

Contextual Constraints based Linear Discriminant Analysis

Zhen Lei*, Stan Li

Center for Biometrics and Security Research & National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences, 95 Zhongguancun Donglu, Beijing 100190, China.

Abstract

Linear feature extraction methods such as LDA have achieved great success in pattern recognition and image processing area. For most existing methods, the image data is usually transformed into a vector representation and the contextual information among pixels is not exploited. However, image data distribute sparsely in high-dimension feature space and the dependence among neighboring pixels is important to represent a natural image. Therefore, in this paper, we propose a novel image contextual constraint based linear discriminant analysis (CCLDA) method by taking into account the pixel dependence of an image in subspace learning process. In this way, a more discriminative subspace could be learned especially in the case of small sample size. Extensive experiments on ORL, Extended Yale-B, PIE and FRGC databases validate the efficacy of the proposed method.

Keywords: Contextual constraint, Discriminant analysis, Face recognition

*Corresponding author. Tel: +86-10-62658306. Fax: +86-10-62632259
Email address: zhen.ray@gmail.com (Zhen Lei)

1. Introduction

Subspace learning, due to its efficacy and efficiency, has achieved great success in pattern recognition and computer vision area, such as face recognition and image retrieval etc. Linear Discriminant Analysis (LDA) (Belhumeur et al. (1997)) is one of the representative methods. The objective of LDA is to find a subspace that maximizes the sample distance from different classes and meanwhile minimizes the sample distance from the same class. Thus, the derived subspace is discriminative to classify different samples correctly. However, in real applications, due to the high dimension of feature (e.g. for one 100×100 image, the number of pixels is 10,000) and usually small number of samples, the classical LDA always fails; this is called small sample size (SSS) problem. To address this problem, a lot of work has been investigated and many LDA-variants have been proposed such as Null-LDA (Chen et al. (2000)), Direct LDA (Yu and Yang (2001)), Regularized LDA (Friedman (1989); Lu et al. (2004)) etc..

All the above LDA methods are based on vector representation. That is, before LDA, the images are transformed into vectors by a certain order. Recently, 2D-LDA (Ye et al. (2004)), which takes operations on image matrix directly was proposed. The most advantage of 2D-LDA is its low computational cost and it implicitly avoids small sample size (SSS) problem since the within class scatter is usually non-singular. Moreover, the 2D-LDA is claimed to retain the image spatial structure so that it is more appropriate than 1D-LDA to address image related problem.

Most of the existing subspace learning methods like PCA (Turk and Pentland (1991)), LDA, LLE (Roweis and Saul (2000)), LPP (He et al. (2005b)),



(a)



(b)

Figure 1: A face image (left) and its re-shuffled one (right). For 1D-LDA classification result, there is no difference between these two forms. However, people almost cannot get any information except some noise from the right one.

26 NPE (He et al. (2005a)), LPDP (Gui et al. (2010)), consider the pixels in im-
27 age independently, not taking into account their spatial relationship. Fig. 1
28 shows an original face image and its relative in which the pixels are re-
29 shuffled. Obviously, we cannot recognize any face information from the right
30 image. However, in vector based discriminant learning methods, if the pix-
31 els in every image are re-shuffled in the same way, it would result in the
32 same classification performance even though there is no semantic sense in
33 shuffled images any more. It is well known that images with certain pattern
34 occupy specific manifold, which is constrained by contextual information, in
35 high-dimension feature space. Therefore, contextual constraint in image is
36 important for image understanding and will provide useful information for
37 classification. One of the most successful work to model the contextual in-
38 formation is the Markov Random Fields (MRFs) (Huang et al. (2004); Dass
39 and Jain (2001); Dass et al. (2002)) which derive the results by maximizing
40 the posterior probability in Bayesian deduction framework. However, the

41 optimization by MRF is somewhat computational expensive and is easy to
42 converge into local minima that limits its application.

43 Wang et al. (Wang et al. (2005)) proposed a novel image matching dis-
44 tance considering the spatial information. However, they didn't demonstrate
45 how to integrate the contextual information into dimensionality reduction
46 problem. In this paper, we propose contextual constraints based linear dis-
47 criminant analysis (CCLDA) which incorporates the contextual information
48 into linear discriminant analysis. The main difference of proposed method
49 from the existing ones is that it takes into account the image contextual
50 constraints during feature dimensionality reduction. Moreover, we study the
51 relationship of CCLDA with LDA and 2D-LDA and indicate that LDA and
52 2D-LDA are essentially certain versions of CCLDA. It should be noted the
53 proposed contextual constraints consideration could be incorporated not only
54 into LDA, but also into other general subspace learning methods.

55 The remainder of this paper is organized as follows. Section 2 briefly
56 reviews some related work of LDA and 2D-LDA. Section 3 details the CCLDA
57 and its relationship with LDA and 2D-LDA. Experiments and discussion
58 on ORL, Extended Yale-B, PIE, and FRGC databases are demonstrated in
59 Section 4 and in Section 5, we conclude the paper.

60 **2. Related Works**

61 Linear discriminant analysis (LDA), due to its simplicity and effective-
62 ness, achieved great success in pattern recognition. The essential idea of
63 LDA is to find the optimal subspace that gathers the samples from the same
64 class and meanwhile disperses the samples from the different ones. Given

Notation	Description
N	number of images in dataset
L	number of classes in dataset
C_i	the i -th class
N_i	number of images in class C_i
\mathbf{X}_i^k	i -th image sample in class C_k
\mathbf{x}_i^k	i -th sample in vector representation in class C_k
r	number of rows in an image
c	number of columns in an image
\mathbf{u}_i	the mean vector of samples in class C_i
\mathbf{u}	the mean vector of samples in dataset
\mathbf{U}_i	the mean matrix of samples in class C_i
\mathbf{U}	the mean matrix of samples in dataset

Figure 2: Notations

65 the data set $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]$, the between class scatter matrix \mathbf{S}_b and
66 within class scatter matrix \mathbf{S}_w are defined as

$$\mathbf{S}_w = \frac{1}{N} \sum_{i=1}^L \sum_{j=1}^{N_i} (\mathbf{x}_j^i - \mathbf{u}_i)(\mathbf{x}_j^i - \mathbf{u}_i)^T$$

$$\mathbf{S}_b = \frac{1}{N} \sum_{i=1}^L N_i (\mathbf{u}_i - \mathbf{u})(\mathbf{u}_i - \mathbf{u})^T$$
(1)

67 LDA aims to learn the projective directions \mathbf{W} which maximize the ratio of
68 between class scatter matrix to within class scatter one as

$$J = \frac{|\tilde{\mathbf{S}}_b|}{|\tilde{\mathbf{S}}_w|} = \frac{|\mathbf{W}^T \mathbf{S}_b \mathbf{W}|}{|\mathbf{W}^T \mathbf{S}_w \mathbf{W}|}$$
(2)

69 The optimal projection matrix \mathbf{W}_{opt} can be obtained by solving the fol-
 70 lowing eigen-value problem

$$\mathbf{S}_w^{-1}\mathbf{S}_b\mathbf{W} = \mathbf{W}\Lambda \quad (3)$$

71 where Λ is the diagonal matrix whose diagonal elements are eigenvalues of
 72 $\mathbf{S}_w^{-1}\mathbf{S}_b$.

73 In real applications, especially the image involving problem, the dimen-
 74 sion of feature is usually much larger than the sample size, so the within
 75 scatter matrix is singular. It is called small sample size or under sample
 76 problem. On the other hand, even if the \mathbf{S}_w is full, due to the high dimen-
 77 sion of feature, the sample size is still too small to reflect the class manifold
 78 in feature space and brings out the over-fitting problem.

79 Later, some researchers pointed out that the vectorization in traditional
 80 LDA disturbed the structure of image which is important for classification.
 81 2D-LDA which takes operation on image matrices directly rather than vectors
 82 is consequently proposed. Let the sample set be $\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n]$, and \mathbf{X}_i
 83 is the i -th image matrix with $r \times c$ size. The within-class scatter matrix \mathbf{S}_w
 84 and the between-class scatter matrix \mathbf{S}_b based on image matrix are defined
 85 as follows:

$$\begin{aligned} \mathbf{S}_w^{2d} &= \frac{1}{N} \sum_{i=1}^L \sum_{j=1}^{N_i} (\mathbf{X}_j^i - \mathbf{U}_i)(\mathbf{X}_j^i - \mathbf{U}_i)^T \\ \mathbf{S}_b^{2d} &= \frac{1}{N} \sum_{i=1}^L N_i (\mathbf{U}_i - \mathbf{U})(\mathbf{U}_i - \mathbf{U})^T \end{aligned} \quad (4)$$

86 The 2D-LDA searches such optimal projections that after projecting the
 87 original data onto these directions, the trace of the ratio of the resulting
 88 between-class scatter matrix to the within-class scatter matrix is maximized.

89 Let \mathbf{W}^{2d} denote a $r \times d$ ($d < r$) projection matrix, and the $r \times c$ image matrix
 90 \mathbf{X} is projected onto \mathbf{W}^{2d} through the following linear transformation:

$$\mathbf{Y} = \mathbf{W}^{2dT} \mathbf{X} \quad (5)$$

91 where the resulting \mathbf{Y} is a $d \times c$ matrix with smaller size than \mathbf{X} . 2D-LDA
 92 chooses \mathbf{W}^{2d} so that the following object function is maximized:

$$J = \frac{|\mathbf{W}^{2dT} \mathbf{S}_b^{2d} \mathbf{W}^{2d}|}{|\mathbf{W}^{2dT} \mathbf{S}_w^{2d} \mathbf{W}^{2d}|} \quad (6)$$

93 The optimal projection matrix \mathbf{W}_{opt}^{2d} can be obtained by solving the fol-
 94 lowing eigen-value problem

$$\mathbf{S}_w^{2d-1} \mathbf{S}_b^{2d} \mathbf{W}^{2d} = \mathbf{W}^{2d} \Lambda \quad (7)$$

95 where Λ is the diagonal matrix whose diagonal elements are eigenvalues of
 96 $\mathbf{S}_w^{2d-1} \mathbf{S}_b^{2d}$.

97 The process described above is called the left unilateral 2D-LDA (Li and
 98 Yuan (2005); Xiong et al. (2005)). There are also other 2D-LDA versions like
 99 right unilateral 2D-LDA and bilateral 2D-LDA (Ye et al. (2004); Yang et al.
 100 (2005); Kong et al. (2005)) etc.

101 **3. Contextual Constraints based Linear Discriminant Analysis (C-** 102 **CLDA)**

103 *3.1. Formulation of CCLDA*

104 As described in Section 1, it is believed that the contextual information
 105 could be utilized to improve the classification performance. As we all know,
 106 the essence of most subspace learning methods is to assign weights to different

107 pixels respectively and then sum them up to project the samples into low-
 108 dimension subspace to be best classified. Direct modeling the contextual
 109 information on pixels seems a little difficult. In this work, we adopt an
 110 alternative way to impose the contextual constraints on the weights instead
 111 of the pixels. The idea is that if the pixels in original image have consistent
 112 properties, the weights corresponding to them should also have similar values.
 113 In this way, we transfer the contextual constraints from the image to the
 114 weight. That is, if we arrange the weights to form a image (called weight
 115 image) according to the order of the pixels in original image, there is also
 116 contextual constraints on the weight image.

117 As stated above, intuitively, if the pixels are of the similar property
 118 or reflect the similar structure, the weights on them would have strong
 119 relationship, otherwise the weights on independent pixels would also be
 120 weakly related. Following this rational, we impose a constraints $J_2(\mathbf{w}) =$
 121 $\frac{1}{2} \sum_{i,j} (w_i - w_j)^2 S_{ij}$ ¹ on traditional LDA to reformulate the object of dis-
 122 criminant analysis as

$$J = \frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{(\mathbf{w}^T \mathbf{S}_w \mathbf{w} + \eta J_2(\mathbf{w}))} \quad (8)$$

123 where S_{ij} describes the similarity of pixels i and j , and η is a coefficient to
 124 balance the trade-off between the training discriminant power and contex-
 125 tual constraints. The constraint function $J_2(\mathbf{w})$ gives a high penalty when
 126 the weights of related pixels differ too much. Due to the symmetry of S_{ij}
 127 in general case, the contextual constraints $J_2(\mathbf{w})$ on weight image can be

¹For ease of interpretation, we discuss the one-dimensional subspace case here and it is easy to be extended to d -dimensional case.

128 formulated using the matrix operations further as follows.

$$\begin{aligned}
J_2(\mathbf{w}) &= \frac{1}{2} \sum_{i,j} (w_i - w_j)^2 S_{ij} \\
&= \frac{1}{2} \left(\sum_{i,j} w_i^2 S_{ij} + \sum_{i,j} w_j^2 S_{ij} - 2 \sum_{i,j} w_i w_j S_{ij} \right) \\
&= \sum_{i,j} w_i^2 S_{ij} - \sum_{i,j} w_i w_j S_{ij} \\
&= \sum_i w_i^2 D_{ii} - \sum_{i,j} w_i w_j S_{ij} \\
&= \mathbf{w}^T \mathbf{D} \mathbf{w} - \mathbf{w}^T \mathbf{S} \mathbf{w} \\
&= \mathbf{w}^T \mathbf{L}^w \mathbf{w}
\end{aligned} \tag{9}$$

129 where $\mathbf{L}^w = \mathbf{D} - \mathbf{S}$ is the Laplacian matrix², and \mathbf{D} is a diagonal matrix where
130 $D_{ii} = \sum_j S_{ij}$. Thus, the objective of CCLDA (Eq. 10) can be reformulated
131 as

$$J = \frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{(\mathbf{w}^T \mathbf{S}_w \mathbf{w} + \eta \mathbf{w}^T \mathbf{L}^w \mathbf{w})} \tag{10}$$

132 The optimal projection \mathbf{w} can be obtained by solving the following general-
133 ized eigen-value problem.

$$\mathbf{S}_b \mathbf{w} = \lambda (\mathbf{S}_w + \eta \mathbf{L}^w) \mathbf{w} \tag{11}$$

134 For a nature image, one assumption often used is that the pixels in local
135 neighboring region have the consistent properties and reflect the similar im-
136 age structure. Therefore, one way to define the contextual matrix \mathbf{S} among

²The formulation of $J_2(\mathbf{w})$ is very similar to LPP. However, they are totally different. LPP describes the neighboring sample structure in high-dimensional feature space while in this paper, we use this formulation to model the weight contextual information among neighboring pixels in one image.

137 different weights of pixels is

$$S_{ij} = \begin{cases} 1 & \text{if } i \text{ and } j \text{ are neighbors} \\ 0 & \text{otherwise} \end{cases} \quad (12)$$

138 If pixel i is in the neighboring region of j or pixel j is in the neighboring
 139 region of i , the corresponding similarity of w_i and w_j is set to be 1, otherwise,
 140 the similarity is set to be 0. In this way, all the pixels in a local adjacent
 141 regions are treated equally and the pixels in different regions are assumed to
 142 be independent. It is known as a rigid constraint. However, sometimes the
 143 pixel values in neighboring regions may change rapidly such as pixels around
 144 the edge. It might be improper to impose the same weight on these pixels.
 145 Therefore, we could define the similarity matrix of weights in another way
 146 named soft constraint as follows.

$$S_{ij} = \begin{cases} e^{-\|f_i - f_j\|^2 / \sigma^2} & \text{if } i \text{ and } j \text{ are neighbors} \\ 0 & \text{otherwise} \end{cases} \quad (13)$$

147 where f_i and f_j are the feature vectors extracted at position i and j respec-
 148 tively to describe the texture and spatial relationship between positions i
 149 and j . It should be noted the definition of weight similarity matrix \mathbf{S} is not
 150 limited to above two versions. Different definitions are possible according to
 151 different problems.

152 3.2. Relationship of CCLDA, LDA and 2DLDA

153 It is obvious that when $\eta = 0$, CCDA will degenerated into classical
 154 LDA. In (Zheng et al. (2008)), it is revealed that the formulations of LDA
 155 and 2DLDA have the following relationship.

Theorem 1. Let $\mathbf{w} = [\tilde{\mathbf{w}}_1^T, \dots, \tilde{\mathbf{w}}_c^T]^T$. If $\tilde{\mathbf{w}}_1, \dots, \tilde{\mathbf{w}}_c$ are imposed to be equivalent, i.e.,

$$\mathbf{w}^{2d} = \tilde{\mathbf{w}}_1 = \dots = \tilde{\mathbf{w}}_c \in \mathbf{R}^{r \times 1}$$

then the following relations are valid:

$$\begin{aligned} \mathbf{w}^T \mathbf{S}_b \mathbf{w} &= \mathbf{w}^{2dT} \mathbf{S}_b^{2d} \mathbf{w}^{2d} + \mathbf{w}^{2dT} \left\{ \sum_{k=1}^L \frac{N_k}{N} \sum_{j=1, h=1, j \neq h}^c (\mathbf{U}_k(j) \right. \\ &\quad \left. - \mathbf{U}(j))(\mathbf{U}_k(h) - \mathbf{U}(h))^T \right\} \mathbf{w}^{2d} \\ \mathbf{w}^T \mathbf{S}_w \mathbf{w} &= \mathbf{w}^{2dT} \mathbf{S}_w^{2d} \mathbf{w}^{2d} + \mathbf{w}^{2dT} \left\{ \frac{1}{N} \sum_{k=1}^L \sum_{i=1}^{N_k} \sum_{i=1, h=1, j \neq h}^c (\mathbf{X}_i^k(j) \right. \\ &\quad \left. - \mathbf{U}_k(j))(\mathbf{X}_i^k(h) - \mathbf{U}_k(h))^T \right\} \mathbf{w}^{2d}. \end{aligned}$$

156 Therefore, if the columns of image are independent and the projection
 157 weights for each column are the same, the result of 2DLDA is equivalent
 158 with that of LDA. However, as pointed in (Zheng et al. (2008)), Theorem 1
 159 indicates that 2DLDA drops the information that characterizes the covari-
 160 ance of columns of image, which may be helpful for classification. Therefore,
 161 the performance of 2DLDA is not always better than LDA.

162 In CCLDA, we impose a constraint on projection to generate the weights
 163 in neighboring regions such that they are close together. Particular, if the
 164 similarity matrix \mathbf{S} is constructed in the way that the similarities of pixels
 165 from the same row are set to be 1 and others are 0, we have the following
 166 theorem.

$$S_{ij} = \begin{cases} 1 & \text{if } i \text{ and } j \text{ are in the same row} \\ 0 & \text{otherwise} \end{cases} \quad (14)$$

Corollary 1. *When η in Eq. 10 is large enough and the \mathbf{S} is defined as Eq. 14, the CCLDA will approach to 2DLDA if the following equations hold:*

$$\sum_{k=1}^L \frac{N_k}{N} \sum_{j=1, h=1, j \neq h}^c (\mathbf{U}_k(j) - \mathbf{U}(j))(\mathbf{U}_k(h) - \mathbf{U}(h))^T = 0$$

$$\frac{1}{N} \sum_{k=1}^L \sum_{i=1}^{N_k} \sum_{i=1, h=1, j \neq h}^c (\mathbf{X}_i^k(j) - \mathbf{U}_k(j))(\mathbf{X}_i^k(h) - \mathbf{U}_k(h))^T = 0$$

167 The proof of this corollary is straightforward. The definition in Eq. 14
 168 assures that the LDA weights of pixels from the same row are as close as
 169 possible. When η is large, more concern will be paid on the contextual
 170 constraint and hence the derived solution will approach to the constraint of
 171 2DLDA if the columns of images are assumed to be independent.

172 Compared with 2DLDA, CCLDA utilizes the contextual information in a
 173 more reasonable and flexible way. Generally, in the formulation of CCLDA,
 174 it models the contextual information as well as preserves the covariance in-
 175 formation among different columns, so it incorporates the advantage of LDA
 176 and 2DLDA and should achieve better performance than LDA and 2DLDA
 177 in practice.

178 Another interesting point is the relationship between CCLDA and reg-
 179 ularized LDA (R-LDA) (Friedman (1989)). The formulations of these two
 180 methods are very similar. The CCLDA utilizes the contextual information by
 181 introducing a laplacian matrix while the R-LDA regularizes the LDA solution
 182 space by adding the identity matrix onto the within class scatter matrix. The
 183 characteristic of CCLDA is that it explicitly models contextual information
 184 in its formulation and hence owns more semantic meanings than R-LDA.

185 4. Experiments

186 Four face databases, ORL, Extended Yale-B, PIE, and FRGC are used to
187 evaluate the performance of CCLDA with Fisher LDA (FLDA) (Belhumeur
188 et al. (1997)), Null LDA (N-LDA) (Chen et al. (2000)), Direct LDA (D-
189 LDA) (Yu and Yang (2001)), R-LDA (Friedman (1989)) and Lu’s R-LDA (R-
190 LDA^L) (Lu et al. (2004)) etc. methods. For CCLDA, we use the soft con-
191 straints defined in Eq. 13 and extract the gray values of images at position
192 i to form the feature vector f_i . The 2DLDA is implemented following the
193 paper (Ye et al. (2004)).

194 4.1. Data Preparation

195 4.1.1. ORL

196 There are 40 persons in ORL database (Samaria and Harter (1994)), each
197 of which contains 10 images. The images are all frontal with slight tilt of
198 the head, including lighting and expression variations. All the images are
199 cropped to 32×32 size.

200 4.1.2. Extended Yale-B

201 The Extended Yale-B database (Georghiades et al. (2001)) contains 2414
202 images from 38 persons each of which contains nearly 64 ones. All the images
203 are rotated, scaled and cropped to 32×32 size.

204 4.1.3. PIE

205 The PIE database (Sim et al. (2003)) consists of 41,368 images from 68
206 people under different poses, illumination conditions and expressions. Five n-
207 ear frontal poses(C05, C07, C09, C27, C29) and all the images under different

208 illuminations and expressions are selected. So there are totally 170 images
209 for each individual. All the images are cropped to 32×32 in experiments.

210 4.1.4. FRGC

211 For FRGC database (Phillips et al. (2005)), we follow the experiment
212 1&4 protocols to evaluate different algorithms. There are 12,776 face images
213 from 222 individuals in training set, including 6360 controlled images and
214 6416 uncontrolled ones. In experiment 1, both target and query images
215 contains 16028 controlled images and in experiment 4, which is considered
216 as the most challenge case, there are 16028 controlled images as the target
217 images and 8014 query images which are uncontrolled ones, from 466 persons.
218 The images are captured over several sessions. All the faces are cropped to
219 64×64 size images.

220 4.2. Parameter Selection

221 There are mainly two parameters in our proposed method that affecting
222 the performance of algorithm. One is the value of σ in Eq. 13, another is the
223 contextual information regularized coefficient η in Eq. 10. In this experiment,
224 the feature vector f in Eq. 13 is extracted by grouping the the pixel values
225 at corresponding position from all the training images. The parameter σ is
226 empirically set to be the average distance among these feature vectors. For
227 the parameter η , we try to examine the impact of different values of η on the
228 performance of algorithm and then choose the best one.

229 Fig. 4 shows the recognition rate of CCLDA on the PIE face database
230 as a function of η . It is easy to see that the CCLDA can achieve better
231 recognition performance than LDA when the value of η is approximately



(a) ORL



(b) Extended Yale-B



(c) PIE



(d) FRGC

Figure 3: Sample examples from (a) ORL, (b) Extended Yale-B, (c) PIE and (d) FRGC databases.

232 between 0.00005 and 0.005, and it achieves its highest recognition rate at
233 the point of $\eta = 0.0005$. Therefore, in all the following experiments, the

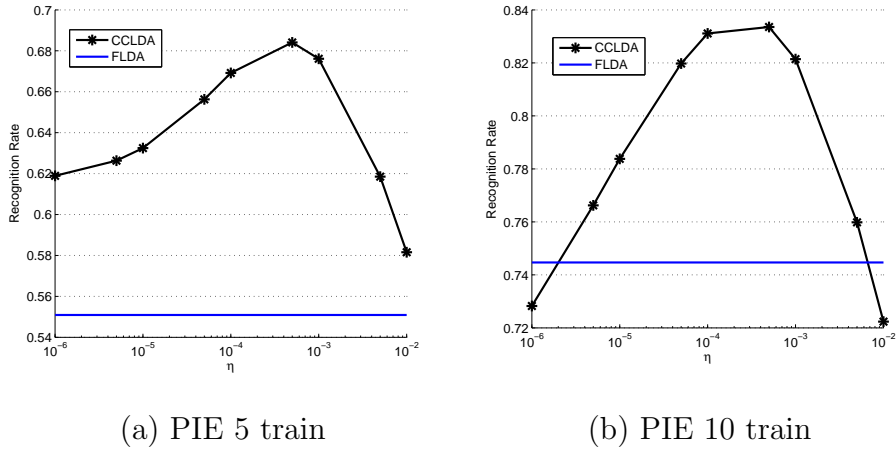


Figure 4: Parameter η selection for CCLDA. The curve shows the recognition rate with respect to η . The solid line shows the recognition rate of FLDA.

234 parameter η is set to 0.0005. For R-LDA and R-LDA^L, the regularized
 235 coefficient is set to 0.01 empirically in the same way.

236 4.3. Recognition Results

237 For the experiments on ORL, Extended Yale-B and PIE databases, the
 238 database is randomly partitioned into gallery and probe sets. For ease of
 239 representation, in following results, G_N denotes that N images per person
 240 are selected to form gallery set and the left ones consist the probe set. The
 241 gallery set is used for training, and in test phase, the images in probe set are
 242 compared with the images in gallery set and the nearest neighboring (NN)
 243 classifier is used to do the classification. All experimental results are reported
 244 as the average of 10 different runs on different splits. For FRGC, we follow
 245 the experiment 1 and 4 protocols and report the verification rate when the
 246 false accept rate is 0.001 in ROC I, ROC II and ROC III.

247 For FLDA, the PCA is first applied by keeping 98% energy, followed
 248 by LDA. The reduced dimension of FLDA, N-LDA, D-LDA, R-LDA and
 249 CCLDA are chosen to be $C - 1$ to retain the most discriminative information,
 250 where C is the class number. For 2DLDA, we finally keep 10×10 dimension
 251 for the classification.

Table 1: ORL: Classification accuracy (mean \pm std%) with different training set size.

Method	G2	G3	G4	G5	G6
FLDA	61.09 \pm 4.88	77.93 \pm 6.13	92.00 \pm 1.88	95.20 \pm 1.55	96.56 \pm 1.82
N-LDA	82.94 \pm 2.59	91.00 \pm 2.33	94.92 \pm 1.16	96.90 \pm 1.31	97.13 \pm 2.40
D-LDA	34.63 \pm 13.01	83.93 \pm 2.65	91.83 \pm 2.18	95.45 \pm 1.30	96.38 \pm 1.90
2DLDA	78.28 \pm 3.50	87.57 \pm 2.21	93.46 \pm 2.21	95.15 \pm 1.73	96.19 \pm 1.90
R-LDA	82.88 \pm 2.74	91.14 \pm 2.23	95.50 \pm 1.49	97.00 \pm 1.03	97.88 \pm 1.89
R-LDA ^L	76.06 \pm 4.78	84.18 \pm 2.76	91.88 \pm 2.18	95.45 \pm 1.30	96.44 \pm 1.93
CCLDA	85.94 \pm 2.59	92.25 \pm 2.08	96.08 \pm 1.31	97.30 \pm 1.03	97.81 \pm 1.62

252 4.4. Analysis and Discussion

253 As mentioned in Section 3.1, the essence of CCLDA is to impose the
 254 smooth constraint on the weight image \mathbf{w} so that the weight values in local
 255 region do not fluctuate too much. This is consistent with the property of
 256 natural images. Fig. 5 shows the first 5 projections of FLDA and CCLDA
 257 in the form of images. It shows that compared to the result of FLDA, the
 258 projections of CCLDA are much smoother; that is, the values of the elements

Table 2: Extended Yale-B: Classification accuracy (mean \pm std%) with different training set size.

Method	G10	G20	G30	G40	G50
FLDA	79.08 \pm 1.38	88.07 \pm 0.58	92.17 \pm 0.64	93.70 \pm 0.59	94.34 \pm 1.16
N-LDA	79.03 \pm 2.31	83.23 \pm 1.39	N/A	N/A	N/A
D-LDA	77.30 \pm 1.16	84.92 \pm 0.61	88.97 \pm 0.77	89.99 \pm 0.88	91.11 \pm 1.28
2DLDA	81.45 \pm 1.66	88.47 \pm 0.89	92.27 \pm 0.57	93.20 \pm 0.57	94.26 \pm 1.31
R-LDA	82.61 \pm 1.53	91.29 \pm 0.64	94.60 \pm 0.53	95.92 \pm 0.62	96.48 \pm 0.86
R-LDA ^L	64.06 \pm 2.43	77.30 \pm 1.17	84.92 \pm 0.61	88.97 \pm 0.77	89.97 \pm 0.87
CCLDA	82.66 \pm 1.47	91.23 \pm 0.77	94.62 \pm 0.63	95.87 \pm 0.53	96.28 \pm 1.04

Table 3: PIE: Classification accuracy (mean \pm std%) with different training set size.

Method	G10	G20	G30	G40	G50
FLDA	74.47 \pm 0.83	87.72 \pm 0.90	91.85 \pm 0.64	93.85 \pm 0.41	95.36 \pm 0.36
N-LDA	69.95 \pm 1.05	N/A	N/A	N/A	N/A
D-LDA	73.96 \pm 0.94	85.11 \pm 0.89	89.68 \pm 0.56	92.07 \pm 0.46	93.77 \pm 0.37
2DLDA	74.44 \pm 0.91	85.08 \pm 0.59	89.32 \pm 0.44	91.69 \pm 0.44	93.19 \pm 0.39
R-LDA	83.17 \pm 0.97	91.89 \pm 0.67	94.35 \pm 0.27	95.56 \pm 0.27	96.54 \pm 0.32
R-LDA ^L	58.95 \pm 1.49	73.94 \pm 0.94	85.10 \pm 0.89	89.68 \pm 0.56	92.07 \pm 0.46
CCLDA	83.35 \pm 0.89	92.04 \pm 0.69	94.50 \pm 0.31	95.71 \pm 0.31	96.70 \pm 0.29

Table 4: FRGC: Verification rates when false accept rate is at 0.001.

Method	Experiment 1			Experiment 4		
	ROC I	ROC II	ROC III	ROC I	ROC II	ROC III
FLDA	0.7932	0.7746	0.7531	0.4514	0.4464	0.4393
D-LDA	0.8714	0.8438	0.8378	0.4240	0.4092	0.3879
2DLDA	0.6384	0.5745	0.5059	0.1735	0.1612	0.1481
R-LDA	0.8715	0.8520	0.8378	0.4844	0.4738	0.4633
R-LDA ^L	0.8720	0.8439	0.8173	0.4241	0.4085	0.3868
CCLDA	0.9222	0.9111	0.9003	0.6022	0.6005	0.6005

259 in neighboring areas are much closer, thus satisfying the motivation behind
 260 this work.

261 Tables 1-4 illustrate the recognition performance of CCLDA and other
 262 popular methods on ORL, Extended Yale-B, PIE and FRGC databases. The
 263 results reveal some interesting points.

- 264 1. In almost all experiments, the performance of CCLDA is always better
 265 than that of LDA, N-LDA, D-LDA, 2DLDA and R-LDA^L. It proves
 266 that the CCLDA method is effective for solving image related pattern
 267 recognition problem such as face recognition.
- 268 2. Compared with FLDA, which is based on vector representation, 2DL-
 269 DA and CCLDA, which consider the image contextual constraints, sig-
 270 nificantly improve the recognition performance when the sample size
 271 is small. It indicates that the contextual information is important and
 272 can be considered as a good prior for image recognition. Moreover, the

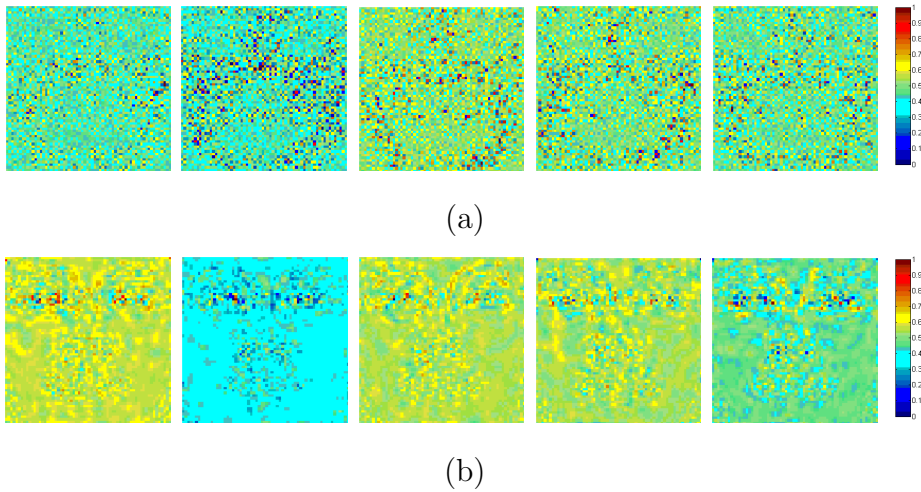


Figure 5: The first five weight images of FLDA (a) and CCLDA (b) on FRGC training set. The weight images of CCLDA are much smoother than those of FLDA.

- 273 performance of CCLDA is always better than 2DLDA, which proves the
 274 CCLDA could better utilize the contextual information than 2DLDA
 275 does.
- 276 3. Compared CCLDA to R-LDA, the performance of these two methods
 277 are very close. Overall, the performance of CCLDA is slightly better
 278 than that of R-LDA, but not significant. The advantage of CCLDA is
 279 that it explicitly models the contextual information in its formulation
 280 and it carries more definite semantic meanings than R-LDA.
- 281 4. On FRGC, where the experiment is considered as the most challenge
 282 one, CCLDA significantly enhances the recognition performance (much
 283 better than R-LDA in this case). It is surprising that the performance
 284 of 2DLDA is much worse than that of FLDA in this case. One possible
 285 explanation is that the training set size of FRGC is large enough for

286 FLDA. Compared with 2DLDA, CCLDA takes into account the con-
287 textual information in a more reasonable and flexible way and achieves
288 better results.

289 **5. Conclusions**

290 In this paper, we propose contextual constraints based linear discriminant
291 method and apply it to image recognition problem. The motivation is that
292 the contextual information is important for image understanding. Therefore,
293 it is possible that contextual information could provide useful clues that help
294 improve image classification performance. Extensive experiments validate
295 our motivation. The proposed contextual constraints could be incorporated
296 into not only linear discriminant analysis, but also other subspace learning
297 methods like LPP, NPE etc. One of our future work is to extend the contex-
298 tual information to the non-linear (kernel) formulation so as to improve the
299 classification performance further.

300 **References**

- 301 Belhumeur, P., Hespanha, J., Kriegman, D., 1997. Eigenfaces vs. fisherfaces: recog-
302 nition using class specific linear projection. *IEEE Trans. PAMI* 19 (7), 711–720.
- 303 Chen, L., Liao, H., Ko, M., Lin, J., Yu, G., October 2000. A new lda-based
304 face recognition system which can solve the small sample size problem. *Pattern*
305 *Recognition* 33 (10), 1713–1726.
- 306 Dass, S. C., Jain, A. K., July 9-12 2001. Markov face models. In: *ICCV*. Vancouver,
307 pp. 112–116.

- 308 Dass, S. C., Jain, A. K., Lu, X., August 2002. Face detection and synthesis using
309 Markov random field models. In: ICPR. Quebec City, pp. 112–116.
- 310 Friedman, J. H., March 1989. Regularized discriminant analysis. *Journal of the*
311 *American Statistical Association* 84 (405), 165.
- 312 Georghiades, A., Belhumeur, P., Kriegman, D., 2001. From few to many: Illumi-
313 nation cone models for face recognition under variable lighting and pose. *IEEE*
314 *Trans. PAMI* 23 (6), 643–660.
- 315 Gui, J., Jia, W., Zhu, L., Wang, S., Huang, D., 2010. Locality preserving discrimi-
316 nant projections for face and palmprint recognition. *Neurocomputing* 73 (13-15),
317 2696–2707.
- 318 He, X., Cai, D., Yan, S., Zhang, H., 2005a. Neighborhood preserving embedding.
319 In: *ICCV*. pp. 1208–1213.
- 320 He, X., Yan, S., Hu, Y., Niyogi, P., Zhang, H., 2005b. Face recognition using
321 laplacianfaces. *IEEE Trans. PAMI* 27 (3), 328–340.
- 322 Huang, R., Pavlovic, V., Metaxas, D., August 2004. A hybrid face recognition
323 method using markov random fields. In: *ICPR*. Cambridge, UK, pp. 157–160.
- 324 Kong, H., Wang, L., Teoh, E. K., Wang, J.-G., Venkateswarlu, R., 2005. A frame-
325 work of 2d fisher discriminant analysis: Application to face recognition with
326 small number of training samples. In: *CVPR*. Vol. 2. pp. 1083–1088.
- 327 Li, M., Yuan, B., 2005. 2d-lda: A novel statistical linear discriminant analysis for
328 image matrix. *Pattern Recognition Letters* 26 (5), 527–532.

- 329 Lu, J., Plataniotis, K. N., Venetsanopoulos, A. N., 2004. Regularization studies
330 on lda for face recognition. In: Proceedings of IEEE International Conference
331 on Image Processing. pp. 63–66.
- 332 Phillips, P. J., Flynn, P. J., Scruggs, W. T., Bowyer, K. W., Chang, J., Hoffman,
333 K., Marques, J., Min, J., Worek, W. J., 2005. Overview of the face recognition
334 grand challenge. In: CVPR. pp. 947–954.
- 335 Roweis, S., Saul, L., 2000. Nonlinear dimensionality reduction by locally linear
336 embedding. *Science* 290 (22), 2323–2326.
- 337 Samaria, F., Harter, A., 1994. Parameterisation of a stochastic model for human
338 face identification. In: Proceedings of 2nd IEEE Workshop on Applications of
339 Computer Vision. Sarasota FL, pp. 138–142.
- 340 Sim, T., Baker, S., Bsat, M., 2003. The CMU pose, illumination, and expression
341 database. *IEEE Trans. PAMI* 25 (12), 1615–1618.
- 342 Turk, M. A., Pentland, A. P., 1991. Face recognition using eigenfaces. In: CVPR.
343 pp. 586–591.
- 344 Wang, L., Zhang, Y., Feng, J., 2005. On the euclidean distance of images. *IEEE*
345 *Trans. PAMI* 27 (8), 1334–1339.
- 346 Xiong, H., Swamy, M. N. S., Ahmad, M. O., July 2005. Two-dimensional fld for
347 face recognition. *Pattern Recognition* 38, 1121–1124.
- 348 Yang, J., Zhang, D., Yong, X., Yang, J.-Y., 2005. Two-dimensional discriminant
349 transform for face recognition. *Pattern Recognition* 38, 1125–1129.
- 350 Ye, J., Janardan, R., Li, Q., 2004. Two-dimensional linear discriminant analysis.
351 In: NIPS. Vol. 17. pp. 1569–1576.

- 352 Yu, H., Yang, J., October 2001. A direct lda algorithm for high-dimensional data
353 with application to face recognition. *Pattern Recognition* 34 (10), 2067–2070.
- 354 Zheng, W.-S., Lai, J. H., Li, S. Z., 2008. 1D-LDA vs. 2D-LDA: When is vector-
355 based linear discriminant analysis better than matrix-based? *Pattern Recogni-
356 tion* 47, 2156–2172.