



HAL
open science

Accommodation to outliers in identification of non linear SISO systems with neural networks

Gérard Bloch, Philippe Thomas, Didier Theilliol

► **To cite this version:**

Gérard Bloch, Philippe Thomas, Didier Theilliol. Accommodation to outliers in identification of non linear SISO systems with neural networks. *Neurocomputing*, 1997, 14 (1), pp.85-99. 10.1016/0925-2312(95)00134-4 . hal-00951969

HAL Id: hal-00951969

<https://hal.science/hal-00951969>

Submitted on 25 Feb 2014

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Accommodation to outliers in identification of non linear SISO systems with neural networks

Gérard Bloch^a, Philippe Thomas^a, Didier Theilliol^a

^aCentre de Recherche en Automatique de Nancy, CNRS URA D821
ESSTIN, Rue Jean Lamour, 54500 Vandoeuvre, France

Abstract

The problem of non-linear Single Input Single Output system identification in the presence of large errors in data is considered. Combining the capabilities of neural networks to solve non-linear problems by learning and a robust recursive prediction error learning rule based on the modeling of the errors, a new algorithm is drawn up. Its potential is illustrated through simulation studies.

Key words: Identification, Neural networks, Robustness, Non-Linear system, Outliers.

1. Introduction

Artificial neural networks have been the focus of a great deal of attention during the last decade, due to their capabilities to solve non-linear problems by learning. Such networks provide a parallel structure with very simple processing elements. Although a broad range of neural networks (NN) architectures and learning rules are available (Grossberg, 1988; Kohonen, 1989; Lippman, 1987; Widrow and Lehr, 1990), the backpropagation algorithm for multilayer feedforward networks (Rumelhart et al., 1986) is the most popular approach for engineering applications. Backpropagation or derived algorithms have been successfully applied for classification and pattern recognition (Rajavelu et al., 1989; Xue et al., 1992), fault detection (Hoskins and Himmelblau, 1988; Kramer and Leonard, 1988), non-linear control (Nahas et al., 1992) and process modelling and identification (Bhat and McAvoy, 1990; Billings et al., 1991; Narendra and Parthasarathy, 1990).

On the other hand, an extensive literature on system identification can be found. Among general textbooks on the subject, Box and Jenkins (1970), Söderström and Stoica (1989) and Ljung (1987) can be mentioned. Particularly, much effort has been devoted to tackle the presence of outliers in experimental input and output data used for identification. Large errors or outliers in data can be for instance caused by offset of sensors, failure of transducers, analog to digital conversion errors or even by malfunctioning of transmission devices. The related works are mainly based on modeling such outliers to produce so called robust identification algorithms. But these works are limited to linear systems.

In this paper, the main feature of neural networks, the ability to identify non-linear systems and a robust recursive prediction error algorithm, based on the modeling of errors due to Huber (1964), are combined. This modeling has been used by Puthenpura and Sinha (1990) for a robust linear recursive least squares type identification algorithm. The convergence of this kind of

algorithms particularly for non-linear systems is very slow. So, a robust feature is introduced in a recursive Gauss-Newton type of algorithm, first employed by Chen et al. (1990a) for neural networks. The goal is to accommodate to outliers in order to eliminate their effects in the identification of non linear SISO systems by an appropriate choice of the criterion to be minimized.

2. Multilayer feedforward neural networks for identification

In this part, the structure of the multilayer feedforward neural network, used for identification of dynamical single input single output (SISO) systems, is presented. The network, shown in Figure 1, is composed of interconnected processing units in three successive layers.

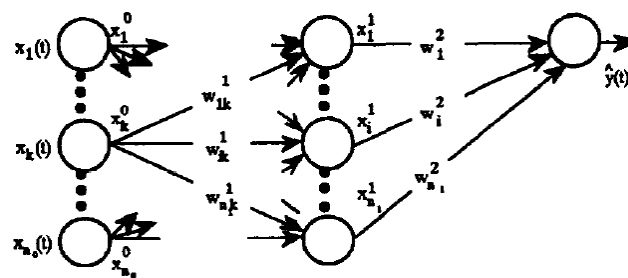


Figure 1: A three layers feedforward neural architecture.

The first or input layer is composed of "transparent" units which do not perform any computation but simply distribute their inputs to all neurons in the next layer called hidden layer ($x_k^0 = x_k(t), \forall k$). For identification purposes, several authors such as Cybenko (1989) or Funahashi (1989) have established that multilayer feedforward neural networks with a single hidden layer are able to approximate continuous functions. The last layer is the output layer, composed of a single neuron and

its output gives the estimated output of the SISO system. Neurons in the hidden and output layers are identical and can be represented as illustrated in Figure 2.

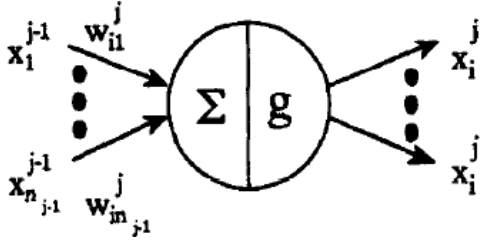


Figure 2: Neuron i in layer j .

The i th neuron in the layer j receives n_{j-1} inputs $\{x_1^{j-1}, \dots, x_{n_{j-1}}^{j-1}\}$ from layer $j-1$ with associated weights $\{w_{i1}^j, \dots, w_{in_{j-1}}^j\}$. This neuron first computes the weighted sum of the n_{j-1} inputs:

$$z_i^j = \sum_{k=1}^{n_{j-1}} w_{ik}^j x_k^{j-1} + b_i^j, \quad (1)$$

where b_i^j is a bias or threshold term. The output of the neuron is a non-linear function of the sum in (1):

$$x_i^j = g(z_i^j), \quad (2)$$

where g denotes the activation function, chosen as often to be a sigmoidal function, here:

$$g(x) = \frac{1}{1 + e^{-x}}, \quad (3)$$

with $\lim_{x \rightarrow -\infty} g(x) = 0$ and $\lim_{x \rightarrow +\infty} g(x) = 1$. So, for the structure and the notations of Figure 1, a network with a single hidden layer can be defined by the following model:

$$\hat{y}(t) = g(z) \quad (4a)$$

$$z = \sum_{k=1}^{n_1} w_k^2 x_k^1 + b^2 \quad (4b)$$

$$x_i^1 = g(z_i^1) \quad (4c)$$

$$z_i^1 = \sum_{k=1}^{n_0} w_{ik}^1 x_k(t) + b_i^1 \quad (4d)$$

$x(t) = [x_1(t) \dots x_{n_0}(t)]^T$ is the n_0 -input vector. In the following, n_d , n_u and n_y refer respectively the input time delay, the numbers of lagged system inputs and outputs to be applied to input layer of the network. In the training phase of the neural network, i.e. the identification step, input $u(t-n_d)$, $u(t-n_d-1)$... and output $y(t-1)$, $y(t-2)$... values of the process are successively applied on the input layer of the network in order to produce an estimated value of the system output:

$$\hat{y}(t, \Theta) = NN(y(t-1), \dots, y(t-n_y), u(t-n_d), \dots, u(t-n_d-n_u+1)) \quad (5)$$

where $\Theta = [\theta_1 \dots \theta_{n_\theta}]^T$ comprises all the unknown weights and biases of the network. The dimension n_θ of the parameter vector Θ is defined as:

$$n_\theta = (n_0 + 1) n_1 + (n_1 + 1),$$

where $n_0 = n_y + n_u$ is the number of the input layer neurons and n_1 is the number of the hidden layer neurons. So the predictor can be noted $NN(n_y, n_u, n_d, n_1)$.

To avoid the saturation of the activation function (3), particularly for the output neuron, contained between 0 and 1, observed input and output data of the system must be scaled between 0 and 1. However, the same notation for original and normalized data is used in the following.

3. Recursive prediction error method

The general framework of the learning rule used in the following is now presented. The backpropagation algorithm (Rumelhart et al., 1986) is the first training method to estimate parameters of multilayer neural network and is a gradient algorithm designed to minimize the mean square error between the output of the network and the desired output.

The recursive prediction error or RPE algorithm, first introduced by Chen et al. (1990a,b) for training neural networks, is a general recursive parameter estimation method which minimizes the prediction error using an approximation of the Gauss-Newton search direction. Only the version of the RPE algorithm introduced in (Chen et al., 1990a) has been considered here. Billings et al. (1991) have shown that the RPE algorithm provides an effective method of learning neural networks. Backpropagation can be viewed as a simplified version of the RPE algorithm. Compared with backpropagation, RPE algorithm involves increased computational load at each iteration, but presents faster convergence, yielding to shorter global computational time. Furthermore, RPE removes the dependence of the estimation algorithm on the user selectable parameters such as learning rate and momentum. Indeed, with an inappropriate combination of these parameters, backpropagation performs badly. In any case, neither backpropagation nor RPE algorithm ensure to reach a final estimation corresponding to a global minimum of the criterion.

The RPE algorithm starts from the general criterion:

$$J(t, \Theta) = \gamma(t) \sum_{k=1}^t \beta(t, k) \ell(\varepsilon(k, \Theta), k), \quad (6)$$

where $\gamma(t)$ is the adaptation gain a time t with $\sum_{k=1}^t \gamma(t) \beta(t, k) = 1$, $\varepsilon(k, \Theta)$ is the scalar prediction error and $\ell(\varepsilon(k, \Theta), k)$ can be chosen as a quadratic function weighted by the innovation variance $\Lambda(k)$:

$$\ell(\varepsilon(k, \Theta), k) = \frac{1}{2} \Lambda^{-1}(k) \varepsilon^2(k, \Theta). \quad (7)$$

The minimization of the criterion (6) can be performed according to:

$$\hat{\Theta}(t) = \hat{\Theta}(t-1) + f_t(\hat{\Theta}(t-1)), \quad (8)$$

where $\hat{\Theta}(t)$ is the estimate of Θ at time t and $f_i(\hat{\Theta}(t-1))$ is a search direction based on information about $J(t, \Theta)$.

The parameter vector Θ is estimated for each $t = 1, \dots, N$, where N is the number of available observations. For off-line estimation, the data set is presented several times and each presentation is called an iteration. The Gauss-Newton search direction is used here and is defined by:

$$f_i(\hat{\Theta}) = -[R(t)]^{-1} \nabla J(t, \Theta), \quad (9)$$

where $R(t)$ and $\nabla J(t, \Theta)$ are respectively the $n_\theta \cdot n_\theta$ approximate Hessian matrix and the $n_\theta \cdot 1$ gradient of $J(t, \Theta)$. The derivation is given by Ljung (1987) and yields the general recursive prediction error algorithm:

$$\varepsilon(t) = y(t) - \hat{y}(t/\hat{\Theta}(t-1)), \quad (10a)$$

$$R(t) = R(t-1) + \gamma(t) [\psi(t) \Lambda^{-1}(t) \psi^T(t) - R(t-1)], \quad (10b)$$

$$\hat{\Theta}(t) = \hat{\Theta}(t-1) + \gamma(t) R^{-1}(t) \psi(t) \Lambda^{-1}(t) \varepsilon(t), \quad (10c)$$

where $\psi(t) = \left[\frac{\partial \hat{y}(t/\Theta)}{\partial \Theta} \right]$ is the $n_\theta \cdot 1$ gradient of \hat{y} with respect to Θ . The elements of $\psi(t)$ must be written depending on the location of the parameters in the network, in the spirit of backpropagation. It can be first shown from (3) that:

$$\frac{\partial g(x)}{\partial x} = g(x) (1 - g(x)). \quad (11)$$

For the parameters of the output layer, Eqs. (11), (4a) and (4b) then yield:

$$\frac{\partial \hat{y}}{\partial w_k^2} = \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w_k^2} = \hat{y} (1 - \hat{y}) x_k^1, \quad (12a)$$

$$\frac{\partial \hat{y}}{\partial b^2} = \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial b^2} = \hat{y} (1 - \hat{y}). \quad (12b)$$

For the parameters of the hidden layer, Eqs. (11) and (4) yield:

$$\frac{\partial \hat{y}}{\partial w_{ik}^1} = \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial x_i^1} \frac{\partial x_i^1}{\partial z_i^1} \frac{\partial z_i^1}{\partial w_{ik}^1} = \hat{y} (1 - \hat{y}) w_i^2 x_i^1 (1 - x_i^1) x_k, \quad (13a)$$

$$\frac{\partial \hat{y}}{\partial b_i^1} = \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial x_i^1} \frac{\partial x_i^1}{\partial z_i^1} \frac{\partial z_i^1}{\partial b_i^1} = \hat{y} (1 - \hat{y}) w_i^2 x_i^1 (1 - x_i^1). \quad (13b)$$

So, the differentiation of \hat{y} with respect to θ_j , $j = 1, \dots, n_\theta$, can be summarized as follows:

$$\psi_j = \begin{cases} \hat{y} (1 - \hat{y}) x_k^1, & \text{if } \theta_j = w_k^2, 1 \leq k \leq n_1, \\ \hat{y} (1 - \hat{y}), & \text{if } \theta_j = b^2, \\ \hat{y} (1 - \hat{y}) w_i^2 x_i^1 (1 - x_i^1) x_k, & \text{if } \theta_j = w_{ik}^1, 1 \leq k \leq n_\theta, \\ \hat{y} (1 - \hat{y}) w_i^2 x_i^1 (1 - x_i^1). & \text{if } \theta_j = b^2. \end{cases} \quad (14)$$

As developed in the appendix, the practical implementation of the algorithm (10) avoids to invert $R(t)$ at each step:

$$\hat{\Theta}(t) = \hat{\Theta}(t-1) + L(t) \varepsilon(t), \quad (15a)$$

$$L(t) = \frac{P(t-1) \psi(t)}{\lambda(t) \Lambda(t) + \psi^T(t) P(t-1) \psi(t)}, \quad (15b)$$

$$P(t) = \frac{1}{\lambda(t)} [P(t-1) - L(t) \psi^T P(t-1)], \quad (15c)$$

where

$$P(t) = \gamma(t) R^{-1}(t), \quad (16)$$

and

$$\lambda(t) = \frac{\gamma(t-1)}{\gamma(t)} (1 - \gamma(t)). \quad (17)$$

The so-called forgetting factor is calculated practically by:

$$\lambda(t) = \lambda_0 (1 - \lambda(t-1)) + (1 - \lambda_0).$$

4. Robustification of the algorithm

Large errors or outliers are quite difficult to be detected and picked out before identification and can cause the parameters to be highly biased. In order to tackle this problem, Puthenpura and Sinha (1990) have developed a robust recursive identification method for linear dynamical systems. This scheme is a weighted least squares algorithm with particular weights and is very similar to the robust Kalman filter obtained by Masreliez and Martin (1977). Based on the modeling of outliers due to Huber (1964), it considers that the measurement noise $e(t)$ which contaminates the noise free output is from the family D_μ , defined by:

$$D_\mu = \{D \mid D = (1 - \mu)G + \mu H, 0 \leq \mu \leq 1\}, \quad (18)$$

where G is the usual normal distribution, H an arbitrary symmetric long-tailed distribution and μ the probability of occurring large errors. In fact H is assumed to be also normal, but with a larger variance compared to G :

$$e(t) \sim (1 - \mu) \mathcal{N}(0, \sigma_1^2) + \mu \mathcal{N}(0, \sigma_2^2), \quad (19)$$

where \mathcal{N} represents a normal distribution, with $\sigma_2^2 > \sigma_1^2$. The probability μ of occurring large errors being unknown, the preceding model is replaced by:

$$e(t) \sim (1 - \delta(t)) \mathcal{N}(0, \sigma_1^2) + \delta(t) \mathcal{N}(0, \sigma_2^2), \quad (20)$$

where $\delta(t) = 0$ for $|\varepsilon(t)| \leq M$ and $\delta(t) = 1$ for $|\varepsilon(t)| > M$, with $\varepsilon(t)$ the prediction error and M a preassigned bound which can be taken as $3\sigma_1$ (Aström, 1980). So the weighting factor appearing in (15) will be chosen as:

$$\Lambda(t) = (1 - \delta(t)) \sigma_1^2 + \delta(t) \sigma_2^2 \quad (21)$$

to reduce the influence of large innovations. Moreover, the variances σ_1^2 and σ_2^2 can be updated as:

$$\begin{aligned} \sigma_1^2(t) &= \sigma_1^2(t-1) + \frac{1}{t-\tau} (\varepsilon^2(t) - \sigma_1^2(t-1)), \text{ for } |\varepsilon(t)| \leq 3\sigma_1(t-1), \\ \sigma_1^2(t) &= \sigma_1^2(t-1), \text{ otherwise,} \end{aligned} \quad (22a)$$

and

$$\begin{aligned} \sigma_2^2(t) &= \sigma_2^2(t-1) + \frac{1}{\tau} (\varepsilon^2(t) - \sigma_2^2(t-1)), \text{ for } |\varepsilon(t)| > 3\sigma_1(t-1), \\ \sigma_2^2(t) &= \sigma_2^2(t-1), \text{ otherwise,} \end{aligned} \quad (22b)$$

with $\tau = 0$, for $t = 1$ and $\tau = \tau + 1$ whenever $|\varepsilon(t)| > 3\sigma_1(t-1)$. As pointed out by Puthenpura and Sinha (1990), τ is the estimated number of outliers.

$\sigma_2^2(0)$ can be chosen as $\sigma_2^2(0) = 3\sigma_1^2(0)$. $\sigma_1^2(0)$ should be chosen much greater than the real value of the noise variance so that in the beginning of the identification no residual $\varepsilon(t)$ appears like outlier. With this choice, σ_1^2 converges to the true value of the noise variance. When σ_1^2 is close to the noise variance, outliers are detected and their influence on identification becomes insignificant. If $\sigma_1^2(0)$ is chosen very small (or zero), all residuals have an absolute value $|\varepsilon(t)|$ greater than $3\sigma_1(t-1)$, only σ_2^2 converges to the noise variance (but with the influence of outliers) and the parameter estimation is biased when outliers are present. If $\sigma_1^2(0)$ is chosen very large, the accommodation to outliers is just delayed.

In the next part, the three following algorithms are applied from (15) to a simulated example:

- NN (Neural Network) with $\Lambda(t) = 1$ and $\lambda(t) = 1$,
- FNN (Neural Network with Forgetting factor) where $\Lambda(t) = 1$ and $\lambda(t)$ is calculated by the recursive relation $\lambda(t) = \lambda_0(1 - \lambda(t-1)) + (1 - \lambda_0)$, with $\lambda_0 = 0.99$ and $\lambda(0) = 0.95$,
- RNN (Robust Neural Network) where $\lambda(t) = 1$ and $\Lambda(t)$ is computed by (21) and (22).

5. Simulation results

A non-linear Hammerstein system example, introduced by Billings et al. (1991), is considered:

$$\begin{aligned} y(t) &= 0.8y(t-1) + 0.4NL(u(t-1)) + e(t) \\ NL(t-1) &= u(t-1) + u(t-1)^2 + u(t-1), \end{aligned} \quad (23)$$

where $u(t)$ is the system input at time t , chosen as a sequence uniformly distributed between -4 and 4, in order to study the system on the whole non-linearity, $y(t)$ is the system output, $e(t)$ is a gaussian noise such as $e \sim \mathcal{N}(0, \sigma_1^2)$ when no outliers are present. In order to show the influence of the outliers on the output, Figure 3 presents the 500 output values, where the noise variance σ_1^2 is equal to 2.56 ($\sigma_1 = 1.61$), contaminated by 25 randomly located outliers. These large errors are simply simulated by multiplying the original values of the noise by a factor f equal to 20.

Figure 4 shows the difference between the preceding series and the corresponding outliers free and noise free simulation. The impact of the outliers filtered by the process dynamics can be noticed. Figure 5 presents the difference between the noisy and outliers free series and the noise free series, which represents the noise filtered by the process, in order to show its variation range, significantly smaller than for the preceding figure.

For the different following experiments, the neural predictor has the same structure $NN(n_y = 2, n_u = 2, n_d = 1, n_1 = 3)$ and the initial weights, randomly chosen between 0 and 1, are kept. The initial value of the covariance matrices P is chosen equal to $100I$ and $\sigma_1^2(0)$ equal to 5 times the variance of the noise.

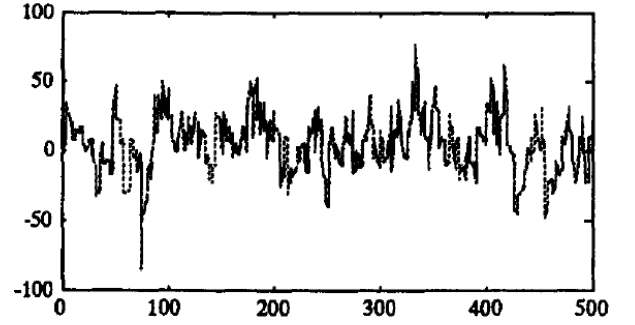


Figure 3: Output signal contaminated by outliers.

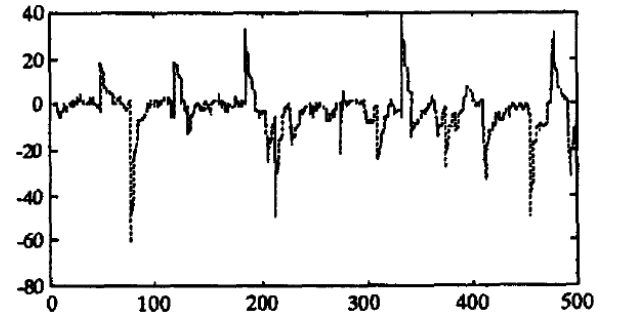


Figure 4: Impact of outliers on the output data.

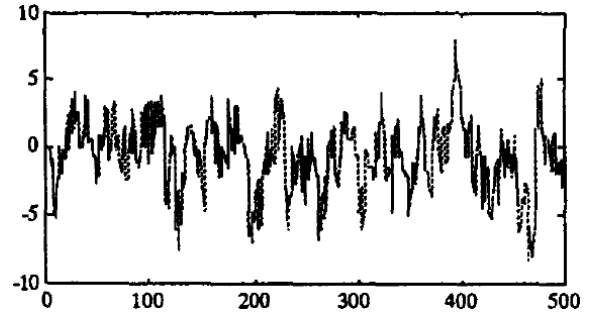


Figure 5: Noise filtered by the process.

For the off-line identification considered in this part, the data set is presented 20 times (iterations). Moreover, a second data set is used for the (cross)-validation of the neural models. This fresh data set is called validation set and has characteristics similar to the original identification set, concerning the input shape, the distribution of the noise and its variance, but is outlier free. The following examples give the values of the residual criterion:

$$V = \frac{1}{N} \sum_{t=1}^N (\hat{y}(t) - y(t))^2,$$

where $y(t)$ is the system output at time t , $\hat{y}(t)$ is the output estimated by the neural model obtained after 20 identification iterations and where $N = 500$.

Figure 6 presents, for $\sigma_1^2 = 0.64$, the variation of the residual criterion when the number of outliers (with multiplicative factor f equal to 25) is varied from 0 (no outliers) to 50. Figure 6(a) concerns the data set used for identification, Figure 6(b) the val-

identification set. NN and RNN algorithms yield very stable results for the validation set contrary to those for FNN. But the residual criterion is reduced by 5 from NN to RNN.

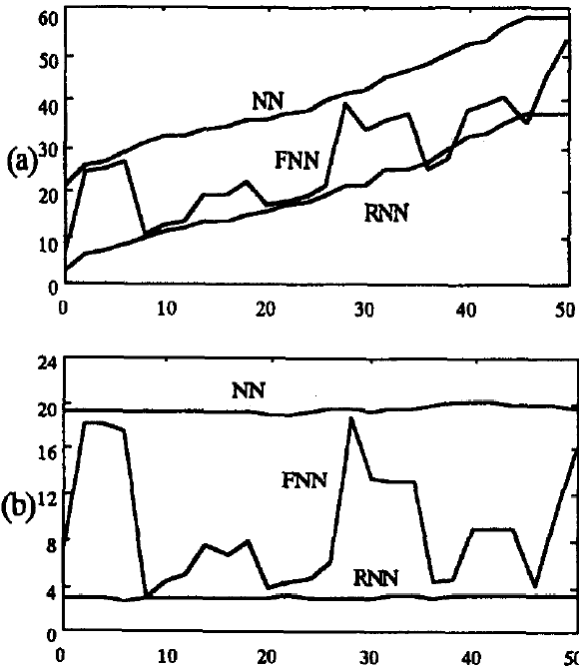


Figure 6: Residual criterion w.r.t. the outliers number. (a) identification set, (b) validation set.

Figures 7(a) and 7(b) show the evolution during learning of the residual criterion for algorithms NN and RNN, respectively, for 0, 25 and 50 outliers. The convergence of the RNN algorithm appears clearly faster. However, as shown in Figure 8, the variation of the criterion computed from the validation set is slower.

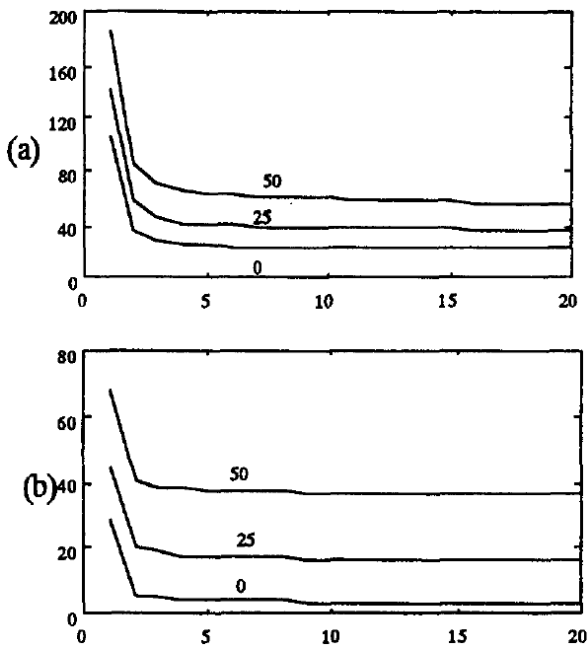


Figure 7: Residual criterion with respect to iterations. (a) NN. (b) RNN.

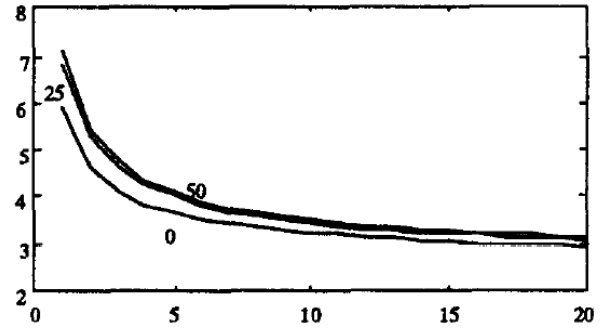


Figure 8: Residual criterion for validation set (RNN).

The second experiment deals with a variation of the multiplicative factor f from 1 (no outliers) to 50, for 50 outliers and for $\sigma_1^2 = 0.64$. The results are similar as those of the preceding experiment, as shown in Figure 9.

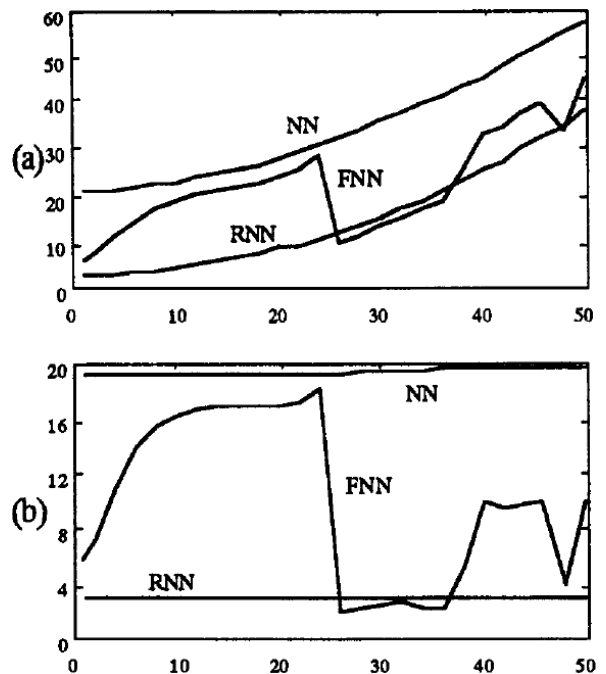


Figure 9: Residual criterion with respect to the factor f . (a) identification set, (b) validation set.

In the third experiment, for 50 outliers with multiplicative factor equal to 25, the noise variance σ_1^2 is varied from 0 (no noise) to 16. Results are given in Figure 10. As for the preceding experiment, the behavior of the FNN algorithm appears disconcerting, such an algorithm being misadapted for contaminated data. Note also on Figure 10(b) that the RNN algorithm, better than the simple NN, produces a good estimation of the noise variance from 2.5.

In the last experiment, σ_1^2 is fixed to 0.64 and a bias varying from 0 (no bias) to 20 is added to the last 50 values of the original outliers free noise. The results given in Figure 11 are very significant. While increasing with RNN for the identification set as the bias increases, the residual criterion is remarkably stable for the validation set, contrary to the other algorithms.

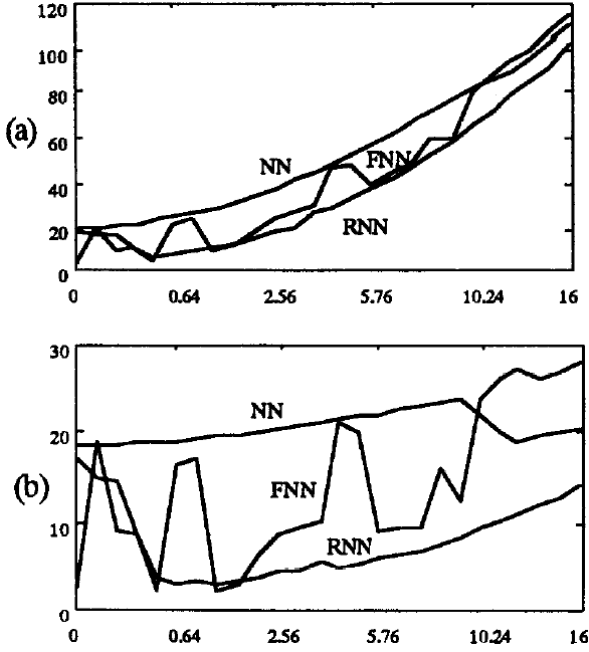


Figure 10: Residual criterion with respect to σ_1^2 . (a) identification set, (b) validation set.

These last results confirm the interest of the presented robust non linear predictor for detection of process changes.

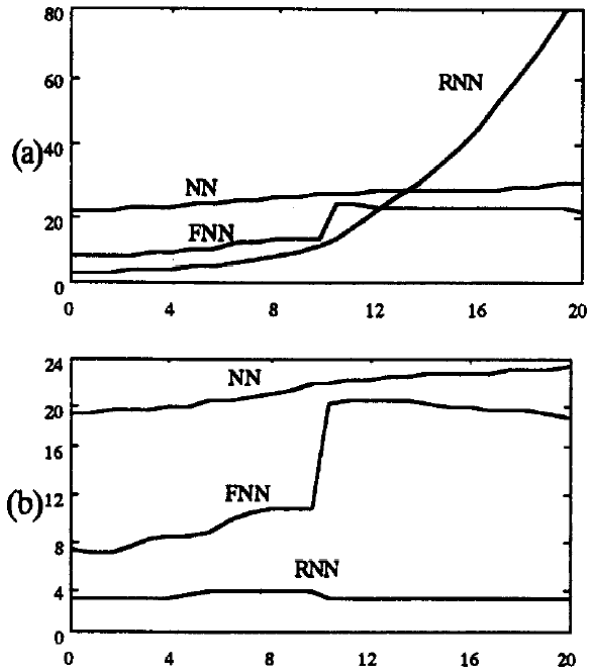


Figure 11: Residual criterion w.r.t. bias magnitude. (a) identification set, (b) validation set.

6. Conclusion

The neural nets have a structure which allows, with the adaptation of the backpropagation, the use of the various and some-

times classical parameter estimation algorithms. The problem of non-linear Single Input Single Output system identification in the presence of outliers in data has been considered. Combining the capabilities of neural networks to solve non-linear problems by learning and a robust recursive prediction error learning rule based on the modeling of the errors, a new algorithm has been drawn up. The results obtained suggest that this algorithm can be employed for identification from contaminated data, but also for failure detection and for robust control. The proposed method can be easily extended to MIMO systems.

Acknowledgements

The authors gratefully acknowledge the reviewers for their instructive comments.

A. Appendix

Let us consider the general prediction error algorithm introduced in (10):

$$\varepsilon(t) = y(t) - \hat{y}(t/\hat{\Theta}(t-1)), \quad (\text{A.1a})$$

$$R(t) = [1 - \gamma(t)] R(t-1) + \gamma(t) \psi(t) \Lambda^{-1}(t) \psi^T(t), \quad (\text{A.1b})$$

$$\hat{\Theta}(t) = \hat{\Theta}(t-1) + \gamma(t) R^{-1}(t) \psi(t) \Lambda^{-1}(t) \varepsilon(t). \quad (\text{A.1c})$$

To avoid inverting $R(t)$ at each step, it is convenient to introduce:

$$P(t) = \gamma(t) R^{-1}(t), \quad (\text{A.2})$$

and apply to (A.1b) the matrix inversion lemma:

$$(A + BCD)^{-1} = A^{-1} - A^{-1}B(C^{-1} + DA^{-1}B)^{-1}DA^{-1}. \quad (\text{A.3})$$

Taking $A = [1 - \gamma(t)] R(t-1)$, $B = D^T = \psi(t)$, and $C = \gamma(t) \Lambda^{-1}(t)$ gives:

$$P(t) = \frac{\gamma(t)}{1 - \gamma(t)} \left[R^{-1}(t-1) - \frac{R^{-1}(t-1)\psi(t)\psi^T(t)\frac{R^{-1}(t-1)}{1 - \gamma(t)}}{\frac{\Lambda(t)}{\gamma(t)} + \psi^T(t)\frac{R^{-1}(t-1)}{1 - \gamma(t)}\psi(t)} \right]. \quad (\text{A.4})$$

Let us recall that the forgetting factor $\lambda(t)$ is linked to $\gamma(t)$ by (17):

$$\lambda(t) = \frac{\gamma(t-1)}{\gamma(t)}(1 - \gamma(t)). \quad (\text{A.5})$$

Equation (A.2) can be rewritten as:

$$P(t-1) = \gamma(t-1) R^{-1}(t-1). \quad (\text{A.6})$$

Combining (A.5) and (A.6) gives:

$$\frac{P(t-1)}{\lambda(t)} = \frac{\gamma(t)}{1 - \gamma(t)} R^{-1}(t-1), \quad (\text{A.7})$$

and introducing (A.7) in (A.4) gives quite directly:

$$P(t) = \frac{1}{\lambda(t)} \left[P(t-1) - \frac{P(t-1)\psi(t)\psi^T(t)P(t-1)}{\lambda(t)\Lambda(t) + \psi^T(t)P(t-1)\psi(t)} \right]. \quad (\text{A.8})$$

Taking:

$$L(t) = \gamma(t) R^{-1}(t) \psi(t) \Lambda^{-1}(t) \quad (\text{A.9})$$

in (A.1c) gives:

$$\hat{\Theta}(t) = \hat{\Theta}(t-1) + L(t)\varepsilon(t). \quad (\text{A.10})$$

Introducing (A.2) and (A.8) in (A.9) yields after some calculations:

$$L(t) = \frac{P(t-1)\psi(t)}{\lambda(t)\Lambda(t) + \psi^T(t)P(t-1)\psi(t)}, \quad (\text{A.11})$$

and substituting (A.11) in (A.8) gives with (A.1a), (A.11) and (A.10) the final algorithm:

$$\varepsilon(t) = y(t) - \hat{y}(t/\hat{\Theta}(t-1)), \quad (\text{A.12a})$$

$$L(t) = \frac{P(t-1)\psi(t)}{\lambda(t)\Lambda(t) + \psi^T(t)P(t-1)\psi(t)}, \quad (\text{A.12b})$$

$$P(t) = \frac{1}{\lambda(t)} \left[P(t-1) - L(t)\psi^T(t)P(t-1) \right], \quad (\text{A.12c})$$

$$\hat{\Theta}(t) = \hat{\Theta}(t-1) + L(t)\varepsilon(t). \quad (\text{A.12d})$$

References

- Aström, K. J., 1980. Maximum likelihood and prediction error methods. *Automatica* 16 (5), 551–574.
- Bhat, N., McAvoy, T. J., 1990. Use of neural nets for dynamic modeling and control of chemical process systems. *Computers & Chemical Engineering* 14 (4-5), 573–582.
- Billings, S. A., Jamaluddin, H. B., Chen, S., 1991. A comparison of the backpropagation and recursive prediction error algorithms for training neural networks. *Mechanical Systems and Signal Processing* 5 (3), 233–255.
- Box, G. E. P., Jenkins, G. M., 1970. *Time Series Analysis: Forecasting and Control*. Holden Day, San Francisco.
- Chen, S., Billings, S. A., Grant, P. M., 1990a. Non-linear system identification using neural networks. *International Journal of Control* 51 (6), 1191–1214.
- Chen, S., Cowan, C. F. N., Billings, S. A., Grant, P. M., 1990b. Parallel recursive prediction error algorithm for training layered neural networks. *International Journal of Control* 51 (6), 1215–1228.
- Cybenko, G., 1989. Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals and Systems* 2 (4), 303–314.
- Funahashi, K.-I., 1989. On the approximate realization of continuous mappings by neural networks. *Neural Networks* 2 (3), 183–192.
- Grossberg, S., 1988. Non-linear neural networks: Principles, mechanisms, and architectures. *Neural Networks* 1 (1), 17–61.
- Hoskins, J. C., Himmelblau, D. M., 1988. Artificial neural network models of knowledge representation in chemical engineering. *Computers & Chemical Engineering* 12 (9-10), 881–890.
- Huber, P. J., 1964. Robust estimation of a location parameter. *Ann. Math. Stat.* 35 (1), 73–101.
- Kohonen, T., 1989. *Self-Organisation and Associative Memory*, 3rd Edition. Springer-Verlag, New-York.
- Kramer, M. A., Leonard, J. A., 1988. Diagnosis using backpropagation neural networks - analysis and criticism. *Computers & Chemical Engineering* 14 (12), 1323–1338.
- Lippman, R. P., 1987. An introduction to computing with neural nets. *IEEE ASSP Magazine* 4 (2), 4–22.
- Ljung, L., 1987. *System Identification - Theory for user*. Prentice Hall, Englewood Cliffs, NJ.
- Masreliez, C. J., Martin, R. D., 1977. Robust Bayesian estimation for the linear model and robustifying the Kalman filter. *IEEE Transactions on Automatic Control* 22 (3), 361–371.
- Nahas, E. P., Henson, M. A., Seborg, D. E., 1992. Non-linear internal model control strategy for neural network models. *Computers & Chemical Engineering* 16 (12), 1039–1057.
- Narendra, K. S., Parthasarathy, K., 1990. Identification and control of dynamical systems using neural networks. *IEEE Transactions on Neural Networks* 1 (1), 4–27.
- Puthenpura, S., Sinha, N. K., 1990. A robust recursive identification method. *Control Theory and Advanced Technology* 6, 683–695.
- Rajavelu, A., Musavi, M., Shirvaikar, M., 1989. A neural network approach to character recognition. *Neural Networks* 2 (5), 387–393.
- Rumelhart, D. E., McClelland, J. L., PDP Research Group, 1986. *Parallel distributed processing*. Vol. 1. MIT Press, Cambridge, MA.
- Söderström, T., Stoica, P., 1989. *System Identification*. Prentice-Hall International, Hemel Hempstead, UK.
- Widrow, B., Lehr, M., 1990. 30 years of adaptive neural networks: perceptron, Madaline, and backpropagation. *Proceedings of the IEEE* 78 (9), 1415–1442.
- Xue, Q., Hu, Y. H., Tompkins, W. J., 1992. Neural-network-based adaptive matched filtering for QRS detection. *IEEE Transactions on Biomedical Engineering* 39 (4), 317–329.