# One-parameter family of nonlinear dynamics
# for online linear discriminant analysis

Kazuyuki Hiraoka, Ken-ichi Hidai,
Hiroshi Mizoguchi, Taketoshi Mishima, and Shuji Yoshizawa

*Dept. of Information and Computer Sciences, Saitama University, 255 Shimo-okubo, Urawa, 338-8570, Japan*
email: hira@me.ics.saitama-u.ac.jp

*Abstract*— A novel variation of online linear discriminant analysis (OLDA) is proposed based on an one-parameter family of nonlinear dynamics. Both previously proposed dynamics for OLDA and Oja's dynamics for principal component analysis are special cases in this family. With regard to this family, fixed-points and their stability are analyzed. The dependence of discrimination performance on the parameter is also studied experimentally.

## I. Introduction

Linear discriminant analysis (LDA) has been applied extensively, e.g. for face identification [8][9]. However, LDA is poor at adaptivity. This is because LDA is a batch learning algorithm. Indeed, we have to re-calculate the discrimination matrix $A$ every time when we add new data to update the identification system. This calculation is heavy for high dimensional data such as face images.

Recently, the authors have proposed an *online* version of LDA [11], which is referred to online LDA (OLDA). By OLDA, the face identification system can be updated with low computational cost when new additional images are presented. Hence OLDA has the ability of adaptation to the change of environment. OLDA also has an advantage that huge matrices never appear in its calculation.

Though iterative algorithms have been proposed[1] for neural network based LDA [6][5], they are not sufficiently "online". Since those algorithms keep $n^2 \times n^2$ matrices when the image size is $n \times n$, they require $O(n^4)$ time for one step updating and $O(n^4)$ memory. In contrast to [6][5], our OLDA algorithm requires only $O(n^2ML)$ time for one step updating and $O(n^2M)$ memory, where $M$ and $L$ are the number of classes and features, respectively. Note that $n^2 \gg M > L$ in typical cases of face identification tasks.

In the present paper, we will propose a novel OLDA algorithm based on an one-parameter family of matrix

---

[1]The original algorithm in [5] is a batch learning. However, we can easily modify it to an iterative learning by replacing the calculation of the standard deviation in [5] with an iterative one.

dynamics. The first algorithm [11] is a special case of it. Moreover, Oja's dynamics for online principal component analysis (PCA) can be viewed as another special case of it.

## II. One-parameter family of matrix dynamics

To determine the discrimination matrix $A$ in LDA, we have to solve a generalized eigenvalue problem. Let $B$ and $W$ be the $N \times N$ variance matrices "between classes" and "within classes" respectively. Then, we have to find an $N \times L$ matrix $A_{sol}$ and an $L \times L$ symmetric matrix $\Gamma$ which satisfy

$$BA_{sol} = WA_{sol}\Gamma,$$
$$A_{sol}^T W A_{sol} = I, \qquad (1)$$

where $I$ is the identity matrix.

In order to obtain the solution $A_{sol}$ iteratively, the following matrix dynamics is discussed in the present study:

$$\frac{d}{dt}A(t) = BA(t) - \alpha BA(t)A(t)^T WA(t)$$
$$- (1-\alpha)WA(t)A(t)^T BA(t), \quad (2)$$

where $\alpha \in \mathbf{R}$ is a parameter.

## III. Main theorems

With regard to the one-parameter family of matrix dynamics (2), the following theorems are proved.

**Theorem 1** *All the solutions of (1) are fixedpoints of (2).*

**Theorem 2** *Suppose that rankB $\geq L$ and $\alpha = 0$. Then all the stable fixedpoints of (2) are solutions of (1).*

**Theorem 3** *Suppose that rankB $\geq L$ and $\alpha > 0$, $\alpha \neq 1$. If $A$ is a stable fixedpoint of (2) and $A$ is not a solutions of (1), the eigenvalues $\omega_1, \cdots, \omega_L$ of $A^T WA$ must satisfy $\omega_k = 1$ or $\omega_k \geq 1/\alpha$ ($k = 1, \cdots, L$).*

Theorem 1 suggests that the dynamics (2) can be used to solve the generalized eigenvalue problem (1). Theorem 2 shows that there is no "spurious" solution when $\alpha = 0$. Theorem 3 shows that the spurious solutions exists only in a region which is "far" from $O$ compared with the "true" solutions, when $0 < \alpha < 1$. This is because the true solutions satisfy $A^T W A = I$.

Thus, as for spurious solutions, smaller $\alpha$ is better. Especially, $\alpha = 0$ is a direct extension of Oja's dynamics for online PCA [1] and this dynamics has no spurious solution.

On the other hand, as for convergence speed, it is experimentally observed that positive $\alpha$ is superior. This will be discussed later.

### IV. Derived OLDA algorithm

From the dynamics (2), following OLDA algorithm is derived.

At every time step $t = 1, 2, 3, \cdots$, a new pair $(x(t), c(t))$ is presented, where $x(t)$ is an $N$-dimensional data vector, $c(t) \in \{1, \cdots, M\}$ is the class of $x(t)$, and $M$ is the number of classes. Based on this pair, auxiliary variables and the $N \times L$ discrimination matrix $A$ are updated as follows:

$$t^c(t) = t^c(t-1) + \delta(c, c(t)), \qquad (3)$$

$$\bar{x}(t) = \left(1 - \frac{1}{t}\right) \bar{x}(t-1) + \frac{1}{t} x(t), \qquad (4)$$

$$\bar{x}^c(t) =$$
$$\begin{cases} \left(1 - \frac{1}{t^c(t)}\right) \bar{x}^c(t-1) + \frac{1}{t^c(t)} x(t) & (c = c(t)), \\ \bar{x}^c(t-1) & (c \neq c(t)), \end{cases} \quad (5)$$

$$v^c(t) = \bar{x}^c(t) - \bar{x}(t), \qquad (6)$$

$$w(t) = x(t) - \bar{x}^{c(t)}(t), \qquad (7)$$

$$y^c(t) = A(t-1)^T v^c(t), \qquad (8)$$

$$z(t) = A(t-1)^T w(t), \qquad (9)$$

$$F(t) = \frac{1}{M} \sum_{c=1}^{M} v^c(t) y^c(t)^T, \qquad (10)$$

$$g(t) = \frac{1}{M} \sum_{c=1}^{M} y^c(t) \left(y^c(t)^T z(t)\right), \qquad (11)$$

$$A(t) = A(t-1) + \eta \Big[ F(t) - \alpha F(t) z(t) z(t)^T$$
$$- (1-\alpha) w(t) g(t)^T - \alpha \epsilon F(t) \left(A(t-1)^T A(t-1)\right)$$
$$- (1-\alpha)\epsilon A(t-1) \left(A(t-1)^T F(t)\right) \Big], \qquad (12)$$

where the parameter $\alpha$ is corresponding to $\alpha$ in (2), the learning coefficient $\eta$ is a small positive number, and the regularization coefficient $\epsilon$ is also a small positive number. The number $L$ of features is less than or

equal to $\min(N, M-1)$. Note that the covariance matrix $W(t) = \frac{1}{t} \sum_{\tau=1}^{t} w(\tau) w(\tau)^T$ within classes is replaced with the instantaneous value $w(t) w(t)^T$. Such replacement is justified by the theory of stochastic approximation [7]. As for the initial values, $t^c(0) = 0$, $\bar{x}(0)$ and $\bar{x}^c(0)$ are arbitrary vectors, and $A(0)$ is an arbitrary matrix which satisfies $\text{rank} A(0) = L$.

### V. Experimental results

We have applied the above OLDA algorithm to face image identification for 4 persons. The performance is evaluated for test samples after $t = 400$ steps of iterations (Fig. 1).
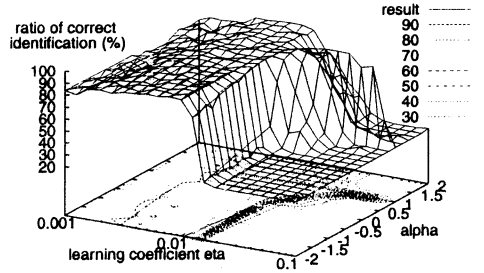


Figure 1: Performance of the proposed method for each $\eta$ and $\alpha$. Stability is improved around $\alpha = 1$.

It is experimentally shown that stability for large learning coefficient is improved around $\alpha = 1$. Moreover, the problem of spurious solutions turns out to be negligible if we take an initial matrix $A(0)$ near the zero matrix $O$.

### VI. Discussion

The reason of the slow convergence for $\alpha \leq 0$ is not clearly understood at now. This point must be studied further. A key will be the fact that $w(t) g(t)^T$ has much larger variance than $F(t) z(t) z(t)^T$ in (12).

Another unclear point of the present study is the question whether the dynamics (2) can be derived from potential functions. The authors have only partial results. For the case $\alpha = 1/2$, (2) is derived as

$$da_{ij}/dt = (1/2)(\partial \phi(A)/\partial a_{ij}), \qquad (13)$$

$$\phi(A) = \text{Tr}\left[A^T B A \left(I - \frac{1}{2} A^T W A\right)\right]. \qquad (14)$$

For the case $\alpha = 0$, (2) is derived similarly by another potential function

$$\psi(A) = \log \det[A^T B A] - \text{Tr}[A^T W A] \qquad (15)$$

together with a metric

$$\langle \Delta A_1, \Delta A_2 \rangle_A = \text{Tr}[\Delta A_1^T \Delta A_2 (A^T B A)^{-1}]. \quad (16)$$

Note that Oja's dynamics [1] is a special case of (2) with $\alpha = 0$ and $W = I$.

## VII. Conclusion

A novel variation of online linear discriminant analysis (OLDA) is proposed based on an one-parameter family of nonlinear dynamics. Both previously proposed dynamics for OLDA and Oja's dynamics for principal component analysis are special cases in this family.

It is theoretically shown that (a) when the parameter $\alpha = 0$, no "spurious solution" exists, and (b) when $0 < \alpha < 1$, spurious solutions exists only in a region which is "far" from $O$ compared with the "true" solutions. On the other hand, it is experimentally shown that stability for large learning coefficient is improved around $\alpha = 1$.

### Acknowledgements

### References

[1] E. Oja, "A simplified neuron model as a principal component analyzer", J. Math. Biol., Vol. 15, pp. 267–273, 1982.

[2] T. D. Sanger, "Optimal unsupervised learning in a single-layer linear feedforward neural network", Neural Networks, Vol. 2, No. 6, pp. 459–473, 1989.

[3] K. Hornik and C.-M. Kuan, "Convergence analysis of local feature extraction algorithms", Neural Networks, Vol. 5, No. 2, pp. 229–240, 1992.

[4] W.-Y. Yan, U. Helmke, and J. B. Moore, "Global analysis of Oja's flow for neural networks", IEEE Tr. NN, Vol. 5, No. 5, pp. 674–683, 1994.

[5] J. Mao and A. K. Jain, "Artificial neural networks for feature extraction and multivariate data projection", IEEE Tr. NN, Vol. 6, No. 2, pp. 296–317, 1995.

[6] C. Chatterjee and V. P. Roychowdhury, "On self-organizing algorithms and networks for class-separability features", IEEE Tr. NN, Vol. 8, No. 3, pp. 663–678, 1997.

[7] M. B. Nevel'son and R. Z. Has'minskiǐ, Stochastic Approximation and Sequential Estimation, Nauka, 1968 (in Russian; translated into Japanese by T. Kitagawa and K. Tajima, 1983).

[8] T. Kurita and S. Hayamizu, "Gesture Recognition using HLAC Features of PARCOR Images and HMM based Recognizer", Proc. 3rd int. conference on automatic face and gesture recognition, April, 1998, pp. 422-427.

[9] W. Zhao, et al., "Discriminant Analysis of Principal Components for Face Recognition", Proc. 3rd int. conference on automatic face and gesture recognition, 1998, pp. 422-427. pp. 336-341.

[10] K. Fukunaga, Statistical Pattern Recognition, Academic Press, New York, 1989.

[11] K. Hiraoka and M. Hamahira, "On Successive Learning Type Algorithm for Linear Discriminant Analysis", IEICE Technical Report, NC99-73, pp. 85–92, 1999 (in Japanese).

[12] Kazuyuki Hiraoka, Shuji Yoshizawa, Ken-ichi Hidai, Masashi Hamahira, Hiroshi Mizoguchi, and Taketoshi Mishima, "Convergence Analysis of Online Linear Discriminant Analysis", International Joint Conference on Neural Networks (IJCNN), 2000, to appear.

### Appendix: Proofs of theorems

**Proof of theorem 1:** Trivial. ∎

**Lemma 1** Let $J_B = A^T B A$ and $J_W = A^T W A$. Then $J_B J_W = J_W J_B = J_B$ if $A$ is a fixedpoint of (2).

**Proof:** Multiplying $A^T$ from left of (2), we obtain

$$\alpha J_B(I - J_W) + (1-\alpha)(I - J_W)J_B = O. \quad (17)$$

Then, by adding (17) and the transpose of (17), we obtain

$$J_B(I - J_W) + (I - J_W)J_B = O. \quad (18)$$

This means $J_B(I - J_W) = (I - J_W)J_B = O$, because both $J_B$ and $I - J_W$ are symmetric. Thus $J_B J_W = J_W J_B = J_B$. ∎

**Corollary 1** By rotating the coordinate system in $\mathbb{R}^L$, we can assume

$$J_B = diag(\beta_1, \cdots, \beta_L), \quad (19)$$
$$J_W = diag(\omega_1, \cdots, \omega_L), \quad (20)$$
$$\beta_1 \geq \cdots \geq \beta_K > 0, \quad (21)$$
$$\beta_{K+1} = \cdots = \beta_L = 0, \quad (22)$$
$$\omega_1 = \cdots = \omega_K = 1, \quad (23)$$
$$\omega_{K+1} \geq \cdots \geq \omega_L \geq 0, \quad (24)$$

without loss of generality.

**Lemma 2** *The following conditions are equivalent when $\alpha \neq 1$:*

1. *$A$ is a fixedpoint of (2).*

2. *$A$ satisfies*

$$BA = WAJ_B, \tag{25}$$

$$J_B = J_W J_B = J_B J_W. \tag{26}$$

**Proof:** ($\Leftarrow$) Trivial. ($\Rightarrow$) Let $A = (a_1, \cdots, a_L)$. When $A$ is a fixedpoint, lemma 1 and its corollary hold. Then, from (2), $(BA - WAJ_B)(I - \alpha J_W) = O$. Thus, $Ba_k = \beta_k W a_k$ if $\omega_k \neq 1/\alpha$. Moreover, $Ba_k = \beta_k W a_k$ also holds if $\omega_k = 1/\alpha \neq 1$, because $\beta_k = 0$ in this case. Note that $Ba_k = 0$ and $\beta_k = a_k^T B a_k = 0$ are equivalent since $B$ is symmetric. ∎

**Proof of theorem 2 and 3:** Let $\breve{A}$ be a "spurious" solution, namely, $\breve{A}$ is a fixedpoint of (2) and $\breve{A}$ is not a solution of (1). Then lemma 1 and its corollary hold. Now we can take a basis $\{e_1, \cdots, e_N\}$ which satisfies $Be_k = \beta_k We_k$, $\breve{A} = (e_1, \cdots, e_L)\text{diag}(\sqrt{\omega_1}, \cdots, \sqrt{\omega_L})$, and $e_k^T We_{k'} = 1(k = k'), 0(k \neq k')$. Suppose that $\omega_1 = \cdots = \omega_K = 1$ and $\omega_{K+1}, \cdots, \omega_L \neq 1$. For simplicity of description, we only show here a discussion for the case that $\beta_1, \cdots, \beta_N$ are different each other except that $\beta_{K+1} = \cdots \beta_L = 0$.

Let $\mathcal{A} = \{\breve{A}C \mid C \text{ is an arbitrary } L \times L \text{ matrix}\}$. Note that there exists a neighborhood $U$ of $\breve{A}$ such that all the fixedpoints in $U$ belong to $\mathcal{A}$. In order to observe the distance between $A(t)$ and $\mathcal{A}$, we define the metric $\langle F, G \rangle \equiv \text{Tr} F^T G$.

Let $\Delta A(t) = A(t) - \breve{A}$ and $\Delta A(0) = (0, \cdots, 0, \epsilon e_H)$, where $H$ satisfies $L + 1 \leq H \leq N$ and $\beta_H > 0$. The existence of such $H$ is guaranteed from assumptions. Since $\Delta \breve{A}^T W A(0) = O$, $\Delta \dot{A} \equiv d\Delta A/dt|_{t=0} = \beta_H(1 - \alpha \omega_L) W \Delta A(0) + o(\epsilon)$. Hence $\langle \breve{A}, \Delta \dot{A} \rangle \approx 0$ for all $\breve{A} \in \mathcal{A}$, namely, $\Delta \dot{A}$ is orthogonal to $\mathcal{A}$.

Moreover, $d\langle \Delta A(t), \Delta A(t) \rangle/dt|_{t=0} \approx \beta_H(1 - \alpha \omega_L)\epsilon^2 > 0$ for sufficiently small $\epsilon$ if $\alpha = 0$ or $\omega_L < 1/\alpha$. Thus the theorems are proved. ∎