

A Deterministic Approach to Regularized Linear Discriminant Analysis

Alok Sharma^{1,2}, Kuldip K. Paliwal¹

¹School of Engineering, Griffith University, Australia

²School of Engineering and Physics, University of the South Pacific, Fiji

Abstract

The regularized linear discriminant analysis (RLDA) technique is one of the popular methods for dimensionality reduction used for small sample size problems. In this technique, regularization parameter is conventionally computed using a cross-validation procedure. In this paper, we propose a deterministic way of computing the regularization parameter in RLDA for small sample size problem. The computational cost of the proposed deterministic RLDA is significantly less than the cross-validation based RLDA technique. The deterministic RLDA technique is also compared with other popular techniques on a number of datasets and favorable results are obtained.

1. Introduction

Linear discriminant analysis (LDA) is a popular technique for dimensionality reduction and feature extraction. Dimensionality reduction is a pre-requisite for many statistical pattern recognition techniques. It is primarily applied for improving generalization capability and reducing computational complexity of a classifier. In LDA the

1 dimensionality is reduced from d -dimensional space to h -dimensional space (where
2 $h < d$) by using a transformation $\mathbf{W} \in \mathbb{R}^{d \times h}$. The transformation (or orientation) matrix
3 \mathbf{W} is found by maximizing the Fisher's criterion: $J(\mathbf{W}) = |\mathbf{W}^T \mathbf{S}_B \mathbf{W}| / |\mathbf{W}^T \mathbf{S}_W \mathbf{W}|$, where
4 $\mathbf{S}_W \in \mathbb{R}^{d \times d}$ is within-class scatter matrix and $\mathbf{S}_B \in \mathbb{R}^{d \times d}$ is between-class scatter matrix.
5 Under this criterion, the transformation of feature vectors from higher dimensional
6 space to lower dimensional space is done in such a manner that the between-class
7 scatter in the lower dimensional space is maximized and within-class scatter is
8 minimized. The orientation matrix \mathbf{W} is computed by the eigenvalue decomposition
9 (EVD) of $\mathbf{S}_W^{-1} \mathbf{S}_B$ [1].

10

11 In many pattern classification applications, the matrix \mathbf{S}_W becomes singular and its
12 inverse computation becomes impossible. This is due to the large dimensionality of
13 feature vectors compared to small number of vectors available for training. This is
14 known as small sample size (SSS) problem [2]. There exist several techniques that can
15 overcome this problem [3]-[11],[19]-[34]. Among these techniques, regularized LDA
16 (RLDA) technique [3] is one of the pioneering methods for solving SSS problem. The
17 RLDA technique has been widely studied in the literature [12]-[14]. It has been applied
18 in areas like face recognition [13],[14] and bioinformatics [15].

1 .

2 In the RLDA technique, the \mathbf{S}_W matrix is regularized to overcome the singularity
3 problem of \mathbf{S}_W . This regularization can be done in various ways. For example, Zhao et al.
4 [12][16][17] have done this by adding a small positive constant α (known as
5 regularization parameter) to the diagonal elements of matrix \mathbf{S}_W ; i.e., the matrix \mathbf{S}_W is
6 approximated by $\mathbf{S}_W + \alpha\mathbf{I}$ and the orientation matrix is computed by EVD of $(\mathbf{S}_W +$
7 $\alpha\mathbf{I})^{-1}\mathbf{S}_B$. The performance of RLDA technique depends on the choice of the
8 regularization parameter α . In the past studies [18], this parameter is chosen rather
9 heuristically, for example, by applying cross-validation procedure on the training data.

10 In the cross-validation based RLDA technique (denoted here as CV-RLDA), the training
11 data is divided into two subsets: training subset and validation subset. The
12 cross-validation procedure searches over a finite range of α values and finds an α
13 value in this range that maximizes the classification accuracy over the validation subset.

14 In the cross-validation procedure, the estimate of α depends on the range over which it
15 is explored. For a given dataset, its classification accuracy can vary depending upon the
16 range of α being explored. Since many values of α have to be searched in this range,
17 the computational cost of this procedure is quite high. In addition, the cross-validation
18 procedure used in the CV-RLDA technique is biased towards the classifier used.

1 In order to address these drawbacks of CV-RLDA technique, we explore a deterministic
2 way for finding the regularization parameter α . This would provide a unique value of
3 the regularization parameter on a given training data. We call this approach as the
4 deterministic RLDA (DRLDA) technique. This technique avoids the use of the heuristic
5 (cross-validation) procedure for parameter estimation and improves the computational
6 efficiency. We show that this deterministic approach computes the regularization
7 parameter by maximizing the Fisher's criterion and its classification performance is
8 quite promising compared to other LDA techniques.

9

10 **2. Related work**

11 In a SSS problem, the within-class scatter matrix \mathbf{S}_W becomes singular and its inverse
12 computation becomes impossible. In order to overcome this problem, generally inverse
13 computation of \mathbf{S}_W is avoided or approximated for the computation of orientation
14 matrix \mathbf{W} . There are several techniques that can overcome this SSS problem. One way
15 to solve this problem is by estimating the inverse of \mathbf{S}_W by its pseudoinverse and then
16 the conventional eigenvalue problem can be solved to compute the orientation matrix \mathbf{W} .
17 This was the basis of Pseudoinverse LDA (PILDA) technique [20]. Some improvements
18 of PILDA have also been presented in [28, 31]. In Fisherface (PCA+LDA) technique,

1 d -dimensional features are firstly reduced to h -dimensional feature space by the
2 application of PCA [2][52][53] and then LDA is applied to further reduce features to k
3 dimensions. There are several ways for determining the value of h in PCA+LDA
4 technique [4][5]. In the Direct LDA (DLDA) technique [7], the dimensionality is reduced
5 in two stages. In the first stage, a transformation matrix is computed to transform the
6 training samples to the range space of \mathbf{S}_B , and in the second stage, the dimensionality
7 of this transformed samples is further transformed by some regulating matrices. The
8 Improved DLDA technique [11], addresses drawbacks of DLDA technique. In the
9 improved DLDA technique, first \mathbf{S}_W is decomposed into its eigenvalues and
10 eigenvectors instead of \mathbf{S}_B matrix as of DLDA technique. Here, both its null space and
11 range space information are utilized by approximating \mathbf{S}_W by a well deterministic
12 substitution. Then \mathbf{S}_B is diagonalized using regulating matrices. For the Null LDA
13 (NLDA) technique [6], the orientation \mathbf{W} is computed in two stages. In the first stage,
14 the data is projected on the null space of \mathbf{S}_W and in the second stage it finds \mathbf{W} that
15 maximizes $|\mathbf{W}^T \mathbf{S}_B \mathbf{W}|$. In orthogonal LDA (OLDA) technique [8], the orientation matrix
16 \mathbf{W} is obtained by simultaneously diagonalizing scatter matrices. It has shown that
17 OLDA is equivalent to NLDA under a mild condition [8]. The Uncorrelated LDA (ULDA)
18 technique [21], is a slight variation of OLDA technique. In ULDA, the orientation

1 matrix \mathbf{W} is made uncorrelated. The fast NLDA (FNLDA) technique [25], is an
2 alternative method of NLDA. In this technique, the orientation matrix is obtained by
3 using the relation $\mathbf{W} = \mathbf{S}_T^+ \mathbf{S}_B \mathbf{Y}$, where \mathbf{Y} is a random matrix of rank $c - 1$, and c is the
4 number of classes. This technique is so far the fastest technique of performing null LDA
5 operation. In extrapolation LDA (ELDA) technique [32], the null space of \mathbf{S}_W matrix is
6 regularized by extrapolating eigenvalues of \mathbf{S}_W using exponential fitting function. This
7 technique utilizes range space information and null space information of \mathbf{S}_W matrix.
8 The two stage LDA (TSLDA) technique [34], exploits all four informative spaces of
9 scatter matrices. These spaces are included in two separate discriminant analyses in
10 parallel. In the first analysis, null space of \mathbf{S}_W and range space of \mathbf{S}_B are retained. In
11 the second analysis, range space of \mathbf{S}_W and null space of \mathbf{S}_B are retained. In
12 eigenfeature regularization (EFR) technique [10], \mathbf{S}_W is regularized by extrapolating
13 its eigenvalues in its null space. The lagging eigenvalues of \mathbf{S}_W is considered as noisy or
14 unreliable which are replaced by an estimation function. The General Tensor
15 Discriminant Analysis (GTDA) technique [48] has been developed for image recognition
16 problems. This work focuses on the representation and pre-processing of
17 appearance-based models for human gait sequences. Two models were presented: gabor
18 gait and tensor gait. In [49], authors proposed a constrained empirical risk

1 minimization framework for distance metric learning (DML) to solve SSS problem. In
 2 Double Shrinking Sparse Dimension Reduction technique [50], the SSS problem is
 3 solved by penalizing the parameter space. A detailed explanation regarding LDA is
 4 given in [51] and an overview regarding SSS based LDA techniques is given in [47].
 5 There are other techniques which can solve SSS problem and applied in various fields of
 6 research [54]-[62]. In this paper, we focus on regularize LDA (RLDA) technique. This
 7 technique overcomes SSS problem by utilizing a small perturbation to the \mathbf{S}_W matrix.
 8 The details of RLDA have been discussed in the next section.

9

10 **3. Regularized linear discriminant techniques for SSS problem**

11 In the RLDA technique, the within-class scatter matrix \mathbf{S}_W is approximated by adding
 12 a regularization parameter to make it a non-singular matrix [3]. There are, however,
 13 different ways to perform regularization (see for details, [3][12]-[14][16][17][30][33]). In
 14 this paper we adopted Zhao's model [12][16][17] to approximate \mathbf{S}_W by adding a
 15 positive constant in the following way $\hat{\mathbf{S}}_W = \mathbf{S}_W + \alpha \mathbf{I}$ ¹. This will make within-class
 16 scatter matrix a non-singular matrix and then its inverse computation would be
 17 possible. The RLDA technique computes the orientation matrix \mathbf{W} by EVD of $\hat{\mathbf{S}}_W^{-1} \mathbf{S}_B$.

¹ In the Friedman's model [3], \mathbf{S}_W is estimated as $\hat{\mathbf{S}}_W = (1 - \alpha) \mathbf{S}_W + \alpha \mathbf{I}$. We have compared Zhao's model and Friedman's model of CV-RLDA and found that Zhao's model exhibits comparatively better generalization capability (see Appendix-I for details). Furthermore, we have considered Zhao's model because it is relatively simpler for establishing deterministic approach of computing α (in DRLDA).

1 Thus, this technique uses null space of \mathbf{S}_W , range space of \mathbf{S}_W and range space of \mathbf{S}_B in
2 one step (i.e., simultaneously).

3

4 In the RLDA technique, a fixed value of regularization parameter can be used, but it
5 may not give the best classification performance as shown in Appendix-II. Therefore,
6 the regularization parameter α is normally computed by the cross-validation procedure.

7 The cross-validation procedure (e.g. leave-one-out or k -fold) employs a particular
8 classifier to estimate α and is conducted on the training set (which is different from the
9 test set). We briefly describe below the leave-one out cross-validation procedure used in

10 the CV-RLDA technique. Let $[a, b]$ be the range of α to be explored and α_0 be any
11 value in this range. Consider a case when n training samples are available. The
12 training set is first subdivided into training subset (consisting of $n - 1$ samples) and

13 validation subset (consisting of 1 sample). For this particular subdivision of training set,
14 the following operations are required: 1) computation of scatter matrices \mathbf{S}_B , \mathbf{S}_W and

15 $\hat{\mathbf{S}}_W = \mathbf{S}_W + \alpha_0 \mathbf{I}$ for $n - 1$ samples in the training subset; 2) EVD of $\hat{\mathbf{S}}_W^{-1} \mathbf{S}_B$ to compute
16 orientation matrix \mathbf{W} ; and 3) classification of the left out sample (from the validation
17 subset) by the classifier to obtain the classification accuracy. These computational

18 operations are carried out for $n - 1$ subdivisions of the training set and the average

1 classification accuracy over the $n - 1$ runs is computed. This average classification
2 accuracy is obtained for a particular value of α (namely α_0). All the above operations
3 will be repeated for other values of α in the range $[a, b]$ to get the highest average
4 classification accuracy. From this description, it is obvious that the cross-validation
5 procedure used in the CV-RLDA technique has the following drawbacks:

6 ● Since the cross-validation procedure repeats the above-mentioned computational
7 operations many times for different values of α , its computation complexity is
8 extremely large.

9 ● Since the cross-validation procedure used in the CV-RLDA technique searches the
10 α parameter over a finite range $[a, b]$, it may not estimate its optimum value. In
11 order to estimate its optimum value, one has to investigate all possible values of α
12 in the range of $(0, \infty)$. However, it is an impossible task (as it will take infinite
13 amount of computation time). Thus, the α value computed by this procedure
14 depends on two factors: 1) the range over which it is searched, and 2) the fineness of
15 the search procedure.

16 ● The cross-validation procedure estimates the regularization parameter in
17 CV-RLDA for a particular classifier. Thus, the estimated value is specific to the
18 classifier and cannot be generalized to other classifiers.

1

2 In our proposed DRLDA technique, we use a deterministic approach to estimate α
3 parameter by maximizing the modified Fisher's criterion. The proposed technique is
4 described in the next section.

5 **4. DRLDA technique**

6 **4.1 Notations**

7 Let $\mathfrak{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ denotes n training samples (or feature vectors) in a
8 d -dimensional space having class labels $\Omega = \{\omega_1, \omega_2, \dots, \omega_n\}$, where $\omega \in \{1, 2, \dots, c\}$ and c
9 is the number of classes. The set \mathfrak{X} can be subdivided into c subsets $\mathfrak{X}_1, \mathfrak{X}_2, \dots, \mathfrak{X}_c$ where
10 \mathfrak{X}_j belongs to class j and consists of n_j number of samples such that:

$$11 \quad n = \sum_{j=1}^c n_j$$

12 and $\mathfrak{X}_j \subset \mathfrak{X}$ and $\mathfrak{X}_1 \cup \mathfrak{X}_2 \cup \dots \cup \mathfrak{X}_c = \mathfrak{X}$.

13

14 If $\boldsymbol{\mu}_j$ is the centroid of \mathfrak{X}_j and $\boldsymbol{\mu}$ is the centroid of \mathfrak{X} , then the total scatter matrix \mathbf{S}_T ,
15 within-class scatter matrix \mathbf{S}_W and between-class scatter matrix \mathbf{S}_B are defined as

16 [1][35][36]

$$17 \quad \mathbf{S}_T = \sum_{\mathbf{x} \in \mathfrak{X}} (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T,$$

$$18 \quad \mathbf{S}_W = \sum_{j=1}^c \sum_{\mathbf{x} \in \mathfrak{X}_j} (\mathbf{x} - \boldsymbol{\mu}_j)(\mathbf{x} - \boldsymbol{\mu}_j)^T,$$

1 and $\mathbf{S}_B = \sum_{j=1}^c n_j (\boldsymbol{\mu}_j - \boldsymbol{\mu})(\boldsymbol{\mu}_j - \boldsymbol{\mu})^T$.

2 Since for SSS problem $d > n$, the scatter matrices will be singular. It is known that the
3 null space of \mathbf{S}_T does not contain any discriminant information [19]. Therefore, the
4 dimensionality can be reduced from d -dimensional space to r_t -dimensional space
5 (where r_t is the rank of \mathbf{S}_T) by applying PCA as a pre-processing step. The range space
6 of \mathbf{S}_T matrix, $\mathbf{U}_1 \in \mathbb{R}^{d \times r_t}$, will be used as a transformation matrix. In the reduced
7 dimensional space the scatter matrices will be given by: $\mathbf{S}_w = \mathbf{U}_1^T \mathbf{S}_w \mathbf{U}_1$ and $\mathbf{S}_b =$
8 $\mathbf{U}_1^T \mathbf{S}_b \mathbf{U}_1$. After this procedure $\mathbf{S}_w \in \mathbb{R}^{r_t \times r_t}$ and $\mathbf{S}_b \in \mathbb{R}^{r_t \times r_t}$ are reduced dimensional
9 within-class scatter matrix and reduced dimensional between-class scatter matrix,
10 respectively.

11 **4.2 Deterministic approach to regularized LDA**

12 In the SSS problem, \mathbf{S}_w matrix becomes singular and its inverse computation becomes
13 impossible. In order to overcome this drawback, the RLDA technique adds a small
14 positive constant α to all the diagonal elements of matrix \mathbf{S}_w to make it non-singular;
15 i.e., \mathbf{S}_w is replaced by $\hat{\mathbf{S}}_w = \mathbf{S}_w + \alpha \mathbf{I}$. In this section, we describe a procedure to compute
16 the regularization parameter α deterministically. In RLDA, the modified Fisher's
17 criterion is given as follows:

$$\hat{J}(\mathbf{w}, \alpha) = \frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{\mathbf{w}^T (\mathbf{S}_w + \alpha \mathbf{I}) \mathbf{w}} \quad (1)$$

1 where $\mathbf{w} \in \mathbb{R}^{r_t \times 1}$ is the orientation vector. Let us denote a function

$$2 \quad f = \mathbf{w}^T \mathbf{S}_b \mathbf{w} \quad (2)$$

3 and a constraint function

$$4 \quad g = \mathbf{w}^T (\mathbf{S}_w + \alpha \mathbf{I}) \mathbf{w} - b = 0 \quad (3)$$

5 where $b > 0$ is any constant. To find the maximum of f under the constraint, let us

6 define a function $F = f - \lambda g$, where λ is Lagrange's multiplier ($\lambda \neq 0$). By setting its

7 derivative to zero, we get

$$\frac{\partial F}{\partial \mathbf{w}} = 2\mathbf{S}_b \mathbf{w} - \lambda(2\mathbf{S}_w \mathbf{w} + 2\alpha \mathbf{w}) = 0$$

$$8 \quad \text{or} \quad \left(\frac{1}{\lambda} \mathbf{S}_b - \mathbf{S}_w\right) \mathbf{w} = \alpha \mathbf{w}. \quad (4)$$

9 Substituting value of $\alpha \mathbf{w}$ from equation (4) into equation (3), we get

$$10 \quad g = \mathbf{w}^T \mathbf{S}_w \mathbf{w} + \mathbf{w}^T \left(\frac{1}{\lambda} \mathbf{S}_b - \mathbf{S}_w\right) \mathbf{w} - b = 0$$

$$11 \quad \text{or} \quad \mathbf{w}^T \mathbf{S}_b \mathbf{w} = \lambda b. \quad (5)$$

12 Also from equation (3), we get

$$13 \quad \mathbf{w}^T (\mathbf{S}_w + \alpha \mathbf{I}) \mathbf{w} = b. \quad (6)$$

14 Dividing equation (5) by equation (6), we get

$$15 \quad \lambda = \frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{\mathbf{w}^T (\mathbf{S}_w + \alpha \mathbf{I}) \mathbf{w}} \quad (7)$$

16 The right-hand side of equation (7) is same as the criterion $\hat{J}(\mathbf{w}, \alpha)$ defined in equation

(1). The left-hand side of equation (7) is the Lagrange's multiplier (in equation (4)).

1 Since our aim is to maximize the modified Fisher's criterion $\hat{J}(\mathbf{w}, \alpha)$, we must set λ
2 equal to maximum of $\hat{J}(\mathbf{w}, \alpha)$. However, it is not possible to find the maximum of $\hat{J}(\mathbf{w}, \alpha)$
3 as α is not known to us. So, as an approximation we set λ equal to the maximum of the
4 original Fisher's criterion $(\mathbf{w}^T \mathbf{S}_b \mathbf{w} / \mathbf{w}^T \mathbf{S}_w \mathbf{w})$. In order to maximize the original Fisher's
5 criterion, we must have eigenvector \mathbf{w} to correspond to the maximum eigenvalue of
6 $\mathbf{S}_w^{-1} \mathbf{S}_b$. Since in SSS problem \mathbf{S}_w is singular and non-invertible, we approximate the
7 inverse of \mathbf{S}_w by its pseudoinverse and carry out the EVD of $\mathbf{S}_w^+ \mathbf{S}_b$ to find the highest
8 (or leading) eigenvalue, where \mathbf{S}_w^+ is the pseudoinverse of \mathbf{S}_w . Thus, if λ_{max} denotes
9 the highest eigenvalue of $\hat{J}(\mathbf{w}, \alpha)$, then

$$\begin{aligned} \lambda_{max} &= \max \left(\frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{\mathbf{w}^T (\mathbf{S}_w + \alpha \mathbf{I}) \mathbf{w}} \right) \\ &\approx \max \left(\frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}} \right) \end{aligned}$$

$$10 \quad \approx \text{largest eigenvalue of } \mathbf{S}_w^+ \mathbf{S}_b \quad (8)$$

11 Thereby, the evaluation of α can be carried out from equation (4) by doing EVD of
12 $(\frac{1}{\lambda} \mathbf{S}_b - \mathbf{S}_w)$, where $\lambda = \lambda_{max}$. The EVD of $(\frac{1}{\lambda} \mathbf{S}_b - \mathbf{S}_w)$ will give $r_b = \text{rank}(\mathbf{S}_b)$
13 eigenvalues. Since the highest eigenvalue will correspond to the most discriminant
14 eigenvector, α is the highest eigenvalue. Therefore, if EVD of $(\frac{1}{\lambda} \mathbf{S}_b - \mathbf{S}_w)$ is given by

$$15 \quad \left(\frac{1}{\lambda} \mathbf{S}_b - \mathbf{S}_w \right) = \mathbf{E} \mathbf{D}_{bw} \mathbf{E}^T \quad (9)$$

16 where $\mathbf{E} \in \mathbb{R}^{r_t \times r_t}$ is a matrix of eigenvectors and $\mathbf{D}_{bw} \in \mathbb{R}^{r_t \times r_t}$ is a diagonal matrix of

1 corresponding eigenvalues. Now the α parameter can be computed as

$$2 \quad \alpha = \max \mathbf{D}_{bw} \quad (10)$$

3 After evaluating α , orientation vector \mathbf{w} can be obtained by performing the EVD of

4 $(\mathbf{S}_w + \alpha\mathbf{I})^{-1}\mathbf{S}_b$; i.e., from

$$5 \quad (\mathbf{S}_w + \alpha\mathbf{I})^{-1}\mathbf{S}_b\mathbf{w} = \gamma\mathbf{w}. \quad (11)$$

6 From the r_b eigenvectors obtained by this EVD, h ($\leq r_b$) eigenvectors corresponding to

7 h highest eigenvalues are used to form the orientation matrix \mathbf{W} .

8

9 It can be shown from Lemma 1 that for DRLDA technique, its maximum eigenvalue is

10 approximately equal to the highest (finite) eigenvalue of Fisher's criterion.

11

12 Lemma 1: *The highest eigenvalue of DRLDA is approximately equivalent to the highest*

13 *(finite) eigenvalue of Fisher's criterion.*

14 Proof 1: From equation 11,

$$15 \quad \mathbf{S}_b\mathbf{w}_j = \gamma_j(\mathbf{S}_w + \alpha\mathbf{I})\mathbf{w}_j, \quad (12)$$

16 where α is the maximum eigenvalue of $(1/\lambda_{max}\mathbf{S}_b - \mathbf{S}_w)$ (from equation 4); $\lambda_{max} \geq 0$

17 is approximately the highest eigenvalue of Fisher's criterion $\mathbf{w}^T\mathbf{S}_b\mathbf{w}/\mathbf{w}^T\mathbf{S}_w\mathbf{w}$ (since

18 λ_{max} is the largest eigenvalue of $\mathbf{S}_w^+\mathbf{S}_b$) [46]; $j = 1 \dots r_b$ and $r_b = \text{rank}(\mathbf{S}_b)$. Substituting

1 $\alpha \mathbf{w} = (1/\lambda_{max} \mathbf{S}_b - \mathbf{S}_w) \mathbf{w}$ (from equation 4, where $\lambda = \lambda_{max}$) into equation 12, we get,

2
$$\mathbf{S}_b \mathbf{w}_m = \gamma_m \mathbf{S}_w \mathbf{w}_m + \gamma_m (1/\lambda_{max} \mathbf{S}_b - \mathbf{S}_w) \mathbf{w}_m,$$

3 or
$$(\lambda_{max} - \gamma_m) \mathbf{S}_b \mathbf{w}_m = 0.$$

4 where $\gamma_m = \max(\gamma_j)$ and \mathbf{w}_m is the corresponding eigenvector. Since $\mathbf{S}_b \mathbf{w}_m \neq 0$ (from
 5 equation 5), $\gamma_m = \lambda_{max}$ and $\gamma_j < \lambda_{max}$, where $j \neq m$. This concludes the proof.

6

7 Corollary 1: The value of regularization parameter is non-negative; i.e., $\alpha \geq 0$ for

8 $r_w \leq r_t$, where $r_t = rank(\mathbf{S}_T)$ and $r_w = rank(\mathbf{S}_w)$.

9 Proof. Please see Appendix-III.

10

11 The summary of the DRLDA technique is given in Table 1².

12 Table 1: DRLDA technique

13 *Step 1.* Pre-processing stage: apply PCA to find the range space $\mathbf{U}_1 \in \mathbb{R}^{d \times r_t}$ of total
 14 scatter matrix \mathbf{S}_T and transform original d -dimensional data space to
 15 r_t -dimensional data space, where $r_t = rank(\mathbf{S}_T)$. Find reduced-dimensional
 16 between-class scatter matrix $\mathbf{S}_b = \mathbf{U}_1^T \mathbf{S}_B \mathbf{U}_1$ and reduced-dimensional
 17 within-class scatter matrix $\mathbf{S}_w = \mathbf{U}_1^T \mathbf{S}_W \mathbf{U}_1$, where $\mathbf{S}_b \in \mathbb{R}^{r_t \times r_t}$ and $\mathbf{S}_w \in \mathbb{R}^{r_t \times r_t}$.

18 *Step 2.* Find the highest eigenvalue λ_{max} by performing EVD of $\mathbf{S}_w^+ \mathbf{S}_b$.

19 *Step 3.* Compute EVD of $(1/\lambda_{max} \mathbf{S}_b - \mathbf{S}_w)$ to find its highest eigenvalue α .

20 *Step 4.* Compute EVD of $(\mathbf{S}_w + \alpha \mathbf{I})^{-1} \mathbf{S}_b$ to find h eigenvectors $\mathbf{w}_j \in \mathbb{R}^{r_t \times 1}$
 21 corresponding to the leading eigenvalues, where $1 \leq h \leq r_b$ and $r_b = rank(\mathbf{S}_b)$.
 22 The eigenvectors \mathbf{w}_j are column vectors of the orientation matrix $\mathbf{W}' \in \mathbb{R}^{r_t \times h}$.

23 *Step 5.* Find orientation matrix $\mathbf{W} \in \mathbb{R}^{d \times h}$ in a d -dimensional space; i.e., $\mathbf{W} = \mathbf{U}_1 \mathbf{W}'$.

² Matlab code will be provided upon acceptance of the paper on our website.

1 The computational requirement for Step 1 of the technique (Table 1) would be $O(dn^2)$;
2 for Step 2 would be $O(n^3)$; for Step 3 would be $O(n^3)$; for Step 4 would be $O(n^3)$; and,
3 for Step 5 would be $O(dn^2)$. Therefore, the total estimated for SSS case ($d \gg n$) would
4 be $O(dn^2)$.

5

6 **5. Experimental setup and results**

7 Experiments are conducted to illustrate the relative performance of the DRLDA
8 technique with respect to other techniques for the following two applications: 1) face
9 recognition and 2) cancer classification. For face recognition, two commonly known
10 datasets namely ORL dataset [37] and AR dataset [38] are utilized. The ORL dataset
11 contains 400 images of 40 people having 10 images per person. We use the
12 dimensionality of the original feature space to be 5152. The AR dataset contains 100
13 classes. We use a subset of AR dataset with 14 face images per class. We use the
14 dimensionality of feature space to be 4980. For cancer classification, 6 commonly
15 available datasets are used. All the datasets used in the experimentation are described
16 in Table 2. For some datasets, number of training samples and test samples are
17 predefined by their donors (Table 2). For these datasets, we use test samples to evaluate
18 the classification performance. For some datasets, the training and test samples are not

1 predefined. For these datasets we carried out k -fold cross-validation procedure³ to
2 compute the classification performance, where $k = 3$.

3

4 The DRLDA technique is compared with the following techniques: Null LDA (NLDA) [6],
5 cross-validation based RLDA (CV-RLDA), Pseudo-inverse LDA (PILDA) [20], Direct
6 LDA (DLDA) [7], Fisherface or PCA+LDA [4][5], Uncorrelated LDA (ULDA) [21] and
7 eigenfeature regularization (EFR) [10]. All the techniques are used to find the
8 orientation matrix $\mathbf{W} \in \mathbb{R}^{d \times c-1}$, thereby, transforming the original space to $c - 1$
9 dimensional space, where c is the number of classes. Then nearest neighbour classifier
10 (NNC) using Euclidean distance measure is used for classifying a test feature vector.

11

12 The setting up of CV-RLDA technique in our experiments is described as follows: the
13 regularization parameter α of CV-RLDA is computed by using leave-one-out
14 cross-validation procedure on the training set. This is done in two steps. In the first step,
15 we perform a coarse search for α by dividing the pre-selected range $[10^{-4}, 1] * \lambda_W$
16 (where λ_W is the maximum eigenvalue of \mathbf{S}_W) into 10 equal intervals and finding the

³ In the k -fold cross-validation procedure [39], we first partition all the available samples randomly into k roughly equal segments. Then hold out one segment as validation data and the remaining $k - 1$ segments as training data. Using the training data, we applied a discriminant technique to obtain orientation matrix and the validation data to compute classification accuracy. This partitioning of samples and computation of classification accuracy is carried out k times to evaluate average classification accuracy.

1 interval whose center value gives the best classification performance on the training set.
2 In the second step, this interval is further divided into 10 subintervals for fine search
3 and the center value of the subinterval that gives the best classification performance is
4 used as the final value of the regularization parameter. Thus, in this procedure, a total
5 of 20 α values are investigated. The regularization parameters computed by CV-RLDA
6 technique on various datasets are shown in Appendix-IV.

7

8 The classification accuracy on all the datasets using the above mentioned techniques
9 are shown in Table 3 (the highest classification accuracies obtained are depicted in bold
10 fonts). It can be seen from Table 3 that out of 8 datasets used, the number of times the
11 highest classification accuracy obtained by NLDA is 2, CV-RLDA is 5, PILDA is 1, DLDA
12 is 1, PCA+LDA is 3, ULDA is 2, EFR is 4 and DRLDA is 6. In particular, DRLDA
13 performs better than CV-RLDA for most of the datasets shown in Table 2 (it
14 outperforms CV-RLDA for 3 out of 8 datasets, shows equal classification accuracy for 3
15 datasets and is inferior to CV-RLDA in the remaining 2 datasets). Note that the
16 CV-RLDA technique when implemented in an ideal form (i.e., when α is searched in
17 the range $(0, \infty)$ with infinitely small step size) should give in principle better results
18 than the DRLDA technique. Since it is not possible for practical reasons (i.e.,

1 computational cost is infinitely large), a finite range is used in CV-RLDA technique. As a
2 result, DRLDA technique is performing here better in terms of classification accuracy
3 for majority of datasets. In addition, the computational cost of CV-RLDA technique
4 (with α being searched in the finite range) is considerably higher than the DRLDA
5 technique as shown in Table 4. Here, we measure the CPU time taken by its ‘Matlab’
6 program on a Sony computer (*core i7 processor at 2.8GHz*).

7 Table 2: Datasets used in the experimentation

Datasets	Class	Dimension	Number of available samples	Number of training samples	Number of test samples
Acute Leukemia [40]	2	7129	72	38	34
ALL subtype [41]	7	12558	327	215	112
GCM [42]	14	16063	198	144	54
Lung Adenocarcinoma [43]	3	7129	96 ([67,19,10])*	-	-
MLL [44]	3	12582	72	57	15
SRBCT [45]	4	2308	83	63	20
Face ORL [37]	40	5152	400 (10/class)	-	-
Face AR [38]	100	4980	1400 (14/class)	-	-

8 * The values in the square parenthesis indicate number of samples per class.
9

10
11 Table 3: Classification accuracy (in percentage) on datasets using various techniques.

Database	NLDA	CV-RLDA	PILDA	DLDA	PCA+LDA	ULDA	EFR	DRLDA
Acute Leukemia	97.1	97.1	73.5	97.1	100.0	97.1	100.0	100.0
ALL subtype	86.6	95.5	62.5	93.8	80.7	82.1	86.6	93.8
GCM	70.4	74.1	46.3	59.3	59.3	66.7	68.5	70.4
Lung Adeno.	81.7	81.7	74.2	72.0	81.7	80.7	83.9	86.0
MLL	100.0	100.0	80.0	100.0	100.0	100.0	100.0	100.0
SRBCT	100.0	100.0	85.0	80.0	100.0	100.0	100.0	100.0
Face ORL	96.9	97.2	96.4	96.7	92.8	92.5	96.7	97.2
Face AR	95.7	96.3	97.3	96.3	94.9	95.8	97.3	97.3

12
13 Table 4: The comparison of cputime (in seconds) of CV-RLDA and DRLDA techniques.

Database	CV-RLDA CPU Time	DRLDA CPU Time
Acute Leukemia	4.68	0.07
ALL subtype	1021.9	1.90
GCM	265.0	1.26
Lung Adeno.	57.9	0.48
MLL	13.6	0.24
SRBCT	17.0	0.08

Face ORL	7396.1	7.41
Face AR	739,380	89.9

1

2 Furthermore, various techniques using artificial data are experimented. For this, we
3 have created a 2-class problem with initial dimensions $d = 10, 25, 30, 50,$ and 100. In
4 order to have ill-posed problem, we generated only 3 samples per class. The
5 dimensionality is reduced from d to 1 for all the techniques and then nearest
6 neighbour classifier is used to evaluate the performance in terms of classification
7 accuracy. For each dimension d , the data is created 100 times to compute average
8 classification accuracy. Table 5 depicts the average classification accuracy over 100 runs.
9 It can be observed from Table 5 that EFR technique is not able to perform because of
10 scarce samples. The DRLDA technique and CV-RLDA technique are performing similar.
11 Pseudoinverse technique (PILDA) is performing the lowest as there is not enough
12 information in the range space of scatter matrices.

13

14 Table 5: Classification accuracy (in percentage) on artificial dataset using various
15 techniques.

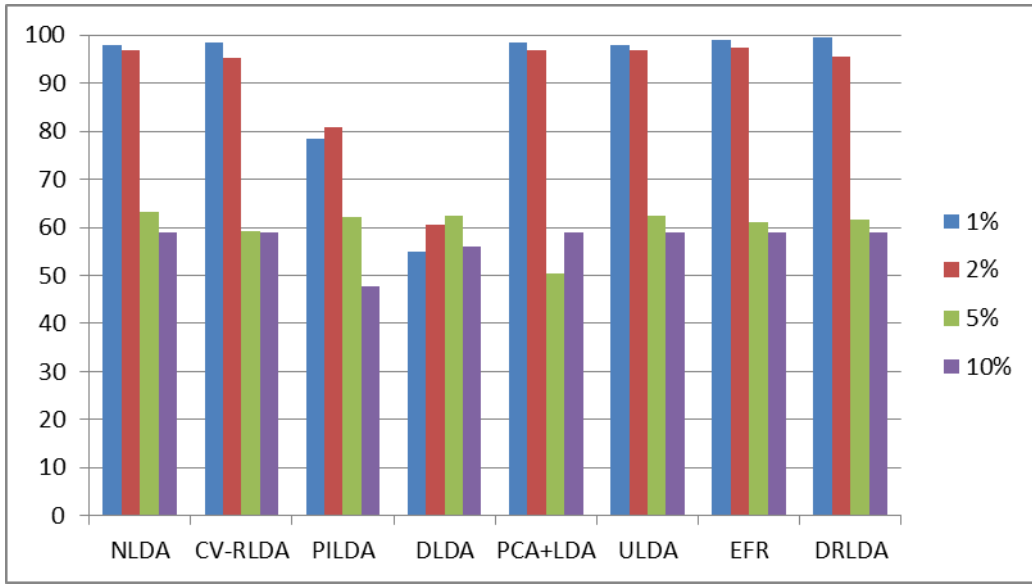
Dimension	NLDA	CV-RLDA	PILDA	DLDA	PCA+LDA	ULDA	EFR	DRLDA
10	84.3	87.2	66.2	87.8	85.7	84.3	-	87.2
25	95.0	96.7	58.2	96.3	93.7	95.0	-	97.2
30	96.0	97.8	52.8	95.8	96.2	96.0	-	98.0
50	98.8	99.2	49.5	99.2	98.7	98.8	-	99.2
100	100	100	50	99.5	99.8	100	-	100

16

17 We have also carried out sensitivity analysis with respect to the classification accuracy.
18 For this purpose, we use Acute Leukemia dataset as a prototype and contaminated the
19 dataset by adding Gaussian noise. We then applied techniques again to evaluate
20 classification performance by using nearest neighbor classifier. The generated noise
21 levels are 1%, 2%, 5% and 10% of the standard deviation of the original feature values.
22 The noisy data has been generated 10 times to compute average classification accuracy.
23 The results are shown in Figure 1. It can be observed from Figure 1 that for low level
24 noise the degradation in classification performance is not enough. But when the noise

1 level increases the classification accuracy deteriorates. The performance of PILDA and
 2 DLDA techniques are lower than other techniques. However, most of the techniques try
 3 to maintain the discriminant information in the noisy environment.

4



5

6 Figure 1: Sensitivity analysis of various techniques on Acute Leukemia dataset at
 7 different noise levels. The y-axis depicts average classification accuracy and x-axis
 8 depicts the techniques used. The noise levels are 1%, 2%, 5% and 10%.

9

10 6. Discussion

11 In order to compare the performance in terms of classification accuracy we compared 7
 12 well known techniques with DRLDA. These techniques compute the orientation matrix

13 \mathbf{W} by utilizing different combinations of informative spaces (i.e., null space of \mathbf{S}_W , range
 14 space of \mathbf{S}_W and range space of \mathbf{S}_B). Each informative space contains a certain level of

15 discriminant information. Theoretically, it is effective to utilize all the informative

1 spaces for the computation of orientation matrix for better generalization capability.
2 How well a technique is combining these spaces would determine its generalization
3 capability. It has been shown that usually the null space of \mathbf{S}_W contains more
4 discriminant information than the range space of \mathbf{S}_W [6][8][22][34]. Therefore, it is
5 likely that a technique that utilizes null space of \mathbf{S}_W effectively, may perform better (in
6 generalization capability) than a technique which does not use the null space of \mathbf{S}_W .
7
8 From the techniques that we have used the NLDA technique employs null space of \mathbf{S}_W
9 and range space of \mathbf{S}_B . Whereas PILDA, DLDA and PCA+LDA techniques employ range
10 space of \mathbf{S}_W and range space of \mathbf{S}_B . Provided the techniques extract the maximum
11 possible information from the spaces they employed then NLDA should beat PILDA,
12 DLDA and PCA+LDA techniques. From Table 3, we can see that NLDA is
13 outperforming PILDA in 7 out of 8 cases. Comparing the classification accuracies of
14 NLDA and DLDA, we can see that NLDA is outperforming DLDA in 4 out of 8 cases and
15 in 2 cases the performance are identical. In a similar way NLDA is surpassing
16 PCA+LDA in 4 out of 8 cases and in 3 cases the performance are identical. On the other
17 hand, the ULDA technique also employs the same spaces as of NLDA technique,
18 however, the classification performance of ULDA is inferior to NLDA (only in 1 out of 8

1 cases ULDA is beating NLDA). This means that orthogonal \mathbf{W} is more effective than
2 uncorrelated \mathbf{W} .
3
4 The other three techniques (CV-RLDA, EFR and DRLDA) employ three spaces; namely,
5 null space of \mathbf{S}_W , range space of \mathbf{S}_W and range space of \mathbf{S}_B . Intuitively, these three
6 techniques contain more discriminant information than above mentioned 5 techniques.
7 However, different strategies of using the three spaces would result in different level of
8 generalization capabilities. In CV-RLDA, the estimation of regularization parameter α
9 depends upon the range of α values being explored (which is restricted due to limited
10 computation time), the cross-validation procedure (e.g. leave-one-out, k -fold) being
11 employed and the classifier used. On the other hand, EFR and DRLDA techniques do
12 not have this problem. The EFR technique utilizes an intuitive model for extrapolating
13 the eigenvalues of range space of \mathbf{S}_W to the null space of \mathbf{S}_W . This way it captures all
14 the spaces. However, the model used for extrapolation is rather arbitrary and it is not
15 necessary that it is an optimum model. The DRLDA technique captures the information
16 from all the spaces by deterministically finding the optimal α parameter from the
17 training samples. From Table 3, it can be observed that EFR is surpassing CV-RLDA in
18 3 out of 8 cases and exhibiting identical classification accuracies in 2 cases. Similarly,

1 DRLDA is outperforming CV-RLDA in 3 out of 8 cases and giving equal results in 3
2 cases. From Table 3 and Table 4, we can also observe that though the classification
3 accuracy of CV-RLDA is high (which depends on the search of the regularization
4 parameter), its computational time is extremely large.

5

6 Thus we have shown that DRLDA technique is performing better than other LDA
7 techniques for the SSS problem. We can intuitively explain its better performance as
8 follows. In the DRLDA technique, we are maximizing the modified Fisher's criterion; i.e.,
9 the ratio of between-class scatter and within-class scatter (see equation 1). To get the α
10 parameter, we are maximizing the difference between the between-class scatter and
11 within-class scatter (see equation 4). Thus, we are combining two different philosophies
12 of LDA mechanism in our DRLDA technique and this is helping us in getting better
13 performance.

14

15 **7. Conclusion**

16 The paper presented a deterministic approach of computing regularized LDA. It avoids
17 the use of the heuristic (cross-validation) procedure for computing the regularization
18 parameter. The technique has been experimented on a number of datasets and

1 compared with several popular techniques. The DRLDA technique exhibits highest
2 classification accuracy for 6 out of 8 datasets and its computational cost is significantly
3 less than CV-RLDA technique.

4 **Appendix-I**

5 In this appendix, the generalization capabilities of Zhao’s model and Friedman’s model
6 of CV-RLDA are demonstrated on several datasets. In order to do this, first we project
7 the original feature vectors onto the range space of total-scatter matrix as a
8 pre-processing step. Then we employ reduced dimensional within-class scatter matrix
9 $\hat{\mathbf{S}}_w$ for the two models of CV-RLDA (see Section 4.1 for details about reduced
10 dimensional matrices). In the first model of CV-RLDA, \mathbf{S}_w is approximated as
11 $\hat{\mathbf{S}}_w = \mathbf{S}_w + \alpha \mathbf{I}$ and in the second model \mathbf{S}_w is approximated as $\hat{\mathbf{S}}_w = (1 - \alpha)\mathbf{S}_w + \alpha \mathbf{I}$. For
12 brevity, we refer the former model of CV-RLDA as CV-RLDA-1 and the latter model as
13 CV-RLDA-2. Table A1 depicts the classification performance of these two models. The
14 details of the datasets and the selection of the regularization parameter α can be found
15 in Section 4.

16

17 It can be seen from Table A1 that CV-RLDA-1 exhibits relatively better classification
18 performance than CV-RLDA-2.

1
2
3
4
5

6 Table A1: Classification accuracy (in percentage) on test set using CV-RLDA-1 and
7 CV-RLDA-2 techniques.

Database	CV-RLDA-1	CV-RLDA-2
Acute Leukemia	97.1	97.1
ALL subtype	95.5	86.6
GCM	74.1	70.4
MLL	100.0	100.0
SRBCT	100.0	100.0

8

9 **Appendix-II**

10 In this appendix, for RLDA technique we show the sensitivity of classification accuracy
11 when selecting the regularization parameter, α . For this purpose we use four values of
12 α . These are $\delta = [0.001, 0.01, 0.1, 1]$, where $\alpha = \delta * \lambda_W$ and λ_W is the maximum
13 eigenvalue of within-class scatter matrix. We applied 3-fold cross-validation procedure
14 on a number of datasets and shown the results in Table A2.

15

16 Table A2: Classification accuracy (in percentage) using 3-fold cross-validation procedure
17 (the highest classification accuracies obtained are depicted in bold fonts).

Database	$\delta = 0.001$	$\delta = 0.01$	$\delta = 0.1$	$\delta = 1$
Acute Leukemia	98.6	98.6	98.6	100
ALL subtype	90.3	90.3	86.0	69.2
GCM	72.7	74.3	76.5	59.0
Lung Adeno.	81.7	80.7	85.0	80.7
MLL	95.7	95.7	95.7	95.7

SRBCT	100.0	100.0	100.0	96.2
Face ORL	96.9	96.9	96.9	96.9
Face AR	95.8	97.9	96.3	81.8

1 It can be observed from the table that the different values of the regularization
2 parameter give different classification accuracies and therefore, the choice of the
3 regularization parameter affects the classification performance. Thus, it is important to
4 select the regularization parameter correctly to get the good classification performance.

5

6 To do this, a cross-validation approach is usually opted. The α parameter is searched in
7 the pre-defined range and the value of α which gives the best classification
8 performance on the training set is selected. It is assumed that the optimum value of α
9 will give the best generalization capability; i.e., the best classification performance on
10 the test set.

11

12 **Appendix III**

13 Corollary 1: The value of regularization parameter is non-negative; i.e., $\alpha \geq 0$ for
14 $r_w \leq r_t$, where $r_t = \text{rank}(\mathbf{S}_T)$ and $r_w = \text{rank}(\mathbf{S}_w)$.

15 Proof 1: From equation 1, we can write

$$16 \quad J = \frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{\mathbf{w}^T (\mathbf{S}_w + \alpha \mathbf{I}) \mathbf{w}}, \quad \text{A1}$$

17 where $\mathbf{S}_b \in \mathbb{R}^{r_t \times r_t}$ and $\mathbf{S}_w \in \mathbb{R}^{r_t \times r_t}$. We can rearrange the above expression as

$$18 \quad \mathbf{w}^T \mathbf{S}_b \mathbf{w} = J \mathbf{w}^T (\mathbf{S}_w + \alpha \mathbf{I}) \mathbf{w} \quad \text{A2}$$

19

1 The eigenvalue decomposition (EVD) of \mathbf{S}_w matrix (assuming $r_w < r_t$) can be given as

2 $\mathbf{S}_w = \mathbf{U}\mathbf{\Lambda}^2\mathbf{U}^T$, where $\mathbf{U} \in \mathbb{R}^{r_t \times r_t}$ is an orthogonal matrix, $\mathbf{\Lambda}^2 = \begin{bmatrix} \mathbf{\Lambda}_w^2 & 0 \\ 0 & 0 \end{bmatrix} \in \mathbb{R}^{r_t \times r_t}$ and

3 $\mathbf{\Lambda}_w = \text{diag}(q_1^2, q_2^2, \dots, q_{r_w}^2) \in \mathbb{R}^{r_w \times r_w}$ are diagonal matrices (as $r_w < r_t$). The eigenvalues
4 $q_k^2 > 0$ for $k = 1, 2, \dots, r_w$. Therefore,

5

6 $\mathbf{S}'_w = (\mathbf{S}_w + \alpha\mathbf{I}) = \mathbf{U}\mathbf{D}\mathbf{U}^T$, where $\mathbf{D} = \mathbf{\Lambda}^2 + \alpha\mathbf{I}$

7 or $\mathbf{D}^{-1/2}\mathbf{U}^T\mathbf{S}'_w\mathbf{U}\mathbf{D}^{-1/2} = \mathbf{I}$ A3

8

9 The between class scatter matrix \mathbf{S}_b can be transformed by multiplying $\mathbf{U}\mathbf{D}^{-1/2}$ on the
10 right side and $\mathbf{D}^{-1/2}\mathbf{U}^T$ on the left side of \mathbf{S}_b as $\mathbf{D}^{-1/2}\mathbf{U}^T\mathbf{S}_b\mathbf{U}\mathbf{D}^{-1/2}$. The EVD of this
11 matrix will give

12 $\mathbf{D}^{-1/2}\mathbf{U}^T\mathbf{S}_b\mathbf{U}\mathbf{D}^{-1/2} = \mathbf{E}\mathbf{D}_b\mathbf{E}^T$, A4

13 where $\mathbf{E} \in \mathbb{R}^{r_t \times r_t}$ is an orthogonal matrix and $\mathbf{D}_b \in \mathbb{R}^{r_t \times r_t}$ is a diagonal matrix.

14 Equation A4 can be rearranged as

15 $\mathbf{E}^T\mathbf{D}^{-1/2}\mathbf{U}^T\mathbf{S}_b\mathbf{U}\mathbf{D}^{-1/2}\mathbf{E} = \mathbf{D}_b$, A5

16 Let the leading eigenvalue of \mathbf{D}_b is γ and its corresponding eigenvector is $\mathbf{e} \in \mathbf{E}$. Then
17 equation A5 can be rewritten as

18 $\mathbf{e}^T\mathbf{D}^{-1/2}\mathbf{U}^T\mathbf{S}_b\mathbf{U}\mathbf{D}^{-1/2}\mathbf{e} = \gamma$, A6

19 The eigenvector \mathbf{e} can be multiplied right side and \mathbf{e}^T on left side of equation A3, we
20 get

21 $\mathbf{e}^T\mathbf{D}^{-1/2}\mathbf{U}^T\mathbf{S}'_w\mathbf{U}\mathbf{D}^{-1/2}\mathbf{e} = 1$ A7

22

23 It can be seen from equations A3 and A5 that matrix $\mathbf{W} = \mathbf{U}\mathbf{D}^{-1/2}\mathbf{E}$ diagonalizes both
24 \mathbf{S}_b and \mathbf{S}'_w , simultaneously. Also vector $\mathbf{w} = \mathbf{U}\mathbf{D}^{-1/2}\mathbf{e}$ simultaneously gives γ and unity
25 eigenvalues in equations A6 and A7. Therefore, \mathbf{w} is a solution of equation A2.

26 Substituting $\mathbf{w} = \mathbf{U}\mathbf{D}^{-1/2}\mathbf{e}$ in equation A2, we get

27

1 $J = \gamma$; i.e., \mathbf{w} is a solution of equation A2.

2

3 From Lemma 1, the maximum eigenvalue of expression $(\mathbf{S}_w + \alpha\mathbf{I})^{-1}\mathbf{S}_b\mathbf{w} = \gamma\mathbf{w}$ is
4 $\gamma_m = \lambda_{max} > 0$ (i.e., real, positive and finite). Therefore, the eigenvectors corresponding
5 to this positive γ_m should also be in real hyperplane (i.e., the components of the vector
6 \mathbf{w} have to have real values). Since $\mathbf{w} = \mathbf{U}\mathbf{D}^{-1/2}\mathbf{e}$ with \mathbf{w} to be in real hyperplane, we
7 must have $\mathbf{D}^{-1/2}$ to be real.

8

9 Since $\mathbf{D} = \mathbf{\Lambda}^2 + \alpha\mathbf{I} = \text{diag}(q_1^2 + \alpha, q_2^2 + \alpha, \dots, q_{r_w}^2 + \alpha, \alpha, \dots, \alpha)$, we have

10
$$\mathbf{D}^{-1/2} = \text{diag}(1/\sqrt{q_1^2 + \alpha}, 1/\sqrt{q_2^2 + \alpha}, \dots, 1/\sqrt{q_{r_w}^2 + \alpha}, 1/\sqrt{\alpha}, \dots, 1/\sqrt{\alpha}).$$

11 Therefore, the elements of $\mathbf{D}^{-1/2}$, must satisfy $1/\sqrt{q_k^2 + \alpha} > 0$ and $1/\sqrt{\alpha} > 0$ for
12 $k = 1, 2, \dots, r_w$ (note $r_w < r_t$); i.e., α cannot be negative or $\alpha > 0$. Furthermore, if $r_w = r_t$
13 then matrix \mathbf{S}_w will be a non-singular matrix and its inverse will exist. In this case,
14 regularization is not required and therefore $\alpha = 0$. Thus, $\alpha \geq 0$ for $r_w \leq r_t$. This
15 concludes the proof.

16

17 **Appendix IV**

18 In this appendix, we show computed value of CV-RLDA technique. The value of α is
19 computed by first doing a coarse search on a predefined range to find a coarse value.
20 After this, a fine search is conducted using this coarse value to get the regularization
21 parameter. In this experiment, we use $\alpha = \delta * \lambda_w$ where $\delta = [10^{-4}, 1]$ and λ_w is the
22 highest eigenvalue of within-class scatter matrix. The values are depicted in Table A3.
23 In addition, we have also shown regularization parameters computed by DRLDA
24 technique as a reference.

25

26

27

1 Table A3: Computed values of regularization parameter for CV-RLDA and DRLDA on
 2 various datasets

Database	CV-RLDA δ	CV-RLDA α	DRLDA α
Acute Leukemia	0.0057	935.3	6.54×10^9
ALL subtype	0.5056	5.17×10^5	1.11×10^{11}
GCM	0.0501	2.42×10^4	1.34×10^9
MLL	0.0057	2621.5	2.98×10^{10}
SRBCT	0.1056	33.01	5715.2

3

4

5 **Reference**

6 [1] R.O. Duda and P.E. Hart, Pattern classification and scene analysis, Wiley, New York, 1973.
 7 [2] K. Fukunaga, Introduction to statistical pattern recognition. Academic Press Inc.,
 8 Hartcourt Brace Jovanovich, Publishers. 1990.
 9 [3] J.H. Friedman, “Regularized discriminant analysis”, *Journal of the American Statistical*
 10 *Association*, vol. 84, no. 405, pp. 165-175, 1989.
 11 [4] D.L. Swets, and J. Weng, “Using discriminative eigenfeatures for image retrieval”, *IEEE*
 12 *Transactions on Pattern Analysis and Machine Intelligence*, 1996, 18, (8), pp. 831-836.
 13 [5] P.N. Belhumeur, J.P. Hespanha and D.J. Kriegman, “Eigenfaces vs. fisherfaces:
 14 recognition using class specific linear projection”, *IEEE Trans. Pattern Analysis and*
 15 *Machine Intelligence.*, vol. 19, no. 7, pp. 711–720, 1997.
 16 [6] L.-F. Chen, H.-Y.M. Liao, M.-T. Ko, J.-C. Lin, and G.-J. Yu, “A new LDA-based face
 17 recognition system which can solve the small sample size problem”, *Pattern Recognition*,
 18 vol. 33, pp. 1713-1726, 2000.
 19 [7] H. Yu and J. Yang, “A direct LDA algorithm for high-dimensional data-with application to
 20 face recognition”, *Pattern Recognition*, vol. 34, pp. 2067-2070, 2001.
 21 [8] J. Ye, “Characterization of a family of algorithms for generalized discriminant analysis on
 22 undersampled problems”, *Journal of Machine Learning Research*, vol. 6, pp. 483-502,
 23 2005.
 24 [9] A. Sharma and K.K. Paliwal, “A gradient linear discriminant analysis for small sample
 25 sized problem”, *Neural Processing Letters*, vol. 27, pp 17-24, 2008.
 26 [10] X. Jiang, B. Mandal and A. Kot, “Eigenfeature regularization and extraction in face
 27 recognition”, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 30,
 28 no. 3, pp. 383-394, 2008.
 29 [11] K.K. Paliwal and A. Sharma, “Improved direct LDA and its application to DNA gene

- 1 microarray data”, *Pattern Recognition Letters*, vol. 31, issue 16, pp. 2489-2492, 2010.
- 2 [12] W. Zhao, R. Chellappa, P.J. Phillips, Subspace linear discriminant analysis for face
3 recognition, Technical Report CAR-TR-914, CS-TR-4009, University of Maryland at
4 College Park, USA, 1999.
- 5 [13] D.Q. Dai and P.C., Yuen, “Regularized discriminant analysis and its application to face
6 recognition”, *Pattern Recognition*, vol. 36, no. 3, pp. 845-847, 2003.
- 7 [14] D.Q. Dai and P.C., Yuen, “Face recognition by regularized discriminant analysis”, *IEEE*
8 *Transactions of SMC Part B*, vol. 37, issue 4, pp. 1080-1085, 2007.
- 9 [15] Y. Guo, T. Hastie, and R. Tibshirani, “Regularized discriminant analysis and its
10 application in microarrays”, *Biostatistics*, vol. 8, no. 1, pp. 86-100, 2007.
- 11 [16] W. Zhao, R. Chellappa, A. Krishnaswamy, “Discriminant analysis of principal components
12 for face recognition”, Proc. Thir Int. Conf. on Automatic Face and Gesture Recognition, pp.
13 336-341, Nara, Japan, 1998.
- 14 [17] W. Zhao, R. Chellappa, and P.J. Phillips, “Face recognition: a literature survey”, *ACM*
15 *Computing Surveys*, vol. 35, no. 4, pp. 399-458, 2003.
- 16 [18] T. Hastie, R. Tibshirani and J. Friedman, *The elements of statistical learning*, Springer, NY,
17 USA, 2001.
- 18 [19] R. Huang, Q. Liu, H. Lu, and S. Ma, “Solving the Small Sample Size Problem of LDA”,
19 *Proceedings of ICPR*, vol. 3, pp. 29-32, 2002.
- 20 [20] Q. Tian, M. Barbero, Z.H. Gu and S.H. Lee, ‘Image classification by the Foley-Sammon
21 transform’, *Optical Engineering*, vol. 25, no. 7, pp. 834-840, 1986.
- 22 [21] J. Ye, R. Janardan, Q. Li, and H. Park, “Feature extraction via generalized uncorrelated
23 linear discriminant analysis” *The Twenty-First International Conference on Machine*
24 *Learning*, pp. 895–902, 2004.
- 25 [22] K.K. Paliwal and A. Sharma, “Approximate LDA technique for dimensionality reduction
26 in the small sample size case”, *Journal of Pattern Recognition Research*, vol. 6, no. 2, pp.
27 298-306, 2011.
- 28 [23] J. Yang, D. Zhang and J.-Y. Yang, “A generesed K-L expansion method which can deal
29 with small samples size and high-dimensional problems”, *Pattern Analysis Application*,
30 vol. 6, pp. 47-54, 2003.
- 31 [24] D. Chu and G.S. Thye, “A new and fast implementation for null space based linear
32 discriminant analysis’, *Pattern Recognition*, vol. 43, pp. 1373-1379, 2010.
- 33 [25] A. Sharma and K.K. Paliwal, “A new perspective to null linear discriminant analysis
34 method and its fast implementation using random matrix multiplication with scatter
35 matrices”, *Pattern Recognition*, vol. 45, issue 6, pp. 2205-2212, 2012.
- 36 [26] J. Ye and T. Xiong, “Computational and theoretical analysis of null space and orthogonal

- 1 linear discriminant analysis”, *Journal of Machine Learning Research*, vol. 7, pp.
2 1183-1204, 2006.
- 3 [27] J. Liu, S.C. Chen, X.Y. Tan, “Efficient Pseudo-inverse Linear Discriminant Analysis and
4 its Nonlinear Form for Face Recognition”, *International Journal of Pattern Recognition
5 and Artificial Intelligence*, vol. 21, no. 8, pp. 1265-1278, 2007.
- 6 [28] K.K. Paliwal and A. Sharma, “Improved Pseudoinverse Linear Discriminant Analysis
7 Method for Dimensionality Reduction”, *International Journal of Pattern Recognition and
8 Artificial Intelligence*, 2011, DOI No: 10.1142/S0218001412500024.
- 9 [29] J. Lu, K. Plataniotis and A. Venetsanopoulos, “Face recognition using kernel direct
10 discriminant analysis algorithms”, *IEEE Transactions on Neural Networks*, vol. 14, no. 1,
11 pp. 117-126, 2003.
- 12 [30] J. Lu, K. Plataniotis and A. Venetsanopoulos, “Regularization studies of linear
13 discriminant analysis in small sample size scenarios with application to face recognition”,
14 *Pattern Recognition Letters*, vol. 26, no. 2, pp. 181-191, 2005.
- 15 [31] F. Song, D. Zhang, J. Wang, H. Liu and Q. Tao, “A parameterized direct LDA and its
16 application to face recognition”, *Neurocomputing*, vol. 71, pp. 191-196, 2007.
- 17 [32] A. Sharma, A. and K.K. Paliwal, “Regularisation of eigenfeatures by extrapolation of
18 scatter-matrix in face-recognition problem”, *Electronics Letters*, IEE, vol. 46, no. 10, pp
19 450-475, 2010.
- 20 [33] J. Lu, K.N. Plataniotis, A.N. Venetsanopoulos, “Regularized discriminant analysis for the
21 small sample”, *Pattern Recognition Letters*, vol. 24, pp. 3079-3087, 2003.
- 22 [34] A. Sharma and K.K. Paliwal, “A two-stage linear discriminant analysis for
23 face-recognition”, *Pattern Recognition Letters*, vol. 33, issue 9, pp. 1157-1162, 2012.
- 24 [35] A. Sharma, and K.K. Paliwal, “Rotational linear discriminant analysis technique for
25 dimensionality reduction”, *IEEE Transactions on Knowledge and Data Engineering*, vol.
26 20, no. 10, pp 1336-1347, 2008.
- 27 [36] A. Sharma, K.K. Paliwal, Onwubolu, G.C., “Class-dependent PCA, LDA and MDC: a
28 combined classifier for pattern classification”, *Pattern Recognition*, vol. 39, issue 7, 2006,
29 pp. 1215-1229.
- 30 [37] F. Samaria and A. Harter, “Parameterization of a stochastic model for human face
31 identification”, *Proc. Second IEEE Workshop Applications of Comp. Vision*, pp. 138-142,
32 1994.
- 33 [38] A.M. Martinez, “Recognizing imprecisely localized, partially occluded, and expression
34 variant faces from a single sample per class”, *IEEE Transactions on Pattern Analysis and
35 Machine Intelligence*, vol. 24, no. 6, pp. 748-763, 2002.
- 36 [39] E. Alpaydin, *Introduction to machine learning*, MIT Press, 2004.

- 1 [40] T.R. Golub, D.K. Slonim, P. Tamayo, C. Huard, M. Gaasenbeek, J.P. Mesirov, H. Coller,
2 M.L. Loh, J.R. Downing, M.A. Caligiuri, C.D. Bloomfield and E.S. Lander, “Molecular
3 classification of cancer: class discovery and class prediction by gene expression
4 monitoring”, *Science*, vol. 286, 531-537, 1999.
- 5 [41] E.J. Yeoh , M.E. Ross, S.A. Shurtleff , W.K. Williams, D. Patel, R. Mahfouz, F.G. Behm,
6 S.C. Raimondi, M.V. Relling, A. Patel, C. Cheng, D. Campana, D. Wilkins, X. Zhou, J. Li,
7 H. Liu, C.H. Pui, W.E. Evans, C. Naeve, L. Wong, J.R. Downing, “Classification, subtype
8 discovery, and prediction of outcome in pediatric acute lymphoblastic leukemia by gene
9 expression profiling”, *Cancer*, vol. 1, no. 2, pp 133-143, 2002.
- 10 [42] S. Ramaswamy, P. Tamayo, R. Rifkin, S. Mukherjee, C.-H. Yeang, M. Angelo, C. Ladd,
11 M. Reich, E. Latulippe, J.P. Mesirov, T. Poggio, W. Gerald, M.. Loda, E.S. Lander and
12 T.R. Golub, “Multiclass cancer diagnosis using tumor gene expression signatures”, *Proc.*
13 *Natl. Acad. Sci. USA*, vol. 98, no. 26, pp 15149-15154, 2001.
- 14 [43] D.G. Beer, S.L.R. Kardia, C.-C. Huang, T.J. Giordano, A.M. Levin, D.E. Misek, L. Lin, G.
15 Chen, T.G. Gharib, D.G. Thomas, M.L. Lizyness, R. Kuick, S. Hayasaka, J.M.G. Taylor,
16 M.D. Iannettoni, M.B. Orringer and S. Hanash, “Gene-expression profiles predict survival
17 of patients with lung adenocarcinoma”, *Nature Medicine*, vol. 8, pp. 816-824, 2002
- 18 [44] S.A. Armstrong, J.E. Staunton, L.B. Silverman, R. Pieters, M.L. den Boer, M.D. Minden,
19 S.E. Sallan, E.S. Lander, T.R. Golub, and S.J. Korsmeyer, “MLL translocations specify a
20 distinct gene expression profile that distinguishes a unique leukemia”, *Nature Genetics*,
21 vol. 30, pp 41-47, 2002.
- 22 [45] J. Khan, J.S. Wei, M. Ringner, L.H. Saal, M. Ladanyi, F. Westermann, F. Berthold, M.
23 Schwab, C.R. Antonescu, C. Peterson and P.S. Meltzer, “Classification and diagnostic
24 prediction of cancers using gene expression profiling and artificial neural network”,
25 *Nature Medicine*, vol. 7, pp. 673-679, 2001.
- 26 [46] Liu, J., Chen, S.C., Tan, X.Y., Efficient pseudo-inverse linear discriminant analysis and its
27 nonlinear form for face recognition, *Int. J. Patt. Recogn. Artif. Intell.* vol. 21, no. 8, pp.
28 1265-1278, 2007.
- 29 [47] Sharma, A., and Paliwal, K.K., Linear discriminant analysis for small sample size
30 problem: an overview, *Int. J. Mach. Learn. & Cyber.*, 2014, DOI
31 10.1007/s13042-013-0226-9.
- 32 [48] Tao, D., Li, X., Wu, X. and Maybank, S., General tensor discriminant analysis and gabor
33 features for gait recognition, *IEEE Trans. Patt. Anal. Mach. Intell.*, vol. 29, issue 10, pp.
34 1700-1715, 2007.

- 1 [49] Bian, W. and Tao, D., Constrained empirical risk minimization framework for distance
2 metric learning, *IEEE Trans. Neural Net. And Learning Sys.*, vol. 23, no. 8, pp. 1194-1205,
3 2012.
- 4 [50] Zhou, T. and Tao, D., Double shrinking sparse dimension reduction, *IEEE Trans. on*
5 *Image. Proc.*, vol. 22, issue 1, pp. 244-257, 2013.
- 6 [51] Tao, D., Li, X., Wu, X. and Maybank, S.J., Geometric mean for subspace selection, *IEEE*
7 *Trans. Patt. Anal. Mach. Learn.*, vol. 31, issue 2, pp. 260-274, 2009.
- 8 [52] Sharma, A., Paliwal, K.K., Imoto, S. and Miyano, S., Principal component analysis using
9 QR decomposition, *Int. Jnr. of Mach. Learn. And Cybernetics*, vol. 4, no. 6, pp. 679-683,
10 2013.
- 11 [53] Sharma, A. and Paliwal, K.K., Fast principal component analysis using fixed-point
12 algorithm, vol. 28, no. 10, pp. 1151-1155, 2007.
- 13 [54] Wang, S.-J., Chen, H.-L., Peng, X.-J. and Zhou, C.-G., Exponential locality preserving
14 projections for small sample size problem, *Neurocomputing*, vol. 74, issue 17, pg.
15 3654-3662, 2011.
- 16 [55] Zhang, L., Zhou, W. and Chang, P.-C., Generalized nonlinear discriminant analysis and its
17 small sample size problems, *Neurocomputing*, vol. 74, issue 4, pp. 568-574, 2011.
- 18 [56] Huang, H., Liu, J., Feng, H. and He, T., Ear recognition based on uncorrelated local Fisher
19 discriminant analysis, *Neurocomputing*, vol. 74, issue 17, pp. 3103-3113, 2011.
- 20 [57] Huerta, E.B., Duval, B. and Hao, J.-K., A hybrid LDA and genetic algorithm for gene
21 selection and classification of microarray data, *Neurocomputing*, vol. 73, issue 13-15, pp.
22 2375-2383, 2010.
- 23 [58] Sharma, A., Imoto, S., Miyano, S., A top-r feature selection algorithm for microarray gene
24 expression data, *IEEE/ACM Trans. on Comp. Biol. and Bioinformatics*, vol. 9, issue 3, pp.
25 754-764, 2012.
- 26 [59] Sharma, A., Imoto, S., Miyano, S., A between-class overlapping filter-based method for
27 transcriptome data analysis, *Journal of Bioinformatics and Computational Biology*, vol. 10,
28 no. 5, pp. 1250010-1 1250010-20, 2012.
- 29 [60] Sharma, A., Imoto, S., Miyano, S., Sharma, V., Null space based feature selection method
30 for gene expression data, *Int. Jnr. of Mach. Learn. and Cybernetics*, vol. 3, issue 4, pp.
31 269-276, 2012.
- 32 [61] Sharma, A. and Paliwal, K.K., Cancer classification by gradient LDA technique using
33 microarray gene expression data, *Data & Knowledge Engineering*, vol. 66, issue 2, pp
34 338-347, 2008.
- 35 [62] Yang, W and Wu, H., Regularized complete linear discriminant analysis, *Neurocomputing*,
36 vol. 137, pp. 185-191, 2014.