



2D-LDA: A statistical linear discriminant analysis for image matrix

Ming Li ^{*}, Baozong Yuan

Institute of Information Science, Beijing Jiaotong University, Beijing 100044, China

Received 19 August 2004

Available online 18 October 2004

Abstract

This paper proposes an innovative algorithm named 2D-LDA, which directly extracts the proper features from image matrices based on Fisher's Linear Discriminant Analysis. We experimentally compare 2D-LDA to other feature extraction methods, such as 2D-PCA, Eigenfaces and Fisherfaces. And 2D-LDA achieves the best performance. © 2004 Elsevier B.V. All rights reserved.

Keywords: Feature extraction; Image representation; Linear discriminant analysis; Subspace techniques; Face recognition

1. Introduction

Feature extraction is the key to face recognition, as it is to any pattern classification task. The aim of feature extraction is reducing the dimensionality of face image so that the extracted features are as representative as possible. The class of image analysis methods called appearance-based approach has been of wide concern, which relies on statistical analysis and machine learning. Turk and Pentland (1991) presented the well-known Eigenfaces method for face recognition, which uses principal component analysis (PCA)

for dimensionality reduction. However, the base physical similarity of the represented images to originals does not provide the best measure of useful information for distinguishing faces from one another (O'Toole, 1993). Belhumeur et al. (1997) proposed Fisherfaces method, which is based on Fisher's Linear Discriminant and produces well separated classes in a low-dimensional subspace. His method is insensitive to large variation in lighting direction and facial expression.

Recently, Yang (2002) investigated the Kernel PCA for learning low dimensional representations for face recognition and found that the Kernel methods provide better representations and achieve lower error rates for face recognition. Bartlett et al. (2002) proposed using ICA for face representation, which is sensitive to the high-order

^{*} Corresponding author. Tel.: +86 10 5168 3149; fax: +51 68 6168 8616.

E-mail address: liming@mail.edu.cn (M. Li).

statistics. This method is superior to representations based on PCA for recognizing faces across days and changes in expression. However, Kernel PCA and ICA are both computationally more expensive than PCA. Weng et al. (2003) presented a fast method, called candid covariance-free IPCA (CCIPCA), to obtain the principal components of high-dimensional image vectors. Moghaddan (2002) compared the Bayesian subspace method with several PCA-related methods (PCA, ICA, and Kernel PCA). The experimental results demonstrated its superiority over PCA, ICA and Kernel PCA.

All the PCA-related methods discussed above are based on the analysis of vectors. When dealing with images, we should firstly transform the image matrixes into image vectors. Then based on these vectors the covariance matrix is calculated and the optimal projection is obtained. However, face images are high-dimensional patterns. For example, an image of 112×92 will form a 10304-dimensional vector. It is difficult to evaluate the covariance matrix in such a high-dimensional vector space. To overcome the drawback, Yang proposed a straightforward image projection technique, named as *image principal component analysis* (IMPCA) (Yang et al., 2004), which is directly based on analysis of original image matrices. Different to traditional PCA, 2D-PCA is based on 2D matrices rather than 1D vectors. This means that the image matrix does not need to be converted into a vector. As a result, 2D-PCA has two advantages: easier to evaluate the covariance matrix accurately and lower time-consuming. Liu et al. (1993) proposed an iterative method to calculate the Foley-Sammon optimal discriminant vectors from image matrixes. And he proposed to substitute $D_t = D_b + D_w$ for D_w to overcome the singularity problem. Liu's method was complicate and didn't resolve the singularity problem well.

In this paper, a statistical linear discriminant analysis for image matrix is discussed. Our method proposes to use Fisher linear projection criterion to find out a good projection. This criterion is based on two parameters: the *between-class scatter matrix* and the *within-class scatter matrix*. Because the dimension of between-class and within-class scatter matrix is much low (compara-

tive to number of training samples). So, the problem, that the within-class scatter matrix maybe singular, will be handled. At the same time, the compute-costing is lower than traditional Fisherfaces. Moreover, we discuss about image reconstruction and conduct a series of experiments on the ORL face database.

The organization of this paper is as follows: In Section 2, we propose the idea and describe the algorithm in detail. In Section 3, we compare 2D-LDA with Eigenfaces, Fisherfaces and 2D-PCA on the ORL face database. Finally, the paper concludes with some discussions in Section 4.

2. Two-dimensional linear discriminant analysis

2.1. Principle: The construction of Fisher projection axis

Let A denotes a $m \times n$ image, and x denotes an n -dimensional column vector. A is projected onto x by the following linear transformation

$$y = Ax. \quad (1)$$

Thus, we get an m -dimensional projected vector y , which is called the feature vector of the image A .

Suppose there are L known pattern classes in the training set, and M denotes the size of the training set. The j th training image is denoted by an $m \times n$ matrix A_j ($j = 1, 2, \dots, M$), and the mean image of all training sample is denoted by \bar{A} and \bar{A}_i ($i = 1, 2, \dots, L$) denoted the mean image of class T_i and N_i is the number of samples in class T_i , the projected class is P_i . After the projection of training image onto x , we get the *projected feature vector*

$$y_j = A_j x, \quad j = 1, 2, \dots, M. \quad (2)$$

How do we judge a projection vector x is good? In fact, the total scatter of the projected samples can be characterized by the trace of the covariance matrix of the projected feature vectors (Turk and Pentland, 1991). From this point of view, we introduced a criterion at first,

$$J(x) = \frac{P_B}{P_W}. \quad (3)$$

There were two parameters

$$P_B = \text{tr}(\mathbf{TS}_B), \quad (4)$$

$$P_W = \text{tr}(\mathbf{TS}_W), \quad (5)$$

where \mathbf{TS}_B denotes the between-class scatter matrix of projected feature vectors of training images, and \mathbf{TS}_W denotes the within-class scatter matrix of projected feature vectors of training images. So,

$$\begin{aligned} \mathbf{TS}_B &= \sum_{i=1}^L N_i (\bar{\mathbf{y}}_i - \bar{\mathbf{y}})(\bar{\mathbf{y}}_i - \bar{\mathbf{y}})^T \\ &= \sum_{i=1}^L N_i [(\bar{\mathbf{A}}_i - \bar{\mathbf{A}})\mathbf{x}][(\bar{\mathbf{A}}_i - \bar{\mathbf{A}})\mathbf{x}]^T, \end{aligned} \quad (6)$$

$$\begin{aligned} \mathbf{TS}_W &= \sum_{i=1}^L \sum_{\mathbf{y}_k \in P_i} (\mathbf{y}_k - \bar{\mathbf{y}}_i)(\mathbf{y}_k - \bar{\mathbf{y}}_i)^T \\ &= \sum_{i=1}^L \sum_{\mathbf{y}_k \in P_i} [(\mathbf{A}_k - \bar{\mathbf{A}}_i)\mathbf{x}][(\mathbf{A}_k - \bar{\mathbf{A}}_i)\mathbf{x}]^T. \end{aligned} \quad (7)$$

So

$$\begin{aligned} \text{tr}(\mathbf{TS}_B) &= \mathbf{x}^T \left(\sum_{i=1}^L N_i (\bar{\mathbf{A}}_i - \bar{\mathbf{A}})^T (\bar{\mathbf{A}}_i - \bar{\mathbf{A}}) \right) \mathbf{x} \\ &= \mathbf{x}^T \mathbf{S}_B \mathbf{x}, \end{aligned} \quad (8)$$

$$\begin{aligned} \text{tr}(\mathbf{TS}_W) &= \mathbf{x}^T \left(\sum_{i=1}^L \sum_{\mathbf{A}_k \in T_i} (\mathbf{A}_k - \bar{\mathbf{A}}_i)^T (\mathbf{A}_k - \bar{\mathbf{A}}_i) \right) \mathbf{x} \\ &= \mathbf{x}^T \mathbf{S}_W \mathbf{x}. \end{aligned} \quad (9)$$

We could evaluate \mathbf{TS}_B and \mathbf{TS}_W directly using the training image samples.

So, the criterion could be expressed by

$$J(\mathbf{x}) = \frac{\mathbf{x}^T \mathbf{S}_B \mathbf{x}}{\mathbf{x}^T \mathbf{S}_W \mathbf{x}}, \quad (10)$$

where \mathbf{x} is a unitary column vector. This criterion is called *Fisher linear projection criterion*. The unitary vector \mathbf{x} that maximizes $J(\mathbf{x})$ is called the *Fisher optimal projection axis*. The optimal projection \mathbf{x}_{opt} is chosen when the criterion is maximized, i.e.,

$$\mathbf{x}_{\text{opt}} = \arg \max_{\mathbf{x}} J(\mathbf{x}). \quad (11)$$

If \mathbf{S}_W is nonsingular, the solution to above optimization problem is to solve the generalized eigenvalue problem (Turk and Pentland, 1991):

$$\mathbf{S}_B \mathbf{x}_{\text{opt}} = \lambda \mathbf{S}_W \mathbf{x}_{\text{opt}}. \quad (12)$$

In the above equation, λ is the maximal eigenvalue of $\mathbf{S}_W^{-1} \mathbf{S}_B$.

The traditional LDA must face to the singularity problem. However, 2D-LDA overcomes this problem successfully. This is because: For each training image, \mathbf{A}_j ($j = 1, 2, \dots, M$), we have $\text{rank}(\mathbf{A}_j) = \min(m, n)$.

From (9), we have

$$\begin{aligned} \text{rank}(\mathbf{S}_W) &= \text{rank} \left(\sum_{i=1}^L \sum_{\mathbf{A}_k \in T_i} (\mathbf{A}_k - \bar{\mathbf{A}}_i)^T (\mathbf{A}_k - \bar{\mathbf{A}}_i) \right) \\ &\leq (M - L) \cdot \min(m, n). \end{aligned} \quad (13)$$

So, in 2D-LDA, \mathbf{S}_W is nonsingular when

$$M \geq L + \frac{n}{\min(m, n)}. \quad (14)$$

In real situation, (14) is always satisfied. So, \mathbf{S}_W is always nonsingular.

In general, it is not enough to have only one Fisher optimal projection axis. We usually need to select a set of projection axis, $\mathbf{x}_1, \dots, \mathbf{x}_d$, subject to the orthonormal constraints. That is,

$$\begin{cases} \{\mathbf{x}_1, \dots, \mathbf{x}_d\} = \arg \max J(\mathbf{x}) \\ \mathbf{x}_i^T \mathbf{x}_j = 0, \quad i \neq j, \quad i, j = 1, \dots, d. \end{cases} \quad (15)$$

In fact, the optimal projection axes, $\mathbf{x}_1, \dots, \mathbf{x}_d$, are the orthonormal eigenvectors of $\mathbf{S}_W^{-1} \mathbf{S}_B$ corresponding to the first d largest eigenvalues.

Using these projection axes, we could form a new *Fisher projection matrix* \mathbf{X} , which is a $n \times d$ matrix,

$$\mathbf{X} = [\mathbf{x}_1 \quad \mathbf{x}_2 \quad \dots \quad \mathbf{x}_d]. \quad (16)$$

2.2. Feature extraction

We will use the optimal projection vectors of 2D-LDA, $\mathbf{x}_1, \dots, \mathbf{x}_d$, for feature extraction. For a given image \mathbf{A} , we have

$$\mathbf{y}_k = \mathbf{A} \mathbf{x}_k, \quad k = 1, 2, \dots, d. \quad (17)$$

Then, we have a family of *Fisher feature vectors* $\mathbf{y}_1, \dots, \mathbf{y}_d$, which formed a $m \times d$ matrix $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_d]$. We called this matrix \mathbf{Y} as the *Fisher feature matrix* of the image \mathbf{A} .

2.3. Reconstruction

In the 2D-LDA method, we can use the Fisher feature matrixes and Fisher optimal projection axes to reconstruct a image by following steps.

For a given image A , the Fisher feature matrix is $Y = [y_1, \dots, y_d]$ and the Fisher optimal projection axes $X = [x_1, \dots, x_d]$, then we have

$$Y = AX. \quad (18)$$

Because x_1, \dots, x_d are orthonormal, it is easy to obtain the reconstructed image of A :

$$\tilde{A} = YX^T = \sum_{k=1}^d y_k x_k^T. \quad (19)$$

We called $\tilde{A}_k = y_k x_k^T$ as a *reconstructed subimage* of A , which have the same size as image A . This means that we use a set of 2D Fisherfaces to reconstruct the original image. If we select $d = n$, then we can completely reconstruct the images in the training set: $\tilde{A} = A$. If $d < n$, the reconstructed image \tilde{A} is an approximation for A .

2.4. Classification

Given two images A_1, A_2 represented by 2D-LDA feature matrix $Y_1 = [y_1^1, \dots, y_d^1]$ and $Y_2 = [y_1^2, \dots, y_d^2]$. So the similarity $d(Y_1, Y_2)$ is defined as

$$d(Y_1, Y_2) = \sum_{k=1}^d \|y_k^1 - y_k^2\|_2, \quad (20)$$

where $\|y_k^1 - y_k^2\|_2$ denotes the Euclidean distance between the two Fisher feature vectors y_k^1 and y_k^2 .

If the Fisher feature matrix of training images are Y_1, Y_2, \dots, Y_M (M is the total number of training images), and each image is assigned to a class T_i . Then, for a given test image Y , if $d(Y, Y_i) =$

$\min_j d(Y_1, Y_j)$ and $Y_i \in T_i$, then the resulting decision is $Y \in T_i$.

3. Experiment and analysis

We evaluated our 2D-LDA algorithm on the ORL face image database. The ORL database (<http://www.cam-orl.co.uk>) contains images of 40 individuals, each person have 10 different images. For some individuals, the images were taken at different times. The facial expression (open or closed eyes, smiling or nonsmiling) and facial details (glasses or no glasses) also vary. All the images were taken against a dark homogeneous background with the subjects in an upright, frontal position (with tolerance for some side movement). The images were taken with a tolerance for some titling and rotation of the face of up to 20° . Moreover, there is also some variation in the scale of up to about 10%. The size of each image is 92×112 pixels, with 256 grey levels per pixel. Five samples of one person in ORL database are shown in Fig. 1. So, we could use the ORL database to evaluate 2D-LDA's performance under conditions where pose and the size of sample are varied.

Using 2D-LDA, we could project the test face image onto the Fisher optimal projection axis, then we could use the Fisher feature vectors set to reconstruct the image. In Fig. 2, some reconstructed images and the original image of one person were given out. In Fig. 2, the variance d denotes the number of dimension used to map and reconstruct the face image. As observed these images, we could find that the reconstructed images are very like obtained by sample the original image on the spacing vertical scanning line. The reconstructed image \tilde{A} is more and more like to the original image A as the value of d increased.

We have done an experiment on the ORL database to evaluate performance of 2D-LDA, 2D-



Fig. 1. Five images in ORL database.

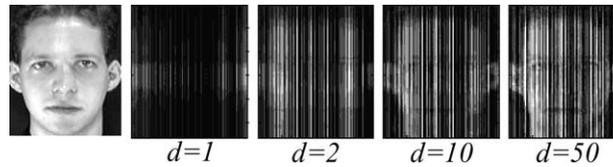


Fig. 2. Some reconstructed images of one person.

PCA (Yang et al., 2004), Eigenfaces (Turk and Pentland, 1991), Fisherfaces (Belhumeur et al., 1997). To evaluate the pure ability of these four in the fair environment, we did not do any preprocess on the face images, and we did not utilize any optimized algorithm. We just realized the algorithm appeared in the literature (Turk and Pentland, 1991; Belhumeur et al., 1997 and Yang et al., 2004) without any modification. In our experiment, we select first five images samples per person for training, and the left five images samples for testing. So, in our experiment, the size of training set and testing set were both 200. So in the 2D-LDA, the size of between-class scatter matrix S_B and within-class scatter matrix S_W are both 92×92 . Fig. 3 shows the classification result. From Fig. 3, we find that the recognition rate of 2D-LDA have achieved the best performance in the four methods. And the best result of 2D-LDA 94.0% is much better than the best result of 2D-PCA 92.5%. And from Fig. 3, we could find that the two 2D feature extraction methods have outstanding performance in the low-dimension

condition, but the conventional ones' ability is very poor.

Table 1 showed out the comparison of the training time of the four algorithms (CPU: Pentium IV 2.66 GHz, RAM: 256M). The four algorithms are realized in the Matlab environment. We could see that the 2D-LDA and 2D-PCA's computing-cost is very low compared with Eigenfaces and Fisherfaces. This is because in the 2D condition, we only need to handle a 92×92 matrix. But using the Eigenfaces and Fisherfaces, we must face to a 10304×10304 matrix. It is a hard work. At last, It must be mentioned that when used the Fisherfaces, we must reduced the dimension of image data to avoid that S_W is singular (Belhumeur et al., 1997). Mapping the original data onto how many dimensions space is a hard problem. We must select the proper number of dimension through experiment. Considering this situation, Fisherfaces is very time-costing.

Table 2 showed that the memory cost of 2D feature extraction is much larger than the 1D ones. This is because 2D methods used a $n \times d$ matrix to present a face image. At the same time the 1D techniques reconstructed face images by a d -dimension vector.

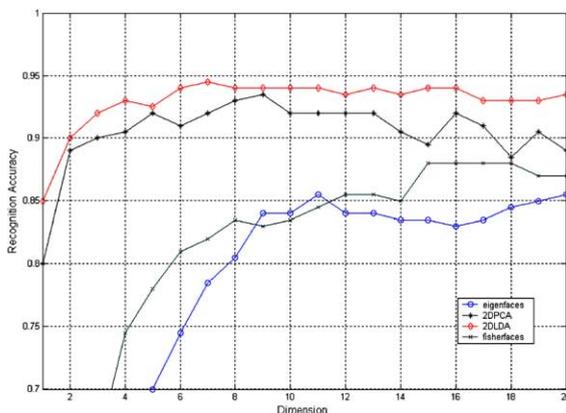


Fig. 3. Comparison of 2D-LDA and 2D-PCA on ORL Database.

Table 1
Comparison of CPU Time (s) for feature extraction using ORL database (15 dimensions)

2D-LDA	2D-PCA	Eigenfaces	Fisherfaces
0.4210	0.4210	28.5000	32.5310

Table 2
Comparison of memory cost (bytes) to present a 92×112 image using different techniques (15 dimensions)

2D-LDA	2D-PCA	Eigenfaces	Fisherfaces
6720	6720	60	60

4. Conclusion

In this paper, a new algorithm for image feature extraction and selection was proposed. This method uses the Fisher Linear Discriminant Analysis to enhance the effect of variation caused by different individuals, other than by illumination, expression, orientation, etc. 2D-LDA uses the image matrix instead of the image vector to compute the between-class scatter matrix and the within-class scatter matrix.

From our experiments, we can see that the 2D-LDA have many advantages over other methods. 2D-LDA achieves the best recognition accuracy in the four algorithms. And this technique's computing cost is very low compared with Eigenfaces and Fisherfaces, and close to 2D-PCA. At the same time, this method shows powerful performance in the low dimension. From Fig. 2, we can see that this new projection method is very like to select spacing vertical scanning lines to present a image. Maybe this is the reason that this algorithm is so effective in image classification.

2D-LDA still has its shortcoming. It needs more memory to store a image than the Eigenfaces and Fisherfaces.

Acknowledgement

This work was supported by the National Natural Science Foundation of China (No. 60441002)

and the University Key Research Project (No. 2003SZ002).

References

- Bartlett, M.S., Movellan, J.R., Sejnowski, T.J., 2002. Face recognition by independent component analysis. *IEEE Trans. Neural Networks* 13 (6), 1450–1464.
- Belhumeur, P.N., Hespanha, J.P., Kriegman, D.J., 1997. Eigenfaces vs. Fisherfaces: Recognition using class specific linear projection. *IEEE Trans. Pattern Anal. Machine Intell.* 19 (7), 711–720.
- Lui, K., Cheng, Y., Yang, J., 1993. Algebraic feature extraction for image recognition based on an optimal discriminant criterion. *Pattern Recognition* 26 (6), 903–911.
- Moghaddan, B., 2002. Principal manifolds and probabilistic subspaces for visual recognition. *IEEE Trans. Pattern Anal. Machine Intell.* 24 (6), 780–788.
- O'Toole, A., 1993. Low-dimensional representation of faces in higher dimensions of the face space. *J. Opt. Soc. Amer.* 10 (3).
- Turk, M., Pentland, A., 1991. Eigenfaces for recognition. *J. Cognitive Neurosci.* 3 (1), 71–86.
- Weng, J., Zhang, Y., Hwang, W., 2003. Candid covariance-free incremental principal component analysis. *IEEE Trans. Pattern Anal. Machine Intell.* 25 (8), 1034–1040.
- Yang, M.H., 2002. Kernel eigenfaces vs. Kernel Fisherfaces: Face recognition using Kernel methods. In: *Proc. 5th Internat. Conf. on Automatic Face and Gesture Recognition (RGR'02)*, pp. 215–220.
- Yang, J., Zhang, D., Frangi, A.F., Yang, J.-y., 2004. Two-dimensional PCA: A new approach to appearance-based face representation and recognition. *IEEE Trans. Pattern Anal. Machine Intell.* 26 (1), 131–137.