

# Coordination algorithm in hierarchical structure of the learning process of Artificial Neural Network

Stanislaw Placzek and Bijaya Adhikari

stanislaw.placzek@wp.pl, bijaya.adhikari1991@gmail.com

Vistula University, Warsaw

## Abstract

While analyzing Artificial Neural Network structures, one usually finds that the first parameter is the number of the ANN layers . Hierarchical structure is an accepted default way to define ANN structure . This structure can be described using different methods, mathematical tools, software and/or hardware realization. In this article, we are proposing ANN decomposition into hidden and output sub networks. To build this kind of learning algorithm, information is exchanged between the first sub networks level and the second coordinator level in every iteration .Learning coefficients are tuned in every iteration. The main coordination task is to choose the coordination parameters in order to minimize both the global target function and all local target functions. In each iteration their values should decrease in asymptotic way to achieve the minimum. In article learning algorithms using forecasting of sub networks connectedness is studied .

## 1 Introduction

Many ANN structures are in practice. The most popular among them is the one with Forward Connections having complete or semi-complete set of weight coefficients. For special needs, ANNs with Forward Cross Connections and Back Connections are used. The full structure of ANN. is depicted on Fig.1. To describe the structure ,independent of the ANN complexity, partition on layers is used: the input layer, one or more hidden layers, and the output layer. Input layer connects ANN with external world ( environment) and performs initial processing , calibration or filtering of input data. The hidden layers are used for main data processing.

In most common structures, hidden layers include more neurons than input layer and they use non-linear activation function. The output layer which sums all signals from hidden layers uses two types of activation functions: linear activation function for classification tasks and non-linear sigmoid or tanth activation functions for approximation tasks . In this paper, to avoid confusion regarding the number of layers, only the hidden layers and the output layer are included. The concept of layers in ANN structures reflects the silent assumption

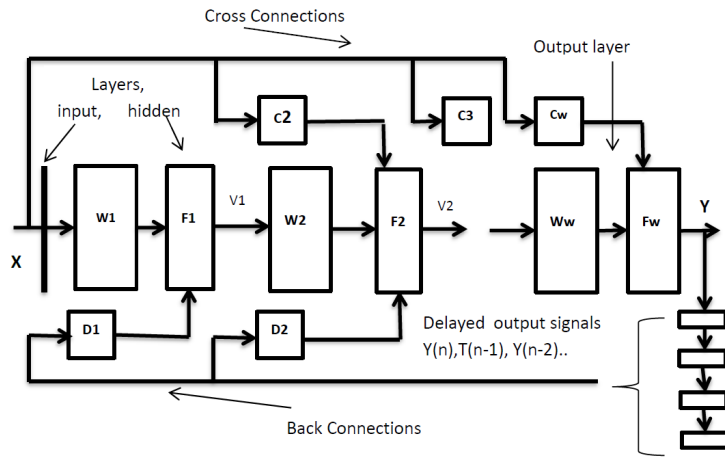


Figure 1: Scheme of the full ANN structure.

that ANN structures are hierarchical. Taking this into account as very important feature of ANN, to describe the network characteristic, a couple of the conceptions can be used.

### 1.1 Abstract description of ANN

To analyze ANN structure, verbal description is used so as to help everybody understand how ANN is built. For more detailed analysis, mathematical description using algebra and/or differential equations is required. Based on these descriptions, ANNs are then implemented by a computer program or an electronic device. So, to achieve complete description of ANN, concepts and models from different fields of science and technology have to be used.

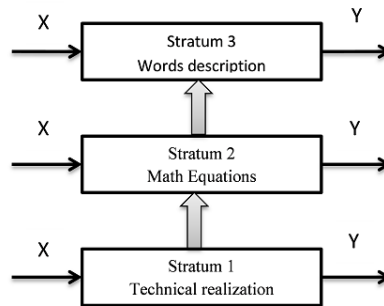


Figure 2: ANN stratification description.

Every model uses its own set of variables and terminology in different abstract level. To describe and understand how a particular ANN is working, some hierarchical set of abstract concepts are used. To separate these concepts from the layer description, a new name is used [15] – delamination of ANN into abstract strata.

## 1.2 Calculation complexity or decision taking.

For multi-layered ANN a lot of hidden layers and output layer can be sectioned off. Every layer has own output vector that is an input vector of the next layer,  $v_i$   $i = 1, 2, \dots, n$ . Both hidden layers and output layer can be described as sub-networks. “ $n$ ” defines the total number of sub-networks. ANN logic decomposition depends on layers separated by establishment of extra output vectors  $v_i$   $i = 1, 2, \dots, n$ . Now the network consists of the set of sub-networks, for each of which local target function is defined by  $\Phi = (\Phi_1, \Phi_2, \dots, \Phi_n)$ .

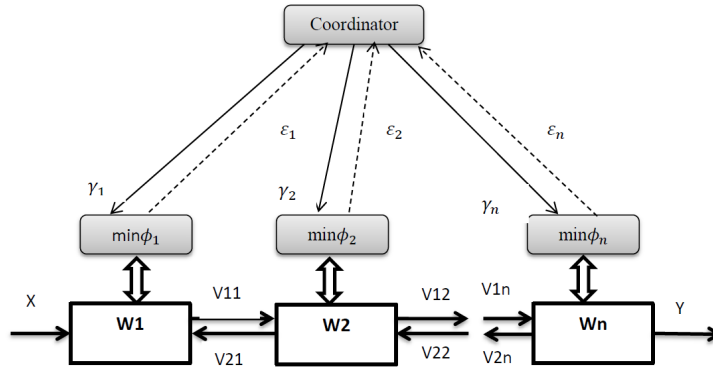


Figure 3: ANN decomposition.

Similar to ANN structure decomposition, learning algorithm using error back propagation can be decomposed too. (Fig.3.). We can sort out:

- - The first level task in which the minimum of the local target functions  $\Phi_i$   $i = 1, 2, \dots, n$  is searched.
- - The second level task which has to coordinate the all first level tasks.

In a learning algorithm constructed this way, there is a set of optimization tasks on the first level. These tasks are searching for the minimum value of target function  $\Phi$ . Unfortunately these are non-linear tasks without constrains. In practice, standard procedures to solve these problems exist. But in two level learning algorithm structure, coordinator is not responsible for solving the global task. Coordinator is obliged to calculate the value of coordination parameters  $\gamma = (\gamma_1, \gamma_2, \dots, \gamma_n)$  for every task on the first level. The first level, searching for the solution of all tasks have to use the coordination parameters value. It is an iterative process. Coordinator in every iteration cycle receives new values of feedback parameters  $\epsilon = (\epsilon_1, \epsilon_2 \dots \epsilon_n)$  from the first level tasks. Using this information coordinator has to make new decisions – calculate the new coordination parameters value. These procedures could be relatively complicated and in the most situations they happen to be non – gradient procedures. In the hierarchical learning algorithm, target functions can be defined as:

- Global target function  $\Psi$ ,
- Set of local target functions for every sub network  $\Phi_i$  where  $i = 1, 2, \dots, n$ ,
- Coordinator target function  $\Omega$ .

According to [15][2], solution of the primary task depends on the minimum global target function  $\Psi$ . The first level tasks should be built in a way that when all the first level tasks are solved, the final solution must be achieved – the minimum of the global target function. This kind of stratified structure is known as level hierarchy [15].

To summarize we conclude:

- Complexity of the problem increases from the first level to the second. Coordinator needs more time to solve its own tasks.
- Coordination tasks could be non – parametric procedures. To study dynamics of changing target functions value  $\Phi$ , coordinator should have the ability to change ( or changing) learning parameters in the first level tasks. As stressed above all the first level task are non – liner and have to be solved using iteration procedures.
- For different tasks, characteristic of ANN learning processes could be different . Coordinator studying feedback information from the first level tasks should have the ability to change all parameters in the both coordinator and the first level procedures.

## 2 Decomposition and coordination of ANN learning algorithm

The two layered ANN with one hidden layer and output layer using full internal forward connections does not have Cross Forward and Back Connections. This kind of networks can be used for both approximation and classification tasks. According to concept introduced above this ANN can be described by using two strata.

### 2.1 Verbal description of Structure. Stratum 2.

ANN with full forward connection contains one hidden layer. In this layer connections between input vector  $X$  and output vector  $V1$  are represented by matrix  $W1$ . All matrix coefficients are defined . Connections in the output layer are defined by matrix  $W2$ . Matrices connect input vector  $V1$  and output vector  $Y$ . In this matrix all weight coefficient are defined , too. Number of input neurons is defined by vector  $X$  which has dimensionality of  $N_0$ . In the same way number of neurons in the output layer is defined by vector  $Y$  which has dimensionality of  $N_2$ . Number of neurons in the hidden layers , $N_1$ , depends on complexity of problem. Usually  $N_1 > N_0$ , so data is not compressed in the first layer. Based on the description introduced above, the ANN can be set off as hierarchical level structure (Fig.4).In the first level, two local target functions,  $\Phi1$  for the first sub-network and  $\Phi2$  for the second sub-network, are defined. On the second level, coordinator is established. Its main goal is to coordinate all the first level tasks and to achieve the minimum of the global target function  $\Psi$ . For coordinator two functions  $G$  and  $H$  are defined which transforms coordination signals ( $V21, V12$ ) and feedback signals ( $V1, V2$ ). At the same time, coordinator should have the ability to change value of learning coefficients  $\alpha_1$  and  $\gamma_2$  by using transformation functions  $h_1(\Phi1, \Phi2)$  and  $h_2(\Phi1, \Phi2)$  (Fig.4.).

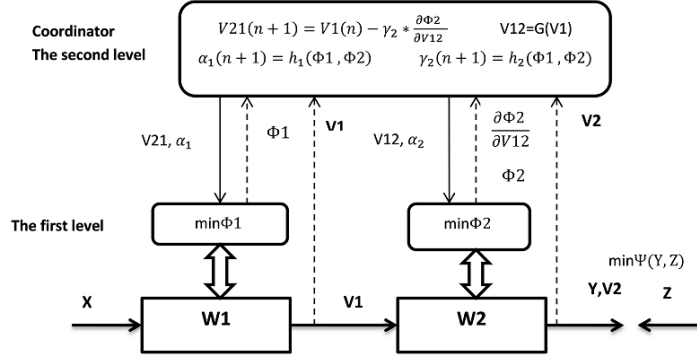


Figure 4: Coordination scheme

## 2.2 Mathematical description structure of ANN. Stratum 1

In the decomposed ANN structure we can defined the next target functions:

- Global target function  $\Psi$ . For all epoch:

$$\Psi(W1, W2, X, Y) = \sum_{k=1}^{N_2} \sum_{p=1}^{N_p} \Psi_k^p = \frac{1}{2} \sum_{k=1}^{N_2} \sum_{p=1}^{N_p} (y_k^p - z_k^p)^2 \quad (1)$$

Where:

$Y[1 : N_2, 1 : N_p]$  - ANN output value,

$Z[1 : N_2, 1 : N_p]$  - teaching data,

$N_2$  - number of output neurons,

$N_p$  - dimensionality of training set.

$\Psi_k^p$  - global target function for "k" output of the second sub-network and p<sup>th</sup> element of training set.

- Local target function s. For all epoch:

$$\epsilon_1 = \Phi1(W1, X, V21) = \sum_{i=1}^{N_1} \sum_{p=1}^{N_p} \phi1_i^p = \frac{1}{2} \sum_{i=1}^{N_1} \sum_{p=1}^{N_p} (v1_i^p - v21_i^p)^2 \quad (2)$$

where:

$V1, V2[1 : N_1, 1 : N_p]$  - coordination matrixes,

$N_1$  - number of hidden neurons,

$N_0$  - number of input neurons.

$\Phi1_i^p$  - local target function for "i" output of the first sub-network and p<sup>th</sup> element of training set.

$$v1_i^p = f(u_i^p) \quad (3)$$

$$u_i^p = \sum_{j=0}^{N_0} W1_{ij} \cdot x_j^p \quad (4)$$

$$\epsilon_2 = \Phi_2(W_2, V_2, Z_2) = \sum_{i=1}^{N_2} \sum_{p=1}^{N_p} \Phi_2^p = \frac{1}{2} \sum_{i=1}^{N_2} \sum_{p=1}^{N_p} (v_2^p - z_2^p)^2 \quad (5)$$

$$v_2^p = f(e_2^p) \quad (6)$$

$$e_2^p = \sum_{i=0}^{N_1} W_{2ki} \cdot v_{1i}^p \quad (7)$$

Where:

f - sigmoid function,

$i = 1, 2, \dots, N_1$   $k = 1, 2, \dots, N_2$ ,

$\Phi_2^p$ - local target function for "k" output of the second sub-network and p<sup>th</sup> element of training set.

On the first level two minimization task  $\Phi_1$  and  $\Phi_2$  have to be solved. These target functions have additive structures. Both could be divided into  $N_1$  and  $N_2$  sub-tasks respectively. This can be used to build programming procedures using appropriate programming language. So, we can formulate  $N_1$  sub-tasks

$$\min \Phi_1 = \min \sum_{i=1}^{N_1} \phi_{1i} = \sum_{i=1}^{N_1} \sum_{p=1}^{N_p} (f[\sum_{j=0}^{N_0} W_{1ij} \cdot x_j^p] - v_{1i}^p)^2 \quad (8)$$

$$\frac{\partial \Phi_1}{\partial W_{1ij}} = \sum_{p=1}^{N_p} (v_{1i}^p - v_{1i}^p) \cdot \frac{\partial f}{\partial u_i} \cdot x_j^p \quad (9)$$

$$W_{1ij}(n+1) = W_{1ij} - \alpha_1 \cdot \frac{\partial \Phi_1}{\partial W_{1ij}} \quad (10)$$

For  $i = 1, 2, \dots, N_1$   $j = 0, 1, 2, \dots, N_0$

$\Phi_{1i}$ -local target function for "i" output of the first sub-network and for whole training set.

In the same way can be formulated  $N_2$  sub-tasks

$$\min \Phi_2 = \min \sum_{k=1}^{N_2} \phi_{2k} = \sum_{k=1}^{N_2} \sum_{p=1}^{N_p} (f[\sum_{i=0}^{N_1} W_{2ki} \cdot v_{1i}^p] - z_k^p)^2 \quad (11)$$

$$\frac{\partial \Phi_2}{\partial W_{2ki}} = \sum_{p=1}^{N_p} (v_2^p - z_k^p) \cdot \frac{\partial f}{\partial e_k} \cdot v_{1i}^p \quad (12)$$

$$\frac{\partial \Phi_2}{\partial v_{1i}^p} = \sum_{k=1}^{N_2} \sum_{p=1}^{N_p} (v_2^p - z_k^p) \cdot \frac{\partial f}{\partial e_k} \cdot W_{2ki} \quad (13)$$

$$W_{2ki}(n+1) = W_{2ki}(n) - \alpha_2 \cdot \frac{\partial \Phi_2}{\partial W_{2ki}} \quad (14)$$

For  $k = 1, 2, \dots, N_2$   $i = 0, 1, 2, \dots, N_1 - 1$

$\Phi_{2k}$ -local target function for "k" output of the second sub-network and for whole training set.

- Coordinator target function

$$\omega = \Omega(\Phi1, \Phi2, V1, V2) \quad (15)$$

The first level tasks calculate control parameters and send them to coordinator. Additionally, in every iteration, coordinator analyzes the local target functions  $\phi_{1_i}(n)$  and  $\phi_{2_k}(n)$ . This information is necessary to calculate the new vector value  $V21$ . At the same time coordinator should have the ability to interfere in learning process by selecting new value of learning parameters  $\alpha_1, \alpha_2, \gamma_2$ . Coordinator can calculate the value of target function by itself using data sent to it by the first level. We should stress that values of the target function change dramatically during the learning process. We observed that values of  $\phi_{1_i}(n)$  and  $\phi_{2_k}(n)$  changed significantly during several hundred iterations. At the same, during learning process the values of the target functions can increase to a big value and then decrease drastically. This process explains that ANN, at the beginning of the learning process, has to attune the weight coefficients of the  $W1$  matrix. In the next step both  $\Phi1$  and  $\Phi2$  target functions change their value in an asymptotic way to achieve their minimum. This means that weight coefficients for both  $W1$  and  $W2$  matrixes are near the stable values and only small corrections are pursued. So, coordinator should study not only the target functions but their dynamic changing process too.

### 3 Example

In an example the main dynamic characteristics of the learning process are shown. The stress is made on the characteristic of the first level local target functions  $\Phi1, \Phi2$ . The structure of ANN is simple and can be described as ANN(3-5-1). This mean that ANN includes, 3 input neurons, 5 neurons in hidden layer and 1 output neuron. Sigmoid activation functions are implemented in both hidden and output layers. Three arguments of XOR function is fed as input data. So, every epoch includes 8 vectors. Changing different learning parameters as  $\alpha_1, \alpha_2, \gamma_2, \beta_1, \beta_2$  dynamic characteristics have been studied.

. In the second part of the test, the simple adaptive coordination algorithm was used. Fig.5. shows how the two target functions  $\Phi1, \Phi2$  changed their value during learning process ( iterations' number). The quality of dynamic processes is different. The function  $\Phi2$ , represent the second local target function ( output one). This process is smooth. This means that at (during) the learning process the value of  $\Phi2$  decreases at a constant rate to the minimum value. Midway through the process, its value decreases very slow. This is correlated with the first target function  $\Phi1$  ( hidden layer). This quality is quite different. From start to 3700 iterations target function  $\Phi1$  increased its value. Two local maximum in 1000 iterations and 3700 iterations are seen. After that, both  $\Phi1$  and  $\Phi1$  functions decreases their value and in the asymptotic way achieves the minimum.

As we stressed in previous sections, hidden layer can be divide into 5 sub-networks. Fig.6. shows the outputs of the three sub-networks ( $\phi_{1_1}, \phi_{1_3}, \phi_{1_4}$ ). The quality of dynamic characteristics are the same, but maximum of the amplitudes are different.

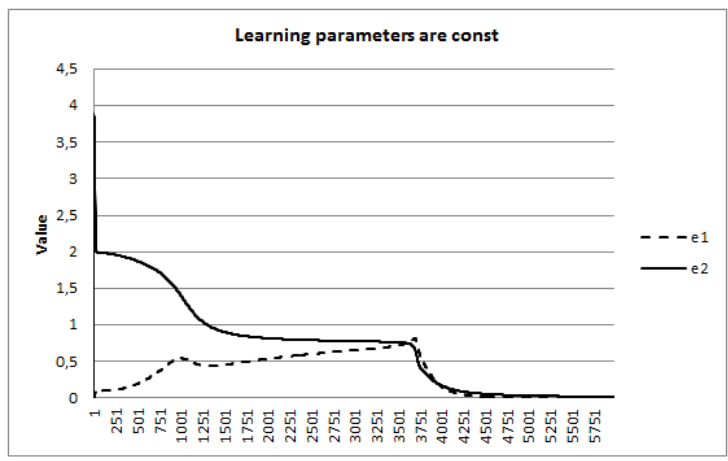


Figure 5:  $\alpha_1 = 0.3$ ,  $\alpha_2 = 0.05$ ,  $\gamma_2 = 0.5$ ,  $\beta_1 = \beta_2 = 0.01$

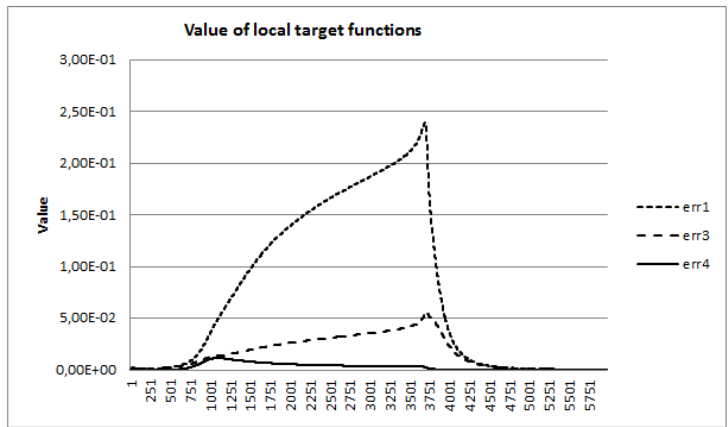


Figure 6: Value of  $\phi_1$ ,  $\phi_3$ ,  $\phi_4$  depending on iteration number



In the next figure (Fig.7), we can see that the quality of learning process depends on  $\gamma_2$  parameter. This parameter is calculated by coordinator and has impact on the forecast of the vector  $V_{21}$  value., For  $\gamma_2 = 0.1$  , that is too small, learning process isn't smooth. small oscillations can be seen. But if  $\gamma_2=0.5$  is too big, the amplitude increase its value more than 5 times . So, coordinator should calculate  $\gamma_2$  using own adaptive algorithm which should achieve from the first level and analyze the target functions  $\Phi_1$  and  $\Phi_2$ .

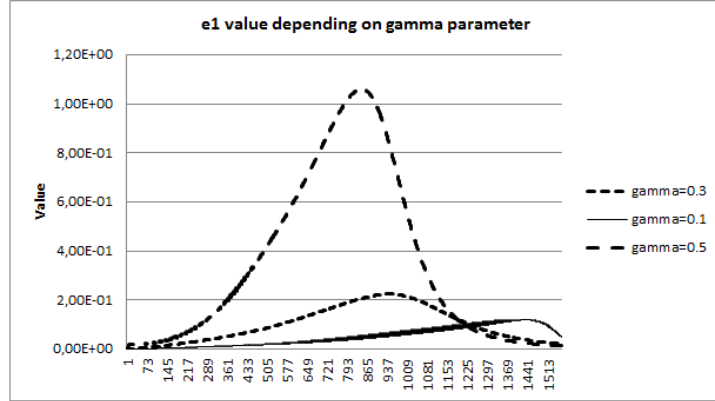


Figure 7: Dynamic characteristic of the learning process for the first layer depending on gamma

To study impact of the coordinator on the quality of learning process , adaptive algorithm changes two parameters  $\alpha_1$  – learning rate for the hidden layer, and  $\gamma_2$  - learning rate for the vector  $V_{21}$  . Vector  $V_{21}$  forecasts the hidden layer's output (Fig.8). When  $\Phi_2$  is greater than  $\Phi_1$  learning rates increases. Their values were increased in very small steps of only 0.05. Learning rates of both  $\alpha_1$  and  $\gamma_2$  shouldn't be extremely large or small.

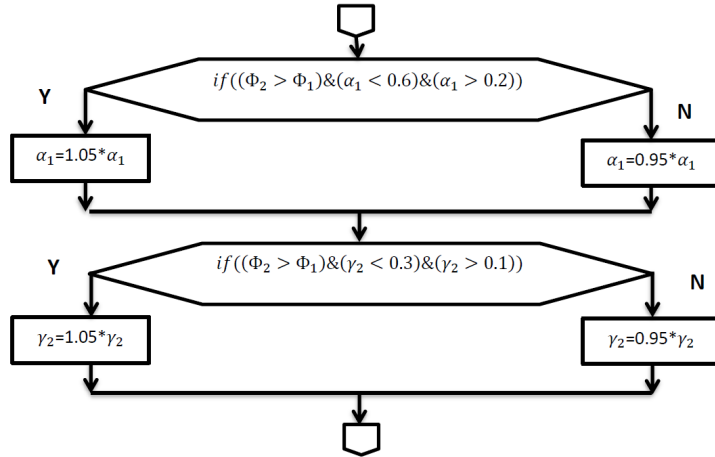


Figure 8: Coordination algorithm

So two extra constraints were used. (Fig.5). shows the coordinator's final

impact on the quality of the learning process. Target function  $\Phi_2$  decreases its value throughout the learning process, but target function  $\Phi_1$  still has the two maximum values. This problem will be studied in future work.

(Fig.10) shows how the value of two learning rates are changed by coordinator.

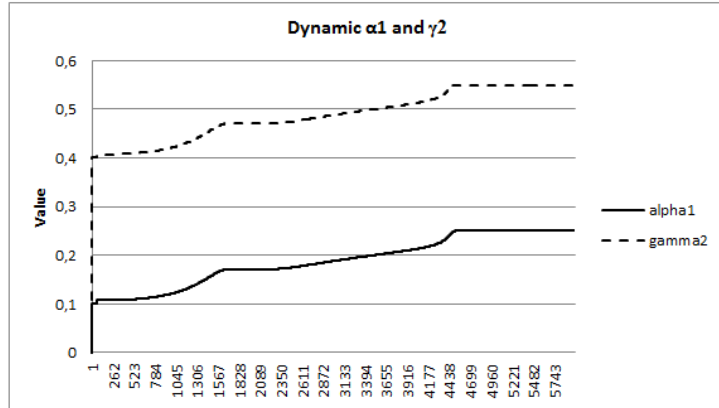


Figure 9: Value of  $\alpha_1$ ,  $\gamma_2$  changing by coordinator

## 4 Conclusion

In [15], few of coordination principles are defined for big hierarchical systems structure. In this article, the following principle is used - the forecast of the connections between sub-networks. In the hierarchical structure of ANN, coordinator should forecast the value of the vector  $V_2$ . This value should be the same as the real value of the hidden layer output  $V_1$ . In this situation global target function  $\Psi$  should achieve its minimum value and then the learning process is finishes.

If the first level of local target functions both  $\Phi_1$  and  $\Phi_2$  meet a couple of conditions [2][15], then convergence is guaranteed. Unfortunately the global target function  $\Psi$  isn't concave and could have a lot of local minimum. Therefore, it is not possible to prove that algorithm is stable and convergent mathematically. But the first local target functions didn't include any constraints and that helps while build learning algorithm. (Fig.10.) shows the final result of the different characteristics of the learning processes.

$\alpha_1$	$\alpha_2$	$\gamma_2$	$\beta_1$	$\beta_2$	Total iterations	Remarks
0,1	0,5	0,3	0,01	0,01	Oscillations	Big error
0,5	0,1	0,3	0,01	0,01	9500	Slow
0,3	0,05	0,3	0,01	0,01	7800	Faster
0,3	0,05	0,5	0,01	0,01	6000	Good
0,3	0,05	0,5	0,1	0,1	5000	The best

Figure 10: Impact of the different learning rates to the quality of the learning process

In the learning processes shown in (Fig.10). all rates were const. Coordinator calculates the new  $V21$  value using  $\gamma_2$ . (Fig.5.) shows that value of target function  $\Phi2$  doesn't change its value between 2000 and 3700 iterations. This is due to the fact that ANN in the first order has to stabilize the  $W1$  matrix weight coefficients. This process depends on  $V21$  vector. When all the  $W1$  weights coefficients are stable, matrix  $W2$  then stabilizes its weight coefficient. In this ANN, the first layer played the most important role. The sub-networks impact on the final value of the first layer's target function  $\Phi1$  is different. There are components in which its impact is very small. This can be explained by the hidden layer structure. The hidden layer includes structural neurons redundancy. Finally, the coordination algorithm is analyzed. Learning rates  $\alpha_1$  and  $\gamma_2$  didn't achieve their maximum value. Probably the value of the learning rate should be calculated using not only the relation between  $\Phi1$  and  $\Phi2$ , but also using their dynamic characteristics as the first difference  $\Delta\Phi1(n) = \Phi1(n) - \Phi1(n-1)$  and  $\Delta\Phi2(n) = \Phi2(n) - \Phi2(n-1)$ . This implies that coordinator should implement the PID controller algorithm.

$$\gamma_2(n+1) = \gamma_2(n) + \lambda_1 \cdot \Phi1(n) + \lambda_2 \cdot (\Phi1(n) - \Phi1(n-1)) \quad (16)$$

This two problems described above should be studied in the future work.

## References

- [1] Ch. M. Bishop, Pattern Recognition and Machine Learning, Springer Science + Business Media, LLC, 2006
- [2] W. Findeisen, J. Szymanowski, A. Wierzbicki, Teoria i metody obliczeniowe optymalizacji. Państwowe Wydawnictwo Naukowe, Warszawa 1977.
- [3] D. J. Montana, L. Davis Training Feed Forward Neural Networks Using Genetic Algorithms. IJCAI Detroit, Michigan 1989.
- [4] S. Osowski, Sieci Neuronowe do Przetwarzania Informacji. Oficyna Wydawnicza Politechniki Warszawskiej, Warsaw 2006.
- [5] S. Osowski, Sieci neuronowe w ujęciu algorytmicznym. WNT, Warszawa 1996.
- [6] Toshinori Munakate, Fundamentals of the New Artificial Intelligence. Second Edition, Springer 2008.
- [7] C. Fyle, Artificial Neural Networks and Information Theory, Department of Computing and information Systems, The University of Paisley, 2000.
- [8] A. Marciniak, J. Korbicz, J. Kus, Wstępne przetwarzanie danych, Sieci Neuronowe tom 6, Akademicka Oficyna Wydawnicza EXIT, Warsaw 2000.
- [9] Z. Mikrut, R. Tadeusiewicz, Sieci neuronowe w przetwarzaniu i rozpoznawaniu obrazów, Sieci Neuronowe tom 6, Akademicka Oficyna Wydawnicza EXIT 2000.
- [10] J. R. Rabunal, J Dorado, Artificial Neural Networks in Real-Life Applications, Idea Group Publishing 2006.

- [11] S.Placzek, B. Adhikari, Analysis of Multilayer Neural Networks with Direct and Cross-Forward Connection , CS&P Conference in the University of Warsaw, Warsaw 2013
- [12] A. Marciniak, J. Korbicz, Neuronowe sieci modularne, Sieci Neuronowe tom 6, Akademicka Oficyna Wydawnicza EXIT 2000.
- [13] Zeng-Guang Hou.Madan M.Gupta, Peter N. Nikiforuk, Min Tan, and Long Cheng, A Recurrent Neural Network for Hierarchical Control of Interconnected Dynamic Systems, IEEE Transactions on Neural Networks, Vol. 18, No. 2, March 2007.
- [14] L. Rutkowski, Metody i techniki sztucznej inteligencji, Wydawnictwo Naukowe PWN, Warsaw 2006.
- [15] M. D. Mesarovic, D. Macko, and Y. Takahara, Theory of hierarchical multilevel systems, Academic Press, New York and London, 1970.