(q, \theta))G_0(q) + W(q, \theta)G(q, \theta)\}u(t) \\ & + (1 - W(q, \theta))v(t) \end{aligned}$$

Now, assume that there exists a value  $\theta_0$ , such that  $G(q, \theta_0) = G_0(q)$ . Then the error of the above prediction becomes

$$\epsilon(t, \theta) = y(t) - \hat{y}(t|\theta_0) = W(q, \theta)v(t) \quad (34)$$

To make this error as small as possible, we must thus match the choice of the filter  $W(q, \theta_0)$  to the properties of the noise  $v(t)$ . Suppose  $v(t)$  can be described as filtered white noise

$$v(t) = H_0(q)e(t) \quad (35)$$

where  $e(t)$  is a sequence of independent random variables. Here, we assume  $H_0(q)$  to be normalized, so that  $H_0(q) = 1 + h_1q^{-1} + \dots$ . Then, it is easy to see from Eq. (34) that no filter  $W(q, \theta_0)$  can do better than  $1/H_0(q)$ , since this makes the prediction error  $\epsilon(t, \theta_0)$  equal to the white noise source  $e(t)$ .

All this leads to the following summarizing conclusion (which is the only thing one needs to understand from this section).

1. In order to distinguish between different predictors, one has to introduce descriptions of the disturbances that act on the process.
2. If the input–output description is assumed to be

$$y(t) = G(q)u(t) + H(q)e(t) \quad (36)$$

where  $\{e(t)\}$  is a white noise source, then the natural predictor of  $y(t)$  given previous observations of inputs and outputs will be

$$\hat{y}(t|\theta) = [1 - H^{-1}(q)]y(t) + H^{-1}(q)G(q)u(t) \quad (37)$$

This predictor gives the smallest possible error if  $\{e(t)\}$  indeed is white noise.

3. Since the dynamics  $G(q)$  and the noise model  $H(q)$  are typically unknown, we will have to work with a parameterized description

$$y(t) = G(q, \theta)u(t) + H(q, \theta)e(t) \quad (38)$$

The corresponding predictor is then obtained from Eq. (37):

$$\hat{y}(t|\theta) = [I - H^{-1}(q, \theta)]y(t) + H^{-1}(q, \theta)G(q, \theta)u(t) \quad (39)$$

We may now return to the question we posed earlier. What is the practical difference between minimizing Eqs. (10) and (25)? Comparing Eq. (23) with Eq. (29), we see that this predictor corresponds to the assumption that  $H = 1$ , i.e., that white measurement noise is added to the output. This also means that minimizing the corresponding prediction error, Eq. (25), will give a clearly better estimate, if this assumption is more or less correct.

### Linear Black-Box Model Parameterization

The model parameterization in Eq. (38) contains a large number of much-used special cases. We have already seen that the ARX model in Eq. (3) corresponds to

$$G(q, \theta) = \frac{B(q)}{A(q)}, \quad H(q, \theta) = \frac{1}{A(q)} \quad (40)$$

That is, we assume the system (plant) dynamics and the noise model to have common poles and no numerator dynamics for the noise. Its main feature is that the predictor  $\hat{y}(t|\theta)$  will be linear in the parameters  $\theta$  according to Eq. (11) or Eq. (7).

We can make Eq. (40) more general by allowing also numerator dynamics. We then obtain the parameterization

$$G(q, \theta) = \frac{B(q)}{A(q)}, \quad H(q, \theta) = \frac{C(q)}{A(q)} \quad (41)$$

The effect of the numerator  $C$  is that the current predicted value of  $y$  will depend upon previous predicted values, not just measured values. This is known as an *ARMAX model*, since the  $C(q)$ -term makes the noise model a moving average of a white noise source. Also, Eq. (41) assumes that the dynamics and the noise model have common poles, and is, therefore, particularly suited for the case where the disturbances enter together with the input, “early in the process,” so to speak.

The output error (OE) model we considered in Eq. (23) corresponds to the case

$$G(q, \theta) = \frac{B(q)}{F(q)}, \quad H(q, \theta) = 1 \quad (42)$$

[We use  $F$  in the denominator to distinguish the case from Eq. (40).] Its unique feature is that the prediction is based on *past inputs only*. It also concentrates on the model dynamics and does not bother about describing the noise.

We can also generalize this model by allowing a general noise model

$$G(q, \theta) = \frac{B(q)}{F(q)}, \quad H(q, \theta) = \frac{C(q)}{D(q)} \quad (43)$$

This particular model parameterization is known as the *Box-Jenkins* (BJ) model, since it as suggested in the well-known book by Box and Jenkins (6).

It differs from the ARMAX-model in Eq. (42) in that it assigns different dynamics (poles) to the noise characteristics from the input–output properties. It is thus better suited for cases where the noise enters “late in the process,” such as measurement noise. See Fig. 3.

One might wonder why we need all these different model parameterizations. As has been mentioned in the text, each has its advantages which can be summarized as follows

- ARX*. Gives a linear regression, very simple to estimate  $\theta$
- ARMAX*. Gives reasonable flexibility to the noise description, assumes that noise enters like the inputs
- OE*. Concentrates on the input–output dynamics
- BJ*. Very flexible, assumes no common characteristics between noise and input–output behavior.

### Physically Parameterized Linear Models

So far, we have treated the parameters  $\theta$  only as vehicles to give reasonable flexibility to the transfer functions in the general linear model in Eq. (38). This model can also be arrived at from other considerations.

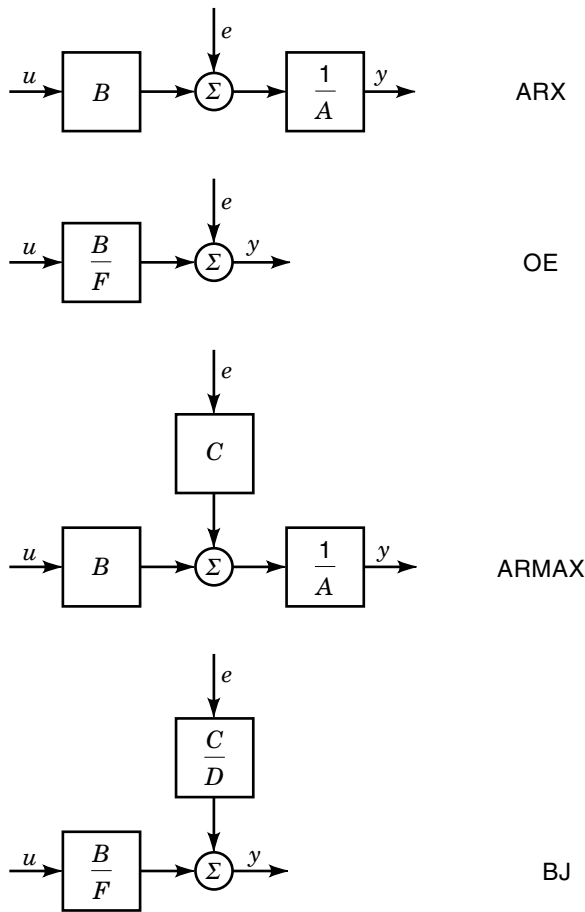


Figure 3. Linear black-box model structures.

Consider a continuous time state space model

$$\dot{x}(t) = A(\theta)x(t) + B(\theta)u(t) \quad (44a)$$

$$y(t) = C(\theta)x(t) + v(t) \quad (44b)$$

Here,  $x(t)$  is the state vector and typically consists of physical variables (such as positions and velocities, etc). The state space matrices  $A$ ,  $B$ , and  $C$  are parameterized by the parameter vector  $\theta$ , reflecting the physical insight we have into the process. The parameters could be physical constants (resistance, heat transfer coefficients, aerodynamical derivatives, etc.) whose values are not known. They could also reflect other types of insights into the system's properties.

### Example 6 An Electric Motor

Consider an electric motor with the input  $u$  being the applied voltage and the output  $y$  being the angular position of the motor shaft.

A first but reasonable approximation of the motor's dynamics is as a first-order system from voltage to angular velocity, followed by an integrator:

$$G(s) = \frac{b}{s(s+a)}$$

If we select the state variables

$$x(t) = \begin{pmatrix} y(t) \\ \dot{y}(t) \end{pmatrix}$$

we obtain the state space form

$$\begin{aligned} \dot{x} &= \begin{pmatrix} 0 & 1 \\ 0 & -a \end{pmatrix} x + \begin{pmatrix} 0 \\ b \end{pmatrix} u \\ y &= (1 \quad 0) x + v \end{aligned} \quad (45)$$

where  $v$  denotes disturbances and noise. In this case, we thus have

$$\begin{aligned} \theta &= \begin{pmatrix} a \\ b \end{pmatrix} \\ A(\theta) &= \begin{pmatrix} 0 & 1 \\ 0 & -a \end{pmatrix} \quad B(\theta) = \begin{pmatrix} 0 \\ b \end{pmatrix} \\ C &= (1 \quad 0) \end{aligned} \quad (46)$$

The parameterization reflects our insight that the system contains an integration, but is in this case not directly derived from detailed physical modeling. Basic physical laws would, in this case, have given us how  $\theta$  depends on physical constants, such as resistance of the wiring, amount of inertia, friction coefficients, and magnetic field constants.

Now, how do we fit a continuous-time model in Eq. (44a) to sampled observed data? If the input  $u(t)$  has been piecewise constant over the sampling interval

$$u(t) = u(kT) \quad kT \leq t < (k+1)T$$

then the states, inputs, and outputs at the sampling instants will be represented by the discrete time model

$$\begin{aligned} x((k+1)T) &= \bar{A}(\theta)x(kT) + \bar{B}(\theta)u(kT) \\ y(kT) &= C(\theta)x(kT) + v(kT) \end{aligned} \quad (47)$$

where

$$\bar{A}(\theta) = e^{A(\theta)T}, \quad \bar{B}(\theta) = \int_0^T e^{A(\theta)\tau} B(\theta) d\tau \quad (48)$$

If the input is not piecewise constant, other sampling rules than Eqs. (47) and (48) will apply. This follows from solving Eq. (44) over one sampling period. We could also further model the added noise term  $v(kT)$  and represent the system in the innovations form

$$\begin{aligned} \bar{x}((k+1)T) &= \bar{A}(\theta)\bar{x}(kT) + \bar{B}(\theta)u(kT) + \bar{K}(\theta)e(kT) \\ y(kT) &= C(\theta)\bar{x}(kT) + e(kT) \end{aligned} \quad (49)$$

where  $\{e(kT)\}$  is white noise. The step from Eqs. (47) to (49) is really a standard Kalman filter step:  $\bar{x}$  will be the one-step ahead predicted Kalman states. A pragmatic way to think about it is as follows: In Eq. (47), the term  $v(kT)$  may not be white noise. If it is colored, we may separate out that part of  $v(kT)$  that cannot be predicted from past values. Denote this part by  $e(kT)$ : it will be the *innovation*. The other part of

$v(kT)$ , the one that can be predicted, can then be described as a combination of earlier innovations,  $e(\ell T)$ ,  $\ell < k$ . Its effect on  $y(kT)$  can then be described via the states by changing them from  $x$  to  $\bar{x}$ , where  $\bar{x}$  contains additional states associated with getting  $v(kT)$  from  $e(\ell T)$ ,  $k \leq \ell$ .

Now, Eq. (49) can be written in input–output form as (let  $T = 1$ )

$$y(t) = G(q, \theta)u(t) + H(q, \theta)e(t) \quad (50)$$

with

$$\begin{aligned} G(q, \theta) &= C(\theta)(qI - \bar{A}(\theta))^{-1}\bar{B}(\theta) \\ H(q, \theta) &= I + C(\theta)(qI - \bar{A}(\theta))^{-1}\bar{K}(\theta) \end{aligned} \quad (51)$$

We are, thus, back at the basic linear model in Eq. (38). The parameterization of  $G$  and  $H$  in terms of  $\theta$  is, however, more complicated than the ones we discussed earlier.

The general estimation techniques, model properties [including the characterization in Eq. (85)], algorithms, etc. apply exactly as described in the section on general parameter estimation techniques.

From these examples, it is also quite clear that non-linear models with unknown parameters can be approached in the same way. We would then typically arrive at a structure

$$\begin{aligned} \dot{x}(t) &= f(x(t), u(t), \theta) \\ y(t) &= h(x(t), u(t), \theta) + v(t) \end{aligned} \quad (52)$$

In this model, all noise effects are collected as additive output disturbances  $v(t)$  which is a restriction, but also a very helpful simplification. If we define  $\hat{y}(t|\theta)$  as the simulated output response to Eq. (52), for a given input, ignoring the noise  $v(t)$ , everything to be said about parameter estimation, model properties, and so on is still applicable.

## MODEL STRUCTURES II: NONLINEAR BLACK-BOX MODELS

In this section, we shall describe the basic ideas behind model structures that have the capability to cover any nonlinear mapping from past data to the predicted value of  $y(t)$ . Recall that we defined a general model structure as a parameterized mapping in Eq. (1):

$$\hat{y}(t|\theta) = g(\theta, Z^{t-1}) \quad (53)$$

We shall consequently allow quite general nonlinear mappings  $g$ . This section will deal with some general principles for how to construct such mappings, and will cover artificial neural networks as a special case. See Refs. 11 and 12 for recent and more comprehensive surveys.

### Nonlinear Black-Box Structures

Now, the model structure family in Eq. (53) is really too general, and it turns out to be useful to write  $g$  as a concatenation of two mappings: one that takes the increasing number of past observations  $Z^{t-1}$  and maps them into a finite dimensional vector  $\varphi(t)$  of fixed dimension and one that takes this vector to the space of the outputs:

$$\hat{y}(t|\theta) = g(\theta, Z^{t-1}) = g(\varphi(t), \theta) \quad (54)$$

where

$$\varphi(t) = \varphi(Z^{t-1}) \quad (55)$$

Let the dimension of  $\varphi$  be  $d$ . As before, we shall call this vector the *regression vector*, and its components will be referred to as the *regressors*. We also allow the more general case that the formation of the regressors is itself parameterized:

$$\varphi(t) = \varphi(Z^{t-1}, \eta) \quad (56)$$

which we for short write  $\varphi(t, \eta)$ . For simplicity, the extra argument  $\eta$  will, however, be used explicitly only when essential for the discussion.

The choice of the nonlinear mapping in Eq. (53) has thus been reduced to two partial problems for dynamical systems:

1. how to choose the nonlinear mapping  $g(\varphi)$  from the regressor space to the output space (i.e., from  $R^d$  to  $R^p$ )
2. how to choose the regressors  $\varphi(t)$  from past inputs and outputs

The second problem is the same for all dynamical systems, and it turns out that the most useful choices of regression vectors are to let them contain past inputs and outputs, and possibly also past predicted/simulated outputs. The regression vector will thus be of the character in Eq. (6). We now turn to the first problem.

### Nonlinear Mappings: Possibilities

Now, let us turn to the nonlinear mapping

$$g(\varphi, \theta) \quad (57)$$

which for any given  $\theta$  maps from  $R^d$  to  $R^p$ . For most of the discussion, we will use  $p = 1$ ; that is, the output is scalar-valued. At this point, it does not matter how the regression vector  $\varphi = (\varphi_1, \dots, \varphi_d)^T$  was constructed. It is just a vector that lives in  $R^d$ .

It is natural to think of the parameterized function family as function expansions:

$$g(\varphi, \theta) = \sum \alpha_k g_k(\varphi) \quad (58)$$

We refer to  $g_k$  as *basis functions*, since the role they play in Eq. (58) is similar to that of a functional space basis. In some particular situations, they do constitute a functional basis. Typical examples are wavelet bases (see below).

We are going to show that expansion in Eq. (58), with different basis functions, plays the role of a unified framework for investigating most known nonlinear black-box model structures.

Now, the key question is: How to choose the basis functions  $g_k$ ? The following facts are essential to understand the connections between most known nonlinear black-box model structures:

- All the  $g_k$  are formed from one “mother basis function,” that we generically denote by  $\kappa(x)$ .
- This function  $\kappa(x)$  is a function of a scalar variable  $x$ .



- Typically,  $g_k$  are dilated (scaled) and translated versions of  $\kappa$ . For the scalar case  $d = 1$ , we may write

$$g_k(\varphi) = g_k(\varphi, \beta_k, \gamma_k) = \kappa(\beta_k(\varphi - \gamma_k)) \quad (59)$$

We, thus, use  $\beta_k$  to denote the dilation parameters and  $\gamma_k$  to denote translation parameters.

**A Scalar Example: Fourier Series** Take  $\kappa(x) = \cos(x)$ . Then Eqs. (58) and (59) will be the Fourier series expansion, with  $\beta_k$  as the frequencies and  $\gamma_k$  as the phases.

**Another Scalar Example: Piecewise Constant Functions** Take  $\kappa$  as the unit interval indicator function:

$$\kappa(x) = \begin{cases} 1 & \text{for } 0 \leq x < 1 \\ 0 & \text{else} \end{cases} \quad (60)$$

and take, for example,  $\gamma_k = k$ ,  $\beta_k = 1/\Delta$ , and  $\alpha_k = f(k\Delta)$ . Then Eqs. (58) and (59) give a piecewise constant approximation of any function  $f$ . Clearly, we would have obtained a quite similar result by a smooth version of the indicator function, for example the Gaussian bell:

$$\kappa(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2} \quad (61)$$

**A Variant of the Piecewise Constant Case** Take  $\kappa$  to be the unit step function

$$\kappa(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases} \quad (62)$$

We, then, just have a variant of Eq. (60), since the indicator function can be obtained as the difference of two steps. A smooth version of the step, like the *sigmoid* function

$$\kappa(x) = \sigma(x) = \frac{1}{1 + e^{-x}} \quad (63)$$

will, of course, give quite similar results.

**Classification of Single-Variable Basis Functions.** Two classes of single-variable basis functions can be distinguished depending on their nature:

- *Local basis functions* are functions having their gradient with bounded support, or at least vanishing rapidly at infinity. Loosely speaking, their variations are concentrated to some interval.
- *Global basis functions* are functions having a nonvanishing gradient for all arguments.

Clearly, the Fourier series is an example of a global basis function, while Eqs. (60), (61), (62) and (63) are all local functions.

**Construction of Multivariable Basis Functions.** In multidimensional case ( $d > 1$ ),  $g_k$  are multivariable functions. In practice, they are often constructed from the single-variable function  $\kappa$  in some simple manner. Let us recall the three most often

used methods for constructing multivariable basis functions from single-variable basis functions.

1. **Tensor Product.** Given  $d$  single-variable functions of the different components  $\varphi_j$  of a  $d$ -dimensional vector  $\varphi$ ,  $h_1(\varphi_1), \dots, h_d(\varphi_d)$  (identical or not), the tensor product construction of the corresponding function from  $R^d$  is then given by their product. In the present case, this means that the basis functions are constructed from the scalar function  $\kappa$  as

$$g_k(\varphi) = \prod_{j=1}^d \kappa(\beta_k^j(\varphi_j - \gamma_k^j)) \quad (64)$$

2. **Radial Construction.** For any single-variable function  $\kappa$ , the radial construction of multivariable basis function of  $\varphi \in R^d$  has the form

$$g_k(\varphi) = g_k f(\varphi, \beta_k, \gamma_k) = \kappa(\|\varphi - \gamma_k\|_{\beta_k}) \quad (65)$$

where  $\|\cdot\|_{\beta_k}$  denotes any chosen norm on the space of the regression vector  $\varphi$ . The norm could typically be a quadratic norm

$$\|\varphi\|_{\beta_k}^2 = \varphi^T \beta_k \varphi \quad (66)$$

with  $\beta_k$  as a possibly  $k$ -dependent positive definite matrix of dilation (scale) parameters. In simple cases,  $\beta_k$  may be just scaled versions of the identity matrix.

3. **Ridge Construction.** Let  $\kappa$  be any single-variable function. Then for all  $\beta_k \in R^d$ ,  $\gamma_k \in R$ , a *ridge* function is given by

$$g_k(\varphi) = g_k(\varphi, \beta_k, \gamma_k) = \kappa(\beta_k^T \varphi + \gamma_k), \varphi \in R^d \quad (67)$$

The ridge function is thus constant for all  $\varphi$  in the subspace  $\{\varphi \in R^d : \beta_k^T \varphi = \text{constant}\}$ . As a consequence, even if the mother basis function  $\kappa$  has local support, the basis functions  $g_k$  will have unbounded support in this subspace. The resulting basis could be said to be *semi-global*, but the term *ridge function* is more precise.

**Approximation Issues.** For any of the described choices, the resulting model becomes

$$g(\varphi, \theta) = \sum_{k=1}^{en} \alpha_k \kappa(\beta_k(\varphi - \gamma_k)) \quad (68)$$

with the different exact interpretations of the argument  $\beta_k(\varphi - \gamma_k)$  just discussed. The expansion is entirely determined by

- the scalar valued function  $\kappa(x)$  of a scalar variable  $x$
- the way the basis functions are expanded to depend on a vector  $\varphi$

The parameterization in terms of  $\theta$  can be characterized by three types of parameters:

- the *coordinates*  $\alpha$
- the *scale* or *dilation* parameters  $\beta$
- the *location* parameters  $\gamma$

A key issue is how well the function expansion is capable of approximating any possible “true system”  $g_0(\varphi)$ . There is rather extensive literature on this subject. For an identification oriented survey, see, for example, Ref. 12.

The bottom line is easy: For almost any choice of  $\kappa(x)$ —except being a polynomial—the expansion in Eq. (68) can approximate any “reasonable” function  $g_0(\varphi)$  arbitrarily well for sufficiently large  $n$ .

It is not difficult to understand this. It is sufficient to check that the delta function—or the indicator function for arbitrarily small areas—can be arbitrarily well approximated within the expansion. Then, clearly, all reasonable functions can also be approximated. For local  $\kappa$  with radial construction, this is immediate: By scaling and location, an arbitrarily small indicator function can be placed anywhere. For the ridge construction, one needs to show that a number of hyperplanes defined by  $\beta$  and  $\gamma$  can be placed and intersect so that any small area in  $R^d$  is cut out.

The question of how *efficient* the expansion is, that is, how large  $n$  is required to achieve a certain degree of approximation is more difficult and has no general answer. We may point to the following aspects:

- If the scale and location parameters  $\beta$  and  $\gamma$  are allowed to depend on the function  $g_0$  to be approximated, then the number of terms  $n$  required for a certain degree of approximation is much less than if  $\beta_k, \gamma_k; k = 1, \dots$  is an a priori fixed sequence.
- For the local, radial approach, the number of terms required to achieve a certain degree of approximation  $\delta$  of a  $p$  times differentiable function is proportional to

$$n \sim \frac{1}{\delta^{(d/p)}} \quad (69)$$

It thus increases exponentially with the number of regressors. This is often referred to as *the curse of dimensionality*.

**Connection to “Named Structures.”** Here we briefly review some popular structures; other structures related to interpolation techniques are discussed in Refs. 11 and 12.

*Wavelets.* The local approach corresponding to Eqs. (58, and 65) has direct connections to wavelet networks and wavelet transforms. The exact relationships are discussed in Ref. 11. Loosely, we note that via the dilation parameters in  $\rho_k$ , we can work with different scales simultaneously to pick up both local and not-so-local variations. With appropriate translations and dilations of a single suitably chosen function  $\kappa$  (the “mother wavelet”), we can make expansion Eq. (58) orthonormal. This is discussed extensively in Ref. 12.

*Wavelet and Radial Basis Networks.* The choice in Eq. (61) without any orthogonalization is found in both wavelet networks (13) and radial basis neural networks (14).

*Neural Networks.* The ridge choice in Eq. (67) with  $\kappa$  given by Eq. (63) gives a much-used neural network structure through the one hidden layer feedforward sigmoidal net.

*Hinging Hyperplanes.* If instead of using the sigmoid  $\sigma$  function we choose “V-shaped” functions (in the form of a higher-dimensional “open book”), Brieman’s *hinging*

*hyperplane* structure is obtained (15).

*Nearest Neighbors or Interpolation.* By selecting  $\kappa$  as in Eq. (60) and the location and scale vector  $\beta_k, \gamma_k$  in structure Eq. (65), such that exactly one observation falls into each “cube,” the nearest neighbor model is obtained: just load the input-output record into a table, and for a given  $\varphi$ , pick the pair  $(\hat{y}, \hat{\varphi})$  for  $\hat{\varphi}$  closest to the given  $\varphi$ ;  $\hat{y}$  is the desired output estimate. If one replaces Eq. (60) by a smoother function and allows some overlapping of the basis functions, we get interpolation-type techniques such as kernel estimators.

*Fuzzy Models.* *Fuzzy models* based on fuzzy set membership belong to the model structures of class Eq. (58). The basis functions  $g_k$  then are constructed from the fuzzy set membership functions and the inference rules using the tensor approach Eq. (64). The exact relationship is described in Ref. 11.

### Estimating Nonlinear Black-Box Models

The model structure is determined by the following choices:

- the regression vector (typically built up from past inputs and outputs)
- the basic function  $\kappa$
- the number of elements (nodes) in the expansion Eq. (58).

Once these choices have been made,  $\hat{y}(t|\theta) = g(\varphi(t), \theta)$  is a well defined function of past data and the parameters  $\theta$ . The parameters are made up of coordinates in expansion Eq. (58) and from location and scale parameters in the different basis functions.

All the algorithms and analytical results of the next section can thus be applied. For neural network applications, these are also the typical estimation algorithms used, often complemented with *regularization*, which means that a term is added to the criterion Eq. (74) that penalizes the norm of  $\theta$ . This will reduce the variance of the model, in that “spurious” parameters are not allowed to take on large, and mostly random, values (11).

For wavelet applications, it is common to distinguish between those parameters that enter linearly in  $\hat{y}(t|\theta)$  (i.e., the coordinates in the function expansion) and those that enter non-linearly (i.e., the location and scale parameters). Often, the latter are seeded to fixed values, and the coordinates are estimated by the linear least squares method. Basis functions that give a small contribution to the fit (corresponding to non-useful values of the scale and location parameters) can then be trimmed away (“pruning” or “shrinking”).

### GENERAL PARAMETER ESTIMATION TECHNIQUES

In this section, we shall deal with issues that are independent of model structure. Principles and algorithms for fitting models to data, as well as the general properties of the estimated models, are all model-structure independent and equally well applicable to, say, ARMAX models and neural network models.

The section is organized as follows. The general principles for parameter estimation are outlined. Then, we deal with the

asymptotic (in the number of observed data) properties of the models, while algorithms are described in the last section.

### Fitting Models to Data

Earlier, we showed several ways to parameterize descriptions of dynamical systems. This is actually the key problem in system identification. No matter how the problem is approached, the bottom line is that such a model parameterization leads to a predictor

$$\hat{y}(t|\theta) = g(\theta, Z^{t-1}) \quad (70)$$

that depends on the unknown parameter vector and past  $Z^{t-1}$  [see Eq. (8)]. This predictor can be linear in  $y$  and  $u$ . This, in turn, contains several special cases both in terms of black-box models and physically parameterized ones, as was discussed earlier. The predictor could also be of general, nonlinear nature, which was also discussed.

In any case, we now need a method to determine a good value of  $\theta$ , based on the information in an observed, sampled data set in Eq. (8). It suggests itself that the basic least-squares like approach in Eqs. (9) through (11) still is a natural approach, even when the predictor  $\hat{y}(t|\theta)$  is a more general function of  $\theta$ .

A procedure with some more degrees of freedom is the following one:

1. From observed data and the predictor  $\hat{y}(t|\theta)$ , form the sequence of prediction errors,

$$\epsilon(t, \theta) = y(t) - \hat{y}(t|\theta), \quad t = 1, 2, \dots, N \quad (71)$$

2. Possibly filter the prediction errors through a linear filter  $L(q)$ ,

$$\epsilon_F(t, \theta) = L(q)\epsilon(t, \theta) \quad (72)$$

so as to enhance or depress interesting or unimportant frequency bands in the signals.

3. Choose a scalar valued, positive function  $\ell(\cdot)$  so as to measure the “size” or “norm” of the prediction error:

$$\ell(\epsilon_F(t, \theta)) \quad (73)$$

4. Minimize the sum of these norms:

$$\hat{\theta}_N = \arg \min_{\theta} V_N(\theta, Z^N) \quad (74)$$

where

$$V_N(\theta, Z^N) = \frac{1}{N} \sum_{t=1}^N \ell(\epsilon_F(t, \theta)) \quad (75)$$

This procedure is natural and pragmatic; we can still think of it as “curve-fitting” between  $y(t)$  and  $\hat{y}(t|\theta)$ . It also has several statistical and information theoretic interpretations. Most importantly, if the noise source in the system [like in Eq. (38)] is supposed to be a sequence of independent random variables  $\{e(t)\}$  each having a probability density function  $f_e(x)$ , then Eq. (74) becomes the maximum likelihood estimate (MLE) if we

choose

$$L(q) = 1 \quad \text{and} \quad \ell(\epsilon) = -\log f_e(\epsilon) \quad (76)$$

The MLE has several nice statistical features and, thus, gives a strong “moral support” for using the outlined method. Another pleasing aspect is that the method is independent of the particular model parameterization used (although this will affect the actual minimization procedure). For example, the method of “back propagation” often used in connection with neural network parameterizations amounts to computing  $\hat{\theta}_N$  in Eq. (74) by a recursive gradient method. We shall deal with these aspects later.

### Model Quality

An essential question is, of course, what properties will the estimate resulting from Eq. (74) have. These will naturally depend on the properties of the data record  $Z^N$  defined by Eq. (8). It is, in general, a difficult problem to characterize the quality of  $\hat{\theta}_N$  exactly. One normally has to be content with the asymptotic properties of  $\hat{\theta}_N$  as the number of data,  $N$ , tends to infinity.

It is an important aspect of the general identification method in Eq. (74) that the asymptotic properties of the resulting estimate can be expressed in general terms for arbitrary model parameterizations.

The first basic result is the following one:

$$\hat{\theta}_N \rightarrow \theta^* \quad \text{as} \quad N \rightarrow \infty \quad (77)$$

where

$$\theta^* = \arg \min_{\theta} E \ell(\epsilon_F(t, \theta)) \quad (78)$$

That is, as more and more data become available, the estimate converges to that value  $\theta^*$  that would minimize the expected value of the “norm” of the filtered prediction errors. This is in a sense the best possible approximation of the true system that is available within the model structure. The expectation  $E$  in Eq. (78) is taken with respect to all random disturbances that affect the data, and it also includes averaging over the input properties. This means, in particular, that  $\theta^*$  will make  $\hat{y}(t|\theta^*)$  a good approximation of  $y(t)$  with respect to those aspects of the system that are enhanced by the input signal used.

The second basic result is the following one. If  $\{\epsilon(t, \theta^*)\}$  is approximately white noise, then the covariance matrix of  $\hat{\theta}_N$  is approximately given by

$$E(\hat{\theta}_N - \theta^*)(\hat{\theta}_N - \theta^*)^T \sim \frac{\lambda}{N} [E\psi(t)\psi^T(t)]^{-1} \quad (79)$$

where

$$\lambda = E\epsilon^2(t, \theta^*) \quad (80)$$

$$\psi(t) = \frac{d}{d\theta} \hat{y}(t|\theta)|_{\theta=\theta^*} \quad (81)$$

Think of  $\psi$  as the sensitivity derivative of the predictor with respect to the parameters. Then Eq. (79) says that the covariance matrix for  $\hat{\theta}_N$  is proportional to the inverse of the covari-

ance matrix of this sensitivity derivative. This is a quite natural result.

Note that for all these results, the expectation operator  $E$  can, under most general conditions, be replaced by the limit of the sample mean; that is,

$$E\psi(t)\psi^T(t) \leftrightarrow \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{t=1}^N \psi(t)\psi^T(t) \quad (82)$$

The results in Eqs. (77) through (81) are general and hold for all model structures, both linear and nonlinear ones, subject only to some regularity and smoothness conditions. They are also fairly natural and will give the guidelines for all user choices involved in the process of identification. See Ref. 4 for more details around this.

**Characterization of the Limiting Model in a General Class of Linear Models.** Let us apply the general limit result in Eqs. (77)–(78) to the linear model structure in Eq. (38). If we choose a quadratic criterion  $\ell(\epsilon) = \epsilon^2$  (in the scalar output case), then this result tells us, in the time domain, that the limiting parameter estimate is the one that minimizes the filtered prediction error variance (for the input used during the experiment.) Suppose that the data actually have been generated by

$$y(t) = G_0(q)u(t) + v(t) \quad (83)$$

Let  $\Phi_u(\omega)$  be the input spectrum and  $\Phi_v(\omega)$  be the spectrum for the additive disturbance  $v$ . Then, the filtered prediction error can be written

$$\begin{aligned} \epsilon_F(t, \theta) &= \frac{L(q)}{H(q, \theta)} [y(t) - G(q, \theta)u(t)] \\ &= \frac{L(q)}{H(q, \theta)} [(G_0(q) - G(q, \theta))u(t) + v(t)] \end{aligned} \quad (84)$$

By Parseval's relation, the prediction error variance can also be written as an integral over the spectrum of the prediction error. This spectrum, in turn, is directly obtained from Eq. (84), so the limit estimate  $\theta^*$  in Eq. (78) can also be defined as

$$\begin{aligned} \theta^* &= \arg \min_{\theta} \left[ \int_{-\pi}^{\pi} |G_0(e^{i\omega}) - G(e^{i\omega}, \theta)|^2 \frac{\Phi_u(\omega) |L(e^{i\omega})|^2}{|H(e^{i\omega}, \theta)|^2} d\omega \right. \\ &\quad \left. + \int_{-\pi}^{\pi} \Phi_v(\omega) |L(e^{i\omega})|^2 / |H(e^{i\omega}, \theta)|^2 d\omega \right] \end{aligned} \quad (85)$$

If the noise model  $H(q, \theta) = H_*(q)$  does not depend on  $\theta$  [as in the output error model Eq. (42)], the expression in Eq. (85) thus shows that the resulting model  $G(e^{i\omega}, \theta^*)$  will give that frequency function in the model set that is closest to the true one, in a quadratic frequency norm with weighting function

$$Q(\omega) = \Phi_u(\omega) |L(e^{i\omega})|^2 / |H_*(e^{i\omega})|^2 \quad (86)$$

This shows clearly that the fit can be affected by the choice of prefilter  $L$ , the input spectrum  $\Phi_u$ , and the noise model  $H_*$ .

### Measures of Model Fit

Some quite general expressions for the expected model fit, that are independent of the model structure, can also be developed.

Let us measure the (average) fit between any model in Eq. (70) and the true system as

$$\bar{V}(\theta) = E|y(t) - \hat{y}(t|\theta)|^2 \quad (87)$$

Here, expectation  $E$  is over the data properties [i.e., expectation over “ $Z$ ” with the notation Eq. (8)]. Recall that expectation also can be interpreted as sample means as in Eq. (82).

Before we continue, let us note the very important aspect that the fit  $\bar{V}$  will depend not only on the model and the true system *but also on data properties*, like input spectra, possible feedback, etc. We shall say that the fit depends on the *experimental conditions*.

The estimated model parameter  $\hat{\theta}_N$  is a random variable because it is constructed from observed data that can be described as random variables. To evaluate the model fit, we then take the expectation of  $\bar{V}(\hat{\theta}_N)$  with respect to the estimation data. That gives our measure

$$F_N = E\bar{V}(\hat{\theta}_N) \quad (88)$$

In general, the measure  $F_N$  depends on a number of things:

- The model structure used
- The number of data points  $N$
- The data properties for which the fit  $\bar{V}$  is defined
- The properties of the data used to estimate  $\hat{\theta}_N$

The rather remarkable fact is that if the two last data properties coincide, then, asymptotically in  $N$  (4, Chap. 16),

$$F_N \approx \bar{V}_N(\theta^*) \left( 1 + \frac{\dim \theta}{N} \right) \quad (89)$$

Here,  $\theta^*$  is the value that minimizes the expected criterion in Eq. (78). The notation  $\dim \theta$  means the number of estimated parameters. The result also assumes that the criterion function  $\ell(\epsilon) = \|\epsilon\|^2$ , and that the model structure is successful in the sense that  $\epsilon_F(t)$  is approximately white noise.

Despite the reservations about the formal validity of Eq. (89), it carries a most important conceptual message: If a model is evaluated on a data set with the same properties as the estimation data, then *the fit will not depend on the data properties*, and it will depend on the model structure only in terms of the number of parameters used and of the best fit offered within the structure.

The expression can be rewritten as follows: Let  $\hat{y}_0(t|t-1)$  denote the “true” one step ahead prediction of  $y(t)$ , and let

$$W(\theta) = E|\hat{y}_0(t|t-1) - \hat{y}(t|\theta)|^2 \quad (90)$$

and let

$$\lambda = E|y(t) - \hat{y}_0(t|t-1)|^2 \quad (91)$$

Then  $\lambda$  is the *innovations* variance, that is, that part of  $y(t)$  that cannot be predicted from the past. Moreover,  $W(\theta^*)$  is the

*bias error*, that is, the discrepancy between the true predictor and the best one available in the model structure. Under the same assumptions as above, Eq. (89) can be rewritten as

$$F_N \approx \lambda + W(\theta^*) + \lambda \frac{\dim \theta}{N} \quad (92)$$

The three terms constituting the model error then have the following interpretations:

- $\lambda$  is the unavoidable error, stemming from the fact that the output cannot be exactly predicted, even with perfect system knowledge.
- $W(\theta^*)$  is the bias error. It depends on the model structure and on the experimental conditions. It will typically decrease as  $\dim \theta$  increases.
- The last term is the variance error. It is proportional to the number of estimated parameters and inversely proportional to the number of data points. It does not depend on the particular model structure or the experimental conditions.

### Algorithmic Aspects

In this section, we shall discuss how to achieve the best fit between observed data and the model, that is, how to carry out the minimization of Eq. (74). For simplicity, we here assume a quadratic criterion and set the prefilter  $L$  to unity:

$$V_N(\theta) = \frac{1}{2N} \sum_{t=1}^N |y(t) - \hat{y}(t|\theta)|^2 \quad (93)$$

No analytic solution to this problem is possible unless the model  $\hat{y}(t|\theta)$  is linear in  $\theta$ , so the minimization has to be done by some numerical search procedure. A classical treatment of the problem of how to minimize the sum of squares is given in Ref. 16.

Most efficient search routines are based on iterative local search in a “down-hill” direction from the current point. We then have an iterative scheme of the following kind:

$$\hat{\theta}^{(i+1)} = \hat{\theta}^{(i)} - \mu_i R_i^{-1} \hat{g}_i \quad (94)$$

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

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

**Search Directions.** The basis for the local search is the gradient

$$V'_N(\theta) = \frac{dV_N(\theta)}{d\theta} = -\frac{1}{N} \sum_{t=1}^N (y(t) - \hat{y}(t|\theta)) \psi(t, \theta) \quad (95)$$

where

$$\psi(t, \theta) = \frac{\partial}{\partial \theta} \hat{y}(t|\theta) \quad (96)$$

The gradient  $\psi$  is, in the general case, a matrix with  $\dim \theta$  rows and  $\dim y$  columns. 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

$$R^{-1}(\theta) V'_N(\theta) \quad (97)$$

where

$$R(\theta) = V''_N(\theta) = \frac{d^2 V_N(\theta)}{d\theta^2} = \frac{1}{N} \sum_{t=1}^N \psi(t, \theta) \psi^T(t, \theta) + \frac{1}{N} \sum_{t=1}^N (y(t) - \hat{y}(t|\theta)) \frac{\partial^2}{\partial \theta^2} \hat{y}(t|\theta) \quad (98)$$

The true Newton direction will, thus, require that the second derivative

$$\frac{\partial^2}{\partial \theta^2} \hat{y}(t|\theta)$$

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

$$R_i = I \quad (99)$$

- *Gauss–Newton Direction.* Use

$$R_i = H_i = \frac{1}{N} \sum_{t=1}^N \psi(t, \hat{\theta}^{(i)}) \psi^T(t, \hat{\theta}^{(i)}) \quad (100)$$

- *Levenberg–Marquard Direction.* Use

$$R_i = H_i + \delta I \quad (101)$$

where  $H_i$  is defined by Eq. (100).

- *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 in Eq. (97) is, however, constructed directly, without explicitly forming and inverting  $V''$ .

It is generally considered (16) that the Gauss–Newton search direction is to be preferred. For ill-conditioned problems, the Levenberg–Marquard modification is recommended.

**Local Minima.** All properties of the estimate that we have discussed really refer to the *global* minimum of the criterion. A fundamental problem with minimization tasks like Eq. (9) 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 worth the effort to find 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, and the genetic algorithm.

### SPECIAL ESTIMATION TECHNIQUES FOR LINEAR BLACK-BOX MODELS

An important feature of a linear, time invariant system is that it is entirely characterized by its *impulse response*. So if we know the system's response to an impulse, we will also know its response to any input. Equivalently, we could study the *frequency response*, which is the Fourier transform of the impulse response.

In this section, we shall consider estimation methods for linear systems that do not use particular model parameterizations. First, we shall consider direct methods to determine the impulse response and the frequency response by simply applying the definitions of these concepts.

Then, spectral analysis for frequency function estimation will be discussed. Finally, a recent method to estimate general linear systems (of given order, by unspecified structure) will be described.

#### Transient and Frequency Analysis

**Transient Analysis.** The first step in modeling is to decide which quantities and variables are important to describe what happens in the system. A simple and common kind of experiment that shows how and in what time span various variables affect each other is called *step-response analysis* or *transient analysis*. In such experiments, the inputs are varied (typically one at a time) as a step:  $u(t) = u_0$ ,  $t < t_0$ ;  $u(t) = u_1$ ,  $t \geq t_0$ . The other measurable variables in the system are recorded during this time. We, thus, study the step response of the system. An alternative would be to study the impulse response of the system by letting the input be a pulse of short duration. From such measurements, information of the following nature can be found:

1. The variables affected by the input in question. This makes it easier to draw block diagrams for the system and to decide which influences can be neglected.
2. The time constants of the system. This also allows us to decide which relationships in the model can be described as static (that is, they have significantly faster time constants than the time scale we are working with).
3. The characteristic (oscillatory, poorly damped, monotone, and the like) of the step responses, as well as the levels of static gains. Such information is useful when studying the behavior of the final model in simulation. Good agreement with the measured step responses should give a certain confidence in the model.

**Frequency Analysis.** If a linear system has the transfer function  $G(q)$ , and the input is

$$u(t) = u_0 \cos \omega kT, \quad (k-1)T \leq t \leq kT \quad (102)$$

then the output after possible transients have faded away will be

$$y(t) = y_0 \cos(\omega t + \varphi), \quad \text{for } t = T, 2T, 3T, \dots \quad (103)$$

where

$$y_0 = |G(e^{i\omega T})| \cdot u_0 \quad (104)$$

$$\varphi = \arg G(e^{i\omega T}) \quad (105)$$

If the system is driven by the input [see Eq. (102)] for a certain  $u_0$  and  $\omega$ , and we measure  $y_0$  and  $\varphi$  from the output signal, it is possible to determine the complex number  $G(e^{i\omega T})$  using Eqs. (104)–(105). By repeating this procedure for a number of different  $\omega$ , we can get a good estimate of the frequency function  $G(e^{i\omega T})$ . This method is called *frequency analysis*. Sometimes, it is possible to see or measure  $u_0$ ,  $y_0$ , and  $\varphi$  directly from graphs of the input and output signals. Most of the time, however, there will be noise and irregularities that make it difficult to determine  $\varphi$  directly. A suitable procedure is then to correlate the output with  $\cos \omega t$  and  $\sin \omega t$ .

#### Estimating the Frequency Response by Spectral Analysis

**Definitions.** The cross spectrum between two (stationary) signals  $u(t)$  and  $y(t)$  is defined as the Fourier transform of their cross covariance function, provided this exists:

$$\Phi_{yu}(\omega) = \sum_{\tau=-\infty}^{\infty} R_{yu}(\tau) e^{i-\omega\tau} \quad (106)$$

where  $R_{yu}(\tau)$  is defined by

$$R_{yu}(\tau) = \mathbf{E}y(t)u(t-\tau) \quad (107)$$

The (auto) spectrum  $\Phi_u(\omega)$  of a signal  $u$  is defined as  $\Phi_{uu}(\omega)$ , that is, as its cross spectrum with itself.

The spectrum describes the frequency contents of the signal. The connection to more explicit Fourier techniques is evident by the following relationship

$$\Phi_u(\omega) = \lim_{N \rightarrow \infty} \frac{1}{N} |U_N(\omega)|^2 \quad (108)$$

where  $U_N$  is the discrete time Fourier transform

$$U_N(\omega) = \sum_{t=1}^N u(t) e^{i\omega t} \quad (109)$$

The relationship in Eq. (108) is shown, for example, in Ref. 3.

Consider now the general linear model:

$$y(t) = G(q)u(t) + v(t) \quad (110)$$

It is straightforward to show that the relationships between the spectra and cross spectra of  $y$  and  $u$  (provided  $u$  and  $v$  are uncorrelated) is given by

$$\Phi_{yu}(\omega) = G(e^{i\omega})\Phi_u(\omega) \quad (111)$$

$$\Phi_y(\omega) = |G(e^{i\omega})|^2\Phi_u(\omega) + \Phi_v(\omega) \quad (112)$$

It is easy to see how the transfer function  $G(e^{i\omega})$  and the noise spectrum  $\psi_v(\omega)$  can be estimated using these expressions, if only we have a method to estimate cross spectra.

**Estimation of Spectra.** The spectrum is defined as the Fourier transform of the correlation function. A natural idea would then be to take the transform of the estimate

$$\hat{R}_{yu}^N(\tau) = \frac{1}{N} \sum_{t=1}^N y(t)u(t-\tau) \quad (113)$$

That will not work in most cases, though. The reason could be described as follows: The estimate  $\hat{R}_{yu}^N(\tau)$  is not reliable for large  $\tau$  since it is based on only a few observations. These “bad” estimates are mixed with good ones in the Fourier transform, thus creating an overall bad estimate. It is better to introduce a weighting, so that correlation estimates for large lags  $\tau$  carry a smaller weight:

$$\hat{\Phi}_{yu}^N(\omega) = \sum_{\ell=-\gamma}^{\gamma} \hat{R}_{yu}^N(\ell) \cdot w_{\gamma}(\ell) e^{-i\ell\omega} \quad (114)$$

This spectral estimation method is known as the *Blackman–Tukey approach*. Here,  $w_{\gamma}(\ell)$  is a window function that decreases with  $|\ell|$ . This function controls the trade-off between frequency resolution and variance of the estimate. A function that gives significant weights to the correlation at large lags will be able to provide finer frequency details (a longer time span is covered). At the same time, it will have to use “bad” estimates, so the statistical quality (the variance) is poorer. We shall return to this trade-off in a moment. How should we choose the shape of the window function  $w_{\gamma}(\ell)$ ? There is no optimal solution to this problem, but the most common window used in spectral analysis is the *Hamming window*:

$$w_{\gamma}(k) = \begin{cases} \frac{1}{2} \left( 1 + \cos \frac{\pi k}{\gamma} \right) & |k| < \gamma \\ 0 & |k| \geq \gamma \end{cases} \quad (115)$$

From the spectral estimates  $\Phi_u$ ,  $\Phi_y$ , and  $\Phi_{yu}$  obtained in this way, we can now use Eq. (111) to obtain a natural estimate of the frequency function  $G(e^{i\omega})$ :

$$\hat{G}_N(e^{i\omega}) = \frac{\hat{\Phi}_{yu}^N(\omega)}{\hat{\Phi}_u^N(\omega)} \quad (116)$$

Furthermore, the disturbance spectrum can be estimated from Eq. (112) as

$$\hat{\Phi}_v^N(\omega) = \hat{\Phi}_y^N(\omega) - \frac{|\hat{\Phi}_{yu}^N(\omega)|^2}{\hat{\Phi}_u^N(\omega)} \quad (117)$$

To compute these estimates, the following steps are performed:

1. Collect data  $y(k)$ ,  $u(k)$ ,  $k = 1, \dots, N$ .
2. Subtract the corresponding sample means from the data. This will avoid bad estimates at very low frequencies.
3. Choose the width of the lag window  $w_{\gamma}(k)$ .
4. Compute  $\hat{R}_y^N(k)$ ,  $\hat{R}_u^N(k)$ , and  $\hat{R}_{yu}^N(k)$  for  $|k| \leq \gamma$  according to Eq. (113).
5. Form the spectral estimates  $\hat{\Phi}_y^N(\omega)$ ,  $\hat{\Phi}_u^N(\omega)$ , and  $\hat{\Phi}_{yu}^N(\omega)$  according to Eq. (114) and analogous expressions.
6. Form Eq. (116) and possibly also Eq. (117).

**Quality of the Estimates.** The estimates  $\hat{G}_N$  and  $\hat{\Phi}_v^N$  are formed entirely from estimates of spectra and cross spectra. Their properties will, therefore, be inherited from the properties of the spectral estimates. For the Hamming window with width  $\gamma$ , it can be shown that the frequency resolution will be

about

$$\frac{\pi}{\gamma\sqrt{2}} \quad \text{radians/time unit} \quad (118)$$

This means that details in the true frequency function that are finer than this expression will be smeared out in the estimate. It is also possible to show that the estimate’s variances satisfy

$$\text{Var } \hat{G}_N(i\omega) \approx 0.7 \cdot \frac{\gamma}{N} \cdot \frac{\Phi_v(\omega)}{\Phi_u(\omega)} \quad (119)$$

and

$$\text{Var } \hat{\Phi}_v^N(\omega) \approx 0.7 \cdot \frac{\gamma}{N} \cdot \Phi_v^2(\omega) \quad (120)$$

[“Variance” here refers to taking expectation over the noise sequence  $v(t)$ .] Note that the relative variance in Eq. (119) typically increases dramatically as  $\omega$  tends to the Nyquist frequency. The reason is that  $|G(i\omega)|$  typically decays rapidly, while the noise-to-signal ratio  $\Phi_v(\omega)/\Phi_u(\omega)$  has a tendency to increase as  $\omega$  increases. In a Bode diagram, the estimates will thus show considerable fluctuations at high frequencies. Moreover, the constant frequency resolution in Eq. (118) will look thinner and thinner at higher frequencies in a Bode diagram due to the logarithmic frequency scale.

See Ref. 3 for a more detailed discussion.

### Choice of Window Size

The choice of  $\gamma$  is a pure trade-off between frequency resolution and variance (variability). For a spectrum with narrow resonance peaks, it is thus necessary to choose a large value of  $\gamma$  and accept a higher variance. For a more flat spectrum, smaller values of  $\gamma$  will do well. In practice, a number of different values of  $\gamma$  are tried. Often, we start with a small value of  $\gamma$  and increase it successively until an estimate is found that balances the trade-off between frequency resolution (true details) and variance (random fluctuations). A typical value for spectra without narrow resonances is  $\gamma = 20 - 30$ .

### Subspace Estimation Techniques for State Space Models

A linear system can always be represented in state space form:

$$\begin{aligned} x(t+1) &= Ax(t) + Bu(t) + w(t) \\ y(t) &= Cx(t) + Du(t) + e(t) \end{aligned} \quad (121)$$

We assume that we have no insight into the particular structure, and we would just estimate any matrices  $A$ ,  $B$ ,  $C$ , and  $D$ , that give a good description of the input-output behavior of the system. This is not without problems, among other things because there are an infinite number of such matrices that describe the same system (the similarity transforms). The coordinate basis of the state-space realization, thus, needs to be fixed.

Let us for a moment assume that not only are  $u$  and  $y$  measured, but also the sequence of state vectors  $x$ . This would, by the way, fix the state-space realization coordinate basis. Now, with known  $u$ ,  $y$ , and  $x$ , the model in Eq. (121) becomes a linear regression: the unknown parameters, all of

the matrix entries in all the matrices, mix with measured signals in linear combinations. To see this clearly, let

$$\begin{aligned} Y(t) &= \begin{pmatrix} x(t+1) \\ y(t) \end{pmatrix} \\ \Theta &= \begin{pmatrix} A & B \\ C & D \end{pmatrix} \\ \Phi(t) &= \begin{pmatrix} x(t) \\ u(t) \end{pmatrix} \\ E(t) &= \begin{pmatrix} w(t) \\ e(t) \end{pmatrix} \end{aligned}$$

Then, Eq. (121) can be rewritten as

$$Y(t) = \Theta\Phi(t) + E(t) \quad (122)$$

From this, all the matrix elements in  $\Theta$  can be estimated by the simple least squares method, as described earlier. The covariance matrix for  $E(t)$  can also be estimated easily as the sample sum of the model residuals. That will give the covariance matrices for  $w$  and  $e$ , as well as the cross covariance matrix between  $w$  and  $e$ . These matrices will, among other things, allow us to compute the Kalman filter for Eq. (121). Note that all of the above holds without changes for multivariable systems, that is, when the output and input signals are vectors.

The only remaining problem is where to get the state vector sequence  $x$ . It has long been known, for example (17,18), that all state vectors  $x(t)$  that can be reconstructed from input-output data in fact are linear combinations of the components of the  $n$   $k$ -step ahead output predictors

$$\hat{y}(t+k|t), \quad k = \{1, 2, \dots, n\} \quad (123)$$

where  $n$  is the model order (the dimension of  $x$ ). See also Appendix 4.A in Ref. 4. We could then form these predictors and select a basis among their components:

$$x(t) = L \begin{pmatrix} \hat{y}(t+1|t) \\ \vdots \\ \hat{y}(t+n|t) \end{pmatrix} \quad (124)$$

The choice of  $L$  will determine the basis for the state-space realization and is done in such a way that it is well conditioned. The predictor  $\hat{y}(t+k|t)$  is a linear function of  $u(s)$ ,  $y(s)$ ,  $1 \leq s \leq t$  and can efficiently be determined by linear projections directly on the input output data. (There is one complication in that  $u(t+1), \dots, u(t+k)$  should not be predicted, even if they affect  $y(t+k)$ .)

What we have described now is the *subspace projection* approach to estimating the matrices of the state-space model in Eq. (121), including the basis for the representation and the noise covariance matrices. There are a number of variants of this approach. See among several references, for example, Refs. 19 and 20.

The approach gives very useful algorithms for model estimation and is particularly well suited for multivariable systems. The algorithms also allow numerically very reliable implementations. At present, the asymptotic properties of the methods are not fully investigated, and the general results

quoted earlier are not directly applicable. Experience has shown, however, that confidence intervals computed according to the general asymptotic theory are good approximations. One may also use the estimates obtained by a subspace method as initial conditions for minimizing the prediction error criterion [see Eq. (74)].

## DATA QUALITY

It is desirable to affect the conditions under which the data are collected. The objective with such *experiment design* is to make the collected data set  $Z^N$  as informative as possible with respect to the models to be built using the data. A considerable amount of theory around this topic can be developed and we shall here just review some basic points.

The first and most important point is the following one

1. The input signal  $u$  must be such that it exposes all the relevant properties of the system.

It must, thus, not be too "simple." For example, a pure sinusoid

$$u(t) = A \cos \omega t$$

will only give information about the system's frequency response at frequency  $\omega$ . This can also be seen from Eq. (85). The rule is that

- the input must contain at least as many different frequencies as the order of the linear model to be built. To be on the safe side, a good choice is to let the input be random (such as filtered white noise). It then contains all frequencies.

Another case where the input is too simple is when it is generated by feedback such as

$$u(t) = -Ky(t) \quad (125)$$

If we would like to build a first-order ARX model

$$y(t) + \alpha y(t-1) = bu(t-1) + e(t)$$

we find that for any given  $\alpha$ , all models such that

$$a + bK = \alpha$$

will give identical input-output data. We can, thus, not distinguish between these models using an experiment with Eq. (125). That is, we can not distinguish between any combinations of " $\alpha$ " and " $b$ " if they satisfy the above condition for a given " $\alpha$ ." The rule is

- If closed-loop experiments have to be performed, the feedback law must not be too simple. It is to be preferred that a set-point in the regulator is being changed in a random fashion.

The second main point in experimental design is

2. Allocate the input power to those frequency bands where a good model is particularly important.

This is also seen from the expression in Eq. (85).



If we let the input be filtered white noise, this gives information on how to choose the filter. In the time domain, it is often useful to think like this:

- Use binary (two-level) inputs if linear models are to be built; this gives maximal variance for amplitude-constrained inputs.
- Check that the changes between the levels are such that the input occasionally stays on one level so long that a step response from the system has time, more or less, to settle. There is no need to let the input signal switch so quickly back and forth that no response in the output is clearly visible.

Note that the second point is really just a reformulation in the time domain of the basic frequency domain advice: let the input energy be concentrated in the important frequency bands.

A third basic piece of advice about experiment design concerns the choice of sampling interval.

3. A typical good sampling frequency is 10 times the bandwidth of the system.

That corresponds roughly to 5–7 samples along the rise time of a step response.

## MODEL VALIDATION AND MODEL SELECTION

### A Pragmatic Viewpoint

The system identification process has, as we have seen, these basic ingredients

- The set of models
- The data
- The selection criterion

Once these have been decided upon, we have, at least implicitly, defined a model: The one in the set that best describes the data according to the criterion. It is thus, in a sense, the best available model in the chosen set. But is it good enough? It is the objective of model validation to answer that question. Often, the answer turns out to be “no,” and we then have to go back and review the choice of model set, or perhaps modify the data set. See Fig. 4.

How do we check the quality of a model? The prime method is to investigate how well it is capable of reproducing the behavior of a new set of data (*the validation data*) that was not used to fit the model. That is, we simulate the obtained model with a new input and compare this simulated output. One may then use one’s eyes or numerical measurements of fit to decide if the fit in question is good enough. Suppose we have obtained several different models in different model structures (say a fourth-order ARX model, a second-order BJ model, a physically parameterized one, and so on) and would like to know which one is best. The simplest and most pragmatic approach to this problem is then to simulate each one of them on validation data, evaluate their performance, and pick the one that shows the most favorable fit to measured data. (This could indeed be a subjective criterion!)

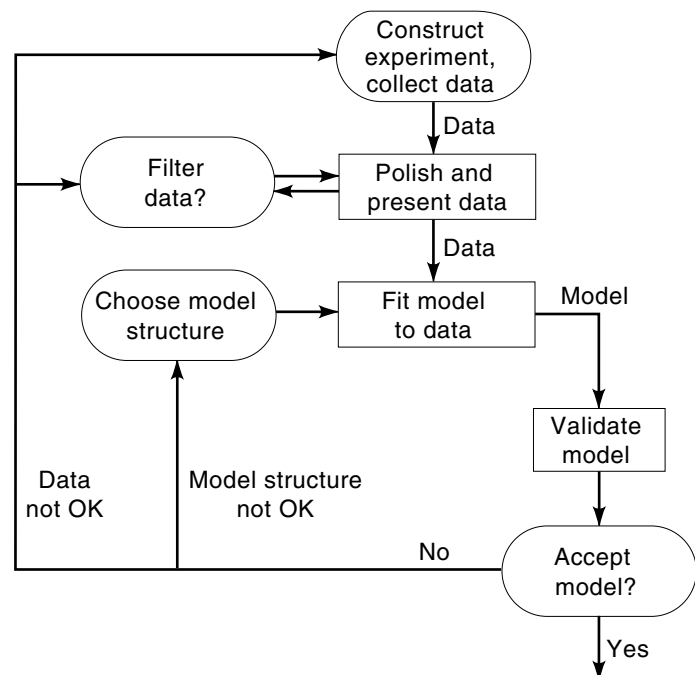


Figure 4. The identification loop.

### The Bias-Variance Trade-off

The heart of the model structure selection process is to handle the trade-off between bias and variance, as formalized by Eq. (92). The “best” model structure is the one that minimizes  $F_N$ , the fit between the model and the data for a fresh data set—one that was not used for estimating the model. Most procedures for choosing the model structures are also aiming at finding this best choice.

**Cross Validation.** A very natural and pragmatic approach is *cross validation*. This means that the available data set is split into two parts, *estimation data*,  $Z_{\text{est}}^{N_1}$  that is used to estimate the models:

$$\hat{\theta}_{N_1} = \arg \min V_{N_1}(\theta, Z_{\text{est}}^{N_1}) \quad (126)$$

and *validation data*,  $Z_{\text{val}}^{N_2}$  for which the criterion is evaluated:

$$\hat{F}_{N_1} = V_{N_2}(\hat{\theta}_{N_1}, Z_{\text{val}}^{N_2}) \quad (127)$$

Here,  $V_N$  is the criterion in Eq. (75). Then,  $\hat{F}_{N_1}$  will be an unbiased estimate of the measure  $F_N$ , defined by Eq. (88), which was discussed at length in the previous section. The procedure would be to try out a number of model structures and choose the one that minimizes  $\hat{F}_{N_1}$ .

Such cross validation techniques to find a good model structure have an immediate intuitive appeal. We simply check if the candidate model is capable of “reproducing” data it hasn’t yet seen. If that works well, we have some confidence in the model, regardless of any probabilistic framework that might be imposed. Such techniques are also the most commonly used ones.

A few comments could be added. In the first place, one could use different splits of the original data into estimation

and validation data. For example, in statistics, there is a common cross validation technique called “leave one out.” This means that the validation data set consists of one data point “at a time” but successively applied to the whole original set. In the second place, the test of the model on the validation data does not have to be in terms of the particular criterion [see Eq. (127)]. In system identification, it is common practice to simulate (or predict several steps ahead) the model using the validation data and then visually inspect the agreement between measured and simulated (predicted) output.

**Estimating the Variance Contribution—Penalizing the Model Complexity.** It is clear that the criterion in Eq. (127) has to be evaluated on the validation data to be of any use; it would be strictly decreasing as a function of model flexibility if evaluated on the estimation data. In other words, the adverse effect of the dimension of  $\theta$  shown in Eq. (92) would be missed. There are a number of criteria, often derived from entirely different viewpoints, that try to capture the influence of this variance error term. The two best known ones at Akaike’s Information Theoretic Criterion, AIC, which has the form (for Gaussian disturbances)

$$\hat{V}_N(\theta, Z^N) = \left(1 + \frac{2 \dim \theta}{N}\right) \frac{1}{N} \sum_{t=1}^N \epsilon^2(t, \theta) \quad (128)$$

and Rissanen’s Minimum Description Length Criterion, MDL, in which  $\dim \theta$  in the expression above is replaced by  $\log N \dim \theta$  (21,22).

The criterion  $\hat{V}_N$  is then to be minimized both with respect to  $\theta$  and to a family of model structures. The relation to the expression in Eq. (89) for  $F_N$  is obvious.

### Residual Analysis

The second basic method for model validation is to examine the residuals (the leftovers) from the identification process. These are the prediction errors

$$\epsilon(t) = \epsilon(t, \hat{\theta}_N) = y(t) - \hat{y}(t|\hat{\theta}_N)$$

that is, what the model could not explain. Ideally, these should be independent of information that was at hand at time  $t - 1$ . For example, if  $\epsilon(t)$  and  $u(t - \tau)$  turn out to be correlated, then there are things in  $y(t)$  that originate from  $u(t - \tau)$  but have not been properly accounted for by  $\hat{y}(t|\hat{\theta}_N)$ . The model has then not squeezed out all relevant information about the system from the data.

It is good practice to always check the residuals for such (and other) dependencies. This is known as residual analysis. A basic reference for how to perform this is Ref. 10.

## BACK TO DATA: THE PRACTICAL SIDE OF IDENTIFICATION

### Software for System Identification

In practice, system identification is characterized by some quite heavy numerical calculations to determine the best model in each given class of models. This is mixed with several user choices, trying different model structures, filtering data, and so on. In practical applications, we will thus need good software support. There are now many different com-

mercial packages for identification available, such as Mathwork’s System Identification Toolbox (2), Matrixx’s System Identification Module (23) and PIM (24). They all have in common that they offer the following routines:

- Handling of data, plotting, and so on: Filtering of data, removal of drift, choice of data segments, etc.
- Non-parametric identification methods: Estimation of covariances, Fourier transforms, correlation- and spectral-analysis, etc.
- Parametric estimation methods: Calculation of parametric estimates in different model structures.
- Presentation of models: Simulation of models, estimation and plotting of poles and zeros, computation of frequency functions, and plotting Bode diagrams, etc.
- Model validation: Computation and analysis of residuals ( $\epsilon(t, \hat{\theta}_N)$ ). Comparison between different models’ properties, etc.

The existing program packages differ mainly in various user interfaces and by different options regarding the choice of model structure according to C above. For example, MATLAB’s Identification Toolbox (2) covers all linear model structures discussed here, including arbitrarily parameterized linear models in continuous time.

Regarding the user interface, there is now a clear trend to make it graphically oriented. This avoids syntax problems and relies more on “click and move,” at the same time as tedious menu-labyrinths are avoided. More aspects of CAD tools for system identification are treated in Ref. 25.

### How to Get to a Good Model?

It follows from our discussion that the most essential element in the process of identification—once the data have been recorded—is to try out various model structures, compute the best model in the structures using Eq. (38), and then validate this model. Typically, this has to be repeated with quite a few different structures before a satisfactory model can be found.

While one should not underestimate the difficulties of this process, the following simple procedure to get started and gain insight into the models could be suggested:

1. Find out a good value for the delay between input and output, for example, by using correlation analysis.
2. Estimate a fourth order linear model with this delay using part of the data, and simulate this model with the input and compare the model’s simulated output with the measured output over the whole data record. In MATLAB language, this is simple

```
z = [y u];
compare(z, arx(z(1:200,:), [4 4 1]));
```

If the model/system is unstable or has integrators, use prediction over a reasonable large time horizon instead of simulation.

Now, either of two things happen:

- The comparison “looks good.” Then, we can be confident that with some extra work—trying out different orders

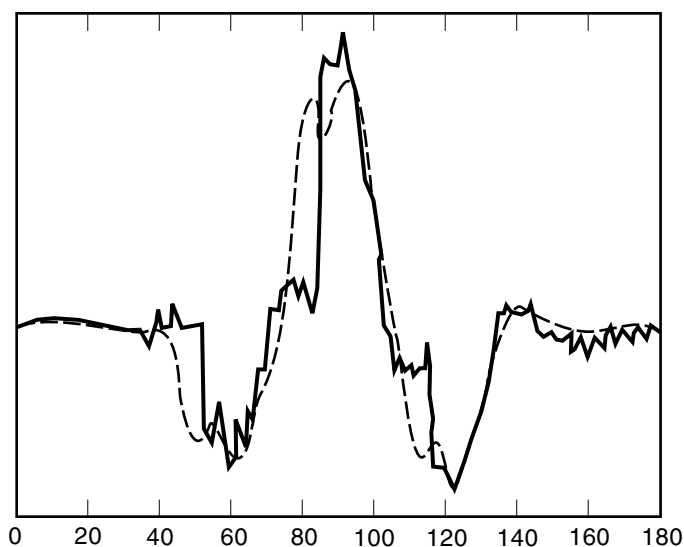
and various noise models—we can fine-tune the model and have an acceptable model quite soon.

- The comparison “does not look good.” Then we must do further work. There are three basic reasons for the failure.
  1. A good description needs higher order linear dynamics. This is actually, in practice, the least likely reason, except for systems with mechanical resonances. One then obviously has to try higher order models or focus on certain frequency bands by band pass filtering.
  2. There are more signals that significantly affect the output. We must then look for what these signals might be, check if they can be measured, and if so, include them among the inputs. Signal sources that cannot be traced or measured are called “disturbances,” and we simply have to live with the fact that they will have an adverse effect on the comparisons.
  3. Some important nonlinearities have been overlooked. We must then resort to semiphysical modeling to find out if some of the measured signals should be subjected to nonlinear transformations. If no such transformations suggest themselves, one might have to try some nonlinear black-box model, like a neural network.

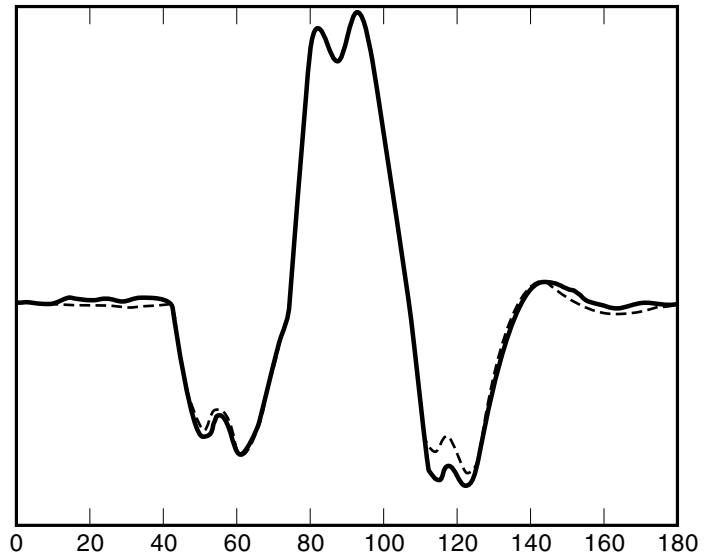
Clearly, this advice does not cover all the art of identification, but it is a reasonable first approximation.

#### Example 7 Aircraft Dynamics

Let us try the recipe on the aircraft data in Fig. 1. Picking the canard angle only as the input, estimating a fourth order model based on the data points 90 to 180, gives Fig. 5. (We use 10-step ahead prediction in this example since the models are unstable, as they should be; JAS has unstable dynamics in this flight case.) It does not “look good.” Let us try alternative 2: More inputs. We repeat the procedure using all three



**Figure 5.** Dashed line: actual pitch rate. Solid line: 10 step ahead predicted pitch rate, based on the fourth order model from canard angle only.



**Figure 6.** As Fig. 5 but using all three inputs.

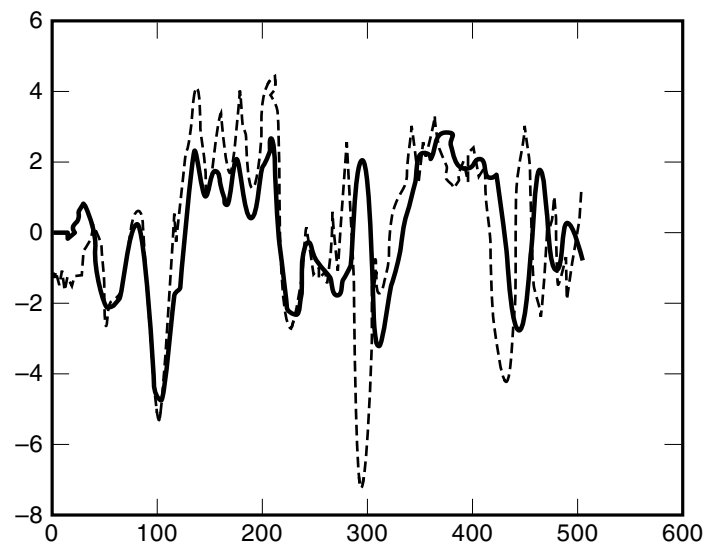
inputs in Fig. 1. That is, the model is computed as

$$\text{arx}([y \ u1 \ u2 \ u3], [4 \ 4 \ 4 \ 4 \ 1 \ 1 \ 1])$$

on the same data set. The comparison is shown in Fig. 6. It “looks good.” By further fine-tuning, as well as using model structures from physical modeling, only slight improvements are obtained.

#### Example 8 Buffer Vessel Dynamics

Let us now consider the pulp process of Fig. 2. We use the  $\kappa$ -number before the vessel as input and the  $\kappa$ -number after the vessel as output. The delay is preliminarily estimated to 12 samples. Our recipe, where a fourth order linear model is estimated using the first 200 samples and then simulated over the whole record gives Fig. 7. It does not look good.



**Figure 7.** Dashed line:  $\kappa$ -number after the vessel, actual measurements. Solid line: simulated  $\kappa$ -number using the input only and a fourth order linear model with delay 12, estimated using the first 200 data points.

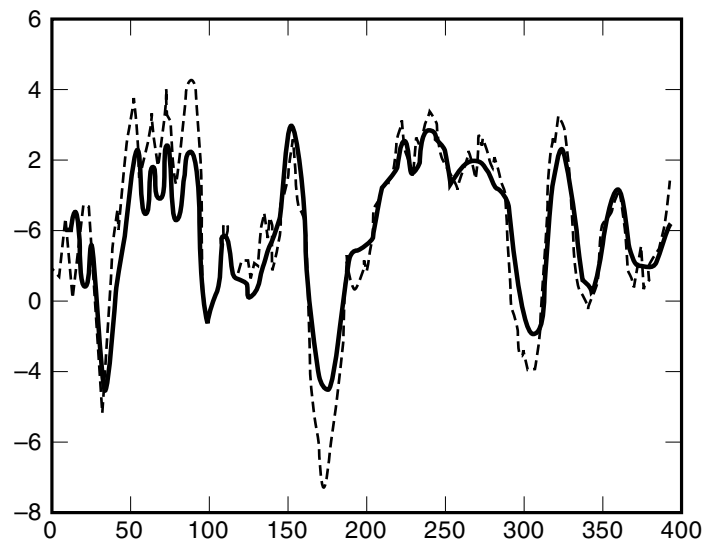


Figure 8. Same as Fig. 7 but applied to resampled data.

Some reflection shows that this process indeed must be non-linear (or time-varying): the flow and the vessel level definitely affect the dynamics. For example, if the flow was a plug flow (no mixing), the vessel would have a dynamics of a pure delay equal to vessel volume divided by flow.

Let us thus resample the data accordingly, that is, so that a new sample is taken (by interpolation from the original measurement) equidistantly in terms of integrated flows divided by volume. In MATLAB terms, this will be

```
z = [y,u]; pf = flow./level;
t = 1:length(z)
newt =
table1([cumsum(pf),t],[pf(1)sum:(pf)]');
newz = table1([t,z],newt);
```

We now apply the same procedure to the resampled data. This gives Fig. 8. This “looks good”. Somewhat better numbers can then be obtained by fine-tuning the orders.

## BIBLIOGRAPHY

1. K. J. Åström and T. Bohlin, Numerical identification of linear dynamic systems from normal operating records, *IFAC Symp. Self-Adapt. Syst.*, Teddington, England, 1965.
2. L. Ljung, *The System Identification Toolbox: The Manual*, Natick, MA: MathWorks Inc., 1995, 4th ed.
3. L. Ljung and T. Glad, *Modeling of Dynamic Systems*, Englewood Cliffs, NJ: Prentice-Hall, 1994.
4. L. Ljung, *System Identification—Theory for the User*, Englewood Cliffs, NJ: Prentice-Hall, 1987.
5. T. Söderström and P. Stoica, *System Identification*, London: Prentice-Hall Int., 1989.
6. G. E. P. Box and D. R. Jenkins, *Time Series Analysis, Forecasting and Control*, San Francisco: Holden-Day, 1970.
7. J. Schoukens and R. Pintelon, *Identification of Linear Systems: A Practical Guideline to Accurate Modeling*, London: Pergamon, 1991.
8. L. Ljung and T. Söderström, *Theory and Practice of Recursive Identification*, Cambridge, MA: MIT Press, 1983.
9. D. Brillinger, *Time Series: Data Analysis and Theory*, San Francisco: Holden-Day, 1981.
10. N. Draper and H. Smith, *Applied Regression Analysis*, 2nd ed, New York: Wiley, 1981, 2nd ed.
11. J. Sjöberg et al., Nonlinear black-box modeling in system identification: A unified overview, *Automatica*, **31** (12): 1691–1724, 1995.
12. A. Juditsky et al., Nonlinear black-box modeling in system identification: Mathematical foundations, *Automatica*, **31** (12): 1724–1750, 1995.
13. Q. Zhang and A. Benveniste, Wavelet networks, *IEEE Trans. Neural Networks*, **3**: 889–898, 1992.
14. T. Poggio and F. Girosi, Networks for approximation and learning, *Proc. IEEE*, **78**: 1481–1497, 1990.
15. L. Breiman, Hinging hyperplanes for regression, classification and function approximation, *IEEE Trans. Inf. Theory*, **39**: 999–1013, 1993.
16. J. E. Dennis and R. B. Schnabel, *Numerical Methods for Unconstrained Optimization and Nonlinear Equations*, Englewood Cliffs, NJ: Prentice-Hall, 1983.
17. J. Rissanen, Basis of invariants and canonical forms for linear dynamic systems, *Automatica*, **10**: 175–182, 1974.
18. H. Akaike, Stochastic theory of minimal realization, *IEEE Trans. Autom. Control*, **AC-19**: 667–674, 1974.
19. P. V. Overschee and B. DeMoor, *Subspace Identification of Linear Systems: Theory, Implementation, Application*, Boston, MA: Kluwer, 1996.
20. W. E. Larimore, System identification, reduced order filtering and modelling via canonical variate analysis, *Proc. 1983 Amer. Control Conf.*, San Francisco, 1983.
21. H. Akaike, A new look at the statistical model identification, *IEEE Trans. Autom. Control*, **AC-19**: 716–723, 1974.
22. J. Rissanen, Modelling by shortest data description, *Automatica*, **14**: 465–471, 1978.
23. *MATRIX<sub>x</sub> Users Guide*, Santa Clara, CA: Integrated Systems Inc., 1991.
24. I. D. Landau, *System Identification and Control Design Using P.I.M. + Software*, Englewood Cliffs, NJ: Prentice-Hall, 1990.
25. L. Ljung, Identification of linear systems, in E. D. Linkens (ed.), *CAD for Control Systems*, New York: Marcel Dekker, 1993, Chap. 6, pp. 147–165.

LENNART LJUNG  
Linköping University

**SYSTEM IDENTIFICATION, MULTIDIMENSIONAL.** See MULTIDIMENSIONAL SIGNAL PROCESSING.