

# A sequential algorithm for feed-forward neural networks with optimal coefficients and interacting frequencies

Enrique Romero\*, René Alquézar

*Departament de Llenguatges i Sistemes Informàtics, Universitat Politècnica de Catalunya, Barcelona, Spain*

Received 17 December 2003; received in revised form 12 July 2005; accepted 13 July 2005

Available online 9 November 2005

Communicated by D. Stork

## Abstract

An algorithm for *sequential approximation with optimal coefficients and interacting frequencies (SAOCIF)* for feed-forward neural networks is presented. SAOCIF combines two key ideas. The first one is the optimization of the coefficients (the linear part of the approximation). The second one is the strategy to choose the frequencies (the non-linear weights), taking into account the interactions with the previously selected ones. The resulting method combines the locality of sequential approximations, where only one frequency is found at every step, with the globality of non-sequential methods, where every frequency interacts with the others. The idea behind SAOCIF can be theoretically extended to general Hilbert spaces. Experimental results show a very satisfactory performance.

© 2005 Elsevier B.V. All rights reserved.

*Keywords:* Bias/variance decomposition; Feed-forward neural networks; Sequential approximations

## 1. Introduction

In terms of the bias/variance decomposition, as far as the number of hidden units of a feed-forward neural network (FNN) grows, bias decreases and variance increases. This happens because the flexibility of the model also grows with the number of hidden units [17,4]. Therefore, choosing an adequate architecture is a very important issue in order to obtain good generalization performance. We will focus on FNNs with one hidden layer of units (two layers of weights), including both multi-layer perceptrons (MLPs) and radial basis function networks (RBFNs).

Sequential approximation algorithms for FNNs (also named incremental or constructive) allow to dynamically construct the network without setting a priori the number of hidden units. They can help to find a proper trade-off between bias and variance by controlling, among other things, the number of hidden units. These methods start

with a small network (usually with no hidden units), and sequentially add new terms (that is, hidden units, each associated with a frequency) until a satisfactory solution is found. We will refer to the weights in the first layer (the non-linear weights) as *frequencies* and to the weights in the second layer (the linear weights) as *coefficients*.

Most of the sequential models found in the literature keep the previously selected frequencies fixed. Among these sequential models, many of them choose the new term so that it matches the previous residue as best as possible (see Section 2). It is well known that, although this strategy leads to approximations convergent towards the target function, it may be far from being the best strategy. This fact can be observed in the example in Fig. 1: when approximating the vector  $f$  with  $v_1$  and  $v_2$  we obtain  $X_2$ ; clearly, this is not the best possible approximation, since  $v_1$  and  $v_2$  form a basis of  $\mathbb{R}^2$ . In this case, optimizing the coefficients of the previously added terms would lead to a much better approximation (exact, in fact) of the target vector. But recalculating the coefficients is not enough, as illustrated in the example in Fig. 2. Suppose that  $X_1$  is a partial approximation of  $f$ , and  $h$  is the vector which best matches the residue  $r_1$ . Since  $h$  does not lie on the plane

\*Corresponding author. Tel.: +34934137796.

*E-mail addresses:* [eromero@lsi.upc.edu](mailto:eromero@lsi.upc.edu) (E. Romero), [alquezar@lsi.upc.edu](mailto:alquezar@lsi.upc.edu) (R. Alquézar).

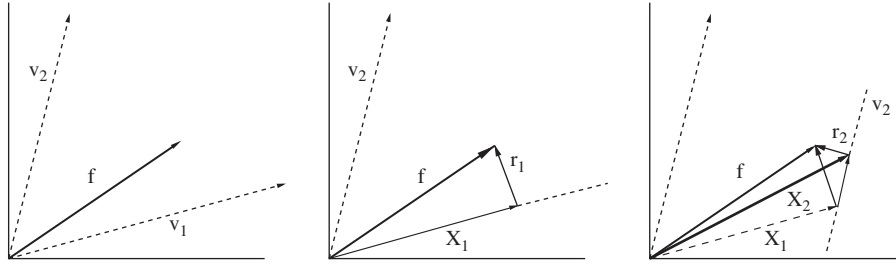


Fig. 1. Sequence of the approximation of a vector  $f$  in  $\mathbb{R}^2$  with  $v_1$  and  $v_2$  matching the previous residue without recalculating the coefficients. In the first step (middle),  $X_1$  is obtained. In the second step (right),  $r_1$  is approximated with  $v_2$ . The resulting vector ( $X_2$ ) is not the best approximation that can be achieved with  $v_1$  and  $v_2$ . In this case, optimizing the coefficients allows to obtain  $f$ .

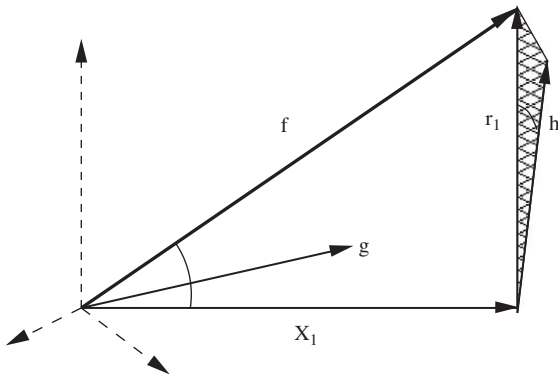


Fig. 2. Approximation of a vector  $f$  in  $\mathbb{R}^3$  matching the previous residue and recalculating the coefficients. Suppose that  $h$  is the vector that best matches the residue  $r_1$ . The vector  $g$ , which lies on the plane that contains  $f$  and  $X_1$ , allows an exact approximation to  $f$  when combined with  $X_1$ . The vector  $h$  (not on this plane) does not have this property. Optimizing the coefficients is not enough if  $h$  is selected.

that contains  $X_1$  and  $f$ , it is not necessarily the vector that, together with  $X_1$ , best approximates the target vector  $f$ . Any vector lying on the plane that contains  $X_1$  and  $f$  ( $g$ , for example) allows an exact approximation of  $f$ . Regardless of the coefficients optimization, matching the residue does not take into account the interactions with the previously selected terms. Any vector lying on the plane that contains  $f$  and a vector of the subspace spanned by the previous terms allows an exact approximation of the target vector. The vector that best matches the residue does not necessarily satisfy this property.

Important exceptions to the idea of matching the residue are the orthogonal least squares learning (OLSL) algorithm [11] and the kernel matching pursuit (KMP) with *pre-fitting* algorithm [48], where an (implicit or explicit) orthogonalization procedure is performed. In these algorithms, every point in the data set is considered as a candidate frequency. After calculating the optimal coefficients for every possible network, the best one (in terms of the minimum squared error) is selected.

In this work, we present an algorithm for *sequential approximation with optimal coefficients and interacting frequencies* (SAOCIF) for FNNs, which combines two key ideas. On the one hand, it optimizes the coefficients, so that the best approximation with the selected vectors is

always achieved, as in Fig. 1. On the other, the frequencies are selected at every step taking into account the interactions with the previously selected terms in a flexible manner. The interactions are discovered by means of their optimal coefficients. In the example in Fig. 2, SAOCIF would select  $g$  (instead of  $h$ ), because it allows a better approximation of  $f$  when combined (interacts) with  $X_1$ . The contribution of the new frequency is measured in terms of its capability of approximation to the target vector together with the previously selected frequencies. There is no explicit intention to match the residue. That is the idea of *interacting frequencies*. Therefore, it may be able to obtain, with the same number of hidden units, better approximations than matching the residue as best as possible. In terms of the bias/variance decomposition, it will be possible to obtain simpler models with the same bias, since the same level of approximation may be achieved with less hidden units.

The proposed algorithm can be seen as an extension and generalization of the OLSL and KMP with *pre-fitting* algorithms in several ways. First, it is not restricted to select the candidate frequencies from the points in the data set. In this sense, a number of different heuristics can be used to select the frequencies. Second, it allows to choose the activation function for every hidden unit. Finally, it is possible to further tune the selected frequencies with any non-linear optimization technique.

The idea behind SAOCIF can be extended to approximation in Hilbert spaces, maintaining orthogonal-like properties. The theoretical results obtained prove that, under reasonable conditions, the residue of the approximation is (in the limit) the best one that can be obtained with any subset of the given set of vectors. The importance of the interacting frequencies lies in the hypothesis that, as it can be seen in Fig. 2, it seems more plausible to find better partial approximations selecting the new frequency taking into account the interactions with the previous frequencies than just matching the residue as best as possible. Therefore, the rate of convergence is expected to increase.

Several experiments were performed in order to test the algorithm. Experimental results show a very satisfactory performance when compared to other sequential approaches. In particular, SAOCIF works better than methods that select the new frequencies based on the idea

of matching the residue, confirming the suitability of the interacting frequencies approach. When the number of hidden units is fixed a priori, better performance is obtained (with a moderate increase in the execution times). When the desired approximation accuracy is fixed a priori, SAOCIF allows to obtain models with less hidden units. These results are in agreement with the bias/variance decomposition.

The rest of the paper is organized as follows. An overview of some sequential approximations can be found in Section 2. The definition of SAOCIF and the algorithm are presented in Section 3. The extension to general Hilbert spaces is discussed in Section 4. The experimental results are described in Section 5. Finally, Section 6 concludes and outlines some directions for further research.

## 2. Sequential approximations

In this section, several sequential approximations for regression are described. Matching the residue is the underlying idea for most of the proposed schemes. Although some of them optimize the coefficients, the frequencies are selected without taking into account the interactions with the previously selected terms. Important exceptions are the OLSL and the KMP with *pre-fitting* algorithms, where an (implicit or explicit) orthogonalization is performed. An extensive review on sequential methods for regression can be found in [26].

### 2.1. Dynamic node creation

The dynamic node creation method [2] is a sequential method where, during the training, a new hidden node is added when the rate of decrease of the average squared error is less than a certain value. After a new node is added, the whole network is trained with standard back-propagation until the solution is satisfactory or another node is needed. Several variants of the dynamic node creation procedure can be found in the literature [7,5,50,43].

### 2.2. Resource-allocating network

The resource-allocating network is presented in [36]. When the network performs well on a presented pattern, the whole network is trained. Otherwise, a new Gaussian RBF hidden unit is added. Therefore, a memorization of training patterns is performed in some sense. The training of the whole network is performed with gradient descent. Several variations of this model can be found in [23,52,42]. All of them are specific for RBFNs.

### 2.3. Projection pursuit

Projection pursuit is a family of optimization methods which can be found in the statistics literature [18]. Projection pursuit regression (PPR) [16], as a particular case of function approximation, estimates the conditional

expectation of a random variable  $Y \in \mathbb{R}$  given  $X \in \mathbb{R}^l$  by means of a sum of ridge functions  $\sum_{j=1}^N g_j(a_j \cdot x)$  as follows (the  $a_j$ 's act as the frequencies). Suppose that the first  $N - 1$  terms of the approximation have been determined. That is, the vectors  $a_j$  and the functions  $g_j$  ( $1 \leq j \leq N - 1$ ) have been calculated. Let  $r_{N-1}(x)$  be the residue at step  $N - 1$ . Find  $a_N$  and  $g_N$  such that  $\|r_{N-1}(x) - g_N(a_N \cdot x)\|$  is the minimum. This process is repeated until the residue is smaller than a user-defined threshold. This procedure may be improved by *back-fitting*: omit some of the earlier summands  $g_j$ , determine better replacements, and then iterate. In [21] it is proved that, under mild smoothness conditions, a minimizing direction  $a_N$  exists at every step such that  $\lim_{N \rightarrow \infty} E[r_N] = 0$ . Later, it was proved that the convergence may be accelerated approximating by an optimal convex combination [22]. The upper bounds for the rate of convergence of approximations with FNNs derived in [3] are based on this result.

#### 2.3.1. Projection pursuit in signal processing

Some methods with the same underlying ideas as projection pursuit can be found in the signal processing area. In [33], matching pursuit (MP) is described, an algorithm that decomposes any signal into a linear expansion of waveforms that are selected from a (possibly redundant) dictionary of functions  $F$ , with  $\|g_\gamma\| = 1$  for every  $g_\gamma \in F$ . Similar to PPR, the MP algorithm works choosing at step  $N$  an element  $g_{\gamma_N} \in F$  which closely matches the residue  $R^N f$ , that is  $|\langle R^N f, g_{\gamma_N} \rangle| \geq \alpha \sup_{g_\gamma \in F} |\langle R^N f, g_\gamma \rangle|$ , where  $0 < \alpha \leq 1$ . The convergence property of MP is proved in [33], based on the results in [21]. After  $N$  steps, a recalculation of the coefficients can be made, named *back-projection*, to approximate  $f$  at best with the selected vectors. This idea was formalized in the orthogonal MP algorithm [35]. Similar results are obtained in [38] with a particular set of functions: the normalized Gaussian functions with adjustable variance and time-frequency center.

#### 2.3.2. Projection pursuit with neural networks

The two layer architecture of a neural network is well suited to construct an approximation with PPR. The projection pursuit learning network (PPLN) [19] is modeled as a one hidden layer MLP that learns unit by unit, and layer by layer cyclically after all the training patterns are presented. Weights are learned while the other ones remain fixed, and no global optimization of the coefficients is performed.

The incremental linear quasiparallel (ILQ) algorithm, a sequential algorithm for neural networks based on the ideas of PPR and MP, is presented in [25]. Every iteration consists of two steps. In the first one, the frequency of the new hidden unit is determined. In the second one, all output weights are recalculated. Given a set of functions  $G$ , the frequencies are determined trying to find  $g_N \in G$  so that nearly maximizes  $|\langle f - f_{N-1}, g \rangle|$  over  $g \in G$ . Therefore, the resulting method is similar to orthogonal MP, where there

exists a *back-projection* procedure at every iteration. Based on the results in [22], upper bounds for the rate of convergence are derived in [25].

With similar ideas, a number of objective functions to obtain the new frequency are explored in [27], all of them with the aim of matching the residue.

#### 2.4. Cascade-correlation

In the neural networks literature, the most used constructive method is, probably, cascade-correlation (CASCOR) [14]. CASCOR combines two key ideas. The first one is the cascade architecture, in which the newly added hidden unit receives inputs from the input layer as well as from the previously added hidden neurons. The second one is the learning algorithm. For each new hidden unit, the algorithm tries to maximize the correlation (or, more precisely, the covariance) between the output of the new unit and the residual error signal of the network. In order to maximize this function, a gradient ascent is performed. The input weights of the hidden units are frozen at the time the new unit is added to the network, and only the output connections are trained. There exist variations of CASCOR where the architecture is not cascaded, and where the learning rule is changed to train directly for minimization of the output errors (see [37] for details). Anyway, the frequencies obtained by these methods are the result of matching the residue at the previous step with only one term. There exist many variants of the original CASCOR algorithm in the literature (see, for example, [20,28,30–32,46]). A comparison of some of them can be found in [29].

#### 2.5. Orthogonal sequential neural networks methods

The orthogonal least squares learning (OLSL) algorithm is proposed in [11], a learning procedure for RBFNs based on the orthogonal least squares method [9]. The procedure starts with a single Gaussian RBF hidden unit and it sequentially increases the number of hidden units, one at a time, until the model error is satisfactory. The frequency of the new hidden unit (the center) is selected among the points in the data set. The classical Gram–Schmidt orthogonalization method is used at each step to form a set of orthogonal vectors for the space spanned by the output vectors of the previously selected hidden units. For every point in the data set, the orthogonal component of its output vector to that space is obtained. The new information introduced by this unit is caused by that part of its output vector which is orthogonal to the space spanned by the output vectors of previously selected hidden units. In this context, an output vector is an element of  $\mathbb{R}^T$ , where  $T$  is the number of patterns, obtained by applying the Gaussian function to every point in the data set. After computing its optimal coefficient (with the squared-error function), the point in the data set maximizing the error reduction ratio is selected. The procedure is

terminated when a predetermined percentage of the total error is reduced. Some extensions of the original procedure can be found in [10,12].

Recently, the kernel matching pursuit (KMP) algorithm was described [48], an extension of MP that can be used to build kernel-based solutions to supervised machine learning problems. The emphasis of the KMP scheme is put on the building of an alternative to support vector machines (SVMs) [47] that controls the sparsity of the solution (that is, the number of support vectors). Whereas good generalization abilities of SVMs are related to margin maximization, KMP is designed to build sparse kernel-based solutions minimizing the squared error function. As in SVMs, the frequencies of the resulting network are a subset of the points in the data set. Given a data set  $D$ , the dictionary (needed in MP) is defined as the set of functions  $F = \{K(x, x_i) : x_i \in D\}$ , where  $K$  is a symmetric positive definite kernel function. Three versions of KMP are defined in [48]: basic KMP (similar to basic MP), KMP with *back-fitting* at every step (similar to orthogonal MP) and KMP with *pre-fitting* (similar to OLSL). The optimization problems posed can be solved exactly because a finite dictionary is used. Experimental comparisons between KMP with *pre-fitting* and SVMs for several classification problems show comparable results with typically much sparser models for KMP with *pre-fitting*. A very similar method, but particular for Gaussian processes, can be found in [44]. The main difference lies in the loss function to be optimized.

### 3. Definition of SAOCIF and algorithm

As previously mentioned, most of the sequential models for FNNs found in the literature choose the new term so that it matches the previous residue as best as possible. This strategy can be far from being the best one, since it does not take into account the interactions with the previously selected terms. In this section SAOCIF is presented, a sequential scheme for FNNs where the contribution of the new frequency is measured in terms of its capability of approximation to the target vector together with the previously selected ones.

#### 3.1. Definition

**Definition.** Let  $H$  be the Hilbert space  $\mathbb{R}^T$ , where  $T$  is the number of patterns in a data set  $D = \{x_1, \dots, x_T\}$ ,  $f = (f_1, \dots, f_T) \in H$  the target vector and  $\Omega$  a space of frequencies. A SAOCIF for FNNs is a sequence of vectors  $\{X_N\}_{N \geq 0}$  in  $H$  whose terms are defined as

- (1)  $X_0 = 0$ .
- (2)  $X_N = \sum_{k=1}^{N-1} \lambda_k^N v_{\omega_k} + \lambda_N^N v_{\omega_N}$ , so that

- (a) The coefficients  $\lambda_1^N, \dots, \lambda_{N-1}^N, \lambda_N^N$  are optimal. That is, the vector  $X_N$  is the best approximation of  $f$  with



vectors  $v_{\omega_1}, \dots, v_{\omega_{N-1}}, v_{\omega_N} \in H$  in terms of minimizing the squared error  $\|f - X_N\|^2$  (that is, with the metric induced by the usual inner product in  $\mathbb{R}^T$ ).

- (b) The frequency  $\omega_N \in \Omega$  is selected taking into account the interactions of  $v_{\omega_N}$  with  $v_{\omega_1}, \dots, v_{\omega_{N-1}}$  in order to minimize  $\|f - X_N\|^2$ .

**Remarks.**

(1) In FNNs terminology, every frequency  $\omega_k \in \Omega$  is associated with a hidden unit  $\varphi_k(\omega_k, x)$ , where  $\varphi_k$  is the activation function. The  $i$ th component of  $v_{\omega_k}$  is the value of the hidden unit  $\varphi_k(\omega_k, x)$  at the  $i$ th point in  $D$ . That is,  $v_{\omega_k} = (\varphi_k(\omega_k, x_1), \dots, \varphi_k(\omega_k, x_T))$ . The output function of the network with  $N$  hidden units is  $X_N(x) = \sum_{k=1}^N \lambda_k^N \varphi_k(\omega_k, x)$ .

(2) At step  $N$ , a new frequency ( $\omega_N$ ) is considered, the number of terms of the approximation is increased by one ( $\lambda_N^N v_{\omega_N}$ ), and the coefficients  $\lambda_1^N, \dots, \lambda_{N-1}^N$  are optimized in order to obtain the best approximation of  $f$  with vectors  $v_{\omega_1}, \dots, v_{\omega_{N-1}}, v_{\omega_N}$ . The frequencies  $\omega_1, \dots, \omega_{N-1}$  are kept fixed. The vectors  $v_{\omega_1}, \dots, v_{\omega_{N-1}}, v_{\omega_N}$  are not necessarily mutually orthogonal.

(3) As it is well known [1], since  $X_N$  is the best approximation of  $f$  with  $v_{\omega_1}, \dots, v_{\omega_{N-1}}, v_{\omega_N}$ , it holds that

$$\forall k : 1 \leq k \leq N \quad \langle f - X_N, v_{\omega_k} \rangle = 0, \tag{1}$$

where  $\langle \cdot, \cdot \rangle$  is the inner product in  $H$ . That is,  $f - X_N$  is orthogonal to the space spanned by  $v_{\omega_1}, \dots, v_{\omega_{N-1}}, v_{\omega_N}$ . By definition of inner product, (1) is equivalent to the following linear equations system:

$$A_N(\lambda_1^N, \dots, \lambda_{N-1}^N, \lambda_N^N)^t = (\langle f, v_{\omega_1} \rangle, \dots, \langle f, v_{\omega_{N-1}} \rangle, \langle f, v_{\omega_N} \rangle)^t, \tag{2}$$

where  $A_N[i, j] = \langle v_{\omega_i}, v_{\omega_j} \rangle$  for  $1 \leq i, j \leq N$ . Thus, once the frequencies  $\omega_1, \dots, \omega_{N-1}, \omega_N \in \Omega$  have been selected, the optimal coefficients  $\lambda_1^N, \dots, \lambda_{N-1}^N, \lambda_N^N$  can be computed by solving (2).

(4) Using (1), it is immediate to verify that

$$\|f - X_N\|^2 = \left\| f - \sum_{k=1, k \neq j}^N \lambda_k^N v_{\omega_k} \right\|^2 - |\lambda_j^N|^2 \|v_{\omega_j}\|^2, \tag{3}$$

$$\|f - X_N\|^2 = \|f\|^2 - \|X_N\|^2, \tag{4}$$

$$\|X_N\|^2 = \sum_{k=1}^N \lambda_k^N \langle f, v_{\omega_k} \rangle. \tag{5}$$

As it can be observed, there is a great parallelism between these properties and those satisfied by an approximation with orthogonal vectors.

**3.2. Algorithm**

A sequential training algorithm for FNNs following the ideas of SAOCIF definition is presented in Fig. 3. Hidden units are added one at a time, choosing the frequencies in a

**Algorithm**

```

repeat
  Increase by 1 the number of hidden units N
  Pick an activation function for the new hidden unit
repeat
  Assign a candidate frequency  $\omega$  to the new hidden unit
  Compute the optimal coefficients  $\{\lambda_k\}_{1 \leq k \leq N}$  by solving (2)
  Set  $\omega_N := \omega$  if  $\|f - X_N\|^2$  is minimized
until the frequency  $\omega_N$  is satisfactory or
  there are no more candidate frequencies (criterion 1)
Optionally, tune the selected frequency  $\omega_N$ 
Fix the frequency  $\omega_N$  in the network
until the network is satisfactory (criterion 2)
end Algorithm
    
```

Fig. 3. An algorithm to construct an FNN following the ideas of SAOCIF.

flexible manner, so as to adjust the network until a satisfactory model is obtained. The algorithm works as follows. Suppose that we are at step  $N$  and we have a certain procedure to generate frequencies. For every candidate frequency, the optimal coefficients of the network are computed with that frequency installed in the new hidden unit, in order to test the real contribution of the frequency (together with the  $N - 1$  previously selected ones) to the approximation to the target vector. There is no explicit intention to match the residue. That is the idea of interacting frequencies. Note that, according to (4), maximizing  $\|X_N\|^2$  is equivalent to minimizing  $\|f - X_N\|^2$ . When the frequency is satisfactory or there are no more candidate frequencies (criterion 1, see below), the selected frequency can be optionally tuned.

Concerning the architecture needed to construct the approximation, it must have the following characteristics:

- (1) It must be a feed-forward architecture with a hidden layer of units (including both MLPs and RBFNs).
- (2) There are no restrictions about the dimension of the input and the output. There will be as many as the target function have. If there are several outputs, the total inner products must be calculated as the summation of the individual inner products for every output.
- (3) There is no restriction about the biases in the hidden units, since they can be treated as part of the frequencies.
- (4) There is no restriction about the activation functions in the hidden units. They can be, for example, sigmoidal, Gaussian, sines, cosines, wavelets, etc. Recent works have shown that the use of non-sigmoidal activation functions for MLPs may lead to very promising results [45]. Obviously, different units may have different activation functions. The output units must have a linear activation function.

As it can be seen, the only real restriction in the architecture is the linear activation function in the output units.

The resulting algorithm combines the locality of sequential approximations, where only one frequency is found at every step, with the globality of non-sequential methods, such as back-propagation [41], where every

frequency interacts with the others. The interactions are discovered by means of the optimal coefficients. The importance of the interacting frequencies lies in the hypothesis that they allow to find better partial approximations, with the same number of hidden units, than frequencies selected just to match the residue as best as possible. In terms of the bias/variance decomposition, it will be possible to obtain simpler models with the same bias, since the same level of approximation may be achieved with less hidden units.

### 3.3. Implementation

The strategy to select the candidate frequency is probably the most important part of the algorithm. In our experiments (see Section 5) three strategies were introduced in order to test the algorithm:

- (1) *Random strategy*: The frequencies are randomly selected within a certain range.
- (2) *Input strategy*: The frequencies are selected from the points in the data set (as is often the case in RBFNs) in a deterministic manner: for every hidden unit to be added, every point in the training set is tested as a candidate frequency.
- (3) *Breeder genetic algorithm (BGA) strategy*: The frequencies are selected with a more sophisticated strategy from the field of evolutionary algorithms, where a population of frequencies evolves driven by BGA [34] with the squared error as the fitness function.

The number of candidate frequencies of every strategy may be very different. Whereas for the Input one it is fixed (given the data set), the random and BGA strategies can be parameterized so as to assign as many candidate frequencies as desired. This is clearly related to the computational cost, as explained below.

Regarding the criterion 1 in Fig. 3, every strategy has its own one. For the random strategy, a fixed number of frequencies is selected. For the input one, the number of points in the data set determines the number of frequencies to test. For the BGA strategy, a maximum number of generations is determined. Regarding the criterion 2 in Fig. 3, many stopping criteria can be used: percentage of learned patterns, early stopping with a validation set, low relative rate of decrease of the error, etc.

SAOCIF satisfies a number of interesting properties to implement it in an efficient fashion, as explained next.

By (4) we have  $\|f - X_N\|^2 = \|f\|^2 - \|X_N\|^2$ , with  $\|f\|^2$  constant. Therefore, the frequency that minimizes the error is such that maximizes  $\|X_N\|^2$ . By (5) we know that  $\|X_N\|^2 = \sum_{k=1}^N \lambda_k^N \langle f, v_{\omega_k} \rangle$ . The values of  $\{\langle f, v_{\omega_k} \rangle\}_{1 \leq k \leq N}$  are the independent vector of the linear equations system (2) just solved to obtain  $\{\lambda_k^N\}_{1 \leq k \leq N}$ , which can be kept stored in memory. Therefore, once the coefficients have been obtained,  $\|X_N\|^2$  can be computed with cost  $O(N)$ , avoiding one pass through the data set (the cost of directly

computing  $\|X_N\|^2$  or  $\|f - X_N\|^2$  is  $O(T \cdot N)$ , where  $T$  is the number of examples). In our implementation,  $\omega_N$  was selected so as to maximize  $\|X_N\|^2$ , computed with (5).

Since the frequencies  $\omega_1, \dots, \omega_{N-1}$  are kept fixed, the proposed linear equations system at step  $N$  is equal to the system solved at step  $N-1$  for the selected frequency, but with a new row and a new column. Therefore, the system (2) posed at step  $N$  can be efficiently solved with bordered systems techniques [13] as follows. Note that (2) can be stated as

$$A_N(\lambda_1^N, \dots, \lambda_N^N)^t = \left( \begin{array}{c|c} A_{N-1} & a \\ \hline a^t & \gamma \end{array} \right) \begin{pmatrix} x \\ \eta \end{pmatrix} = \begin{pmatrix} b \\ \beta \end{pmatrix} \\ = (\langle f, v_{\omega_1} \rangle, \dots, \langle f, v_{\omega_N} \rangle)^t,$$

where  $x, a, b$  are vectors with  $N-1$  components, and  $\eta, \beta, \gamma$  are scalars. It can be easily verified that

$$\eta = \frac{\beta - a^t \cdot A_{N-1}^{-1} \cdot b}{\gamma - a^t \cdot A_{N-1}^{-1} \cdot a} \quad \text{and} \\ x = A_{N-1}^{-1} (b - \eta \cdot a). \quad (6)$$

Since  $A_{N-1}$  and  $b$  are equal for every candidate frequency, matrix  $A_{N-1}$  was inverted only once at every step, prior to the selection of the first frequency, and kept stored in memory. For every candidate frequency, its associated linear equations system was solved using (6). Whereas the inversion of a matrix has cost  $O(N^3)$ , the computational cost of multiplying it by a vector is  $O(N^2)$ , where  $N$  is the matrix dimension (in our case, the number of hidden units).

Regarding the computational cost, let  $N_{\text{SAOCIF}}$  be the number of hidden units of the resulting network,  $F$  the number of candidate frequencies for every added hidden unit (assuming that it is equal for every one),  $T$  the number of patterns and  $I$  the input dimension. Using the aforementioned implementation properties, and assuming  $N_{\text{SAOCIF}} \leq F$  and  $N_{\text{SAOCIF}} \leq T$ , the computational cost of the algorithm in Fig. 3 can be bounded by  $C_1 \cdot N_{\text{SAOCIF}}^2 \cdot F \cdot h_1(T \cdot I)$ , where  $C_1$  is constant and  $h_1(T \cdot I) \in O(T \cdot I)$ . The conditions  $N_{\text{SAOCIF}} \leq T$  and  $N_{\text{SAOCIF}} \leq F$  are not real restrictions, since it makes no sense, for example, constructing a network with more hidden units than examples in the data set. Similarly, the number of candidate frequencies is usually larger than the number of hidden units. In addition, these conditions make that the quadratic factor  $N_{\text{SAOCIF}}^2$  is not so relevant, in practice, as one could think (see Section 5). The computational cost mainly comes from the construction of the linear equations system for every candidate frequency rather than from solving it.

### 3.4. Comparison with other sequential schemes

The algorithm in Fig. 3 can be seen as an extension and generalization of the OLSL and KMP with *pre-fitting* algorithms (see Section 2) in several ways. First, it is not restricted to select the candidate frequencies from the

points in the data set. In this sense, a number of different heuristics can be used. Second, it allows to choose the activation function for every hidden unit. Finally, it is possible to further tune the selected frequency with any non-linear optimization technique. Note that the OLSL and KMP with *pre-fitting* algorithms can be simulated with the proposed algorithm by selecting the frequency with the input strategy without tuning.

#### 4. Extension to Hilbert spaces

In this section, SAOCIF is extended to general Hilbert spaces. Consider the following approximation problem: “Let  $H$  be a Hilbert space with inner product  $\langle \cdot, \cdot \rangle : H \times H \rightarrow \mathbb{C}$ , a space of parameters  $\Omega$ , and  $f \in H$  a vector to approximate with vectors  $v_\omega = v(\omega)$ ,  $v : \Omega \rightarrow H$ ,  $\omega \in \Omega$ , such that  $\forall \omega \in \Omega \|v_\omega\| \neq 0$ . Find  $\omega_1, \omega_2, \dots \in \Omega$  and  $\lambda_1, \lambda_2, \dots \in \mathbb{C}$  such that  $\lim_{N \rightarrow \infty} \|f - \sum_{k=1}^N \lambda_k v_{\omega_k}\| = 0$ ”. This definition is, in essence, the traditional one in approximation of vectors in Hilbert spaces [1]. The condition  $\|v_\omega\| \neq 0$  is equivalent, by the inner product properties, to that of  $v_\omega \neq 0$ . Observe that every vector  $v_\omega \in H$  depends on a parameter  $\omega \in \Omega$ .

In this setting, the same definition of SAOCIF stated in Section 3.1 is also valid for a general Hilbert space  $H$ , with the metric induced by its inner product. As previously, the term *frequency* refers to every  $\omega_1, \omega_2, \dots \in \Omega$ , and *coefficient* to every  $\lambda_1, \lambda_2, \dots \in \mathbb{C}$ . The following results show that SAOCIF converges towards the target vector in this general setting under reasonable conditions.

**Proposition 1.** *Let  $H$  be a Hilbert space, and  $f \in H$ . Any SAOCIF  $\{X_N\}_{N \geq 0}$  satisfies the following properties:*

- (Pa)  $\forall N \geq 0 \|f - X_{N+1}\|^2 \leq \|f - X_N\|^2$ .
- (Pb) *If  $M \geq N$ , then  $\langle f - X_M, f - X_N \rangle = \|f - X_M\|^2$ .*
- (Pc) *Suppose that  $\lambda_N^N \neq 0$ . The vector  $v_{\omega_N}$  is orthogonal to the space spanned by  $\{v_{\omega_1}, \dots, v_{\omega_{N-1}}\}$  if and only if the coefficients do not change between steps  $N - 1$  and  $N$ .*
- (Pd)  $\{X_N\}_{N \geq 0}$  *is convergent in  $H$ . That is,  $\exists g \in H \lim_{N \rightarrow \infty} \|g - X_N\| = 0$ .*

The proof can be found in Appendix A.

Observe that by (Pc), the only directions that guarantee that, without recalculating the coefficients, the approximation is optimal are the orthogonal directions. Therefore, if vectors are not mutually orthogonal, coefficients must be recalculated.

**Theorem 1.** *Let  $H$  be a Hilbert space,  $f \in H$  and a SAOCIF  $\{X_N\}_{N \geq 0}$ . Let  $g$  be such that  $\lim_{N \rightarrow \infty} \|g - X_N\| = 0$  (Pd), and suppose that for every  $\mu \in \mathbb{C}$  and every  $\omega_0 \in \Omega$  we have*

$$\|f - X_{N+1}\|^2 \leq \|f - (X_N + \mu v_{\omega_0})\|^2 + \alpha_N \quad (7)$$

*for every  $N \geq 0$ . That is, the approximation of  $f$  with  $X_{N+1}$  is better (up to  $\alpha_N$ ) than the best approximation of the*

*residue  $f - X_N$  that one could achieve with only one vector  $v_{\omega_0} \in v(\Omega)$ . If  $\limsup_{N \rightarrow \infty} \alpha_N \leq 0$ , then:*

(Ta) *The vector  $g$  satisfies:*

$$(Ta1) \forall \omega_0 \in \Omega \lim_{N \rightarrow \infty} \langle g - X_N, v_{\omega_0} \rangle = 0.$$

$$(Ta2) \forall \omega_0 \in \Omega \langle f, v_{\omega_0} \rangle = \langle g, v_{\omega_0} \rangle.$$

$$(Ta3) \langle f - g, g \rangle = 0.$$

(Ta4) *There is no subset of vectors in  $v(\Omega)$  that approximate  $f$  more than  $g$ . That is,*

$$\|f - g\| = \inf_{\mu_k \in \mathbb{C}, \psi_k \in \Omega} \left\| f - \sum_k \mu_k v_{\psi_k} \right\|.$$

(Tb) *If there exists  $A \subseteq \Omega$  such that the set of vectors  $\{v_\psi : \psi \in A\}$  spans  $H$ , then  $\{X_N\}_{N \geq 0}$  converges towards  $f$ .*

The proof can be found in Appendix B.

Observe that these results are not very restrictive, since the universal approximation capability is a necessary condition for the convergence property. Hence, SAOCIF allows (by selecting  $\Omega$  and  $v(\Omega)$ ) to choose any (or some) of the multiple vector families satisfying this property. In the particular case of  $L^2$ , it can be used with approximations by algebraic polynomials, Fourier series, wavelets and FNNs (including both MLPs and RBFNs), among others.

The hypothesis about the tolerance  $\alpha_N$  is, in essence, the same as in [22,3,25,24]. Imposing conditions on  $\alpha_N$ , as in [3,25], the rate of convergence can be proved to be  $O(1/\sqrt{N})$ , which has the same order as the optimal ones that can be found in the literature. However, the actual rate of convergence is expected to increase, since it seems more plausible to find better approximations (or equivalently, smaller  $\alpha_N$ ) selecting the new frequency taking into account the interactions with the previously selected ones than, for example, matching the residue as best as possible.

#### 5. Experiments

We now comment some generalities about the experiments performed in order to illustrate and test the algorithm for SAOCIF presented in Section 3.

In the tables of results, MPR means “matching the previous residue”, and it is a sequential scheme where the new frequency is selected so as to match the previous residue as best as possible and the previous coefficients are not recalculated. MPR can be implemented with a slight change in the algorithm in Fig. 3: the selected frequency maximizes  $\langle f - X_{N-1}, v_\omega \rangle^2 / \|v_\omega\|^2$  and the coefficient of the new hidden unit is computed as  $\lambda_N = \langle f - X_{N-1}, v_\omega \rangle / \|v_\omega\|^2$ . The idea behind MPR is exactly the same as in PPR, PPLN or MP without *back-projection* (see Section 2). As previously stated (see Section 3.3), let  $F$  be the number of candidate frequencies for every hidden unit,  $T$  the number of patterns and  $I$  the input dimension. In this setting, the computational cost of MPR can be bounded by  $C_2 \cdot N_{\text{MPR}} \cdot F \cdot h_2(T \cdot I)$ , where

Table 1  
NMSE in the training and test sets after the addition of 30 hidden units for the MPR, OCMPR and SAOCIF schemes

| Data set | Training |          |          | Test     |          |          |
|----------|----------|----------|----------|----------|----------|----------|
|          | MPR      | OCMPR    | SAOCIF   | MPR      | OCMPR    | SAOCIF   |
| HEA1     | 0.001937 | 0.000017 | 0.000000 | 0.002219 | 0.000073 | 0.000000 |
| HEA2     | 0.002146 | 0.000175 | 0.000028 | 0.005095 | 0.000710 | 0.000133 |
| HEA3     | 0.015259 | 0.006857 | 0.003222 | 0.059817 | 0.019668 | 0.011035 |
| HEA4     | 0.002875 | 0.002993 | 0.000406 | 0.009441 | 0.008655 | 0.000994 |
| HEA5     | 0.002725 | 0.001107 | 0.000506 | 0.007358 | 0.003459 | 0.001706 |

For every model, the frequencies were selected with the BGA strategy without tuning.

$C_2$  is constant,  $N_{\text{MPR}}$  is the number of hidden units of the network obtained with MPR and  $h_2(T \cdot I) \in O(T \cdot I)$ . When *back-projection* is present (that is, the coefficients are optimized after the selection of the frequency), we will refer to this method as “optimal coefficients after matching the previous residue” (OCMPR), equivalent to Orthogonal MP or ILQ. In this case, and assuming that  $N_{\text{OCMPR}} \leq F$ , the computational cost of OCMPR can be bounded by  $C_3 \cdot N_{\text{OCMPR}} \cdot F \cdot h_3(T \cdot I) + C_4 \cdot N_{\text{OCMPR}}^4$ , where  $C_3$  and  $C_4$  are constant,  $N_{\text{OCMPR}}$  is the number of hidden units of the resulting network and  $h_3(T \cdot I) \in O(T \cdot I)$ . As it can be seen, the respective computational costs of SAOCIF (see Section 3.3), MPR and OCMPR strongly depend on the number of hidden units of the final models ( $N_{\text{SAOCIF}}$ ,  $N_{\text{MPR}}$  and  $N_{\text{OCMPR}}$ ). This is a very important issue in order to assess the execution times found for each method.

Several activation functions (“AF” in the tables) have been tested, such as linear (lin), logistic (lgt), cosine (cos) in the MLP model and Gaussian (gau) in the RBFN model. Output units did not have biases, and in the hidden units they were only present for the logistic activation function. Regarding the strategies to select the frequencies, the column “WR” indicates the range of weights to look for candidate frequencies. When the value “input” is present, it means that the candidate frequencies are selected from the points in the training set, as explained in Section 3. Standard parameters were used for the BGA (see [6], for example). The average number of hidden units in the resulting networks is shown in column “NH”. An “NP” value means “not possible”, indicating that the learning of the training set was unsatisfactory.

### 5.1. HEA data sets

The data sets described in [19] were used to compare the approximation capabilities of MPR, OCMPR and SAOCIF. These data sets were constructed from five non-linear two-dimensional functions with different levels of complexity. Every training set contains 225 points generated from the uniform distribution  $U[0, 1]^2$  with their corresponding function value. The test set comprises 10,000 points of a regularly spaced grid on  $[0, 1]^2$ . No noise was added. We will refer to those data sets as HEA $_n$ , where  $n$  varies from 1

to 5. These data sets have been widely used in the literature (see, for example, [27,32,46]). Different from [19], 10 independent training sets were constructed and trained for every function. In some preliminary experiments, we observed that overfitting was not present with these data sets. Results are shown in Table 1 as the average, over 10 runs with the respective training sets for every function, of the normalized mean squared error

$$\text{NMSE} = \frac{\sum_{i=1}^T (\text{HEA}_n(x_i) - X_N(x_i))^2}{\sum_{i=1}^T (\text{HEA}_n(x_i) - \overline{\text{HEA}_n})^2}$$

after the addition of 30 hidden units. For every model (MPR, OCMPR and SAOCIF), the frequencies were selected with the BGA strategy (100 generations of an initial random population of 100 individuals) without tuning. The activation function was also selected for every unit, in addition to the frequency as follows: for a predefined set of activation functions, the process of selecting the frequency is repeated for every activation function in the set. Different units may have different activation functions. The set of activation functions contained the hyperbolic tangent, Gaussian function, sine, cosine or the polynomial kernel  $k(x) = (x + 1)^2$ , with MLP units.

As it can be observed in Table 1, the approximation capability of SAOCIF is superior to that of OCMPR, which in turn compares favorably with MPR. Generalization results are also better for SAOCIF than OCMPR and MPR. Since overfitting was not observed during the learning process with these data sets, the best results are obtained by those models that are able to fit more accurately the data. Therefore, the suitability of the idea of interacting frequencies is confirmed. Taking the mean execution time for MPR as 1, the relative execution times for OCMPR and SAOCIF in these experiments were 1.01 and 1.25, respectively (the execution time for MPR was 2'29" on a Pentium IV processor at 2.7 GHz). Note that the relative improvement of SAOCIF with respect to MPR and OCMPR (measured as the ratio between the respective NMSE in the test set) is always greater than 1.78. Therefore, a better performance is obtained with a moderate increase in the execution time.



Table 2  
Results for the *two spirals* problem with the random strategy (MPR, OCMPR, SAOCIF) and CASCOR

| Method | AF  | WR         | NH    | Training (%) | Test (%) |
|--------|-----|------------|-------|--------------|----------|
| MPR    | lgt | [−16, +16] | 500   | 98.97        | 99.48    |
| MPR    | lgt | [−8, +8]   | NP    | –            | –        |
| OCMPR  | lgt | [−8, +8]   | 105.2 | 100          | 100      |
| SAOCIF | lgt | [−8, +8]   | 91.7  | 100          | 100      |
| CASCOR | lgt | [−1, +1]   | 13.3  | 100          | 99.48    |

Columns “AF”, “WR” and “NH” indicate the activation function used, the range of weights to look for candidate frequencies and the average number of hidden units in the resulting networks, respectively. The reduced number of hidden units in CASCOR is due to its cascade architecture, but its average number of frequencies was 101.75.

### 5.2. The two spirals data set

The well-known *two spirals* problem consists in identifying the points of two interlocking spirals. Both training and test set comprise 194 points with balanced classes. This problem was tested with MPR, OCMPR and SAOCIF for the logistic activation function. The frequencies were selected randomly (200 attempts) within a certain range of weights. There was no further training to tune the frequency after a new hidden unit was added. The maximum number of hidden units added was 500, and no more hidden units were added when the whole training set was learned.

Results are shown in Table 2 as the average of 10 runs. The CASCOR algorithm (see Section 2) was also tested. Column “Train” indicates the percentage of the training set which has been learned. Column “Test” indicates the generalization performance obtained by an average-output committee of the resulting networks. As already known, this is a very hard problem for the logistic function because of its intrinsic high non-linearity and radial symmetry, but it could be learned with SAOCIF and an adequate (and very large) range of weights. MPR could not solve the problem with the same ranges, at least with 500 hidden units. Generalization in the obtained solutions was very good. Fig. 4 compares the generalization obtained by the CASCOR and SAOCIF algorithms.

In these experiments, OCMPR took less execution time than SAOCIF, which in turn was faster than MPR. The execution time for MPR was very high because of the large number of needed hidden units. Taking the mean execution time for OCMPR as 1, the relative execution times for SAOCIF and MPR in these experiments were 1.52 and 10.98, respectively (the execution time for OCMPR was 57" on a Pentium IV processor at 2.7 GHz). However, as expected from its better approximation capability, the number of hidden units in the obtained solutions is smaller for SAOCIF than for OCMPR and MPR.

### 5.3. The Pima Indians Diabetes data set

Another comparison was performed with a widely used medical diagnosis data set, namely the “Pima Indians Diabetes” database produced in the Applied Physics



Fig. 4. Generalization obtained by CASCOR (left) and SAOCIF (right), respectively with logistic functions for the *two spirals* problem.

Laboratory, Johns Hopkins University, 1988. The same data set as in [37] was used, with input values in [0, 1].

This data set was tested with MPR, OCMPR and SAOCIF for different activation functions. The input and the BGA strategies were tested to select the candidate frequencies. The BGA constructed, for every frequency, 20 generations of an initial random population of 100 individuals. Every selected frequency was further tuned with back-propagation. After every epoch, the coefficients were computed as in the selection step. Learning rates of the frequencies were automatically set following the *bold driver* technique proposed in [49], and the momentum was set to 0. The tuning was stopped when the relative rate of decrease of the error was less than 0.0005. In some of the experiments, linear activation functions were combined with non-linear activation functions in the hidden layer (“lin-fun” in the tables). When linear activation functions were present, their optimal frequencies can be calculated analytically, solving a linear equations system similar to (2) and setting the coefficients to 1. For the SAOCIF scheme, the coefficients of the linear hidden units are kept fixed when new non-linear units are added, in order to approximate only the non-linear component of the function. Different from the *two spirals* data set, previous experience with this data set suggested that it would be necessary to control the complexity of the model in order to obtain a good performance. A double 5–4-fold cross-validation (CV) was performed as follows [40]. First, a 5-fold CV (the outer CV) was performed to obtain 5 folds (4 folds to “learn” and 1 fold to test). Then, the 4 folds of the

Table 3

Results for the *Diabetes* data set with the input (left) and BGA (right) strategy for the MPR, OCMPR and SAOCIF schemes

| Method | AF      | Input |       |       | BGA         |       |       |
|--------|---------|-------|-------|-------|-------------|-------|-------|
|        |         | WR    | NH    | Test% | WR          | NH    | Test% |
| MPR    | lgt     | Input | 5.65  | 76.60 | [−0.5, 0.5] | 3.29  | 76.24 |
| MPR    | lin-lgt | Input | 14.25 | 77.12 | [−0.5, 0.5] | 11.16 | 77.12 |
| MPR    | cos     | Input | 3.80  | 76.47 | [−1.5, 1.5] | 1.90  | 76.99 |
| MPR    | lin-cos | Input | 3.65  | 76.47 | [−1.5, 1.5] | 3.27  | 77.23 |
| MPR    | gau     | Input | 9.10  | 76.73 | [0, 1]      | 7.66  | 76.78 |
| MPR    | lin-gau | Input | 11.90 | 76.99 | [0, 1]      | 6.65  | 76.97 |
| OCMPR  | lgt     | Input | 5.10  | 76.47 | [−0.5, 0.5] | 3.13  | 76.18 |
| OCMPR  | lin-lgt | Input | 3.90  | 76.47 | [−0.5, 0.5] | 4.13  | 76.63 |
| OCMPR  | cos     | Input | 2.85  | 76.99 | [−1.5, 1.5] | 1.97  | 77.02 |
| OCMPR  | lin-cos | Input | 4.90  | 76.73 | [−1.5, 1.5] | 3.93  | 76.86 |
| OCMPR  | gau     | Input | 6.95  | 76.73 | [0, 1]      | 4.52  | 76.05 |
| OCMPR  | lin-gau | Input | 7.75  | 76.60 | [0, 1]      | 5.43  | 76.44 |
| SAOCIF | lgt     | Input | 3.00  | 76.47 | [−0.5, 0.5] | 2.30  | 76.05 |
| SAOCIF | lin-lgt | Input | 4.30  | 77.25 | [−0.5, 0.5] | 4.90  | 77.07 |
| SAOCIF | cos     | Input | 2.40  | 76.34 | [−1.5, 1.5] | 2.47  | 77.18 |
| SAOCIF | lin-cos | Input | 4.20  | 77.39 | [−1.5, 1.5] | 3.12  | 77.39 |
| SAOCIF | gau     | Input | 4.80  | 78.04 | [0, 1]      | 5.42  | 77.41 |
| SAOCIF | lin-gau | Input | 4.90  | 77.12 | [0, 1]      | 4.70  | 77.36 |

Columns “AF”, “WR” and “NH” indicate the activation function used, the range of weights to look for candidate frequencies and the average number of hidden units in the resulting networks, respectively.

“learning set” of the outer CV were used as follows: 3 folds to train and 1 fold to validate, as in a 4-fold CV (the inner CV). Therefore, the number of trained models in a double 5–4-fold CV is 20. For the BGA strategy, this procedure was repeated 5 times, in order to alleviate the effect of random initial populations. Previous to every CV, the examples in the data set were randomly shuffled. For every training, no more hidden units were added when the error on the validation set did not improve for 5 consecutive hidden units.

The results are shown in Table 3, as the average in the network with minimum validation set error. The column “Test” indicates the mean generalization performance obtained by the average-output committees of the resulting networks in the inner CV. Although these results seem very similar for the different parameter configurations, there are some regularities which can be observed for this problem:

- (1) Non-linear activation functions different from the classical sigmoidal and Gaussian (such as cosines) may be satisfactorily used. Linear hidden units have mostly a positive influence on the results, and selecting the frequencies from the points in the data set seems well suited not only for RBFNs, as commonly used, but also for MLPs.
- (2) Interestingly, the performance of the input strategy was quite similar to that of the BGA strategy, with a much smaller computational cost. Although more experiments are needed, it may be an interesting idea to

obtain good and cheap models for data sets of moderate size.

- (3) The number of hidden units of the solutions obtained with SAOCIF is usually less than the number of hidden units of the solutions obtained with OCMPR and MPR. The mean number of hidden units for SAOCIF, OCMPR and MPR in Table 3 were 3.88, 4.55 and 6.86, respectively.
- (4) For every non-linear activation function tested, the best results are obtained with the SAOCIF scheme, when compared to OCMPR and MPR. This fact is related to the number of hidden units of the obtained solutions, since SAOCIF obtains simpler models (in terms of number of hidden units) with the same level of approximation.
- (5) In this case, the execution times were very similar for SAOCIF, OCMPR and MPR. The larger computational cost of every step for SAOCIF was compensated with the lower number of steps (number of hidden units) in the obtained solutions.

It is not easy to make a direct comparison of these results with previous studies on this benchmark. Since the results are very sensitive to the partitions made on the data set, a realistic comparison can only be made when similar experimental techniques are used. In Table 4 our results are compared to several methods found in the literature that have been applied to this data set with similar resampling techniques. In the same aforementioned conditions, the CASCOR algorithm was tested for more than 70

Table 4  
Comparison of results for the *Diabetes* data set

| Source    | Method                  | Test% |
|-----------|-------------------------|-------|
| [15]      | Boosting                | 75.60 |
| [51]      | Heterogeneous RBFNs     | 76.30 |
| [39]      | RBFNs                   | 75.90 |
| [39]      | AdaBoost                | 76.10 |
| [39]      | Support Vector Machines | 76.50 |
| [48]      | Kernel Matching Pursuit | 76.10 |
| This work | Cascade-Correlation     | 76.76 |
| This work | MPR (BGA/lin-cos)       | 77.23 |
| This work | OCMPR (BGA/cos)         | 77.02 |
| This work | SAOCIF (Input/gau)      | 78.04 |

different sets of parameters. Results are also shown in Table 4.

## 6. Conclusions and future work

An algorithm for SAOCIF for FNNs has been presented. Coefficients are optimized for every candidate frequency, so that the approximations maintain orthogonal-like properties. The new frequency is not chosen in order to match the previous residue as best as possible. It is selected with global approximation of the target vector criteria (interacting frequencies). The algorithm can be seen as an extension and generalization of the OLSL [11] and KMP with *pre-fitting* [48] algorithms in several ways. The idea behind SAOCIF can be extended to general Hilbert spaces. Theoretical results prove that, under reasonable conditions, the residue of the obtained approximation is (in the limit) the best one that can be obtained with any subset of the given set of vectors.

Experimental results show a very satisfactory performance for SAOCIF and several suggesting ideas for future experiments, such as the selection of the frequencies from the data set, as in OLSL and KMP, or the combination of linear and non-linear activation functions in the hidden units. The candidate frequencies can be selected with heuristics different from current strategies. In principle, a more intelligent selection could lead to better approximations. Likewise, the selection of the activation function for the new hidden unit admits any number of heuristics.

The theoretical results (see the proofs for details) are a consequence of the optimality of the coefficients in SAOCIF. The importance of the interacting frequencies lies in the hypothesis that it seems more plausible to find better partial approximations selecting the new frequency taking into account the interactions with the previous frequencies than, for example, matching the residue as best as possible. Therefore, the actual rate of convergence is expected to increase with this strategy. Experimental results confirmed this hypothesis. In order to theoretically prove this claim it would be necessary to formalize the concept of “interacting frequencies”, a non-trivial task in our opinion, since it is not an existential property but a procedural one.

## Acknowledgements

The authors want to thank the anonymous reviewers for their valuable comments and suggestions in order to prepare the final version of the paper.

This work has been supported by the Consejo Interministerial de Ciencia y Tecnología (CICYT), under projects DPI2002-03225 and CGL2004-04702-C02-02.

## Appendix A

- (Pa) Evident, since the coefficients of  $X_{N+1}$  are optimal.  
(Pb) Expressing  $f - X_N$  as  $f - X_M + X_M - X_N$  we have

$$\begin{aligned} \langle f - X_M, f - X_N \rangle \\ = \|f - X_M\|^2 + \langle f - X_M, X_M - X_N \rangle. \end{aligned}$$

By (1),  $\langle f - X_M, X_M - X_N \rangle = 0$  holds.

- (Pc) The necessity is clear by (1). To prove the sufficiency, suppose that  $X_N = X_{N-1} + \lambda_N^N v_{\omega_N}$ . By (1), for every  $1 \leq j \leq N-1$  we have  $\langle f - X_N, v_{\omega_j} \rangle = 0$  and  $\langle f - X_{N-1}, v_{\omega_j} \rangle = 0$ . Therefore, for every  $1 \leq j \leq N-1$  we have

$$\begin{aligned} 0 = \langle f - X_N, v_{\omega_j} \rangle &= \langle f - X_{N-1} - \lambda_N^N v_{\omega_N}, v_{\omega_j} \rangle \\ &= \langle -\lambda_N^N v_{\omega_N}, v_{\omega_j} \rangle, \end{aligned}$$

which finishes the proof, since  $\lambda_N^N \neq 0$ .

- (Pd) Since  $H$  is complete, it suffices to prove that  $\lim_{N,M \rightarrow \infty} \|X_M - X_N\|^2 = 0$ . Suppose that  $M > N$ . Expressing  $X_M - X_N$  as  $(X_M - f) + (f - X_N)$ , and using (Pb) we easily obtain

$$\|X_M - X_N\|^2 = \|f - X_N\|^2 - \|f - X_M\|^2.$$

Since the sequence  $\{\|f - X_N\|^2\}_{N \geq 0}$  is decreasing and positive (Pa), it is convergent. Hence,

$$\lim_{N,M \rightarrow \infty} (\|f - X_N\|^2 - \|f - X_M\|^2) = 0.$$

## Appendix B

- (Ta1) From Schwartz inequality, for every  $\omega_0 \in \Omega$  we have

$$|\langle g - X_N, v_{\omega_0} \rangle| \leq \|g - X_N\| \|v_{\omega_0}\|.$$

The proof finishes using that  $\lim_{N \rightarrow \infty} \|g - X_N\| = 0$ .

- (Ta2) Let  $\omega_0 \in \Omega$ . By hypothesis, for every  $N \geq 0$  and every  $\mu \in \mathbb{C}$

$$\begin{aligned} \|f - X_{N+1}\|^2 &\leq \|f - (X_N + \mu v_{\omega_0})\|^2 + \alpha_N \\ &= \|f - X_N\|^2 - 2 \operatorname{Re}(\langle f - X_N, \mu v_{\omega_0} \rangle) \\ &\quad + |\mu|^2 \|v_{\omega_0}\|^2 + \alpha_N. \end{aligned}$$

Expressing  $f - X_N$  as  $f - g + g - X_N$  we have

$$\begin{aligned} \|f - X_{N+1}\|^2 - \|f - X_N\|^2 \\ \leq |\mu|^2 \|v_{\omega_0}\|^2 - 2 \operatorname{Re}(\langle f - g, \mu v_{\omega_0} \rangle) \\ - 2 \operatorname{Re}(\langle g - X_N, \mu v_{\omega_0} \rangle) + \alpha_N. \end{aligned}$$

Hence, renaming  $D_N = \|f - X_N\|^2 - \|f - X_{N+1}\|^2$ ,

$$\begin{aligned} & 2 \operatorname{Re}(\langle f - g, \mu v_{\omega_0} \rangle) - |\mu|^2 \|v_{\omega_0}\|^2 \\ & \leq D_N - 2 \operatorname{Re}(\langle g - X_N, \mu v_{\omega_0} \rangle) + \alpha_N \\ & \leq D_N + 2|\mu| |\langle g - X_N, v_{\omega_0} \rangle| + \alpha_N \\ & = D_N + 2|\mu| |\langle g - X_N, v_{\omega_0} \rangle| + \alpha_N. \end{aligned}$$

for every  $\mu \in \mathbb{C}$ . Let  $\mu_0 = \langle f - g, v_{\omega_0} \rangle / \|v_{\omega_0}\|^2$ , and  $\varepsilon > 0$ .

Using (Pa), (Ta1), and the hypothesis, there exists  $N_0$  such that for every  $N \geq N_0$ ,

$$\|f - X_N\|^2 - \|f - X_{N+1}\|^2 \leq \varepsilon/3,$$

$$2|\mu_0| |\langle g - X_N, v_{\omega_0} \rangle| \leq \varepsilon/3,$$

$$\alpha_N \leq \varepsilon/3.$$

Thus, we have  $2 \operatorname{Re}(\langle f - g, \mu_0 v_{\omega_0} \rangle) - |\mu_0|^2 \|v_{\omega_0}\|^2 \leq \varepsilon$ . Since

$$\langle f - g, \mu_0 v_{\omega_0} \rangle = \overline{\mu_0} \langle f - g, v_{\omega_0} \rangle = |\mu_0|^2 \|v_{\omega_0}\|^2,$$

$2 \operatorname{Re}(\langle f - g, \mu_0 v_{\omega_0} \rangle) = 2|\mu_0|^2 \|v_{\omega_0}\|^2$  holds. Therefore, for every  $\varepsilon \geq 0$

$$|\mu_0|^2 \|v_{\omega_0}\|^2 = 2 \operatorname{Re}(\langle f - g, \mu_0 v_{\omega_0} \rangle) - |\mu_0|^2 \|v_{\omega_0}\|^2 \leq \varepsilon.$$

Hence,  $|\mu_0|^2 \|v_{\omega_0}\|^2 = 0$ . Since  $\|v_{\omega_0}\|^2 \neq 0$ , we derive  $\mu_0 = 0$ . Therefore,  $\langle f - g, v_{\omega_0} \rangle = 0$  for every  $\omega_0 \in \Omega$ .

(Ta3) Using (Ta2) and Schwartz inequality we have

$$\begin{aligned} |\langle f - g, g \rangle| &= |\langle f - g, g - X_N \rangle| \\ &\leq \|f - g\| \|g - X_N\|. \end{aligned}$$

The proof finishes using that  $\lim_{N \rightarrow \infty} \|g - X_N\| = 0$ .

(Ta4) According to (Ta2), any vector combination  $\sum_k \mu_k v_{\psi_k}$  in  $v(\Omega)$  satisfies  $\langle f - g, \sum_k \mu_k v_{\psi_k} \rangle = 0$ . Hence we have

$$\|f - g\|^2 = \left\langle f - g, f - \sum_k \mu_k v_{\psi_k} \right\rangle - \langle f - g, g \rangle.$$

Using (Ta3) and Schwartz inequality we have

$$\|f - g\|^2 \leq \|f - g\| \left\| f - \sum_k \mu_k v_{\psi_k} \right\|.$$

Therefore,  $\|f - g\| \leq \|f - \sum_k \mu_k v_{\psi_k}\|$  for any vector combination. The other inequality is clear, since for every  $N \geq 0$

$$\begin{aligned} \inf_{\mu_k \in \mathbb{C}, \psi_k \in \Omega} \left\| f - \sum_k \mu_k v_{\psi_k} \right\| &\leq \|f - X_N\| \\ &\leq \|f - g\| + \|g - X_N\|. \end{aligned}$$

The proof finishes using that  $\lim_{N \rightarrow \infty} \|g - X_N\| = 0$ .

(Tb) It is immediately derived from (Ta4).

## References

[1] N.I. Achieser, Theory of Approximation, Frederick Ungar, New York, 1956.

- [2] T. Ash, Dynamic node creation in backpropagation networks, *Connect. Sci.* 1 (4) (1989) 365–375.
- [3] A.R. Barron, Universal approximation bounds for superposition of a sigmoidal function, *IEEE Trans. Inform. Theory* 39 (3) (1993) 930–945.
- [4] A.R. Barron, Approximation and estimation bounds for artificial neural networks, *Machine Learn.* 14 (1) (1994) 115–133.
- [5] E.B. Bartlett, Dynamic node architecture learning: an information theoretic approach, *Neural Networks* 7 (1) (1994) 129–140.
- [6] L. Belanche, A case study in neural network training with the breeder genetic algorithm, Technical Report LSI-00-7-R, Departament de Llenguatges i Sistemes Informàtics, Universitat Politècnica de Catalunya, Spain, 2000.
- [7] M.G. Bello, Enhanced training algorithms, and integrated training/architecture selection for multilayer perceptron networks, *IEEE Trans. Neural Networks* 3 (6) (1992) 864–875.
- [9] S. Chen, S.A. Billings, W. Luo, Orthogonal least squares methods and their applications to non-linear system identification, *Int. J. Control* 50 (5) (1989) 1873–1896.
- [10] S. Chen, E.S. Chng, K. Alkadhimi, Regularized orthogonal least squares algorithm for constructing radial basis function networks, *Int. J. Control* 64 (5) (1996) 829–837.
- [11] S. Chen, C.F.N. Cowan, P.M. Grant, Orthogonal least squares learning algorithm for radial basis function networks, *IEEE Trans. Neural Networks* 2 (2) (1991) 302–309.
- [12] S. Chen, Y. Wu, B.L. Luk, Combined genetic algorithm and regularized orthogonal least squares learning for radial basis function networks, *IEEE Trans. Neural Networks* 10 (5) (1999) 1239–1243.
- [13] V.N. Faddeeva, Computational Methods of Linear Algebra, Dover Publications Inc., New York, 1959.
- [14] S.E. Fahlman, C. Lebiere, The cascade-correlation learning architecture, in: D.S. Touretzky (Ed.), *Advances in Neural Information Processing Systems*, vol. 2, Morgan Kaufmann, Los Altos, CA, 1990, pp. 524–532.
- [15] Y. Freund, R.E. Schapire, Experiments with a new boosting algorithm, in: 13th International Conference on Machine Learning, 1996, pp. 148–156.
- [16] J.H. Friedman, W. Stuetzle, Projection pursuit regression, *J. Am. Stat. Assoc.* 76 (1981) 817–823.
- [17] S. Geman, E. Bienenstock, R. Doursat, Neural networks and the bias/variance dilemma, *Neural Comput.* 4 (1) (1992) 1–58.
- [18] P.J. Huber, Projection pursuit, *Ann. Stat.* 13 (2) (1985) 435–475.
- [19] J.N. Hwang, S.R. Ray, M. Maechler, D. Martin, J. Schimert, Regression modelling in back-propagation and projection pursuit learning, *IEEE Trans. Neural Networks* 5 (3) (1994) 342–353.
- [20] M.M. Islam, K. Murase, A new algorithm to design compact two-hidden-layer artificial neural networks, *Neural Networks* 14 (9) (2001) 1265–1278.
- [21] L.K. Jones, On a conjecture of Huber concerning the convergence of projection pursuit regression, *Ann. Stat.* 15 (2) (1987) 880–882.
- [22] L.K. Jones, A simple lemma on greedy approximation in Hilbert space and convergence rates for projection pursuit regression and neural network training, *Ann. Stat.* 20 (1) (1992) 608–613.
- [23] V. Kadirkamanathan, M. Niranjan, A function estimation approach to sequential learning with neural networks, *Neural Comput.* 5 (6) (1993) 954–975.
- [24] V. Kůrková, Incremental approximation by neural networks, in: M. Karny, K. Warwick, V. Kůrková (Eds.), *Dealing With Complexity: A Neural Network Approach*, Springer, London, 1998, pp. 177–188.
- [25] V. Kůrková, B. Beliczynski, Incremental approximation by one-hidden-layer neural networks, in: *International Conference on Artificial Neural Networks*, vol. 1, 1995, pp. 505–510.
- [26] T.Y. Kwok, D.Y. Yeung, Constructive algorithms for structure learning in feedforward neural networks for regression problems, *IEEE Trans. Neural Networks* 8 (3) (1997) 630–645.
- [27] T.Y. Kwok, D.Y. Yeung, Objective functions for training new hidden units in constructive neural networks, *IEEE Trans. Neural Networks* 8 (5) (1997) 1131–1148.



- [28] J.J.T. Lahnajärvi, M.I. Lehtokangas, J.P.P. Saarinen, Fixed cascade error—a novel constructive neural network for structure learning, in: International Conference on Artificial Neural Networks in Engineering, 1999, pp. 25–30.
- [29] J.J.T. Lahnajärvi, M.I. Lehtokangas, J.P.P. Saarinen, Evaluation of constructive neural networks with cascaded architectures, *Neurocomputing* 48 (1–4) (2002) 573–607.
- [30] M.I. Lehtokangas, Modelling with constructive backpropagation, *Neural Networks* 12 (4–5) (1999) 707–716.
- [31] E. Littmann, H. Ritter, Cascade LLM networks, in: International Conference on Artificial Neural Networks, vol. 1, 1992, pp. 253–257.
- [32] L. Ma, K. Khorasani, New training strategies for constructive neural networks with application to regression problems, *Neural Networks* 17 (4) (2004) 589–609.
- [33] S.G. Mallat, Z. Zhang, Matching pursuits with time-frequency dictionaries, *IEEE Trans. Signal Process.* 41 (12) (1993) 3397–3415.
- [34] H. Mühlenbein, D. Schlierkamp-Voosen, Predictive models for the breeder genetic algorithm I. Continuous parameter optimization, *Evol. Comput.* 1 (1) (1993) 25–49.
- [35] Y.C. Pati, R. Rezaifar, P.S. Krishnaprasad, Orthogonal matching pursuit: recursive function approximation with application to wavelet decomposition, in: 27th Asilomar Conference on Signals, Systems and Computers, vol. 1, 1993, pp. 40–44.
- [36] J. Platt, A resource-allocating network for function interpolation, *Neural Comput.* 3 (2) (1991) 213–225.
- [37] L. Prechelt, Investigation of the CasCor family of learning algorithms, *Neural Networks* 10 (5) (1997) 885–896.
- [38] S. Qian, D. Chen, Signal representation using adaptive normalized Gaussian functions, *Signal Process.* 36 (1) (1994) 1–11.
- [39] G. Rätsch, T. Onoda, K.R. Müller, An improvement of AdaBoost to avoid overfitting, in: International Conference on Neural Information Processing, 1998, pp. 506–509.
- [40] B.D. Ripley, Statistical ideas for selecting network architectures, in: B. Kappen, S. Gielen (Eds.), *Neural Networks: Artificial Intelligence and Industrial Applications*, Springer, London, 1995, pp. 183–190.
- [41] D.E. Rumelhart, G.E. Hinton, R.J. Williams, Learning internal representations by error propagation, in: D.E. Rumelhart, J.L. McClelland (Eds.), *Parallel Distributed Processing: Explorations in the Microstructure of Cognition*, vol. 1, MIT Press, Cambridge, MA, 1986, pp. 318–362.
- [42] M. Salmerón, J. Ortega, C.G. Puntonet, A. Prieto, Improved RAN sequential prediction using orthogonal techniques, *Neurocomputing* 41 (1–4) (2001) 153–172.
- [43] Y. Shin, J. Ghosh, Ridge polynomial networks, *IEEE Trans. Neural Networks* 6 (3) (1995) 610–622.
- [44] A.J. Smola, P. Bartlett, Sparse greedy Gaussian process regression, in: *Advances in Neural Information Processing Systems*, vol. 13, MIT Press, Cambridge, MA, 2001, pp. 619–625.
- [45] J.M. Sopena, E. Romero, R. Alquézar, Neural networks with periodic and monotonic activation functions: a comparative study in classification problems, in: *Proceedings of the Ninth International Conference on Artificial Neural Networks*, vol. 1, 1999, pp. 323–328.
- [46] N.K. Treadgold, T.D. Gedeon, Exploring constructive cascade networks, *IEEE Trans. Neural Networks* 10 (6) (1999) 1335–1350.
- [47] V.N. Vapnik, *The Nature of Statistical Learning Theory*, Springer, New York, 1995.
- [48] P. Vincent, Y. Bengio, Kernel matching pursuit, *Machine Learn.* 48 (1–3) (2002) 165–187 (Special issue on new methods for model combination and model selection).
- [49] T.P. Vogl, J.K. Mangis, A.K. Rigler, W.T. Zink, D.L. Alkon, Accelerating the convergence of the back-propagation method, *Biol. Cybernet.* 59 (1988) 257–263.
- [50] Z. Wang, C. Di Massimo, M.T. Tham, J. Morris, A procedure for determining the topology of multilayer feedforward neural networks, *Neural Networks* 7 (2) (1994) 291–300.
- [51] D.R. Wilson, T.R. Martinez, Heterogeneous radial basis function networks, in: *International Conference on Neural Networks*, vol. 2, 1996, pp. 1263–1267.
- [52] L. Yingwei, N. Sundararajan, P. Saratchandran, A sequential learning scheme for function approximation using minimal radial basis function neural networks, *Neural Comput.* 9 (2) (1997) 461–478.



**Enrique Romero** was born in Barcelona, Spain, in 1966. He received the licenciate degree in Mathematics in 1989 from the *Universitat Autònoma de Barcelona*. In 1994, he received the licenciate degree in Computer Science from the *Universitat Politècnica de Catalunya* (UPC). In 1996, he joined the Department of *Llenguatges i Sistemes Informàtics*, UPC, as an assistant professor. He received the M.Sc. degree in Artificial Intelligence and the Ph.D. degree in Computer

Science from the UPC in 2000 and 2004, respectively. His research interests include Neural Networks, Support Vector Machines and Feature Selection.



**René Alquézar** was born in Barcelona, Spain, in 1962. He received the licenciate and Ph.D. degrees in Computer Science from the *Universitat Politècnica de Catalunya* (UPC), Barcelona, in 1986 and 1997, respectively. From 1987 to 1991 he was with the Spanish company NTE as technical manager of R&D projects on image processing applications for the European Space Agency (ESA). From 1991 to 1994 he was with the *Institut de Cibernètica*, Barcelona, holding a pre-doctoral research

grant from the Government of Catalonia and completing the doctoral courses of the UPC on Artificial Intelligence. He joined the Department of *Llenguatges i Sistemes Informàtics*, UPC, in 1994 as an assistant professor, and since March 2001 he has been associate professor in the same department. His current research interests include Neural Networks, Structural and Syntactic Pattern Recognition and Computer vision.