# Construction of confidence intervals for neural networks based on least squares estimation

I. Rivals[*], L. Personnaz

*Laboratoire d'Électronique, École Supérieure de Physique et de Chimie Industrielles, 10 rue Vauquelin, 75231 Paris Cedex 05, France*

Received 14 October 1997; accepted 5 July 1999

## Abstract

We present the theoretical results about the construction of confidence intervals for a nonlinear regression based on least squares estimation and using the linear Taylor expansion of the nonlinear model output. We stress the assumptions on which these results are based, in order to derive an appropriate methodology for neural black-box modeling; the latter is then analyzed and illustrated on simulated and real processes. We show that the linear Taylor expansion of a nonlinear model output also gives a tool to detect the possible ill-conditioning of neural network candidates, and to estimate their performance. Finally, we show that the least squares and linear Taylor expansion based approach compares favorably with other analytic approaches, and that it is an efficient and economic alternative to the nonanalytic and computationally intensive bootstrap methods. © 2000 Elsevier Science Ltd. All rights reserved.

*Keywords*: Nonlinear regression; Neural networks; Least squares estimation; Linear Taylor expansion; Confidence intervals; Ill-conditioning detection; Model selection; Approximate leave-one-out score

## 1. Introduction

For any modeling problem, it is very important to be able to estimate the reliability of a given model. This problem has been investigated to a great extent in the framework of linear regression theory, leading to well-established results and commonly used methods to build confidence intervals (CIs) for the regression, that is the process output expectation (Seber, 1977); more recently, these results have been extended to nonlinear models (Bates & Watts, 1988; Seber & Wild, 1989). In the neural network modeling studies however, these results are seldom exploited, and generally only an average estimate of the neural model reliability is given through the mean square model error on a test set; but in an application, one often wishes to know a CI at any input value of interest. Nevertheless, thanks to the increase of computer power, the use of bootstrap methods has increased in the past years (Efron & Tibshirani, 1993). These nonanalytic methods have been proposed to build CIs for neural networks (Heskes, 1997; Paass, 1993; Tibshirani, 1996), but with the shortcoming of requiring a large number of trainings.

This paper presents an economic alternative to the construction of CIs using neural networks. This approach being built on the *linear* least squares (LS) theory applied to the linear Taylor expansion (LTE) of the output of nonlinear models, we first recall how to establish CIs for linear models in Section 2, and then approximate CIs for nonlinear models in Section 3. In Section 4, we exploit these known theoretical results for practical modeling problems involving neural models. We show that the LTE of a nonlinear model output not only provides a CI at any input value of interest, but also gives a tool to detect the possible ill-conditioning of the model, and as in Monari (1999) and Monari and Dreyfus (submitted), to estimate its performance through an approximate leave-one-out (LOO) score. A real-world illustration is given through an industrial application, the modeling of the elasticity of a complex material from some of its structural descriptors. Section 5 compares the LS LTE approach to other analytic approaches, and discusses its advantages with respect to bootstrap approaches.

We consider single-output models, since each output of a multi-output model can be handled separately. We deal with

* Corresponding author. Tel.: +33-1-40-79-45-45; fax: +33-1-40-79-44-25.

*E-mail address:* isabelle.rivals@espci.fr (I. Rivals).

Current address: Équipe de Statistique Appliquée, École Supérieure de Physique et de Chimie Industrielles, 10 rue Vauquelin, 75231 Paris Cedex 05, France.

*Abbreviations*: CI: confidence interval; LOO: leave-one-out; LS: least squares; LTE: linear Taylor expansion; SISO: single input—single output; MISO: multi input—single output; MSTE: mean square training error: MSPE: mean square performance error.

## Nomenclature

We distinguish between random variables and their values (or realizations) by using upper- and lowercase letters; all vectors are column vectors, and are denoted by boldface letters; nonrandom matrices are denoted by light lowercase letters

$\boldsymbol{x}$       Nonrandom $n$-input vector

$Y_p = Y_p|\boldsymbol{x}$   Random scalar output depending on $\boldsymbol{x}$

$E(Y_p|\boldsymbol{x})$   Mathematical expectation, or regression function, of $Y_p$ given $\boldsymbol{x}$

$W$       Random variable with zero expectation denoting additive noise

$\sigma^2$       Variance of $W$

$\{\boldsymbol{x}^k, y_p^k\}_{k=1 \text{ to } N}$ Data set of $N$ input–output pairs, where the $\{\boldsymbol{x}^k\}$ are nonrandom $n$-vectors, and the $\{y_p^k\}$ are the corresponding realizations of the random outputs $\{Y_p^k = Y_p|\boldsymbol{x}^k\}$

$\{(\boldsymbol{x}^k)^{\mathrm{T}}\boldsymbol{\theta}, \boldsymbol{\theta} \in \mathbb{R}^n\}$ Family of linear functions of $\boldsymbol{x}$ parameterized by $\boldsymbol{\theta}$

$\boldsymbol{\theta}_p$       Unknown true $q$-parameter vector ($q = n$ in the linear case)

$x = [\boldsymbol{x}^1\ \boldsymbol{x}^2\ ...\ \boldsymbol{x}^N]^{\mathrm{T}}$ Nonrandom $(N,n)$ input matrix

$\boldsymbol{Y}_p = [Y_p^1\ Y_p^2\ ...\ Y_p^N]^{\mathrm{T}}$ Random $N$-vector of the outputs of the data set

$\boldsymbol{W} = [W^1\ W^2\ ...\ W^N]^{\mathrm{T}}$ Random $N$-vector with $E(\boldsymbol{W}) = 0$

$J(\boldsymbol{\theta})$       Value of the least squares cost function

$\boldsymbol{\Theta}_{\mathrm{LS}}$       Least squares estimator of $\boldsymbol{\theta}_p$

$\boldsymbol{\theta}_{\mathrm{LS}}$       Least squares estimate of $\boldsymbol{\theta}_p$

$\boldsymbol{R} = \boldsymbol{Y}_p - x\boldsymbol{\Theta}_{\mathrm{LS}}$ Least squares residual random $N$-vector in the linear case

$\boldsymbol{r}$       Value of $\boldsymbol{R}$

$\mathscr{m}(x)$       Range of $x$ (linear manifold)

$p_x$       Orthogonal projection matrix on $\mathscr{m}(x)$

$S^2$       Estimator of $\sigma^2$

$s^2$       Value of $S^2$

$\{f(\boldsymbol{x}, \boldsymbol{\theta}), \boldsymbol{\theta} \in \mathbb{R}^q\}$ Family of nonlinear functions of $\boldsymbol{x}$ parameterized by $\boldsymbol{\theta}$

$\boldsymbol{f}(x, \boldsymbol{\theta})$   $N$-vector $[f(\boldsymbol{x}^1, \boldsymbol{\theta})...f(\boldsymbol{x}^k, \boldsymbol{\theta})...f(\boldsymbol{x}^N, \boldsymbol{\theta})]^{\mathrm{T}}$

$\boldsymbol{R} = \boldsymbol{Y}_p - \boldsymbol{f}(x, \boldsymbol{\Theta}_{\mathrm{LS}})$ Least squares residual random $N$-vector

$\xi = [\boldsymbol{\xi}^1\ \boldsymbol{\xi}^2\ ...\ \boldsymbol{\xi}^N]^{\mathrm{T}}$ Unknown nonrandom $(N,q)$ matrix with $\boldsymbol{\xi}^k = \partial f(\boldsymbol{x}^k, \boldsymbol{\theta})/\partial \boldsymbol{\theta}|_{\boldsymbol{\theta}=\boldsymbol{\theta}_p}$

$\mathscr{m}(\xi)$       Range of $\xi$

$p_\xi$       Orthogonal projection matrix on $\mathscr{m}(\xi)$

$z = [z^1\ z^2\ ...\ z^N]^{\mathrm{T}}$ Matrix approximating $\xi$ with $z^k = \partial f(\boldsymbol{x}^k, \boldsymbol{\theta})/\partial \boldsymbol{\theta}|_{\boldsymbol{\theta}=\boldsymbol{\theta}_{\mathrm{LS}}}$

$\mathscr{m}(z)$       Range of $z$

$p_z$       Orthogonal projection matrix on $\mathscr{m}(z)$

$I_N$       $(N,N)$ identity matrix

$\boldsymbol{\theta}_{\mathrm{LS}}^{(k)}$       Leave-one-out (the $k$th example) least squares estimate of $\boldsymbol{\theta}_p$

$\{e^k\}_{k=1 \text{ to } N}$ Leave-one-out errors

$n_{\mathrm{h}}$       Number of hidden neurons of a neural network

$H$       Random Hessian matrix of the cost function

$h$       Value of the Hessian matrix of the cost function

$\mathrm{var}(\widehat{f(\boldsymbol{x}_1 \boldsymbol{\Theta}_{\mathrm{LS}})})_{\mathrm{ref}}$ Reference variance estimate

$\mathrm{var}(\widehat{f(\boldsymbol{x}, \boldsymbol{\Theta}_{\mathrm{LS}})})_{\mathrm{LTE}}$ LTE estimate of the variance of a nonlinear model output

$\mathrm{var}(\widehat{f(\boldsymbol{x}, \boldsymbol{\Theta}_{\mathrm{LS}})})_{\mathrm{Hessian}}$ Hessian estimate of the variance of a nonlinear model output

$\mathrm{var}(\widehat{f(\boldsymbol{x}, \boldsymbol{\Theta}_{\mathrm{LS}})})_{\mathrm{sandwich}}$ Sandwich estimate of the variance of a nonlinear model output

Abbreviations

CI       confidence interval

LOO       leave-one-out

LS       least squares

LTE       linear Taylor expansion

SISO       single input - single output

MISO       multi input - single output

MSTE       mean square training error

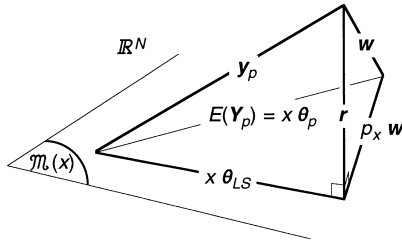MSPE       mean square performance error

Fig. 1. Geometric representation of the linear LS solution (true assumed model).

static modeling problems for the case of a nonrandom (noise free) $n$-input vector $\boldsymbol{x} = [x_1\ x_2\ ...\ x_n]^{\mathrm{T}}$, and a noisy measured output $y_p$ which is considered as the actual value of a random variable $Y_p = Y_p | \boldsymbol{x}$ depending on $\boldsymbol{x}$. We assume that there exists an unknown regression function $E(Y_p | \boldsymbol{x})$ such that for any fixed value of $\boldsymbol{x}$:

$$Y_p | \boldsymbol{x} = E(Y_p | \boldsymbol{x}) + W | \boldsymbol{x} \tag{1}$$

where $W | \boldsymbol{x}$ is thus a random variable with zero expectation. A family of parameterized functions $\{f(\boldsymbol{x}, \boldsymbol{\theta}),\ \boldsymbol{x} \in \mathbb{R}^n,\ \boldsymbol{\theta} \in \mathbb{R}^q\}$ is called an *assumed* model. This assumed model is said to be *true* if there exists a value $\boldsymbol{\theta}_p$ of $\boldsymbol{\theta}$ such that, $\forall \boldsymbol{x}$ in the input domain of interest, $f(\boldsymbol{x}, \boldsymbol{\theta}_p) = E(Y_p | \boldsymbol{x})$. In the following, a data set of $N$ input–output pairs $\{\boldsymbol{x}^k, y_p^k\}_{k=1 \text{ to } N}$ is available, where the $\boldsymbol{x}^k = [x_1^k\ x_2^k\ ...\ x_n^k]^{\mathrm{T}}$ are nonrandom $n$-vectors, and the $\{y_p^k\}$ are the corresponding realizations of the random variables $\{Y_p^k = Y_p | \boldsymbol{x}^k\}$.[1] The goal of the modeling procedure is not only to estimate the regression $E(Y_p | \boldsymbol{x})$ in the input domain of interest with the output of a model, but also to compute the value of a CI for the regression, that is the value of a random interval with a chosen probability to contain the regression. For the presentation of the results of linear and nonlinear regression estimation, we deal with the true model (a model which is linear in the parameters in Section 2, a nonlinear one in Section 3), i.e. we consider that a family of functions containing the regression is known. In Section 4, we consider the general realistic case of neural black-box modeling where a preliminary selection among candidate neural models is first performed because a true model is not known a priori.

## 2. Confidence intervals for linear models

We consider a true linear assumed model, that is the associated family of linear functions $\{\boldsymbol{x}^{\mathrm{T}} \boldsymbol{\theta},\ \boldsymbol{x} \in \mathbb{R}^n,\ \boldsymbol{\theta} \in \mathbb{R}^n\}$ contains the regression; Eq. (1) can thus be uniquely rewritten as:

$$Y_p | \boldsymbol{x} = \boldsymbol{x}^{\mathrm{T}} \boldsymbol{\theta}_p + W | \boldsymbol{x} \tag{2}$$

where $\boldsymbol{\theta}_p$ is an unknown $n$-parameter vector. Model (2) associated to the data set leads to:

$$\boldsymbol{Y}_p = x \boldsymbol{\theta}_p + \boldsymbol{W} \tag{3}$$

where $x = [\boldsymbol{x}^1\ \boldsymbol{x}^2\ ...\ \boldsymbol{x}^N]^{\mathrm{T}}$ is the nonrandom $(N,n)$ input matrix, $\boldsymbol{Y}_p = [Y_p^1\ Y_p^2\ ...\ Y_p^N]^{\mathrm{T}}$ and $\boldsymbol{W} = [W^1\ W^2\ ...\ W^N]^{\mathrm{T}}$ are random $N$-vectors, with $E(\boldsymbol{W}) = 0$. Geometrically, this means that $E(\boldsymbol{Y}_p | \boldsymbol{x}) = x \boldsymbol{\theta}_p$ belongs to the solution surface, the linear manifold $\mathcal{m}(x)$ of the observation space $\mathbb{R}^N$ spanned by the columns of the input matrix[2] (the range of $x$). We assume that $\mathcal{m}(x)$ is of dimension $n$, that is rank$(x) = n$. In other words, the model is identifiable, i.e. the data set is appropriately chosen, possibly using experimental design.

### 2.1. The linear least squares solution

The LS estimate $\boldsymbol{\theta}_{\mathrm{LS}}$ of $\boldsymbol{\theta}_p$ minimizes the empirical quadratic cost function:

$$J(\boldsymbol{\theta}) = \tfrac{1}{2} \sum_{k=1}^{N} (y_p^k - (\boldsymbol{x}^k)^{\mathrm{T}} \boldsymbol{\theta})^2 = \tfrac{1}{2} (\boldsymbol{y}_p - x\boldsymbol{\theta})^{\mathrm{T}} (\boldsymbol{y}_p - x\boldsymbol{\theta}) \tag{4}$$

The estimate $\boldsymbol{\theta}_{\mathrm{LS}}$ is a realization of the LS estimator $\boldsymbol{\Theta}_{\mathrm{LS}}$ whose expression is:

$$\boldsymbol{\Theta}_{\mathrm{LS}} = (x^{\mathrm{T}}x)^{-1} x^{\mathrm{T}} \boldsymbol{Y}_p = \boldsymbol{\theta}_p + (x^{\mathrm{T}}x)^{-1} x^{\mathrm{T}} \boldsymbol{W} \tag{5}$$

As the assumed model is true, this estimator is unbiased. The orthogonal projection matrix on $\mathcal{m}(x)$ is $p_x = x(x^{\mathrm{T}}x)^{-1} x^{\mathrm{T}}$. It follows from Eq. (5) that the unbiased LS estimator of $E(\boldsymbol{Y}_p | \boldsymbol{x})$ is:

$$x \boldsymbol{\Theta}_{\mathrm{LS}} = x \boldsymbol{\theta}_p + p_x \boldsymbol{W} \tag{6}$$

that is the sum of $E(\boldsymbol{Y}_p | \boldsymbol{x})$ and of the projection of $\boldsymbol{W}$ on $\mathcal{m}(x)$, as shown in Fig. 1. Let $\boldsymbol{R}$ denote the residual random $N$-vector $\boldsymbol{R} = \boldsymbol{Y}_p - x \boldsymbol{\Theta}_{\mathrm{LS}}$, that is the vector of the errors on the data set, then:

$$\boldsymbol{R} = (I_N - p_x) \boldsymbol{W} \tag{7}$$

Under the assumption that the $\{W^k\}$ are identically distributed and uncorrelated (homoscedastic), i.e. the noise covariance matrix is $K(\boldsymbol{W}) = \sigma^2 I_N$, it follows from Eq. (5) that the variance of the LS estimator of the regression for any input $\boldsymbol{x}$ of interest is:[3]

$$\mathrm{var}(\boldsymbol{x}^{\mathrm{T}} \boldsymbol{\Theta}_{\mathrm{LS}}) = \sigma^2 \boldsymbol{x}^{\mathrm{T}} (x^{\mathrm{T}}x)^{-1} \boldsymbol{x} \tag{8}$$

Using Eq. (7), we obtain the unbiased estimator

$$S^2 = \frac{\boldsymbol{R}^{\mathrm{T}} \boldsymbol{R}}{N - n}$$

of $\sigma^2$; the corresponding (unbiased) estimate of the variance

---

[1] We recall that we distinguish between random variables and their values (or realizations) by using upper and lowercase letters, e.g. $Y_p^k$ and $y_p^k$; all vectors are column vectors, and are denoted by boldface letters, e.g. the $n$-vectors $\boldsymbol{x}$ and $\{\boldsymbol{x}^k\}$; nonrandom matrices are denoted by light lowercase letters (except the unambiguous identity matrix).

[2] $\mathcal{m}(x)$ is sometimes called the "expectation surface" (Seber & Wild, 1989); as a matter of fact, the solution surface coincides with the expectation surface only when the assumed model is true.

[3] We recall that $\boldsymbol{x}$ (boldface) is the $(n, 1)$ input vector of interest, and that $x$ is the experimental $(N, n)$ input matrix.
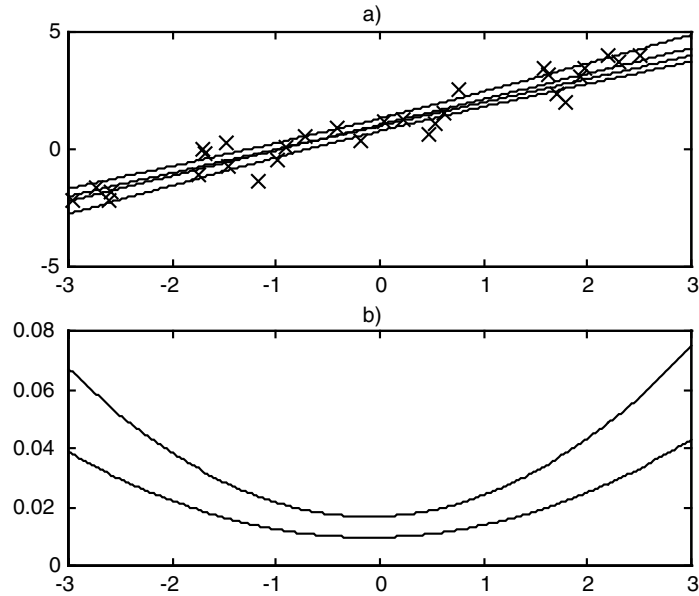
Fig. 2. CI for process #1, a simulated linear SISO process (true assumed model with $n = 2$ parameters): (a) regression (thin line), data set (crosses), model output and 99% CI (thick lines); and (b) true variance (thin line) and LS estimate of the variance (thick line) of $\boldsymbol{x}^{\mathrm{T}}\boldsymbol{\Theta}_{\mathrm{LS}}$.

of $\boldsymbol{x}^{\mathrm{T}}\boldsymbol{\Theta}_{\mathrm{LS}}$ is thus:

$$\widehat{\mathrm{var}(\boldsymbol{x}^{\mathrm{T}}\boldsymbol{\Theta}_{\mathrm{LS}})} = s^2 \boldsymbol{x}^{\mathrm{T}}(x^{\mathrm{T}}x)^{-1}\boldsymbol{x} \qquad (9)$$

where $s$ is the value of $S$.

### 2.2. Confidence intervals for a linear regression

If the $\{W^k\}$ are homoscedastic gaussian variables, that is $W \to N_N(\mathbf{0}, \sigma^2 I_N)$ :

**Theorem 1.**

$$\boldsymbol{\Theta}_{\mathrm{LS}} - \boldsymbol{\theta}_p \to N_n(\mathbf{0}, \sigma^2(x^{\mathrm{T}}x)^{-1}) \qquad (10)$$

**Theorem 2.**

$$\frac{\boldsymbol{R}^{\mathrm{T}}\boldsymbol{R}}{\sigma^2} \to \chi^2_{N-n} \qquad (11)$$

**Theorem 3.** $\boldsymbol{\Theta}_{\mathrm{LS}}$ *is statistically independent from* $\boldsymbol{R}^{\mathrm{T}}\boldsymbol{R}$.

The proof of the above theorems follows from Fig. 1 and from the Fisher–Cochrane theorem (Goodwin & Payne, 1977), see for instance (Seber, 1977).

The goal is to build a CI for the regression value $E(Y_p|\boldsymbol{x}) = \boldsymbol{x}^{\mathrm{T}}\boldsymbol{\theta}_p$, for any input vector $\boldsymbol{x}$ of interest. The variance of the measurements $\sigma^2$ being unknown, let us build a normalized centered gaussian variable where both

$E(Y_p|\boldsymbol{x})$ and $\sigma$ appear:

$$\frac{\boldsymbol{x}^{\mathrm{T}}\boldsymbol{\Theta}_{\mathrm{LS}} - E(Y_p|\boldsymbol{x})}{\sigma\sqrt{\boldsymbol{x}^{\mathrm{T}}(x^{\mathrm{T}}x)^{-1}\boldsymbol{x}}} \to N(0, 1) \qquad (12)$$

Thus, using the Pearson variable (11), which is independent from Eq. (12) according to Theorem 3, we obtain the following Student variable:

$$\frac{\boldsymbol{x}^{\mathrm{T}}\boldsymbol{\Theta}_{\mathrm{LS}} - E(Y_p|\boldsymbol{x})}{S\sqrt{\boldsymbol{x}^{\mathrm{T}}(x^{\mathrm{T}}x)^{-1}\boldsymbol{x}}} \to \text{Student } (N - n) \qquad (13)$$

A $100(1 - \alpha)\%$ CI for $E(Y_p|\boldsymbol{x})$ is thus:

$$\boldsymbol{x}^{\mathrm{T}}\boldsymbol{\theta}_{\mathrm{LS}} \pm t_{N-n}\left(1 - \frac{\alpha}{2}\right)s\sqrt{\boldsymbol{x}^{\mathrm{T}}(x^{\mathrm{T}}x)^{-1}\boldsymbol{x}} \qquad (14)$$

where $t_{N-n}$ is the inverse of the Student$(N - n)$ cumulative distribution.

Note that Eq. (14) allows to compute a CI corresponding to any input vector, and that it is much more informative than average values such as that the mean square error on the data set, or the mean of the variance estimate over the data set;[4] as a matter of fact, the latter invariably equals $s^2 n/N$.

---

[4] The mean of the variance estimate over the training data set is: $\frac{1}{N}\sum_{k=1}^{N} s^2(\boldsymbol{x}^k)^{\mathrm{T}}(x^{\mathrm{T}}x)^{-1}\boldsymbol{x}^k = \frac{s^2}{N}\sum_{k=1}^{N}[p_x]_{kk} = \frac{s^2}{N}trace(p_x)$. As $p_x$ is the orthogonal projection matrix on a $n$-dimensional subspace, $trace(p_x) = n$.
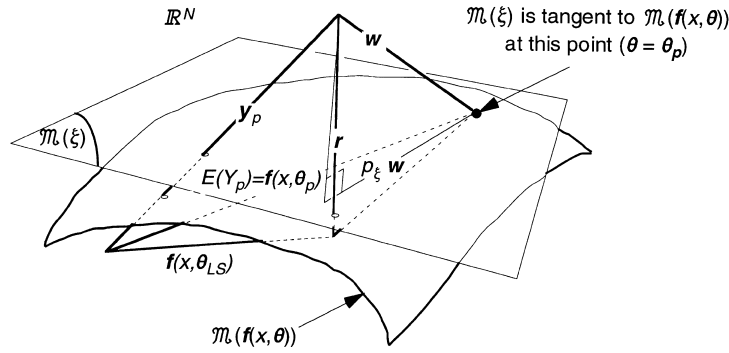
Fig. 3. Geometric representation of the nonlinear LS solution and of its LTE approximation (true assumed model).

## 2.3. Example of a simulated linear SISO process (process #1)

We consider a simulated single input-single output (SISO) linear process:

$$y_p^k = \theta_{p_1} + \theta_{p_2} x^k + w^k \quad k = 1 \text{ to } N \tag{15}$$

We take $\theta_{p_1} = 1$, $\theta_{p_2} = 1$, $\sigma^2 = 0.5$, $N = 30$. The inputs $\{x^k\}$ of the data set are uniformly distributed in $[-3; 3]$, as shown in Fig. 2a. The family of functions $\{\theta_1 + \theta_2 x, \; \boldsymbol{\theta} \in \mathbb{R}^2\}$ is considered, that is the assumed model is true, and we choose a confidence level of 99% ($t_{28}(1\%) = 2.76$). The LS estimation leads to $s^2 = 0.29$, i.e. underestimates the noise variance. Fig. 2b displays the estimate (9) of the variance of $\boldsymbol{x}^{\mathrm{T}} \boldsymbol{\Theta}_{\mathrm{LS}}$, and the true variance (8). The estimator $S^2$ of the noise variance being unbiased, the difference between the estimated variance (9) and the (usually unknown) true variance (8) is only due to the particular values of the measurement noise. Fig. 2(a) shows the regression $E(Y_p|x)$, the data set, the model output and the 99% CI for the regression computed with Eq. (14).

## 3. Approximate confidence intervals for nonlinear models

We consider a family of nonlinear functions $\{f(\boldsymbol{x}, \boldsymbol{\theta}), \; \boldsymbol{x} \in \mathbb{R}^n, \; \boldsymbol{\theta} \in \mathbb{R}^q\}$ which contains the regression, that is the assumed model is true; Eq. (1) can thus be rewritten as:

$$Y_p|\boldsymbol{x} = f(\boldsymbol{x}, \boldsymbol{\theta}_p) + W|\boldsymbol{x} \tag{16}$$

where $\boldsymbol{\theta}_p$ is an unknown $q$-parameter vector. We denote by $\boldsymbol{f}(x, \boldsymbol{\theta}_p)$ the unknown vector $[f(\boldsymbol{x}^1, \boldsymbol{\theta}_p)...f(\boldsymbol{x}^k, \boldsymbol{\theta}_p)... f(\boldsymbol{x}^N, \boldsymbol{\theta}_p)]^{\mathrm{T}}$ defined on the data set, thus:

$$\boldsymbol{Y}_p = \boldsymbol{f}(x, \boldsymbol{\theta}_p) + \boldsymbol{W} \tag{17}$$

As in Section 2, $x$ denotes the $(N,n)$ input matrix,[5] and $\boldsymbol{Y}_p$ and $\boldsymbol{W}$ are random $N$-vectors with $E(\boldsymbol{W}) = 0$. Geometrically, this means that $E(\boldsymbol{Y}_p|x)$ belongs to the solution surface, the manifold $\mathcal{m}(f(\boldsymbol{x}, \boldsymbol{\theta})) = \{f(\boldsymbol{x}, \boldsymbol{\theta}), \; \boldsymbol{\theta} \in \mathbb{R}^q\}$ of $\mathbb{R}^N$.

## 3.1. The linear Taylor expansion of the nonlinear least squares solution

A LS estimate $\boldsymbol{\theta}_{\mathrm{LS}}$ of $\boldsymbol{\theta}_p$ minimizes the empirical cost function:[6]

$$J(\boldsymbol{\theta}) = \tfrac{1}{2} \sum_{k=1}^N (y_p^k - f(\boldsymbol{x}^k, \boldsymbol{\theta}))^2 = \tfrac{1}{2}(\boldsymbol{y}_p - \boldsymbol{f}(x, \boldsymbol{\theta}))^{\mathrm{T}}(\boldsymbol{y}_p - \boldsymbol{f}(x, \boldsymbol{\theta})) \tag{18}$$

The estimate $\boldsymbol{\theta}_{\mathrm{LS}}$ is a realization of the LS estimator $\boldsymbol{\Theta}_{\mathrm{LS}}$. Efficient algorithms are at our disposal for the minimization of the cost function (18): they can lead to an absolute minimum, but they do not give an analytic expression of the estimator that could be used to build CIs. In order to take advantage of the results concerning linear models, it is worthwhile considering a linear approximation of $\boldsymbol{\Theta}_{\mathrm{LS}}$ which is obtained by writing the LTE of $f(\boldsymbol{x}, \boldsymbol{\theta})$ around $f(\boldsymbol{x}, \boldsymbol{\theta}_p)$:

$$f(\boldsymbol{x}, \boldsymbol{\theta}) \approx f(\boldsymbol{x}, \boldsymbol{\theta}_p) + \sum_{r=1}^q \left( \left. \frac{\partial f(\boldsymbol{x}, \boldsymbol{\theta})}{\partial \theta_r} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}_p} (\theta_r - \theta_{p_r}) \right)$$

$$= f(\boldsymbol{x}, \boldsymbol{\theta}_p) + \boldsymbol{\xi}^{\mathrm{T}}(\boldsymbol{\theta} - \boldsymbol{\theta}_p) \tag{19}$$

where

$$\xi = \left. \frac{\partial f(\boldsymbol{x}, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}_p}$$

Thus, with the matrix notation:

$$\boldsymbol{f}(x, \boldsymbol{\theta}) \approx \boldsymbol{f}(x, \boldsymbol{\theta}_p) + \xi(\boldsymbol{\theta} - \boldsymbol{\theta}_p) \tag{20}$$

where

$$\xi = [\boldsymbol{\xi}^1 \boldsymbol{\xi}^2 ... \boldsymbol{\xi}^N]^{\mathrm{T}}$$

---

[5] In the case of a nonlinear model, $x$ is merely a two-dimensional array.

[6] In the case of a multilayer neural network, the minimum value of the cost function can be obtained for several values of the parameter vector; but, since the only function-preserving transformations are neuron exchanges, as well as sign flips for odd activation functions like the hyperbolic tangent [Sussman 1992], we will legitimately consider the neighborhood of one of these values only.

and

$$\boldsymbol{\xi}^k = \frac{\partial f(\boldsymbol{x}^k, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}\bigg|_{\boldsymbol{\theta}=\boldsymbol{\theta}_p}$$

The $(N, q)$ matrix $\xi$ is the nonrandom and unknown (since $\boldsymbol{\theta}_p$ is unknown) Jacobian matrix of $f$. Using Eq. (20), one obtains, similarly to the linear case, the following approximation of $\boldsymbol{\Theta}_{LS}$ (see Appendix A.1 for a detailed derivation of Eqs. (21) and (23)):

$$\boldsymbol{\Theta}_{LS} \approx \boldsymbol{\theta}_p + (\xi^T \xi)^{-1} \xi^T \boldsymbol{W} \tag{21}$$

The range $m(\xi)$ of $\xi$ is tangent to the manifold $m(f(x, \boldsymbol{\theta}))$ at $\boldsymbol{\theta} = \boldsymbol{\theta}_p$; this manifold is assumed to be of dimension $q$, hence $\text{rank}(\xi) = q$. The matrix $p_\xi = \xi(\xi^T \xi)^{-1} \xi^T$ is the orthogonal projection matrix on $m(\xi)$. From Eqs. (20) and (21), the LS estimator of $E(\boldsymbol{Y}_p|x)$ can be approximately expressed by:

$$f(x, \boldsymbol{\Theta}_{LS}) \approx f(x, \boldsymbol{\theta}_p) + p_\xi \boldsymbol{W} \tag{22}$$

i.e. it is approximately the sum of $E(\boldsymbol{Y}_p|x)$ and of the projection of $\boldsymbol{W}$ on $m(\xi)$, as illustrated in Fig. 3. If $K(\boldsymbol{W}) = \sigma^2 I_N$ (homoscedasticity), the variance of the model output, that is the LS estimator of the regression, for an input $\boldsymbol{x}$ is approximately:

$$\text{var}(f(\boldsymbol{x}, \boldsymbol{\Theta}_{LS})) \approx \sigma^2 \boldsymbol{\xi}^T (\xi^T \xi)^{-1} \boldsymbol{\xi} \tag{23}$$

In the following, approximation (23) will be termed "the LTE approximation" of the model output variance. Let $\boldsymbol{R}$ denote the LS residual vector $\boldsymbol{R} = \boldsymbol{Y}_p - f(x, \boldsymbol{\theta}_{LS})$, thus:

$$\boldsymbol{R} \approx (I_N - p_\xi)\boldsymbol{W} \tag{24}$$

Under the assumption of appropriate regularity conditions on $f$, and for large $N$, the curvature of the solution surface[7] $m(f(x, \boldsymbol{\theta}))$ is small; thus, using Eq. (24), one obtains the asymptotically (i.e. when $N \to \infty$) unbiased estimator $S^2 = (\boldsymbol{R}^T \boldsymbol{R})/(N - q)$ of $\sigma^2$. In Eq. (23), the matrix $\xi$ takes the place of matrix $x$ in the linear case. But, as opposed to $x$, $\xi$ is unknown since it is a function of the unknown parameter $\boldsymbol{\theta}_p$. The $(N, q)$ matrix $\xi$ may be approximated by

$$z = [z^1 z^2 ... z^N]^T$$

where

$$z^k = \frac{\partial f(\boldsymbol{x}^k, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}\bigg|_{\boldsymbol{\theta}=\boldsymbol{\theta}_{LS}},$$

that is

$$z_r^k = \frac{\partial f(\boldsymbol{x}^k, \boldsymbol{\theta})}{\partial \theta_r}\bigg|_{\boldsymbol{\theta}=\boldsymbol{\theta}_{LS}} \tag{25}$$

In the following, we assume that $\text{rank}(z) = q$. Like the

---

[7] The curvature is usually decomposed in two components: (i) the intrinsic curvature, which measures the degree of bending and twisting of the solution surface $m(f(x, \boldsymbol{\theta}))$, and (ii) the parameter-effects curvature, which describes the degree of curvature induced by the choice of the parameters $\boldsymbol{\theta}$.
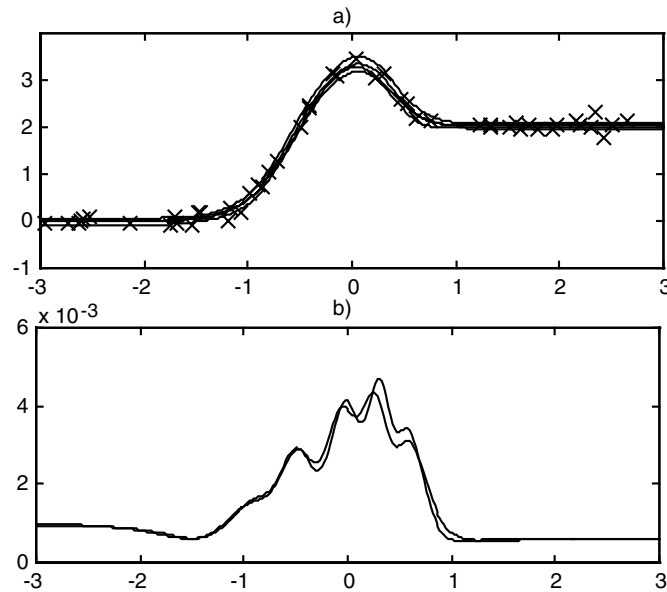
matrix $\xi$, the vector $\boldsymbol{\xi}$ is not available, and its value may be approximated by:

$$z = \frac{\partial f(\boldsymbol{x}, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}\bigg|_{\boldsymbol{\theta}=\boldsymbol{\theta}_{LS}} \tag{26}$$

From Eqs. (23), (25) and (26), the LTE estimate of the variance of the LS estimator of the regression for an input $\boldsymbol{x}$ is thus:

$$\widehat{\text{var}(f(\boldsymbol{x}, \boldsymbol{\Theta}_{LS}))}_{LTE} = s^2 z^T (z^T z)^{-1} z \tag{27}$$

*3.2. Approximate confidence intervals for a nonlinear regression*

If $\boldsymbol{W} \to N_N(\boldsymbol{0}, \sigma^2 I_N)$, and for large $N$, it follows from the above relations and from the linear LS theory (Seber & Wild, 1989) that:

**Theorem 4.**

$$\boldsymbol{\Theta}_{LS} \sim\!\!\longrightarrow N_q(\boldsymbol{\theta}_p, \sigma^2 (\xi^T \xi)^{-1}) \tag{28}$$

**Theorem 5.**

$$\frac{\boldsymbol{R}^T \boldsymbol{R}}{\sigma^2} \sim\!\!\longrightarrow \chi^2_{N-q} \tag{29}$$

**Theorem 6.** $\boldsymbol{\Theta}_{LS}$ *is approximately statistically independent from* $\boldsymbol{R}^T \boldsymbol{R}$.

Using Eqs. (23) and (28), let us again build a quasi-normalized and centered gaussian variable where both $E(Y_p|\boldsymbol{x})$ and $\sigma$ appear:

$$\frac{f(\boldsymbol{x}, \boldsymbol{\Theta}_{LS}) - E(Y_p|\boldsymbol{x})}{\sigma\sqrt{\boldsymbol{\xi}^T (\xi^T \xi)^{-1} \boldsymbol{\xi}}} \sim\!\!\longrightarrow N(0, 1) \tag{30}$$

Thus, the variable (29) being approximately independent from Eq. (30) according to Theorem 6, we have:

$$\frac{f(\boldsymbol{x}, \boldsymbol{\Theta}_{LS}) - E(Y_p|\boldsymbol{x})}{S\sqrt{\boldsymbol{\xi}^T (\xi^T \xi)^{-1} \boldsymbol{\xi}}} \sim\!\!\longrightarrow \text{Student } (N - q) \tag{31}$$

A $100(1 - \alpha)\%$ *approximate* CI for $E(Y_p|\boldsymbol{x})$ is thus:

$$f(\boldsymbol{x}, \boldsymbol{\theta}_{LS}) \pm t_{N-q}\left(1 - \frac{\alpha}{2}\right)s\sqrt{z^T (z^T z)^{-1} z} \tag{32}$$

Note that, when $N$ is large, the Student distribution is close to the normal distribution, and thus $t_{N-q}(1 - \frac{\alpha}{2}) \approx n(1 - \frac{\alpha}{2})$, where $n$ is the inverse of the normal cumulative distribution.

Like in the linear case, Eq. (32) allows to compute a CI at any input $\boldsymbol{x}$ of interest, which gives much more information than the value of the mean variance estimate over the data

Fig. 4. CI for process #2, a simulated "neural" SISO process (the assumed model, a two hidden neurons network with $q = 7$ parameters, is true): (a) regression (thin line), the $N = 50$ examples of the data set (crosses), model output and 99% approximate CI (thick lines); and (b) reference (thin line) and LTE (thick line) estimates of the variance of $f(x, \Theta_{LS})$.

set: as a matter of fact, the latter always approximately equals $s^2(q/N)$.

From a practical point of view, the construction of a CI for a neural model output at any input $x$ of interest involves once and for all the computation of the matrix $z$, that is the $N \times q$ partial derivatives of the model output with respect to the parameters evaluated at $\theta_{LS}$ for the data inputs $\{x^k\}_{k=1 \text{ to } N}$, and, for each $x$ value, that of $z$, i.e. the derivatives at $x$. In the case of a neural model, these derivatives are easily obtained with the backpropagation algorithm.

All the previous results and the above considerations are valid provided an absolute minimum of the cost function (18) is reached. In real-life, several estimations of the parameters must thus be made starting from different initial values, the estimate corresponding to the lowest minimum being kept in order to have a high probability to obtain an absolute minimum. In the examples of this work, the algorithm used for the minimization of the cost function is the Levenberg algorithm, as described for example in Bates and Watts (1988), and several trainings are performed. The probability of getting trapped in a relative minimum increasing with the number of parameters of the network and decreasing when the size of the data set increases, the number of trainings is chosen accordingly.

### 3.3. Quality of the approximate confidence intervals

#### 3.3.1. Theoretical analysis

The quality of the approximate CI depends essentially on the curvature of the solution surface $m(f(x, \theta))$, which depends on the regularity of $f$ and on the value of $N$. In practice, $f$ is often regular enough for a first-order approximation to be satisfactory, provided that $N$ is large enough.

Thus, if $N$ is large: (i) as in the linear case, the estimator of the noise variance $S^2$ is unbiased, and the difference between $s^2$ and $\sigma^2$ is only due to the particular values of the noise; and (ii) the variance of $f(x, \Theta_{LS})$ is small, and $\theta_{LS}$ is likely to be close to $\theta_p$: $z$ and $z$ are thus likely to be good approximations of, respectively, $\xi$ and $\boldsymbol{\xi}$. A reliable estimate of a CI may thus be obtained from the LTE variance estimate (27). On the other hand, if $N$ is too small: (i) as opposed to the linear case, even if the assumed model is true, the estimator of the noise variance $S^2$ is biased; and (ii) the variance of $f(x, \Theta_{LS})$ is large, and $\theta_{LS}$ is likely to differ from $\theta_p$: $z$ and $z$ risk to be poor approximations of $\xi$ and $\boldsymbol{\xi}$. Thus, if $N$ is diagnosed as too small, one cannot "have confidence" in the confidence intervals (32), and additional data should be gathered.

#### 3.3.2. Quantitative analysis

As detailed for example by Antoniadis, Berruyer and Carmona (1992), Bates and Watts (1988) and Seber and Wild (1989) different measures of the curvature can be computed, and can be used in each particular case to evaluate the accuracy of the LTE. In Section 4.3 dealing with neural network modeling, we give indications on how to judge if $N$ is large enough for the approximate CI to be accurate.

In order to evaluate the accuracy of the LTE variance estimate (27) when dealing with simulated processes, we introduce an estimate of the unknown true variance of $f(x, \Theta_{LS})$ that is not biased by curvature effects, *the reference variance estimate* $\mathrm{var}(\widehat{f(\mathbf{x}_1 \Theta_{LS})})_{\mathrm{ref}}$. This estimate is computed on a large number $M$ of other sets of $N$ outputs corresponding to the inputs of the training set, and whose values are obtained with different realizations (simulated
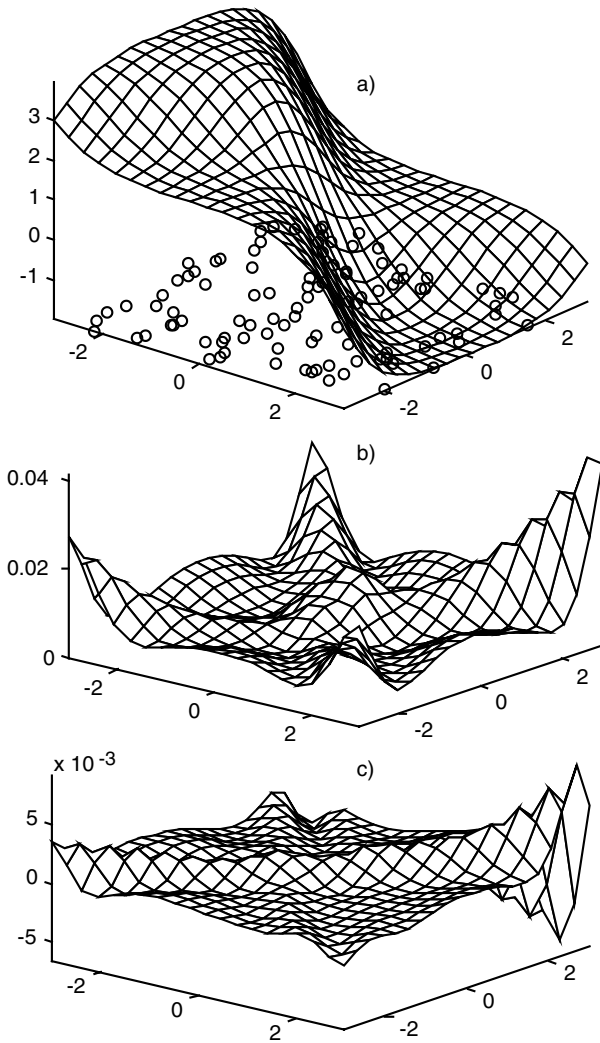
Fig. 5. CI for process #3, a simulated "neural" MISO process (the assumed model, a two hidden neurons network with $q = 9$ parameters, is true): (a) the $N = 100$ inputs of the data set (circles) and regression; (b) LTE estimate of the variance of $f(x, \Theta_{\mathrm{LS}})$; and (c) difference between the reference and the LTE estimates of the variance of $f(x, \Theta_{\mathrm{LS}})$.

values) of the noise $W$. The $i$th LS estimate $f(x, \theta_{\mathrm{LS}}^{(i)})$ of $E(Y_p|x)$ is computed with data set $i$ ($i = 1$ to $M$), and the reference estimate of the variance at input $x$ is computed as:

$$\mathrm{var}(\widehat{f(\mathbf{x}_1 \Theta_{\mathrm{LS}})})_{\mathrm{ref}} = \frac{1}{M} \sum_{i=1}^{M} (f(x, \theta_{\mathrm{LS}}^{(i)}) - \langle f(x) \rangle)^2,$$

where $\langle f(x) \rangle = \dfrac{1}{M} \displaystyle\sum_{i=1}^{M} f(x, \theta_{\mathrm{LS}}^{(i)})$ \hfill (33)

In the nonlinear case, we thus use three notions related to the (true) variance of $f(x, \Theta_{\mathrm{LS}})$: (1) the LTE variance approximation (23), which is a good approximation of the true variance if the curvature is small, as we show in Section 3.4.3, and which can be computed only when the process is

simulated; (2) the LTE variance estimate (27), which is the estimate that can be computed in real-life; and (3) the reference variance estimate (33), which tends to the true variance when $M$ tends to infinity because it is not biased by curvature effects, and which can be computed only when the process is simulated.

### 3.4. Illustrative examples

As we are concerned with neural models, and since, in this section, the assumed model is true, the following examples bring into play "neural" processes, that is to say processes whose regression function is the output of a neural network; the more realistic case of arbitrary processes for whom a family of nonlinear functions (a neural network with a given architecture) containing the regression is unknown is tackled in the next section.

### 3.4.1. Example of a simulated "neural" SISO process (process #2)

We consider a SISO process simulated by a neural network with one hidden layer of two hidden neurons with hyperbolic tangent activation function and a linear output neuron:

$$y_p^k = \theta_{p_1} + \theta_{p_2} \tanh(\theta_{p_3} + \theta_{p_4} x^k)$$

$$+ \theta_{p_5} \tanh(\theta_{p_6} + \theta_{p_7} x^k) + w^k \; k = 1 \text{ to } N \quad (34)$$

We take $\boldsymbol{\theta}_p = [1;\; 2;\; 1;\; 2;\; -1;\; -1;\; 3]^{\mathrm{T}}$, $\sigma^2 = 10^{-2}$, $N = 50$. The inputs $\{x^k\}$ of the data set are uniformly distributed in $[-3; 3]$, as shown in Fig. 4a. The family of functions $\{\theta_1 + \theta_2 \tanh(\theta_3 + \theta_4 x) + \theta_5 \tanh(\theta_6 + \theta_7 x), \boldsymbol{\theta} \in \mathbb{R}^7\}$ is considered, that is the assumed model is true, and we choose a confidence level of 99% ($t_{43}(1\%) = 2.58$). The minimization of the cost function with the Levenberg algorithm leads to $s^2 = 1.02 \times 10^{-2}$. Fig. 4b displays the LTE estimate of the variance of $f(x,\ \Theta_{\mathrm{LS}})$ (27), and the reference estimate (33) computed over $M = 10\,000$ sets. Fig. 4a shows the regression, the data set used for the LS estimation, and the corresponding model output and 99% CI (32). The model being true and the size $N = 50$ of the data set being relatively large with respect to the number of parameters and to the noise variance, we observe that: (i) $s^2 \approx \sigma^2$; (ii) the model output is close to the regression, leading to good approximations of $\xi$ and of $\boldsymbol{\xi}$ by $z$ and $z$. Thus, (i) and (ii) lead to an accurate LTE estimate of the variance, and hence of the CI.

### 3.4.2. Example of a simulated "neural" MISO process (process #3)

We consider a MISO process simulated by a neural network with two inputs, one hidden layer of two "tanh"
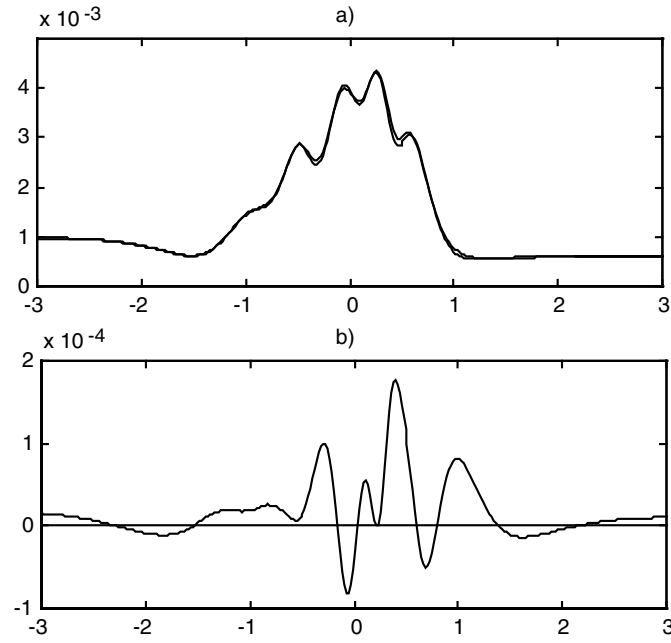
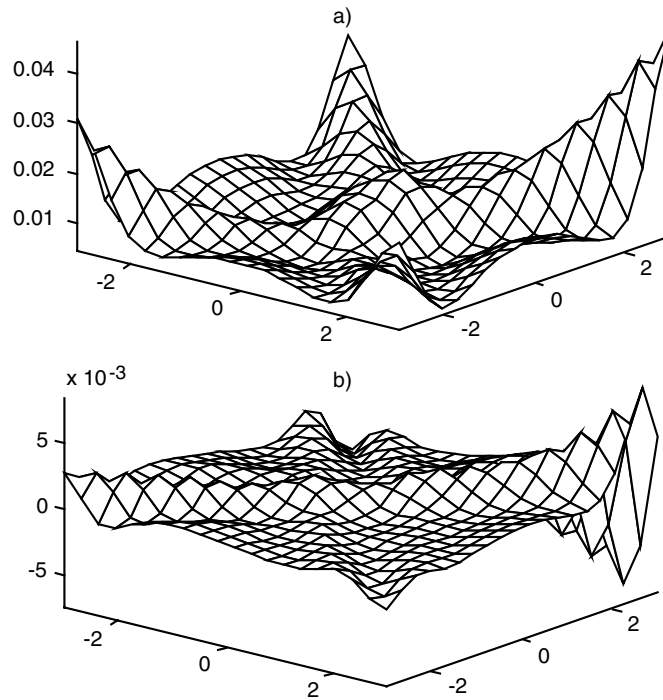Fig. 6. Accuracy of the LTE approximation of the variance for process #2: (a) reference estimate of the variance of $f(x, \Theta_{LS})$ (thin line), LTE approximation obtained with the true values $\theta_p$ and $\sigma^2$ (thick line); and (b) difference between the reference estimate of the variance of $f(x, \Theta_{LS})$ and the LTE approximation obtained with $\theta_p$ and $\sigma^2$.

hidden neurons and a linear output neuron:

$$y_p^k = \theta_{p_1} + \theta_{p_2} \tanh(\theta_{p_3} + \theta_{p_4} x_1^k + \theta_{p_5} x_2^k)$$

$$+ \theta_{p_6} \tanh(\theta_{p_7} + \theta_{p_8} x_1^k + \theta_{p_9} x_2^k) + w^k \ k = 1 \text{ to } N \quad (35)$$

We take $\theta_p = [1; 1; 0; 1; -1; -2; 0; 1; 1]^T$, $\sigma^2 = 10^{-1}$, $N = 100$. The inputs $\{x_1^k\}$ and $\{x_2^k\}$ of the data set are uniformly distributed in $[-3; 3]$. As for process #2, the assumed model is true, i.e. the neural network associated to Eq. (35) is used; the minimization of the cost function with



Fig. 7. Accuracy of the LTE approximation of the variance for process #3: (a) reference estimate of the variance of $f(x, \Theta_{LS})$; and (b) difference between the reference estimate of the variance of $f(x, \Theta_{LS})$ and the LTE approximation obtained with $\theta_p$ and $\sigma^2$.

the Levenberg algorithm leads to $s^2 = 9.73 \times 10^{-2}$. Fig. 5a shows the inputs of the data set and the regression; Fig. 5b displays the LTE estimate of the variance of $f(\boldsymbol{x}, \boldsymbol{\Theta}_{LS})$ (27); Fig. 5c displays the difference between the reference variance estimate (33) computed over $M = 10\,000$ sets, and the LTE variance estimate (27). As $s^2$ is slightly smaller than $\sigma^2$, the variance is globally slightly underestimated. However except in the domain around the corner (3, 3) where the density of the inputs is lower, and where the slope of the output surface is steep, the LTE variance estimate is satisfactory. A reliable estimate of the CI may thus be obtained.

### 3.4.3. Accuracy of the LTE variance approximation (processes #2 and #3)

Let us now show on the example of processes #2 and #3 that the curvature of the solution surface is small enough for the LTE approximation of the variance (23) to be satisfactory. For both processes, we have computed approximation (23), using the values of $\xi$ and of $\boldsymbol{\xi}$ (at $\boldsymbol{\theta}_p$) and the value of $\sigma^2$ used for the noise simulation. As shown in Fig. 6a and b for process #2, the LTE approximation of the variance (23) is very close to the reference variance estimate. As a matter of fact, the difference between them (Fig. 6b) is only due to the curvature, which is small as $N = 50$ is large with respect to the complexity of the regression. Expression (23) also leads to satisfactory results in the case of process #3 as shown in Fig. 7a and b ($N = 100$, two inputs). This tends to show that a first-order approximation of the variance is often sufficient, and that it is not worth to bother with a higher-order approximation. Seber and Wild (1989) introduces a quadratic approximation of the LS estimator using the curvature of the solution surface. This approximation uses arrays of projected second derivatives, the intrinsic and parameter-effects curvatures arrays; but their presentation is beyond the scope of this paper.

## 4. Confidence intervals for neural networks

In the previous sections, the model used for the construction of CIs is true. For real-world black-box modeling problems however, a family of functions that contains the regression is not known a priori. The first task is thus to select the less complex family of functions which contains a function approximating the regression to a certain degree of accuracy in the input domain delimited by the data set. In practice, several families of increasing complexity (for example neural networks with one layer of an increasing number $n_h$ of hidden units, and a linear output neuron) are considered, and the data set is used both to estimate their parameters, and to perform the selection between the candidates. In order to retain the less complex family containing a good approximation of the regression, it is of interest to perform the selection only between neural candidates which are not unnecessarily large, and which are (that is their matrix $z$ is) sufficiently well-conditioned to allow the

computation of the approximate CI (32). We propose to discard too large models by a systematic detection of ill-conditioning, and to perform the selection among the approved, i.e. well-conditioned models using an approximate value of their LOO score whose computation does not require further training. Both the ill-conditioning detection and the estimation of the LOO score of a neural candidate are based on the LTE of its output.

### 4.1. Ill-conditioning detection for model approval

A too large neural model, trained up to convergence with a simple LS cost-function, will generally overfit. Overfitting is often avoided by using early stopping of the training algorithm or by adding regularization terms in the cost function, e.g. "weight decay" (Bishop, 1995). Unfortunately, as only the estimator whose value corresponds to an absolute minimum of the quadratic cost function (18) without weight decay terms is unbiased, both early stopping and weight decay introduce a bias in the estimation of the regression: the corresponding estimates thus lead to questionable CIs for the regression.

To detect and discard too large networks, we propose, after the training of each candidate up to a (hopefully) absolute minimum of the cost function (18), to check the conditioning of their matrix $z$ (see Rivals & Personnaz, 1998). The fact that $z$ be ill-conditioned is the symptom that some parameters are useless, since the elements of $z$ represent the sensibility of the model output with respect to the parameters. A typical situation is the saturation of a "tanh" hidden neuron, a situation which generates in the matrix $z$ a column of $+1$ or $-1$ that corresponds to the parameter between the output of the saturated hidden neuron and the linear output neuron, and columns of zeros that correspond to the parameters between the network inputs and the saturated hidden neuron.[8] In practice, we propose to perform a singular value factorization of $z$, and to compute its condition number, that is the ratio of its largest to its smallest singular value, see, e.g. Golub & Van Loan, 1983. The matrix $z$ can be considered as very ill-conditioned when its condition number reaches the inverse of the computer precision, which is of the order of $10^{-16}$.

Further, in order to be able to compute the approximate CI (32) which involve $(z^T z)^{-1}$, the cross-product Jacobian matrix $z^T z$ must also be well conditioned. As the condition number of $z^T z$ is the square of the condition number of $z$, the networks whose matrix $z$ has a condition number much larger than $10^8$ cannot be approved.

There are other studies of the ill-conditioning of neural networks, but they deal with their training rather than with their approval, like in the work by Zhou and Si (1998) where an algorithm avoiding the Jacobian rank deficiency is

---

[8] Such a situation might also correspond to a relative minimum; to check the conditioning of $z$ helps thus also to discard neural networks trapped in relative minima, and leads to retrain the neural candidate starting from different initial weights.
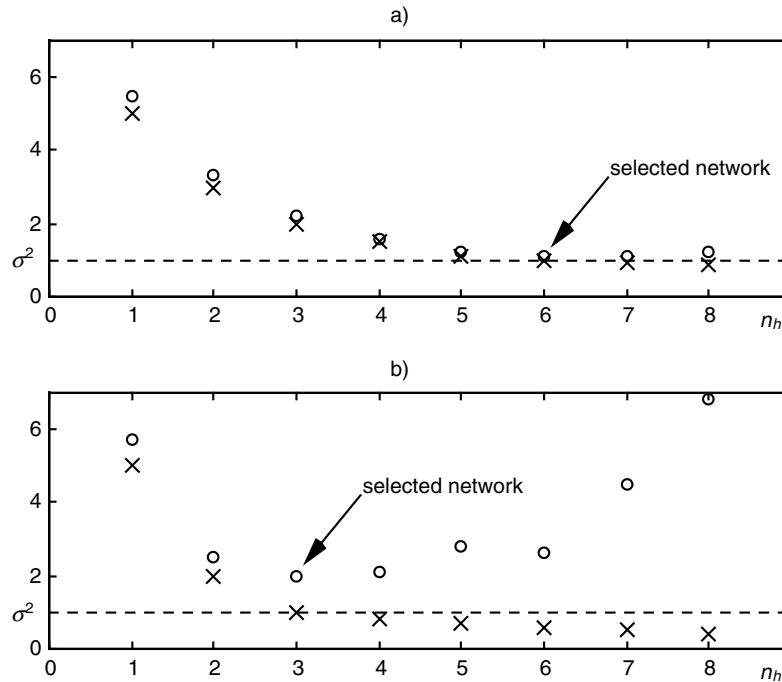
Fig. 8. Schematic evolution of the *MSTE* (crosses) and *MSPE* (circles) as a function of the number of hidden neurons of the neural network candidates, the network with the smallest *MSPE* being selected: (a) large data set: the ratio *MSPE/MSTE* of the selected network (six hidden neurons) is roughly equal to 1, hint that the data set size *N* is large; and (b) small data set: the ratio *MSPE/MSTE* of the selected network (three hidden neurons) is roughly equal to 2, hint that the data set size *N* is small.

presented, or by Saarinen, Bramley and Cybenko (1993) where the Hessian rank deficiency is studied during the training. In our view, rank deficiency is not very relevant *during the training* because with a Levenberg algorithm, the matrix to be "inverted" is made well conditioned by the addition of a scalar matrix $\lambda \, l_q$, $\lambda > 0$, to the cross-product Jacobian.

### 4.2. Approximate leave-one-out scores for model selection

The selection among the networks which have been approved can be performed with statistical tests (Rivals & Personnaz, 1998; Urbani, Roussel-Ragot, Personnaz & Dreyfus, 1994). Another approach, cross validation, consists in partitioning the data set in training and test sets, and in selecting the smallest network leading to the smallest mean square error on the test sets.[9] One of the drawbacks of cross validation is to require a successful training of the candidate models on many test sets, that is $N$ successful trainings in the case of LOO cross validation. Let us denote by $e^k$ the error obtained on the left out example $k$ with the model trained on the $N - 1$ remaining examples ($k$th test set). In this section, we derive an approximate expression of $e^k$, expression which allows an economic estimation of the LOO score

without performing these $N$ time-consuming trainings of each candidate network, as proposed in Monari (1999) and Monari and Dreyfus (submitted)).

In the case of a linear model, it is well known (Efron & Tibshirani, 1993) that the $k$th LOO error $e^k$ can be directly derived from the corresponding residual $r^k$:

$$e^k = \frac{r^k}{1 - [p_x]_{kk}} \quad k = 1 \text{ to } N \qquad (36)$$

where, we recall, $p_x$ denotes the orthogonal projection matrix on the range of $x$. Expression (36) holds irrespective of whether or not the assumed model is true.

In the case of a nonlinear model, we show (see Appendix B) that a useful approximation of the $k$th LOO error can be obtained using the LTE of the model output at $\boldsymbol{\theta}_{LS}$:

$$e^k \approx \frac{r^k}{1 - [p_z]_{kk}} \quad k = 1 \text{ to } N \qquad (37)$$

where $p_z$ denotes the orthogonal projection matrix on the range of $z$. The approximation (37) is thus similar to Eq. (36).[10] Like in the linear case, expression (37) holds independently on the assumed model being true or not. Hence the LOO score:

$$\text{LOO score} = \frac{1}{N} \sum_{k=1}^{N} (e^k)^2 \qquad (38)$$

---

[9] Note that statistical tests may advantageously be used complementarily to cross validation in order to take a decision (Rivals and Personnaz, 1999); these tests can also be established by applying LS theory to the LTE of nonlinear models (Bates and Watts, 1988), but this exceeds the scope of this paper.

[10] An expression similar to (37) is proposed in Hansen and Larsen (1993), but unfortunately, it is not valid even in the linear case.
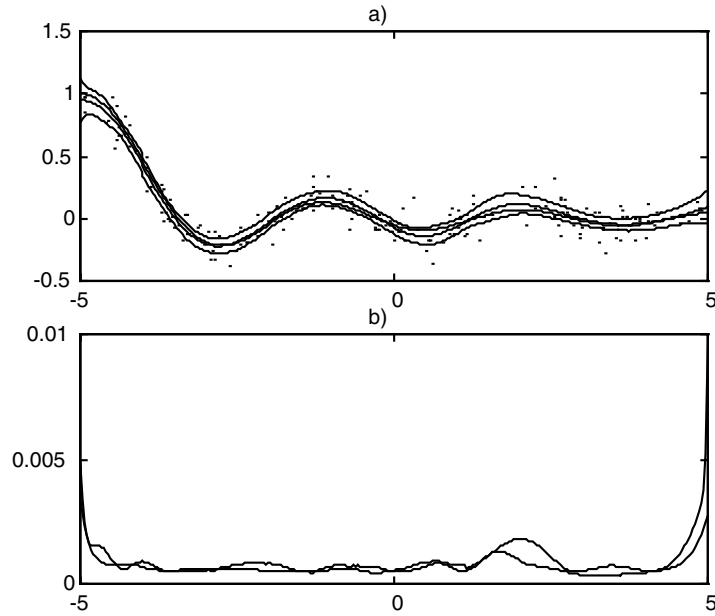
Fig. 9. CI for process #4, a simulated nonlinear SISO process, in the case of a data set of size $N = 200$ (the selected model is a four hidden neurons network with $q = 13$ parameters): (a) regression (thin line), data set (small points), model output and 99% approximate CI (thick lines); and (b) reference (thin line) and LTE (thick line) estimates of the variance of $f(x, \theta_{LS})$.

This LOO score can be used as an estimate of the mean square performance error, and we thus denote it as MSPE, as opposed to $(2/N)J(\theta_{LS})$, the mean square training error (MSTE). The interested reader will find in the work of Monari (1999) and Monari and Dreyfus (submitted) a systematic model selection procedure based on both the approximate LOO score and the distribution of the values of the $\{[p_z]_{kk}\}$. Nevertheless, another performance measure could be chosen as well (a 10-fold cross validation score, a mean square error on an independent set, etc.).

### 4.3. Accuracy of the approximate confidence intervals

The quality of the selected model $f(x, \theta_{LS})$, and thus of the associated approximate CI, depends essentially on the size $N$ of the available data set with respect to the complexity of the unknown regression function and to the noise variance $\sigma^2$.

1. $N$ is large: it is likely that the selected family $\{f(x, \theta),\ \theta \in \mathbb{R}^q\}$ contains the regression $E(Y_p|x)$, i.e. that the LS estimator is asymptotically unbiased, that the model $f(x, \theta_{LS})$ is a good approximation of $E(Y_p|x)$ in the domain delimited by the dataset, and that the curvature is small. In this case, reliable CIs can be computed with Eq. (32).
2. $N$ is small: it is likely that the selected family $\{f(x, \theta),\ \theta \in \mathbb{R}^q\}$ is too small[11] to contain $E(Y_p|x)$, i.e. that the LS estimator is biased, and that the model

---

[11] It will generally not be too large since the approval procedure proposed in Section 4.1 prevents from selecting a neural network with useless parameters.

$f(x, \theta_{LS})$ thus underfits. The approximate CIs are thus questionable, and additional data should be gathered.

A good indicator of whether the data set size $N$ is large enough or not is the ratio *MSPE/MSTE* of the selected candidate: if its value is close to 1, then $N$ is probably large enough, whereas a large value is the symptom of a too small data set size, as shown in Fig. 8 (and as illustrated numerically in the following examples).

### 4.4. Example of a simulated nonlinear SISO process (process #4)

This first example is based on a simulated process. Like in the previous sections, a reference estimate of the variance of the output of a neural network is made, using $M = 1000$ other sets; to ensure that an absolute minimum is reached on each of the $M$ sets, 5–30 trainings (depending on the network size) with the Levenberg algorithm for different initializations of the weights are performed, and the weights giving the smallest value of the cost function (18) are kept. We consider the SISO process simulated with:

$$y_p^k = \text{sinc}\,(2(x^k + 5)) + w^k \quad k = 1 \text{ to } N \tag{39}$$

where sinc denotes the cardinal sine function; we take $\sigma^2 = 10^{-2}$.

First, a data set of $N = 200$ input–output pairs is computed, with input values uniformly distributed in $[-5; 5]$. As a family of nonlinear functions (a neural network with a given architecture) containing the regression is not known a priori, neural networks with a linear output neuron and a layer of $n_h$ "tanh" hidden neurons are trained. The numerical
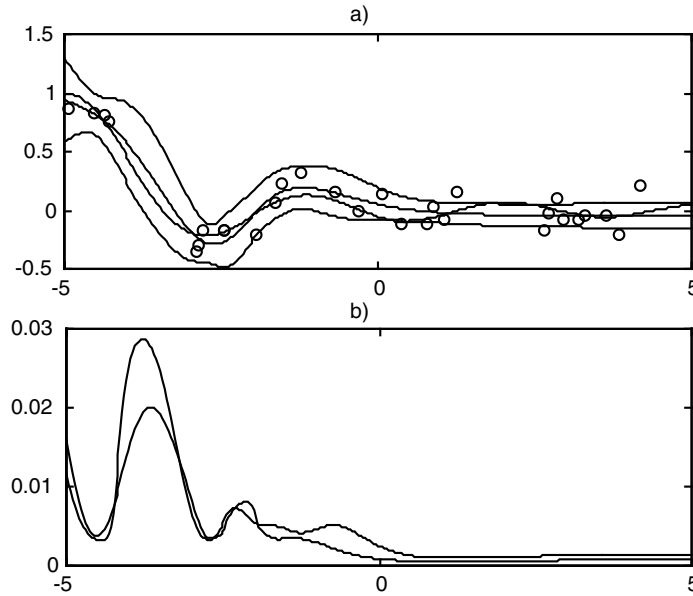
Fig. 10. CI for process #4, a simulated nonlinear SISO process, in the case of a data set of size $N = 30$ (the selected model is a two hidden neurons network with $q = 7$ parameters): (a) regression (thin line), data set (circles), model output and 99% approximate CI (thick lines); and (b) reference (thin line) and LTE (thick line) estimates of the variance of $f(\boldsymbol{x}, \boldsymbol{\Theta}_{\mathrm{LS}})$.

results are summarized in Table 1. We list the number of parameters $q$, the MSTE (i.e. the smallest MSTE obtained with the network for its different random weight initializations), the condition number of $z$, and, if the latter is not too large, the MSPE (corresponding approximate LOO score computed with Eqs. (37) and (38)) and the ratio MSPE/MSTE. The candidates with more than six hidden neurons cannot be approved, because $\mathrm{cond}(z) \gg 10^8$: for $n_{\mathrm{h}} = 7$, $\mathrm{cond}(z) = 10^{11}$. The optimal number of neurons $n_{\mathrm{h}}^{\mathrm{opt}}(200) = 4$ is selected on the basis of the MSPE. The fact that the corresponding ratio MSPE/MSTE is close to 1 is the symptom that $N$ is large enough, so that the selected family of networks contains a good approximation of the regression, and that the curvature is small (case 1 of Section 4.3). The results obtained for the selected neural network are shown in Fig. 9. The model output is close to the regression, the LTE variance estimate (27) is close to the reference variance estimate (33), and the CI is thus accurate.

Second, a data set of $N = 30$ input–output pairs is computed, the numerical results being summarized in Table 2. The data set being much smaller, the candidates cannot be approved as soon as $n_{\mathrm{h}} > 4$: for $n_{\mathrm{h}} = 5$, $\mathrm{cond}(z) = 10^{15}$. The optimal number of neurons $n_{\mathrm{h}}^{\mathrm{opt}}(30) = 2$ is selected on the basis of the MSPE. The ratio MSPE/MSTE of the selected network equals 2.1, symptom that $N$ is relatively small, and that the selected family of networks is likely not to contain the regression (case 2 of Section 4.3). The results obtained for the selected neural network are shown in Fig. 10. The family of functions implemented by a network with two hidden units is obviously too small to contain a good approximation of the regression, and though the estimate of the output variance is good (it is close to the reference variance estimate), since the output of the neural network differs from the regression, the CI is less accurate than in the case where $N = 200$. Note that in the input domain [0, 5] where the model underfits, the variance remains constant and low. This is due to the fact that, in this domain, the model output is insensitive to most parameters of the network (this is usually the case when,

Table 1
Results obtained on the modeling of the simulated SISO process #4 using neural networks, in the case $N = 200$

| $n_{\mathrm{h}}$ | $q$ | MSTE | Cond($z$) | MSPE | MSPE/MSTE |
|---|---|---|---|---|---|
| 1 | 4 | $1.4 \times 10^{-2}$ | 10 | $1.4 \times 10^{-2}$ | 1.0 |
| 2 | 7 | $1.2 \times 10^{-2}$ | $10^3$ | $1.3 \times 10^{-2}$ | 1.1 |
| 3 | 10 | $9.7 \times 10^{-3}$ | $10^6$ | $1.1 \times 10^{-2}$ | 1.1 |
| **4** | **13** | $\mathbf{8.5 \times 10^{-3}}$ | $\mathbf{10^2}$ | $\mathbf{9.8 \times 10^{-3}}$ | **1.1** |
| 5 | 16 | $8.4 \times 10^{-3}$ | $10^6$ | $9.9 \times 10^{-3}$ | 1.2 |
| 6 | 19 | $8.2 \times 10^{-3}$ | $10^7$ | $1.0 \times 10^{-2}$ | 1.2 |
| 7 | 22 | $7.9 \times 10^{-3}$ | $10^{11}$ | – | – |

Table 2
Results obtained on the modeling of the simulated SISO process #4 using neural networks, in the case $N = 30$

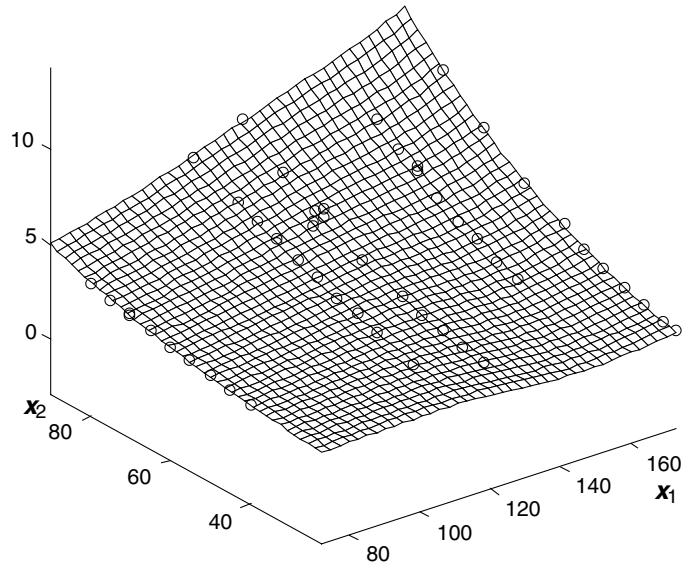| $n_{\mathrm{h}}$ | $q$ | MSTE | Cond($z$) | MSPE | MSPE/MSTE |
|---|---|---|---|---|---|
| 1 | 4 | $2.4 \times 10^{-2}$ | $10^1$ | $2.7 \times 10^{-2}$ | 1.1 |
| **2** | **7** | $\mathbf{1.1 \times 10^{-2}}$ | $\mathbf{10^6}$ | $\mathbf{2.3 \times 10^{-2}}$ | **2.1** |
| 3 | 10 | $8.1 \times 10^{-3}$ | $10^3$ | $2.4 \times 10^{-2}$ | 3.0 |
| 4 | 13 | $7.1 \times 10^{-3}$ | $10^4$ | $4.3 \times 10^1$ | $6.1 \times 10^3$ |
| 5 | 16 | $5.0 \times 10^{-3}$ | $10^{15}$ | – | – |

Fig. 11. Industrial modeling problem (the selected model is a two hidden neurons network with $q = 9$ parameters): model output, and the $N = 69$ examples of the data set (circles).
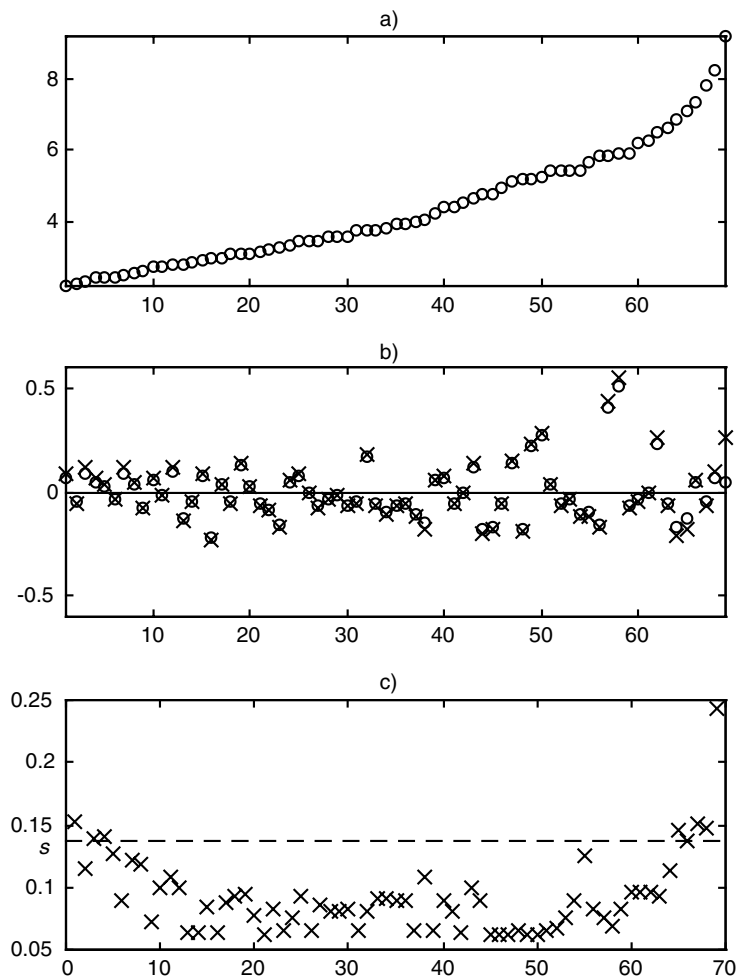


Fig. 12. Industrial modeling problem (a) the $N = 69$ outputs of the data set presented in increasing order of their values; (b) the corresponding residuals (circles) and approximate LOO errors (crosses); and (c) half width of the 95% approximate Cl at the $N = 69$ examples of the data set, and LS estimate $s$ of the noise standard-deviation $s$ (dotted line).
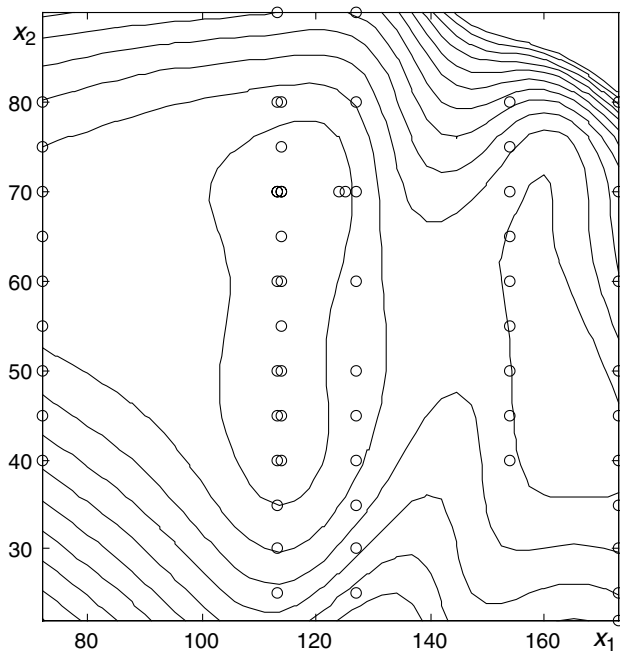
Fig. 13. Industrial modeling problem: isocontours of the LTE estimate of the standard deviation of $f(\mathbf{x}, \boldsymbol{\Theta}_{\mathrm{LS}})$, and the $N = 69$ inputs of the data set (circles).

Table 3
Results obtained on the modeling of the industrial process using neural networks

| $n_{\mathrm{h}}$ | $q$ | MSTE | Cond($z$) | MSPE | MSPE/MSTE |
|---|---|---|---|---|---|
| 1 | 5 | $5.2 \times 10^{-2}$ | $10^4$ | $6.6 \times 10^{-2}$ | 1.3 |
| **2** | **9** | $\mathbf{1.6 \times 10^{-2}}$ | $\mathbf{10^5}$ | $\mathbf{2.1 \times 10^{-2}}$ | **1.3** |
| 3 | 13 | $1.5 \times 10^{-2}$ | $10^4$ | $1.7 \times 10^{-1}$ | $1.1 \times 10^1$ |
| *4* | *17* | *$1.2 \times 10-2$* | *$10^{12}$* | *–* | *–* |

like here, the output of a network does not vary[12]): the elements of the $z$s in this domain are thus constant and small, hence a small and constant variance at the corresponding $\mathbf{x}$s.

### 4.5. Industrial modeling problem

We apply here the presented methodology (LS parameter estimation, model approval, model selection, CI construction) to an industrial example first tackled in (Rivals & Personnaz, 1998), that is the modeling of a mechanical property of a complex material from three structural descriptors. We have been provided with a data set of $N = 69$ examples; the inputs and outputs are normalized for the LS estimations. Thanks to repetitions in the data, and assuming homoscedasticity, the "mean square pure error" (Draper & Smith, 1998) gives a good estimate of the noise variance: $\widehat{\sigma^2} = 3.38 \times 10^{-2}$. Using this reliable estimate, statistical tests establish the significance of two inputs. An affine model with these two inputs gives the estimate $s^2 = 2.38 \times 10^{-1}$ of the variance, hence the necessity of nonlinear modeling.

Neural networks with a linear output neuron and a layer of $n_{\mathrm{h}}$ "tanh" hidden neurons are trained. The numerical results are summarized in Table 3. It shows that the candi-

dates with more than three hidden neurons cannot be approved: for $n_{\mathrm{h}} = 4$, cond($z$) $= 10^{12}$. The optimal number of neurons $n_{\mathrm{h}}^{\mathrm{opt}}(69) = 2$ is selected on the basis of the MSPE. The fact that the corresponding ratio MSPE/MSTE equals 1.3 indicates that $N$ is large enough, so that the selected family of networks contains probably a good approximation of the regression, and that the curvature is small (case 1 of Section 4.3). The function implemented by the selected network is shown in Fig. 11.

The $N = 69$ output values of the training set are presented in the increasing order in Fig. 12a, and the corresponding residuals and approximate LOO errors in Fig. 12b: both appear quite uncorrelated and homoscedastic. A CI with a level of significance of 95% is then computed with Eq. (32); the half width of the 95% CI on the $N = 69$ examples of the data set is shown in Fig. 12c. In order to check the confidence, which can be attached to the model, the variance of its output must be examined in the whole input domain of interest. Fig. 13 shows the isocontours of the LTE estimate of the standard deviation of the model output $s\sqrt{z^{\mathrm{T}}(z^{\mathrm{T}}z)^{-1}z}$ in the input domain defined by the training set. The computation of the LTE variance estimate thus allows not only to construct a CI at any input of interest, but also to diagnose that, at the top right corner of the input domain, the model standard deviation is larger than that of the noise itself (the highest isocontour value equals that of the estimate of the noise standard deviation $s = 1.39 \times 10^{-1}$). Little confidence can thus be attached to the model output in this input domain, where more data should be gathered. On the contrary, there is a large region on the left of the diagram where there are very few training examples, but where the LTE estimate of the standard deviation is surprisingly rather small; like for the modeling of process #4 in Section 4.4, this is due to the fact that the model output is less sensitive to most parameters of the network in this region (the model output varies very little, see Fig. 11).

## 5. Comparisons

In this section, we discuss the advantages of the LS LTE approach to the construction of confidence intervals for neural networks with respect to other analytic approaches and to the bootstrap methods, and compare them on simulated examples.

---

[12] The output of a neural network with one layer of tanh hidden units remains constant in a given domain of its inputs when the "tanh" activation functions of all hidden units are saturated in this domain: the output of the network is thus insensitive to all the parameters of the hidden units.
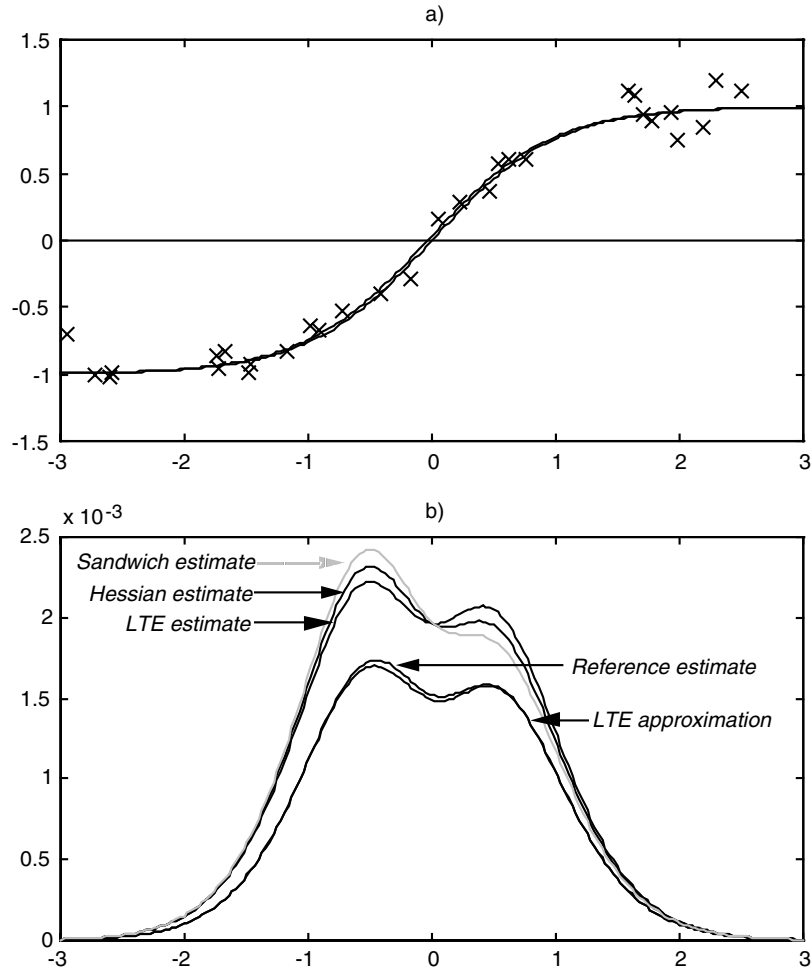
Fig. 14. Comparison of different estimates of the variance of a nonlinear model output for process #5, a simulated "neural" SISO process (the assumed model, a single nonlinear neuron with $q = 2$ parameters, is true): (a) regression with a "gentle" slope (thin line), the $N = 30$ examples of the data set (crosses), model output (thick line); and (b) LTE approximation and estimates of the variance of $f(\mathbf{x}, \boldsymbol{\Theta}_{LS})$.

## 5.1. Comparison to other analytic approaches

### 5.1.1. Maximum likelihood approach

In the case of gaussian homoscedastic data, likelihood theory leads to the same approximate variance (23), but two different estimators of it are commonly encountered (see Appendix A.2):

$$var\,(\widehat{f(\mathbf{x}, \boldsymbol{\Theta}_{LS})})_{LTE} = \widehat{\sigma^2}\, z^T (z^T z)^{-1} z$$

i.e. the same estimate as Eq. (27), and also:

$$var\,(\widehat{f(\mathbf{x}, \boldsymbol{\Theta}_{LS})})_{Hessian} = \widehat{\sigma^2}\, z^T [h(\boldsymbol{\theta}_{LS})]^{-1} z \qquad (40)$$

which necessitates the computation of the Hessian. Efficient methods for computing the Hessian are presented by Buntine and Weigend (1994).

### 5.1.2. Bayesian approach

The Bayesian approach is an alternative approach to the sampling theory (or the frequentist approach) for modeling problems, and also leads to the design of CIs. These two approaches are conceptually very different: the Bayesian approach treats the unknown parameters as random variables, whereas they are considered as certain in the frequentist approach. Nevertheless, as outlined by Bishop (1995), MacKay (1992a,b), the Bayesian approach leads to a posterior distribution of the parameters with a covariance matrix whose expression is very similar to that of the covariance matrix of the least-squares estimator of the parameters, and thus to CIs which are similar to those presented in this paper. We thus make a brief comparison between the CIs these two approaches lead to.

The most important difference is that the estimator which is considered here is the one whose estimate minimizes the cost function (18), whereas in the Bayesian approach, a cost-function with an additional weight-decay regularization term is minimized; the presence of this weight-decay term stems from the assumption of a gaussian prior for the parameters.

Nevertheless, the least squares cost function (18) can be seen as the limit where the regularization term is zero, which corresponds to an uninformative prior for the parameters. In
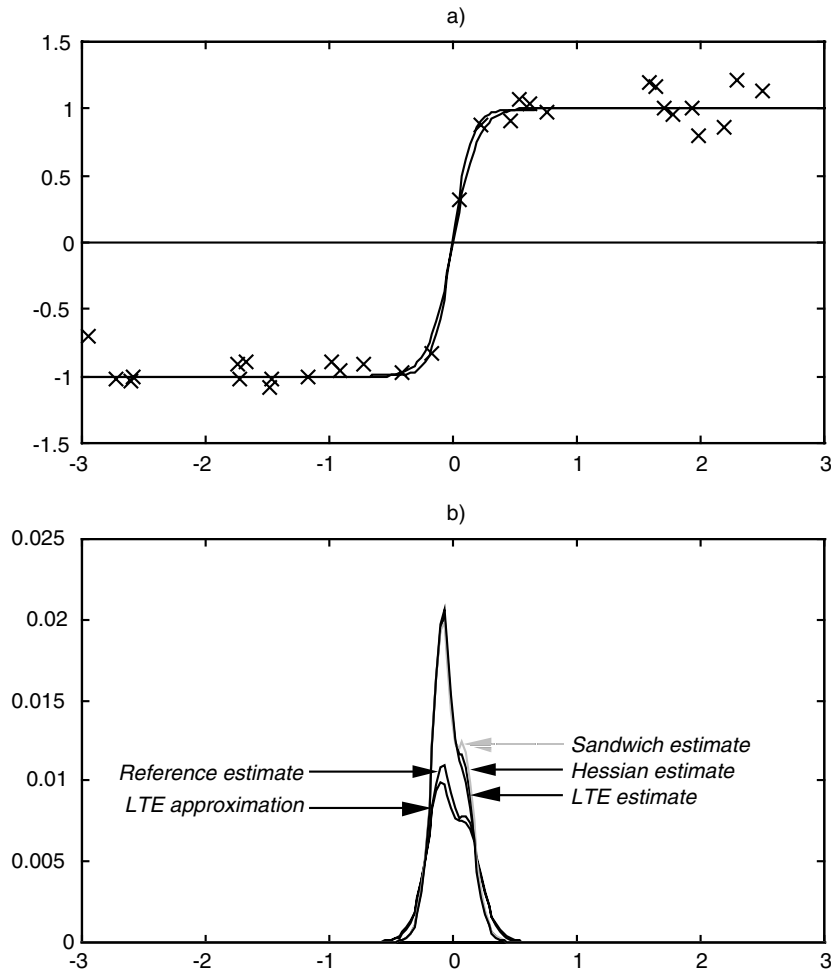
Fig. 15. Comparison of different estimates of the variance of a nonlinear model output for process #6, a simulated "neural" SISO process (the assumed model, a single nonlinear neuron with $q = 2$ parameters, is true): (a) regression with a "steep" slope (thin line), the $N = 30$ examples of the data set (crosses), model output (thick line); and (b) LTE approximation and estimates of the variance of $f(\mathbf{x}, \boldsymbol{\Theta}_{\mathrm{LS}})$.

this case (that is Eq. (18) is minimized as in this paper), there is another small difference in the Bayesian approach as outlined by Bishop (1995) and MacKay (1992a,b). Under hypotheses which we cannot detail here, the Bayesian approach leads to a posterior parameter distribution with the approximate covariance matrix $\sigma^2[h(\boldsymbol{\theta}_{\mathrm{LS}})]^{-1}$, $h(\boldsymbol{\theta}_{\mathrm{LS}})$ being the Hessian of the cost function evaluated at the most probable value of the parameter, that is here $\boldsymbol{\theta}_{\mathrm{LS}}$. A LTE of the estimator output leads then to the following estimate of its variance at input $\mathbf{x}$:

$$\mathrm{var}\,(\widehat{f(\mathbf{x}, \boldsymbol{\Theta}_{\mathrm{LS}})})_{\mathrm{Hessian}} = \widehat{\sigma^2}\, z^{\mathrm{T}}[h(\boldsymbol{\theta}_{\mathrm{LS}})]^{-1}z$$
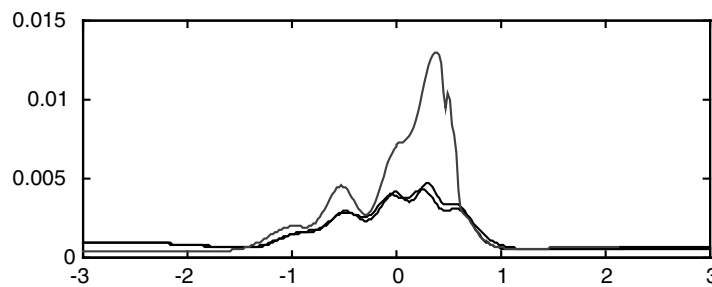
i.e. it also leads to estimate (40).



Fig. 16. Comparison of the LS LTE and bootstrap pairs approach estimates of the variance for process #2. reference (thin line), LTE (thick line), and bootstrap pairs (dotted line) estimates of the variance of $f(\mathbf{x}, \boldsymbol{\Theta}_{\mathrm{LS}})$.

### 5.1.3. Sandwich estimator

The sandwich estimate of the variance of a nonlinear model output can be derived in various frameworks (a possible derivation in the frequentist approach is given in Appendix A.3):

$$\text{var }(\widehat{f(\boldsymbol{x}, \boldsymbol{\theta}_{\text{LS}})})_{\text{sandwich}} = \widehat{\sigma^2}\boldsymbol{z}^{\text{T}}[h(\boldsymbol{\theta}_{\text{LS}})]^{-1}\boldsymbol{z}^{\text{T}}\boldsymbol{z}[h(\boldsymbol{\theta}_{\text{LS}})]^{-1}\boldsymbol{z} \quad (41)$$

The sandwich estimator is known to be robust to model incorrectness, i.e. the considered family of functions is too small (see, e.g. Efron & Tibshirani, 1993; Ripley, 1995).

### 5.1.4. Numerical comparison (processes #5 and #6)

Here, we perform a numerical comparison of the three variance estimates considered above on a very simple example. We consider a SISO process simulated by a single "tanh" neuron:

$$y_p^k = \tanh(\theta_{p_1} + \theta_{p_2}x^k) + w^k \quad k = 1 \text{ to } N \quad (42)$$

with $\sigma^2 = 0.01$, $N = 30$. For this comparison, the noise variance $\sigma^2$ is estimated with $s^2$ in the three (LTE, Hessian, and sandwich) output variance estimates.

We first simulate a process with $\theta_{p_1} = 0$, $\theta_{p_2} = 1$ (process #5). The corresponding results are shown in Fig. 14. The variance reference estimate is computed on $M = 10\,000$ data sets. The LTE approximation (23) of the variance is almost perfect. The LTE (27), Hessian (40), and sandwich (41) estimates are comparable: the parameter estimates being accurate ($\theta_{\text{LS}_1} = 3.63\ 10^{-2}$, $\theta_{\text{LS}_2} = 0.996$), the fact that they are overestimated is almost only due to the noise variance estimate $s^2 = 1.32 \times 10^{-2}$. Nevertheless, the shape of the LTE estimate is closer to the reference estimate than that of the two others.

We then simulate a process with: $\theta_{p_1} = 0$, $\theta_{p_2} = 5$ (process #6). The corresponding results are shown in Fig. 15. The function being steeper, the curvature is larger, and the LTE approximation (23) of the variance is a little less accurate. The three estimates are still very similar but, here, their overestimation is due not only to the noise variance estimate $s^2 = 1.25 \times 10^{-2}$, but also to the bias of the parameter estimates ($\theta_{\text{LS}_1} = 3.79\ 10^{-2}$, $\theta_{\text{LS}_2} = 6.58$).

The computational cost of the LTE estimate being lower (is does not necessitate the computation of the Hessian matrix), there is no reason to prefer one of the two other estimates. As a matter of fact, since the Hessian depends on the data set, it is the realization of a random matrix. Thus, in the maximum likelihood as well as in the Bayesian approach, it is often recommended to take the expectation of the Hessian, and to evaluate it at the available $\boldsymbol{\theta}_{\text{LS}}$, i.e. to replace it by the cross-product Jacobian $\boldsymbol{z}^{\text{T}}\boldsymbol{z}$ (Seber & Wild, 1989): estimates (40) and (41) then reduce to estimate (27). As mentioned above, the sandwich variance estimator is known to be robust to model incorrectness, a property which is not tested with this simple setting, but this is beyond the scope of this paper.

### 5.2. Comparison to bootstrap approaches

The bootstrap works by creating many pseudo replicates of the data set, the bootstrap sets, and reestimating the LS solution (retraining the neural network) on each bootstrap set; the variance of the neural model output, and the associated CI, are then computed over the trained networks, typically *a hundred* (Efron & Tibshirani, 1993). In the "bootstrap pairs approach" for example, a bootstrap set is created by sampling with replacement from the data set (Efron & Tibshirani, 1993). The first advantage of the LS LTE approach is to require only *one* successful training of the network on the data set to compute the LTE estimate of the variance of its output, whereas the bootstrap methods require a hundred successful trainings of the network on the different bootstrap sets.

Studies on bootstrap where only one training with a random initialization of the weights was performed for each bootstrap set show a pathological overestimation of the variance. This can be seen in Tibshirani (1996), examples 2 and 3; but the overestimation of the bootstrap is not detected in this work because the reference estimate is also overestimated for the same reasons (one single training per set). As pointed out by Refenes, Zapranis and Utans (1997), a way to reduce this overestimation is to start each training on a bootstrap set with the weights giving the smallest value of the cost function (18) (that is on the original data set); but even so, the bootstrap method becomes untractable for large networks, and/or for multi input processes.

The claim that bootstrap methods are especially efficient for problems with small data sets (see, e.g. Heskes, 1997) may be subject to criticism. As an illustration, the variance was estimated for process #2 using the bootstrap pairs approach on 300 bootstrap sets, the network weights being initialized twice for each training, once with the true ones, and once with those obtained by training the network on the whole data set. As shown in Fig. 16, though the size of the data set is not very small ($N = 50$), the bootstrap variance estimate is far away from the reference estimate. Increasing the number of bootstrap sets up to 1000 did not improve the variance estimate.

In fact, the bootstrap is especially suited to the estimation of the variance of estimators defined by a formula, like for example an estimator of a correlation coefficient (Efron & Tibshirani, 1993): for each bootstrap set, an estimate is computed using the formula, and the estimate of the variance is easily obtained. However the bootstrap is definitely not the best method if each estimation associated to a bootstrap set involves an iterative algorithm like the training of a neural network, which is the case for the construction of a CI with a neural model. However, if the data set is large enough, and if the training time is considered unimportant, the bootstrap pairs approach is a solution in the case of heteroscedasticity (that is $K(\boldsymbol{W})$ is not scalar anymore), whereas the LS LTE approach, as well as the "bootstrap residuals" approach (Efron & Tibshirani, 1993), are no longer valid.

## 6. Conclusion

We have given a thorough analysis of the LS LTE approach to the construction of CIs for a nonlinear regression using neural network models, and put emphasis on its enlightening geometric interpretation. We have stressed the underlying assumptions, in particular the fact that the approval and selection procedures must have led to a parsimonious, well-conditioned model containing a good approximation of the regression. Our whole methodology (LS parameter estimation, model approval, model selection, CI construction) has been illustrated on representative examples, bringing into play simulated processes and an industrial one.

We have also shown that, as opposed to the computationally intensive bootstrap methods, the LS LTE approach to the estimation of CIs is both accurate and economical in terms of computer power, and that it leads to CIs which are comparable to those obtained by other analytic approaches under similar assumptions, at a lower computational cost.

A rigorous assessment of the accuracy of the results obtained with the LS LTE approach, as well as with any statistical approach dealing with nonlinear models and assuming the local planarity of the solution surface, remains an open problem: it could be enlightened by a specific study of the curvature of the solution surface of neural networks.

## Acknowledgements

We thank an anonymous reviewer whose constructive comments contributed to improve the quality of this paper. We also thank Howard Gutowitz whose advice we greatly appreciated.

## Appendix A. Estimates of a nonlinear model output variance

In order to make this paper self-contained, we provide derivations of the different variance estimates.

### A.1. LTE variance estimate in sampling theory

The well-known approximation (Seber & Wild, 1989) we use in this paper is based on a single expansion, the LTE of the nonlinear model output for an input $x$ at the true parameter value $\theta_p$:

$$f(x, \theta) \approx f(x, \theta_p) + \boldsymbol{\xi}^{\mathrm{T}}(\theta - \theta_p) \tag{A1}$$

This expansion leads, for the data set, to

$$\boldsymbol{f}(x, \theta) \approx \boldsymbol{f}(x, \theta_p) + \xi(\theta - \theta_p) \tag{A2}$$

We now use Eq. (A2) in the expression of the cost-function

$$J(\theta) = \tfrac{1}{2}(y_p - \boldsymbol{f}(x, \theta))^{\mathrm{T}}(y_p - \boldsymbol{f}(x, \theta)) \tag{A3}$$

This leads to

$$
\begin{aligned}
J(\theta) &\approx \tfrac{1}{2}(y_p - \boldsymbol{f}(x, \theta_p) - \xi(\theta - \theta_p))^{\mathrm{T}} \\
&\quad \times (y_p - \boldsymbol{f}(x, \theta_p) - \xi(\theta - \theta_p)) \\
&= \tfrac{1}{2}(y_p - \boldsymbol{f}(x, \theta_p) + \xi\theta_p)^{\mathrm{T}}(y_p - \boldsymbol{f}(x, \theta_p) + \xi\theta_p) \\
&\quad - \theta^{\mathrm{T}}\xi^{\mathrm{T}}(y_p - \boldsymbol{f}(x, \theta_p) + \xi\theta_p) + \tfrac{1}{2}\theta^{\mathrm{T}}\xi^{\mathrm{T}}\xi\theta
\end{aligned}
$$

An approximate expression of the gradient of the cost-function follows

$$\frac{\partial J}{\partial \theta} \approx -\xi^{\mathrm{T}}(y_p - \boldsymbol{f}(x, \theta_p) + \xi\theta_p) + \xi^{\mathrm{T}}\xi\theta \tag{A4}$$

Hence an approximate expression of the least-squares estimate of the parameters

$$\theta_{\mathrm{LS}} \approx \theta_p + (\xi^{\mathrm{T}}\xi)^{-1}\xi^{\mathrm{T}}(y_p - \boldsymbol{f}(x, \theta_p)) \tag{A5}$$

And hence the corresponding approximation of the least-squares estimator (i.e. the random vector $\boldsymbol{\Theta}_{\mathrm{LS}}$) of the parameters (expression (21) in the main text)

$$\boldsymbol{\Theta}_{\mathrm{LS}} \approx \theta_p + (\xi^{\mathrm{T}}\xi)^{-1}\xi^{\mathrm{T}}(Y_p - \boldsymbol{f}(x, \theta_p)) = \theta_p + (\xi^{\mathrm{T}}\xi)^{-1}\xi^{\mathrm{T}}W \tag{A6}$$

Using the linear Taylor expansion (A1), we obtain an approximation of the variance of the LS estimator of the regression for any input $x$ (expression (23) in the main text):

$$\mathrm{var}(f(x, \boldsymbol{\Theta}_{\mathrm{LS}})) \approx \sigma^2 \boldsymbol{\xi}^{\mathrm{T}}(\xi^{\mathrm{T}}\xi)^{-1}\boldsymbol{\xi} \tag{A7}$$

The derivatives involved in $\boldsymbol{\xi}$ and $\xi$ being performed at the unknown $\theta = \theta_p$, they may be estimated by the derivatives at $\theta = \theta_{\mathrm{LS}}$, that is by replacing $\boldsymbol{\xi}$ by $z$ and $\xi$ by $z$. Hence the LTE variance estimate presented in the paper:

$$\widehat{\mathrm{var}(f(x, \boldsymbol{\Theta}_{\mathrm{LS}}))}_{\mathrm{LTE}} = s^2 z^{\mathrm{T}}(z^{\mathrm{T}}z)^{-1}z = \frac{r^{\mathrm{T}}r}{N - q}z^{\mathrm{T}}(z^{\mathrm{T}}z)^{-1}z \tag{A8}$$

### A.2. LTE variance estimates in maximum likelihood theory

For comparison, we sum up the results obtained with maximum likelihood theory (see, e.g. Efron & Tibshirani, 1993; Tibshirani, 1996). We make the same assumptions as for sampling theory, i.e. that the nonlinear assumed model is true and that $K(W) = \sigma^2 I_N$ (homoscedasticity), and we consider a gaussian distributed noise. In this case, the log likelihood function is:

$$L(\theta) = -\frac{1}{2\sigma^2}(y_p - \boldsymbol{f}(x, \theta))^{\mathrm{T}}(y_p - \boldsymbol{f}(x, \theta)) + \mathrm{cte} \tag{A9}$$

The parameters that maximize Eq. (A9) are those that minimize Eq. (A3), i.e. $\theta_{\mathrm{ML}} = \theta_{\mathrm{LS}}$.

It can be shown (Seber & Wild, 1989) that the covariance matrix of $\boldsymbol{\Theta}_{\mathrm{ML}}$ is given asymptotically by the inverse of the Fisher information matrix evaluated at $\theta_p$. The Fisher information matrix being the mathematical expectation of the random matrix $M(\theta)$ of the second derivatives of the log

likelihood function, we have

$$[M(\boldsymbol{\theta})]_{ij} = -\frac{\partial^2 L}{\partial \theta_i \partial \theta_j}$$

$$= \frac{1}{\sigma^2} \sum_{k=1}^{N} \left( \frac{\partial f(\boldsymbol{x}^k, \boldsymbol{\theta})}{\partial \theta_i} \frac{\partial f(\boldsymbol{x}^k, \boldsymbol{\theta})}{\partial \theta_j} + (Y_p^k - f(\boldsymbol{x}^k, \boldsymbol{\theta})) \frac{\partial^2 f(\boldsymbol{x}^k, \boldsymbol{\theta})}{\partial \theta_i \partial \theta_j} \right) \tag{A10}$$

The assumed model being true, i.e. $E(Y_p^k - f(\boldsymbol{x}^k, \boldsymbol{\theta}_p)) = E(W^k) = 0$, the Fisher information matrix evaluated at $\boldsymbol{\theta}_p$ is given by:

$$[E(M(\boldsymbol{\theta}_p))]_{ij} = -\frac{\partial^2 L}{\partial \theta_i \partial \theta_j}\bigg|_{\boldsymbol{\theta}=\boldsymbol{\theta}_p}$$

$$= \frac{1}{\sigma^2} \sum_{k=1}^{N} \left[ \frac{\partial f(\boldsymbol{x}^k, \boldsymbol{\theta})}{\partial \theta_i}\bigg|_{\boldsymbol{\theta}=\boldsymbol{\theta}_p} \frac{\partial f(\boldsymbol{x}^k, \boldsymbol{\theta})}{\partial \theta_j}\bigg|_{\boldsymbol{\theta}=\boldsymbol{\theta}_p} \right] \tag{A11}$$

$$E(M(\boldsymbol{\theta}_p)) = \frac{1}{\sigma^2} \xi^{\mathrm{T}} \xi$$

Thus, the covariance matrix of $\boldsymbol{\Theta}_{\mathrm{ML}} = \boldsymbol{\Theta}_{\mathrm{LS}}$ is approximately given by:

$$K(\boldsymbol{\Theta}_{\mathrm{ML}}) \approx [E(M(\boldsymbol{\theta}_p))]^{-1} = \sigma^2 (\xi^{\mathrm{T}} \xi)^{-1} \tag{A12}$$

Using the linear Taylor expansion (A1), the maximum likelihood approximation of the variance of the output in the gaussian case is obtained

$$\mathrm{var}(f(\boldsymbol{x}, \boldsymbol{\Theta}_{\mathrm{LS}})) \approx \sigma^2 \boldsymbol{\xi}^{\mathrm{T}} (\xi^{\mathrm{T}} \xi)^{-1} \boldsymbol{\xi} \tag{A13}$$

Hence, the maximum likelihood approximate variance (A13) is identical to the sampling theory approximate variance (A7).

**Remark.** *The Hessian matrix h is the value of the random matrix H with elements:*

$$[H(\boldsymbol{\theta})]_{ij} = \frac{\partial^2 J}{\partial \theta_i \partial \theta_j}$$

$$= \sum_{k=1}^{N} \left( \frac{\partial f(\boldsymbol{x}^k, \boldsymbol{\theta})}{\partial \theta_i} \frac{\partial f(\boldsymbol{x}^k, \boldsymbol{\theta})}{\partial \theta_j} + (Y_p^k - f(\boldsymbol{x}^k, \boldsymbol{\theta})) \frac{\partial^2 f(\boldsymbol{x}^k, \boldsymbol{\theta})}{\partial \theta_i \partial \theta_j} \right) \tag{A14}$$

Thus

$$E(M(\boldsymbol{\theta}_p)) = \frac{1}{\sigma^2} E(H(\boldsymbol{\theta}_p)) = \frac{1}{\sigma^2} \xi^{\mathrm{T}} \xi \tag{A15}$$

We can thus estimate the variance with

$$\mathrm{var}(\widehat{f(\boldsymbol{x}, \boldsymbol{\Theta}_{\mathrm{LS}})})_{\mathrm{LTE}} = \widehat{\sigma^2} z^{\mathrm{T}} (z^{\mathrm{T}} z)^{-1} z \tag{A16}$$

In likelihood theory, the variance of the noise is estimated with

$$\frac{R^{\mathrm{T}} R}{N} = \frac{N-q}{N} s^2 \approx s^2,$$

but we will skip over this minor difference; Eq. (A16) is thus identical to Eq. (A8).

It is also proposed to estimate the Fisher information matrix $E(M(\boldsymbol{\theta}_p))$ with the "observed information matrix" $m(\boldsymbol{\theta}_{\mathrm{LS}})$; this leads to estimate the variance with

$$\mathrm{var}(\widehat{f(\boldsymbol{x}, \boldsymbol{\Theta}_{\mathrm{LS}})})_{\mathrm{Hessian}} = \widehat{\sigma^2} z^{\mathrm{T}} [h(\boldsymbol{\theta}_{\mathrm{LS}})]^{-1} z \tag{A17}$$

In contrary to estimate (A16), estimate (A17) necessitates the computation of the Hessian.

### A.3. Sandwich variance estimate

Let us propose a derivation of this estimate in the sampling theory. A second expansion is needed, the LTE of the gradient at the true parameter value $\boldsymbol{\theta}_p$:

$$\frac{\partial J}{\partial \boldsymbol{\theta}}\bigg|_{\boldsymbol{\theta}=\boldsymbol{\theta}_{\mathrm{LS}}} \approx \frac{\partial J}{\partial \boldsymbol{\theta}}\bigg|_{\boldsymbol{\theta}=\boldsymbol{\theta}_p} + \frac{\partial^2 J}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^{\mathrm{T}}}\bigg|_{\boldsymbol{\theta}=\boldsymbol{\theta}_p} (\boldsymbol{\theta}_{\mathrm{LS}} - \boldsymbol{\theta}_p)$$

$$= \frac{\partial J}{\partial \boldsymbol{\theta}}\bigg|_{\boldsymbol{\theta}=\boldsymbol{\theta}_p} + h(\boldsymbol{\theta}_p)(\boldsymbol{\theta}_{\mathrm{LS}} - \boldsymbol{\theta}_p) \tag{A18}$$

where $h(\boldsymbol{\theta}_p)$ is the value of the random Hessian matrix (see Section A.2) evaluated at $\boldsymbol{\theta}_p$. Hence an approximate expression of the LS estimate of the parameters

$$\boldsymbol{\theta}_{\mathrm{LS}} \approx \boldsymbol{\theta}_p - [h(\boldsymbol{\theta}_p)]^{-1} \frac{\partial J}{\partial \boldsymbol{\theta}}\bigg|_{\boldsymbol{\theta}=\boldsymbol{\theta}_p} \tag{A19}$$

In Eq. (A19), we can replace the gradient by its expression

$$\frac{\partial J}{\partial \boldsymbol{\theta}}\bigg|_{\boldsymbol{\theta}=\boldsymbol{\theta}_p} = -\xi^{\mathrm{T}} (y_p - f(x, \boldsymbol{\theta}_p)) = -\xi^{\mathrm{T}} W \tag{A20}$$

Hence the corresponding approximation of the least-squares estimator (random vector $\boldsymbol{\Theta}_{\mathrm{LS}}$) of the parameters

$$\boldsymbol{\Theta}_{\mathrm{LS}} \approx \boldsymbol{\theta}_p + [H(\boldsymbol{\theta}_p)]^{-1} \xi^{\mathrm{T}} W \tag{A21}$$

Using the LTE of the model output (A1), we obtain:

$$f(\boldsymbol{x}, \boldsymbol{\Theta}_{\mathrm{LS}}) \approx f(\boldsymbol{x}, \boldsymbol{\theta}_p) + \boldsymbol{\xi}^{\mathrm{T}} [H(\boldsymbol{\theta}_p)]^{-1} \xi^{\mathrm{T}} W \tag{A22}$$

Neglecting the random character of $H$ ($H$ being replaced by $h$), the output variance can be approximated by

$$\mathrm{var}(f(\boldsymbol{x}, \boldsymbol{\Theta}_{\mathrm{LS}})) \approx \sigma^2 \boldsymbol{\xi}^{\mathrm{T}} [h(\boldsymbol{\theta}_p)]^{-1} \xi^{\mathrm{T}} \xi [h(\boldsymbol{\theta}_p)]^{-1} \boldsymbol{\xi} \tag{A23}$$

This leads to propose the sandwich estimate

$$\mathrm{var}(\widehat{f(\boldsymbol{x}, \boldsymbol{\Theta}_{\mathrm{LS}})})_{\mathrm{sandwich}} = s^2 z^{\mathrm{T}} [h(\boldsymbol{\theta}_{\mathrm{LS}})]^{-1} z^{\mathrm{T}} z [h(\boldsymbol{\theta}_{\mathrm{LS}})]^{-1} z \tag{A24}$$

This estimate also necessitates the computation of the Hessian of the cost-function.

## Appendix B. Derivation of an approximate LOO error

The following derivation is inspired from the work of Antoniadis et al. (1992) and is valid irrespective of whether or not the assumed model is true. We denote by $\boldsymbol{\theta}_{\text{LS}}^{(k)}$ the LS estimate on the $k$th LOO set $\{\boldsymbol{x}^i, y_p^i\}_{i=1\text{to }N, i \neq k}$. We have the $k$th residual $r^k$ and the $k$th LOO error $e^k$:

$$\begin{cases} r^k = y_p^k - f(\boldsymbol{x}^k, \boldsymbol{\theta}_{\text{LS}}) \\ e^k = y_p^k - f(\boldsymbol{x}^k, \boldsymbol{\theta}_{\text{LS}}^{(k)}) \end{cases} \tag{B1}$$

Let us denote by $\boldsymbol{y}_p^{(k)}$ the $(N-1)$-vector obtained by deletion of the $k$th component of the measured output vector $\boldsymbol{y}_p$, by $z^{(k)}$ the $(N-1, q)$ matrix obtained by deletion of the $k$th row of $z$, by $x^{(k)}$ the $(N-1, q)$ matrix obtained by deletion of the $k$th row of $x$. The LOO estimate $\boldsymbol{\theta}_{\text{LS}}^{(k)}$ minimizes the cost-function

$$J^{(k)}(\boldsymbol{\theta}) = \tfrac{1}{2}(\boldsymbol{y}_p^{(k)} - \boldsymbol{f}(x^{(k)}, \boldsymbol{\theta}))^{\text{T}}(\boldsymbol{y}_p^{(k)} - \boldsymbol{f}(x^{(k)}, \boldsymbol{\theta})) \tag{B2}$$

We first approximate $\boldsymbol{f}(x^{(k)}, \boldsymbol{\theta})$ by its LTE at $\boldsymbol{\theta}_{\text{LS}}$:

$$\boldsymbol{f}(x^{(k)}, \boldsymbol{\theta}) \approx \boldsymbol{f}(x^{(k)}, \boldsymbol{\theta}_{\text{LS}}) + z^{(k)}(\boldsymbol{\theta} - \boldsymbol{\theta}_{\text{LS}}) \tag{B3}$$

Hence the approximation of $\boldsymbol{\theta}_{\text{LS}}^{(k)}$

$$\boldsymbol{\theta}_{\text{LS}}^{(k)} \approx \boldsymbol{\theta}_{\text{LS}} + \left(z^{(k)\text{T}} z^{(k)}\right)^{-1} z^{(k)\text{T}}(\boldsymbol{y}_p^{(k)} - \boldsymbol{f}(x^{(k)}, \boldsymbol{\theta}_{\text{LS}})) \tag{B4}$$

In the previous expression, we have

$$z^{(k)\text{T}}(\boldsymbol{y}_p^{(k)} - \boldsymbol{f}(x^{(k)}, \boldsymbol{\theta}_{\text{LS}})) = z^{\text{T}}(\boldsymbol{y}_p - \boldsymbol{f}(x, \boldsymbol{\theta}_{\text{LS}})) - \boldsymbol{z}^k r^k$$

$$= z^{\text{T}}\boldsymbol{r} - \boldsymbol{z}^k r^k = -\boldsymbol{z}^k r^k \tag{B5}$$

because the columns of $z$ are orthogonal to the residual vector $\boldsymbol{r}$. Using the matrix inversion lemma, we can express $(z^{(k)\text{T}} z^{(k)})^{-1}$ in Eq. (B4) in terms of $(z^{\text{T}} z)^{-1}$

$$\left(z^{(k)\text{T}} z^{(k)}\right)^{-1} = (z^{\text{T}} z)^{-1} + \frac{(z^{\text{T}} z)^{-1} \boldsymbol{z}^{(k)} \boldsymbol{z}^{(k)\text{T}} (z^{\text{T}} z)^{-1}}{1 - \boldsymbol{z}^{(k)\text{T}} (z^{\text{T}} z)^{-1} \boldsymbol{z}^{(k)}}$$

$$= (z^{\text{T}} z)^{-1} + \frac{(z^{\text{T}} z)^{-1} \boldsymbol{z}^{(k)} \boldsymbol{z}^{(k)\text{T}} (z^{\text{T}} z)^{-1}}{1 - [p_z]_{kk}} \tag{B6}$$

where $p_z$ denotes the orthogonal projection matrix on the range of $z$.

Replacing Eqs. (B5) and (B6) into Eq. (B4), we finally obtain

$$\boldsymbol{\theta}_{\text{LS}}^{(k)} \approx \boldsymbol{\theta}_{\text{LS}} - (z^{\text{T}} z)^{-1} \boldsymbol{z}^k \frac{r^k}{1 - [p_z]_{kk}} \tag{B7}$$

Expanding $e^k$ at $\boldsymbol{\theta}_{\text{LS}}$ and replacing Eq. (B7) into this expansion, we obtain an approximate expression of the LOO error

which is similar to the expression of the linear LOO error (36)

$$e^k \approx \frac{r^k}{1 - [p_z]_{kk}} \qquad k = 1 \text{ to } N \tag{B8}$$

In practice, the diagonal terms of $p_z$ are computed using the singular value factorization of $z = u\Sigma v^{\text{T}}$, where $u$ is an orthogonal $(N,N)$ matrix, $\Sigma$ is a diagonal $(N,q)$ matrix, and $v$ is an orthogonal $(q,q)$ matrix (see, Golub & Van Loan, 1983). Then:

$$[p_z]_{kk} = \sum_{i=1}^{q} [u]_{ki}^2 \qquad k = 1 \text{ to } N \tag{B9}$$

The diagonal elements of $p_z$ that differ from 1 by a threshold consistent with the computer precision are considered as equal to 1 (theoretically, the values of the $[p_z]_{kk}$ are comprised between $1/N$ and 1).

## References

Antoniadis, A., Berruyer, J., & Carmona, R. (1992). *Régression non linéaire et applications*. Paris: Economica.

Bates, D. M., & Watts, D. G. (1988). *Nonlinear regression analysis and its applications*. New York: Wiley.

Bishop, M. (1995). *Neural networks for pattern recognition*. Oxford: Clarendon Press.

Buntine, W., & Weigend, A. (1994). Computing second derivatives in feedforward neural networks: a review. *IEEE Transactions on Neural Networks*, 5 (3), 480–488.

Draper, N. R., & Smith, H. (1998). *Applied regression analysis*. New York: Wiley.

Efron, B., & Tibshirani, R. J. (1993). *An introduction to the bootstrap*. New York: Chapman.

Golub, G. H., & Van Loan, C. F. (1983). *Matrix computations*. Baltimore: John Hopkins University Press.

Goodwin, G. C., & Payne, R. L. (1977). *Dynamic system identification; experiment design and data analysis*. New York: Academic Press.

Hansen, L. K., & Larsen, J. (1993). Linear unlearning for cross validation. *Advances in Computational Mathematics*, 5, 286–290.

Heskes, T. (1997). Practical confidence and prediction intervals. In M. Mozer, M. Jordan & T. Petsche, *Advances in neural information processing systems*, vol. 9. Cambridge, MA: MIT Press, 176–182.

MacKay, D. J. C. (1992a). Bayesian interpolation. *Neural Computation*, 4, 415–447.

MacKay, D. J. C. (1992b). A practical Bayesian framework for backprop networks. *Neural Computation*, 4, 448–472.

Monari, G. (1999). *Sélection de modèles non linéaires par leave-one-out; étude théorique et application des réseaux de neurones au procédé de soudage par points*. Thèse de Doctorat de l'Université Paris 6.

Monari, G., & Dreyfus, G. Local linear least squares: performing leave-one-out without leaving anything out. Submitted for publication.

Paass, G. (1993). Assessing and improving neural network predictions by the bootstrap algorithm. In S. J. Hanson, J. D. Cowan & C. L. Giles, *Advances in neural information processing systems* (pp. 186–203), 5. Cambridge, MA: MIT Press.

Refenes, A.-P. N., Zapranis, A. D., & Utans, J. (1997). Neural model identification, variable selection and model adequacy. In A. S. Weigend, Y. Abu-Mostafa & A.-P. N. Refenes, *Decision technologies for financial engineering* (pp. 243–261). Singapore: World Scientific.

Ripley, B. D. (1995). *Pattern recognition and neural networks*. Cambridge: Cambridge University Press.

Rivals, I., & Personnaz, L. (1998). Construction of confidence intervals in neural modeling using a linear Taylor expansion. *Proceedings of the International Workshop on Advanced Black-Box Techniques for Nonlinear Modeling* (pp. 17–22), 8–10 July, Leuwen.

Rivals, I., & Personnaz, L. (1999). On cross-validation for model selection. *Neural Computation*, *11*, 863–870.

Saarinen, S., Bramley, R., & Cybenko, G. (1993). Ill-conditioning in neural network training problems. *SIAM Journal on Scientific and Statistical Computing*, *14*, 693–714.

Seber, G. A. F. (1977). *Linear regression analysis*. New York: Wiley.

Seber, G. A. F., & Wild, C. (1989). *Nonlinear regression*. New York: Wiley.

Sussman, H. J. (1992). Uniqueness of the weights for minimal feedforward nets with a given input–output map. *Neural Networks*, *5*, 589–593.

Tibshirani, R. J. (1996). A comparison of some error estimates for neural models. *Neural Computation*, *8*, 152–163.

Urbani, D., Roussel-Ragot, P., Personnaz, L., & Dreyfus, G. (1994). The selection of neural models of non-linear dynamical systems by statistical tests. *Neural Networks for Signal Processing, Proceedings of the 1994 IEEE Workshop*.

Zhou, G., & Si, J. (1998). A systematic and effective supervised learning mechanism based on Jacobian rank deficiency. *Neural Computation*, *10*, 1031–1045.