

GROWING NEURAL NETWORKS USING NONCONVENTIONAL ACTIVATION FUNCTIONS

Yevgeniy Bodyanskiy, Iryna Pliss, Oleksandr Slipchenko

Abstract: *In the paper, an ontogenic artificial neural network (ANNs) is proposed. The network uses orthogonal activation functions that allow significant reducing of computational complexity. Another advantage is numerical stability, because the system of activation functions is linearly independent by definition. A learning procedure for proposed ANN with guaranteed convergence to the global minimum of error function in the parameter space is developed. An algorithm for structure network structure adaptation is proposed. The algorithm allows adding or deleting a node in real-time without retraining of the network. Simulation results confirm the efficiency of the proposed approach.*

Keywords: *ontogenic artificial neural network, orthogonal activation functions, time-series forecasting.*

ACM Classification Keywords: *I.2.6 Learning – Connectionism and neural nets*

Introduction

Artificial neural networks (ANNs) are widely applied to solving a variety of problems such as information processing, data analysis, system identification, control etc. under structural and parametric uncertainty [1, 2].

One of the most attractive properties of ANNs is the possibility to adapt their behavior to the changing characteristics of the modeled system. By adaptivity we understand not only the adjustment of parameters (synaptic weights), but also the possibility to adjust the architecture (the number of nodes). The goal of the present paper is the development of an algorithm for structural and synaptic adaptation of ANNs for nonlinear system modeling, capable of online operation, i.e. sequential information processing without re-training after structure modification.

The problem of optimization of neural network architecture has been studied for quite a long time. The algorithms that start their operation with simple architecture and gradually add new nodes during learning, are called 'constructive algorithms'. In contrast, destructive algorithms start their operation with an initially redundant network, and simplify it as learning proceeds. This process is called 'pruning'.

Radial basis function network (RBFN) is one of the most popular neural network architectures [3]. One of the first constructive algorithms for such networks was proposed by Platt and named 'resource allocation' [4]. By present time, a number of modifications of this procedure is known [5, 6]. One of the most known is the cascade-correlation architecture developed by Fahlman and Lebiere [7].

Among the destructive algorithms, the most popular are the 'optimal brain damage' [8] and 'optimal brain surgeon' [9]. In these methods, the significance of a node or a connection between nodes is determined by the change in error function that its deletion incurs. For this purpose, the matrix of second derivatives of the optimized function with respect to the tunable parameters is analyzed. Both procedures are quite complex computationally. Besides that, an essential disadvantage is the need for re-training after the deletion of non-significant nodes. This, in turn, makes the real-time operation of these algorithms impossible. Other algorithms such as [10] are heuristic and lack universality.

It should be noted that there is no universal and convenient algorithm, which could be used for the manipulation of the number of nodes and suitable for most problems and architectures. Many of the algorithms proposed so far lack theoretical justification as well as the predictability of the results of their application and the ability to operate in real time.

Network Architecture

Let's consider the network architecture, that implements the following nonlinear mapping

$$\hat{y}(k) = \hat{f}(x(k)) = \sum_{i=1}^n \sum_{j=1}^{h_i} w_{ji} \phi_{ji}(x_i(k)) \quad (4)$$

where $k = 1, 2, \dots$ – discrete time or ordinal number of sample in training set, w_{ji} – tunable synaptic weights, $\phi_{ji}(\bullet)$ – j -th activation function for i -th input variable, h_i – number of activation functions for appropriate input variable, $x_i(k)$ – value of i -th input signal at time moment k (or for k -th training sample).

This architecture contains $h = \sum_{i=1}^n h_i$ tunable parameters and it can be readily seen that the this number is

between the scatter-partitioned and grid-partitioned systems.

We propose the use of orthogonal polynomials of one variable for the activation functions. Particular system of functions can be chosen according to the specificity of the solved problem. If the input data are normalized on the hypercube $[-1, 1]^n$, the system of Legendre polynomials orthogonal on the interval $[-1, 1]$ with weight $\gamma(x) \equiv 1$ [17] can be used:

$$P_n(x) = 2^{-n} \sum_{m=0}^{[n/2]} (-1)^m \frac{(2n-2m)!}{m!(n-m)!(n-2m)!} x^{n-2m}, \quad (5)$$

where $[\bullet]$ is the integer part of a number.

System of Legendre polynomials is best suited for the case when we know exact interval of data changes before network construction. This is quite a common situation as well as an opposite one. For the latter case the following system of Hermite orthogonal polynomials can be used:

$$H_n(x) = n! \sum_{m=0}^{[n/2]} (-1)^m \frac{(2x)^{n-2m}}{m!(n-2m)!}. \quad (6)$$

This system is orthogonal on $(-\infty; +\infty)$ with weight function $h(x) = e^{-x^2}$ and gives us a possibility to decrease influence of the data lying far from the point of origin.

Normalized Hermite polynomials usually denoted by $\hat{H}_n(\bullet)$ (i.e. those with $\|\hat{H}_n(\bullet)\| = 1$) can be obtained from (6):

$$\hat{H}_n(x) = \sqrt{\frac{n!}{2^n \sqrt{\pi}}} \sum_{m=0}^{[n/2]} (-1)^m \frac{(2x)^{n-2m}}{m!(n-2m)!}. \quad (7)$$

Among other possible choices for activation functions we should mention Chebyshev [15, 16] and Hermite [18] functions as well as non-sinusoidal orthogonal systems proposed by Haar and Walsh.

Synaptic Adaptation

The sum of squared errors will be used as the learning criterion:

$$E(k) = \sum_{p=1}^k e^2(p) = \sum_{p=1}^k (y(p) - \sum_{i=1}^n \sum_{j=1}^{h_i} w_{ji} \phi_{ji}(x_i(p)))^2, \quad (8)$$

where k is the ordinal number of an element in the learning sequence or the discrete time when the data is processed in the order of its arrival, $y(p)$ – value of learning signal at time moment p (or for p -th training sample).

For the convenience of further notation, let us re-write the expression for the output of the neural network (4) in the form

$$\hat{y}(k+1) = \phi^T(k+1)W(k), \quad (9)$$

where $\phi(k) = (\phi_{11}(x(k)), \phi_{21}(x(k)), \dots, \phi_{h,n}(x(k)))^T$ is a $(h \times 1)$ vector of the values of the basis functions for the k -th element of the training set (or at the instant k for sequential processing), $W(k) = (w_{11}(k), w_{21}(k), \dots, w_{h,n}(k))^T$ is a $(h \times 1)$ vector of synaptic weights estimates at the iteration k .

Since the output of the proposed neural network linearly depends on the tuned parameters, we can use the least squares procedure to estimate them. For sequential processing, e.g. in the case of online learning, we can use the recursive least squares method:

$$\begin{cases} W(k+1) = W(k) + \frac{P(k)(y(k+1) - W^T(k)\phi(k+1))\phi(k+1)}{1 + \phi^T(k+1)P(k)\phi(k+1)}, \\ P(k+1) = P(k) - \frac{P(k)\phi(k+1)\phi^T(k+1)P(k)}{1 + \phi^T(k+1)P(k)\phi(k+1)}. \end{cases} \quad (10)$$

Because of the orthogonality of the basis functions, the matrix $P(k)$ will tend to diagonal form as $k \rightarrow \infty$. If the activation functions are orthonormal, $P(k)$ will tend to the unity matrix. Due to this property, the learning procedure will retain numerical stability with the increase of the number of samples in the training sequence.

Structure Adaptation

Let's consider sequential learning that minimizes (8) and leads to the estimate

$$W_h(k) = R_h^{-1}(k)F_h(k), \quad (11)$$

$$R_h^{-1}(k) = R_h^{-1}(k-1) - \frac{R_h^{-1}(k-1)\phi(k)\phi(k)^T R_h^{-1}(k-1)}{1 + \phi(k)^T R_h^{-1}(k-1)\phi(k)}, \quad (12)$$

$$F_h(k) = F_h(k-1) + \phi(k)y(k). \quad (13)$$

The use of the recursive least squares (RLS) method and its modifications allows to obtain an accurate and well-interpretable measure of significance of each function in the mapping (4). This mapping can be considered as an expansion of an unknown reconstructed function in the basis $\{\phi_{ji}(\cdot)\}$. Obviously, if the absolute value of any of the coefficients in this expansion is small, then the corresponding function can be excluded from the basis without significant loss of accuracy. The remaining synaptic weights does not need to be retrained if the weight of the excluded node is close to zero. Otherwise, the network should be retrained.

Assume that a vector of synaptic weights $W_h(k)$ of a network comprising h nodes was obtained at the instant k using the formula (11), where the index h determines the number of basis functions (the dimension of $\phi(k)$). Also assume that the absolute value of the considered parameter $w_h(k)$ is small, and we want to exclude corresponding unit function from the expansion (4). The assumption about the insignificance of the activation h is not restrictive, because we always can re-number the basis functions. This will result only in the rearrangement of the rows and columns in the matrix $R_h(k)$ and in the change of ordering of the elements of the vector $F_h(k)$. However, the rearrangement of columns and/or rows of a matrix does not influence the subsequent matrix operations.

Taking into account the fact that the matrix $R_h(k)$ is symmetric, we obtain:

$$W_h(k) = R_h^{-1}(k)F_h(k) = \begin{pmatrix} R_{h-1}(k) & \beta_{h-1}(k) \\ \beta_{h-1}^T(k) & r_{hh}(k) \end{pmatrix}^{-1} \begin{pmatrix} F_{h-1}(k) \\ f_h(k) \end{pmatrix}, \quad (14)$$

where $r_{ij}(k)$ is the element of the i -th row and j -th column of the matrix $R_h(k)$,

$\beta_{h-1}(k) = (r_{1h}(k), \dots, r_{hh-1}(k))^T = (r_{h1}(k), \dots, r_{hh-1}(k))^T$, $f_i(k)$ is the i -th element of vector $F_h(k)$.

After simple transformations of (14) we obtain the expression

$$W_h(k) = \begin{pmatrix} W_{h-1}(k) - R_{h-1}^{-1}(k)\beta_{h-1}(k)w_h(k) \\ w_h(k) \end{pmatrix} \tag{15}$$

that enables us to exclude the function from (4) and obtain the corrected estimates of the remaining parameters of the ANN. For this operation, we use only the information accumulated in the matrix $R_h(k)$ and vector $F_h(k)$.

Using the same technique as above, it is possible to write a procedure that can be used to add a new function to the existing basis. Direct application of the Frobenius formula [12] leads to the algorithm

$$W_{h+1}(k) = R_{h+1}^{-1}(k)F_{h+1}(k) = \begin{pmatrix} R_h(k) & \beta_h(k) \\ \beta_h^T(k) & r_{h+1,h+1}(k) \end{pmatrix}^{-1} \begin{pmatrix} F_h(k) \\ f_{h+1}(k) \end{pmatrix} = \begin{pmatrix} W_h(k) + R_h^{-1}(k)\beta_h(k) \frac{\beta_h^T(k)W_h(k) - f_{h+1}(k)}{r_{h+1,h+1}(k) - \beta_h^T(k)R_h^{-1}(k)\beta_h(k)} \\ \frac{-\beta_h^T(k)W_h(k) + f_{h+1}(k)}{r_{h+1,h+1}(k) - \beta_h^T(k)R_h^{-1}(k)\beta_h(k)} \end{pmatrix} \tag{16}$$

where $\beta_h(k) = (r_{1h+1}(k), \dots, r_{hh+1}(k))^T = (r_{h+11}(k), \dots, r_{hh+1}(k))^T$.

Thus, with the help of equation (16) we can add a new function (neuron) to the model (4), and exclude an existing function using the formula (15) without retraining remaining weights. In order to perform these operations in real time, it is necessary to accumulate the information about a larger number of basis functions than currently being used. E.g., we can initially introduce a redundant number of basis functions H and accumulate information in the matrix $R_H(k)$ and vector $F_H(k)$ as new data arrive, with only $h < H$ basis functions being used for the description of the unknown mapping. The complexity of the model can be either reduced or increased as required.

The analysis of equations (11), (15), and (16) shows that the efficiency of the proposed learning algorithm is directly related to the condition number of the matrix $R_h(k)$. This matrix will be non-singular if the functions $\{\varphi_i(\cdot)\}_{i=1}^h$ used in the expansion (4) are linear-independent. The best situation is when the function system $\{\varphi_i(\cdot)\}_{i=1}^h$ is orthogonal. In this case, the matrix $R_h(k)$ becomes diagonal, the formulas (11), (15), and (16) being greatly simplified because

$$diag(a_1, \dots, a_n)^{-1} = diag\left(\frac{1}{a_1}, \dots, \frac{1}{a_n}\right), \tag{17}$$

where $diag(a_1, \dots, a_n)$ is an $(n \times n)$ matrix with non-zero elements a_1, \dots, a_n only on the main diagonal.

Simulation Results

We have applied the proposed ontogenic network with orthogonal activation functions to online identification of a rat's (*Ratus Norvegicus Vistar*) brain activity during sleeping phase.

The signal was measured with frequency of 64 Hz. We took a fragment of signal containing 3200 points (50 second of measuring), that was typical for sleeping phase of rat's life activity. Two neural networks of type (4) were trained in real-time. Each network had 10 inputs – delayed signal values ($y(k), y(k-1), \dots, y(k-9)$) and was trained to output one-step ahead value of the process – $y(k+1)$. First network utilized synaptic adaptation algorithm (11) while second one also involved the structure adaptation technique (15), (16). Initially both ANNs had 5 activation functions per input, the one with synaptic adaptation only retained all 50 tunable

parameters during it's work while ANN with structure adaptation mechanism had only 25 fired functions (the most significant ones chosen in real-time). For the results comparing purpose we also trained multilayer perceptron (further referred as MLP) with the same structure of inputs and training signal, having 5 units in the 1st and 4 in the 2nd hidden layers (that totals to 74 tunable parameters). As MLP is not capable of real-time data processing, all samples are used as training set and test criteria are calculated on the same data points. MLP was trained during 250 epochs with Levenberg-Marquardt algorithm. Our research showed that this is enough to achieve precision comparable to proposed ontogenic neural network with orthogonal activation functions.

For visual presentation of processed data see Fig. 1 which shows the results of identification using proposed neural network together with original time series.

Results of identification can be found in table 1. We used some different measures of identification quality. First, we analyse normalized root mean squared error, which is closely related to the learning criterion. Two other criteria used: "Wegstrecke" [19] characterizes the quality of the model for prediction/identification (+1 means perfect one), "Trefferquote" [20] is percent value of correctly predicted direction changes.

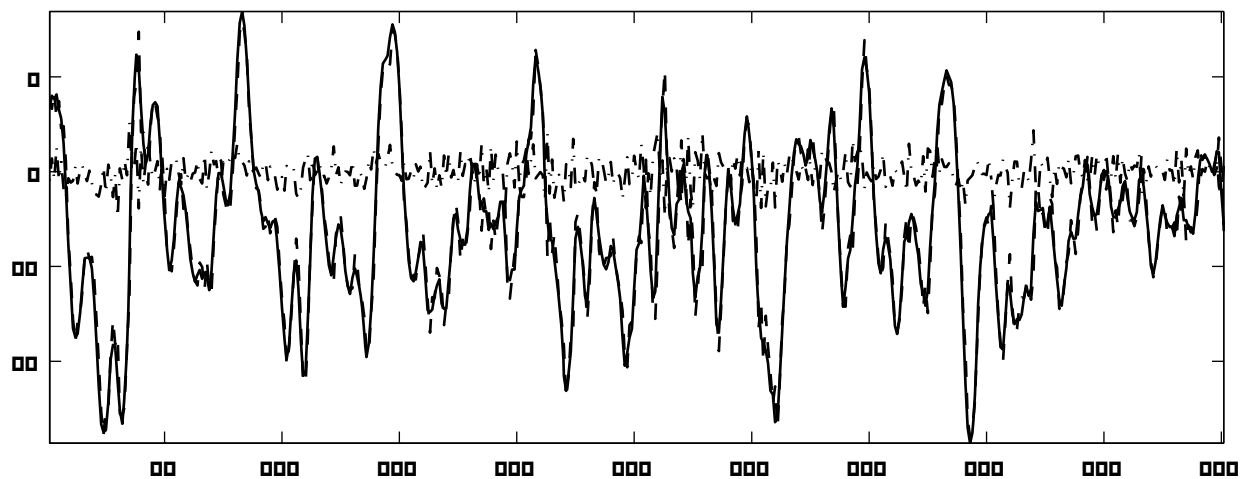


Figure 1. Identification of a rat's brain activity during sleeping phase using proposed neural network with orthogonal activation functions – brain activity signal (*solid line*), network output (*dashed line*), and identification error (*dash-dot line*)

Table 1 – Identification results for different architectures

<i>Description</i>	<i>NRMSE</i>	<i>Trefferquote</i>	<i>Wegstrecke</i>
<i>OrthoNN, real-time processing</i>	<i>0.1852</i>	<i>82.2847</i>	<i>0.85312</i>
<i>OrthoNN, real-time processing, variable number of nodes</i>	<i>0.2175</i>	<i>77.6357</i>	<i>0.74625</i>
<i>MLP, offline learning (250 epochs), error on the training set</i>	<i>0.1685</i>	<i>83.9533</i>	<i>0.87192</i>

We can see that utilizing structure adaptation technique leads to somewhat worth results. This is the tradeoff for having less tunable parameters and possibility to process non-stationary signals.

Adaptation of neural network in real time benefits us in a number of ways. First, as noted earlier, it can reduce computational complexity. Second and perhaps more important benefit is in using adapting neuromodel as a basis for some higher level system of data processing (e.g. time-series classification, diagnostics system etc.).

After obtaining promising results of online identification we used proposed neural network architecture to monitor rat's state in real time. Second level of monitoring system was built with the help of expert which initially provided recorded activity of rat's brain together with animal state for each moment of time. We processed the data and analyzed neural network's set of states. The analysis showed that states are dividable in a space of synaptic weights. A slightly modified Bayes estimator for the state of observed object was synthesized and trained. Simulation showed that developed automated monitoring system is capable of online data processing and gives correct state in 94,5% of cases. The response of the systems in form of object's state was later verified by the expert and found reliable enough to be used for data preprocessing in day to day activity.

Conclusion

A new computationally efficient neural network with orthogonal activation functions was proposed. It has a simple and compact architecture not affected by the curse of dimensionality, and provides high precision of nonlinear dynamic system identification. An apparent advantage is much easier implementation and lower computational load as compared to the conventional neural network architectures.

The approach presented in the paper can be used for nonlinear system modeling, control, and time series prediction. An interesting direction of further work is the use of the network with orthogonal activation functions as a part of hybrid multilayer architecture. Another possible application of proposed ontogenic neural network is its use as a basis for diagnostic systems.

References

1. Handbook of Neural Computation. IOP Publishing and Oxford University Press, 1997.
2. Nelles O. Nonlinear System Identification. Berlin, Springer, 2001.
3. Poggio T. and Girosi F. A Theory of Networks for Approximation and Learning. A.I. Memo No. 1140, Artificial Intelligence Laboratory, Massachusetts Institute of Technology, 1989.
4. Platt J. A resource allocating network for function interpolation. *Neural Computation*, 3, 1991, p. 213-225.
5. Nag A. and Ghosh J. Flexible resource allocating network for noisy data. In: Proc. SPIE Conf. on Applications and Science of Computational Intelligence, SPIE Proc. Vol. 3390, Orlando, FL., April 1998, p. 551-559.
6. Yingwei L., Sundararajan N. and Saratchandran P. Performance evaluation of a sequential minimal radial basis function (RBF) neural network learning algorithm. *IEEE Trans. on Neural Networks*, 9, 1998, p. 308-318.
7. Fahlman S. E. and Lebiere C. The cascade-correlation learning architecture. Technical Report CMU-CS-90-100, School of Computer Science, Carnegie Mellon University, Pittsburgh, PA, 1990.
8. Cun Y. L., Denker J. S., Solla S. A. Optimal Brain Damage. *Advances in Neural Information Processing Systems*, 2, 1990, p. 598-605.
9. Hassibi B. and Stork D. G. Second-order derivatives for network pruning: Optimal brain surgeon. In: *Advances in Neural Information Processing Systems*, Hanson et al. (Eds), 1993, p. 164-171.
10. Prechelt L. Connection pruning with static and adaptive pruning schedules. *Neurocomputing*, 16, 1997, p. 49-61.
11. Takagi T. and Sugeno M. Fuzzy identification of systems and its application to modeling and control. *IEEE Trans. on System, Man and Cybernetics*. 15, 1985, p. 116-132.
12. Gantmacher F. R. *The Theory of Matrices*. Chelsea Publ. Comp., New York, 1977
13. Narendra K. S. and Parthasarathy K. Identification and control of dynamic systems using neural networks. *IEEE Trans. on Neural Networks*, 1, 1990, p. 4-26.
14. Scott I. and Mulgrew B. "Orthonormal function neural network for nonlinear system modeling". In: *Proceedings of the International Conference on Neural Networks (ICNN-96)*, June, 1996.
15. Patra J.C. and Kot A.C. Nonlinear dynamic system identification using Chebyshev functional link artificial neural network. *IEEE Trans. on System, Man and Cybernetics – Part B*, 32, 2002, p. 505-511.
16. Bodyanskiy Ye.V., Kolodyazhniy V.V., and Slipchenko O.M. "Forecasting neural network with orthogonal activation functions" In: Proc. of 1st Int. conf. "Intelligent decision-making systems and information technologies", Chernivtsi, Ukraine, 2004, p. 57. (in Russian)
17. Bateman, H., Erdelyi, A.: *Higher Transcendental Functions*. Vol.2. McGraw-Hill (1953)

18. Liying M., Khorasani K. Constructive Feedforward Neural Network Using Hermite Polinomial Activation Functions. IEEE Trans. On Neural Networks, 16, No. 4, 2005, p.821–833.
19. Baumann M. Nutzung neuronale Netze zur Prognose von Aktionkursen. – Report Nr. 2/96, TU Ilmenau, 1996.
20. Fueser K. Neuronale Neteze in der Finanzwirtschaft. – Wiesbaden: Gabler, 1995.

Authors' Information

Yevgeniy Bodyanskiy - Dr. Sc., Prof., Head of Control Systems Research Laboratory, Kharkiv National University of Radio Electronics, Lenin Av., 14, Kharkiv, 61166, Ukraine, e-mail: bodya@kture.kharkov.ua

Iryna Pliss - Ph.D., Senior research scientist, Control Systems Research Laboratory, Kharkiv National University of Radio Electronics, Lenin Av., 14, Kharkiv, 61166, Ukraine, e-mail: pliss@kture.kharkov.ua

Oleksandr Slipchenko - Ph.D., Senior research scientist, Control Systems Research Laboratory, Kharkiv National University of Radio Electronics, Lenin Av., 14, Kharkiv, 61166, Ukraine, e-mail: aslipchenko@kture.kharkov.ua

DOUBLE-WAVELET NEURON BASED ON ANALYTICAL ACTIVATION FUNCTIONS

Yevgeniy Bodyanskiy, Nataliya Lamonova, Olena Vynokurova

Abstract: *In this paper a new double-wavelet neuron architecture obtained by modification of standard wavelet neuron, and its learning algorithm are proposed. The offered architecture allows to improve the approximation properties of wavelet neuron. Double-wavelet neuron and its learning algorithm are examined for forecasting non-stationary chaotic time series.*

Keywords: *wavelet, double-wavelet neuron, recurrent learning algorithm, forecasting, emulation, analytical activation function.*

ACM Classification Keywords: *I.2.6 Learning – Connectionism and neural nets*

Introduction

Recently, in the analysis tasks and the non-stationary series processing under the uncertainty conditions computational intelligence techniques particularly hybrid neural networks are widely used. The most important tasks related to signal processing are forecasting and emulation of dynamic non-stationary states of systems in the future.

For solving such kind of forecasting problems a variety of neural network architectures including hybrid architectures are used. However they are either bulky because of their architecture (for instance multilayer perceptron) or poorly adjusted to learning process in real time. In most cases the activation functions for these neural networks are sigmoidal functions, splines, polynomials and radial basis functions.

In addition wavelet theory is widespread [1-3] and allows to recognize the local characteristics of the non-stationary signals with high accuracy. At the confluence of the two approaches, hybrid neural networks and wavelet theory, have evolved the so-called wavelet neural networks [4-18] that have good approximating properties and sensitivity to the characteristics changes of the analyzed processes.

Previous studies have proposed and described [19-21] attractive features of wavelet neuron such as technical realization, ensured accuracy and learning simplicity. At the same time the wavelet functions are incarnated either at the level of synaptic weights or the neuron output, and as a learning algorithm the gradient learning algorithm with constant step is used. For the improvement of approximation abilities and the acceleration of the learning