

A Multipurpose Linear Component Analysis Method Based on Modulated Hebb Oja Learning Rule

Marko V. Jankovic, *Senior Member, IEEE*, and Masashi Sugiyama

Department of Computer Science, Tokyo Institute of Technology,
Japan

Abstract –This letter presents a Hebb-type learning algorithm for on-line linear calculation of principal components. The proposed method is based on a recently proposed cooperative-competitive concept, named the time-oriented hierarchical method. The algorithm performs deflation on the signal power rather than on the signal itself. It will be also shown when, or how, this algorithm can be used as a blind signal separation algorithm. The proposed synaptic efficacy learning rule does not need the explicit information about the value of the other efficacies to make individual efficacy modification. The number of necessary global calculation circuits is one.

Index terms - adaptive algorithm, principal/independent component analysis

EDICS: [SAS-STAT](#) Detection, estimation and classification theory and methods, statistical signal processing
[SAS-ADAP](#) Adaptive systems and adaptive filtering

Corresponding author:

Marko Jankovic,

Tokyo Institute of Technology, Department of Computer Science,

Ookayama 2-12-1, Meguro-ku, Tokyo, 152-8552, Japan;

Tel: +81-3- 5734-2699, FAX: +81-3-5734-2699

e-mail: marko@sg.cs.titech.ac.jp

1. Introduction

Neural networks provide a way for parallel on-line computations of principal component analysis (PCA) or principal subspace analysis (PSA) [1]. The goal of PSA or PCA is to extract from a stationary, random, zero-mean process $\mathbf{x}(i) \in \mathfrak{R}^K$, with covariance matrix $\mathbf{C} = E\{\mathbf{x}(i)\mathbf{x}(i)^T\}$, the subspace spanned by the N principal eigenvectors of \mathbf{C} or principal vectors themselves, respectively. Within last years various PCA and PSA learning algorithms have been proposed and mathematically investigated (see e.g. [1]). Most of the proposed algorithms are based on local Hebbian learning - due to locality it has been argued that these algorithms are biologically plausible. However, it is not frequently analyzed how those algorithms could be related to the known biological neural architectures or could they have multiple purposes.

In [2], biologically inspired PSA methods, named Modulated Hebbian (MH) and Modulated Hebb-Oja (MHO) learning rules have been introduced. Major objectives for the methods' derivation were to obtain a network which has a learning rule for individual synaptic efficacy that requires the least possible amount of explicit information about the other synaptic efficacies, especially those related to other neurons and to minimize the neural hardware that is necessary for implementation of the proposed learning rule.

In this paper modification of the MHO algorithm, named the Multiurpose Linear Component Analysis MHO (MILICA-MHO) algorithm is proposed and analyzed. Generally speaking the MILICA-MHO algorithm performs PCA. If the input signal is prewhitened the same algorithm can perform blind signal separation (BSS). The new algorithm is obtained by implementation of the modified (generalized) time-oriented hierarchical method [5,6] on the MHO method. Comparing to bigradient or related nonlinear PCA algorithms [7] which can be used in a similar way, or Generalized Hebbian Learning (GHA) [9] the proposed algorithm is much more based on local calculations.

2. Theory

Let $\mathbf{x} \in \mathfrak{R}^K$ denote the input random variables with mean zero, and let $\mathbf{y} = \mathbf{W}^T \mathbf{x} \in \mathfrak{R}^N$ denotes the output, where $\mathbf{W} \in \mathfrak{R}^{K \times N}$ denote the synaptic weight matrix. In scalar form, for the n -th output, we have $y_n = \mathbf{w}_n^T \mathbf{x}$ ($n = 1, \dots, N$), where $\mathbf{w}_n \in \mathfrak{R}^K$ denotes the column vector of \mathbf{W} .

The Modulated Hebbian (Oja) (MHO) learning rule is introduced in [2] and analyzed in [3, 4]. The MH(O) rule can be derived as a gradient descent learning rule for minimization of the following cost function:

$$\begin{aligned} J(\mathbf{W}) &:= \mathbb{E}\left(\|\mathbf{x} - \mathbf{W}\mathbf{y}\|^4\right) = \mathbb{E}\left(\left(\|\mathbf{x} - \mathbf{W}\mathbf{W}^T\mathbf{x}\|^2\right)^2\right) \\ &= \mathbb{E}\left(\left(\text{tr}(\mathbf{x}\mathbf{x}^T) - 2\text{tr}(\mathbf{W}^T\mathbf{x}\mathbf{x}^T\mathbf{W}) + \text{tr}(\mathbf{W}^T\mathbf{x}\mathbf{x}^T\mathbf{W}\mathbf{W}^T\mathbf{W})\right)^2\right), \quad \|\mathbf{w}_n\| = 1 \end{aligned} \quad (1)$$

or, under assumption $\mathbf{W}^T\mathbf{W}=\mathbf{I}$ (\mathbf{I} is identity matrix), in compact notation

$$\begin{aligned} J_N^{MHP\text{SA}}(\mathbf{W}) &:= \mathbb{E}\left\{\left(\mathbf{x}^T\mathbf{x} - \mathbf{y}^T\mathbf{y}\right)^2\right\} = \mathbb{E}\left\{\left(\mathbf{x}^T(\mathbf{I} - \mathbf{W}\mathbf{W}^T)\mathbf{x}\right)^2\right\} \\ &= \mathbb{E}\left\{\left(\sum_{k=1}^K x_k^2 - \sum_{n=1}^N y_n^2\right)^2\right\}, \quad \|\mathbf{w}_n\| = 1. \end{aligned} \quad (2)$$

As it was explained in [2-4] MHO learning algorithm represents the learning rule in which individual synaptic efficacy modification does not require the explicit information about the other synaptic efficacies, especially those related to other neurons. In order to obtain a new PCA algorithm with the same feature, we will introduce the following set of cost functions:

$$\begin{aligned} J_n^{MHP\text{CA}}(\mathbf{W}) &:= J_N^{MHP\text{SA}}(\mathbf{W}) + f(\mathbf{a})\mathbb{E}\left\{\left(\sum_{k=1}^K x_k^2 - \sum_{j=1}^n y_j^2\right)^2\right\}, \quad \|\mathbf{w}_n\| = 1, \\ f(\mathbf{a}) &= \begin{cases} \mathbf{a}, & n \neq N, \\ 0, & n = N. \end{cases} \end{aligned}$$

It can be noticed that for every weight vector \mathbf{w}_n we have a different learning rule which consists of one common part and one part which is specific for the n -th output neuron. By minimization of the proposed cost function the new learning rule for the n -th weight vector can be written in the form

$$\begin{aligned} \mathbf{w}_n(i+1) &= \mathbf{w}_n(i) + g(i)\left(\mathbf{x}(i)y_n(i) - \mathbf{w}_n(i)y_n^2(i)\right)\left(\|\mathbf{x}(i)\|^2 - \|\mathbf{y}(i)\|^2\right) \\ &\quad + f(\mathbf{a})g(i)\left(\mathbf{x}(i)y_n(i) - \mathbf{w}_n(i)y_n^2(i)\right)\left(\|\mathbf{x}(i)\|^2 - \sum_{j=1}^n y_j^2(i)\right), \quad n = 1, \dots, N, \end{aligned} \quad (3)$$

where $\gamma(i)$ is the learning factor, and i defines steps in time (t). We actually have a system of equations that have the same family (common) part of the learning rule, while all individual (specific) parts of

the learning rules are different. While the GHA algorithm [9] for modification of synaptic vector w_n (related to output neuron n) requires explicit information about all synaptic efficacies related to output neurons 1 to $n-1$, the proposed algorithm does not need any explicit information about the synaptic efficacies related to any other neurons.

Now we can state the following theorem (proof is diverted to Appendix):

Theorem 1: Columns of W evolve to span the principal subspace of C under the assumption that the input is bounded and a is properly chosen (in a sense that it is possible to achieve stable solution). If $\text{rank}[W] = N$, then the space of locally-asymptotically-stable stationary points satisfy $W^T W = I$ and columns of W will be equal to principal eigenvectors of C .

3. Numerical experiments

Here, we will examine the small scale numerical simulations results. The number of inputs was set at $K = 5$ and the number of output neurons was $N = 3$. Artificial zero-mean vectors with uncorrelated elements were generated by the following equations:

$$\begin{aligned} x(1, i) &= .45 \sin(i/2); \\ x(2, i) &= .45 ((\text{rem}(i,23) - 11)/9).^5; \\ x(3, i) &= .35 \sin(i/17.8); \\ x(4, i) &= .145 ((\text{rand}(1,1) < .5) * 2 - 1) \cdot \log(\text{rand}(1,1) + .5); \\ x(5, i) &= .18 \text{ randn}(1,1), \end{aligned}$$

where rem , rand and randn represent standard MATLAB functions (rem - remainder after integer division, rand – random number generator from uniform distribution and randn – random number generator from normal distribution). The input signal is constructed as $s = 0.47 * \text{mix} * x$, where the mixing matrix mix is defined as $\text{mix} = -.5 + \text{rand}(K)$.

In Fig.1 cosine of the angles between column vectors of matrix W and numerically calculated eigenvectors of the input signal covariance matrix were used as illustration of the algorithm efficiency. The learning rate was chosen constant, $g_i = 3.45$, for the first 15000 iterations, and then was set at $g_i = 0.115$. The individual part was taken as $\alpha=0.5$. The initial value for matrix W was selected as $W_{\text{init}} = -.5 + \text{rand}(5,3)$.

In Fig. 2 we can see that proposed method can be efficiently used for extraction of the deterministic input signals (under the assumption that mixing matrix is orthogonal). The first column in Fig. 2 represents original signals, the second column represents the input signals obtained from the original signals by multiplication with the orthonormal matrix *mix*, and the last column represents the output obtained by implementation of MILICA-MHO.

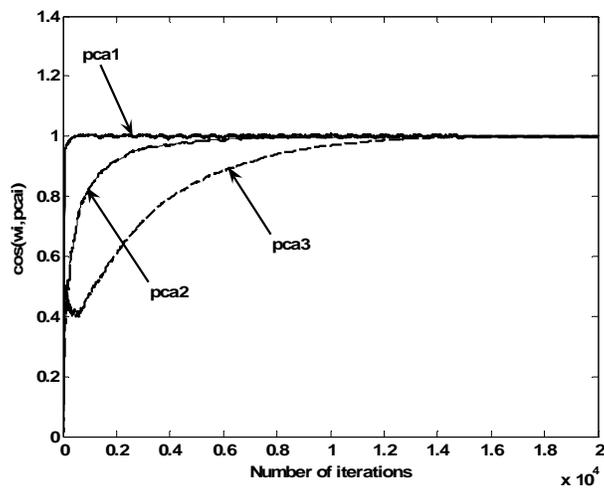


Fig. 1. Cosine of angles between column vectors of W and three principal eigenvectors of input covariance matrix versus the number of iterations

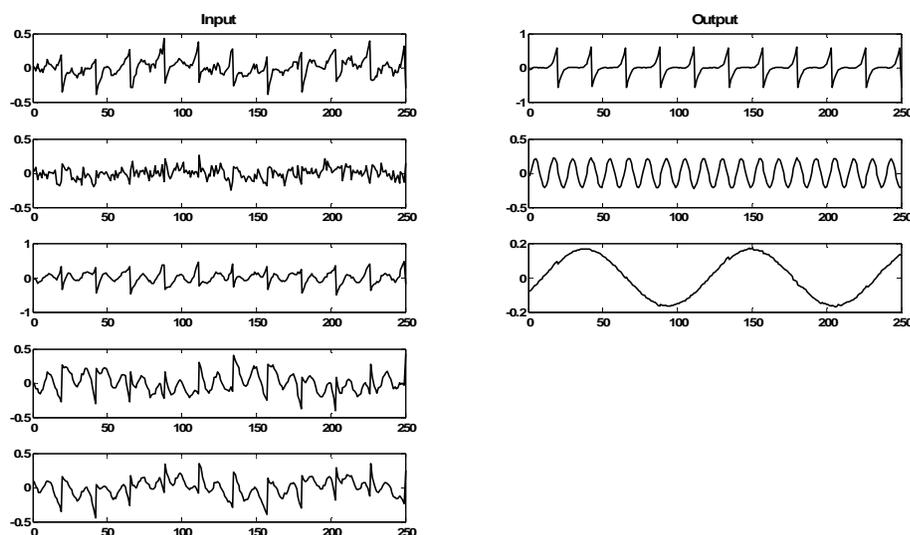


Fig. 2. Blind signal extraction of deterministic components (mixing matrix is orthonormal)

4. Can we use the proposed method for Blind Signal Separation?

In order to show how the algorithm proposed in this paper can be used as a BSS algorithm in some cases and in order to explain results obtained in Fig. 2, we will consider the simplest case $k = n = 2$. Also, we will assume that the input data is prewhitened. Then, the MILICA-MHO algorithm can be represented by the following equations:

$$\begin{aligned} \mathbf{w}_1(i+1) &= \mathbf{w}_1(i) + \mathbf{g}(i) \left(\mathbf{x}(i)y_1(i) - \mathbf{w}_1(i)y_1^2(i) \right) \left(\|\mathbf{x}(i)\|^2 - \|\mathbf{y}(i)\|^2 \right) \\ &\quad + \mathbf{a}\mathbf{g}(i) \left(\mathbf{x}(i)y_2(i) - \mathbf{w}_2(i)y_2^2(i) \right) \left(\|\mathbf{x}(i)\|^2 - y_1^2(i) \right) \\ \mathbf{w}_2(i+1) &= \mathbf{w}_2(i) + \mathbf{g}(i) \left(\mathbf{x}(i)y_2(i) - \mathbf{w}_2(i)y_2^2(i) \right) \left(\|\mathbf{x}(i)\|^2 - \|\mathbf{y}(i)\|^2 \right) \end{aligned} \quad (4)$$

that are obtained by minimization of the following two cost functions

$$\begin{aligned} J_1^{MHPCA}(\mathbf{W}) &= J_2^{MHPSA}(\mathbf{W}) + \mathbf{a}\mathbf{E} \left\{ \left((x_1^2 + x_2^2) - y_1^2 \right)^2 \right\} \\ J_2^{MHPSA}(\mathbf{W}) &= \mathbf{E} \left\{ \left(\sum_{k=1}^2 x_k^2 - \sum_{j=1}^2 y_j^2 \right)^2 \right\} \end{aligned} \quad (5)$$

We can see that both cost functions contain a part which is used for symmetric orthogonalization while J_1^{MHPCA} contains also an additional part which directs the solution in some particular direction. Since we know that the input data is priwhitened and that the final solution will be an orthogonal matrix, we can approximately rewrite (5) as

$$\begin{aligned} J_1^{MHPCA}(\mathbf{W}) &\approx J_2^{MHPSA}(\mathbf{W}) + \mathbf{a}\mathbf{E} \left\{ \left(\sum_{k=1}^2 y_k^2 - y_1^2 \right)^2 \right\} \\ J_2^{MHPSA}(\mathbf{W}) &= \mathbf{E} \left\{ \left(\sum_{k=1}^2 x_k^2 - \sum_{j=1}^2 y_j^2 \right)^2 \right\} \end{aligned} \quad (6)$$

or

$$\begin{aligned} J_1^{MHPCA}(\mathbf{W}) &\approx J_2^{MHPSA}(\mathbf{W}) + \mathbf{a}\mathbf{E}(y_2^4) \\ J_2^{MHPSA}(\mathbf{W}) &= \mathbf{E} \left\{ \left(\sum_{k=1}^2 x_k^2 - \sum_{j=1}^2 y_j^2 \right)^2 \right\} \end{aligned} \quad (7)$$

since, in the vicinity of the solution it holds

$$\sum_{k=1}^2 x_k^2 \cong \sum_{j=1}^2 y_j^2. \quad (8)$$

So, if at most one of the input signals is supergaussian our method will perform signal separation from their prewhitened mixture, since the algorithm (7) is going to minimize the 4th moment of the output [7]. This can be easily extended to dimensions higher than 2.

If all signals, except one, are supergaussian, then the algorithm can perform BSS, given that the output is calculated as $y_n = -\mathbf{w}_n^T \mathbf{x}$ ($n = 1, \dots, N$), where $\mathbf{w}_n \in \mathfrak{R}^K$ denotes the column vector of \mathbf{W} . This can also be extended to higher dimensions.

5. Conclusion

In this letter, we have proposed a novel PCA/BSS algorithm based on modification of the PSA MH(O) algorithm. The algorithm employs a deflation technique that is implemented on the signal power, which could be seen as a more natural solution (especially in the PCA case). Locality of calculations could be seen as a desirable property for possible implementation in parallel hardware. If we know in advance the number of supergaussian and subgaussian signals, the algorithm can be easily adjusted to perform BSS of any kind of signals. If it is not the case, the algorithm and network structure have to be additionally modified. Theory and simulations indicate that the algorithm works well.

Appendix

Proof of Theorem 1:

Equation (3) can be written in the following form

$$\begin{aligned} \mathbf{w}_n(i+1) = & \mathbf{w}_n(i) \\ & + \mathbf{g}(i) \left(\mathbf{C}(i) \left(\|\mathbf{x}(i)\|^2 - \|y(i)\|^2 + f(\mathbf{a}) \left(\|\mathbf{x}(i)\|^2 - \sum_{j=1}^n y_j^2(i) \right) \right) \mathbf{w}_n(i) - \mathbf{w}_n(i) \mathbf{q}(i) \right), \end{aligned} \quad (9)$$

where $\mathbf{C}(i) = \mathbf{x}(i)\mathbf{x}(i)^T$, and $\theta(i)$ is defined as

$$\mathbf{q}(i) = y_n^2 \left(\|\mathbf{x}(i)\|^2 - \|\mathbf{y}(i)\|^2 + f(\mathbf{a}) \left(\|\mathbf{x}(i)\|^2 - \sum_{j=1}^n y_j^2(i) \right) \right). \quad (10)$$

Differential counterpart of this difference equation can be written in the form [8, 9]

$$\frac{d\mathbf{w}_n}{dt} = \mathbf{F}_n \mathbf{w}_n - \mathbf{w}_n \Theta_n \quad n = 1, 2, \dots, N, \quad (11)$$

where \mathbf{F}_n is defined as

$$\begin{aligned} \mathbf{F}_n &= \mathbb{E} \left(\mathbf{C}(i) \left(\|\mathbf{x}(i)\|^2 - \|\mathbf{y}(i)\|^2 + f(\mathbf{a}) \left(\|\mathbf{x}(i)\|^2 - \sum_{j=1}^n y_j^2(i) \right) \right) \right) \\ &= \mathbb{E} \left(\mathbf{C}(i) \left(\|\mathbf{x}(i)\|^2 - \|\mathbf{y}(i)\|^2 \right) \right) + f(\mathbf{a}) \mathbb{E} \left(\mathbf{C}(i) \left(\left(\|\mathbf{x}(i)\|^2 - \sum_{j=1}^n y_j^2(i) \right) \right) \right) \\ &= \mathbf{F}f + \mathbf{F}i_n, \end{aligned} \quad (12)$$

where expectation is over \mathbf{x} and Θ_n is defined as

$$\begin{aligned} \Theta_n &= \mathbb{E} \left(y_n^2 \left(\|\mathbf{x}(i)\|^2 - \|\mathbf{y}(i)\|^2 + f(\mathbf{a}) \left(\|\mathbf{x}(i)\|^2 - \sum_{j=1}^n y_j^2(i) \right) \right) \right) \\ &= \mathbb{E} \left(y_n^2 \left(\|\mathbf{x}(i)\|^2 - \|\mathbf{y}(i)\|^2 \right) \right) + f(\mathbf{a}) \mathbb{E} \left(y_n^2 \left(\left(\|\mathbf{x}(i)\|^2 - \sum_{j=1}^n y_j^2(i) \right) \right) \right) \\ &= \Theta f_n + \Theta i_n. \end{aligned} \quad (13)$$

It is not difficult to conclude that, at stationary points, equation (11) represents system of equations for calculations of eigenvector of matrices \mathbf{F}_n . It is possible to prove that all matrices \mathbf{F}_n have same eigenvectors as matrix $\mathbf{C} = \mathbb{E}(\mathbf{x}\mathbf{x}^T)$ – it can be done by straight calculation if we use empirical expectations and than show that matrices \mathbf{F}_n commute with matrix \mathbf{C} . Also, it can easily be shown that all matrices \mathbf{F}_n are symmetric – they represent expectation of the symmetric matrix multiplied by some scalar. Since \mathbf{F}_n are real and symmetric then, if matrix \mathbf{W} is of maximum rank for all i (for all time t), it is easy to see that matrix \mathbf{W} is an orthonormal matrix.

In the analysis that follows, stable points of equation (11) will be investigated. In order to prove that original algorithm (9) will converge toward the same stable points, it should be proved that outputs are bounded and that (9) visits infinitely often, almost surely, a compact subset of the domain

of attraction of the asymptotically stable solution of (11) [8]. By implementation of the same method that was proposed in [4], it can be shown that synaptic vectors have bounded norms under some reasonable assumptions. This, together with assumptions of the boundness of the input signals, means that outputs are bounded. However, it must be said that proving that (9) visits infinitely often, almost surely, a compact subset of the domain of attraction of the asymptotically stable solution of (11) is very difficult. Such kind of proof does not exist for many of the known algorithms. Based on experience, generally can be said that if the desired limit of the discrete algorithm is not an asymptotically stable point of the averaged differential equation, then the convergence will not take place and the behavior of the discrete algorithm is not good [4]. Here will be said that in all simulations in which (9) converged it had the same stable points as (11).

Now, we will rewrite the equation (11) in the following form

$$\begin{aligned}
\frac{d\mathbf{w}_n}{dt} &= (\mathbf{F}\mathbf{f} + \mathbf{F}\mathbf{i}_n)\mathbf{w}_n - \mathbf{w}_n(\Theta\mathbf{f}_n + \Theta\mathbf{i}_n) \\
&= (\mathbf{F}\mathbf{f}\mathbf{w}_n - \mathbf{w}_n\Theta\mathbf{f}_n) + (\mathbf{F}\mathbf{i}_n\mathbf{w}_n - \mathbf{w}_n\Theta\mathbf{i}_n) \\
&= \frac{d\mathbf{w}\mathbf{f}_n}{dt} + \frac{d\mathbf{w}\mathbf{i}_n}{dt} \quad n=1,2,\dots,N.
\end{aligned} \tag{14}$$

Without a loss of generality, we can assume that coordinate system coincides with eigenvectors of matrix C . Then, the input data is decorrelated and covariance matrix of input data is diagonal. In that case analysis that was carried on in [4] could be applied. From [4, Theorem 4.3] we know that if \mathbf{W} spans any subspace defined by some of the eigenvectors of matrix C then family part of the learning equations (11), denoted by $d\mathbf{w}\mathbf{f}_n/dt$, is zero, so equation (14) can be written as

$$\frac{d\mathbf{w}_n}{dt} = \frac{d\mathbf{w}\mathbf{i}_n}{dt} = \mathbf{F}\mathbf{i}_n\mathbf{w}_n - \mathbf{w}_n\Theta\mathbf{i}_n \quad n=1,2,\dots,N. \tag{15}$$

In the case $n=1$ equation (15) represents the equation for calculation of the principal eigenvector of matrix C (having in mind [4]). Then, for $n=2$, this equation represents the equation for calculation of the eigenvector which is orthonormal to the first principal eigenvector \mathbf{w}_1 of C (\mathbf{W} is an orthonormal matrix) and corresponds to the maximum possible eigenvalue. So, it represents equation for calculation of eigenvector that corresponds to the second largest eigenvalue of the matrix C . This can be seen as a deflation procedure in which the data are not back projected and then subtracted from the original

data, but rather the energy of input signal is decreased by the amount that corresponds to the energy contained in the direction of the principal component. So, for $n > 2$ we can continue the same way of reasoning and conclude that w_n represents the eigenvector of matrix C which corresponds to n -th biggest eigenvalue of the matrix C . This concludes the proof. Proper range of the parameter a can be assessed through stability analysis.

REFERENCES

- [1] A. Chichocki, S.-I. Amari, (2003) "Adaptive Blind Signal and Image Processing – Learning Algorithms and Applications", John Wiley and Sons, New York.
- [2] M. Jankovic (2001) "A new modulated Hebbian learning rule – Method for local computation of a principal subspace", ICONIP 2001: 470-475
- [3] M. Jankovic and H. Ogawa (2003) "A new modulated Hebb learning rule – Biologically plausible method for local computation of principal subspace", Int. J. Neural Systems 13: 215-224.
- [4] M. Jankovic and H. Ogawa (2006) "Modulated Hebb-Oja learning rule – A method for principal subspace analysis", IEEE Trans. on Neural Networks 17: 345-356.
- [5] M. Jankovic and B. Reljin (2005) "Neural learning on Grassman/Stiefel principal/minor submanifold", EUROCON 2005: 249-252.
- [6] M. Jankovic and H. Ogawa (2004) "Time-Orineted Hierarachical Method for Computation of Principal Components using Subspace Learning Algorithm", Int. J. Neural Systems 14: 313-324.
- [7] J. Karhunen, E. Oja, L. Wang, R. Vigario, J. Joutsalo (1997) "A Class of Neural Networks for Independent Component Analysis", IEEE Trans. on Neural Networks 8: 486-504.
- [8] L. Ljung (1977) "Analysis of recursive stochastic algorithms", IEEE Trans. Automat. Contr. 22: 551-575.
- [9] T. Sanger (1989) "Optimal Unsupervised Learning in a Single-Layer Linear Feedforward Neural Network", Neural Networks 2: 459-473.