

PAPER • OPEN ACCESS

## Extreme learning machine adapted to noise based on optimization algorithms

To cite this article: A Vásquez *et al* 2020 *J. Phys.: Conf. Ser.* **1514** 012006

View the [article online](#) for updates and enhancements.



**IOP | ebooks™**

Bringing together innovative digital publishing with leading authors from the global scientific community.

Start exploring the collection—download the first chapter of every title for free.

# Extreme learning machine adapted to noise based on optimization algorithms

A Vásquez<sup>1</sup>, M Mora<sup>1</sup>, E Salazar<sup>1</sup>, and E Gelvez<sup>2</sup>

<sup>1</sup> Laboratorio de Investigaciones Tecnológicas en Reconocimiento de Patrones, Universidad Católica del Maule, Talca, Chile

<sup>2</sup> Facultad de Ciencias Básicas y Biomédicas, Universidad Simón Bolívar, San José de Cúcuta, Colombia

E-mail: marcomoracofre@gmail.com

**Abstract.** The extreme learning machine for neural networks of feedforward of a single hidden layer randomly assigns the weights of entry and analytically determines the weights the output by means the Moore-Penrose inverse, this algorithm tends to provide an extremely fast learning speed preserving the adjustment levels achieved by classifiers such as multilayer perception and support vector machine. However, the Moore-Penrose inverse loses precision when using data with additive noise in training. That is why in this paper a method to robustness of extreme learning machine to additive noise proposed. The method consists in computing the weights of the output layer using non-linear optimization algorithms without restrictions. Tests are performed with the gradient descent optimization algorithm and with the Levenberg-Marquardt algorithm. From the implementation it is observed that through the use of these algorithms, smaller errors are achieved than those obtained with the Moore-Penrose inverse.

## 1. Introduction

Artificial neural networks have been widely used to solve problems of medical diagnosis [1], voice recognition [2], image processing [3], facial recognition [4], between many other applications [1, 2, 5–8]. Despite this, they present clear inconveniences as the high time in the training and the convergence to local minimums. Many works have been developed with the finality of solve these problems and one of the most recent is the method known as extreme learning machine (ELM).

In general, the supervised learning algorithms of the neural networks (multilayer perceptron, convolutional neural networks) are based on non-linear optimization algorithms without restrictions, such as gradient descent or higher speed variants. Extreme learning machine was initially proposed for single-hidden layer feedforward neural networks in [9], and in its training algorithm the input weights and biases of the hidden layer can be chosen randomly, and the weights of the output layer is calculated using a least squares method based on the application of the generalized Moore-Penrose inverse. Unlike other learning algorithms based on gradient descent, the ELM algorithm does not require iterative techniques to adjust the weights and biases of the hidden layer during the training process, so it becomes in a simple learning method, and with an extremely fast learning speed.

Recently, it has been shown in [10] that ELM with weights and biases in the hidden layer assigned arbitrarily and with almost any nonzero activation function correspond to a universal



approximator of functions. This result makes that ELM networks compete with the well known multilayer perceptron and support vector machine.

Noisy data is a common factor in the processing of data that comes of the real-world. Sources of noise include physical measurement limitations, the use of stochastic simulation models, incomplete sampling of large spaces, and human-computer interaction. Many of the search methods in use nowadays for data with noise can be traced back to the least squares methodology approach, which consists in determine a curve that best fits to the data set [11–13]. A known fact is that in this type of methods the approximation error is degraded as the Noisy data increases. This paper proposes the use of nonlinear optimization algorithms without restrictions to improve the robustness of the estimate that provided the ELM network in presence of highly noisy data. The algorithms considered are the gradient descent and the Levenberg-Marquardt method, which allow to estimate the model parameters optimally in the sense of the mean square error. The results obtained show that the optimization methods improve the estimation of the parameters of the ELM network than the provided by the Pseudo-Inverse in a acceptable search time.

In section 2, presents the Moore-Penrose generalized inverse solution, the least squares solution for a general linear system, two optimization algorithms that allow to determine the minimum value of a function and a brief introduction to the ELM algorithm. Section 3 explains the materials and methods to address the proposed approaches. Performance evaluation is presented in section 4, while the conclusions are given in section 5.

## 2. Preliminaries

This section gives a brief introduction to the Moore-Penrose generalized inverse, the least squares solution, the gradient descent, the Levenberg-Marquardt method and the extreme learning machine algorithm.

### 2.1. Moore-Penrose generalized inverse

The solution of a linear system  $\mathbf{Ax} = \mathbf{y}$  can be calculated very simply by using the generalized inverse of Moore-Penrose [14], where the matrix  $\mathbf{A}$  can be singular and even rectangular.

A matrix  $\mathbf{G}$  of order  $n \times m$  is the Moore-Penrose generalized inverse of matrix  $\mathbf{A}$  of order  $m \times n$ , if it satisfies the conditions of Equation (1).

$$\mathbf{AGA} = \mathbf{A}, \quad \mathbf{GAG} = \mathbf{G}, \quad (\mathbf{AG})^T = \mathbf{AG} \quad \text{y} \quad (\mathbf{GA})^T = \mathbf{GA}. \quad (1)$$

The Moore-Penrose generalized inverse of matrix  $\mathbf{A}$  will be denoted by  $\mathbf{G} = \mathbf{A}^\dagger$ .

### 2.2. Least squares solution

We say that  $\mathbf{x}_0$  is a least squares solution for a linear system  $\mathbf{Ax} = \mathbf{y}$ , if for any  $\mathbf{y} \in \mathbb{R}^m$  satisfies the Equation (2).

$$\|\mathbf{x}_0\| \leq \|\mathbf{x}\|, \quad \forall \mathbf{x} \in \{\mathbf{x} : \|\mathbf{Ax} - \mathbf{y}\| \leq \|\mathbf{Az} - \mathbf{y}\|, \quad \forall \mathbf{z} \in \mathbb{R}^n\}, \quad (2)$$

that is, a solution  $\mathbf{x}_0$  is said to be a minimum norm least squares solution of a linear system  $\mathbf{Ax} = \mathbf{y}$  if it has the smallest norm among all the least squares solutions.

### 2.3. Optimization algorithms gradient descent and Levenberg-Marquardt

The gradient descent algorithm is described by the Equation (3), this is most intuitive technique to find minimum of a function [15,16]. The search process of the minimum is carried out through an iterative process following the opposite direction to the gradient of a function  $E$  with a convergence rate  $\alpha$ .

$$\mathbf{W}_{k+1} = \mathbf{W}_k - \alpha \nabla E. \quad (3)$$

It can be observe of [15–17] that the gradient descent method and the of Gauss-Newton are two complementary methods in the advantages they provide. In order to make sure that the approximated Hessian matrix  $\mathbf{H} \approx \mathbf{J}^T \mathbf{J}$  is invertible [16], Levenberg-Marquardt algorithm introduces another approximation to Hessian matrix by the Equation (4).

$$\mathbf{H} \approx \mathbf{J}^T \mathbf{J} + \mu \mathbf{I}, \quad (4)$$

where  $\mu$  is the combination coefficient,  $\mathbf{I}$  is the identity matrix and  $\mathbf{J}$  is the Jacobian matrix. Levenberg and Marquardt proposed a new algorithm based on these observations, whose update rule is a combination of the algorithms mentioned above, such as be shown in Equation (5).

$$\mathbf{W}_{k+1} = \mathbf{W}_k - (\mathbf{J}_k^T \mathbf{J}_k + \mu \mathbf{I})^{-1} \mathbf{J}_k. \quad (5)$$

The combination coefficient  $\mu$  in the previous equation is a parameter that be adjusted in each cycle of according a the evolution of the error. If  $\mu$  is very small, the  $\mathbf{J}_k^T \mathbf{J}_k$  matrix becomes an approximation to the Hessian and Gauss-Newton algorithm is used. If  $\mu \gg 1$ , the method becomes analogous to gradient descent.

#### 2.4. Extreme learning machine

For  $N$  arbitrary training samples  $\{(\mathbf{x}_i, \mathbf{t}_i)\}_{i=1}^N$ , where  $\mathbf{x}_i = [x_{i1}, x_{i2}, \dots, x_{in}]^T \in \mathbb{R}^n$  and  $\mathbf{t}_i = [t_{i1}, t_{i2}, \dots, t_{im}]^T \in \mathbb{R}^m$ . The ELM learning algorithm for single hidden layer neural networks with  $L$  hidden neurons and activation function  $g(x)$  are mathematically modeled by Equation (6).

$$\sum_{i=1}^L \beta_i g(\mathbf{w}_i \cdot \mathbf{x}_j + b_i) = \mathbf{o}_j, \quad j = 1, \dots, N, \quad (6)$$

where  $\mathbf{w}_i = [w_{i1}, w_{i2}, \dots, w_{in}]^T$  is the weight vector connecting the  $i$ th hidden neuron and the input neurons,  $\beta_i = [\beta_{i1}, \beta_{i2}, \dots, \beta_{im}]^T$  is the weight vector connecting the  $i$ th hidden neuron and the output neurons,  $b_i$  is the threshold of the  $i$ th hidden neuron and  $\mathbf{w}_i \cdot \mathbf{x}_j$  denotes the inner product of  $\mathbf{w}_i$  and  $\mathbf{x}_j$  [9].

The standard ELM algorithm with  $L$  hidden neurons and activation function  $g(x)$  can approximate these  $N$  samples with zero error means that  $\sum_{j=1}^N \|\mathbf{o}_j - \mathbf{t}_j\| = 0$ , that is, there exist  $\beta_i$ ,  $\mathbf{w}_i$  and  $b_i$  such as Equation (7).

$$\sum_{i=1}^L \beta_i g(\mathbf{w}_i \cdot \mathbf{x}_j + b_i) = \mathbf{t}_j, \quad j = 1, \dots, N. \quad (7)$$

The Equation (7) can be written in matrix form of Equation (8).

$$\mathbf{H} \beta = \mathbf{T}, \quad (8)$$

where the elements of Equation (8) are given by Equation (9).

$$\mathbf{H} = \begin{bmatrix} g(\mathbf{w}_1 \cdot \mathbf{x}_1 + b_1) & \cdots & g(\mathbf{w}_L \cdot \mathbf{x}_1 + b_L) \\ \vdots & \ddots & \vdots \\ g(\mathbf{w}_1 \cdot \mathbf{x}_N + b_1) & \cdots & g(\mathbf{w}_L \cdot \mathbf{x}_N + b_L) \end{bmatrix}, \quad \beta = \begin{bmatrix} \beta_1^T \\ \vdots \\ \beta_L^T \end{bmatrix} \quad \text{and} \quad \mathbf{T} = \begin{bmatrix} \mathbf{t}_1^T \\ \vdots \\ \mathbf{t}_N^T \end{bmatrix}. \quad (9)$$

As mentioned in [18, 19], the matrix  $\mathbf{H}$  defined in Equation (9) is called the hidden layer output matrix of the neural network; the  $i$ th column of  $\mathbf{H}$  is the  $i$ th hidden neuron's output vector with respect to inputs  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$ .

If the number of hidden neurons is equal to the number of training samples  $L = N$ , matrix  $\mathbf{H}$  is square and invertible, and SLFNs can approximate these training samples with zero error [18–20]. However, in most cases the number of hidden neurons is much less than the number of training samples  $L \ll N$ ,  $\mathbf{H}$  is a nonsquare matrix and there may not exist  $\mathbf{w}_i, b_i, \beta_i$  ( $i = 1, \dots, L$ ) such that  $\mathbf{H}\boldsymbol{\beta} = \mathbf{T}$ . Thus, instead one may need to find specific  $\hat{\mathbf{w}}_i, \hat{b}_i, \hat{\boldsymbol{\beta}}$  ( $i = 1, \dots, L$ ) such that Equation (10).

$$\|\mathbf{H}(\hat{\mathbf{w}}_1, \dots, \hat{\mathbf{w}}_L, \hat{b}_1, \dots, \hat{b}_L)\hat{\boldsymbol{\beta}} - \mathbf{T}\| = \min_{\mathbf{w}_i, b_i, \beta} \|\mathbf{H}(\mathbf{w}_1, \dots, \mathbf{w}_L, b_1, \dots, b_L)\boldsymbol{\beta} - \mathbf{T}\|. \quad (10)$$

The Equation (10) is equivalent to minimizing the cost function.

The unknown  $\mathbf{W}$  in Equations (3) and Equation (5) corresponds to the set of weights and biases  $\{\mathbf{w}_i, \beta_i, b_i\}$  of the ELM algorithm. In addition, the comparison of results will be carried out minimizing the objective function  $E$  defined in Equation (11) using the Levenberg-Marquardt algorithm and the gradient descent.

$$E = \frac{1}{N} \sum_{j=1}^N \left( \sum_{i=1}^L \beta_i g(\mathbf{w}_i \cdot \mathbf{x}_j + b_i) - \mathbf{t}_j \right)^2. \quad (11)$$

### 3. Materials and methods

The training and testing data for the ELM algorithm have been generated around of a  $n$  grade polynomial function through additive noise, that is Equation (12).

$$\text{data} = f(\mathbf{x}) + r \times \text{rand}(N, 1), \quad (12)$$

where  $f(\mathbf{x}) = a_0 \times \text{ones}(N, 1) + a_1 \mathbf{x} + a_2 \mathbf{x}^2 + \dots + a_n \mathbf{x}^n$ ,  $N$  is the set of training samples,  $\mathbf{x}$  is a vector of dimension  $N$  and  $r$  is the amplitude of the noise. The analysis was done taking different scenarios with respect to the degree of the polynomial and the number of hidden neurons in the neural network.

In this paper be proposes the gradient descent method (GD-ELM) and the of Levenberg-Marquardt (LM-ELM) for determine the output weights of the ELM algorithm through iterative processes defined in the Equation (3) and Equation (5). The procedure consists in estimate the model parameters optimally in the sense of the mean square error adjusting the noise  $r$  of the data obtained from Equation (12). In addition, we choose the sigmoid function (*i.e.*  $g(\mathbf{w}, \mathbf{x}, \mathbf{b}) = 1/(1 + \exp(-(\mathbf{w}\mathbf{x} + \mathbf{b})))$ ) as the activation function of GD-ELM, LM-ELM and ELM.

The input weights  $\mathbf{w}$  are randomly generated from the range of  $(-1, 1)$  and the hidden biases  $\mathbf{b}$  are generated randomly from the range of  $(0, 1)$  using a uniform distribution. For each regression problem, the average results over 50 trials are obtained for each algorithm. In this study, the algorithms have been implemented through the MATLAB R2014a environment.

### 4. Experimental results and analysis

In this section, a series of experiments are carried out to demonstrate the effectiveness of the proposed GD-ELM and LM-ELM approaches. To illustrate the good behavior of GD-ELM and LM-ELM, are compared with the basic ELM. The training problem that serves as example is to reconstruct a polynomial through the network. The results obtained are compared via the root mean square error (RMSE), such as is shown in the following Equation (13).

$$\text{RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2}, \quad (13)$$

where  $\hat{y}_i$  is the estimated value and  $y_i$  is the true value.

The experiments are conducted on three sets ( $S_1$ ,  $S_2$  and  $S_3$ ) specified in Table 1. The coefficients of the polynomials are random constants between -1 and 1. Training and testing samples were taken, where the input are distributed evenly in  $[0, 1]$ . Furthermore, the number of maximum iterations for gradient descent and the Levenberg-Marquardt algorithm are set in 15000 and 70, respectively.

**Table 1.** Specification of three datasets.

Polynomial degree (PD)	Training data	Testing data	Noise amplitude
Set 1 ( $S_1$ )	90	322	$r = 0.5$
Set 2 ( $S_2$ )	50	254	$r = 0.4$
Set 3 ( $S_3$ )	20	189	$r = 0.3$

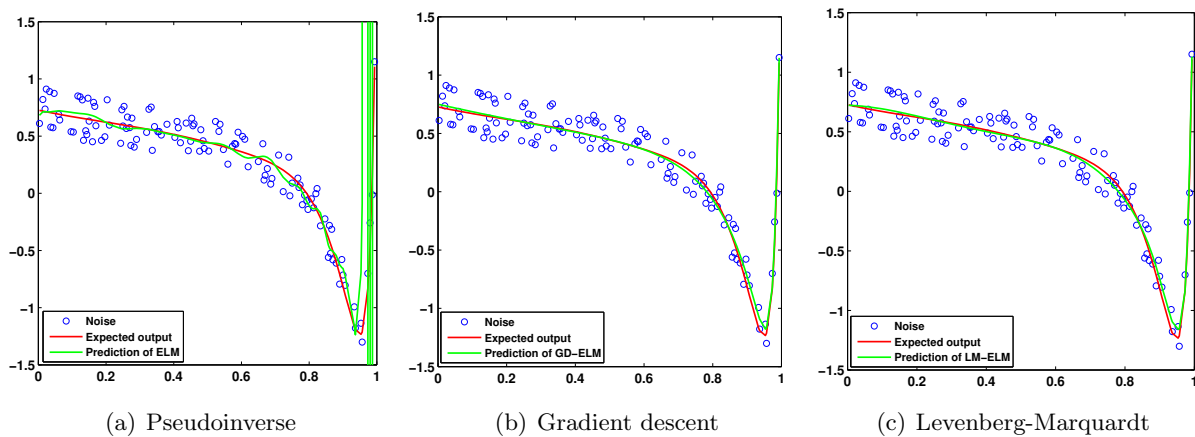
The testing error, the training error, the standard deviation (SD) and learning time are selected as the indicators for performance testing. The smaller testing RMSE denotes the better generalization performance of the algorithm and the smaller SD indicates the better stability of the algorithm. The comparison of the three algorithms is shown in Table 2.

**Table 2.** Performance comparison of the GD-ELM, LM-ELM and ELM.

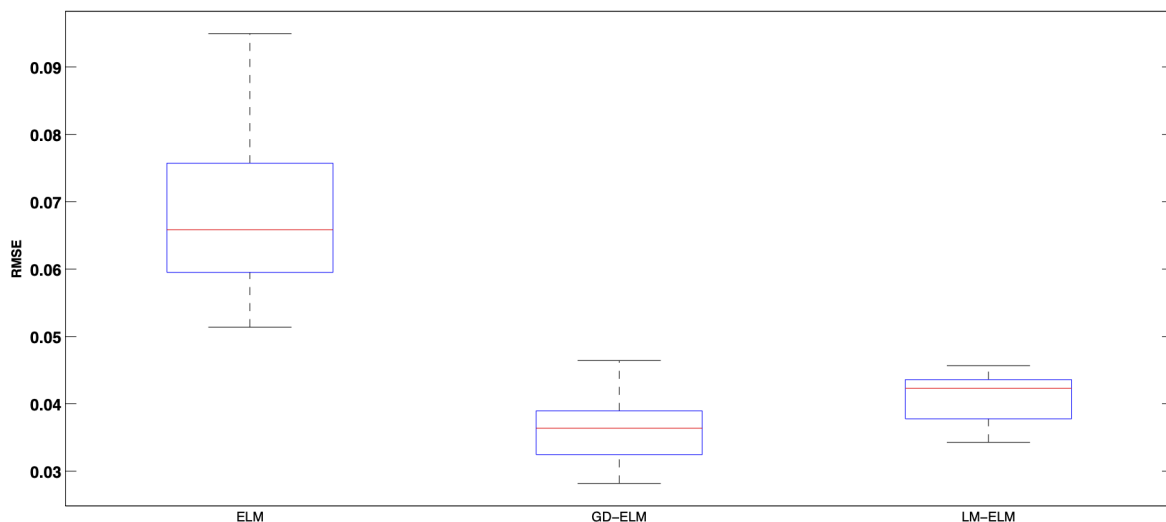
Datasets	Algorithms	Testing RMSE	Testing SD	Time	Hidden nodes
$S_1$	ELM	0.4857	0.2334	0.0092	30
	GD-ELM	0.0886	0.0192	3.2729	30
	LM-ELM	0.2096	0.0897	0.9012	30
$S_3$	ELM	0.2803	0.2803	0.0018	15
	GD-ELM	0.0843	0.0251	2.4390	15
	LM-ELM	0.1907	0.1128	0.0874	15
$S_3$	ELM	0.1096	0.0396	0.0030	20
	GD-ELM	0.0331	0.0023	2.6522	20
	LM-ELM	0.0366	0.0019	0.0585	20

From Table 2, we observe that the proposed GD-ELM and LM-ELM algorithms has the RMSE and the SD of testing smallest for the three data sets, which means that GD-ELM and LM-ELM can achieve better generalization performance and stability than ELM.

Figure 1 shows the behavior of the exact and approximate curves obtained by the optimization algorithms for the set  $S_1$ . Then, testing samples were used to see the stability of the proposed algorithms compared with the basic ELM. Figure 2 shows the test results obtained to the executing the code 50 times maintaining the same number of hidden neurons. It is observed that GD-ELM and LM-ELM always behave better compared with ELM.



**Figure 1.** Comparison of the actual curve and the output predicted by ELM, GD-ELM and LM-ELM for a noise of radius 0.2. The GD-ELM and LM-ELM algorithms present good approximations to the real model as shown in (b) and (c), however the ELM algorithm begins to oscillate as shown in (a) obtaining bad models.



**Figure 2.** Comparison of the RMSE of ELM, GD-ELM and LM-ELM considering 50 initializations of the weights.

## 5. Conclusion

In this study, two optimization algorithms were presented for the calculation of ELM parameters for regression problems. Experimental results show that the GD-ELM and LM-ELM algorithms can achieve better test results in compared with the pseudoinverse to determine the output weights of the neural network when working with for highly noisy data. Such as literature indicates, the LM-ELM algorithm is faster than the GD-ELM algorithm, which allows the solution of the problem in a reasonable time.

## References

- [1] Al-Shayea Q K 2011 Artificial neural networks in medical diagnosis *International Journal of Computer Science Issues* **8(2)** 150
- [2] Melin P, Urias J, Solano D, Soto M, Lopez M and Castillo O 2006 Voice Recognition with Neural Networks, Type-2 Fuzzy Logic and Genetic Algorithms *Engineering Letters* **13(2)** 108

- [3] Suresh S, Babu R V and Kim H J 2009 No-reference image quality assessment using modified extreme learning machine classifier *Applied Soft Computing* **9(2)** 541
- [4] Mohammed A A, Minhas R, Wu Q J and Sid-Ahmed M A 2011 Human face recognition based on multidimensional PCA and extreme learning machine *Pattern Recognition* **44(10-11)** 2588
- [5] Bai Z, Huang G B, Wang D, Wang H and Westover M B 2014 Sparse extreme learning machine for classification *IEEE Transactions on Cybernetics* **44(10)** 1858
- [6] Huang G, Song S, Gupta J N and Wu C 2014 Semi-supervised and unsupervised extreme learning machines *IEEE Transactions on Cybernetics* **44(12)** 2405
- [7] Huang G B, Zhou H, Ding X and Zhang R 2011 Extreme learning machine for regression and multiclass classification *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)* **42(2)** 513
- [8] Huang Z, Yu Y, Gu J and Liu H 2016 An efficient method for traffic sign recognition based on extreme learning machine *IEEE Transactions on Cybernetics* **47(4)** 920
- [9] Huang G B, Zhu Q Y and Siew C K 2004 Extreme learning machine: a new learning scheme of feedforward neural networks *Neural Networks* **2** 985
- [10] Huang G B, Chen L and Siew C K 2006 Universal approximation using incremental constructive feedforward networks with random hidden nodes *IEEE Trans. Neural Networks* **17(4)** 879
- [11] Arnold D V and Beyer H G 2003 A comparison of evolution strategies with other direct search methods in the presence of noise *Computational Optimization and Applications* **24(1)** 135
- [12] Cantú-Paz E 2004 Adaptive sampling for noisy problems *Genetic and Evolutionary Computation-GECCO 2004* **3102** 947
- [13] Ridout D and Judd K 2002 Convergence properties of gradient descent noise reduction *Physica D: Nonlinear Phenomena* **165(1-2)** 26
- [14] Courrieu P 2008 Fast computation of Moore-Penrose inverse matrices *Neural Information Processing - Letters and Reviews* **8(2)** 25
- [15] Ranganathan A 2004 The levenberg-marquardt algorithm *Tutorial on LM algorithm* **11(1)** 101
- [16] Yu H and Wilamowski B M 2011 Levenberg-marquardt training *Industrial electronics handbook* **5(12)** 1
- [17] Gavin H 2013 *The Levenberg-Marquardt method for nonlinear least squares curve-fitting problems* (Durham: Duke University) p 1
- [18] Huang G B and Babri H A 1998 Upper bounds on the number of hidden neurons in feedforward networks with arbitrary bounded nonlinear activation functions *IEEE Transactions on Neural Networks* **9(1)** 224
- [19] Huang G B 2003 Learning capability and storage capacity of two-hidden-layer feedforward networks *IEEE Transactions on Neural Networks* **14(2)** 274
- [20] Tamura S I and Tateishi M 1997 Capabilities of a four-layered feedforward neural network: four layers versus three *IEEE Transactions on Neural Networks* **8(2)** 251