

# GENERALIZED SUBSPACE RULES FOR ON-LINE PCA AND THEIR APPLICATION IN SIGNAL AND IMAGE COMPRESSION

Toshihisa Tanaka

Department of Electrical and Electronic Engineering,  
Tokyo University of Agriculture and Technology (TUAT)  
Nakacho, Koganei-shi, Tokyo, 184-8588, Japan. tanakat@cc.tuat.ac.jp

## ABSTRACT

Weighted subspace (WS) algorithms developed by Oja and Xu for PCA are unified into generalized forms and theoretically analyzed. It is then proved that the generalized rules are stable at only the fixed point where the principal components are extracted. We moreover find the optimal parameter in terms of the preservation of orthogonality of estimated principal components during tracking. To understand the theoretical behavior, then, toy numerical examples are shown. Moreover, a possibility for the application of adaptive data compression is discussed, by showing examples of backward adaptation image coding.

## 1. INTRODUCTION

Principal component analysis (PCA), or equivalently called the Karhunen-Loève transform, is important and extensively used in the fields of signal compression, neural networks pattern recognition, and so on [1]. The principal components (PCs) can be obtained by solving an eigenvalue problem of the correlation matrix  $\mathbf{R} = E_{\mathbf{x}}[\mathbf{x}\mathbf{x}^T]$ , where  $E_{\mathbf{x}}[\cdot]$  denotes the ensemble average corresponding to a stochastic vector  $\mathbf{x}$  in the Euclidean space  $\mathbb{R}^N$ , and  $\cdot^T$  denotes the transposition. We presume in this paper that  $\mathbf{R}$  has full rank and distinct eigenvalues. Then, there exist the  $N$  eigenvectors  $\phi_i$ ,  $i = 0, \dots, N-1$  and the corresponding eigenvalues  $\lambda_0 > \dots > \lambda_{N-1} > 0$ . The first  $K$  eigenvectors correspond to the  $K$  PCs and span the principal subspace of dimension  $K$ .

Estimation and tracking of the principal subspace/components by the gradient descent method can be accomplished by the update as

$$\mathbf{W}(n+1) = \mathbf{W}(n) + \beta_n \Delta \mathbf{W}(n) \quad (1)$$

where  $\Delta \mathbf{W}(n)$  is the updating matrix,  $\mathbf{W}(n)$  is the principal subspace/components estimate of size  $N \times K$  at the  $n$ th iteration step, and  $\beta_n$  is a positive learning parameter. Several updating algorithms have been proposed. The well-known method for principal subspace tracking proposed by Oja [2] is given by

$$\Delta \mathbf{W}(n) = \mathbf{x}(n)\mathbf{y}^T(n) - \mathbf{W}(n)\mathbf{y}(n)\mathbf{y}^T(n), \quad (2)$$

where  $\mathbf{y}(n) = \mathbf{W}^T(n)\mathbf{x}(n)$ . By applying the direct differentiation to the cost function

$$J[\mathbf{W}] = E_{\mathbf{x}} \|\mathbf{x} - \mathbf{W}\mathbf{W}^T \mathbf{x}\|^2, \quad (3)$$

This work was partially supported by Japan Society for the Promotion of Science (JSPS), Grant-in-Aid for Scientific Research (C) (2), 15500101, 2003. He is also with Laboratory for Advanced Brain Signal Processing, RIKEN Brain Science Institute, Japan.

where  $\mathbf{W}$  is of  $N \times K$ , we obtain the modified algorithm given as

$$\Delta \mathbf{W}(n) = 2\mathbf{x}(n)\mathbf{y}^T(n) - \mathbf{W}(n)\mathbf{y}(n)\mathbf{y}^T(n) - \mathbf{x}(n)\mathbf{y}^T(n)\mathbf{W}^T(n)\mathbf{W}(n), \quad (4)$$

which is proposed by Yang [3] and Xu [4], independently. This rule is sometimes referred to as normalized Oja's rule [5]. This rule also preserves orthogonality of  $\mathbf{W}$  during learning. In the sense of the principal subspace tracking, it is reported in [5] that the learning rule given in (4) works better compared to Oja's rule as in (2). In many applications, however, adaptive rules which track not only a principal subspace but also PCs are desirable. For this purpose, Oja *et al* proposed in [6] a weighted version of the subspace rule given as

$$\Delta \mathbf{W}(n) = \mathbf{x}(n)\mathbf{y}^T(n) - \mathbf{W}(n)\mathbf{y}(n)\mathbf{y}^T(n)\mathbf{D}, \quad (5)$$

where  $\mathbf{D}$  is a diagonal matrix with positive entries in increasing order  $0 < d_0 < \dots < d_{K-1}$ . Similarly, Xu proposed the weighted subspace rule given as

$$\Delta \mathbf{W}(n) = \mathbf{x}(n)\mathbf{y}^T(n)\mathbf{D}^{-1} - \mathbf{W}(n)\mathbf{D}^{-1}\mathbf{y}(n)\mathbf{y}^T(n), \quad (6)$$

and its normalized version given as

$$\begin{aligned} \Delta \mathbf{W}(n) = & 2\mathbf{x}(n)\mathbf{y}^T(n)\mathbf{D}^{-1} - \mathbf{W}(n)\mathbf{D}^{-1}\mathbf{y}(n)\mathbf{y}^T(n) \\ & - \mathbf{x}(n)\mathbf{y}^T(n)\mathbf{D}^{-1}\mathbf{W}^T(n)\mathbf{W}(n). \end{aligned} \quad (7)$$

As we have seen, the weighted learning rules for PCA summarized here can be categorized into two groups: weighted subspace (WS) rules, and the normalized weighted subspace (NWS) rules. Based on these observation, in this paper, we propose a unified approach of WS rules that allow individual PCs to be extracted. In the rest of the paper, firstly, we propose generalized forms for the two groups of weighted rules. Secondly, it is proved that the generalized rules are stable at only the point where the PCs are extracted. It is finally shown that in special cases, orthogonality of  $\mathbf{W}$  is preserved during the learning. It is also proved that in these cases, the normalized rule preserves orthogonality more strongly than the non-normalized rules. Simple numerical examples are provided for illustrating theoretically justified properties of the proposed rules. Moreover, we provide examples of the application to backward adaptation of transform image coding. We would understand theoretical results through the numerical simulations.

## 2. GENERALIZED WEIGHTED SUBSPACE RULES AND THEIR STABILITY

The generalized forms for non-normalized and normalized weighted subspace rules are derived in this section. By replacing  $\mathbf{W}\mathbf{D}^{-1/2}$  by

$\mathbf{W}$  in the Xu's weighted rule as in (6), we have the modified form given as

$$\Delta \mathbf{W}(n) = \mathbf{x}(n)\mathbf{y}^T(n)\mathbf{D}^{-1} - \mathbf{W}(n)\mathbf{y}(n)\mathbf{y}^T(n) \quad (8)$$

In the same way as the above, the normalized version as in (7) becomes

$$\Delta \mathbf{W}(n) = \mathbf{x}(n)\mathbf{y}^T(n)\mathbf{D}^{-1} - \mathbf{W}(n)\mathbf{y}(n)\mathbf{y}^T(n) - \mathbf{x}(n)\mathbf{y}^T(n)\mathbf{W}^T\mathbf{W}. \quad (9)$$

Comparing Xu's rules (8) and (9) with Oja's rule (2), we may unify them and generalize into the following unified rules:

- Generalized WS (GWS) rule:

$$\Delta \mathbf{W}(n) = \mathbf{x}(n)\mathbf{y}^T(n)\mathbf{D}^{-p} - \mathbf{W}(n)\mathbf{y}(n)\mathbf{y}^T(n)\mathbf{D}^{1-p}, \quad (10)$$

- Generalized NWS (GNWS) rule:

$$\Delta \mathbf{W}(n) = \mathbf{x}(n)\mathbf{y}^T(n)\mathbf{D}^{-p} - \mathbf{W}(n)\mathbf{y}(n)\mathbf{y}^T(n)\mathbf{D}^{1-p} - \mathbf{x}(n)\mathbf{y}^T(n)\mathbf{W}^T\mathbf{W}\mathbf{D}^{1-p}. \quad (11)$$

where  $p$  is a scalar constant, which will change the learning behavior.

If we set in the GWS rule that  $p = 0$  and  $p = 1$ , we obtain the Oja's WS rule (5) and the Xu's WS rule (6), respectively. Furthermore, if we set in the GNWS rule that  $p = 1$ , we obtain the Xu's NWS rule (9).

In order to analyze the stochastic behaviors of the algorithms, as done, for instance, in [1, 4, 6], we introduce the corresponding ordinary differential equations (ODE) given as

$$\frac{d\mathbf{W}}{dt} = \mathbf{R}\mathbf{W}\mathbf{D}^{-p} - \mathbf{W}\mathbf{W}^T\mathbf{R}\mathbf{W}\mathbf{D}^{1-p}, \quad (12)$$

for the GWS rule and

$$\frac{d\mathbf{W}}{dt} = 2\mathbf{R}\mathbf{W}\mathbf{D}^{-p} - \mathbf{W}\mathbf{W}^T\mathbf{R}\mathbf{W}\mathbf{D}^{1-p} - \mathbf{R}\mathbf{W}\mathbf{W}^T\mathbf{W}\mathbf{D}^{1-p}, \quad (13)$$

for the GNWS rule.

First of all, we find the fixed points of the gradients (10) and (11). Assume that  $\mathbf{R}$  has nonzero distinct eigenvalues  $\lambda_0 > \lambda_1 > \dots > \lambda_{N-1}$ . Let  $\mathbf{\Lambda}$  be a diagonal matrix given as  $\mathbf{\Lambda} = \text{diag}[\lambda_{\pi(0)}, \dots, \lambda_{\pi(N-1)}]$  where  $\pi$  is a permutation of  $\{0, \dots, N-1\}$ , that is,  $\{\pi(0), \dots, \pi(N-1)\} = \{0, \dots, N-1\}$ , and  $\text{diag}[\cdot]$  denotes a diagonal matrix. Then, an eigenvalue decomposition may be given as  $\mathbf{R} = \mathbf{P}\mathbf{\Lambda}\mathbf{P}^T$ , where  $\mathbf{\Lambda} = \text{diag}[\lambda_{\pi(0)}, \dots, \lambda_{\pi(N-1)}]$ . Now, we present the following result.

**Theorem 1** *Let  $\bar{\mathbf{P}}$  be an  $N \times K$  matrix whose columns correspond to eigenvectors of  $\mathbf{R}$  in arbitrary order without duplication. The fixed point is then given by  $\mathbf{W} = \bar{\mathbf{P}}\mathbf{D}^{-1/2}$ .*

Proofs in this paper are omitted due to lack of space. All the proofs are provided in [7].

Keep in mind that the theorem means that at any fixed point,  $\mathbf{W}$  is in the manifold  $S$ .

## 2.1. Stability of the Fixed Points

We investigate in this section the stability of the fixed points. Firstly, we show the main result as follows.

**Theorem 2** *The rules (10) and (11) are stable if and only if  $\mathbf{W}$  is at the fixed point  $\mathbf{P}^*\mathbf{D}^{-1/2}$ , where  $\mathbf{P}^*$  is the  $M \times K$  matrix whose columns are the  $K$  PCs in decreasing order of the eigenvalues, that is,  $\mathbf{P}^* = [\phi_0, \dots, \phi_{K-1}]$ .*

## 2.2. Stability on the Manifold

We investigate the behavior of the dynamics of  $(\mathbf{W}^T\mathbf{W})(t)$  in this subsection. It is shown that in some particular cases, the manifold  $S$  is stable under the rules (10) and (11). The stability of the manifold is a desirable property for adaptive algorithms, since no additional normalization, which would make the system more complicated, is needed.

Before the stability analysis on the manifold  $S$ , we find the special cases where the manifold  $S$  is stable. Assume that at the  $n$ th iteration,  $\mathbf{W}(n)$  is in the manifold  $S$ , that is,  $\mathbf{W}^T(n)\mathbf{W}(n) = \mathbf{D}^{-1}$ . Then, the GNWS rule (11) becomes identical to the GWS rule (10). From (1) and (10), we have

$$\begin{aligned} \mathbf{W}^T(n+1)\mathbf{W}(n+1) &= \mathbf{D}^{-1} + \beta \left[ (\mathbf{D}^{-p}\mathbf{W}^T\mathbf{R}\mathbf{W} + \mathbf{W}^T\mathbf{R}\mathbf{W}\mathbf{D}^{-p}) \right. \\ &\quad \left. - (\mathbf{D}^{1-p}\mathbf{W}^T\mathbf{R}\mathbf{W}\mathbf{D}^{-1} + \mathbf{D}^{-1}\mathbf{W}^T\mathbf{R}\mathbf{W}\mathbf{D}^{1-p}) \right] + O(\beta^2), \end{aligned} \quad (14)$$

where we omit the iteration index  $n$  in the right hand side. It should be here noted that the signal is not necessarily stationary. In order that  $\mathbf{W}^T(n+1)\mathbf{W}(n+1)$  is still in the manifold  $S$  under the assumption that  $\beta^2$  is negligible small, it must hold that

$$\mathbf{D}^{-p}\mathbf{W}^T\mathbf{R}\mathbf{W} + \mathbf{W}^T\mathbf{R}\mathbf{W}\mathbf{D}^{-p} = \mathbf{D}^{1-p}\mathbf{W}^T\mathbf{R}\mathbf{W}\mathbf{D}^{-1} + \mathbf{D}^{-1}\mathbf{W}^T\mathbf{R}\mathbf{W}\mathbf{D}^{1-p}. \quad (15)$$

If  $\mathbf{W}^T\mathbf{R}\mathbf{W}$  is diagonal, the equality holds; however, the diagonality is guaranteed at the fixed points of  $\mathbf{W}$ . We have now the following result of the preservation of orthogonality:

**Theorem 3** *By the update with a small parameter  $\beta$ , it holds that*

$$\mathbf{W}^T(n+1)\mathbf{W}(n+1) = \mathbf{D}^{-1} + O(\beta^2), \quad (16)$$

*if and only if  $K = 1$  or  $p = 1$ .*

If we choose  $K = 1$ ,  $\mathbf{W}^T\mathbf{R}\mathbf{W}$  is always a scalar, and then the equality holds for arbitrary  $p$ . If  $K \neq 1$ , which is a more general case,  $p$  must be unity so that the equality holds.

If a signal is stationary, we can provide a precise discussion and a stronger result of the stability analysis on the manifold in this special case  $K = 1$  or  $p = 1$ .

**Theorem 4** *In the learning rules (10) and (11), the manifold is stable under perturbation when  $K = 1$  or  $p = 1$ .*

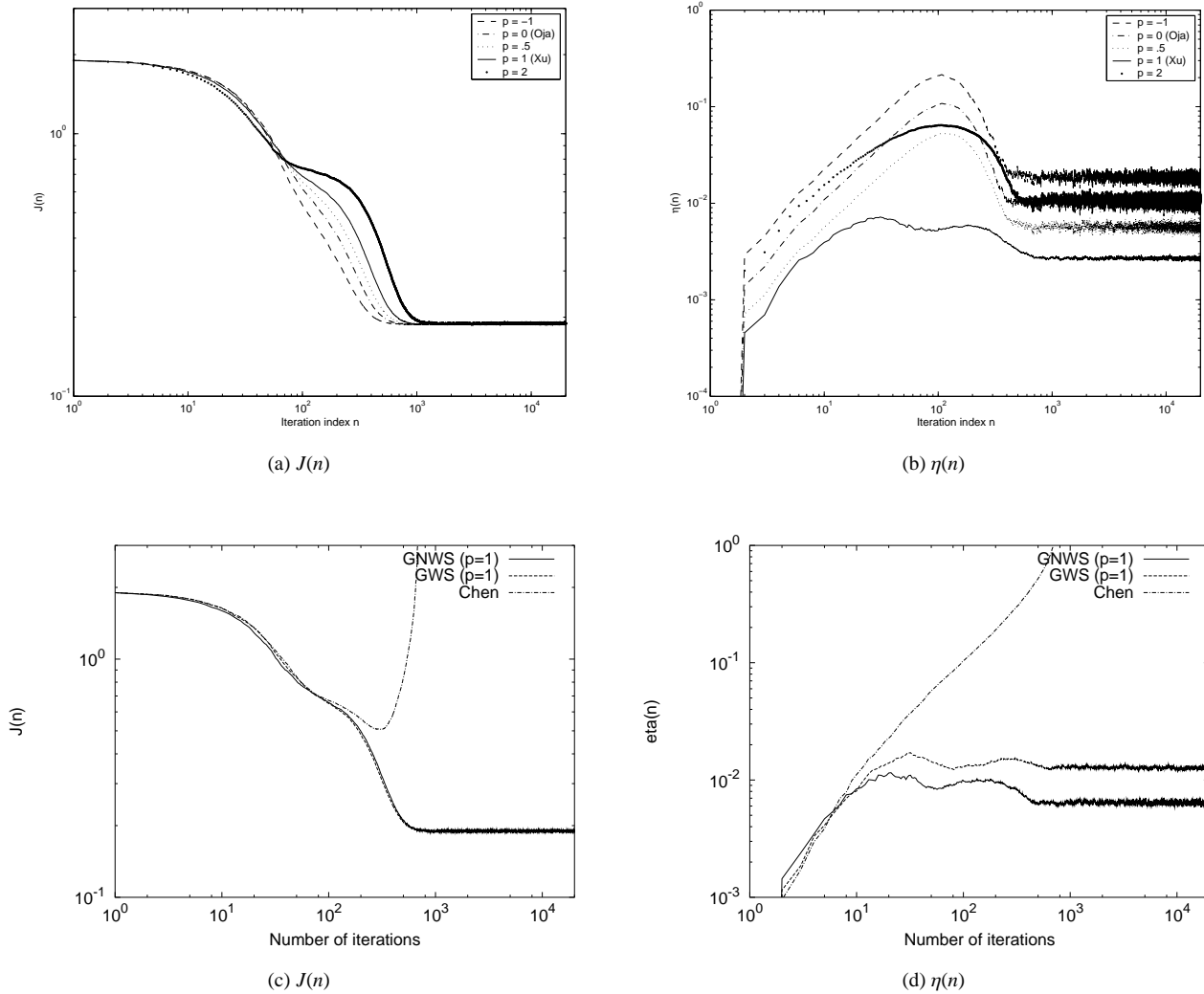
*Furthermore, the perturbation in (11) decays to zero faster than that in (10).*

## 3. NUMERICAL EXAMPLES

### 3.1. Random Vectors Generated from One Source

In this section, we show a simulation study of the proposed rules for a better understanding of the analytical results. We choose  $\mathbf{x}(n)$  to be a sequence of independent jointly-Gaussian random vectors with the correlation matrix

$$\mathbf{R} = \begin{bmatrix} 0.9 & 0.4 & 0.7 & 0.3 \\ 0.4 & 0.3 & 0.5 & 0.4 \\ 0.7 & 0.5 & 1.0 & 0.6 \\ 0.3 & 0.4 & 0.6 & 0.9 \end{bmatrix}. \quad (17)$$



**Fig. 1.** Average behaviors for different values of  $p$  under the GNWS rule ((a) and (b)). Average behaviors of the GWS rule ( $p = 1$ ), the GNWS rule ( $p = 1$ ), and the Chen's rule [8] ((c) and (d)).

We choose  $K = 2$ ,  $\beta = 0.01$ , and  $\mathbf{D} = \text{diag}[d_0, d_1] = \text{diag}[0.51, 1.01]$ . For evaluation of the algorithms, we compute the ensemble averages of the performance factors as follows:

$$J(n) = \text{tr}[(\mathbf{I} - \mathbf{P}(n)\mathbf{P}^T(n))\mathbf{R}(\mathbf{I} - \mathbf{P}(n)\mathbf{P}^T(n))^T], \quad (18)$$

$$\eta(n) = \|\mathbf{P}^T(n)\mathbf{P}(n) - \mathbf{I}\|_F^2 \quad (19)$$

for each algorithm using 100 independent runs, where  $\mathbf{P}(n) = \mathbf{W}(n)\mathbf{D}^{1/2}$ ,  $\|\cdot\|_F$  denotes the Frobenius norm,  $\text{tr}[\cdot]$  is the trace, and  $\mathbf{D}^{1/2}$  is a diagonal matrix with the entries  $\sqrt{d_0}, \dots, \sqrt{d_{K-1}}$ .

Figures 1(a) and 1(b) illustrate the average behaviors of the cost functions  $J(n)$  and  $\eta(n)$  for different values  $p$  under the GNWS rule, respectively. We examine the cases where  $p = -1, 0, 0.5, 1$ , and  $2$ . The GNWS rule when  $p = 0.5$  is identical to the Xu's weighted rule (8). The difference of behavior is observed both in Figs. 1(a) and 1(b). Interestingly, it can also be observed from Fig. 1(b) that the GNWS rule with  $p = 1$  preserves orthogonality more strongly than the GNWS rules with  $p \neq 1$ . Theorem 4 can

theoretically justify this fact.

Then, we show the difference between the GWS and the GNWS rules in Fig. 1(c) and 1(d). The behavior of the Chen's rule [8] is also depicted. It can be observed that the Chen's rule diverges rapidly. The difference of behavior between the GWS and the GNWS rules is small; however, it should be noted that we can observe the distinct difference in the behavior of  $\eta$  as seen in Fig. 1(d). The GNWS rule keeps orthogonality more strongly than the GWS rule. This observation coincides with the theoretical result of Theorem 4.

### 3.2. Backward Adaptation in the Presence of Quantization

We next show an application of those rules to a backward adaptive transform system [9]. Block diagram of transform coding system with backward adaptive transform updates is illustrated in Fig. 2. In this adaptive system, the transform matrix  $\mathbf{W}(n-1)$  at the  $(n-1)$ th iteration is updated with the quantized vector  $\hat{\mathbf{y}}(n)$ ; therefore,

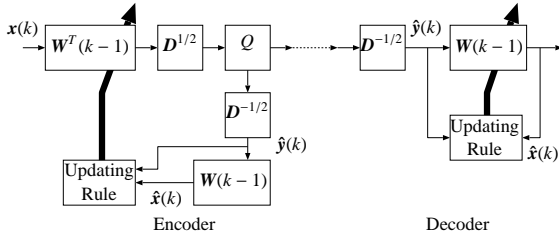


Fig. 2. Backward adaptive transform system.

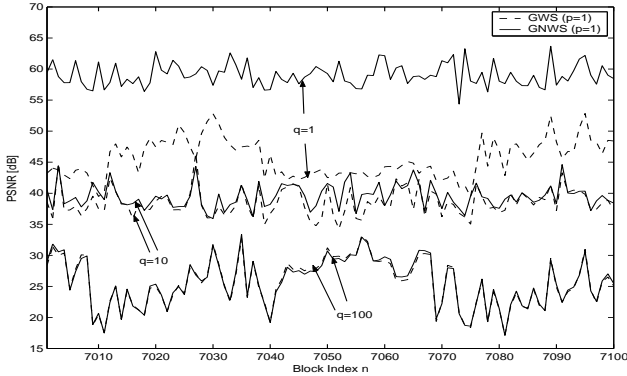


Fig. 3. PSNR at  $n$ th blocks in dB for “Lena.”

the decoder can update the transform matrix in the same way as the encoder does. Specifically, the quantized vector  $\hat{y}(n)$  is obtained by the quantization operator (non-linear) defined as

$$\hat{y} = \mathbf{D}^{-1/2} Q[\mathbf{D}^{1/2} \mathbf{y}] = q \mathbf{D}^{-p/2} \text{round}(1/q \mathbf{D}^{1/2} \mathbf{y}), \quad (20)$$

where  $\text{round}(\cdot)$  is the rounding operator for vector elements and  $q$  is the step size of quantization. The update in (10) and (11) is performed by  $\hat{y}(n)$  and  $\hat{x}(n) = \mathbf{W}(n-1)\hat{y}(n)$  instead of  $\mathbf{y}(n)$  and  $\mathbf{x}(n)$ , respectively.

We compare two rules in the backward adaptation system. For simplicity, we examine the one-dimensional case where an input image is raster-scanned horizontally and segmented into blocks of  $M$  samples. Figure 3 depicts the comparison of PSNRs for the test image “Lena,” which is of size  $512 \times 512$ . In this figure, PSNRs of 7001st to 7100th blocks, i.e.,  $\text{PSNR}(n) = -20 \log_{10} \frac{\|\mathbf{x}(n) - \hat{\mathbf{x}}(n)\|}{255 \times M}$  for  $7001 \leq n \leq 7100$ , are plotted. It is clearly observed that the GNWS provides better performance in PSNR than the GWS for  $q = 1$  and  $q = 10$ . The difference is significant in the case  $q = 1$ . To understand the reason for this effect, we illustrate the plots of  $\eta(n)$  for  $1 \leq n \leq 9000$  and  $q = 1$  in Fig. 4, which shows that the GWS is less orthogonal than the GNWS during the iteration. This fact may result in the difference of PSNRs between two rules.

#### 4. CONCLUSIONS

In this paper, the framework for weighted subspace (WS) rules for PCA have been introduced. The proposed form includes the conventional WS rules. The stability of the algorithms have been

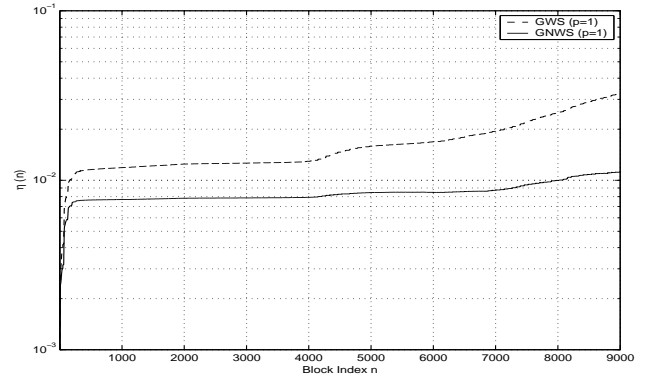


Fig. 4. Stability of orthogonality in the “Lena” image.

analyzed. The special case of the proposed form such that orthogonality of the estimated matrix is preserved has been found. We conclude that the GNWS rule with the parameter  $p = 1$  (Xu’s normalized rule) is the most promising among the rules dealt with in the paper. Future work may involve the development of fast and orthogonal algorithms as done in [5, 10]

#### 5. REFERENCES

- [1] K. I. Diamantaras and S. Y. Kung, *Principal Component Neural Networks: Theory and Applications*. New York: Wiley, 1996.
- [2] E. Oja, “Neural networks, principal components, and subspaces,” *Int. J. Neural Systems*, vol. 1, pp. 61–68, Apr. 1989.
- [3] B. Yang, “Projection approximation subspace tracking,” *IEEE Trans. Signal Processing*, vol. 43, pp. 95–107, Jan. 1995.
- [4] L. Xu, “Least mean square error reconstruction principle for self-organizing neural-nets,” *Neural Networks*, vol. 6, pp. 627–648, 1993.
- [5] S. Attallah and K. Abed-Meraim, “Fast algorithms for subspace tracking,” *IEEE Signal Processing Letters*, vol. 8, pp. 203–206, July 2001.
- [6] E. Oja, H. Ogawa, and J. Wangviwattana, “Principal component analysis by homogeneous neural networks, Part I: The weighted subspace criterion,” *IEICE Trans. Inf. & Syst.*, vol. E75-D, pp. 366–375, May 1992.
- [7] T. Tanaka, “Generalized weighted rules for principal components tracking,” *to appear in IEEE Trans. Signal Processing*, 2004.
- [8] T. Chen, S. Amari, and Q. Lin, “A unified algorithm for principal and minor components extraction,” *Neural Networks*, vol. 11, no. 3, pp. 385–390, 1998.
- [9] V. K. Goyal, J. Zhuang, and M. Vetterli, “Transform coding with backward adaptive updates,” *IEEE Trans. Information Theory*, vol. 46, pp. 1623–1633, Feb. 1997.
- [10] K. Abed-Meraim, S. Attallah, A. Chkeif, and Y. Hua, “Orthogonal Oja algorithm,” *IEEE Signal Processing Letters*, vol. 7, pp. 116–119, May 2000.