

# Speed Up Kernel Discriminant Analysis

Deng Cai · Xiaofei He · Jiawei Han

Received: 13 December 2009 / Revised: 7 April 2010 / Accepted: 16 April 2010

**Abstract** Linear Discriminant Analysis (LDA) has been a popular method for dimensionality reduction which preserves class separability. The projection vectors are commonly obtained by maximizing the between class covariance and simultaneously minimizing the within class covariance. LDA can be performed either in the original input space or in the reproducing kernel Hilbert space (RKHS) into which data points are mapped, which leads to Kernel Discriminant Analysis (KDA). When the data are highly nonlinear distributed, KDA can achieve better performance than LDA. However, computing the projective functions in KDA involves eigen-decomposition of kernel matrix, which is very expensive when a large number of training samples exist. In this paper, we present a new algorithm for kernel discriminant analysis, called *Spectral Regression Kernel Discriminant Analysis* (SRKDA). By using spectral graph analysis, SRKDA casts discriminant analysis into a regression framework which facilitates both efficient computation and the use of regularization techniques. Specifically, SRKDA only needs to solve a set of regularized regression problems and there is no eigenvector computation involved, which is a huge save of computational cost. The new formulation makes it very easy to develop incremental version of the algorithm which can

fully utilize the computational results of the existing training samples. Moreover, it is easy to produce sparse projections (Sparse KDA) with a  $L_1$ -norm regularizer. Extensive experiments on spoken letter, handwritten digit image and face image data demonstrate the effectiveness and efficiency of the proposed algorithm.

**Keywords** Kernel Discriminant Analysis · Regression · Subspace Learning · Dimensionality Reduction

## 1 Introduction

Dimensionality reduction has been a key problem in many fields of information processing, such as data mining, information retrieval, and pattern recognition. When data is represented as points in a high-dimensional space, one is often confronted with tasks like nearest neighbor search. Many methods have been proposed to index the data for fast query response, such as  $K$ - $D$  tree,  $R$  tree,  $R^*$  tree, etc [16]. However, these methods can only operate with small dimensionality, typically less than 100. The effectiveness and efficiency of these methods drop exponentially as the dimensionality increases, which is commonly referred to as the “curse of dimensionality”.

During the last decade, with the advances in computer technologies and the advent of the World Wide Web, there has been an explosion in the amount and complexity of digital data being generated, stored, analyzed, and accessed. Much of this information is multimedia in nature, including text, image, and video data. The multimedia data are typically of very high dimensionality, ranging from several thousands to several hundreds of thousands. Learning with such high dimensionality in many cases is almost infeasible. Thus, learnability necessitates dimensionality reduction [8, 30, 33, 34].

---

Deng Cai · Xiaofei He  
State Key Lab of CAD&CG,  
College of Computer Science,  
Zhejiang University  
Tel.: +86-571-88206681  
Fax: +86-571-88206680  
E-mail: {dengcai,xiaofeihe}@cad.zju.edu.cn

Jiawei Han  
Department of Computer Science,  
University of Illinois at Urbana Champaign  
E-mail: hanj@cs.uiuc.edu

Once the high-dimensional data is mapped into lower-dimensional space, conventional indexing schemes can then be applied [20, 28, 29].

One of the most popular dimensionality reduction algorithms might be Linear Discriminant Analysis (LDA). LDA is a supervised method that has proved successful on classification problems [9, 15]. The projection vectors are commonly obtained by maximizing the between class covariance and simultaneously minimizing the within class covariance. The classical LDA is a linear method and fails for nonlinear problems. To deal with this limitation, nonlinear extensions of LDA through “kernel trick” have been proposed. The main idea of kernel-based methods is to map the input data to a feature space through a nonlinear mapping, where the inner products in the feature space can be computed by a kernel function without knowing the nonlinear mapping explicitly [27]. Kernel Fisher Discriminant Analysis (KFD) in [22] and Generalized Discriminant Analysis (GDA) in [1] are two independently developed approaches for kernel-based nonlinear extensions of LDA. They are essentially equivalent. To avoid confusion, we will refer this approach as Kernel Discriminant Analysis (KDA) hereafter.

When solving the optimization problem of KDA, we need to handle the possible singularity problem of the total scatter matrix. There are two approaches trying to address this issue either by using regularization techniques [22] or by applying singular value decomposition [1, 26][35]. Both of these two approaches for solving optimization problem of KDA involve the eigen-decomposition of the kernel matrix which is computationally expensive. Moreover, due to the difficulty of designing an incremental solution for the eigen-decomposition on the kernel matrix, there has been little work on designing incremental KDA algorithms that can efficiently incorporate new data examples as they become available.

In [23], S. Mika et al. made a first attempt to speed up KDA through a greedy approximation technique. However, their algorithm was developed to handle the binary classification problem. For a multi-class problem, the authors suggested the one against the rest scheme by considering all two-class problems. Recent studies [4–7, 9] show that various linear dimensionality reduction algorithms can be formulated as regression problems and thus have efficient computational solutions. Particularly, our previous work [6] has demonstrated that LDA can be formulated as a regression problem and be efficiently solved. The similar idea has been applied to unsupervised dimensionality reduction algorithms [7] and semi-supervised dimensionality reduction algorithms [5]. However, it is not clear that how

these similar techniques can be applied on non-linear dimensionality reduction algorithms which use kernel techniques.

In this paper, we propose a new algorithm for kernel discriminant analysis, called *Spectral Regression Kernel Discriminant Analysis* (SRKDA). Our analysis essentially follows our previous idea for speeding up LDA [6]. By using spectral graph analysis, SRKDA casts discriminant analysis into a regression framework which facilitates both efficient computation and the use of regularization techniques. Specifically, SRKDA only needs to solve a set of regularized regression problems and there is no eigenvector computation involved, which is a huge save of computational cost. Moreover, the new formulation makes it very easy to develop incremental version of the algorithm which can fully utilize the previous computational results on the existing training samples.

The points below highlight the contributions of this paper:

- KDA in the binary-class case has been shown to be equivalent to regularized kernel regression with the class label as the output [27]. Our paper extends this relation to multi-class case.
- We provides a new formulation of KDA optimization problem. With this new formulation, the KDA optimization problem can be efficiently solved by avoiding the eigen-decomposition of the kernel matrix. Theoretical analysis shows that the new approach can achieve 27-times speedup over the ordinary KDA approaches.
- Moreover, SRKDA can be naturally performed in the incremental manner. The computational results on the existing training samples can be fully utilized when new training samples are injected into the system. Theoretical analysis shows that SRKDA in the incremental mode has only quadratic-time complexity, which is a huge improvement comparing to the cubic-time complexity of the ordinary KDA approaches.
- Since SRKDA uses regression as a building block, various kinds of regularization techniques can be easily incorporated (*e.g.*,  $L_1$ -norm regularizer to produce sparse projections). Our approach provides a huge possibility to develop new variations of kernel discriminant analysis.
- A short version of this work has been published in ICDM [3]. In this journal version, we provide two new sections on theoretical analysis (Section 3.1) and sparse KDA (Section 5). Moreover, we have added significant amount of the experimental results.

The remainder of the paper is organized as follows. In Section 2, we provide a brief review of LDA and KDA, plus a detailed computational analysis of KDA. Section 3 introduces our proposed *Spectral Regression Kernel Discriminant Analysis* algorithm. The incremental version of SRKDA is introduced in Section 4 and the extensive experimental results are presented in Section 5. Finally, we provide some concluding remarks in Section 6.

## 2 A Brief Review of LDA and KDA

Linear Discriminant Analysis (LDA) seeks directions on which the data points of different classes are far from each other while requiring data points of the same class to be close to each other [15]. Suppose we have a set of  $m$  samples  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m \in \mathbb{R}^n$ , belonging to  $c$  classes. The objective function of LDA is as follows:

$$\mathbf{a}_{opt} = \arg \max \frac{\mathbf{a}^T S_b \mathbf{a}}{\mathbf{a}^T S_w \mathbf{a}}, \quad (1)$$

$$S_b = \sum_{k=1}^c m_k (\boldsymbol{\mu}^{(k)} - \boldsymbol{\mu})(\boldsymbol{\mu}^{(k)} - \boldsymbol{\mu})^T,$$

$$S_w = \sum_{k=1}^c \left( \sum_{i=1}^{m_k} (\mathbf{x}_i^{(k)} - \boldsymbol{\mu}^{(k)})(\mathbf{x}_i^{(k)} - \boldsymbol{\mu}^{(k)})^T \right),$$

where  $\boldsymbol{\mu}$  is the global centroid,  $m_k$  is the number of samples in the  $k$ -th class,  $\boldsymbol{\mu}^{(k)}$  is the centroid of the  $k$ -th class, and  $\mathbf{x}_i^{(k)}$  is the  $i$ -th sample in the  $k$ -th class. We call  $S_w$  the within-class scatter matrix and  $S_b$  the between-class scatter matrix.

Define the total scatter matrix  $S_t = \sum_{i=1}^m (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T$ , we have  $S_t = S_b + S_w$  [15]. The objective function of LDA in Eqn. (1) is equivalent to

$$\mathbf{a}_{opt} = \arg \max \frac{\mathbf{a}^T S_b \mathbf{a}}{\mathbf{a}^T S_t \mathbf{a}}. \quad (2)$$

The optimal  $\mathbf{a}$ 's are the eigenvectors corresponding to the non-zero eigenvalue of eigen-problem:

$$S_b \mathbf{a} = \lambda S_t \mathbf{a}. \quad (3)$$

Since the rank of  $S_b$  is bounded by  $c - 1$ , there are at most  $c - 1$  eigenvectors corresponding to non-zero eigenvalues [15].

For the sake of clarity, we provide in a summary for notations usage in Table 1. To extend LDA to the nonlinear case, we consider the problem in a feature space  $\mathcal{F}$  induced by some nonlinear mapping

$$\phi : \mathbb{R}^n \rightarrow \mathcal{F}$$

**Table 1** Notations

Notations	Descriptions
$m$	the number of total training data points
$n$	the number of features
$c$	the number of classes
$m_k$	the number of data points in $k$ -th class
$\mathbf{x}_i$	the $i$ -th data point
$\mathbf{x}_i^{(k)}$	the $i$ -th data point in the $k$ -th class
$\boldsymbol{\mu}$	the total sample mean vector
$\boldsymbol{\mu}^{(k)}$	the mean vector of the $k$ -th class
$X$	the data matrix
$S_b$	the between-class scatter matrix
$S_w$	the within-class scatter matrix
$S_t$	the total scatter matrix
$\boldsymbol{\alpha}$	the transformation vector

For a proper chosen  $\phi$ , an inner product  $\langle \cdot, \cdot \rangle$  can be defined on  $\mathcal{F}$  which makes for a so-called reproducing kernel Hilbert space (RKHS). More specifically,

$$\langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle = \mathcal{K}(\mathbf{x}, \mathbf{y})$$

holds where  $\mathcal{K}(\cdot, \cdot)$  is a positive semi-definite kernel function. Several popular kernel functions are: Gaussian kernel  $\mathcal{K}(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|^2 / 2\sigma^2)$ ; polynomial kernel  $\mathcal{K}(\mathbf{x}, \mathbf{y}) = (1 + \mathbf{x}^T \mathbf{y})^d$ ; Sigmoid kernel  $\mathcal{K}(\mathbf{x}, \mathbf{y}) = \tanh(\mathbf{x}^T \mathbf{y} + \alpha)$ .

Let  $S_b^\phi$ ,  $S_w^\phi$  and  $S_t^\phi$  denote the between-class, within-class and total scatter matrices in the feature space, respectively. We have

$$S_b^\phi = \sum_{k=1}^c m_k (\boldsymbol{\mu}_\phi^{(k)} - \boldsymbol{\mu}_\phi)(\boldsymbol{\mu}_\phi^{(k)} - \boldsymbol{\mu}_\phi)^T,$$

$$S_w^\phi = \sum_{k=1}^c \left( \sum_{i=1}^{m_k} (\phi(\mathbf{x}_i^{(k)}) - \boldsymbol{\mu}_\phi^{(k)})(\phi(\mathbf{x}_i^{(k)}) - \boldsymbol{\mu}_\phi^{(k)})^T \right),$$

$$S_t^\phi = \sum_{i=1}^m (\phi(\mathbf{x}_i) - \boldsymbol{\mu}_\phi)(\phi(\mathbf{x}_i) - \boldsymbol{\mu}_\phi)^T,$$

where  $\boldsymbol{\mu}_\phi^{(k)}$  and  $\boldsymbol{\mu}_\phi$  are the centroids of the  $k$ -th class and the global centroid, respectively, in the feature space.

Let  $\boldsymbol{\nu}$  denote the projective function in the feature space, the corresponding objective function (2) in the feature space is

$$\boldsymbol{\nu}_{opt} = \arg \max \frac{\boldsymbol{\nu}^T S_b^\phi \boldsymbol{\nu}}{\boldsymbol{\nu}^T S_t^\phi \boldsymbol{\nu}}, \quad (4)$$

which can be solved by the eigen-problem:

$$S_b^\phi \boldsymbol{\nu} = \lambda S_t^\phi \boldsymbol{\nu}.$$

Because the eigenvectors are linear combinations of  $\phi(\mathbf{x}_i)$  [1][27], there exist coefficients  $\alpha_i$  such that

$$\boldsymbol{\nu} = \sum_{i=1}^m \alpha_i \phi(\mathbf{x}_i).$$

Let  $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_m]^T$ , it can be proved [1] that Eqn. (4) is equivalent to:

$$\boldsymbol{\alpha}_{opt} = \arg \max \frac{\boldsymbol{\alpha}^T K W K \boldsymbol{\alpha}}{\boldsymbol{\alpha}^T K K \boldsymbol{\alpha}}, \quad (5)$$

and the corresponding eigen-problem is:

$$K W K \boldsymbol{\alpha} = \lambda K K \boldsymbol{\alpha}. \quad (6)$$

where  $K$  is the kernel matrix ( $K_{ij} = \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)$ ) and  $W$  is defined as:

$$W_{ij} = \begin{cases} 1/m_k, & \text{if } \mathbf{x}_i \text{ and } \mathbf{x}_j \text{ both belong to} \\ & \text{the } k\text{-th class;} \\ 0, & \text{otherwise.} \end{cases} \quad (7)$$

Each eigenvector  $\boldsymbol{\alpha}$  gives a projective function  $\boldsymbol{\nu}$  in the feature space. For a data example  $\mathbf{x}$ , we have

$$\begin{aligned} \langle \boldsymbol{\nu}, \phi(\mathbf{x}) \rangle &= \sum_{i=1}^m \alpha_i \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}) \rangle \\ &= \sum_{i=1}^m \alpha_i \mathcal{K}(\mathbf{x}_i, \mathbf{x}) \\ &= \boldsymbol{\alpha}^T K(:, \mathbf{x}) \end{aligned}$$

where  $K(:, \mathbf{x}) \doteq [K(\mathbf{x}_1, \mathbf{x}), \dots, K(\mathbf{x}_m, \mathbf{x})]^T$ . Let  $\{\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_{c-1}\}$  be the  $c-1$  eigenvectors of the eigen-problem in Eqn. (6) with respect to the non-zero eigenvalues. The transformation matrix  $\Theta = [\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_{c-1}]$  is an  $m \times (c-1)$  matrix and a data sample  $\mathbf{x}$  can be embedded into  $c-1$  dimensional subspace by

$$\mathbf{x} \rightarrow \mathbf{z} = \Theta^T K(:, \mathbf{x}).$$

The above approach extends LDA into RKHS by using ‘‘kernel trick’’ is independently developed by Mika *et al.* [22] and Baudat *et al.* [1]. This algorithm was named as Kernel Fisher Discriminant (KFD) in [22] and Generalized Discriminant Analysis (GDA) in [1].

## 2.1 Computational Analysis of KDA

To get a stable solution of the eigen-problem in Eqn. (6), the matrix  $KK^T$  is required to be non-singular [17]. When  $K$  is singular, there are two methods to solve this problem. The first method is by using eigen-decomposition of  $K$ , which was proposed in [1].

Suppose the rank of  $K$  is  $r (r \leq m)$  and the eigen-decomposition of  $K$  is as follows:

$$K = U \Sigma U^T = U_r \Sigma_r U_r^T$$

where  $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_m)$  is the diagonal matrix of sorted eigenvalues ( $\sigma_1 \geq \dots \geq \sigma_m \geq 0$ ) and  $U$  is the matrix of normalized eigenvectors associated to  $\Sigma$ .

$\Sigma_r = \text{diag}(\sigma_1, \dots, \sigma_r)$  is the diagonal matrix of nonzero eigenvalues and  $U_r$  is the first  $r$  columns of  $U$ . Thus  $\Sigma_r^{-1}$  exists and  $U_r^T U_r = I$ , where  $I$  is the identity matrix.

Substituting  $K$  in Eqn. (5), we get

$$\boldsymbol{\alpha}_{opt} = \arg \max \frac{(\Sigma_r U_r^T \boldsymbol{\alpha})^T U_r^T W U_r (\Sigma_r U_r^T \boldsymbol{\alpha})}{(\Sigma_r U_r^T \boldsymbol{\alpha})^T U_r^T U_r (\Sigma_r U_r^T \boldsymbol{\alpha})}.$$

We proceed to variable modification using  $\boldsymbol{\beta} = \Sigma_r U_r^T \boldsymbol{\alpha}$  and get:

$$\boldsymbol{\beta}_{opt} = \arg \max \frac{\boldsymbol{\beta}^T U_r^T W U_r \boldsymbol{\beta}}{\boldsymbol{\beta}^T \boldsymbol{\beta}},$$

Thus, the optimal  $\boldsymbol{\beta}$ 's are the leading eigenvectors of matrix  $U_r^T W U_r$ . Once  $\boldsymbol{\beta}$ 's are calculated,  $\boldsymbol{\alpha}$  can be computed as  $\boldsymbol{\alpha} = U_r \Sigma_r^{-1} \boldsymbol{\beta}$ .

The second method is using the idea of regularization, by adding constant values to the diagonal elements of  $KK$ , as  $KK + \gamma I$ , for  $\gamma > 0$ . It is easy to see that  $KK + \gamma I$  is nonsingular. This method is used in [22]. By noticing that

$$KK + \gamma I = U \Sigma U^T U \Sigma U^T + \gamma I = U (\Sigma^2 + \gamma I) U^T,$$

we define  $\tilde{\Sigma} = (\Sigma^2 + \gamma I)^{1/2}$ , the objective function of regularized KDA can be written as:

$$\begin{aligned} & \max \frac{\boldsymbol{\alpha}^T K W K \boldsymbol{\alpha}}{\boldsymbol{\alpha}^T (K K + \gamma I) \boldsymbol{\alpha}} \\ &= \max \frac{\boldsymbol{\alpha}^T U \Sigma U^T W U \Sigma U^T \boldsymbol{\alpha}}{\boldsymbol{\alpha}^T U \tilde{\Sigma} \tilde{\Sigma} U^T \boldsymbol{\alpha}} \\ &= \max \frac{\boldsymbol{\beta}^T \tilde{\Sigma}^{-1} \Sigma U^T W U \Sigma \tilde{\Sigma}^{-1} \boldsymbol{\beta}}{\boldsymbol{\beta}^T \boldsymbol{\beta}} \end{aligned}$$

where  $\boldsymbol{\beta} = \tilde{\Sigma} U^T \boldsymbol{\alpha}$ . The optimal  $\boldsymbol{\beta}$ 's are the leading eigenvectors of matrix  $\tilde{\Sigma}^{-1} \Sigma U^T W U \Sigma \tilde{\Sigma}^{-1}$ . With this formulation, the above two methods can be computed in exactly the same way.

To reduce the computation in calculating  $\boldsymbol{\beta}$ , we shall exploit the special structure of  $W$ . Without loss of generality, we assume that the data points are ordered according to their labels. It is easy to check that the matrix  $W$  has a block-diagonal structure

$$W = \begin{bmatrix} W^{(1)} & 0 & \dots & 0 \\ 0 & W^{(2)} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & W^{(c)} \end{bmatrix} \quad (8)$$

where  $\{W^{(k)}\}_{k=1}^c$  is an  $m_k \times m_k$  matrix with all the elements equal to  $1/m_k$ .

We partition the  $m \times r$  matrix  $U_r$  as  $[U_r^{(1)}, \dots, U_r^{(c)}]^T$ , where  $U_r^{(k)} \in \mathbb{R}^{r \times m_k}$ . Let  $\mathbf{v}_i^{(k)}$  be the  $i$ -th column vector of  $U_r^{(k)}$ , we have:

$$\begin{aligned} U_r^T W U_r &= \sum_{k=1}^c U_r^{(k)} W^{(k)} (U_r^{(k)})^T \\ &= \sum_{k=1}^c \frac{1}{m_k} \left( \sum_{i=1}^{m_k} \mathbf{v}_i^{(k)} \sum_{i=1}^{m_k} (\mathbf{v}_i^{(k)})^T \right) \\ &= \sum_{k=1}^c m_k \bar{\mathbf{v}}^{(k)} (\bar{\mathbf{v}}^{(k)})^T \\ &= H H^T \end{aligned}$$

where  $H = [\sqrt{m_1} \bar{\mathbf{v}}^{(1)}, \dots, \sqrt{m_c} \bar{\mathbf{v}}^{(c)}] \in \mathbb{R}^{r \times c}$  and  $\bar{\mathbf{v}}^{(k)}$  is the average vector of  $\mathbf{v}_i^{(k)}$ .

To calculate the  $c$  leading eigenvectors of  $H H^T$ , it is not necessary to work on matrix  $H H^T$  which is of size  $r \times r$ . We can use a much more efficient algorithm. Suppose the Singular Value Decomposition of  $H$  is

$$H = P \Gamma Q^T,$$

it is easy to check that the column vectors of  $P$  are the eigenvectors of  $H H^T$  and the column vectors of  $Q$  are the eigenvectors of  $H^T H$  [32]. Moreover, if  $P$  or  $Q$  is given, we can recover the other via the formula  $H Q = P \Gamma$  and  $P^T H = \Gamma Q^T$ . Since  $c \ll r$ , we can calculate the  $c$  eigenvectors of  $H^T H$  and then recover the eigenvectors of  $H H^T$ , which are  $\beta$ 's.

We use the term *flam* [31], a compound operation consisting of one addition and one multiplication, to measure the operation counts. All the kernel methods need to compute the kernel matrix  $K$  which requires  $O(m^2 n)$  flam, where  $n$  is the number of features. The eigen-decomposition of  $K$  requires  $\frac{9}{2} m^3$  flam [32, 17]; Calculating the  $c - 1$  eigenvectors  $\beta$ 's requires  $\frac{9}{2} c^3 + \frac{3}{2} m c^2$  flam; Computing  $\alpha$ 's from  $\beta$ 's requires  $m^2 c$  flam. Finally, we conclude the time complexity of KDA measured by flam is

$$\frac{9}{2} m^3 + m^2 c + O(m^2 n) + \frac{3}{2} m c^2 + \frac{9}{2} c^3.$$

Considering  $m \gg c$ , the above time complexity can be simplified as

$$\frac{9}{2} m^3 + m^2 c + O(m^2 n). \quad (9)$$

For a large scale problem, we have  $m \gg n$ . Thus, the time complexity of KDA is determined by  $\frac{9}{2} m^3$ , which is the cost of eigen-decomposition of size  $m \times m$  kernel matrix  $K$ .

### 3 Efficient KDA via Spectral Regression

In order to solve the KDA eigen-problem in Eqn. (6) efficiently, we use the following theorem:

**Theorem 1** Let  $\mathbf{y}$  be the eigenvector of eigen-problem  $W \mathbf{y} = \lambda \mathbf{y}$  (10)

with eigenvalue  $\lambda$ . If  $K \alpha = \mathbf{y}$ , then  $\alpha$  is the eigenvector of eigen-problem in Eqn. (6) with the same eigenvalue  $\lambda$ .

*Proof* We have  $W \mathbf{y} = \lambda \mathbf{y}$ . At the left side of Eqn. (6), replace  $K \alpha$  by  $\mathbf{y}$ , we have

$$K W K \alpha = K W \mathbf{y} = K \lambda \mathbf{y} = \lambda K \mathbf{y} = \lambda K K \alpha$$

Thus,  $\alpha$  is the eigenvector of eigen-problem Eqn. (6) with the same eigenvalue  $\lambda$ .

Theorem 1 shows that instead of solving the eigen-problem Eqn. (6), the KDA projective functions can be obtained through two steps:

1. Solve the eigen-problem in Eqn. (10) to get  $\mathbf{y}$ .
2. Find  $\alpha$  which satisfies  $K \alpha = \mathbf{y}$ . The kernel matrix  $K$  is positive semi-definite. When  $K$  is non-singular (positive definite), for any given  $\mathbf{y}$ , we have a unique  $\alpha = K^{-1} \mathbf{y}$  which satisfies the above linear equations system. When  $K$  is singular, the system may have no solution or have infinite many solutions (the linear equations system is underdetermined) [17]. A possible way is to approximate  $\alpha$  by solving the following linear equations:

$$(K + \delta I) \alpha = \mathbf{y} \quad (11)$$

where  $I$  is the identity matrix and  $\delta \geq 0$  is the regularization parameter.

The advantages of this two-step approach are as follows:

1. We will show later how the eigen-problem in Eqn. (10) is *trivial* and we can directly get those eigenvectors  $\mathbf{y}$ .
2. The eigen-decomposition of  $K$  is avoided. Since the matrix  $K + \delta I$  is positive definite, the Cholesky decomposition can be used to efficiently solve the linear equations in Eqn. (11) [17], [31]. The computational complexity analysis will be provided in the later section.

The linear equations system in Eqn. (11) has a close connection with regularized regression [36]. We denote the projective function in the feature space as:

$$f(\mathbf{x}) = \langle \boldsymbol{\nu}, \phi(\mathbf{x}) \rangle = \sum_{i=1}^m \alpha_i K(\mathbf{x}, \mathbf{x}_i)$$

It can be easily verified that the solution  $\mathbf{\alpha}^* = (K + \delta I)^{-1} \mathbf{y}$  given by equations in Eqn. (11) is the optimal solution of the following regularized regression problem [36]:

$$\min_{f \in \mathcal{F}} \sum_{i=1}^m (f(\mathbf{x}_i) - y_i)^2 + \delta \|f\|_K^2 \quad (12)$$

where  $y_i$  is the  $i$ -th element of  $\mathbf{y}$ ,  $\mathcal{F}$  is the RKHS associated with Mercer kernel  $\mathcal{K}$  and  $\|\cdot\|_K$  is the corresponding norm.

Now let us analyze the eigenvectors of  $W$  which is defined in Eqn. (7) and (8). The  $W$  is block-diagonal, thus, its eigenvalues and eigenvectors are the union of the eigenvalues and eigenvectors of its blocks (the latter padded appropriately with zeros). It is straightforward to show that  $W^{(k)}$  has eigenvector  $\mathbf{e}^{(k)} \in \mathbb{R}^{m_k}$  associated with eigenvalue 1, where  $\mathbf{e}^{(k)} = [1, 1, \dots, 1]^T$ . Also there is only one non-zero eigenvalue of  $W^{(k)}$  because the rank of  $W^{(k)}$  is 1. Thus, there are exactly  $c$  eigenvectors of  $W$  with the same eigenvalue 1. These eigenvectors are

$$\mathbf{y}_k = \left[ \underbrace{0, \dots, 0}_{\sum_{i=1}^{k-1} m_i}, \underbrace{1, \dots, 1}_{m_k}, \underbrace{0, \dots, 0}_{\sum_{i=k+1}^c m_i} \right]^T \quad k = 1, \dots, c \quad (13)$$

Since 1 is a repeated eigenvalue of  $W$ , we can just pick any other  $c$  orthogonal vectors in the space spanned by  $\{\mathbf{y}_k\}$ , and define them to be our  $c$  eigenvectors. The vector of all ones  $\mathbf{e}$  is naturally in the spanned space. This vector is useless since the corresponding projective function will embed all the samples to the same point. Therefore, we pick  $\mathbf{e}$  as our first eigenvector of  $W$  and use Gram-Schmidt process to orthogonalize the remaining eigenvectors. The vector  $\mathbf{e}$  can then be removed, which leaves us exactly  $c - 1$  eigenvectors of  $W$ . We denote them as:

$$\{\bar{\mathbf{y}}_k\}_{k=1}^{c-1}, (\bar{\mathbf{y}}_k^T \mathbf{e} = 0, \bar{\mathbf{y}}_i^T \bar{\mathbf{y}}_j = 0, i \neq j) \quad (14)$$

The above two-step approach essentially combines the spectral analysis of the matrix  $W$  and regression techniques. Therefore, we name this new approach as *Spectral Regression Kernel Discriminant Analysis* (SRKDA) [8]:

In the following several subsections, we will provide the theoretical and computational analysis on SRKDA. Please see [6] for applying the similar technique on Linear Discriminant Analysis to obtain an efficient algorithm.

### 3.1 Theoretical Analysis

SRKDA calculates the projective functions through the linear equations system in Eqn. (11). When the kernel

matrix  $K$  is positive definite and the  $\delta = 0$ , Theorem 1 shows that the  $c - 1$  solutions  $\mathbf{\alpha}_k = K^{-1} \mathbf{y}_k$  are exactly the eigenvectors of the KDA eigen-problem in Eqn. (6) with respect to the eigenvalue 1. In this case, SRKDA is equivalent to ordinary KDA. Thus, it is interesting and important to see when the positive semi-definite kernel matrix  $K$  will be positive definite.

One of the most popular kernels is the Gaussian RBF kernel,  $\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / 2\sigma^2)$ . Our discussion in this section will only focus on Gaussian kernel. Regarding the Gaussian kernel, we have the following lemma:

**Lemma 1 (Full Rank of Gaussian RBF Gram Matrices [21])** *Suppose that  $\mathbf{x}_1, \dots, \mathbf{x}_m$  are distinct points, and  $\sigma \neq 0$ . The matrix  $K$  given by*

$$K_{ij} = \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / 2\sigma^2)$$

*has full rank.*

*Proof* See [21] and Theorem 2.18 in [27].

In other words, the kernel matrix  $K$  is positive definite (provided no two  $\mathbf{x}_i$  are the same).

Thus, we have the following theorem:

**Theorem 2** *If all the sample vectors are different and the Gaussian RBF kernel is used, all  $c - 1$  projective functions in SRKDA are eigenvectors of eigen-problem in Eqn. (6) with respect to eigenvalue 1 when  $\delta = 0$ . In other words, the SRKDA and ordinary KDA are equivalent.*

*Proof* This theorem can be easily proofed by combining Lemma 1 and Theorem 1.

It is easy to check that the values of the  $i$ -th and  $j$ -th entries of any vector  $\mathbf{y}$  in the space spanned by  $\{\mathbf{y}_k\}$  in Eqn. (13) are the same as long as  $\mathbf{x}_i$  and  $\mathbf{x}_j$  belong to the same class. Thus the  $i$ -th and  $j$ -th rows of  $\bar{Y}$  are the same, where  $\bar{Y} = [\bar{\mathbf{y}}_1, \dots, \bar{\mathbf{y}}_{c-1}]$ . Theorem (1) shows that when the kernel matrix is positive definite, the  $c - 1$  projective functions of KDA are exactly the solutions of the  $c - 1$  linear equations systems  $K \mathbf{\alpha}_k = \bar{\mathbf{y}}_k$ . Let  $\Theta = [\mathbf{\alpha}_1, \dots, \mathbf{\alpha}_{c-1}]$  be the KDA transformation matrix which embeds the data points into the KDA subspace

$$\Theta^T [K(:, \mathbf{x}_1), \dots, K(:, \mathbf{x}_m)] = \bar{Y}^T.$$

The columns of matrix  $\bar{Y}^T$  are the embedding results of data samples in the KDA subspace. Thus, the data points with the same label are corresponding to the same point in the KDA subspace when the kernel matrix  $K$  is positive definite.

These projective functions are optimal in the sense of separating training samples with different labels. However, they usually overfit the training set thus may not be able to perform well for the test samples.

### 3.2 Computational Analysis

The computation of SRKDA involves two steps: responses ( $\bar{y}_k$  in Eqn. 14) generation and regularized regression. The cost of the first step is mainly the cost of Gram-Schmidt method, which requires  $(mc^2 - \frac{1}{3}c^3)$  flam [31].

To solve the  $c - 1$  linear equations systems in Eqn. (11), we can use the Cholesky decomposition, which uniquely factorizes the positive definite matrix  $K + \delta I$  in the form  $K + \delta I = R^T R$ , where  $R$  is upper triangular with positive diagonal elements. The Cholesky decomposition requires  $\frac{1}{6}m^3$  flam [31]. With this Cholesky decomposition, the  $c - 1$  linear equations can be solved within  $m^2c$  flam [31]. Besides solving the SRKDA optimization problem, we also need to compute the kernel matrix  $K$  which requires  $O(m^2n)$  flam, where  $n$  is the number of features. Thus, the computational cost of SRKDA is

$$\frac{1}{6}m^3 + m^2c + O(m^2n) + mc^2 - \frac{1}{3}c^3,$$

which can be simplified as

$$\frac{1}{6}m^3 + m^2c + O(m^2n).$$

Comparing to the computational cost of ordinary KDA in Eqn. (9), SRKDA reduces the dominant part, which is  $\frac{9}{2}m^3$  of ordinary KDA, to  $\frac{1}{6}m^3$ ; achieves a 27-times speedup.

### 4 Incremental KDA via Spectral Regression

Due to the difficulty of designing an incremental solution for the eigen-decomposition on the kernel matrix in KDA, there has been little work on designing incremental KDA algorithms that can efficiently incorporate new data examples as they become available. The SRKDA algorithm uses regression instead of eigen-decomposition to solve the optimization problem, which provides us the chance to develop incremental version of SRKDA.

The major cost in SRKDA computation is the step of Cholesky decomposition which requires  $\frac{1}{6}m^3$  flam. Fortunately, the Cholesky decomposition can be easily implemented in the incremental manner [31]. Actually, *Sherman's march*, one of the most popular Cholesky decomposition algorithms, is implemented in the incremental manner [31].

The procedure of Sherman's march is illustrated graphically in Figure 1. The gray area represents the part of the Cholesky decomposition that has already

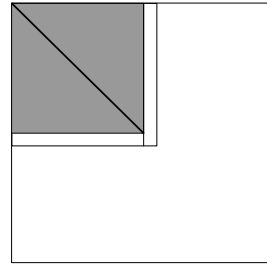


Fig. 1 Sherman's march (Cholesky decomposition)

been computed with  $R$  and  $R^T$  separated by a diagonal line<sup>1</sup>. The white area represents untouched elements of the original matrix. The thin vertical box represents the column of  $R$  about to be computed. The algorithm is easy to derive. We show how to proceed from  $(m - 1) \times (m - 1)$  submatrix to a  $m \times m$  matrix. We have

$$\begin{aligned} K_m &= \begin{pmatrix} K_{m-1} & \mathbf{k}_{1m} \\ \mathbf{k}_{1m}^T & k_{mm} \end{pmatrix} \\ &= \begin{pmatrix} R_{m-1}^T & \mathbf{0} \\ \mathbf{r}_{1m}^T & r_{mm} \end{pmatrix} \begin{pmatrix} R_{m-1} & \mathbf{r}_{1m} \\ \mathbf{0} & r_{mm} \end{pmatrix}, \end{aligned}$$

which leads to

$$\begin{aligned} K_{m-1} &= R_{m-1}^T R_{m-1} \\ \mathbf{k}_{1m} &= R_{m-1}^T \mathbf{r}_{1m} \\ k_{mm} &= \mathbf{r}_{1m}^T \mathbf{r}_{1m} + r_{mm}^2 \end{aligned}$$

When the Cholesky decomposition of the  $(m - 1) \times (m - 1)$  submatrix  $K_{m-1}$  is known, it is easy to get the Cholesky decomposition of the  $m \times m$   $K_m$ . For detailed derivation, please see [31].

Now, let us consider the additional computational cost of incremental SRKDA when  $\Delta m$  new data samples are injected to the system which already has  $m$  samples. Comparing to the batch mode of SRKDA, we can get computational saving on two steps:

1. We only need to calculate the additional part of kernel matrix which requires  $O(nm\Delta m + n\Delta m^2)$  flam;
2. The incremental Cholesky decomposition requires  $\frac{1}{6}(m + \Delta m)^3 - \frac{1}{6}m^3$  flam [31].

Thus, the computation cost of incremental SRKDA measured by flam is

$$\begin{aligned} &\frac{1}{2}m^2\Delta m + \frac{1}{2}m\Delta m^2 + \frac{1}{6}\Delta m^3 + (m + \Delta m)^2c \\ &+ O(nm\Delta m + n\Delta m^2) + (m + \Delta m)c^2 - \frac{1}{3}c^3. \end{aligned}$$

When  $\Delta m \ll m$  and  $c \ll m$ , the above cost can be simplified as

$$\left(\frac{\Delta m}{2} + c\right)m^2 + O(nm\Delta m).$$

**Table 2** Computational complexity of KDA and SRKDA

Algorithm		operation counts ( <i>flam</i> [31])
Batch mode	KDA	$\frac{9}{2}m^3 + cm^2 + O(nm^2)$
	SRKDA	$\frac{1}{6}m^3 + cm^2 + O(nm^2)$
Incremental mode	KDA	$\frac{9}{2}m^3 + cm^2 + O(nm\Delta m)$
	SRKDA	$(\frac{\Delta m}{2} + c)m^2 + O(nm\Delta m)$

$m$ : the number of data samples

$n$ : the number of features

$c$ : the number of classes

$\Delta m$ : the number of new data samples

We summarize our complexity analysis results in Table 2. The main conclusions include:

- The ordinary KDA needs to perform eigen-decomposition on the kernel matrix, which is very computationally expensive. Moreover, it is difficult to develop incremental algorithm based on the ordinary KDA formulation. In both batch and incremental modes, ordinary KDA has the dominant part of the cost as  $\frac{9}{2}m^3$ .
- SRKDA performs regression instead of eigen-decomposition. In the batch mode, it only has the dominant part of the cost as  $\frac{1}{6}m^3$ , which is a 27-times speedup of ordinary KDA. Moreover, it is easy to develop incremental version of SRKDA which only has quadratic-time complexity with respect to  $m$ . This computational advantage makes SRKDA much more practical in real world applications.

## 5 Sparse KDA via Spectral Regression

Since SRKDA uses regression as a building block, various kinds of regularization techniques can be easily incorporated which makes SRKDA more flexible. In this section, we will discuss the usage of  $L_1$ -norm regularizer to produce a sparse KDA solution.

Recently, there are considerable interests on developing sparse subspace learning algorithms, *i.e.*, the projective vectors are sparse. While the traditional linear subspace learning algorithms (*e.g.*, PCA, LDA) learn a set of combined features which are linear combinations of all the original features, sparse linear subspace learning algorithms can learn the combined features which are linear combinations of part of the original features (*important* ones). Such parsimony not only produces a set of projective functions that are easy to interpret but also leads to better performance [37,25].

Zou *et al.* [37] proposed an elegant sparse PCA algorithm (SPCA) using their “Elastic Net” framework for  $L_1$ -penalized regression on regular principle components, solved very efficiently using *least angle regression* (LARS) [13]. Subsequently, d’Aspremont *et al.* [11] relaxed the hard cardinality constraint and solved for a convex approximation using semi-definite programming. In [24,25], Moghaddam *et al.* proposed a spectral bounds framework for sparse subspace learning. Particularly, they proposed both exact and greedy algorithms for sparse PCA and sparse LDA.

The projective function of a kernel subspace learning algorithm can be written as

$$f(\mathbf{x}) = \boldsymbol{\alpha}^T K(:, \mathbf{x}) = \sum_{i=1}^m \alpha_i \mathcal{K}(\mathbf{x}_i, \mathbf{x}). \quad (15)$$

In the ordinary kernel subspace learning algorithms,  $\alpha_i$  are usually nonzero and the projective function is dependent on all the samples in the training set. When we aim at learning a sparse function (sparse  $\boldsymbol{\alpha}$ ), many  $\alpha_i$  will equal to zero. Thus, the projective function will only depend on part of the training samples. From this sense, the sparse kernel subspace learning algorithms share the similar idea of Support Vector Machines [36]. Those samples with non-zero  $\alpha_i$  can also be called as support vectors. One advantage of this parsimony is that it requires less storage for the model and less computational time in the testing phase.

Following [25], the objective function of Sparse Kernel Discriminant Analysis (SparseKDA) can be defined as the following cardinality-constrained optimization:

$$\begin{aligned} \max \quad & \frac{\boldsymbol{\alpha}^T K W K \boldsymbol{\alpha}}{\boldsymbol{\alpha}^T K K \boldsymbol{\alpha}} \\ \text{subject to} \quad & \text{card}(\boldsymbol{\alpha}) = k \end{aligned} \quad (16)$$

The feasible set is all sparse  $\boldsymbol{\alpha} \in \mathbb{R}^m$  with  $k$  non-zero elements and  $\text{card}(\boldsymbol{\alpha})$  as their  $L_0$ -norm. Unfortunately, this optimization problem is NP-hard and generally intractable.

In [24,25], Moghaddam *et al.* proposed a spectral bounds framework for sparse subspace learning. Particularly, they proposed both exact and greedy algorithms for sparse PCA and sparse LDA. Their spectral bounds framework is based on the following optimal condition of the sparse solution.

For simplicity, we define  $A = K W K$  and  $B = K K$ . A sparse vector  $\boldsymbol{\alpha} \in \mathbb{R}^m$  with cardinality  $k$  yielding the maximum objective value in Eqn. (16) would necessarily imply that

$$\lambda_{max} = \frac{\boldsymbol{\alpha}^T A \boldsymbol{\alpha}}{\boldsymbol{\alpha}^T B \boldsymbol{\alpha}} = \frac{\boldsymbol{\beta}^T A_k \boldsymbol{\beta}}{\boldsymbol{\beta}^T B_k \boldsymbol{\beta}}$$

where  $\boldsymbol{\beta} \in \mathbb{R}^k$  contains the  $k$  non-zero elements in  $\boldsymbol{\alpha}$  and the  $k \times k$  principle sub-matrices of  $A$  and  $B$  obtained

<sup>1</sup> Actually, we only need to store  $R$ .



by deleting the rows and columns corresponding to the zero indices of  $\alpha$ . The  $k$ -dimensional quadratic form in  $\beta$  is equivalent to a standard unconstrained generalized Rayleigh quotient, which can be solved by a generalized eigen-problem.

The above observation gives the exact algorithm for sparse subspace learning: a discrete search for the  $k$  indices which maximize  $\lambda_{max}$  of the subproblem  $(A_k, B_k)$ . However, such observation does not suggest an efficient algorithm because an exhaustive search is still NP-hard. To solve this problem, Moghaddam *et al.* proposed an efficient greedy algorithm which combines *backward elimination* and *forward selection* [24,25]. However, there are two major drawbacks of their approach:

1. Even their algorithm is a greedy one, the cost of backward elimination is with complexity  $O(m^4)$ [25].
2. In reality, more than one projective functions are usually necessary for subspace learning. However, the optimal condition of the sparse solution only gives the guide to find ONE sparse “eigenvector”, which is the first projective function. It is unclear how to find the following projective functions. Although [24] suggests to use recursive deflation, the sparseness of the following projective functions is not guaranteed.

Our SRKDA algorithm uses the regression instead of eigen-decomposition to solve the optimization problem. Thus, it provides us the chance to develop more efficient sparse KDA algorithm.

Recall the second step of SRKDA, which is solving the linear equations system  $K\alpha = \mathbf{y}$ . Essentially, we try to solve a regression problem:

$$\min_{\alpha} \sum_{i=1}^m \left( K(:, \mathbf{x}_i)^T \alpha - y_i \right)^2$$

where  $K(:, \mathbf{x}_i)$  is the  $i$ -th column of  $K$  and  $y_i$  is the  $i$ -th element of  $\mathbf{y}$ . We can add different regularizers to get different solutions with desired properties. The SRKDA algorithm we described in the previous section essentially adds a  $L_2$ -norm regularizer. We can also use a  $L_1$ -norm regularizer:

$$\min_{\alpha} \left( \sum_{i=1}^m \left( K(:, \mathbf{x}_i)^T \alpha - y_i \right)^2 + \delta \sum_{i=1}^m |\alpha_i| \right) \quad (17)$$

which is usually referred as *lasso regression* [18]. Due to the nature of the  $L_1$  penalty, some coefficients  $\alpha_i$  will be shrunk to exact zero if  $\delta$  is large enough. Therefore, the lasso produces a sparse projective function, which is exactly what we want.

By using the *Least Angel Regression* (LARS) algorithm [13], the entire solution path (the solutions with

**Table 3** Statistics of the three data sets

dataset	dim ( $n$ )	train size ( $m$ )	test size	# of classes ( $c$ )
Isolet	617	6238	1559	26
USPS	256	7291	2007	10
PIE	1024	8000	3554	68

all the possible cardinality on  $\alpha$ ) of the regression problem in Eqn. (17) can be computed in  $O(m^3)$ . Thus, SRKDA with a  $L_1$ -norm regularizer provides us an efficient algorithm to compute the sparse KDA solution.

## 6 Experimental Results

In this section, we investigate the performance of our proposed SRKDA algorithm in batch mode, incremental mode and sparse mode. All of our experiments have been performed on a P4 3.20GHz Windows XP machine with 2GB memory. For the purpose of reproducibility, we provide all the algorithms used in these experiments at:

<http://www.zjucadcg.cn/dengcai/Data/data.html>

### 6.1 Datasets

Three datasets are used in our experimental study, including spoken letter, handwritten digit image, and face image data sets. The important statistics of three datasets are summarized below (see also Table 3):

- The Isolet spoken letter recognition database<sup>2</sup> was first used in [14]. It contains 150 subjects who spoke the name of each letter of the alphabet twice. The speakers are grouped into sets of 30 speakers each, and are referred to as isolet1 through isolet5. In the past usage [14][12], isolet1&2&3&4 were used as the training set and isolet5 was used as the test set. For the purposes of our experiment, we also choose isolet5 as the test set and perform several runs with isolet1, isolet1&2, isolet1&2&3, and isolet1&2&3&4 as the training set respectively.
- The USPS handwritten digit database is described in [19]. A popular subset<sup>3</sup> contains 9298  $16 \times 16$  handwritten digit images in total, which is then split into 7291 training images and 2007 test images. In our experiment, we train all the algorithms on the first 1500 (3000, 4500, 6000, and 7291) images in the training set and test on the 2007 test images.

<sup>2</sup> <http://www.ics.uci.edu/~mllearn/MLSummary.html>

<sup>3</sup> <http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/multiclass.html#usps>

- The CMU PIE face database<sup>4</sup> contains 68 subjects with 41,368 face images as a whole. The face images were captured under varying pose, illumination and expression. In our experiment, the five near frontal poses (C05, C07, C09, C27, C29) under different illuminations and expressions are used which leaves us 11,554 face images. All the images are manually aligned and cropped. The cropped images are  $32 \times 32$  pixels, with 256 gray levels per pixel<sup>5</sup>. Among the 11,554 images, 8,000 images are used as the training set and the remaining 3,554 images are used for testing. We also run several cases by training all the algorithms on the first 2000, 3000,  $\dots$ , 8000 images in the training set.

## 6.2 Compared algorithms

Four algorithms which are compared in our experiments are listed below:

1. Linear Discriminant Analysis (LDA) [15], which provides us a baseline performance of linear algorithms. We can examine the usefulness of kernel approaches by comparing the performance of KDA and LDA.
2. Kernel Discriminant Analysis (KDA) as discussed in Section 2. We test the regularized version and choose the regularization parameter  $\delta$  by five fold cross-validation on the training set.
3. Spectral Regression Kernel Discriminant Analysis (SRKDA), our approach proposed in this paper. The regularization parameter  $\delta$  is also chosen by five fold cross-validation on the training set.
4. Support Vector Machine (SVM) [36], which is believed as one of the state-of-the-art classification algorithms. Specifically, we use the LibSVM system [10] which implements the multi-class classification with one versus one strategy. SVM is used to get the sense that how good the performance of KDA is.

We use the Gaussian RBF kernel for all the kernel-based methods. We tune the kernel width parameter  $\sigma$  and large margin parameter  $C$  in SVM to achieve best testing performance for SVM. Then, the same kernel width parameter  $\sigma$  is used in all the other kernel-based algorithms.

## 6.3 Results

The classification error rate as well as the training time (second) for each method on the three data sets are reported on the Table (4 ~ 6) respectively.

The main observations from the performance comparisons include:

- The Kernel Discriminant Analysis model is very effective in classification. SRKDA has the best performance for almost all the cases in all the three data sets (even better than SVM). For Isolet data set, previous study [12] reported the minimum error rate training on Isolet1+2+3+4 by OPT<sup>6</sup> with 30 bit ECOC is 3.27%. KDA (SRKDA) achieved better performance in our experiment for this train/test split. For USPS data set, previous studies [27] reported error rate 3.7% for KDA and 4.0% for SVM, slightly better than the results in our experiment. For all the cases, KDA (SRKDA) achieved significantly better performance than LDA, which suggests the effectiveness of kernel approaches.
- Since the eigen-decomposition of the kernel matrix is involved, the ordinary KDA is computationally expensive in training. SRKDA uses regression instead of eigen-decomposition to solve the optimization problem, and thus achieves significant speedup comparing to ordinary KDA. The empirical results are consistent with the theoretical estimation of the efficiency. The time of training SRKDA is comparable with that of training SVM. SRKDA is faster than SVM on Isolet and PIE data sets, while slower than SVM on USPS data set. This is because the time of training SVM is dependant with the number of support vectors [2]. For some data sets with lots of noise (*e.g.*, USPS), the number of support vectors is far less than the number of samples. In this case, SVM can be trained very fast.

## 6.4 Experiments on Incremental KDA

In this experiment, we study the computational cost of SRKDA performing in the incremental manner. The USPS and PIE data sets are used. We start from the training set with the size of 1000 (the first 1000 samples in the whole training set) and increase the training size by 200 for each step. SRKDA is then performed in the incremental manner. It is important to note that SRKDA in the incremental manner gives the exactly same projective functions as the SRKDA in the batch mode. Thus, we only care about the computational costs in this experiment.

Figure 2 and 3 shows log-log plots of how CPU-time of KDA (SRKDA, incremental SRKDA) increases with the size of the training set on USPS and PIE data set respectively. Lines in a log-log plot correspond to polynomial growth  $O(m^d)$ , where  $d$  corresponds to the slope

<sup>4</sup> [http://www.ri.cmu.edu/projects/project\\_418.html](http://www.ri.cmu.edu/projects/project_418.html)

<sup>5</sup> <http://www.zjucadcg.cn/dengcai/Data/FaceData.html>

<sup>6</sup> Conjugate-gradient implementation of back-propagation

**Table 4** Performance comparisons on Isolet dataset

Training Set	Error (%)				Time (s)				Speedup
	LDA	KDA	SRKDA	SVM	LDA	KDA	SRKDA	SVM	
Isolet1	15.27	11.74	12.89	12.51	1.93	18.86	1.21	4.75	15.6
Isolet1+2	6.61	3.79	3.85	4.11	2.14	134.6	5.51	13.79	24.4
Isolet1+2+3	5.90	2.99	3.08	3.34	2.37	451.6	14.09	23.84	32.1
Isolet1+2+3+4	5.71	2.82	2.89	3.27	2.56	991.2	27.86	34.82	35.6

\*Column labeled “Speedup” shows how many times faster the SRKDA is (comparing to ordinary KDA).

**Table 5** Performance comparisons on USPS dataset

Training Set	Error (%)				Time (s)				Speedup
	LDA	KDA	SRKDA	SVM	LDA	KDA	SRKDA	SVM	
1500	10.61	6.58	5.88	6.85	0.21	14.97	0.92	0.78	16.3
3000	9.77	5.53	5.38	5.58	0.27	111.9	4.35	2.20	25.7
4500	9.52	5.53	4.88	5.13	0.34	354.3	11.29	4.06	31.4
6000	9.92	5.03	4.43	5.08	0.40	825.3	22.74	6.22	36.3
7291	10.26	4.83	4.04	4.83	0.47	1553.6	37.59	8.18	41.3

**Table 6** Performance comparisons on PIE dataset

Training Set	Error (%)				Time (s)				Speedup
	LDA	KDA	SRKDA	SVM	LDA	KDA	SRKDA	SVM	
2000	5.29	5.18	4.81	6.30	8.77	36.51	2.47	24.13	14.8
3000	4.61	4.25	3.94	4.70	9.06	116.9	5.39	43.99	21.7
4000	4.14	5.53	3.24	3.74	9.42	256.6	10.35	68.43	24.8
5000	3.85	3.23	2.90	3.29	9.73	502.3	17.40	96.26	28.9
6000	3.57	2.91	2.53	2.84	10.06	830.7	27.21	125.6	30.5
7000	3.40	2.65	2.19	2.64	10.39	1340.9	38.65	155.6	34.7
8000	3.35	2.41	2.17	2.34	10.79	1908.1	53.75	186.7	35.5

of the line. The ordinary KDA scales roughly  $O(m^{2.9})$ , which is slightly better than the theoretical estimation. SRKDA in the batch mode has better scaling, which is also better than theoretical estimation with roughly  $O(m^{2.6})$  over much of the range. This explains why SRKDA can be more than 27 times faster than ordinary KDA in the previous experiments. The SRKDA in the incremental mode has the best scaling, which is (to some surprise) better than quadratic with roughly  $O(m^{1.8})$  over much of the range.

## 6.5 Experiments on Sparse KDA

In this experiment, we study the performance of SRKDA performing in the sparse mode, *i.e.*, the SRKDA with  $L_1$ -norm regularizer to produce the sparse KDA solution. To the best of our knowledge, there is no other published method to generate a sparse KDA solution. Moghaddam’s sparse LDA approach [25] can be modified to generate the sparse KDA solution. However, as we pointed out in the last section, their approach can only generate ONE sparse projective function and is only suitable for binary class problem, while all the three data sets studied in this paper are multi-class data sets.

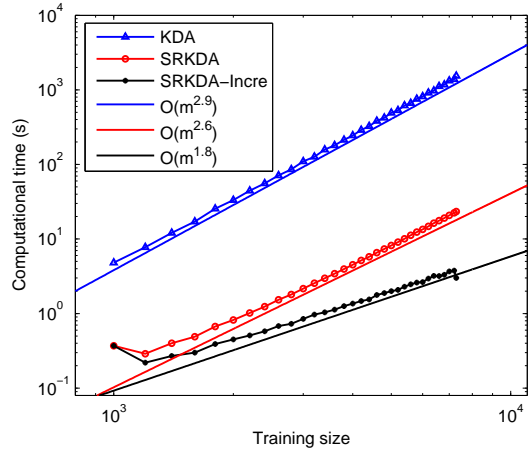
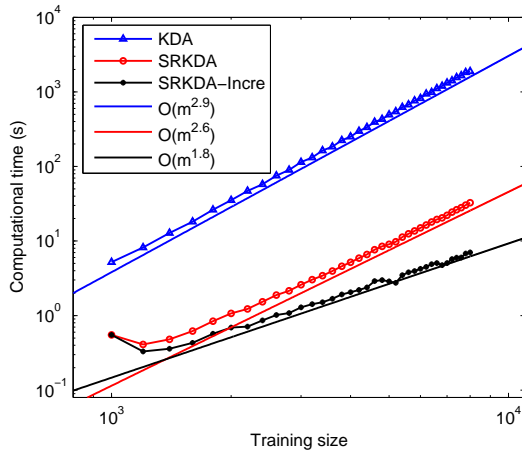
**Fig. 2** Computational cost of KDA, batch SRKDA and incremental SRKDA on the USPS data set.

Table (7), (8) and (9) show the classification error rate of SRKDA in sparse mode on the three data sets respectively. By using the *Least Angel Regression* (LARS) algorithm [13], the entire solution path (the solutions with all the possible cardinality on the projective function  $\alpha$ ) can be computed. After this, we use cross validation to select the optimal cardinality of the projective function in the experiment. We also show the sparsity of the projective function of SRKDA(sparse) in the ta-



**Fig. 3** Computational cost of KDA, batch SRKDA and incremental SRKDA on the PIE data set.

**Table 7** Classification error on Isolet dataset

Training Set	Error (%)			Sparsity
	KDA	SRKDA	SRKDA(Sparse)	
Isolet1	11.74	12.89	11.74	60%
Isolet1+2	3.79	3.85	3.59	60%
Isolet1+2+3	2.99	3.08	2.82	60%
Isolet1+2+3+4	2.82	2.89	2.82	60%

**Table 8** Classification error on USPS dataset

Training Set	Error (%)			Sparsity
	KDA	SRKDA	SRKDA(Sparse)	
1500	6.58	5.88	5.83	60%
3000	5.53	5.38	5.13	60%
4500	5.53	4.88	4.73	60%
6000	5.03	4.43	4.04	60%
7291	4.83	4.04	3.94	60%

**Table 9** Classification error on PIE dataset

Training Set	Error (%)			Sparsity
	KDA	SRKDA	SRKDA(Sparse)	
2000	5.18	4.81	4.73	60%
3000	4.25	3.94	3.71	60%
4000	5.53	3.24	3.12	60%
5000	3.23	2.90	2.81	60%
6000	2.91	2.53	2.44	60%
7000	2.65	2.19	2.17	60%
8000	2.41	2.17	2.14	60%

bles. The sparsity is defined as the percentage of zero entries in a projective vector. For ordinary KDA and SRKDA, the projective functions (vectors) are dense and the sparsity is zero.

As can be seen, the SRKDA(sparse) generates much more parsimonious model. The sparsity of the projective function in SRKDA(sparse) is 60%, which means the number of the “support vectors” is less than half of the total training samples. Moreover, such parsimony leads to better performance. In all the cases, the per-

formance of SRKDA(sparse) is better than that of the ordinary KDA and SRKDA.

## 7 Conclusions

In this paper, we propose a novel algorithm for kernel discriminant analysis, called *Spectral Regression Kernel Discriminant Analysis* (SRKDA). Our algorithm is developed from a graph embedding viewpoint of KDA problem. It combines the spectral graph analysis and regression to provide an efficient approach for kernel discriminant analysis. Specifically, SRKDA only needs to solve a set of regularized regression problems and there is no eigenvector computation involved, which is a huge save of computational cost. The theoretical analysis shows that SRKDA can achieve 27-times speedup over the ordinary KDA. Moreover, the new formulation makes it very easy to develop incremental version of the algorithm which can fully utilize the computational results of the existing training samples. With incremental implementation, the computational cost of SRKDA reduces to quadratic-time complexity. Since SRKDA uses regression as a building block, various kinds of regularization techniques can be easily incorporated (*e.g.*,  $L_1$ -norm regularizer to produce sparse projections). Our approach provides a huge possibility to develop new variations of kernel discriminant analysis. Extensive experimental results show that our method consistently outperforms the other state-of-the-art KDA extensions considering both effectiveness and efficiency.

**Acknowledgements** This work was supported in part by National Natural Science Foundation of China under Grants 60905001 and 90920303, National Key Basic Research Foundation of China under Grant 2009CB320801. Any opinions, findings, and conclusions expressed here are those of the authors and do not necessarily reflect the views of the funding agencies.

## References

1. G. Baudat and F. Anouar. Generalized discriminant analysis using a kernel approach. *Neural Computation*, (12):2385–2404, 2000.
2. C. J. C. Burges. A tutorial on support vector machines for pattern recognition. *Data Mining and Knowledge Discovery*, 2(2):121–167, 1998.
3. D. Cai, X. He, and J. Han. Efficient kernel discriminant analysis via spectral regression. In *Proc. Int. Conf. on Data Mining (ICDM’07)*, 2007.
4. D. Cai, X. He, and J. Han. Spectral regression: A unified approach for sparse subspace learning. In *Proc. Int. Conf. on Data Mining (ICDM’07)*, 2007.
5. D. Cai, X. He, and J. Han. Spectral regression: A unified subspace learning framework for content-based image retrieval. In *Proceedings of the 15th ACM International Conference on Multimedia*, Augsburg, Germany, 2007.

6. D. Cai, X. He, and J. Han. SRDA: An efficient algorithm for large scale discriminant analysis. *IEEE Transactions on Knowledge and Data Engineering*, 20(1):1–12, January 2008.
7. D. Cai, X. He, W. V. Zhang, and J. Han. Regularized locality preserving indexing via spectral regression. In *Proceedings of the 16th ACM conference on Conference on information and knowledge management (CIKM'07)*, pages 741–750, 2007.
8. K. Chakrabarti and S. Mehrotra. Local dimensionality reduction: A new approach to indexing high dimensional spaces. In *Proc. 2000 Int. Conf. Very Large Data Bases (VLDB'00)*, 2000.
9. S. Chakrabarti, S. Roy, and M. V. Soundalgekar. Fast and accurate text classification via multiple linear discriminant projections. *The VLDB Journal*, 12(2):170–185, 2003.
10. C.-C. Chang and C.-J. Lin. *LIBSVM: a library for support vector machines*, 2001. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
11. A. d'Aspremont, L. E. Chaoui, M. I. Jordan, and G. R. G. Lanckriet. A direct formulation for sparse PCA using semidefinite programming. In *Advances in Neural Information Processing Systems 17*, 2004.
12. T. G. Dietterich and G. Bakiri. Solving multiclass learning problems via error-correcting output codes. *Journal of Artificial Intelligence Research*, (2):263–286, 1995.
13. B. Efron, T. Hastie, I. Johnstone, and R. Tibshirani. Least angle regression. *Annals of Statistics*, 32(2):407–499, 2004.
14. M. A. Fandy and R. Cole. Spoken letter recognition. In *Advances in Neural Information Processing Systems 3*, 1990.
15. K. Fukunaga. *Introduction to Statistical Pattern Recognition*. Academic Press, 2nd edition, 1990.
16. V. Gaede and O. Günther. Multidimensional access methods. *ACM Comput. Surv.*, 30(2):170–231, 1998.
17. G. H. Golub and C. F. V. Loan. *Matrix computations*. Johns Hopkins University Press, 3rd edition, 1996.
18. T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. New York: Springer-Verlag, 2001.
19. J. J. Hull. A database for handwritten text recognition research. *IEEE Trans. Pattern Anal. Mach. Intell.*, 16(5), 1994.
20. H. Jin, B. C. Ooi, H. T. Shen, C. Yu, and A. Zhou. An adaptive and efficient dimensionality reduction algorithm for high-dimensional indexing. In *Proc. 2003 Int. Conf. on Data Engineering (ICDE'03)*, 2003.
21. C. A. Micchelli. Algebraic aspects of interpolation. In *Proceedings of Symposia in Applied Mathematics*, volume 36, pages 81–102, 1986.
22. S. Mika, G. Rätsch, J. Weston, B. Schölkopf, and K.-R. Müller. Fisher discriminant analysis with kernels. In *Proc. of IEEE Neural Networks for Signal Processing Workshop (NNSP)*, 1999.
23. S. Mika, A. Smola, and B. Schölkopf. An improved training algorithm for kernel fisher discriminants. In *Proceedings AISTATS 2001*. Morgan Kaufmann, 2001.
24. B. Moghaddam, Y. Weiss, and S. Avidan. Spectral bounds for sparse PCA: Exact and greedy algorithms. In *Advances in Neural Information Processing Systems 18*, 2005.
25. B. Moghaddam, Y. Weiss, and S. Avidan. Generalized spectral bounds for sparse LDA. In *ICML '06: Proceedings of the 23rd international conference on Machine learning*, pages 641–648, 2006.
26. C. H. Park and H. Park. Nonlinear discriminant analysis using kernel functions and the generalized singular value decomposition. *SIAM J. Matrix Anal. Appl.*, 27(1):87–102, 2005.
27. B. Schölkopf and A. J. Smola. *Learning with Kernels*. MIT Press, 2002.
28. H. T. Shen, B. C. Ooi, and X. Zhou. Towards effective indexing for very large video sequence database. In *Proceedings of the 2005 ACM SIGMOD international conference on Management of data*, pages 730–741, 2005.
29. H. T. Shen, X. Zhou, and B. Cui. Indexing text and visual features for www images. In *7th Asia Pacific Web Conference (APWeb2005)*, 2005.
30. H. T. Shen, X. Zhou, and A. Zhou. An adaptive and dynamic dimensionality reduction method for high-dimensional indexing. *The VLDB Journal*, 16(2):219–234, 2007.
31. G. W. Stewart. *Matrix Algorithms Volume I: Basic Decompositions*. SIAM, 1998.
32. G. W. Stewart. *Matrix Algorithms Volume II: Eigensystems*. SIAM, 2001.
33. D. Tao, X. Li, X. Wu, and S. J. Maybank. General tensor discriminant analysis and gabor features for gait recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 29(10):1700–1715, 2007.
34. D. Tao, X. Li, X. Wu, and S. J. Maybank. Geometric mean for subspace selection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31(2):260–274, 2009.
35. D. Tao, X. Tang, X. Li, and Y. Rui. Kernel direct biased discriminant analysis: A new content-based image retrieval relevance feedback algorithm. *IEEE Transactions on Multimedia*, 8(4):716–727, 2006.
36. V. N. Vapnik. *Statistical learning theory*. John Wiley & Sons, 1998.
37. H. Zou, T. Hastie, and R. Tibshirani. Sparse principle component analysis. Technical report, Statistics Department, Stanford University, 2004.