Rapid and Brief Communication

# A direct LDA algorithm for high-dimensional data — with application to face recognition

## Hua Yu*, Jie Yang

*Interactive System Labs, School of Computer Science, Carnegie Mellon University, 5000 Forbes Avenue, Pittsburgh, PA 15213, USA*

## 1. Introduction

Linear discriminant analysis (LDA) has been successfully used as a dimensionality reduction technique to many classification problems, such as speech recognition, face recognition, and multimedia information retrieval. The objective is to find a projection $A$ that maximizes the ratio of between-class scatter $S_b$ against within-class scatter $S_w$ (Fisher's criterion):

$$\arg \max_A \frac{|A S_b A^{\mathrm{T}}|}{|A S_w A^{\mathrm{T}}|}.$$

However, for a task with very high-dimensional data such as images, the traditional LDA algorithm encounters several difficulties. Consider face recognition for example. A low-definition face image of size $64 \times 64$ implies a feature space of $64 \times 64 = 4096$ dimensions, and therefore scatter matrices of size $4096 \times 4096 = 16M$. First, it is computationally challenging to handle big matrices (such as computing eigenvalues). Second, those matrices are almost always singular, as the number of training images needs to be at least 16M for them to be non-degenerate.

Due to these difficulties, it is commonly believed that a direct LDA solution for such high-dimensional data is infeasible. Thus, ironically, before LDA can be used to reduce dimensionality, another procedure has to be first applied for dimensionality reduction.

In face recognition, many techniques have been proposed (for a good review, see Ref. [1]). Among them, the most notable is a *two-stage* PCA + LDA approach [2,3]:

$$A = A_{\mathrm{LDA}} A_{\mathrm{PCA}}.$$

Principal component analysis (PCA) is used to project images from the original *image space* into a *face-subspace*, where dimensionality is reduced and $S_w$ is no longer degenerate, so that LDA can proceed without trouble. A potential problem is that the PCA criterion may not be compatible with the LDA criterion, thus the PCA step may discard dimensions that contain important discriminative information.

Chen et al. have recently proved that the null space of $S_w$ contains the most discriminative information [1]. But, their approach fell short of making use of any information outside of that null space. In addition, heuristics are needed to extract a small number of features for image representation, so as to avoid computational problems associated with large scatter matrices.

In this paper, we present a direct, exact LDA algorithm for high-dimensional data set. It accepts high-dimensional data (such as raw images) an input, and optimizes Fisher's criterion directly, without any feature extraction or dimensionality reduction steps.

## 2. Direct LDA solution

At the core of the direct LDA algorithm lies the idea of simultaneous diagonalization, the same as in the traditional LDA algorithm. As the name suggests, it tries

---

*Corresponding author. Tel.: + 1-412-268-5479; fax: + 1-412-268-6298.

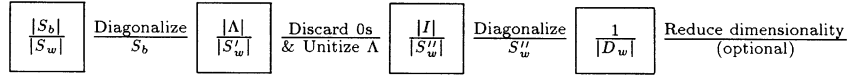*E-mail address:* hyu@cs.cmu.edu (H. Yu).

Fig. 1. Thumbnail of the direct LDA algorithm.

to find a matrix that simultaneously diagonalizes both $S_w$ and $S_b$:

$$AS_wA^T = I, \quad AS_bA^T = \Lambda,$$

where $\Lambda$ is a diagonal matrix with diagonal elements sorted in decreasing order. To reduce dimensionality to $m$, we simply pick the top $m$ rows of $A$, which corresponds to the largest $m$ diagonal elements in $\Lambda$. Details of the algorithm can be found in Ref. [4].

The key idea of our new algorithm is to discard the null space of $S_b$ — which contains no useful information — rather than discarding the null space of $S_w$, which contains the most discriminative information. This can be achieved by diagonalizing $S_b$ first and then diagonalizing $S_w$. The traditional procedure takes the reverse order. While both approaches produce the same result when $S_w$ is not singular, the reversal in order makes a drastic difference for high-dimensional data, where $S_w$ is likely to be singular.

The new algorithm is outlined below. Fig. 1 provides a conceptual overview of this algorithm. Computational issues will be discussed shortly after.

(1) Diagonalize $S_b$: find matrix $V$ such that

$$V^TS_bV = \Lambda, \quad \text{eigen decomposition / SVD}$$

where $V^TV = I$. $\Lambda$ is a diagonal matrix sorted in decreasing order.

This can be done using the traditional eigenanalysis, i.e. each column of $V$ is an eigenvector of $S_b$, and $\Lambda$ contains all the eigenvalues. As $S_b$ might be singular, some of the eigenvalues will be 0 (or close to 0). It is necessary to discard those eigenvalues and eigenvectors, as projection directions with a total scatter of 0 do not carry any discriminative power at all. at most c - 1 non-zero

Let $Y$ be the first $m$ columns of $V$ (an $n \times m$ matrix, $n$ being the feature space dimensionality), now

$$Y^TS_bY = D_b > 0,$$

where $D_b$ is the $m \times m$ principal sub-matrix of $\Lambda$.

(2) Let $Z = YD_b^{-1/2}$,

$$(YD_b^{-1/2})^T S_b(YD_b^{-1/2}) = I \Rightarrow Z^T S_bZ = I.$$

Thus, $Z$ unitizes $S_b$, and reduces dimensionality from $n$ to $m$. m-by-n

Diagonalize $Z^TS_wZ$ by eigenanalysis:

$$U^TZ^TS_wZU = D_w, \quad \text{Dw is also m-by-m}$$

where $U^TU = I$. $D_w$ may contain zeros in its diagonal. we can simply keep all the m and do this finally.

Since the objective is to maximize the ratio of total-scatter against within-class scatter, we can sort the diagonal elements of $D_w$ and discard some eigenvalues in the high end, together with the corresponding eigenvectors. It is important not to keep the dimensions with the smallest eigenvalues, especially zeros. This is exactly the reason why we started by diagonalizing $S_b$, rather than $S_w$. See Section 2.2 for more discussion.

(3) Let the LDA matrix

$$A = U^TZ^T.$$

$A$ diagonalizes both the numerator and the denominator in Fisher's criterion

$$AS_wA^T = D_w, \quad AS_bA^T = I.$$

(4) For classification purpose, notice that $A$ already diagonalizes $S_w$; therefore the final transformation that spheres the data should be

$$x^* \leftarrow D_w^{-1/2} Ax.$$

### 2.1. Computational considerations

Although the scheme above gives an exact solution for Fisher's criterion, we have not addressed the computational difficulty that both scatter matrices are too big to be held in memory, let alone their eigenanalysis.

Fortunately, the method presented by Turk and Pentland [5] for the eigenface problem is still applicable. The key observation is that scatter matrices can be represented in a way that both saves memory, and facilitates eigenanalysis. For example,

$$S_b = \sum_{i=1}^{J} n_i(\mu_i - \mu)(\mu_i - \mu)^T = \Phi_b\Phi_b^T \quad (n \times n),$$

where

$$\Phi_b = [\sqrt{n_1}(\mu_1 - \mu), \sqrt{n_2}(\mu_2 - \mu), \dots] \quad (n \times J)$$

with $J$ being the number of classes and $n_i$ the number of training images for class $i$. Thus, instead of storing an $n \times n$ matrix, we need only to store $\Phi_b$ which is $n \times J$. The eigenanalysis is simplified by virtue of the following lemma:

**Lemma 1.** *For any $n \times m$ matrix $L$, mapping $x \to Lx$ is a one-to-one mapping that maps eigenvectors of $L^TL$ ($m \times m$) onto those of $LL^T$ ($n \times n$).*

As $\Phi_b^{\mathrm{T}}\Phi_b$ is a $J \times J$ matrix, eigenanalysis is affordable. In Step 2 of our algorithm, to compute eigenvalues for $Z^{\mathrm{T}}S_wZ$, simply notice

$$S_w = \sum_i (x_i - \mu_{k_i})(x_i - \mu_{k_i})^{\mathrm{T}} = \Phi_w\Phi_w^{\mathrm{T}},$$

where

$$\Phi_w = [x_1 - \mu_{k_1}, x_2 - \mu_{k_2}, \dots ] \quad (n \times n_t),$$

with $n_t$ being the total number of images in the training set. Thus

$$Z^{\mathrm{T}}S_wZ = Z^{\mathrm{T}}\Phi_w\Phi_w^{\mathrm{T}}Z = (\Phi_w^{\mathrm{T}}Z)^{\mathrm{T}}\Phi_w^{\mathrm{T}}Z.$$

We can again use Lemma 1 to compute eigenvalues.

### 2.2. Discussions

#### 2.2.1. Null space of $S_w$

The traditional simultaneous diagonalization begins by diagonalizing $S_w$. If $S_w$ is not degenerate, it gives the same result as our approach. If $S_w$ is singular, however, the traditional approach runs into a dilemma: to proceed, it has to discard those eigenvalues equal to 0; but those discarded eigenvectors are the most important dimensions!

As Chen et al. pointed out [1], the null space of $S_w$[1] carries most of the discriminative information. More precisely, for a projection direction $a$, if $S_wa = 0$, and $S_ba \neq 0$, $aS_ba^{\mathrm{T}}/aS_wa^{\mathrm{T}}$ is maximized. The intuitive explanation is that, when projected onto direction $a$, within-class scatter is 0 but between-class scatter is not. Obviously, perfect classification can be achieved in this direction.

Different from the algorithm proposed in Ref. [1], which operates solely in the null space, our algorithm can take advantage of all the information, both within and outside of $S_w$'s null space. Our algorithm can still be used in cases where $S_w$ is not singular, which is common in tasks like speech recognition.

#### 2.2.2. Equivalence to PCA + LDA

As Fukunaga pointed out [4], there are other variants of Fisher's criterion

$$\arg\max_A \frac{|A^{\mathrm{T}}S_tA|}{|A^{\mathrm{T}}S_wA|} \quad \text{or} \quad \arg\max_A \frac{|A^{\mathrm{T}}S_bA|}{|A^{\mathrm{T}}S_tA|},$$

where $S_t = S_b + S_w$ is the *total scatter matrix*.

Interestingly, if we use the first variant (with $S_t$ in the numerator), Step 1 of our algorithm becomes exactly PCA. Discarding $S_t$'s eigenvectors with 0 eigenvalues reduces dimensionality, just as Belhumeur et al. proposed in their two-stage PCA + LDA method [3]. If their LDA step handled $S_w$'s null space properly, the two approaches would give the same performance. In a sense our method can be called "unified PCA + LDA", since there is no separate PCA step. It not only leads to a clean presentation, but also results in an efficient implementation.

### 3. Face recognition experiments

We tested the direct LDA algorithm on face images from Olivetti-Oracle Research Lab (ORL, http://www.cam-orl.co.uk). The ORL data set consists of 400 frontal faces: 10 tightly, cropped images of 40 individuals with variations in pose, illumination, facial expression (open/closed eyes, smiling /not smiling) and facial details (glasses/no glasses). The size of each image is $92 \times 112$ pixels, with 256 grey levels per pixel.

Three sets of experiments are conducted. In all cases we randomly choose five images per person for training, the other five for testing. To reduce variation, each experiment is repeated at least 10 times.

Without dimensionality reduction in Step 2, average recognition accuracy is 90.8%. With dimensionality reduction, where everything outside of $S_w$'s null space is discarded, average recognition accuracy becomes 86.6%. This verifies that while $S_w$'s null space is important, discriminative information does exist outside of it.

### 4. Conclusions

In this paper, we proposed a direct LDA algorithm for high-dimensional data classification, with application to face recognition in particular. Since the number of samples is typically smaller than the dimensionality of the samples, both $S_b$ and $S_w$ are singular. By modifying the simultaneous diagonalization procedure, we are able to discard the null space of $S_b$ — which carries no discriminative information — and to keep the null space of $S_w$, which is very important for classification. In addition, computational techniques are introduced to handle large scatter matrices efficiently. The result is a unified LDA algorithm that gives an exact solution to Fisher's criterion whether or not $S_w$ is singular.

### References

[1] L. Chen, H. Liao, M. Ko, J. Lin, G. Yu, A new LDA-based face recognition system which can solve the small sample size problem, Pattern Recognition 33 (10) (2000) 1713–1726.

---

[1] Null space of $S_w = \{x \mid S_wx = 0, x \in \mathbf{R}^n\}$.

[2] D. Swets, J. Weng, Using discriminant eigenfeatures for image retrieval, Pattern Anal. Mach. Intell. 18 (8) (1996) 831–836.

[3] P.N. Belhumeur, J.P. Hespanha, D.J. Kriegman, Eigenfaces vs. fisherface: recognition using class specific linear projection, Pattern Anal. Mach. Intell. 19 (7) (1997) 711–720.

[4] K. Fukunaga, Introduction to Statistical Pattern Recogniton (2nd Edition), Academic Press, New York, 1990.

[5] M. Turk, A. Pentland, Eigenfaces for recognition, J. Cognitive Neurosci. 3 (1) (1991) 72–86.