

# Fast Algorithm for Updating Discriminant Functions in Linear Discriminant Analysis

Kyi Min Htut, Hironori Tamaki, Atsushi Nakajima and Takaomi Shigehara  
Department of Information and Computer Sciences,  
Faculty of Engineering, Saitama University,  
255 Shimo-Okubo, Saitama, Saitama 338-8570, Japan  
Tel.: +81-48-858-9035, Fax: +81-48-858-3716  
E-mail: sigehara@ics.saitama-u.ac.jp

**Abstract:** We propose a new numerical algorithm for linear discriminant analysis which makes it possible to *update* the discriminant functions with very low computational cost.

## 1. Introduction

To clarify the point of the issue, we first give a brief summary for the framework of linear discriminant analysis (LDA). Discriminant analysis is a standard method for statistical pattern recognition [1]. In pattern recognition, we use the term *pattern* to denote the  $n$ -dimensional data vector  $\mathbf{x} \in \mathbf{R}^n$  whose components are measurements of features of an object under consideration. The main purpose in pattern recognition is to determine the *class*  $c \in \Gamma \equiv \{1, 2, \dots, M\}$  to which a given pattern vector  $\mathbf{x}$  belongs, where  $M$  is the number of classes. The problem is reduced to finding the *discriminant functions*  $g_c(\mathbf{x})$ ,  $c \in \Gamma$ , such that

$$g_c(\mathbf{x}) = \min_{c' \in \Gamma} g_{c'}(\mathbf{x}) \implies \mathbf{x} \text{ belongs to the class } c. \quad (1)$$

In LDA, the discriminant functions are determined by the following procedure. For given  $d$  sample data of pattern  $\mathbf{x}(t) \in \mathbf{R}^n$  and class  $c(t) \in \Gamma$  to which  $\mathbf{x}(t)$  belongs ( $t = 1, 2, \dots, d$ ), we first solve the *generalized eigenvalue problem* (GEP) of  $n$ -th order;

$$A_d \mathbf{p} = \lambda G_d \mathbf{p}, \quad (2)$$

where

$$A_d = \frac{1}{M} \sum_{c=1}^M (\bar{\mathbf{x}}_d^c - \bar{\mathbf{x}}_d) (\bar{\mathbf{x}}_d^c - \bar{\mathbf{x}}_d)^T \quad (3)$$

and

$$G_d = \frac{1}{d} \sum_{t=1}^d (\mathbf{x}(t) - \bar{\mathbf{x}}_d^{c(t)}) (\mathbf{x}(t) - \bar{\mathbf{x}}_d^{c(t)})^T \quad (4)$$

are the between-class and within-class variance-covariance matrices, respectively. Here

$$\bar{\mathbf{x}}_d = \frac{1}{d} \sum_{t=1}^d \mathbf{x}(t) \quad \text{and} \quad \bar{\mathbf{x}}_d^c = \frac{\sum_{t=1}^d \delta(c, c(t)) \mathbf{x}(t)}{\sum_{t=1}^d \delta(c, c(t))} \quad (5)$$

are the mean vector of the whole sample patterns and the mean vector of the sample patterns which belong to the class  $c$ , respectively. (In Eq.(5),  $\delta(c, c') = 1$  for  $c = c'$  and 0 for  $c \neq c'$ .) It is important to notice that  $M \ll n \ll d$  holds in a generic case; For face identification for example,  $M \simeq 10$  is the number of persons to be discriminated,  $n \simeq 10^2 \sim 10^3$  is the number of pixels of image and  $d \simeq 10^3 \sim 10^4$  is the number of the sample image data. The GEP (2) has at most  $M - 1$  nonzero eigenvalues and there indeed exist  $M - 1$  nonzero eigenvalues  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{M-1} > 0$  in a generic case. Let  $\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_{M-1}$  be the corresponding normalized eigenvectors;  $\mathbf{p}_i^T G_d \mathbf{p}_j = \delta(i, j)$ . Then the discriminant functions in LDA are given by

$$g_c(\mathbf{x}) = |P^T (\mathbf{x} - \bar{\mathbf{x}}_d^c)|^2 \quad (6)$$

with  $P = (\mathbf{p}_1 \mathbf{p}_2 \dots \mathbf{p}_{M-1})$ . Obviously, the most expensive part for determining the discriminant functions (6) is in solving the GEP (2) and it is highly desired to develop an efficient numerical algorithm for the GEP (2). This is exactly our main subject. In this paper, we concentrate on *updating* the discriminant functions when some new sample data are added. In most applications of discriminant analysis, it is very hard to construct the *good* discriminant functions which keep a high rate for correct identification from the sample data given at first and it is necessary to reconstruct the discriminant functions by adding some new sample data. For this purpose, an iterative method based on a nonlinear matrix dynamical system has been proposed in [2]. In this paper, we propose an alternative new numerical algorithm. Our method is based on a direct method and makes it possible to find the *exact* updated discriminant functions only with  $O(n^2)$  arithmetic operations, contrary to with  $O(n^3)$  arithmetic operations which are required to solve the GEP (2) in a direct manner.

## 2. Main Results

The standard numerical algorithms for the GEP with a generic constraint such that  $A_d$  is symmetric and  $G_d$  is positive definite symmetric in Eq.(2) follow the two steps:

- (I) Calculate the Cholesky decomposition  $G_d$ .
- (II) By using the Cholesky decomposition  $G_d$ , make the problem reduce to a symmetric eigenvalue problem of  $n$ -th order and solve it.

Both of the steps (I) and (II) require  $O(n^3)$  arithmetic operations. The first step of our strategy is to make the symmetric eigenvalue problem of  $n$ -th order in the step (II) reduce to a mathematically equivalent symmetric eigenvalue problem of  $M$ -th order by using a specific feature of the between-class variance-covariance matrix  $A_d$  in Eq.(3). Indeed, we can show the following proposition;

**Proposition 1** Assume that the Cholesky decomposition of  $G_d$  is given;  $G_d = \Lambda_d \Lambda_d^T$  with an  $n$ -dimensional lower triangular matrix  $\Lambda_d$ . Then the GEP (2) with Eqs.(3) and (4) reduces to a symmetric eigenvalue problem of  $M$ -th order which can be solved with  $O(n^2)$  arithmetic operations.

**Proof** Let us define an  $n \times M$  matrix  $U_d$  by

$$U_d = \frac{1}{\sqrt{M}} (\bar{x}_d^1 - \bar{x}_d \quad \bar{x}_d^2 - \bar{x}_d \quad \cdots \quad \bar{x}_d^M - \bar{x}_d). \quad (7)$$

Then the between-class variance-covariance matrix can be written as

$$A_d = U_d U_d^T. \quad (8)$$

By multiplying  $U_d^T G_d^{-1} = (\Lambda_d^{-1} U_d)^T \Lambda_d^{-1}$  from the left on both sides, Eq.(2) reduces to an  $M$ -th order eigenvalue equation

$$\tilde{A}_d \mathbf{q} = \lambda \mathbf{q}, \quad (9)$$

where

$$\mathbf{q} = U_d^T \mathbf{p} \in \mathbf{R}^M \quad (10)$$

and the  $M$ -dimensional matrix  $\tilde{A}_d$  is defined by

$$\tilde{A}_d = \tilde{U}_d^T \tilde{U}_d \quad (11)$$

with

$$\tilde{U}_d = \Lambda_d^{-1} U_d. \quad (12)$$

A theorem of singular value decomposition (with a metric matrix  $G$ ) indicates that the nonzero eigenvalues of Eq.(9) coincide with those of Eq.(2). After finding a nonzero eigenvalue  $\lambda \neq 0$  and the associated eigenvector  $\mathbf{q}$  in Eq.(9), the corresponding eigenvector  $\mathbf{p}$  in Eq.(2) is obtained by the solution of a linear system

$$G_d \mathbf{p} = U_d \mathbf{q}. \quad (13)$$

Since  $n \gg M$ , the reduction of Eq.(2) to Eq.(9) brings about a substantial speed up for numerical computation. If the Cholesky decomposition of  $G_d$  is available, we can solve the eigenvalue problem (9) with  $O(n^2)$  arithmetic operations. Note that  $\tilde{U}_d$  in Eq.(12) can be determined by solving  $M$  linear systems with a common  $n$ -dimensional lower triangular coefficient matrix;  $\Lambda_d \tilde{U}_d = U_d$ . Also Eq.(13) is solved with  $O(n^2)$  arithmetic operations if the Cholesky decomposition of  $G_d$  is completed.

The second step of our strategy is to update the Cholesky decomposition  $G_d = \Lambda_d \Lambda_d^T$  with a cheap numerical cost when a new sample data of pattern  $\mathbf{x}(d+1) \in \mathbf{R}^n$  and class  $c(d+1) \in \Gamma$  to which  $\mathbf{x}(d+1)$  belongs is added. This is indeed possible by the following proposition;

**Proposition 2** Let  $G_d = \Lambda_d \Lambda_d^T$  be the Cholesky decomposition of  $G_d$ . Assume that a new sample data of pattern  $\mathbf{x}(d+1) \in \mathbf{R}^n$  and class  $c(d+1) \in \Gamma$  is added. Then the Cholesky decomposition of the new within-class variance-covariance matrix

$$G_{d+1} = \frac{1}{d+1} \sum_{t=1}^{d+1} (\mathbf{x}(t) - \bar{\mathbf{x}}_{d+1}^{c(t)}) (\mathbf{x}(t) - \bar{\mathbf{x}}_{d+1}^{c(t)})^T \quad (14)$$

can be obtained with  $O(n^2)$  arithmetic operations.

**Proof** After a somewhat lengthy but straightforward calculation, we reach

$$\begin{aligned} G_{d+1} &= \frac{d}{d+1} G_d \\ &+ \frac{1}{d+1} (\mathbf{x}(d+1) - \bar{\mathbf{x}}_d^{c(d+1)}) \\ &\quad \times (\mathbf{x}(d+1) - \bar{\mathbf{x}}_d^{c(d+1)})^T \\ &- \frac{\sum_{t=1}^{d+1} \delta(c(d+1), c(t))}{d+1} (\bar{\mathbf{x}}_{d+1}^{c(d+1)} - \bar{\mathbf{x}}_d^{c(d+1)}) \\ &\quad \times (\bar{\mathbf{x}}_{d+1}^{c(d+1)} - \bar{\mathbf{x}}_d^{c(d+1)})^T. \end{aligned} \quad (15)$$

Eq.(15) shows that  $G_{d+1}$  is obtained by adding a rank-one perturbation to  $G_d$  and subtracting a rank-one perturbation successively. Thus we can use the technique of the *Cholesky updating/downdating* [3] (see Appendix for details), which makes it possible to obtain the Cholesky decomposition of  $G_{d+1}$  from  $\Lambda_d$  with  $O(n^2)$  arithmetic operations. ■

Since the updating of the between-class variance-covariance matrix  $A_d$  (actually the updating of  $U_d$  in Eq.(7)) is carried out with  $O(n)$  arithmetic operations, we can update the discriminant functions only with  $O(n^2)$  arithmetic operations by combining the methods described in the proof of Propositions 1 and 2.

### 3. Numerical Experiment

Numerical experiment has been performed on Sun microsystems workstation (OS: Solaris 2.6, CPU: Micro Sparc 100 MHz×2, Main Memory: 128 MB, Compiler: g77 ver.2.95.2 with option '-O3'). In Table 1, we show the execution time for updating the Cholesky decomposition of the within-class variance-covariance matrix  $G_d$  when a new sample data of pattern  $\mathbf{x}(d+1) \in \mathbf{R}^n$  and class  $c(d+1) \in \Gamma$  is added. The second line (PM) is the execution time for the proposed method in Proposition 2, while the third line (CM) is the execution time for the conventional method. In the latter case, an additional  $O(n^3)$  procedure (matrix product and Cholesky decomposition) is required. We has utilized `dpotrf` routine

Table 1. Execution time for updating Cholesky decomposition of  $G_d$  for  $d = 1000$ ,  $M = 10$ . PM and CM are for proposed and conventional methods, respectively.

data dim. $n$	100	300	500	700	900
PM (sec)	0.01	0.05	0.13	0.28	0.41
CM (sec)	0.94	8.93	25.68	52.87	89.93
speed up	94	178	197	188	219

in LAPACK ver.3.0 for the Cholesky decomposition of the new within-class variance-covariance matrix  $G_{d+1}$ . Table 1 shows that the proposed method is quite satisfactory; Compared to the standard  $O(n^3)$  procedure, the speed up for updating the within-class variance-covariance matrix  $G_d$  by the proposed method is about 100 ~ 200 for a wide range of the matrix size covering  $n \simeq 100 \sim 900$ . Together with the method in Proposition 1, we expect to update the discriminant functions within at most a few seconds when a new sample data is added.

One of the authors (T.S.) thanks Dr. Kazuyuki Hiraoka for fruitful discussions and helpful comments.

## References

- [1] A. Webb, "Statistical Pattern Recognition", Arnold, London, 1999.
- [2] K. Hiraoka, M. Hamahira, K. Hidai, H. Mizoguchi, T. Mishima and S. Yoshizawa, "Fast Algorithm for Online Linear Discriminant Analysis", Proceedings of the 2000 International Technical Conference on Circuits/Systems, Computers, and Communications, pp. 274-277, 2000; IEICE Transactions on Fundamentals of Electronics, Communications and Computer Sciences, Vol. E84-A, pp. 1431-1434, 2001.
- [3] G. W. Stewart, "Matrix Algorithms - Volume 1: Basic Decompositions", SIAM, Philadelphia, 1998.

## Appendix Cholesky Updating/Downdating

Let  $G$  be an  $n \times n$  positive definite real symmetric matrix and let

$$G = \Lambda \Lambda^T \quad (16)$$

be the Cholesky decomposition of  $G$ , where  $\Lambda = (l_1 l_2 \cdots l_n)$  is a lower triangular matrix. We consider how the Cholesky decomposition (16) of  $G$  is perturbed if a rank-one perturbation  $\mathbf{x}\mathbf{x}^T$  ( $\mathbf{x} \in \mathbf{R}^n$ ) is added to  $G$ , or subtracted from  $G$ . We will see that the new Cholesky decomposition after the rank-one perturbation can be calculated with  $O(n^2)$  arithmetic operations by using  $\Lambda$  and  $\mathbf{x}$  in a direct manner.

### A1 Cholesky Updating

We first consider the case that a rank-one perturbation is added to  $G$ . Let us define

$$\tilde{G} \equiv G + \mathbf{x}\mathbf{x}^T. \quad (17)$$

Note that  $G$  is symmetric and still positive definite, and hence  $\tilde{G}$  has the Cholesky decomposition.

A first note is that Eq.(17) is rewritten as

$$\tilde{G} = \Lambda \Lambda^T + \mathbf{x}\mathbf{x}^T = \Omega_0 \Omega_0^T \quad (18)$$

with an  $n \times (n+1)$  matrix

$$\Omega_0 = (\Lambda \ \mathbf{x}_0) = (l_1 l_2 \cdots l_n \ \mathbf{x}_0), \quad (19)$$

where we set  $\mathbf{x}_0 = \mathbf{x}$ . So, in order to obtain the Cholesky decomposition of  $\tilde{G}$ , we have only to find an  $(n+1) \times (n+1)$  orthogonal matrix  $P$  such that

$$\Omega_0 P = (\Lambda \ \mathbf{x}_0) P = (\tilde{\Lambda} \ \mathbf{0}), \quad (20)$$

where  $\tilde{\Lambda}$  is a lower triangular matrix, giving the Cholesky decomposition of  $\tilde{G}$ ;  $\tilde{G} = \tilde{\Lambda} \tilde{\Lambda}^T$ . Note that the matrix  $\Omega_0$  in Eq.(19) has such a form as

$$\Omega_0 = \begin{pmatrix} * & & & & \bullet \\ * & * & & & \bullet \\ \vdots & \vdots & \ddots & & \vdots \\ * & * & * & * & \bullet \\ * & * & * & * & * & \bullet \end{pmatrix}, \quad (21)$$

where asterisks and black dots indicate the non-zero entries. To find an orthogonal matrix  $P$  in (20), we first multiply an  $(n+1) \times (n+1)$  rotation matrix

$$R_1(\theta_1) = \begin{pmatrix} \cos \theta_1 & & & & -\sin \theta_1 \\ & 1 & & & \\ & & \ddots & & \\ & & & 1 & \\ \sin \theta_1 & & & & \cos \theta_1 \end{pmatrix} \quad (22)$$

with

$$\cos \theta_1 = \frac{l_{1,1}}{\sqrt{l_{1,1}^2 + x_{0,1}^2}}, \quad \sin \theta_1 = \frac{x_{0,1}}{\sqrt{l_{1,1}^2 + x_{0,1}^2}} \quad (23)$$

to  $\Omega_0$  from the right. Then we can eliminate the first component (top black dot) in the last column of  $\Omega_0$ ;

$$\begin{aligned} \Omega_1 \equiv \Omega_0 R_1(\theta_1) &= (\tilde{l}_1 l_2 \cdots l_n \ \mathbf{x}_1) \\ &= \begin{pmatrix} * & & & & 0 \\ * & * & & & \bullet \\ \vdots & \vdots & \ddots & & \vdots \\ * & * & * & * & \bullet \\ * & * & * & * & * & \bullet \end{pmatrix} \end{aligned} \quad (24)$$

A similar procedure can be applied to eliminate the  $k$ -th component in the last column of  $\Omega_0$  successively ( $k = 2, 3, \dots, n$ ). If we define

$$\Omega_k \equiv \Omega_{k-1} R_k(\theta_k) \quad (25)$$

