



# A comparison of generalized linear discriminant analysis algorithms

Cheong Hee Park<sup>a,\*</sup>, Haesun Park<sup>b2</sup>

<sup>a</sup>Department of Computer Science and Engineering, Chungnam National University, 220 Gung-dong, Yuseong-gu, Daejeon 305 763, Republic of Korea

<sup>b</sup>College of Computing, Georgia Institute of Technology, 801 Atlantic Drive, Atlanta, GA 30332, USA

Received 26 April 2006; received in revised form 27 June 2007; accepted 30 July 2007

## Abstract

Linear discriminant analysis (LDA) is a dimension reduction method which finds an optimal linear transformation that maximizes the class separability. However, in undersampled problems where the number of data samples is smaller than the dimension of data space, it is difficult to apply the LDA due to the singularity of scatter matrices caused by high dimensionality. In order to make the LDA applicable, several generalizations of the LDA have been proposed recently. In this paper, we present theoretical and algorithmic relationships among several generalized LDA algorithms and compare their computational complexities and performances in text classification and face recognition. Towards a practical dimension reduction method for high dimensional data, an efficient algorithm is proposed, which reduces the computational complexity greatly while achieving competitive prediction accuracies. We also present nonlinear extensions of these LDA algorithms based on kernel methods. It is shown that a generalized eigenvalue problem can be formulated in the kernel-based feature space, and generalized LDA algorithms are applied to solve the generalized eigenvalue problem, resulting in nonlinear discriminant analysis. Performances of these linear and nonlinear discriminant analysis algorithms are compared extensively.  
 © 2007 Published by Elsevier Ltd on behalf of Pattern Recognition Society.

**Keywords:** Dimension reduction; Feature extraction; Generalized linear discriminant analysis; Kernel methods; Nonlinear discriminant analysis; Undersampled problems

## 1. Introduction

Linear discriminant analysis (LDA) seeks an optimal linear transformation by which the original data is transformed to a much lower dimensional space. The goal of LDA is to find a linear transformation that maximizes class separability in the reduced dimensional space. Hence the criteria for dimension reduction in LDA are formulated to maximize the between-class scatter and minimize the within-class scatter. The scatters are measured by using scatter matrices such as the

between-class scatter matrix ( $S_b$ ), within-class scatter matrix ( $S_w$ ) and total scatter matrix ( $S_t$ ). Let us denote a data set  $A$  as

$$A = [a_1, \dots, a_n] = [A_1, A_2, \dots, A_r] \in \mathbb{R}^{m \times n}, \quad (1)$$

where a collection of data items in the class  $i$  ( $1 \leq i \leq r$ ) is represented as a block matrix  $A_i \in \mathbb{R}^{m \times n_i}$  and  $N_i$  is the index set of data items in the class  $i$ . Each class  $i$  has  $n_i$  elements and the total number of data is  $n = \sum_{i=1}^r n_i$ . The between-class scatter matrix  $S_b$ , within-class scatter matrix  $S_w$  and total scatter matrix  $S_t$  are defined as

$$S_b = \sum_{i=1}^r n_i (c_i - c)(c_i - c)^T, \quad (2)$$

$$S_w = \sum_{i=1}^r \sum_{j \in N_i} (a_j - c_i)(a_j - c_i)^T, \quad (3)$$

$$S_t = \sum_{j=1}^n (a_j - c)(a_j - c)^T, \quad (4)$$

\* Corresponding author. Tel.: +82 42 821 6293.

E-mail addresses: [cheonghee@cnu.ac.kr](mailto:cheonghee@cnu.ac.kr) (C.H. Park), [hpark@cc.gatech.edu](mailto:hpark@cc.gatech.edu) (H. Park).

<sup>1</sup> This study was financially supported by research funded by Chungnam National University in 2005.

<sup>2</sup> This work was supported in part by the National Science Foundation grant CCF-0621889. Any opinions, findings and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation (NSF).

where  $c_i = (1/n_i)\sum_{j \in N_i} a_j$  and  $c = (1/n)\sum_{j=1}^n a_j$  are class centroids and the global centroid, respectively.

The optimal dimension reducing transformation  $G^T \in \mathbb{R}^{l \times m}$  ( $l < m$ ) for LDA is the one that maximizes the between-class scatter and minimizes the within-class scatter in a reduced dimensional space. Common optimization criteria for LDA are formulated as the maximization problem of objective functions

$$J_1(G) = \frac{\text{trace}(G^T S_b G)}{\text{trace}(G^T S_w G)},$$

$$J_2(G) = \text{trace}((G^T S_w G)^{-1} (G^T S_b G)),$$

$$J_3(G) = \frac{|G^T S_b G|}{|G^T S_w G|}, \quad (2)$$

where  $\tilde{S}_i = G^T S_i G$  for  $i = b, w$  are scatter matrices in the space transformed by  $G^T$ . It is well known [1,2] that when  $S_w$  is nonsingular, the transformation matrix  $G$  is obtained by the eigenvectors corresponding to the  $r - 1$  largest eigenvalues of

$$S_w^{-1} S_b g = \lambda g. \quad (3)$$

However, for undersampled problems such as text classification and face recognition where the number of data items is smaller than the data dimension, scatter matrices become singular and their inverses are not defined. In order to overcome the problems caused by the singularity of the scatter matrices, several methods have been proposed [3–8]. In this paper, we present theoretical relationships among several generalized LDA algorithms and compare computational complexities and performances of them.

While linear dimension reduction has been used in many application areas due to its simple concept and easiness in computation, it is difficult to capture a nonlinear relationship in the data by a linear function. Recently kernel methods have been widely used for nonlinear extension of linear algorithms [9]. The original data space is transformed to a feature space by an implicit nonlinear mapping through kernel methods. As long as an algorithm can be formulated with inner product computations, without knowing the explicit representation of a nonlinear mapping we can apply the algorithm in the transformed feature space, obtaining nonlinear extension of the original algorithm. We present nonlinear extensions of generalized LDA algorithms through the formulation of a generalized eigenvalue problem in the kernel-based feature space.

The rest of the paper is organized as follows. In Section 2, a theoretical comparison of generalized LDA algorithms is presented. We study theoretical and algorithmic relationships among several generalized LDA algorithms and compare their computational complexities and performances. Computationally efficient algorithm is also proposed which computes the exactly same solution as that in Refs. [4,10] but saves computational complexities greatly. In Section 3, nonlinear extensions of these generalized LDA algorithms are presented. A generalized eigenvalue problem is formulated in the nonlinearly transformed feature space for which all the generalized LDA algorithms can be applied resulting in nonlinear dimension reduction methods. Extensive comparisons of these linear

Table 1  
Summary of the notations used

Notations	Description
$m$	Data dimension
$n$	Number of data items
$r$	Number of classes
$n_i$	Number of data items in class $i$
$c, c_i$	The global and class centroids
$A$	Data matrix of size $m \times n$
$S_b, S_w, S_t$	Scatter matrices of size $m \times m$
$H_b$	The matrix of size $m \times r$ such that $S_b = H_b H_b^T$
$H_w$	The matrix of size $m \times n$ such that $S_w = H_w H_w^T$
$H_t$	The matrix of size $m \times n$ such that $S_t = H_t H_t^T$
$s$	Rank of the matrix $[H_b H_w]$
$I_\tau, 0_\tau$	Identity and zero matrices of size $\tau \times \tau$

and nonlinear discriminant analysis algorithms are conducted. Conclusion follows in Section 4.

For convenience, important notations used throughout the rest of the paper are listed in Table 1.

## 2. A comparison of generalized LDA algorithms for undersampled problems

### 2.1. Regularized LDA

In the regularized LDA (RLDA) [3], when  $S_w$  is singular or ill-conditioned, a diagonal matrix  $\alpha I$  with  $\alpha > 0$  is added to  $S_w$ . Since  $S_w$  is symmetric positive semidefinite,  $S_w + \alpha I$  is nonsingular with any  $\alpha > 0$ . Therefore, we can apply the algorithm for the classical LDA to solve the eigenvalue problem

$$S_b g = \lambda (S_w + \alpha I) g. \quad (4)$$

#### 2.1.1. Two-class problem

We now consider a simple case when the data set has two classes, since in that case a comparison of generalized LDA algorithms is easy to illustrate. The two-class problem in LDA is known as Fisher discriminant analysis (FDA) [2]. In a two-class case,  $S_b$  can be expressed as

$$S_b = \frac{n_1 n_2}{n} (c_1 - c_2)(c_1 - c_2)^T, \quad (5)$$

and the eigenvalue problem (3) is simplified to

$$S_w^{-1} (c_1 - c_2)(c_1 - c_2)^T g = \lambda g, \quad (6)$$

when  $S_w$  is nonsingular. The solution for Eq. (6) is a nonzero multiple of  $g = S_w^{-1} (c_1 - c_2)$ , and the 1-dimensional representation of any data item  $z \in \mathbb{R}^{m \times 1}$  by LDA is obtained as

$$g^T z = (c_1 - c_2)^T S_w^{-1} z = (c_1 - c_2)^T U_w \Sigma_w^{-1} U_w^T z, \quad (7)$$

where  $S_w = U_w \Sigma_w U_w^T$  is the eigenvalue decomposition (EVD) of  $S_w$ . Since  $S_w + \alpha I = U_w (\Sigma_w + \alpha I) U_w^T$ , the regularized LDA

1 gives the solution

$$g^T z = (c_1 - c_2)^T U_w (\Sigma_w + \alpha I)^{-1} U_w^T z,$$

3 and the regularization parameter  $\alpha$  affects the scales of the principal components of  $S_w$ .

5 In the regularized LDA, the parameter  $\alpha$  is to be optimized experimentally since no theoretical procedure for choosing an optimal parameter is easily available. Recently, a generalization of LDA through simultaneous diagonalization of  $S_b$  and  $S_w$  using the generalized singular value decomposition (GSVD) has been developed [4]. This LDA/GSVD, summarized in the next section, does not require any parameter optimization.

## 2.2. LDA based on the GSVD

13 Howland et al. [4,10] applied the generalized singular value decomposition (GSVD) due to Paige and Saunders [11] to overcome the limitation of the classical LDA. When the GSVD is applied to two matrices  $Z_1$  and  $Z_2$  with the same number of columns,  $p$ , we obtain

$$U_1^T Z_1 X = \begin{bmatrix} \Gamma_1 & 0 \\ \gamma & p-\gamma \end{bmatrix} \quad \text{and} \quad U_2^T Z_2 X = \begin{bmatrix} \Gamma_2 & 0 \\ \gamma & p-\gamma \end{bmatrix}$$

$$\text{for } \gamma = \text{rank} \left( \begin{bmatrix} Z_1 \\ Z_2 \end{bmatrix} \right),$$

19 where  $U_1$  and  $U_2$  are orthogonal and  $X$  is nonsingular,  $\Gamma_1^T \Gamma_1 + \Gamma_2^T \Gamma_2 = I_\gamma$  and  $\Gamma_1^T \Gamma_1$  and  $\Gamma_2^T \Gamma_2$  are diagonal matrices with nonincreasing and nondecreasing diagonal components, respectively.

23 The method in Ref. [4] utilized the representations of the scatter matrices

$$25 \quad S_b = H_b H_b^T, \quad S_w = H_w H_w^T, \quad \text{and} \quad S_t = H_t H_t^T. \quad (7)$$

where

$$27 \quad H_b = [\sqrt{n_1}(c_1 - c), \dots, \sqrt{n_r}(c_r - c)] \in \mathbb{R}^{m \times r}, \quad (8)$$

$$H_w = [A_1 - c_1 e_1, \dots, A_r - c_r e_r] \in \mathbb{R}^{m \times n}, \quad (9)$$

$$29 \quad H_t = [a_1 - c, \dots, a_n - c] \in \mathbb{R}^{m \times n}, \quad (10)$$

31 and  $e_i = [1, \dots, 1] \in \mathbb{R}^{1 \times n_i}$ . Suppose the GSVD is applied to the matrix pair  $(H_b^T, H_w^T)$  and we obtain

$$U_b^T H_b^T X = [\Gamma_b \ 0] \quad \text{and} \quad U_w^T H_w^T X = [\Gamma_w \ 0], \quad (11)$$

33 where  $U_b \in \mathbb{R}^{r \times r}$  and  $U_w \in \mathbb{R}^{n \times n}$  are orthogonal,  $X \in \mathbb{R}^{m \times m}$  is nonsingular, and

$$35 \quad \Gamma_b^T \Gamma_b + \Gamma_w^T \Gamma_w = I_s \quad \text{for } s = \text{rank} \left( \begin{bmatrix} H_b^T \\ H_w^T \end{bmatrix} \right).$$

Then Eqs. in (11) give

$$37 \quad \begin{aligned} X^T S_b X &= X^T (H_b H_b^T) X = (X^T H_b U_b) (U_b^T H_b^T X) \\ &= [\Gamma_b \ 0]^T [\Gamma_b \ 0], \end{aligned}$$

$$\begin{aligned} X^T S_w X &= X^T (H_w H_w^T) X = (X^T H_w U_w) (U_w^T H_w^T X) \\ &= [\Gamma_w \ 0]^T [\Gamma_w \ 0]. \end{aligned} \quad (12)$$

From (12) and  $\Gamma_b^T \Gamma_b + \Gamma_w^T \Gamma_w = I_s$ , we have 39

$$X^T S_b X = \begin{bmatrix} \Gamma_b^T \Gamma_b & \\ & 0_{m-s} \end{bmatrix} \equiv \begin{bmatrix} I_\mu & & \\ & D_\tau & \\ & & 0_{s-\mu-\tau} \\ & & & 0_{m-s} \end{bmatrix} \quad (13)$$

and 41

$$X^T S_w X = \begin{bmatrix} \Gamma_w^T \Gamma_w & \\ & 0_{m-s} \end{bmatrix} \equiv \begin{bmatrix} 0_\mu & & \\ & E_\tau & \\ & & I_{s-\mu-\tau} \\ & & & 0_{m-s} \end{bmatrix}, \quad (14)$$

43 where  $D_\tau + E_\tau = I_\tau$  and the subscripts in  $I$  and  $0$  denote the size of square identity and zero matrices. Denoting the diagonal elements in  $\Gamma_b^T \Gamma_b$  as  $\eta_i$ 's and the diagonal elements in  $\Gamma_w^T \Gamma_w$  as  $\zeta_i$ 's, we have 45

$$\zeta_i S_b x_i = \eta_i S_w x_i, \quad i = 1, \dots, m, \quad (15) \quad 47$$

49 where  $x_i$  is the column vectors of  $X$ . Note that  $x_i$ ,  $i = s + 1, \dots, m$ , belong to  $\text{null}(S_b) \cap \text{null}(S_w)$ . Hence  $\eta_i$  and  $\zeta_i$  for  $i = s + 1, \dots, m$  in Eq. (15) can be any arbitrary numbers.

By partitioning  $X$  in Eqs. (13)–(14) as 51

$$X = \begin{bmatrix} X_1 & X_2 & X_3 & X_4 \end{bmatrix} \in \mathbb{R}^{m \times m}, \quad (16)$$

$\mu \quad \tau \quad s-\mu-\tau \quad m-s$

53 the generalized eigenvalues and eigenvectors obtained by the GSVD can be classified as shown in Table 2. For the last  $m - s$  vectors  $x$  belonging to  $\text{null}(S_w) \cap \text{null}(S_b)$ , 55

$$0 = x^T S_b x = (x^T H_b) (H_b^T x) = \|x^T H_b\|^2 = \sum_{i=1}^r n_i |x^T c_i - x^T c|^2$$

and 57

$$0 = x^T S_w x = \sum_{j=1}^n |x^T a_j - x^T c_j|^2,$$

where  $a_j$  belongs to a class  $i$ . 59

Hence

$$\begin{cases} x^T c_i = x^T c & \text{for } i = 1, \dots, r, \\ x^T a_j = x^T c_i & \text{for all } a_j \text{ in a class } i, \end{cases} \quad (17) \quad 61$$

therefore

$$X_4^T z = X_4^T c \quad (18) \quad 63$$

65 for any given data item  $z = a_i$ . This implies that the vectors  $x_i$ ,  $i = s + 1, \dots, m$ , belonging to  $\text{null}(S_b) \cap \text{null}(S_w)$  do not convey discriminative information among the classes, even though 67

Table 2  
Generalized eigenvalues  $\lambda_i$ 's and eigenvectors  $x_i$ 's from the GSVD

	$\eta_i$	$\zeta_i$	$\lambda_i = \frac{\eta_i}{\zeta_i}$	$x_i$ belongs to
$1 \leq i \leq \mu$	1	0	$\infty$	$\text{null}(S_w) \cap \text{null}(S_b)^c$
$\mu + 1 \leq i \leq \mu + \tau$	$1 > \eta_i > 0$	$0 < \zeta_i < 1$	$\infty > \lambda_i > 0$	$\text{null}(S_w)^c \cap \text{null}(S_b)^c$
$\mu + \tau + 1 \leq i \leq s$	0	1	0	$\text{null}(S_w) \cap \text{null}(S_b)$
$s + 1 \leq i \leq m$	Any value	Any value	Any value	$\text{null}(S_w) \cap \text{null}(S_b)$

The superscript  $c$  denotes the complement.

1 the corresponding eigenvalues are not necessarily zeros. Since  
rank( $S_b$ )  $\leq r - 1$ , from Eqs. (13)–(14), we have

$$3 \quad x_i^T S_b x_i = 0 \quad \text{and} \quad x_i^T S_w x_i = 1 \quad \text{for} \quad r \leq i \leq s,$$

and the between-class scatter becomes zero by the projection  
5 onto the vector  $x_i$ . Hence  $r - 1$  leftmost columns of  $X$  gives  
an optimal transformation  $G_h^T$  for LDA. This method is called  
7 LDA/GSVD.

### 2.2.1. An efficient algorithm for LDA/GSVD

9 The algorithm to compute the GSVD for the pair ( $H_b^T$ ,  $H_w^T$ )  
was presented in Ref. [4] as follows:

11 (1) Compute the singular value decomposition (SVD) of

$$Z = \begin{bmatrix} H_b^T \\ H_w^T \end{bmatrix} \in \mathbb{R}^{(r+n) \times m}; \quad Z = P \begin{bmatrix} A & 0 \\ 0 & 0 \end{bmatrix} U^T$$

13 where  $s = \text{rank}(Z)$  and  $P \in \mathbb{R}^{(r+n) \times (r+n)}$  and  $U \in \mathbb{R}^{m \times m}$   
are orthogonal and the diagonal components of  $A \in \mathbb{R}^{s \times s}$   
15 is nonincreasing.

17 (2) Compute  $V$  from the SVD of  $P(1 : r, 1 : s)$ ,<sup>3</sup> which is  
 $P(1 : r, 1 : s) = W F V^T$ .

19 (3) Compute the first  $r - 1$  columns of  $X = U \begin{bmatrix} A^{-1} V & 0 \\ 0 & I \end{bmatrix}$ , and  
assign them to the transformation matrix  $G_h$ .

21 Now we show that this algorithm can be computed rather  
simply, producing an efficient and intuitive approach for  
LDA/GSVD. Since  $\Gamma_b^T \Gamma_b + \Gamma_w^T \Gamma_w = I_s$ , from Eqs. (13)–(14),  
23 we have

$$X^T S_t X = X^T S_b X + X^T S_w X = \begin{bmatrix} I_s & 0 \\ 0 & 0 \end{bmatrix}, \quad (19)$$

25 where  $s = \text{rank}(Z)$ . Eq. (19) implies  $s = \text{rank}(S_t)$  and from step  
3 in the LDA/GSVD algorithm

$$27 \quad S_t = X^{-T} \begin{bmatrix} I_s & 0 \\ 0 & 0 \end{bmatrix} X^{-1} = U \begin{bmatrix} \Sigma_1 & 0 \\ 0 & 0 \end{bmatrix} U^T, \quad \Sigma_1 = A^T A, \quad (20)$$

<sup>3</sup> The notation  $P(1 : r, 1 : s)$  which may appear as a MATLAB shorthand  
denotes a submatrix of  $P$  composed of the components from the first to the  
 $r$ th row and from the first to  $s$ th column.

which results in the EVD of  $S_t$ . Partitioning  $U$  as  $U =$   
 $\begin{bmatrix} U_1 & U_2 \\ \underbrace{\hspace{1cm}}_s & \underbrace{\hspace{1cm}}_{m-s} \end{bmatrix}$ , we have

$$X = U \begin{bmatrix} A^{-1} V & 0 \\ 0 & I \end{bmatrix} = [U_1 \Sigma_1^{-1/2} V \quad U_2]. \quad (21)$$

By substituting  $X$  in Eq. (13) with Eq. (21),

$$\Sigma_1^{-1/2} U_1^T S_b U_1 \Sigma_1^{-1/2} = V \Gamma_b^T \Gamma_b V^T. \quad (22)$$

Note that the optimal transformation matrix  $G_h$  by LDA/GSVD  
is obtained by the leftmost  $r - 1$  columns of  $X$ , which are the  
leftmost  $r - 1$  columns of  $U_1 \Sigma_1^{-1/2} V$ . Eqs. (20) and (22) show  
that  $U_1$  and  $\Sigma_1$  can be computed from the EVD of  $S_t$  and  $V$   
from the EVD of  $\Sigma_1^{-1/2} U_1^T S_b U_1 \Sigma_1^{-1/2}$ . This new approach for  
LDA/GSVD is summarized in Algorithm 1.

**Algorithm 1.** An efficient algorithm for LDA/GSVD.

(1) Compute the EVD of  $S_t$ :  $S_t = [U_1 \quad U_2] \begin{bmatrix} \Sigma_1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} U_1^T \\ U_2^T \end{bmatrix}$ .

(2) Compute  $V$  from the EVD of  $\tilde{S}_b \equiv \Sigma_1^{-1/2} U_1^T S_b U_1 \Sigma_1^{-1/2}$ :  
 $\tilde{S}_b = V \Gamma_b^T \Gamma_b V^T$ .

(3) Assign the first  $r - 1$  columns of  $U_1 \Sigma_1^{-1/2} V$  to  $G_h$ .

In Algorithm 1, the matrices  $U_1$  and  $\Sigma_1$  in the EVD of  $S_t \in$   
 $\mathbb{R}^{m \times m}$  can be obtained by the EVD of  $H_t^T H_t \in \mathbb{R}^{n \times n}$  instead  
of  $H_t H_t^T \in \mathbb{R}^{m \times m}$  [1] by which computational complexity can  
be reduced from  $O(m^3)$  to  $O(n^3)$ . Especially when  $m$  is much  
bigger than  $n$ , computational savings become great. Let the  
EVD of  $H_t^T H_t$  be

$$H_t^T H_t = \underbrace{\begin{bmatrix} J_1 & \\ & J_2 \end{bmatrix}}_s \underbrace{\begin{bmatrix} & \\ & n-s \end{bmatrix}}_{n-s} \begin{bmatrix} D_1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} J_1^T \\ J_2^T \end{bmatrix}, \quad (23)$$

where  $s = \text{rank}(H_t) = \text{rank}(S_t)$ . From Eq. (23)

$$S_t (H_t J_1) = H_t (H_t^T H_t) J_1 = (H_t J_1) D_1,$$

and therefore the columns in  $H_t J_1$  are eigenvectors of  $S_t$  corre-  
sponding to nonzero eigenvalues in the diagonal of  $D_1$ . Since  
 $(H_t J_1)^T (H_t J_1) = D_1$ , we obtain the orthonormal eigenvectors  
and corresponding nonzero eigenvalues of  $S_t$  by  $H_t J_1 D_1^{-1/2}$   
and  $D_1$ , which are  $U_1$  and  $\Sigma_1$ , respectively. In this new ap-  
proach, we just need to compute the EVD of a much smaller  
 $n \times n$  matrix  $H_t^T H_t$  instead of  $m \times m$  matrix  $S_t = H_t H_t^T$  when  
 $m \gg n$ . However, in the regularized LDA or the method by Chen  
et al. which is presented next, we cannot resort to this approach.

1 The regularized LDA needs the entire  $m$  eigenvectors of  $S_w$  and  
 2 the method based on the projection to  $\text{null}(S_w)$  needs to com-  
 3 pute a basis of  $\text{null}(S_w)$  which are eigenvectors corresponding  
 4 to zero eigenvalues.

### 5 2.2.2. Two-class problem

6 Now we consider the two-class problem in LDA/GSVD. By  
 7 Eq. (5), we have

$$\begin{aligned} \Sigma_1^{-1/2} U_1^T S_b U_1 \Sigma_1^{-1/2} &= \Sigma_1^{-1/2} U_1^T \rho (c_1 - c_2) (c_1 - c_2)^T U_1 \Sigma_1^{-1/2} \\ &= \left( \frac{w}{\|w\|_2} \right) \rho \|w\|_2^2 \left( \frac{w}{\|w\|_2} \right)^T, \end{aligned}$$

9 where  $\rho = n_1 n_2 / n$  and  $w = \Sigma_1^{-1/2} U_1^T (c_1 - c_2)$ . Hence the  
 10 transformation matrix  $g \in \mathbb{R}^{m \times 1}$  is given by

$$11 \quad g = v U_1 \Sigma_1^{-1/2} w = v U_1 \Sigma_1^{-1} U_1^T (c_1 - c_2)$$

12 for some scalar  $v$ , and the dimension reduced representation of  
 13 any data item  $z$  is given by

$$g^T z = v (c_1 - c_2)^T U_1 \Sigma_1^{-1} U_1^T z = v (c_1 - c_2)^T S_t^+ z,$$

15 where  $S_t^+$  denotes the pseudoinverse of  $S_t$ . When  $S_w$  is non-  
 16 singular, by applying the Sherman–Morrison formula [12] to  
 17  $S_t = S_w + S_b$ , we have

$$\begin{aligned} S_t^{-1} &= (S_w + \rho (c_1 - c_2) (c_1 - c_2)^T)^{-1} = S_w^{-1} \\ &\quad - \frac{S_w^{-1} \rho (c_1 - c_2) (c_1 - c_2)^T S_w^{-1}}{1 + \rho (c_1 - c_2)^T S_w^{-1} (c_1 - c_2)} \end{aligned}$$

19 and

$$g^T z = v (c_1 - c_2)^T S_t^{-1} z = v_1 (c_1 - c_2)^T S_w^{-1} z \quad (24)$$

21 for a scalar  $v_1 = v / (1 + \rho (c_1 - c_2)^T S_w^{-1} (c_1 - c_2))$ . Eq. (24)  
 22 shows that LDA/GSVD is equal to the classical LDA when  $S_w$   
 23 is nonsingular.

### 24 2.3. A method based on the projection onto $\text{null}(S_w)$

25 In face recognition, in the efforts to overcome the singularity  
 26 of scatter matrices caused by high dimensionality, some meth-  
 27 ods have been proposed [5,6]. The basic principle of the algo-  
 28 rithms proposed in Refs. [5,6] is that the transformation using  
 29 a basis of either  $\text{range}(S_b)$  or  $\text{null}(S_w)$  is performed in the first  
 30 stage and then in the transformed space the second projective  
 31 directions are searched. These methods are summarized in this  
 32 and next section where we also present their algebraic relation-  
 33 ships.

34 Chen et al. [5] proposed a generalized method of LDA which  
 35 solves undersampled problems and applied it for face recogni-  
 36 tion. The method projects the original space onto the null space  
 37 of  $S_w$  using an orthonormal basis of  $\text{null}(S_w)$ , and then in the  
 38 projected space, a transformation that maximizes the between-  
 39 class scatter is computed.

Consider the SVD of  $S_w \in \mathbb{R}^{m \times m}$ ,

$$S_w = U_w \Sigma_w U_w^T. \quad 41$$

Partitioning  $U_w$  as  $U_w = \begin{bmatrix} U_{w1} & U_{w2} \end{bmatrix}$  where  $s_1 = \text{rank}(S_w)$ ,

$$\text{null}(S_w) = \text{span}(U_{w2}). \quad (25) \quad 43$$

44 First, the transformation by  $U_{w2} U_{w2}^T$  projects the original  
 45 data to  $\text{null}(S_w)$ . Then, the eigenvectors corresponding to  
 46 the largest eigenvalues of the between-class scatter matrix  
 47  $\tilde{S}_b$  in the projected space are found. Let the EVD of  $\tilde{S}_b \equiv$   
 $U_{w2} U_{w2}^T S_b U_{w2} U_{w2}^T$  be

$$\tilde{S}_b = \tilde{U}_b \tilde{\Sigma}_b \tilde{U}_b^T = \begin{bmatrix} \tilde{U}_{b1} & \tilde{U}_{b2} \end{bmatrix} \begin{bmatrix} \tilde{\Sigma}_{b1} & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \tilde{U}_{b1}^T \\ \tilde{U}_{b2}^T \end{bmatrix}, \quad (26) \quad 49$$

50 where  $\tilde{U}_b^T \tilde{U}_b = I$ ,  $s_2 = \text{rank}(\tilde{S}_b)$  and  $\tilde{\Sigma}_{b1} \in \mathbb{R}^{s_2 \times s_2}$ . Then, the  
 51 transformation matrix  $G_e$  is obtained by

$$G_e = U_{w2} U_{w2}^T \tilde{U}_{b1}. \quad (27) \quad 51$$

Let us call this method  $To-N(S_w)$  as an abbreviation. 53

#### 52 2.3.1. Two-class problem

In the two-class problem,  $S_b$  is expressed as in Eq. (5) and 55

$$\begin{aligned} \tilde{S}_b &= U_{w2} U_{w2}^T \rho (c_1 - c_2) (c_1 - c_2)^T U_{w2} U_{w2}^T \\ &= \left( \frac{w}{\|w\|_2} \right) \rho \|w\|_2^2 \left( \frac{w}{\|w\|_2} \right)^T, \end{aligned}$$

56 where  $\rho = n_1 n_2 / n$  and  $w = U_{w2} U_{w2}^T (c_1 - c_2) \in \mathbb{R}^{m \times 1}$ . Hence  
 57 the transformation matrix  $g \in \mathbb{R}^{m \times 1}$  is obtained by

$$g = U_{w2} U_{w2}^T \frac{w}{\|w\|_2} = v U_{w2} U_{w2}^T (c_1 - c_2), \quad 59$$

60 with  $v = 1 / \|w\|_2$ . For any data item  $z \in \mathbb{R}^{m \times 1}$ , the dimension  
 61 reduced representation is given by

$$g^T z = v (c_1 - c_2)^T U_{w2} U_{w2}^T z.$$

#### 62 2.3.2. Relationship with LDA/GSVD 63

From (26), we have

$$\begin{bmatrix} \tilde{U}_{b1}^T \\ \tilde{U}_{b2}^T \end{bmatrix} U_{w2} U_{w2}^T S_b U_{w2} U_{w2}^T \begin{bmatrix} \tilde{U}_{b1} \\ \tilde{U}_{b2} \end{bmatrix} = \begin{bmatrix} \tilde{\Sigma}_{b1} & 0 \\ 0 & 0 \end{bmatrix}, \quad (28) \quad 65$$

$$\begin{bmatrix} \tilde{U}_{b1}^T \\ \tilde{U}_{b2}^T \end{bmatrix} U_{w2} U_{w2}^T S_w U_{w2} U_{w2}^T \begin{bmatrix} \tilde{U}_{b1} \\ \tilde{U}_{b2} \end{bmatrix} = 0. \quad (29) \quad 67$$

68 The second equation holds due to Eq. (25). Eqs. (28)–(29)  
 69 imply that the column vectors of  $G_e$  given in Eq. (27) belong  
 70 to  $\text{null}(S_w) \cap \text{null}(S_b)^c$  and they are discriminative vectors,  
 71 since the transformation by these vectors minimizes the within-  
 72 class scatter to zero and increases the between-class scatter.  
 73 The top row of Table 2 shows that the LDA/GSVD solution  
 also includes the vectors from  $\text{null}(S_w) \cap \text{null}(S_b)^c$ . Based on

this observation, this method  $To-N(S_w)$  can be compared with LDA/GSVD. By denoting  $X$  in LDA/GSVD as

$$X = \left[ \underbrace{X_1}_{\mu} \underbrace{X_2}_{\tau} \underbrace{X_3}_{s-\mu-\tau} \underbrace{X_4}_{m-s} \right], \quad (30)$$

we find a relationship between  $X_1$  and  $G_e = U_{w2}U_{w2}^T\tilde{U}_{b1}$ .

Eq. (14) implies that  $[X_1 \ X_4]$  is a basis of  $\text{null}(S_w)$ . Hence any vector in  $\text{null}(S_w)$  can be represented as a linear combination of column vectors in  $[X_1 \ X_4]$ . The following Theorem shows the condition for any vector in  $\text{null}(S_w)$  to belong to  $\text{null}(S_w) \cap \text{null}(S_b)^c$ .

**Theorem I.** Any vector  $x$  belongs to  $\text{null}(S_w) \cap \text{null}(S_b)^c$  if and only if  $x$  is represented as  $X_1h + X_4k$  where  $h \neq 0 \in \mathbb{R}^{\mu \times 1}$  and  $k \in \mathbb{R}^{(m-s) \times 1}$ .

**Proof.** Let  $x \in \text{null}(S_w) \cap \text{null}(S_b)^c$ . Since  $[X_1 \ X_4]$  is a basis of  $\text{null}(S_w)$ ,  $x = X_1h + X_4k$  for some  $h \in \mathbb{R}^{\mu \times 1}$  and  $k \in \mathbb{R}^{(m-s) \times 1}$ . Suppose  $h = 0$ . Then  $x = X_4k \in \text{null}(S_w) \cap \text{null}(S_b)$ , which contradicts to  $x \in \text{null}(S_w) \cap \text{null}(S_b)^c$ . Hence  $h \neq 0$ .

Now let us prove that if  $h \neq 0$  then  $x = X_1h + X_4k$  belongs to  $\text{null}(S_w) \cap \text{null}(S_b)^c$ . Since  $x = X_1h + X_4k \in \text{null}(S_w)$ , it is enough to show  $x \notin \text{null}(S_b)$ . From Eq. (13),

$$\begin{aligned} x^T S_b x &= (X_1 h)^T S_b (X_1 h) = h^T (X_1^T S_b X_1) h \\ &= h^T I_{\mu} h = \|h\|_2^2 \neq 0. \quad \square \end{aligned}$$

By Theorem I,

$$U_{w2}U_{w2}^T\tilde{U}_{b1} = X_1 H + X_4 K$$

for some matrices  $H \in \mathbb{R}^{\mu \times s_2}$  and  $K \in \mathbb{R}^{(m-s) \times s_2}$  with  $s_2 = \text{rank}(\tilde{S}_b)$ , where each column of  $H$  is nonzero. Hence for any data item  $z \in \mathbb{R}^{m \times 1}$ , the reduced dimensional representation by  $G_e = U_{w2}U_{w2}^T\tilde{U}_{b1}$  is given as

$$G_e^T z = H^T X_1^T z + K^T X_4^T z. \quad (31)$$

As explained in Eq. (17) of Section 2.2, since all data items are transformed to one point by  $x^T$  for  $x \in \text{null}(S_w) \cap \text{null}(S_b)$ , the second part  $K^T X_4^T z$  in (31) corresponds to the translation which does not affect the classification performance.

While the transformation matrix  $G_e = U_{w2}U_{w2}^T\tilde{U}_{b1}$  by the method  $To-N(S_w)$  is related to  $X_1$  of LDA/GSVD as in Eq. (31), the main difference between the two methods is due to the eigenvectors in  $\text{null}(S_w)^c \cap \text{null}(S_b)^c$ , which correspond to the second row in Table 2. The projection to  $\text{null}(S_w)$  by  $U_{w2}U_{w2}^T$  excludes vectors in  $\text{null}(S_w)^c$ , and therefore  $\text{null}(S_w)^c \cap \text{null}(S_b)^c$ . When

$$\text{rank}(\tilde{S}_b) < \text{rank}(S_b) \leq r - 1,$$

where  $r$  is the number of classes, the reduced dimension by  $G_e = U_{w2}U_{w2}^T\tilde{U}_{b1}$  is  $\text{rank}(\tilde{S}_b)$ , therefore less than  $r - 1$ , while LDA/GSVD includes  $r - 1$  vectors from both  $\text{null}(S_w) \cap \text{null}(S_b)^c$  and  $\text{null}(S_w)^c \cap \text{null}(S_b)^c$ . In order to

demonstrate this case, we conducted an experiment using data in text classification, of which characteristics will be discussed in detail in the section for experiments. The data was collected from Reuters-21578 database and contains four classes. Each class has 80 samples and the data dimension is 2412. After splitting the data set randomly to training data and test data with a ratio of 4:1, the linear transformations by LDA/GSVD and the method  $To-N(S_w)$  were computed by using training data. While the rank of  $S_b$  was 3, the rank of  $\tilde{S}_b$  was 2 in this data set. Hence the reduced dimension by the method  $To-N(S_w)$  due to Chen et al. was 2. On the other hand, LDA/GSVD produced two eigenvectors from  $\text{null}(S_w) \cap \text{null}(S_b)^c$  and one eigenvector from  $\text{null}(S_w)^c \cap \text{null}(S_b)^c$ , resulting in the reduced dimension 3. Fig. 1 illustrates the reduced dimensional spaces by both methods. The top three figures were generated by LDA/GSVD. For the visualization, the data reduced to 3-dimensional space by LDA/GSVD was projected to 2-dimensional spaces,  $x$ - $y$ ,  $x$ - $z$  and  $y$ - $z$  spaces, respectively. In  $x$ - $y$  space, two classes ( $\Delta$  and  $*$ ) are well separated, while two other classes ( $O$  and  $+$ ) are mixed together. However, as shown in the second and third figures, two classes mixed in  $x$ - $y$  space are separated in  $x$ - $z$  and  $y$ - $z$  spaces along  $z$  axis. This shows the third eigenvector from  $\text{null}(S_w)^c \cap \text{null}(S_b)^c$  improves the separation of classes. The bottom three figures were generated by the method based on the projection to  $\text{null}(S_w)$ . Since  $\text{rank}(\tilde{S}_b) = 2$ , the reduced dimension by that method was 2 and the first figure illustrates the reduced dimensional space. The second and third figures show that adding one more column vector from  $U_{w2}U_{w2}^T\tilde{U}_{b2}$  and increasing the reduced dimension to 3 does not improve the separation of classes mixed in  $x$ - $y$  space, since the one extra dimension comes from  $\text{null}(S_w) \cap \text{null}(S_b)$ . On the other hand, when

$$\text{rank}(\tilde{S}_b) = \text{rank}(S_b) = r - 1,$$

both LDA/GSVD and the method  $To-N(S_w)$  obtain transformation matrices  $G_h$  and  $G_e$  from  $\text{null}(S_w) \cap \text{null}(S_b)^c$ . Then the difference between two methods comes from the diagonal components of  $I_{r-1}$  and  $\tilde{S}_{b1}$  in

$$G_h^T S_b G_h = I_{r-1} \quad \text{and} \quad G_e^T S_b G_e = \tilde{S}_{b1},$$

where  $\tilde{S}_{b1}$  has nonincreasing diagonal components. As shown in the experimental results of Section 2.7, the effects of different scaling in the diagonal components may depend on the characteristics of data.

#### 2.4. A method based on the transformation by a basis of $\text{range}(S_b)$

In this section, we review another two-step approach by Yu and Yang [6] proposed to handle undersampled problems, and illustrate its relationship to other methods. Contrary to the method discussed in Section 2.3, the method presented in this section first transforms the original space by using a basis of  $\text{range}(S_b)$ , and then in the transformed space the minimization of within-class scatter is pursued.

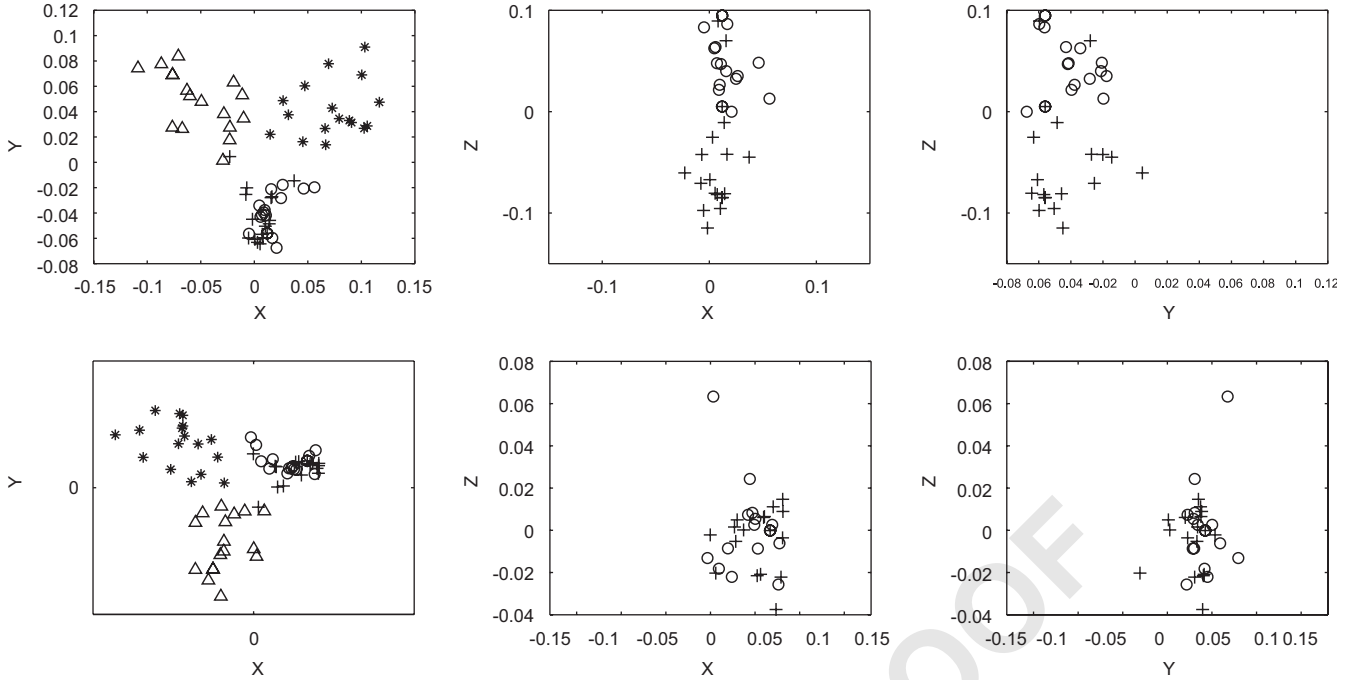


Fig. 1. The visualization of the data in the reduced dimensional spaces by LDA/GSVD (figures in the first row) and the method  $To-N(S_w)$  (figures in the second row).

1 Consider the EVD of  $S_b$ ,

$$S_b = U_b \Sigma_b U_b^T = \begin{bmatrix} U_{b1} & U_{b2} \end{bmatrix} \begin{bmatrix} \Sigma_{b1} & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} U_{b1}^T \\ U_{b2}^T \end{bmatrix},$$

$s_1$     $m-s_1$

3 where  $U_b$  is orthogonal,  $\text{rank}(S_b) = s_1$  and  $\Sigma_{b1}$  is a diagonal  
 5 matrix with nonincreasing positive diagonal components. Then  
 7  $\text{range}(S_b) = \text{span}(U_{b1})$ . In the method by Yu and Yang, the  
 original data is first transformed to an  $s_1$ -dimensional space by  
 $V_y = U_{b1} \Sigma_{b1}^{-1/2}$ . Then the between-class scatter matrix  $\tilde{S}_b$  in  
 the transformed space becomes

$$9 \quad \tilde{S}_b \equiv V_y^T S_b V_y = I_{s_1}.$$

Now consider the EVD of  $\tilde{S}_w \equiv V_y^T S_w V_y$ ,

$$11 \quad \tilde{S}_w = \tilde{U}_w \tilde{\Sigma}_w \tilde{U}_w^T, \quad (32)$$

13 where  $\tilde{U}_w \in \mathbb{R}^{s_1 \times s_1}$  is orthogonal and  $\tilde{\Sigma}_w \in \mathbb{R}^{s_1 \times s_1}$  is a diagonal matrix. Then

$$\tilde{U}_w^T V_y^T S_b V_y \tilde{U}_w = I_{s_1} \quad \text{and} \quad \tilde{U}_w^T V_y^T S_w V_y \tilde{U}_w = \tilde{\Sigma}_w. \quad (33)$$

15 In most applications,  $\text{rank}(S_w)$  is greater than  $\text{rank}(S_b)$ , and  
 $\tilde{\Sigma}_w$  is nonsingular since

$$\begin{aligned} \text{rank}(\tilde{U}_w^T V_y^T S_w V_y \tilde{U}_w) &= \text{rank}(S_w) \geq \text{rank}(S_b) \\ 17 \quad &= \text{rank}(\tilde{U}_w^T V_y^T S_b V_y \tilde{U}_w) = s_1. \end{aligned}$$

Scaling (33) by  $\tilde{\Sigma}_w^{-1/2}$ , we have

$$\begin{aligned} (\tilde{\Sigma}_w^{-1/2} \tilde{U}_w^T V_y^T) S_b (V_y \tilde{U}_w \tilde{\Sigma}_w^{-1/2}) &= \tilde{\Sigma}_w^{-1}, \\ (\tilde{\Sigma}_w^{-1/2} \tilde{U}_w^T V_y^T) S_w (V_y \tilde{U}_w \tilde{\Sigma}_w^{-1/2}) &= I_{s_1}. \end{aligned} \quad (34) \quad 19$$

The authors in Ref. [6] proposed the transformation matrix

$$G_y = V_y \tilde{U}_w \tilde{\Sigma}_w^{-1/2}. \quad 21$$

Eqs. (34) imply that each column of  $G_y$  belongs to  $\text{null}(S_w)^c \cap \text{null}(S_b)^c$ . We call this method  $To-R(S_b)$  for short. 23

#### 2.4.1. Two-class problem

In a two-class problem, since 25

$$\begin{aligned} S_b &= \rho(c_1 - c_2)(c_1 - c_2)^T \\ &= \left( \frac{c_1 - c_2}{\|c_1 - c_2\|_2} \right) \rho \|c_1 - c_2\|_2^2 \left( \frac{c_1 - c_2}{\|c_1 - c_2\|_2} \right)^T, \end{aligned}$$

27 where  $\rho = n_1 n_2 / n$ , a data item is transformed to the 1-  
 dimensional space by  $g = (c_1 - c_2) / (\sqrt{\rho} \|c_1 - c_2\|_2)$ . The  
 29 dimension reduced representation of any data item  $z$  is given by  
 $g^T z = v(c_1 - c_2)^T z$  for some scalar  $v$ . Note that no minimization  
 of within-class scatter in the transformed space is possible. 31

The optimization criteria by  $J_2$  and  $J_3$  in (2) are invariant  
 33 under any nonsingular linear transformation, i.e. for any non-  
 singular matrix  $F$  whose order is the same as that of the column  
 35 dimension of  $G$ ,

$$J_i(G) = J_i(GF), \quad i = 2, 3, \quad (35)$$

Table 3  
The prediction accuracies (%)

Face data	Transformation matrix		
	$G_y = V_y$	$G_y = V_y \tilde{U}_w$	$G_y = V_y \tilde{U}_w \tilde{\Sigma}_w^{-1/2}$
AT&T	94.3	94.3	99.0
Yale	80.6	80.6	89.7

while the objective function  $J_1$  is not. Hence in the transformation matrix  $G_y = V_y \tilde{U}_w \tilde{\Sigma}_w^{-1/2}$  obtained by the method  $To-R(S_b)$ , none of the components  $\tilde{\Sigma}_w^{-1/2}$  and  $\tilde{U}_w \tilde{\Sigma}_w^{-1/2}$  involved in the second step (those in Eqs. (32)–(34)) improves the optimization criteria by  $J_2$  and  $J_3$ . However, the following experimental results show that the scaling by  $\tilde{\Sigma}_w^{-1/2}$  can make dramatic effects on the classification performances. Postponing the detailed explanation on the data sets and experimental setting until Section 2.7, experimental results on the face recognition data sets are shown in Table 3. After dimension reduction, 1-NN classifier was used in the reduced dimensional space.

### 2.5. A method of PCA plus transformations to range( $S_w$ ) and null( $S_w$ )

As shown in the analysis of the compared methods, they search for discriminative vectors in  $\text{null}(S_w) \cap \text{null}(S_b)^c$  and  $\text{null}(S_w)^c \cap \text{null}(S_b)^c$ . The method  $To-N(S_w)$  by Chen et al. finds solution vectors in  $\text{null}(S_w) \cap \text{null}(S_b)^c$  and  $To-R(S_b)$  by Yu et al. restricts the search space to  $\text{null}(S_w)^c \cap \text{null}(S_b)^c$ . LDA/GSVD by Howland et al. finds solution from both spaces, however, the number of possible discriminative vectors cannot be greater than  $\text{rank}(S_b)$ , possibly resulting in solution vectors only from  $\text{null}(S_w) \cap \text{null}(S_b)^c$  in the case of high dimensional data. Recently Yang et al. [7] have proposed a method to obtain solution vectors in both spaces, which we will call  $To-NR(S_w)$ .

In the method by Yang et al., first, the transformation by the orthonormal basis of range( $S_t$ ), as in PCA, is performed. Let the SVD of  $S_t$  be

$$S_t = U_t \Sigma_t U_t^T = \underbrace{[U_{t1}]}_s \underbrace{[U_{t2}]}_{m-s} \begin{bmatrix} \Sigma_{t1} & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} U_{t1}^T \\ U_{t2}^T \end{bmatrix},$$

where  $s = \text{rank}(S_t)$ . In the transformed space by  $U_{t1}$ , let the within-scatter matrix be  $\tilde{S}_w = U_{t1}^T S_w U_{t1}$ . Then the basis of  $\text{null}(\tilde{S}_w)$  and  $\text{range}(\tilde{S}_w)$  can be found by the EVD of  $\tilde{S}_w$  as

$$\tilde{S}_w = \tilde{U}_w \tilde{\Sigma}_w \tilde{U}_w^T = [\tilde{U}_{w1} \ \tilde{U}_{w2}] \begin{bmatrix} \tilde{\Sigma}_{w1} & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \tilde{U}_{w1}^T \\ \tilde{U}_{w2}^T \end{bmatrix}. \quad (36)$$

In the transformed space by the basis  $\tilde{U}_{w2}$  of  $\text{null}(\tilde{S}_w)$ , let  $Y$  be the matrix whose columns are the eigenvectors corresponding to nonzero eigenvalues of

$$\tilde{S}_b \equiv \tilde{U}_{w2}^T U_{t1}^T S_b U_{t1} \tilde{U}_{w2}. \quad (37)$$

On the other hand, in the transformed space by the basis  $\tilde{U}_{w1}$  of  $\text{range}(\tilde{S}_w)$ , let  $Z$  be the matrix whose columns are the

eigenvectors<sup>4</sup> with the  $k$  largest nonzero eigenvalues of  $\hat{S}_t^{-1} \hat{S}_b$  where  $\hat{S}_b \equiv \tilde{U}_{w1}^T U_{t1}^T S_b U_{t1} \tilde{U}_{w1}$  and  $\hat{S}_t \equiv \tilde{U}_{w1}^T U_{t1}^T S_t U_{t1} \tilde{U}_{w1}$ . Then the transformation matrix by the method  $To-NR(S_w)$  is constructed as

$$G_d = [U_{t1} \tilde{U}_{w2} Y \quad U_{t1} \tilde{U}_{w1} Z]. \quad (38)$$

When two parts  $U_{t1} \tilde{U}_{w2} Y$  and  $U_{t1} \tilde{U}_{w1} Z$  are used for transformation matrix  $G_d$ , it will be better to normalize the columns in  $U_{t1} \tilde{U}_{w1} Z$  so that effects of both parts can be balanced.

#### 2.5.1. Relationship with the method $To-N(S_w)$

Recall from Section 2.3 that the method  $To-N(S_w)$  projects the original space onto the null space of  $S_w$  using an orthonormal basis of  $\text{null}(S_w)$ , and then in the projected space, a transformation that maximizes the between-class scatter is computed.

Since  $U_{t2}$  is a basis of  $\text{null}(S_t)$  and  $\text{null}(S_t) \subset \text{null}(S_w)$ , from (36)

$$\begin{bmatrix} U_{t1}^T \\ U_{t2}^T \end{bmatrix} S_w [U_{t1} \quad U_{t2}] = \begin{bmatrix} \tilde{U}_w \tilde{\Sigma}_w \tilde{U}_w^T & 0 \\ 0 & 0 \end{bmatrix}. \quad (39)$$

By Eq. (39), we can obtain the EVD of  $S_w$  as

$$\begin{aligned} S_w &= [U_{t1} \tilde{U}_w \quad U_{t2}] \begin{bmatrix} \tilde{\Sigma}_w & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \tilde{U}_w^T U_{t1}^T \\ U_{t2}^T \end{bmatrix} \\ &= [U_{t1} \tilde{U}_{w1} \quad U_{t1} \tilde{U}_{w2} \quad U_{t2}] \begin{bmatrix} \tilde{\Sigma}_{w1} & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \tilde{U}_{w1}^T U_{t1}^T \\ \tilde{U}_{w2}^T U_{t1}^T \\ U_{t2}^T \end{bmatrix}. \end{aligned} \quad (40)$$

Eq. (40) shows that the columns of  $V \equiv [U_{t1} \tilde{U}_{w2} \quad U_{t2}]$  is an orthonormal basis of  $\text{null}(S_w)$ . Hence the transformation by  $V V^T$  gives the projection onto the null space of  $S_w$ .

Now by notation (37) and  $\text{span}(U_{t2}) = \text{null}(S_t) \subset \text{null}(S_b)$ ,

$$\begin{aligned} [U_{t1} \tilde{U}_{w2} \quad U_{t2}] \begin{bmatrix} (U_{t1} \tilde{U}_{w2})^T \\ U_{t2}^T \end{bmatrix} S_b [U_{t1} \tilde{U}_{w2} \quad U_{t2}] \begin{bmatrix} (U_{t1} \tilde{U}_{w2})^T \\ U_{t2}^T \end{bmatrix} \\ = U_{t1} \tilde{U}_{w2} \tilde{S}_b \tilde{U}_{w2}^T U_{t1}^T, \end{aligned}$$

which is the between-class scatter matrix in the projected space by  $V V^T$ . Let the EVD of  $\tilde{S}_b$  be

$$\tilde{S}_b = [\tilde{U}_{b1} \ \tilde{U}_{b2}] \begin{bmatrix} \tilde{\Sigma}_{b1} & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \tilde{U}_{b1}^T \\ \tilde{U}_{b2}^T \end{bmatrix}. \quad (65)$$

<sup>4</sup> In Ref. [7], it was claimed that the orthonormal eigenvectors of  $\hat{S}_t^{-1} \hat{S}_b$  should be used. However,  $\hat{S}_t^{-1} \hat{S}_b$  may not be symmetric therefore it is not guaranteed that there exist orthonormal eigenvectors of  $\hat{S}_t^{-1} \hat{S}_b$ .



1 Then we have the transformation matrix  $G_e$  by the method  $To-N(S_w)$  as

$$G_e = [U_{t1} \tilde{U}_{w2} \ U_{t2}] \begin{bmatrix} (U_{t1} \tilde{U}_{w2})^T \\ U_{t2}^T \end{bmatrix} U_{t1} \tilde{U}_{w2} \bar{U}_{b1} \\ = U_{t1} \tilde{U}_{w2} \bar{U}_{b1}, \quad (41)$$

which is exactly same as  $U_{t1} \tilde{U}_{w2} Y$  in  $G_d$  of (38).

## 5 2.6. Other approaches for generalized LDA

### 6 2.6.1. PCA plus LDA

7 Using PCA as a preprocessing step before applying LDA has  
8 been a traditional technique for undersampled problems and  
9 successfully applied for face recognition [13]. In this approach,  
10 data dimension is reduced by PCA so that in the reduced dimen-  
11 sional space the within-class scatter matrix becomes nonsing-  
12 ular and classical LDA can be performed. However, choosing  
13 optimal dimensions reduced by PCA is not easy and experi-  
14 mental process for it can be expensive. In Section 2.7 where we  
15 present experimental comparison of the discussed algorithms,  
16 we demonstrate the difficulty with choosing the optimal dimen-  
17 sion reduced by PCA.

### 18 2.6.2. GSLDA

19 Zheng et al. claimed that the most discriminant vectors for  
20 LDA can be chosen from

$$21 \text{null}(S_t)^\perp \cap \text{null}(S_w), \quad (42)$$

22 where  $\text{null}(S_t)^\perp$  denotes the orthogonal complement of  $\text{null}(S_t)$   
23 [8]. They also proposed a computationally efficient method  
24 called GSLDA [14] which uses the modified Gram–Schmidt  
25 orthogonalization (MGS) in order to obtain an orthogonal basis  
26 of  $\text{null}(S_t)^\perp \cap \text{null}(S_w)$ . In Ref. [14], under the assumption that  
27 the given data items are independent, MGS is applied to

$$[H_w^*, H_b^*], \quad (43)$$

28 obtaining an orthogonal basis  $Q$  of Eq. (43), where  $H_w^*$  is con-  
29 structed by deleting one column from each subblock  $A_i - c_i e_1$ ,  
30  $1 \leq i \leq r$ , in  $H_w$  and  $H_b^* = [c_1 - c, \dots, c_{r-1} - c]$ . Then the last  
31  $r - 1$  columns of  $Q$  give an orthogonal basis of Eq. (42). When  
32 applying  $L_2$ -norm as a similarity measure, using any orthog-  
33 onal basis of  $\text{null}(S_t)^\perp \cap \text{null}(S_w)$  as a transformation matrix  
34 gives the same classification performances [14].

35 In Section 2.5, it was shown that a transformation matrix  $G_e$   
36 by the method  $To-N(S_w)$  is same as the first part  $U_{t1} \tilde{U}_{w2} Y$   
37 in the transformation matrix  $G_d$  by the method  $To-NR(S_w)$ . In  
38 fact, it is not difficult to prove that under the assumption of the  
39 independence of data items,  $U_{t1} \tilde{U}_{w2} Y$  is an orthogonal basis  
40 of Eq. (42), and therefore prediction accuracies by the method  
41  $To-N(S_w)$  and GSLDA should be same.

### 42 2.6.3. Uncorrelated LDA

43 Instead of the orthogonality of the columns  $\{g_i\}$  in the trans-  
44 formation matrix  $G$ , i.e.,  $g_i^T g_j = 0$  for  $i \neq j$ , uncorrelated  
45 LDA (ULDA) imposes the  $S_t$ -orthogonal constraint,  $g_i^T S_t g_j = 0$

for  $i \neq j$  [15]. In Ref. [16], it was shown that discriminant  
46 vectors obtained by the LDA/GSVD solve the  $S_t$ -orthogonal  
47 constraint. Hence the proposed algorithm 1 can also give solu-  
48 tions for ULDA more efficiently.

## 49 2.7. Experimental comparisons of generalized LDA algorithms

50 In order to compare the discussed methods, we conducted  
51 extensive experiments using two types of data sets in text clas-  
52 sification and face recognition.

53 Text classification is a task to assign a class label to a new  
54 document based on the information from pre-classified docu-  
55 ments. A collection of documents are assumed to be repre-  
56 sented as a term-document matrix, where each document is  
57 represented as a column vector and the components of the col-  
58 umn vector denote frequencies of words appeared in the docu-  
59 ment. The term-document matrix is obtained after preprocess-  
60 ing with common words and rare term removal, stemming, term  
61 frequency and inverse term frequency weighting and normal-  
62 ization [17]. The term-document matrix representation often  
63 makes the high dimensionality inevitable.

64 For all text data sets,<sup>5</sup> they were randomly split to the train-  
65 ing set and the test set with the ratio of 4:1. Experiments are re-  
66 peated 10 times to obtain mean prediction accuracies and stan-  
67 dard deviation as a performance measure. Detailed description  
68 of text data sets is given in Table 4. After computing a trans-  
69 formation matrix using training data, both training data and test  
70 data were represented in the reduced dimensional space. In the  
71 transformed space, the nearest neighbor classifier was applied  
72 to compute the prediction accuracies for classification. For each  
73 data item in test set, it finds the nearest neighbor from the train-  
74 ing data set and predicts a class label for the test data accord-  
75 ing to the class label of the nearest neighbor. Table 5 reports the  
76 mean prediction accuracies from 10 times random splitting to  
77 training and test sets.

78 The second experiment, face recognition, is a task to iden-  
79 tify a person based on given face images with different facial  
80 expressions, illumination and poses. Since the number of pic-  
81 tures for each subject is limited and the data dimension is the  
82 number of pixels of a face image, face recognition data sets are  
83 typically severely undersampled.

84 Our experiments used two data sets, AT&T (formerly ORL)  
85 face database and Yale face database. The AT&T database has  
86 400 images, which consists of 10 images of 40 subjects. All the  
87 images were taken against a dark homogeneous background,  
88 with slightly varying lighting, facial expressions (open/closed  
89 eyes, smiling/nonsmiling), and facial details (glasses/no-  
90 glasses). The subjects are in up-right, frontal positions with  
91 tolerance for some side movement [18]. For the manageable  
92 data sizes, the images have been downsampled from the size  
93  $92 \times 112$  to  $46 \times 56$  by averaging the grey level values on  $2 \times 2$   
94 blocks. Yale face database contains 165 images, 11 images  
95

<sup>5</sup> The text data sets were downloaded from <http://www-users.cs.umn.edu/~karypis/cluto/download.html> which were collected from Reuter-21578 and TREC-5, TREC-6, TREC-7 database and preprocessed to reduce to manageable data size.

Table 4  
The description of data sets

Data	Re1	Tr12	Tr23	Tr31	Tr41	Tr45	AT&T	Yale
Dim.	3094	5896	5825	8104	7362	8175	2576	8586
No. data	490	210	187	841	757	575	400	165
Classes	5	7	4	4	5	6	40	15

Table 5  
Prediction accuracies (%)

Data	RLDA	LDA/GSVD	$To-N(S_w)$	$To-R(S_b)$	$To-NR(S_w)$
<i>Text classification</i>					
Re1	<b>95.8</b>	95.1	94.5	94.2	94.7
Tr12	95.7	<b>98.3</b>	98.1	96.7	97.6
Tr23	87.9	90.3	91.5	88.2	<b>91.8</b>
Tr31	98.6	98.4	98.6	97.7	<b>98.7</b>
Tr41	<b>98.0</b>	97.3	97.0	96.3	97.1
Tr45	93.6	93.3	94.2	94.1	<b>94.4</b>
<i>Face recognition</i>					
AT&T	98.0	93.5	98.0	<b>99.0</b>	98.8
Yale	97.6	<b>98.8</b>	97.6	89.7	98.2

For RLDA, the best accuracy among  $\alpha = 0.5, 1, 1.5$  is reported. For each data set, the best prediction accuracy is shown in boldface.

of 15 subjects. The 11 images per subject were taken under various facial expressions or configurations: center-light, with glasses, happy, left-light, without glasses, normal, right-light, sad, sleepy, surprised, and wink [19]. In our experiment, each image has been downsampled from  $320 \times 243$  to  $106 \times 81$  by averaging the grey values on  $3 \times 3$  blocks. Detailed description of face data sets is also given in Table 4. Since the number of images for each subject is small, leave-one-out method was performed where it takes one image for test set and the remaining images are used as a training set. Each image serves as a test datum by turns and the ratio of the number of correctly classified cases and the total number of data is considered as a prediction accuracy.

Table 5 summarizes the prediction accuracies from both experiments. For the regularized LDA, we report the best among the accuracies obtained with the regularization parameter  $\alpha = 0.5, 1, 1.5$ . The method based on the transformation to  $\text{range}(S_b)$ ,  $To-R(S_b)$ , gives relatively low prediction accuracies compared with the methods utilizing the null space of the within-class scatter matrix  $S_w$ . While no single method works the best in all situations, computational complexities can be dramatically different among the compared methods as we will discuss in the next section.

When PCA is performed as a preprocessing step for LDA, it is not easy to determine the dimension obtained by PCA. In the next experiment we compare *PCA plus LDA* with the generalized LDA methods discussed. Varying the dimensions reduced by PCA, LDA was applied to reduce the data dimension further to  $r - 1$  where  $r$  is the number of classes. Three figures in Fig. 2 show prediction accuracies for three data sets Tr12, Tr31 and Yale face data, respectively. The values on the horizontal axis denote the intermediate data dimensions obtained by PCA. They demonstrate the difficulty in choosing the op-

timal dimension in applying PCA as a preprocessing step for LDA, although the best prediction accuracies indicated by the peak points on the graphs are comparable with those in Table 5.

## 2.8. Analysis of computational complexities

In this section we analyze computational complexities for the discussed methods. The computational complexity for the SVD decomposition depends on what parts need to be explicitly computed. We use flop counts for the analysis of computational complexities where one flop (floating point operation) represents roughly what is required to do one addition/subtraction or one multiplication/division [12]. For the SVD of a matrix  $H \in \mathbb{R}^{p \times q}$  when  $p \gg q$ ,

$$H = U \Sigma V^T = \begin{bmatrix} U_1 & U_2 \end{bmatrix} \Sigma V^T, \\ \begin{matrix} q & p-q \end{matrix}$$

where  $U \in \mathbb{R}^{p \times p}$ ,  $\Sigma \in \mathbb{R}^{p \times q}$  and  $V \in \mathbb{R}^{q \times q}$ , the complexities (flops) can be roughly estimated as follows [12, p. 254].

Need to be computed explicitly	Complexities
$U_1, \Sigma$	$6pq^2 + 11q^3$
$U, \Sigma$	$4p^2q + 13q^3$
$U, \Sigma, V$	$4p^2q + 22q^3$

For the multiplication of the  $p_1 \times p_2$  matrix and the  $p_2 \times p_3$  matrix,  $2p_1 p_2 p_3$  flops can be counted.

For simplicity, cost for constructing  $H_b \in \mathbb{R}^{m \times r}$ ,  $H_w \in \mathbb{R}^{m \times n}$  and  $H_t \in \mathbb{R}^{m \times n}$  in Eqs. (8)–(10) was not included for the comparison, since the construction of scatter matrices is required in all the methods. For  $H \in \mathbb{R}^{p \times q}$  and  $p \gg q$ , when

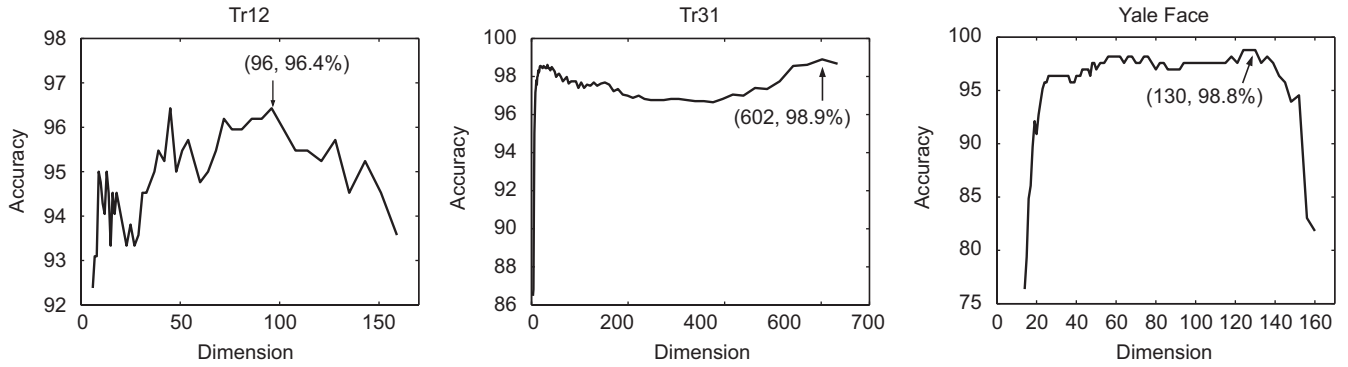


Fig. 2. The effects of the dimensions reduced by PCA on the prediction accuracies. The values on the horizontal axis denote data dimensions reduced by PCA before LDA is applied.

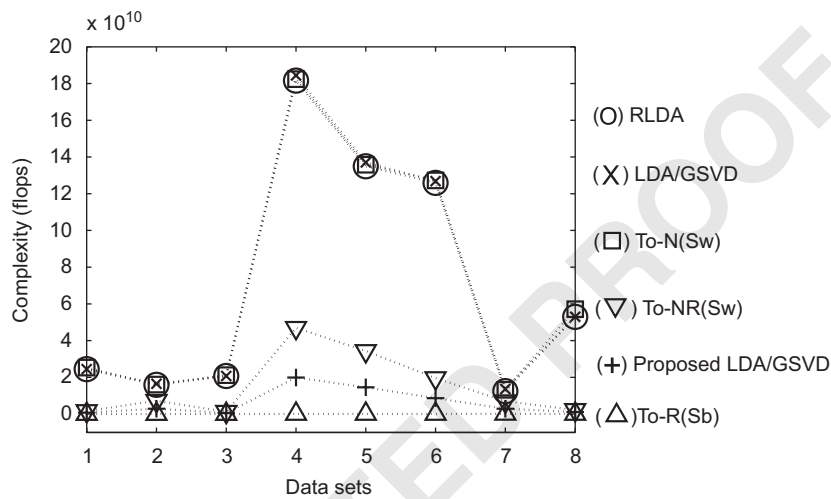


Fig. 3. Comparison of computational complexities of the generalized LDA methods using the sizes of training data used in experiments. From the left on x-axis, the data sets, Tr12, Re1, Tr23, Tr31, Tr41, Tr45, AT&T and Yale, are corresponded.

only eigenvectors corresponding to the nonzero eigenvalues of  $HH^T \in \mathbb{R}^{p \times p}$  are needed, the approach of computing the EVD of  $H^T H$  instead of  $HH^T$  as explained in Section 2.2 was utilized.

Fig. 3 compares computational complexities of the discussed methods by using specific sizes of training data sets used in the experiments. As shown in Fig. 3, regularized LDA, LDA/GSVD [4] and the method  $To-N(S_w)$  [5] have high computational complexities overall. The method  $To-R(S_b)$  [6] obtained the lowest computational costs compared with other methods while its performance cannot be ranked highly. The proposed algorithm for LDA/GSVD reduced the complexity of the original algorithm dramatically while it achieves competitive prediction accuracies as shown in Section 2.7. This new algorithm can save computational complexities even more when the number of terms is much greater than the number of documents.

### 3. Nonlinear discriminant analysis based on kernel methods

Linear dimension reduction is conceptually simple and has been used in many application areas. However, it has a limita-

tion for the data which is not linearly separable since it is difficult to capture a nonlinear relationship with a linear mapping. In order to overcome such a limitation, nonlinear extensions of linear dimension reduction methods using kernel methods have been proposed [20–25]. The main idea of kernel methods is that without knowing the nonlinear feature mapping or the mapped feature space explicitly, we can work on the nonlinearly transformed feature space through kernel functions. It is based on the fact that for any kernel function  $\kappa$  satisfying Mercer's condition, there exists a reproducing kernel Hilbert space  $H$  and a feature map  $\Phi$  such that

$$\kappa(x, y) = \langle \Phi(x), \Phi(y) \rangle, \quad (44)$$

where  $\langle \cdot, \cdot \rangle$  is an inner product in  $H$  [9,26,27].

Suppose that given a kernel function  $\kappa$  original data space is mapped to a feature space (possibly an infinite dimensional space) through a nonlinear feature mapping  $\Phi: \mathcal{A} \subset \mathbb{R}^m \rightarrow \mathcal{F} \subset \mathbb{R}^N$  satisfying Eq. (44). As long as the problem formulation depends only on the inner products between data points in  $\mathcal{F}$  and not on the data points themselves, without explicit representation of the feature mapping  $\Phi$  or the feature space  $\mathcal{F}$ , we can work on the feature space  $\mathcal{F}$  through relation (44). As

positive definite kernel functions satisfying Mercer's condition, polynomial kernel and Gaussian kernel

$$\kappa(x, y) = (\gamma_1(x \cdot y) + \gamma_2)^d, \quad d > 0 \text{ and } \gamma_1, \gamma_2 \in \mathbb{R},$$

$$\kappa(x, y) = \exp(-\|x - y\|^2/2\sigma^2), \quad \sigma \in \mathbb{R}$$

are in wide use.

In this section, we present the formulation of a generalized eigenvalue problem in the kernel-based feature space and apply the generalized LDA algorithms, obtaining nonlinear discriminant analysis. Given a kernel function  $\kappa$ , let  $\mathcal{S}_b$  and  $\mathcal{S}_w$  be the between-class and within-class scatter matrices in the feature space  $\mathcal{F} \subset \mathbb{R}^N$  which has been transformed by a mapping  $\Phi$  satisfying Eq. (44). Then the LDA in  $\mathcal{F}$  finds a linear transformation  $\mathcal{G} = [\varphi_1, \dots, \varphi_l] \in \mathbb{R}^{N \times l}$ , where the columns of  $\mathcal{G}$  are the generalized eigenvectors corresponding to the  $l$  largest eigenvalues of

$$\mathcal{S}_b \varphi = \lambda \mathcal{S}_w \varphi. \quad (45)$$

As in Eq. (7),  $\mathcal{S}_b$  and  $\mathcal{S}_w$  can be expressed as

$$\mathcal{S}_b = \mathcal{H}_b \mathcal{H}_b^T \quad \text{and} \quad \mathcal{S}_w = \mathcal{H}_w \mathcal{H}_w^T,$$

where

$$\mathcal{H}_b = [\sqrt{n_1}(\tilde{c}_1 - \tilde{c}), \dots, \sqrt{n_r}(\tilde{c}_r - \tilde{c})] \in \mathbb{R}^{N \times r}, \quad (46)$$

$$\mathcal{H}_w = [\Phi(A_1) - \tilde{c}_1 e_1, \dots, \Phi(A_r) - \tilde{c}_r e_r] \in \mathbb{R}^{N \times n},$$

$$\tilde{c}_i = \frac{1}{n_i} \sum_{j \in N_i} \Phi(a_j), \quad \tilde{c} = \frac{1}{n} \sum_{i=1}^n \Phi(a_i) \quad \text{and}$$

$$e_i = [1, \dots, 1] \in \mathbb{R}^{1 \times n_i}. \quad (47)$$

The notation  $\Phi(A_i)$  is used to denote  $\Phi(A_i) = \Phi([a_j, \dots, a_k]) = [\Phi(a_j), \dots, \Phi(a_k)]$ .

Let  $\varphi$  be represented as a linear combination of  $\Phi(a_i)$ 's such as  $\varphi = \sum_{i=1}^n u_i \Phi(a_i)$ , and define

$$u = [u_1, \dots, u_n]^T,$$

$$\mathcal{H}_b = [b_{ij}]_{(1 \leq i \leq n, 1 \leq j \leq r)},$$

$$b_{ij} = \sqrt{n_j} \left( \frac{1}{n_j} \sum_{p \in N_j} \kappa(a_i, a_p) - \frac{1}{n} \sum_{p=1}^n \kappa(a_i, a_p) \right). \quad (48)$$

Then we have

$$\mathcal{H}_b^T \varphi = \mathcal{H}_b^T u, \quad (49)$$

since

$$\begin{aligned} \mathcal{H}_b^T \varphi &= \begin{bmatrix} \sqrt{n_1}(\tilde{c}_1 - \tilde{c})^T \\ \vdots \\ \sqrt{n_r}(\tilde{c}_r - \tilde{c})^T \end{bmatrix} \left( \sum_{i=1}^n u_i \Phi(a_i) \right) \\ &= \begin{bmatrix} \sqrt{n_1} \left( \frac{1}{n_1} \sum_{p \in N_1} \Phi(a_p) - \frac{1}{n} \sum_{p=1}^n \Phi(a_p) \right)^T \\ \vdots \\ \sqrt{n_r} \left( \frac{1}{n_r} \sum_{p \in N_r} \Phi(a_p) - \frac{1}{n} \sum_{p=1}^n \Phi(a_p) \right)^T \end{bmatrix} \end{aligned}$$

$$[\Phi(a_1), \dots, \Phi(a_n)] \begin{bmatrix} u_1 \\ \vdots \\ u_n \end{bmatrix}$$

$$= \mathcal{H}_b^T u. \quad (50) \quad 35$$

Similarly, we can obtain

$$\mathcal{H}_w^T \varphi = \mathcal{H}_w^T u, \quad (50) \quad 37$$

where

$$\mathcal{H}_w = [w_{ij}]_{(1 \leq i \leq n, 1 \leq j \leq n)}, \quad (51) \quad 39$$

$$w_{ij} = \kappa(a_i, a_j) - \frac{1}{n_\delta} \sum_{p \in N_\delta} \kappa(a_i, a_p) \quad (51) \quad 41$$

when  $a_j$  belongs to the class  $\delta$ .

From Eqs. (49) and (50), for any  $\varphi = \sum_{i=1}^n u_i \Phi(a_i)$  and  $\psi = \sum_{i=1}^n v_i \Phi(a_i)$  we have

$$\begin{aligned} \mathcal{S}_b \varphi &= \lambda \mathcal{S}_w \varphi \\ \Leftrightarrow \psi^T \mathcal{H}_b \mathcal{H}_b^T \varphi &= \lambda \psi^T \mathcal{H}_w \mathcal{H}_w^T \varphi \\ \Leftrightarrow v^T \mathcal{H}_b \mathcal{H}_b^T u &= \lambda v^T \mathcal{H}_w \mathcal{H}_w^T u \\ \text{for } u &= [u_1, \dots, u_n]^T, \quad v = [v_1, \dots, v_n]^T \\ \Leftrightarrow \mathcal{H}_b \mathcal{H}_b^T u &= \lambda \mathcal{H}_w \mathcal{H}_w^T u. \end{aligned} \quad (52) \quad 45$$

Therefore, the generalized eigenvalue problem  $\mathcal{S}_b \varphi = \lambda \mathcal{S}_w \varphi$  becomes

$$\mathcal{H}_b \mathcal{H}_b^T u = \lambda \mathcal{H}_w \mathcal{H}_w^T u. \quad (53) \quad 47$$

### Algorithm 2. Nonlinear discriminant analysis.

Given a data matrix  $A = [a_1, \dots, a_n] \in \mathbb{R}^{m \times n}$  with  $r$  classes and a kernel function  $\kappa$ , it computes the  $l$  dimensional representation of any input vector  $z \in \mathbb{R}^{m \times 1}$  by applying the generalized LDA algorithm in the kernel-based feature space composed of the columns of  $K = [\kappa(a_i, a_j)]_{(1 \leq i \leq n, 1 \leq j \leq n)}$ .

Table 6

Prediction accuracies (%) by the classical LDA in the original space and the generalized LDA algorithms in the nonlinearly transformed feature space

Data	Dim.	No. data	Classes	Linear	Nonlinear methods				
				LDA	RLDA	LDA/GSVD	$To-N(S_w)$	$To-R(S_b)$	$To-NR(S_w)$
Musk	166	6599	2	91.2	97.6	<b>99.4</b>	<b>99.4</b>	89.2	99.3
Isolet	617	7797	26	93.9	95.8	96.8	97.0	89.7	<b>97.1</b>
Car	6	1728	4	88.2	94.7	94.1	94.9	84.5	<b>95.2</b>
Mfeature	649	2000	10	–	94.4	98.1	<b>98.3</b>	94.0	<b>98.3</b>
Bcancer	9	699	2	95.3	95.2	<b>96.4</b>	93.5	92.8	94.3
Bscale	4	625	3	87.0	<b>94.1</b>	86.5	86.5	86.5	86.1

In the Mfeature data set, the classical LDA was not applicable in the original space due to the singularity of the within-class scatter matrix.

- (1) Compute  $\mathcal{K}_b \in \mathbb{R}^{m \times r}$ ,  $\mathcal{K}_w \in \mathbb{R}^{n \times n}$  and  $\mathcal{K}_t \in \mathbb{R}^{n \times n}$  according to (48), (51) and (55).
- (2) Compute transformation matrix  $\mathcal{G}$  by applying the generalized LDA algorithms discussed in Section 2.
- (3) For any input vector  $z \in \mathbb{R}^{m \times 1}$ , a dimension reduced representation is computed by Eq. (57).

Note that  $\mathcal{K}_b \mathcal{K}_b^T$  and  $\mathcal{K}_w \mathcal{K}_w^T$  can be viewed as the between-class scatter matrix and within-class scatter matrix of the kernel matrix

$$K = [\kappa(a_i, a_j)]_{(1 \leq i \leq n, 1 \leq j \leq n)}, \quad (54)$$

when each column  $[\kappa(a_1, a_j), \dots, \kappa(a_n, a_j)]^T$  in  $K$  is considered as a data point in the  $n$ -dimensional space. It can be observed by comparing the structures of  $\mathcal{K}_b$  and  $\mathcal{K}_w$  with those of  $\mathcal{K}_b$  and  $\mathcal{K}_w$  in Eqs. (46)–(47). As in  $\mathcal{K}_b$  and  $\mathcal{K}_w$  of Eqs. (48) and (51),  $\mathcal{K}_t$  can be computed as

$$\mathcal{K}_t = [t_{ij}]_{(1 \leq i \leq n, 1 \leq j \leq n)},$$

$$t_{ij} = \kappa(a_i, a_j) - \frac{1}{n} \sum_{p=1}^n \kappa(a_i, a_p). \quad (55)$$

Since  $\mathcal{K}_b \mathcal{K}_b^T$  and  $\mathcal{K}_w \mathcal{K}_w^T$  are both singular in the feature space, the classical LDA cannot be applied for the generalized eigenvalue problem (53). Now we apply the generalized LDA algorithms discussed in Section 2 to solve Eq. (53), obtaining nonlinear discriminant analysis. Let

$$\mathcal{G} = [u^{(1)}, \dots, u^{(l)}] \in \mathbb{R}^{n \times l} \quad (56)$$

be the transformation matrix obtained by applying any generalized LDA algorithm in the feature space. Then the dimension reduced representation of any data item  $z \in \mathbb{R}^{m \times 1}$  is given by

$$\mathcal{G}^T \begin{bmatrix} \kappa(a_1, z) \\ \vdots \\ \kappa(a_n, z) \end{bmatrix} \in \mathbb{R}^{l \times 1}. \quad (57)$$

Algorithm 2 summarizes nonlinear extension of generalized LDA algorithms by kernel methods. 23

### 3.1. Experimental comparisons of nonlinear discriminant analysis algorithms 25

For this experiment, six data sets from UCI machine learning repository were used. By randomly splitting the data to the training and test set of equal size and repeating it 10 times, 10 pairs of training and test sets were constructed for each data. For the Bcancer and Bscale data sets, the ratio of training and test set was set as 4:1. Using the training set of the first pair among ten pairs and the nearest-neighbor classifier, five cross-validation was used in order to determine the optimal value for  $\sigma$  in the Gaussian kernel function  $\kappa(x, y) = \exp(-(\|x - y\|^2)/2\sigma^2)$ . After finding the optimal  $\sigma$  values, mean prediction accuracies from ten pairs of training and test sets were calculated and they are reported in Table 6. In the regularization method, while the regularization parameter was set as 1, the optimal  $\sigma$  value was searched by the cross-validation. Table 6 also reports the prediction accuracies by the classical LDA in the original data space and it demonstrates that nonlinear discriminant analysis can improve prediction accuracies compared with LDA. 27 29 31 33 35 37 39 41

Fig. 4 illustrates the computational complexities using the specific sizes of the training data used in Table 6. As in the comparison of the generalized LDA algorithms, the method  $To-R(S_b)$  [5] gives the lowest computational complexities among the compared methods. However, combining  $To-R(S_b)$  with kernel methods does not make effective nonlinear dimension reduction method as shown in Table 6. In the generalized eigenvalue problem, 43 45 47 49

$$\mathcal{K}_b \mathcal{K}_b^T u = \lambda \mathcal{K}_w \mathcal{K}_w^T u, \quad (51)$$

where  $\mathcal{K}_b \mathcal{K}_b^T, \mathcal{K}_w \mathcal{K}_w^T \in \mathbb{R}^{n \times n}$ . The data dimension is equal to the number of data and the rank of  $\mathcal{K}_w \mathcal{K}_w^T$  is not severely smaller than the data dimension. However, poor performances by  $To-R(S_b)$  demonstrate that the null space of  $\mathcal{K}_w \mathcal{K}_w^T$  contains discriminative information. Figs. 3 and 4 show that the proposed LDA/GSVD method can reduce greatly the computational cost of the original LDA/GSVD in both the original space and the feature space. 53 55 57 59

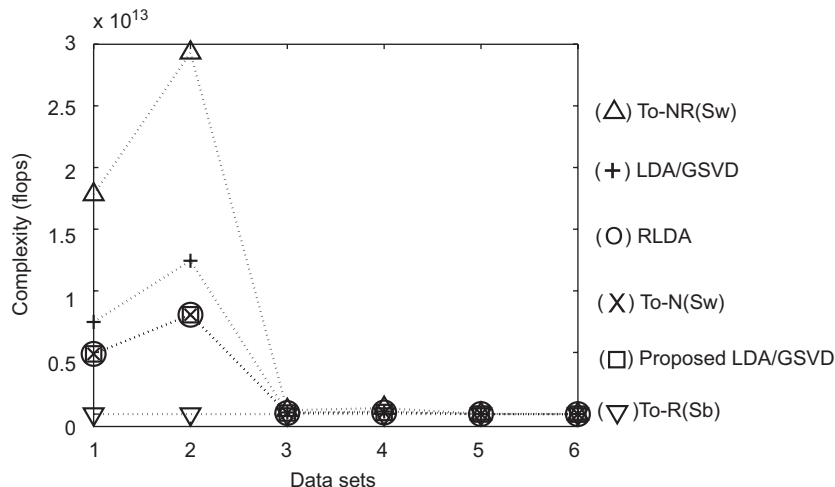


Fig. 4. The figures compare complexities required for the generalized LDA algorithms in the feature space for specific problem sizes of training data used in Table 6. From the left on  $x$ -axis, the data sets, Musk, Isolet, Car, Mfeature, Bcancer and Bscale are corresponded.

#### 4. Conclusions/Discussions

We presented the relationships among several generalized LDA algorithms developed for handling undersampled problems and compared their computational complexities and performances. As discussed in the theoretical comparison, many algorithms are closely related, and experimental results indicate that computational complexities are important issues in addition to classification performances. The LDA/GSVD showed competitive performances throughout the experiments, but the computational complexities can be expensive especially for high dimensional data. An efficient algorithm has been proposed, which produces the same solution as LDA/GSVD. The computational savings are remarkable especially for high dimensional data.

Nonlinear extensions of the generalized LDA algorithms by the formulation of generalized eigenvalue problem in the kernel-based feature space were presented. Experimental results using data sets from UCI database demonstrate that nonlinear discriminant analysis can improve prediction accuracies compared with LDA.

#### References

[1] K. Fukunaga, Introduction to Statistical Pattern Recognition, second ed., Academic Press, New York, 1990.

[2] R.O. Duda, P.E. Hart, D.G. Stork, Pattern Classification, Wiley-Interscience, New York, 2001.

[3] J.H. Friedman, Regularized discriminant analysis, *J. Am. Stat. Assoc.* 84 (405) (1989) 165–175.

[4] P. Howland, M. Jeon, H. Park, Structure preserving dimension reduction for clustered text data based on the generalized singular value decomposition, *SIAM J. Matrix Anal. Appl.* 25 (1) (2003) 165–179.

[5] L. Chen, H.M. Liao, M. Ko, J. Lin, G. Yu, A new LDA-based face recognition system which can solve the small sample size problem, *Pattern Recognition* 33 (2000) 1713–1726.

[6] H. Yu, J. Yang, A direct LDA algorithm for high-dimensional data— with application to face recognition, *Pattern Recognition* 34 (2001) 2067–2070.

[7] J. Yang, J.-Y. Yang, Why can LDA be performed in PCA transformed space?, *Pattern Recognition* 36 (2003) 563–566.

[8] W. Zheng, L. Zhao, C. Zou, An efficient algorithm to solve the small sample size problem for LDA, *Pattern Recognition* 37 (2004) 1077–1079.

[9] N. Cristianini, J. Shawe-Taylor, An Introduction to Support Vector Machines and other Kernel-based Learning Methods, Cambridge, 2000.

[10] P. Howland, H. Park, Generalizing discriminant analysis using the generalized singular value decomposition, *IEEE Trans. Pattern Anal. Mach. Intell.* 26 (8) (2004) 995–1006.

[11] C.C. Paige, M.A. Saunders, Towards a generalized singular value decomposition, *SIAM J. Numer. Anal.* 18 (1981) 398–405.

[12] G.H. Golub, C.F. Van Loan, Matrix Computations, third ed., Johns Hopkins University Press, Baltimore, MD, 1996.

[13] P.N. Belhumeur, J.P. Hespanha, D.J. Kriegman, Eigenfaces vs. fisherfaces: recognition using class specific linear projection, *IEEE Trans. Pattern Anal. Mach. Learn.* 19 (7) (1997) 711–720.

[14] W. Zheng, C. Zou, L. Zhao, Real-time face recognition using Gram–Schmidt orthogonalization for LDA, in: The Proceedings of the 17th International Conference on Pattern Recognition, 2004.

[15] Z. Jin, J.-Y. Yang, Z.-S. Hu, Z. Lou, Face recognition based on the uncorrelated discriminant transformation, *Pattern Recognition* 34 (2001) 1405–1416.

[16] J. Ye, R. Janardan, Q. Li, H. Park, Feature extraction via generalized uncorrelated linear discriminant analysis, in: The Proceedings of the 21st International Conference on Machine Learning, 2004.

[17] T.G. Kolda, D.P. O’Leary, A semidiscrete matrix decomposition for latent semantic indexing in information retrieval, *ACM Trans. Inf. Systems* 16 (4) (1998) 322–346.

[18] (<http://www.uk.research.att.com/facedatabase.html>).

[19] (<http://cvc.yale.edu/projects/yalefaces/yalefaces.html>).

[20] B. Scholkopf, A.J. Smola, K.-R. Muller, Nonlinear component analysis as a kernel eigenvalue problem, *Neural Comput.* 10 (1998) 1299–1319.

[21] S. Mika, G. Ratsch, J. Weston, B. Scholkopf, K.-R. Muller, Fisher discriminant analysis with kernels, in: E. Wilson, J. Larsen, S. Douglas (Eds.), *Neural Networks for Signal Processing*, vol. IX, IEEE, New York, 1999, p. 41–48.

[22] G. Baudat, F. Anouar, Generalized discriminant analysis using a kernel approach, *Neural Comput.* 12 (2000) 2385–2404.

[23] V. Roth, V. Steinhage, Nonlinear discriminant analysis using kernel functions, *Adv. Neural Inf. Process. Systems* 12 (2000) 568–574.

[24] S.A. Billings, K.L. Lee, Nonlinear fisher discriminant analysis using a minimum squared error cost function and the orthogonal least squares algorithm, *Neural Networks* 15 (2) (2002) 263–270.

- 1 [25] C.H. Park, H. Park, Nonlinear discriminant analysis using kernel  
3 functions and the generalized singular value decomposition, *SIAM J.*  
4 *Matrix Anal. Appl.* 27 (1) (2005) 87–102.
- 5 [26] B. Scholkopf, S. Mika, C.J.C. Burges, P. Knirsch, K.-R. Muller, G.  
Ratsch, A.J. Smola, Input space versus feature space in kernel-based  
methods, *IEEE Trans. Neural Networks* 10 (5) (1999) 1000–1017.
- [27] C.J.C. Burges, A tutorial on support vector machines for pattern  
recognition, *Data Min. Knowl. Discovery* 2 (2) (1998) 121–167.

**About the Author**—CHEONG HEE PARK received her Ph.D. in Mathematics from Yonsei University, Korea in 1998. She received the M.S. and Ph.D. degrees in Computer Science at the Department of Computer Science and Engineering, University of Minnesota in 2002 and 2004, respectively. She is currently in the Department of Computer Science and Engineering, Chungnam National University, Korea as an Assistant Professor. Her research interests include pattern recognition, data mining, bioinformatics and machine learning.

**About the Author**—HAESUN PARK received her B.S. degree in Mathematics from Seoul National University, Seoul Korea, in 1981 with summa cum laude and the University President's Medal for the top graduate, and her M.S. and Ph.D. degrees in Computer Science from Cornell University, Ithaca, NY, in 1985 and 1987, respectively. She was on the faculty of the Department of Computer Science and Engineering, University of Minnesota, Twin Cities, from 1987 to 2005. From 2003 to 2005, she served as a Program Director, the Computing and Communication Foundations Division at the National Science Foundation, Arlington, VA, USA. Since July 2005, she has been a Professor in Division of Computational Science and Engineering, College of Computing, Georgia Institute of Technology, Atlanta, GA. Her current research interests include numerical algorithms, pattern recognition, bioinformatics, information retrieval, and data mining. She has published over 100 research papers in these areas. Prof. Park served on numerous conference committees and editorial boards of journals. Currently she is on the editorial board of *BIT Numerical Mathematics*, *SIAM Journal on Matrix Analysis and Applications*, and *International Journal of Bioinformatics Research and Applications*.

UNCORRECTED PROOF