

Applications of Interval Analysis to Statistical Problems

Vincent Vigneron

¹ IBISC CNRS FRE 2494

Université d'Evry

91020 Evry Cedex, France

{vincent.vigneron}@ibisc.univ-evry.fr

² CES CNRS-UMR 8174 Equipe SAMOS-MATISSE

75634 Paris cedex 13, France

vigneron@univ-paris1.fr

Abstract. This paper contribution is about *guaranteed* numerical methods based on interval analysis and the application of these methods to vast classes of statistical problems. 'Guaranteed' means here the inner and outer approximations of the sets of interest are obtained, which can be made as precise as desired, at the cost of increasing the computational effort. It thus becomes possible to achieve tasks still thought by many to be out of the reach of numerical methods, such as finding *all* solutions of non-linear equations and inequalities or a global optimizer of possible multi-modal criteria.

1 Introduction

Before using interval analysis as a basic tool in the following section, we shall now introduce its main concepts.

Interval computation is a special case of computation on sets, and set theory provides the formulations for interval analysis. The reader interested in interval analysis may found more developments in [1,2]. Set and interval mathematics come from the same general theory developed during the 30's by the french School of Topology: if a number a and a bound b of a approximate the value of some number x such that $|x - a| \leq b$, then *interval mathematics* tells us that x is in the interval $[a - b, a + b]$. Hence, the type of an interval is dual: at the same time *number* and *set*, with evident implications in set arithmetics. For instance, suppose a real axis provided with an order relation \leq , the interval $[x]$, bounded by \underline{x} and \overline{x} , is a closed connected subset of real numbers $\{x | \underline{x} \leq x \leq \overline{x}\}$. Interval arithmetic follows from order properties and basic operations on real numbers or vectors extend in a natural way to intervals *e.g.*

$$\begin{aligned}[x] + [y] &= \{x + y \mid x \in [x], y \in [y]\} = [\underline{x} + \underline{y}; \overline{x} + \overline{y}] \\ -[y] &= \{-x \mid x \in [y]\} = [-\overline{y}; -\underline{y}] \\ [x] - [y] &= \{x - y \mid x \in [x], y \in [y]\} = [x] + (-[y]) \\ [x] \times [y] &= \{x \times y \mid x \in [x], y \in [y]\} \\ &= [\min(\underline{x} \cdot \underline{y}, \underline{x} \cdot \overline{y}, \overline{x} \cdot \underline{y}, \overline{x} \cdot \overline{y}), \\ &\quad \max(\underline{x} \cdot \underline{y}, \underline{x} \cdot \overline{y}, \overline{x} \cdot \underline{y}, \overline{x} \cdot \overline{y})] \\ \frac{[x]}{[y]} &= [x] \times \frac{1}{[y]}.\end{aligned}$$

However, the arithmetical rules for intervals differ from those for real numbers. For instance, $x^2 + x + 100 = (x + \frac{1}{2})^2 + \frac{399}{4}$ whereas $[x]^2 + [x] + 100$ differs from $([x] + \frac{1}{2})^2 + \frac{399}{4}$, as illustrated in the following example.

Example 1. At $[x] = [0, 1]$, $[x]^2 + [x] + 100 = [0, 1]^2 + [0, 1] + 100 = [100, 102]$, and $([0, 1] + \frac{1}{2}) + \frac{399}{4} = [\frac{1597}{16}, \frac{1621}{16}]$. The first result is a pessimistic approximation of the image set of $x^2 + x + 100$ at $[0, 1]$ whereas the second (and better one) is equal to this image set. \square

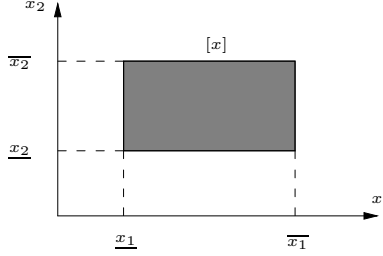


Fig. 1. A two-dimensional box.

Algebraic structures and topological structures have been developed by Moore in the 60's [3]. Note that there are areas of interval mathematics and fuzzy set theory that have been developed in parallel, in particular the extension principle of Zadeh [4].

A box or *vector* of intervals $[x]$ (see *e.g.* Fig. 1) is the cartesian product of n intervals and is noted $[x_1] \times \dots \times [x_n]$, with $[x_i] = [\underline{x}_i, \overline{x}_i]$, $i = 1, \dots, n$. The i th interval component $[x_i]$ is the projection of $[x]$ onto the i th axis. Many of the notions introduced for intervals extend without difficulty to boxes. The set of all boxes of \mathbb{R}^n is denoted \mathbb{IR}^n . Note that $]-\infty, \infty[\times \dots \times]-\infty, \infty[$ is an element of \mathbb{IR}^n .

Consider the situation where we have a model which acts as a function f , mapping (inputs) x to (outputs) y . This model f might be quite complex, with multiple input parameters and with different kinds of uncertainty represented on them: information available on inputs may be rich or sparse, so-called “aleatory” and may be made known through objective measurements. Mathematically inputs might be represented by probability or possibility distributions, by strong or sparse collections of data points, by simple intervals, or even by non-quantified linguistic expressions.

Given f , how can we propagate the uncertainty on x to y through f ? Moreover, how can we do so in a way which respects all the original uncertainty quantifications as provided, making no unnecessary assumptions? How can we do such in a way which uses *only*, but *all of* what we are given?

In this paper, we propose an approach for solving such problems: we assume that experimental points (both inputs and outputs) are modeled as intervals and provide exact solutions. The contribution of this article is three-fold: firstly, it is shown that interval analysis can directly be applied to perform optimization, yielding to close form expressions of results. Secondly, a C++ library which implements the techniques is introduced. The program and the library are open-source may be obtained from the author on simple request. Third, short illustrative examples including a nonlinear dynamic process and a blind source separation problem are given which show the significative improvement of the approach in comparison with standard linear identification.

2 A refresher on parameter estimation

Consider some function $f : \mathcal{X} \subset \mathbb{R}^n \rightarrow \mathcal{Y} \subset \mathbb{R}^s$, where \mathcal{Y} is a bounded set and suppose we wish to construct a model $g : \mathbb{X} \subset \mathcal{X} \rightarrow \mathbb{Y} \subset \mathcal{Y}$, where \mathbb{X} and \mathbb{Y} are some domains

of interest, by choosing a parameter vector $\mathbf{p} \in \mathbb{R}^p$ so that, mathematically speaking,

$$\mathbf{y} = f(\mathbf{x}) = g(\mathbf{x}, \mathbf{p}) + e(\mathbf{p}), \quad (1)$$

for all $\mathbf{x} \in \mathbb{X}$ and $\mathbf{y} \in \mathbb{Y}$, where the error in approximation, $e(\mathbf{p})$, is as small as possible.

We suppose that all that is available to choose the parameters \mathbf{p} in g is some part of the *unknown* function f in the form of the input-output data pair associations. The i th input-output pair for the system f is denoted by $(\mathbf{x}_i, \mathbf{y}_i)$, where $\mathbf{x}_i \in \mathbb{X}$, $\mathbf{y}_i \in \mathbb{Y}$ and $\mathbf{y}_i = f(\mathbf{x}_i)$. This may correspond for instance to $n + s$ scalar measurements corresponding to various experimental conditions on a static process or on a dynamical one. The row vector $\mathbf{z}_i = (\mathbf{x}_i^T \mathbf{y}_i^T) \in \mathbb{R}^{n+s}$ denotes one particular data sample. Stacking N consecutive samples on top of each other gives the data matrix

$$\mathbf{Z} = \begin{bmatrix} \mathbf{z}_1 \\ \mathbf{z}_2 \\ \vdots \\ \mathbf{z}_N \end{bmatrix} \in \mathbb{X} \times \mathbb{Y}. \quad (2)$$

The purpose of parameter estimation is, for instance, to find \mathbf{p} such that $g(\mathbf{x}, \mathbf{p})$ best fits \mathbf{y} in a sense to be specified. In [5], the parameters are considered admissible if the error $e(\mathbf{p})$ belongs to some prior compact set of admissible error $\mathbb{E} \subset \mathbb{R}^s$. For instance, \mathbb{E} may be the box defined as

$$\mathbb{E} = \{e | e^- \leq e \leq e^+\}, \quad (3)$$

where e^+ and e^- are some prior bounds. One is then interested in finding the set \mathbb{S} of all values of \mathbf{p} such that the error is admissible, i.e. $\mathbb{S} = \{\mathbf{p} | e(\mathbf{p}) \in \mathbb{E}\}$. This set has been called *membership set*, *likelihood set* and *posterior feasible set*. If the data were generated by a statistical model $g(\mathbf{x}, \mathbf{p}^*)$, where \mathbf{p}^* is some true value of the parameters and if $e(\mathbf{p}^*) \in \mathbb{E}$, then \mathbb{S} contains \mathbf{p}^* . Thus, \mathbb{S} provides an accurate description of the uncertainty with which \mathbf{p}^* is estimated [5].

If the reciprocal of g exists and is denoted g^{-1} , \mathbb{S} is defined as $\mathbb{S} = g^{-1}(\mathbf{y} - \mathbb{E}) = g^{-1}(\mathbb{Y})$, where $\mathbb{Y} = \mathbf{y} - \mathbb{E}$ is the *measurement set*. In other words, for any $\mathbf{p} \in \mathbb{S}$ there exists $e \in \mathbb{E}$ such that $\mathbf{y} = g(\mathbf{x}, \mathbf{p}) + e$.

System identification (i.e. function approximation) amounts to adjusting \mathbf{p} using information from \mathbf{Z} so that $g(\mathbf{x}, \mathbf{p}) \approx f(\mathbf{x})$, $\forall \mathbf{x} \in \mathbb{X}$. Measured and model outputs never match perfectly in practice, but differ as $e(\mathbf{p})$. An obvious modeling goal must be that this discrepancy is “small” in some sense that is achieved by the value of the *approximation error* we wish to bound. Such a bound is for example

$$\sup_{\mathbf{x} \in \mathbb{X}} \|f(\mathbf{x}) - g(\mathbf{x}, \mathbf{p})\|, \quad (4)$$

which requires that f is known everywhere. The problem is that we only know the part of f given by \mathbf{Z} , and it is the only evaluation we can make based on known information. When g is not affine in \mathbf{p} , one may linearize it around some value of the parameters estimated beforehand and then use any method for linear models. Scanning the parametric space using random search offers no guarantee as to the global nature of the obtained results [6].

A thorough examination of $e(\mathbf{p})$ reveals that *imprecision* play a key role in the parameter value formulation: indeed, it is reasonable to assume that a large part of estimation deviations of \mathbf{p} may come from imperfection in input assessment, together with the model imprecision.

3 Identification method

3.1 Core of the modeling

The structure of the model, *i.e.* the form of g and the matrix \mathbf{Z} , are determined by the statistician on the basis of prior knowledge and/or by comparing several candidate models in terms of the prediction error. Once a model is chosen, the parameters of the model may be estimated. Modeling possible values of variables by means of real intervals accounts for some imprecision.

Let us consider a system with imprecise input and output x and y resp., which are readings from unreliable sensor (“noisy data”). To simplify exposition, we shall only consider output errors :

$$e_i(\mathbf{p}) = y_i - g(\mathbf{x}_i; \mathbf{p}), \quad i = 1, \dots, s \quad (5)$$

but other types of errors could be considered as well. We assume that these errors should satisfy $\underline{e}_i \leq e_i(\mathbf{p}) \leq \bar{e}_i$, $i = 1, \dots, s$ to be “admissible”, where \underline{e}_i and \bar{e}_i are known lower and upper prior bounds of the approximation error resulting from technical specifications or pointing out how far we can go in accepting discrepancies between our data and model outputs. Note that $e_i(\mathbf{p})$ is an interval.

Let \mathbf{y} be the vector of all data $y_i, i = 1, \dots, s$ collected on a given system, and let $\mathbf{g}(\mathbf{x}; \mathbf{p})$ be the vector of all corresponding model outputs $g(\mathbf{x}_i; \mathbf{p}), i = 1, \dots, s$. Equation (5) can then be rewritten as

$$\mathbf{e}(\mathbf{p}) = \mathbf{y} - \mathbf{g}(\mathbf{x}; \mathbf{p}), \quad (6)$$

Picked in intervals, equation (6) gives:

$$[\mathbf{e}(\mathbf{p})] = [\mathbf{y}] - \mathbf{g}([\mathbf{x}]; [\mathbf{p}]), \quad (7)$$

The model (7) is admissible if \mathbf{p} is such that $\mathbf{e} \in \mathbb{E}$, *i.e.* the set of all error vectors \mathbf{e} satisfies $\underline{\mathbf{e}} \leq \mathbf{e} \leq \bar{\mathbf{e}}$ where $\underline{\mathbf{e}}$ and $\bar{\mathbf{e}}$ are known. Estimating interval $[\mathbf{p}]$ consists on looking for the set of all admissible values of P^α that are *consistent* with (7), *i.e.* errors should satisfy

$$[\mathbf{y}] - \mathbf{g}([\mathbf{x}]; [\mathbf{p}]) \in [\mathbf{e}](\mathbf{p}). \quad (8)$$

Equation (8) relies on the assumption that there exists an *inclusion function* $[g]$ of g that returns an *enveloping box* guaranteed to contain the image by g of any given box $[\mathbf{x}]$ included in the domain of g . $[g]([\mathbf{x}])$ is a box such that

$$\mathbf{g}([\mathbf{x}]) \subset [g]([\mathbf{x}]), \quad (9)$$

The including function $[g]$ is easy to compute for usual elementary functions that can be obtained by composition of elementary operations such as $+$, $-$, \times , $/$, \exp , \tan , \sin , \dots by replacing each of these elementary operations by their inclusion function in the formal expression of g .

When no efficient algorithm exists for its computation, $[g]$ can be approximated by an inclusion function $\mathbb{G} : \mathbb{R}^n \rightarrow \mathbb{R}^p$ satisfying the equation (9) and such that:

$$w([\mathbf{x}]) \rightarrow 0 \Rightarrow w([\mathbb{G}([\mathbf{x}])]) \rightarrow 0. \quad (10)$$

where $w([\mathbf{x}])$ is the width of the box $[\mathbf{x}]$, defined as the length of its largest side(s). It must be noted that the algorithm presented in this paper is guaranteed to converge only if (10) is valid.

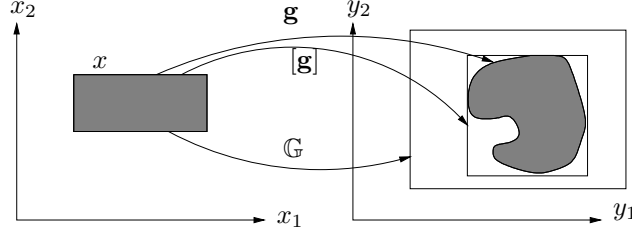


Fig. 2. Minimal inclusion function $[g]$ and inclusion function G of a function g .

Figure 2 illustrates conditions (9) and (10).

Let $\hat{\mathbb{S}} = \{\mathbf{p} | g([\mathbf{x}]; [\mathbf{p}]) \in \mathbb{Y}\}$. Then

$$g([\mathbf{x}]; [\mathbf{p}]) \in \mathbb{Y} \Leftrightarrow \mathbf{p} \in g^{-1}(\mathbb{Y}) \quad \text{and} \quad \mathbf{p} \in [\mathbf{p}](0) \quad (11)$$

$$\Leftrightarrow \mathbf{p} \in [\mathbf{p}](0) \cap g^{-1}(\mathbb{Y}) \quad (12)$$

where $[\mathbf{p}](0)$ is the search domain. Thus

$$\hat{\mathbb{S}} = [\mathbf{p}](0) \cap g^{-1}(\mathbb{Y}), \quad (13)$$

and characterizing $\hat{\mathbb{S}}$ is a *set inversion problem*.

We shall say that $[\mathbf{p}]$ is *feasible* if $[\mathbf{p}] \subset \mathbb{S}$, *unfeasible* if $[\mathbf{p}] \cap \mathbb{S} = \emptyset$, else $[\mathbf{p}]$ is *ambiguous*.

To perform computation of \mathbb{S} in an approximate but guaranteed way, an interval analysis algorithm is applied to compute the possible interval range $[\mathbf{p}]$ of \mathbf{p} . The interval computation stage is detailed in the following.

3.2 Set characterization

Methods allowing to implement interval analysis are relatively few and date from the nineties and among them, one may quote the Moore's algorithm [3] and SIVIA proposed by Jaulin [7]. Most of the methods for estimating parameters are based on computations performed at point values of the parameter vector. The main interest in the notion of paving is to replace point values by subsets of the parameter space. For simplicity, we will use pavings based upon boxes.

A *paving* of a compact subset $\{\mathbb{P}\} \subset \mathbb{R}^n$ is a set of non overlapping boxes with nonzero width such that the union of these boxes corresponds to $\{\mathbb{P}\}$. A subpaving \mathbb{K} of \mathbb{P} is a subset of \mathbb{P} . Upon completion, the algorithm encloses \mathbb{S} between two compact sets corresponding to 2 subpavings (Fig. 4). Let \mathbb{E} be the feasible error set. Initialisation is performed by setting $\mathbb{Y} = \mathbf{y} - \mathbb{E}$. The principle is as follows:

- Define an initial box of interest $[\mathbf{p}](0)$ within which the search will be performed
- Compute a paving $\{\mathbb{P}\}$ of $[\mathbf{p}](0)$
- Compute the image $g([\mathbf{x}]; [\mathbf{p}])$ for each box of this paving. Three situations must then be considered (see figure 3).

$$\begin{cases} [g]([\mathbf{x}]; [\mathbf{p}]) \subset [\mathbb{Y}] \Rightarrow [\mathbf{p}] \subset \mathbb{S} & \text{so that } [\mathbf{p}] \text{ is feasible} \\ [g]([\mathbf{x}]; [\mathbf{p}]) \cap [\mathbb{Y}] = \emptyset \Rightarrow [\mathbf{p}] \cap \mathbb{S} = \emptyset & \text{so that } [\mathbf{p}] \\ & \text{is unfeasible} \\ \text{otherwise, } [\mathbf{p}] & \text{is indetermined.} \end{cases}$$

The exploration algorithm performs a recursive implementation of the principle that has just been described: a *bisection algorithm* splits each box of the subpaving into smaller boxes whenever needed until the width of the box becomes smaller than some tolerance parameter ϵ to be specified by the user. Cutting is carried out again as long as the boxes contain solutions or stops if the boxes do not contain any.

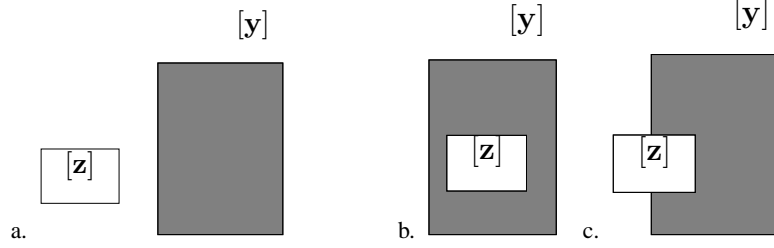


Fig. 3. Feasibility of boxes: a test function allows to distinguish the cases (a-c) represented in this figure. Let $[z]$ and $[y]$ be two boxes. Suppose that $[n] = [z] \cap [y]$. a) $[n] = \emptyset$ means that $[z]$ has an empty intersection with $[y]$. b) $[n] = [z]$, so that $[z]$ is included inside $[y]$. c) $[n] \subset [z]$ and $[n] \neq [z]$. $[z]$ intersects $[y]$, so that we cannot conclude. The box $[z]$ can be split again.

We shall split \mathbb{P} iteratively into three subpavings $\underline{\mathbb{S}}$, \mathbb{S} and $\overline{\mathbb{S}}$ corresponding to the sets of all feasible, unfeasible and indetermined boxes respectively, as plotted in figure 4. These subpavings satisfy the following relations :

1. $\{\underline{\mathbb{S}}\} \subset \mathbb{S} \subset \{\underline{\mathbb{S}}\} \cup \{\overline{\mathbb{S}}\}$
2. $\text{vol}(\{\underline{\mathbb{S}}\}) \leq \text{vol}(\mathbb{S}) \leq \text{vol}(\{\underline{\mathbb{S}}\}) + \text{vol}(\{\overline{\mathbb{S}}\})$
3. $\{\{\underline{\mathbb{S}}\}\} \subset \{\mathbb{S}\} \subset \{\{\underline{\mathbb{S}}\}\} \cup \{\{\overline{\mathbb{S}}\}\}$

Provided that \mathbb{S} is full, this means that the pair $\{\underline{\mathbb{S}}; \overline{\mathbb{S}}\}$ defines a neighbourhood $\partial\mathbb{S} \triangleq \overline{\mathbb{S}} \setminus \underline{\mathbb{S}}$ of \mathbb{S} with a diameter that can be chosen arbitrarily small.

The previous algorithm makes an extensive use of a *stack* L of boxes, *i.e.* a dynamical structure on which only 3 operations are possible: at any time, one may put an element on top of the list, remove the top element or test the stack for emptiness. We define the required accuracy ϵ for the paving \mathbb{P} as the maximum width that an indetermined box can have. In the following, the principal plane of a box is the symmetry plane of this box that is orthogonal to the axis $i \in \{j | w([p]) = w([p_j])\}$, where the operator “width” $w([\cdot])$ of a box is the length of its largest side.

Let $[p](0)$ be the box considered at iteration k . Initialisation is performed by setting $k = 0$, $L = \emptyset$, $\underline{\mathbb{S}} = \overline{\mathbb{S}} = \emptyset$. The recursive algorithm can be described as follows:

The union of all the boxes in the list L returned by the program contains \mathbb{S} ; the partition \mathbb{P} consisting of feasible, unfeasible and indetermined boxes can be plotted in the parameter space in the case the space dimension is less than 4 (see fig. 4).

4 Discussion

The method of *set characterization* introduced in section 3.2 appeals to some comments.

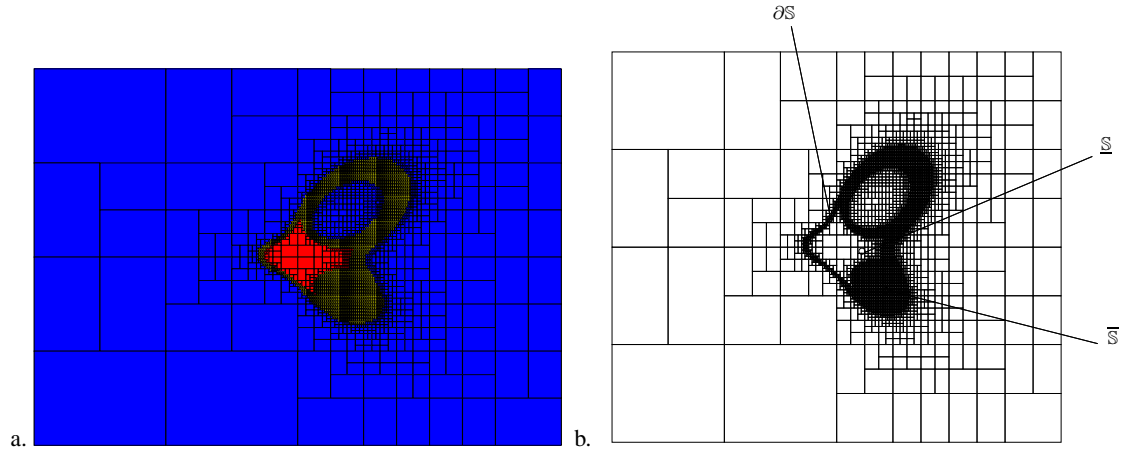


Fig. 4. a. Regular paving of a box: accepted, rejected and indetermined subpavings are respectively coloured in red, blue and yellow. b. $\{\underline{S}$ and $\bar{S}\}$ brackets the portion of S contained in $[p](0)$.

Algorithm BISECT
<p>INPUTS data: \mathbf{y} inclusion function: $[g](\cdot)$ feasible error set: \mathbb{E} prior feasible box: $[p](0)$ accuracy for the paving: ϵ</p> <p>INITIALIZATION $\mathbb{Y} = \mathbf{y} - \mathbb{E}$; stack: $L = \emptyset$ iteration: $k = 0$ $[p] = [p](0)$;</p> <p>ITERATION k Begin</p> <p>step 1: if $[g]([x]; [p](k)) \subset [\mathbb{Y}]$ then $\underline{S} := \underline{S} \cup [p]$ and $\bar{S} := \bar{S} \cup [p]$; step 2: else if $[g]([x]; [p](k)) \cap [\mathbb{Y}] = \emptyset$, then $\underline{S} := \underline{S} \cup [p]$; then unstack $[p](k)$ as unfeasible; step 3: else if $w([p](k)) \leq \epsilon$, then $\bar{S} := \bar{S} \cup [p](k)$; else cut $[p](k)$ along the principal plane and stack the resulting boxes in L. step 4: if the stack is not empty, then unstack and store the resulting box in $[p](k+1)$; $k = k + 1$; go to step 1; End</p>

Table 1. Recursive implementation of the bisection algorithm.

Upon completion, this approach encompasses *all* the acceptable values of the parameter vector in a set that is fully characterized by BISECT: $\underline{\mathbb{S}}$ and $\overline{\mathbb{S}}$ will tend to \mathbb{S} from inside and outside when $\epsilon \rightarrow 0$. Since \mathbb{S} is a finite union of boxes guaranteed to contain the portion of \mathbb{S} of interest, it is very convenient for implementing set-theoretic manipulations [8,5].

The advantages of this approach are threefold:

- (i) no assumption is made on the image function g ,
- (ii) no statistical assumption on the modeling error is required,
- (iii) any bounded error can be treated independently from its origin (modeling and/or measurement error).

An other advantage of the proposed approach is that the input-output roles of the variables x and y can be reversed since the linking function $\mathbf{g} : x \rightarrow y$ can be run *forward* as well as *backward* when using interval analysis. Subpavings form a useful class of objects for manipulating statistical estimations.

The algorithm requires a possibly very large search box $[\mathbf{p}](0)$ to which \mathbb{S} is guaranteed to belong. Solvers split the search box into an union of boxes (the *paving*) with guaranteed error bounds (*i.e.* mathematically valid) [6] (see section 3.2). The paving is built by the solver itself. A computer program can represent a set of (eventually disjoint) intervals as a *list* L . The precision of the solver is controlled by coefficients specifying, for example, the width ϵ of the smallest boxes of the paving, or the accuracy in the localization of a global optimum. The computing time of the solver can increase quickly with the dimension and size of the list L .

Special care must be taken in avoiding memorizing unnecessary information, otherwise the quantity of memory required to store the paving of \mathbb{S} will increase linearly at each iteration, which may result into a memory overflow even for problems of modest dimension.

One can observe that the parameter space is not isotropic because the sensitivities of g with respect to the various components of \mathbf{p} are not of the same order of magnitude. The basic bisection technique suggested in Tab. 1 may not be efficient enough. The problem is then how to choose the fastest bisection policy that results in a convergence as fast as possible. Jaulin *et al.* [5] suggests the bisections of $[\mathbf{p}]$ into boxes $[\mathbf{p}_1]$ and $[\mathbf{p}_2]$ that minimize $\text{vol}(g([\mathbf{x}]; [\mathbf{p}_1])) + \text{vol}(g([\mathbf{x}]; [\mathbf{p}_2]))$. If $[\mathbf{p}]$ is not ambiguous, this policy will tend to avoid classifying as indetermined. Experiments seem to show that this can improve the efficiency of the solver when the anisotropy is severe.

5 Application to statistical problems

5.1 Parameter estimation

In this example, a simplified version of the problem explored by Jaulin and Walter³ [5] is given. The vector comprising all available data \mathbf{y} is:

$$\mathbf{y} = (1.59, 1.44, 1.30, 1.18, 1.07, 0.96, 0.87, 0.79, 0.71, 0.64)^T.$$

The numerical values of the corresponding interval data are given in Tab. 5.

The set $\hat{\mathbb{P}}$ to be characterized consists of the set of variable vectors $([p_1], [p_2])^T$ such that the graph of the function :

$$f(\mathbf{p}, t) = p_2 \exp(-p_1 t), \quad (14)$$

³ The extension of the method to multiple-output problems is straightforward.

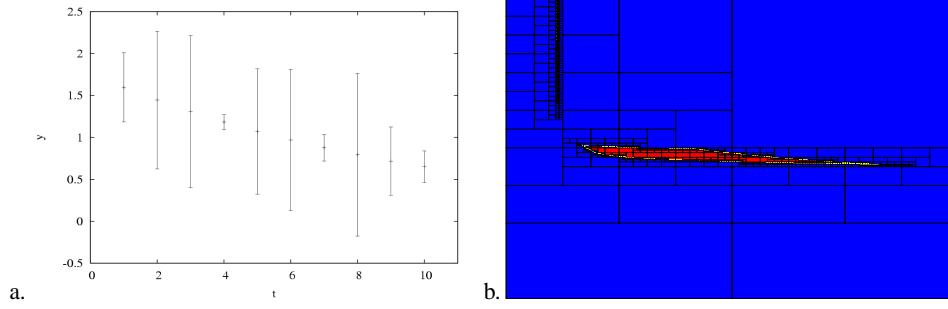


Fig. 5. a) Measurement times and corresponding interval data. b) Top-view of the paving generated by BISECT to bracket the solution set (in red) in the parameter space. The outer frame corresponds to the box $[-1, 5] \times [-5, 5]$.

crosses all the data bars of Figure 5. In this simulated example, the $[y_i]$ were computed by adding a random error interval with radius $\rho_i = 0.5|y_i| + 1$ to the y_i . The initial box domains for the parameters p_1 and p_2 may be arbitrarily large, for example

$$[p_1] = [-10, 000; 10, 000] \text{ and } [p_2] = [-10, 000; 10, 000], \quad (15)$$

i.e. no prior information is available on the parameters. The feasible fuzzy set for the parameters is given by (13), where the search domain $[\mathbf{p}](0)$ is $[-1, 5] \times [-5, 5]$. In less than 1s, on a PENTIUM IV, BISECT generates the pavings of figure 5, thus bracketing the posterior feasible set for \mathbf{p}_α between the inner and the outer approximations. Figures 5 gives a top-view representation of \mathbf{p} , *i.e.* the bounded support of the variables p_1 and p_2 that is consistent with the equation (14) and the domains (15).

5.2 Curve estimation

The curvature $\kappa(t)$ of an arbitrary twisted curve \mathcal{C} measures the rate of change of the tangent when moving along the curve. It measures, so to speak, the deviation of the curve from a straight line in the neighbourhood of any of its points. It is quite easy to derive an analytic expression of the curvature which is valid when \mathcal{C} is represented by an allowable parametric representation $\mathbf{x}(t)$:

$$\kappa(t) = \frac{|\mathbf{x}' \times \mathbf{x}''|}{|\mathbf{x}'|^3}, \quad (16)$$

where \times denotes the vector product. Derivatives with respect to time are denoted by primes, *e.g.* $\mathbf{x}' = \frac{d\mathbf{x}}{dt}$ and $\mathbf{x}'' = \frac{d^2\mathbf{x}}{dt^2}$. (16) is equivalent to (see [9] and references therein):

$$\kappa(t) = \frac{\sqrt{(\mathbf{x}' \cdot \mathbf{x}')(\mathbf{x}'' \cdot \mathbf{x}'') - (\mathbf{x}' \cdot \mathbf{x}'')^2}}{(\mathbf{x}' \cdot \mathbf{x}')^{\frac{3}{2}}}, \quad (17)$$

When \mathcal{C} is a *clelia* curve (see figure 5.2.a), then it has the following cartesian representation⁴:

$$\mathbf{x}(t) = (R \cos nt \cos t, R \cos nt \sin t, R \sin nt) \quad (18)$$

⁴ *Clelia* curve was studied by Guido Grandi in 1728.

with t as parameter, it is easy to obtain from (17) a simpler representation of the curvature:

$$\kappa(t) = \frac{\sqrt{(4n^2 - n^4) \cos^2(nt) + n^6 + 4n^4 + (1 - n^2) \cos^4(nt)}}{R \sqrt[3]{n^2 + \cos^2(nt)}}. \quad (19)$$

Figure 5.2.b plots the domain of $\kappa(t)$.

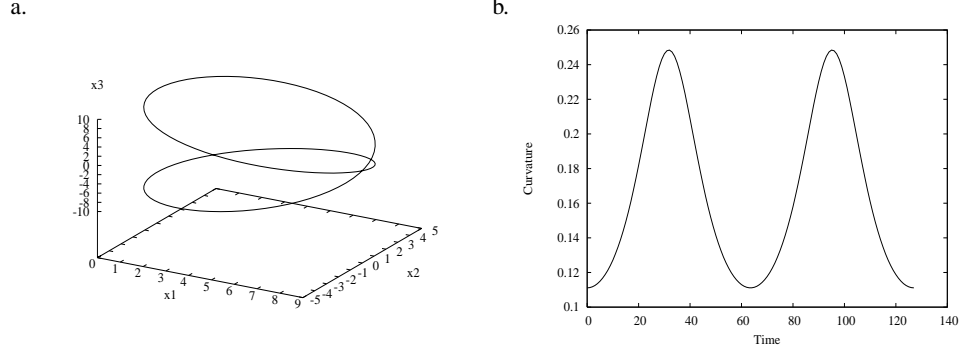


Fig. 6. a. 3-dimensional representation of the Clelia curve. b. Graph of the function $\kappa(t)$ given by equation (19).

A statistical approach assume that a set of N input-output data pairs $(\mathbf{x}_i, \kappa_i)_{i=1}^N$ is available. Recall that $\mathbf{x}_i \in \mathbb{R}^3$ are vectors and κ_i is scalar. The data set is split into a training and a validation data set. A neural network such as represented in Fig. ?? is designed for supervised learning and prediction task from a particular \mathbf{x} . Such architecture is called *multi-layer perceptron* (MLP). Hidden units are placed between the features units and the predictions units. A deterministic MLP trained by a method such as backpropagation [10], can implement any input-output function provided that the number of hidden neurons n_h is sufficiently large.

Fig. 5.2.a depicts the *training error* performed on the training data set. With increased number of hidden neurons, the accuracy of the resulting neural system with respect to the training data set is improved, but the ability of the model to generalize for inputs (test set) may be degraded as illustrated in our results (Fig. 5.2.b). In a similar way, decreasing the output error improves the accuracy with respect to the training data set, but accuracy in the presence of inputs different from the training data set is degraded. These tests illustrate some effects in parameter choice on the resulting neural model.

A second approach consists in using interval analysis tools. In this simulated example, we assume that the values taken by κ are imprecise and sampled at the rate $n_{\frac{\pi}{11}}$: the prior intervals $[\kappa_i]$ are computed by adding a centered error interval to the associated measurement κ_i . The solution set $\hat{\mathbb{S}}$ to be characterized consists of all the values of $\mathbf{p} = (R, n)^T$ such that the graph of the function $\kappa(t)$ crosses all the data bars of figure 5.2. The dataset is made of 22 data, so very few for a learning model.

For $\epsilon = 0.03$, BISECT generates the subpaving represented in Figures 9 in 7s on a PENTIUM IV, consistent with the equation (19). The prior box for the parameters is taken as $[\mathbf{p}] = [-10, 10] \times [-10, 10]$.

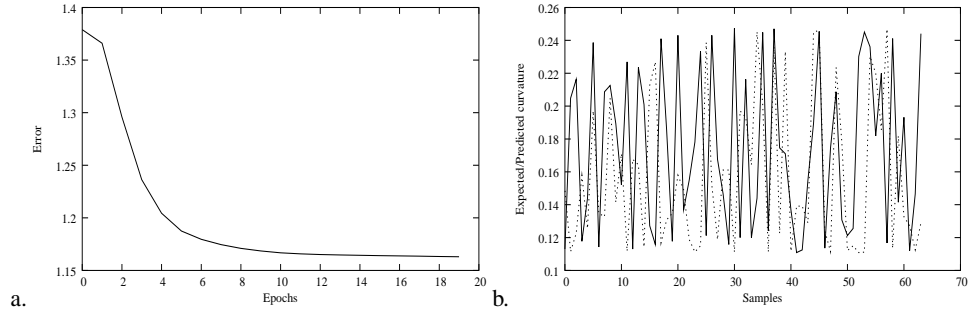


Fig. 7. a. Measured outputs (–) and simulated outputs (..) computed by the neural network.

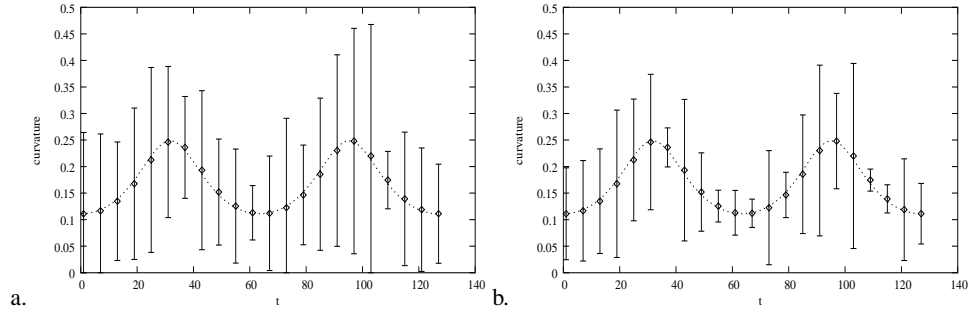


Fig. 8. a) Experimental data (♦) together with their uncertainty intervals and graph of the function $\kappa(t)$ (see Eq. (19)). b) Posterior feasible intervals for the κ_i superposed on the graph of the function $\kappa(t)$.

More interesting, one may want to generate the posterior feasible set for the κ_i (posterior estimates) from the formula 19 as can be seen in figure 5.2.

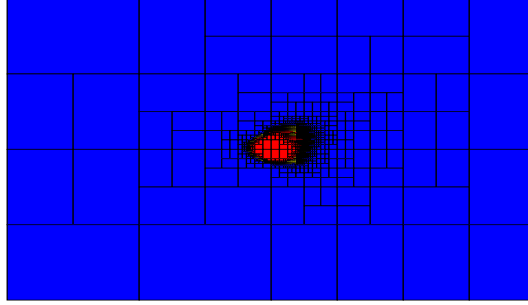


Fig. 9. Top-view of the paving generated by BISECT to bracket the solution set (in red) in the parameter space. The outer frame corresponds to the box $[-10, 10] \times [-5, 5]$.

5.3 Blind source separation

The problem addressed here is the recovery of n unknown independent sources $s_i(t)$ from the observation of a n linear mixtures x_i . In matrix and vector notations, this model reads $\mathbf{x} = A\mathbf{s}$, where $\mathbf{s} = (s_1, s_2)^T$ is the vector of sources, $\mathbf{x} = (x_1, x_2)^T$ the vector of observations and $A = \{a_{ij}\}$ the mixture matrix. It is a noise-free, time-free model. To recover a vector \mathbf{y} close to the source vector \mathbf{s} knowing the vector \mathbf{x} only, one should estimate some inverse of A , denoted as B . The corresponding estimate of \mathbf{s} is $\mathbf{y} = B\mathbf{x}$. It should be noted, however that the matrix A (or its inverse) is not identifiable from the observations (see e.g.[11]): even if we can extract n independent components, we do not know their ordering. This implies that there exists a *freedom of permutations* of the original signals. The magnitudes of the original signals s_i are also not recoverable, because a scalar multiple of s_i , cs_i of s_i by a constant c cannot be distinguished from multiplication of the i th column of A by the same constant c . Therefore, therefore, we can recover only a permuted and rescaled version of the sources, *i.e.* we can obtain at best PDA^{-1} , where P is a permutation matrix and D is nonsingular scaling matrix.

The important question is: does the independence of the components of \mathbf{y} imply necessarily the *separability* of the mixing model? The answer to this question is positive in the linear instantaneous domain: the transformation which maps a non-Gaussian random vector with independent components to a random vector with independent components is unique, up to some trivial transformation. This property is a direct result of the Darmois-Skitovich theorem [12]. For sake of simplicity, the problem here is limited to the case of 2 sources s_1 and s_2 observable through 2 signals x_1 and x_2 (omitting time index t). A solution to the problem was first addressed by Herault and Jutten [13]. They proposed a fully interconnected network (see figure 10) which computes

$$\begin{cases} y_1 &= x_1 + b_{12}y_2 \\ y_2 &= x_2 + b_{21}y_1. \end{cases} \quad (20)$$

where b_{12} and b_{21} are adaptive weights adjusted by means of an adaptation law based on the product of 2 nonlinear functions f and g . The sources are assumed to be zero-mean (i.e. $\mathbb{E}[s_1] = \mathbb{E}[s_2] = 0$), stationary and independent.

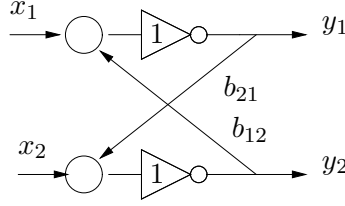


Fig. 10. The two-channels Herault-Jutten network.

The independence of the signals means that these sources must have zero covariance, i.e. $\mathbb{E}[s_1 s_2] = 0$ and therefore $\mathbb{E}[y_1 y_2] = 0$. The following adaptation rule is considered in [13] to learn the coefficients b_{ij}

$$\frac{db_{12}}{dt} = \mu y_1^3 y_2, \quad \frac{db_{21}}{dt} = \mu y_2^3 y_1, \quad (21)$$

where μ is a positive constant. It is well-known that the equilibrium points are solutions of

$$\begin{cases} \mathbb{E}[y_1^3 y_2] = 0 \\ \mathbb{E}[y_1 y_2^3] = 0 \end{cases} \quad (22)$$

However this does not *guarantee* that $\begin{pmatrix} 1 & b_{12} \\ b_{21} & 1 \end{pmatrix}$ converges to the inverse of the mixing equation, even locally.

Example 2. Let $\begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}$ be the (unknown) separating matrix. The solutions of the following system:

$$\begin{pmatrix} 1 & b_{12} \\ b_{21} & 1 \end{pmatrix}^{-1} \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \text{ or } \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}. \quad (23)$$

reads $b_{12} = \frac{a_{12}}{a_{22}}$ and $b_{22} = \frac{a_{21}}{a_{11}}$, or $b_{12} = \frac{a_{11}}{a_{21}}$ and $b_{22} = \frac{a_{22}}{a_{12}}$. \square

Comon *et al.* [14] and Sorouchyari [15] investigate the convergence properties of the algorithm and perform a stability analysis for such a network. They demonstrate that there are *exactly* 4 paired equilibrium points (see figure 12): indeed, if the point (b_{12}^*, b_{21}^*) is a equilibrium point, then the point $(\frac{1}{b_{12}^*}, \frac{1}{b_{21}^*})$ is also a solution (see e.g. [15]). One can show that such a pair of equilibrium points is on a line passing through the origine of the (b_{12}, b_{21}) plane. But only one of these stationary points will be a *stable separating solution* [16]. A stationary point (b_{12}^*, b_{21}^*) verifies (22) and is said to be (*locally asymptotically stable*) if the eigenvalues of matrix Γ , defined as the partial derivatives of (22)

$$\Gamma = \frac{-\mu}{1 - b_{12} b_{21}} \begin{pmatrix} 3\mathbb{E}[y_1^2 y_2^2] & \mathbb{E}[y_1^4] \\ \mathbb{E}[y_2^4] & 3\mathbb{E}[y_1^2 y_2^2] \end{pmatrix} \quad (24)$$

have negative real parts. *Stability conditions* are therefore

$$9(\mathbb{E}[y_1^2 y_2^2])^2 > \mathbb{E}[y_2^4] \mathbb{E}[y_1^4] \quad (25)$$

$$b_{12} b_{21} < 1 \quad (26)$$

Condition (25) follows from the requirements that the loop again of the network must be less than zero and depends only of the statistics of the input sources, whereas the condition (26) is of greater significance: it indicates that the source signals require to have negative kurtosis for the separating point to be stable when the cubic and identity functions are used.

The set \mathbb{S} to be characterized consists of all the values of \mathbf{b} such that the graphs of $B\mathbf{x}$ coincides with the graphs of \mathbf{s} . \mathbb{S} can be defined in vector form as

$$\mathbb{S} \triangleq \{\mathbf{b} \in \mathbb{R}^2 | \mathbf{f}(\mathbf{b}) \in \mathbf{x}\}, \quad (27)$$

where $\mathbf{f}(\cdot)$ is the vector function

$$\begin{pmatrix} 9(E[y_1^2 y_2^2])^2 - E[y_2^4] \mathbb{E}[y_1^4] \\ 1 - b_{12} b_{21} \\ E[y_1^3 y_2^2] \\ E[y_1 y_2^3] \end{pmatrix}, \quad (28)$$

and $\mathbf{x} =]0; +\infty[^2 \times]-\epsilon, \epsilon[^2 \times [\tilde{b}_1] \times [\tilde{b}_2]$ for any small $\gamma > 0$. SIVIA can be used to characterise \mathbb{S} provided that the following test

$$(\mathbf{b} \in [\mathbf{b}](0)) \wedge (\mathbf{f}(\mathbf{b}) \in \mathbf{x}) \quad (29)$$

is true, where $[\mathbf{b}](0)$ is the prior search box. The considered problem is then much more general than the HJ-learning rule.

As an illustration, consider the discrete time model in which the data have been generated by simulating for $k = 1, \dots, 500$: $\begin{pmatrix} x_1(k) \\ x_2(k) \end{pmatrix} = \begin{pmatrix} 1 & 0 & 6 \\ 0 & 3 & 1 \end{pmatrix} \begin{pmatrix} s_1(k) \\ s_2(k) \end{pmatrix}$, where $s_1(k) = \sin(7, 3kT_e)$, $s_2(k) = \sin(4kT_e)$; $T_e = 0, 2$. Prior interval $([x_1(k)], [x_2(k)])$ are obtained by adding a random white noise $n(k) \sim \mathcal{U}_{[-\epsilon, \epsilon]}$, with $\epsilon = 0.01$. Fig. 11.a plots the sources and Fig. 11.b the observed data.

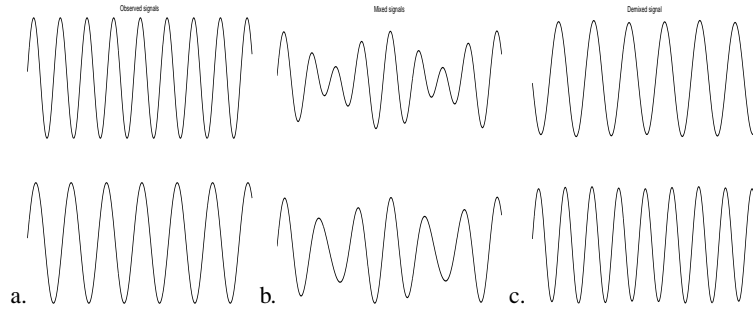


Fig. 11. a) The unknown sources b) the observations c) The separated sources.

For $\epsilon_r = 0.01$, the algorithm generates in the (b_{12}, b_{21}) plane (Fig. 12.a) the subpaving of the parameters satisfying the independence constraint (22) in less than 1s on a PEN-TIUM III: the similarity with Fig. 12 is relevant.

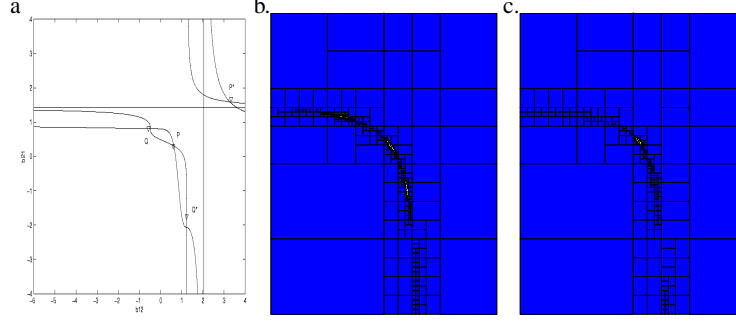


Fig. 12. a) Theoretical equilibrium points of the system b) stable equilibrium state using (25) c) separating state concatenating (25) and (26) in the (b_{12}, b_{21}) plane.

The yellow boxes have been proven to be included in \mathcal{S} and the blue grey boxes have been proved to have an empty intersection with \mathcal{S} . No conclusion has been reached for the red boxes. After completion, the contracted intervals in the figures 12.b and c include the true values of the parameters when constraints (25) and (26) respectively are satisfied. Large portion of the prior parameter space can be quickly eliminated, before concentrating on the indeterminate region. The required memory for the stack remains surprisingly limited. As could be expected, the number of boxes in the subpavings increases quickly when the number of parameters increases or when the tolerance parameter ϵ_r is decreased.

6 Conclusion

The problem of parameter estimation of a (non)linear model from prior knowledge, experimental data and collateral constraints is viewed in this article as one of *set inversion*, which is solved in an approximate but guaranteed way with the tools of interval analysis. It is possible to characterize the set of all parameter vectors that are *consistent* with the data in the sense that the errors between the data and corresponding model outputs fall within known prior bounds. This has been illustrated on simple simulated examples for time-invariant models whose outputs are linear in their inputs, even if nonlinear in their parameters. It is worth stressing that the scheme of section 3.2 is robust in the sense that the fit of the tuned model is at least as good as what is obtained with the classical optimization approach.

Upon completion of the algorithm, a paving bracketing the contours of the solution membership functions is found (or not) with a precision controlled by the solver.

This computation process has drawbacks: (i) its complexity is exponential in the number of parameters which restricts its use to low-dimensional problems, (ii) the algorithm presented here is far from optimal from the viewpoint of computational time and significant improvements can be expected in the near future, (iii) efficient functions are needed which are available only when an explicit solution for the equations defining the model can be found.

Realistic advantages can be found compared to the statistical approach:

1. The error structure is quite simple and similar information is usually available in most practical cases, not assuming *any a priori* statistical information about the error.

2. The computation of the parameter domain is conceptually simple and is practically feasible even if the number of data is not large.
3. The algorithm is *deterministic*.
4. *guaranteed* results are available for (non)linear models even when the parameters are not identifiable. Nonlinear constraints are easily handled.
5. the solver characteristics are different from optimization approaches that requires a (large) set of data points.

Least-Square estimation suffers from the fact that the cost function to be minimized is a sum of terms involving the same parameters, so multioccurrence of these parameters is unavoidable and tends to make inclusion functions for the cost function very pessimistic, which complicates the elimination of interesting parts of the search domain.

References

1. L. Jaulin, M. Kieffer, O. Didrit, and E. Walter. *Applied interval analysis*. Springer, Paris, 2001.
2. V. Vigneron, M. Olteanu, and H. Maaref. Identification of fuzzy functions via interval analysis. *Fuzzy Sets and Systems*, 2007. to be published.
3. R.E. Moore. *Interval analysis*. Prentice Hall, New York, 1966.
4. P. Loslever and S. Bouilland. Marriage of fuzzy sets and multiple correspondence analysis. examples with subjective interval data and biomedical signals. *Fuzzy Sets and Systems*, 107:255–275, 1999.
5. L. Jaulin and E. Walter. Guaranteed nonlinear parameter estimation from bound-error data via interval analysis. *Mathematics and Computers in Simulation*, 35(2):123–137, 1993.
6. E. Walter, editor. *Special issue on parameter identification with error bounds*, volume 32, 1990.
7. L. Jaulin. *Solution globale et garantie de problèmes ensemblistes. Application à l'estimation non-linéaire et à la commande robuste*. Phd thesis, Université Paris-Sud, Orsay, France, 1994.
8. J.P. Norton. Special issue on bounded-error estimation: Issue 1. In *International Symposium Control and Signal Processing*, pages 1–118, 1990.
9. E. Kreyszig. *Differential geometry*. Dover, 1991.
10. D.E. Rumelhart, G.E. Hinton, and R.J. Williams. Nature. *Learning representations by back-propagating errors*, 323:533–536, 1986.
11. V. Vigneron and J. Jutten. Bounded approximation for score fonction selection. In *Proceedings ICA'2003*, pages 119–124, April 2003.
12. G. Darmon. Analyse générale des liaisons stochastiques. *Rev. Inst. Intern. Stat.*, 21:2–8, 1953. (in french).
13. J. Herault and C. C. Jutten. Space ot time adaptative signal processing by neural network models. In *Neural networks for computing*, Snowbird, 1986.
14. O. Macchi and E. Moreau. Self-adaptive source separation, part i: Covergence analysis of a direct linear network controlled by the herault-jutten algorithm. *IEEE Trans. on Signal Processing*, 45(4), April 1997.
15. E. Sorouchyari. Blind separation of sources, part 3: stability analysis. *Signal Processing*, 24:21–30, 1991.
16. N. Minorsky. *Nonlinear oscillations*. Reinhold, New York, 1962.