

Ranking Principal Components in Face Spaces Through AdaBoost.M2 Linear Ensemble

Tiene A. Filisbino, Gilson A. Giraldi
National Laboratory for Scientific Computing- LNCC
Petrópolis, Brazil
{tiene,gilson}@lncc.br

Carlos Eduardo Thomaz
Department of Electrical Engineering, FEI
São Bernardo do Campo, Brazil
cet@fei.edu.br

Abstract—Despite the success of Principal Component Analysis (PCA) for dimensionality reduction, it is known that its most expressive components do not necessarily represent important discriminant features for pattern recognition. In this paper, the problem of ranking PCA components, computed from multi-class databases, is addressed by building multiple linear learners that are combined through the AdaBoost.M2 in order to determine the discriminant contribution of each PCA feature. In our implementation, each learner is a weakened version of a linear support vector machine (SVM). The strong learner built by the ensemble technique is processed following a strategy to get the global discriminant vector to sort PCA components according to their relevance for classification tasks. Also, we show how the proposed methodology to compute the global discriminant vector can be applied to other multi-class approaches, like the linear discriminant analysis (LDA). In the computational experiments we compare the obtained approaches with counterpart ones using facial expression experiments. Our experimental results have shown that the principal components selected by the proposed technique allows higher recognition rates using less linear features.

Keywords-PCA; Ranking PCA Components; Separating Hyperplanes; Ensemble Methods; AdaBoost; Face Image Analysis

I. INTRODUCTION

Nowadays, increasingly large amount of high dimensional image databases are being generated, leading to a strong demand for dimensionality reduction for discarding redundancy, and features selection techniques to reduce the feature space for discriminating sample groups before executing classification tasks [1].

In this avenue, we follow statistical learning approaches whose basic pipeline can be described as follows [2]: (a) Linear subspace learning for dimensionality reduction; (b) Among the linear components obtained, select the most discriminant ones; (c) Solve the classification problem; (d) Reconstruction problem, that is, visualize the information captured by the discriminant linear components.

The step (a) can be accomplished through classical works on linear dimensionality reduction including the principal component analysis (PCA), factor analysis (FA) [3], multi-dimensional scaling (MDS) [1] and projection pursuit (PP) [4], [3]. The determination of discriminant features (step (b) above) is very known in the context of PCA. In this case, it was observed that, since PCA explains the covariance structure of all the data its most expressive components, that

is, the first principal components with the largest eigenvalues, do not necessarily represent the most important discriminant directions to separate sample groups [5], [6]. This observation motivates the development of specific techniques to compute discriminant subspaces which, in general, depend on the incorporation of prior information based on labeled data. The Fisher's linear discriminant analysis (LDA) [1], discriminant principal components analysis (DPCA) [5] and its extension to multi-class problems, named Multi-Class DPCA [7], Zhu and Martinez [8] criterion, are techniques reported in the literature for discriminant features selection.

In this work we focus on discriminant analysis on multi-class problems. In this case, given an N -class database, the Multi-Class DPCA builds a linear support vector machine (SVM) ensemble, composed of N SVM machines, to get the discriminant weights that are combined through the AdaBoost technique in order to determine the discriminant contribution of each feature.

Contributions: In this paper we keep the Multi-Class DPCA methodology, but we replace the AdaBoost by the AdaBoost.M2 algorithm and combine the separating SVM hyperplanes through a simple strategy to compute the global discriminant weights. In this way, we get a new ranking method for the principal components, called Multi-Class.M2 DPCA algorithm, given by the group-differences extracted by a linear ensemble based on the AdaBoost.M2 technique. The computational experiments demonstrate that the new discriminant technique improves the Multi-Class DPCA for both reconstruction and recognition. Also, we show that the proposed methodology to compute discriminant weights can be applied to other multi-class approaches, like the linear discriminant analysis (LDA).

It is important to highlight that we do not deal with the problem of computing general discriminant directions that are not principal components. Rather, we apply the idea of using a set of linear classifiers and an ensemble method (AdaBoost.M2, in this case) to compute a matrix of discriminant weights that is processed to select among the principal components the most discriminant ones. We have focused here on the SVM [9] method but any other separating hyperplane could be used.

To evaluate the Multi-Class.M2 DPCA algorithm, we perform group separation tasks in facial expression experiments involving neutral, happiness, sad, fear, and anger face images.

The experiments show that the SVM can be used as an effective component classifier to generate the discriminant weights for the multi-class discriminant principal components analysis. Furthermore, the computational experiments demonstrate the benefits of sorting principal components using the Multi-Class.M2 DPCA if compared with the traditional PCA, and the Multi-Class DPCA methodologies for selecting PCA components.

The paper is organized as follows. In section I-A we survey related works for discriminant analysis. Next, section I-B presents the main stages of the proposed method. Then, in section II we review the theory behind DPCA approach. Section III presents the Multi-Class.M2 DPCA algorithm. The computational experiments are described in section IV. Finally, in Section V, we conclude the paper, summarizing its main contributions and describing further developments.

A. Related work

Given a feature space, a key question is "how can we determine (or compute) the most important discriminant features for a pattern recognition task, like classification?" Discriminant analysis techniques address this question, which is very known in the context of PCA.

The Figure 1 is a simple example that pictures the limitation of PCA for discriminant features extraction. Both Figures 1.(a) and 1.(b) represent the same data set. Figure 1.(a) just shows the PCA directions (\tilde{x} and \tilde{y}) and the distribution of the samples over the space. However, in Figure 1.(b) we distinguish two patterns: plus (+) and triangle (▼). We observe that the principal PCA direction \tilde{x} can not discriminate samples of the considered groups because the projection of the data points over direction \tilde{x} will mix the patterns in the corresponding one-dimensional subspace.

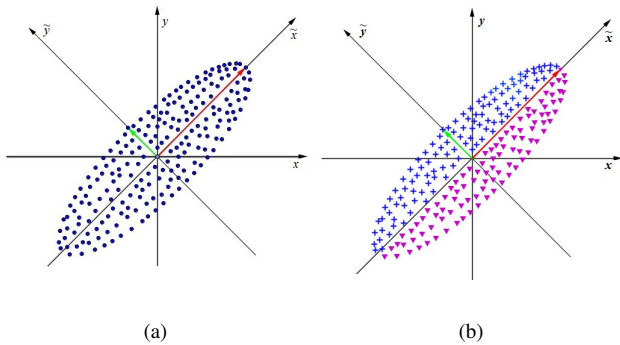


Fig. 1. (a) Scatter plot and PCA directions. (b) The same population but distinguishing patterns plus (+) and triangle (▼).

In general, Fisher's linear discriminant analysis (LDA) is used to identify the most important linear directions for separating sample groups rather than PCA [1]. This method, as well as the weighted pairwise variant of the well-known multi-class Fisher criterion introduced in [10] has the limitation of finding number of groups - 1 meaningful discriminant directions.

In [5] it is proposed the DPCA technique, based on the idea of using the discriminant weights obtained by separating hyperplanes to select among the principal components the most discriminant ones. In [7] the DPCA was extended for multi-class problems and, the so called multi-class discriminant principal components analysis (Multi-Class DPCA), consists of the following steps: (a) apply PCA technique for dimensionality reduction in order to eliminate redundancy. (b) Compute a linear ensemble, based on the one-against-all SVM multi-class approach. (c) Combine the discriminant weights computed through the separating SVM hyperplanes in order to determine the discriminant contribution of each feature. So, given a N -class database, the step (b) builds N SVM machines in the PCA space. The step (c) is implemented by adapting an ensemble technique, the AdaBoost one [11], to yield a global discriminant vector. The proposed solution was evaluated in group separation tasks involving facial expression experiments and achieves higher recognition rates using less PCA features. However, the Multi-Class DPCA is not efficient for reconstruction. Also, the number of iterations (step (b) above) is equal to the number of classes in the main loop of Multi-Class DPCA. Such characteristic may limits the ability of the method to select discriminant features. These drawbacks have motivated the current work that is described next.

B. Technique overview

The whole Multi-Class.M2 DPCA methodology is schematized in Figure 2. We follow [7] and keep the application of PCA technique for dimensionality reduction in the step (1) of the pipeline. Then, in step (2), we compute a set of linear SVM hyperplanes, based on the one-against-all SVM multi-class approach. We also apply an ensemble technique, the AdaBoost.M2 algorithm, to combine the linear classifiers in order to compute the global discriminant vector. The key idea of this step is based on the fact that AdaBoost.M2 linearly combines weak classifiers to get the strong hypothesis. So, it is straightforward to obtain the global discriminant weights from the expression that defines the strong classifier by using a simple scheme, that corresponds to step (3) of Figure 2. This strategy can be also used to combine discriminant directions computed by other multi-class approaches, like linear discriminant analysis (LDA).

However, it is known that a strong learner like SVM does not work well as the base component for Adaboost [12]. Therefore, we follow [12] and implement a strategy to compute a weakened version of SVM that is useful as an Adaboost.M2 component [13].

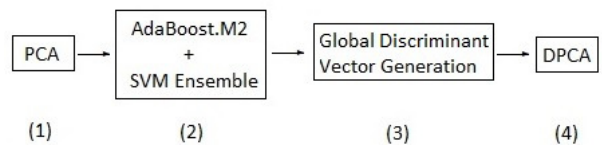


Fig. 2. Flowchart with main steps of the proposed technique.

Finally, in the stage (4) of Figure 2, we follow the traditional DPCA proposal and sort PCA components in the decreasing order of the global discriminant weights. The result of Multi-Class.M2 DPCA algorithm is the PCA components arranged according to the discriminant weights. The method is not restricted to any application or particular probability density function of the sample groups and the number of meaningful discriminant directions is not limited to the number of groups.

II. TECHNICAL BACKGROUND

The Multi-Class.M2 DPCA technique is based on the DPCA [5], the weakened SVM proposed in [12], the AdaBoost.M2 algorithm described in [13], and the nonseparable linear SVM [9]. Following, we describe the DPCA methodology. The reader can find a summary of all the other techniques in [14].

Let the training observations $\mathbf{x}_i \in \mathfrak{R}^n$, $i = 1, \dots, M$ that generate a $M \times n$ training set matrix $\tilde{\Theta}$ centered respect to the global mean $\tilde{\mathbf{x}}$. Hence, the PCA algorithm computes a transformation matrix $P_{pca} = [\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_{m'}]$ whose columns \mathbf{p}_i , $i = 1, \dots, m'$ minimize the mean square reconstruction error, being the $m' \leq n$ eigenvectors of the covariance matrix Ω of $\tilde{\Theta}$ that correspond to the m' largest eigenvalues [15].

If to each training sample \mathbf{x}_i it is associated a label $y_i \in \{-1, 1\}$, then we have a labeled training set:

$$X = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2) \dots (\mathbf{x}_M, y_M)\}, \quad (1)$$

and, we can apply the DPCA technique to select the most discriminant principal components to separate sample groups.

The original DPCA is implemented taking as input a training set X , like in expression (1). Firstly, for discarding redundancies, the PCA transformation matrix $P_{pca} = [\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_{m'}]$ is computed and each zero mean data vector $\tilde{\mathbf{x}}_i$ is projected generating a vector $\bar{\mathbf{x}}_i = (P_{pca})^T \tilde{\mathbf{x}}_i$. Afterwards, the obtained $M \times m'$ data matrix and their corresponding labels are used as input to calculate the separating hyperplane. In the following we focus on the SVM technique, although any other linear classifier could be used.

Since DPCA assumes only two classes to separate, there are only one discriminant vector $\phi_{svm} = (w_1, w_2, \dots, w_{m'})$ given by the SVM hyperplane. If we multiply the $M \times m'$ most expressive features matrix by the $m' \times 1$ discriminant SVM vector:

$$\begin{aligned} c_1 &= \bar{\mathbf{x}}_{11}w_1 + \bar{\mathbf{x}}_{12}w_2 + \dots + \bar{\mathbf{x}}_{1m'}w_{m'}, \\ c_2 &= \bar{\mathbf{x}}_{21}w_1 + \bar{\mathbf{x}}_{22}w_2 + \dots + \bar{\mathbf{x}}_{2m'}w_{m'}, \\ &\dots \\ c_M &= \bar{\mathbf{x}}_{N1}w_1 + \bar{\mathbf{x}}_{N2}w_2 + \dots + \bar{\mathbf{x}}_{Nm'}w_{m'}. \end{aligned} \quad (2)$$

we get the most discriminant feature $c_i \in \mathbb{R}$ of each one of the m' -dimensional vectors $\bar{\mathbf{x}}_i$. Therefore, we can determine the discriminant contribution of each feature by investigating the weights $[w_1, w_2, \dots, w_{m'}]$. In fact, weights that are estimated to be 0 or approximately 0 have negligible contribution on the discriminant scores c_i described in equation (2), indicating

that the corresponding features are not significant to separate the sample groups. In contrast, largest weights (in absolute values) indicate that the corresponding features contribute more to the discriminant score and consequently are important to characterize the differences between the groups.

Therefore, instead of sorting these features by selecting the corresponding principal components in decreasing order of eigenvalues, as PCA does, DPCA selects as the most important features for classification the ones with the highest discriminant weights, that is, $|w_1| \geq |w_2| \geq \dots \geq |w_{m'}|$.

III. MULTI-CLASS.M2 DISCRIMINANT ANALYSIS

The Multi-Class.M2 DPCA procedure is described by the Algorithm 1. At the input of the procedure, the training instances in the database $X \subset \mathfrak{R}^n$ are supposed independently and identically distributed from an uniform distribution D . Following the pipeline in Figure 2, the first stage of Multi-Class.M2 DPCA applies the PCA, for dimensionality reduction in order to eliminate redundancy (line 2).

The labeled projected data set is built in line 3 and composes the input to generate the weak SVM classifiers. In line 9 of Multi-Class.M2 DPCA algorithm, each weak learner generates an hypotheses, which has the form $h : X \times Y \rightarrow [0, 1]$, and can be interpreted as the probability that y is the correct label associated with instance \mathbf{x} . So, given a sample \mathbf{x}_i , the probability of choosing an incorrect label y is [13]: $P_r = \frac{1}{2} (1 - h(\mathbf{x}_i, y_i) + h(\mathbf{x}_i, y))$.

However, we have $|Y|-1$ possibilities to obtain the incorrect answer. So, we can define the loss of the hypothesis through a weighted average according to some $q_{i,y}$, called the label weighting function, that assigns to each example i in the training set a load, with $\sum_{y \neq y_i} q_{i,y} = 1$. The resulting formula is called the pseudo-loss of h on training instance i with respect to q [13]:

$$ploss_q(h, i) = \frac{1}{2} \left(1 - h(\mathbf{x}_i, y_i) + \sum_{y \neq y_i} q_{i,y} h(\mathbf{x}_i, y) \right). \quad (3)$$

So, following the AdaBoost.M2 strategy [13], in each iteration t of the Algorithm 1, the weak learner's goal is to minimize the expected pseudo-loss, computed in line 10 of the Algorithm 1, for a distribution D^t and weighting function q^t . The algorithm uses a second weight vector whose values at time t are denoted by $w_{i,y}^t$, $i = 1, \dots, M$, $y \in Y - \{y_i\}$, which is initialized in line 1, based on the initial distribution D . The main loop of the algorithm aims to update these weights in order to minimize the expected pseudo-loss. So, the weighting function q^t and the distribution D^t are computed using the $w_{i,y}^t$ (line 5 of procedure 1).

Next, the Multi-Class.M2 DPCA computes a set of SVM hyperplanes, based on the one-against-all SVM multi-class approach presented in [16]. Hence, as we have N classes, the internal loop in the Algorithm 1 (line 6 to 9) constructs N weakened SVMs, in the PCA subspace, using the Algorithm 2. To do this, in line 7 of Algorithm 1 we build the $\tilde{\Theta}^y$ set by taking all k_y projected samples from class y and label them as 1. Then, using random sampling we choose $(2k_y)/(N-1)$

projected samples from classes other than y and label them as -1 . The obtained set of feature vectors $\bar{\mathbf{x}}_m^y \in \mathbb{R}^{m'}$ and corresponding labels $y_m \in \{-1, 1\}$:

$$\bar{\Theta}^y = \left\{ (\bar{\mathbf{x}}_1^y, l_1), (\bar{\mathbf{x}}_2^y, l_2), \dots, (\bar{\mathbf{x}}_{3k_y}^y, l_{3k_y}) \right\}, \quad (4)$$

are the input to call the Algorithm 2 which construct the weak SVM (WSVM) model y , represented by a hyperplane direction (ϕ_y^t) and a linear coefficient (b_y^t) . Each hypothesis h^t , in line 9 of Algorithm 1, is generated through a WSVM and the following normalization function:

$$f(z) = \frac{z - z_{min,y}^t}{z_{max,y}^t - z_{min,y}^t}, \quad (5)$$

where $f: [z_{min,y}^t, z_{max,y}^t] \rightarrow [0, 1]$, with $z_{min,y}^t$ and $z_{max,y}^t$ being the minimum and maximum values, respectively, of the set $\{ \langle \mathbf{x}_i, \phi_y^t \rangle + b_y^t, i = 1, 2, \dots, M \}$.

The lines 16-18 of the Algorithm 1 are based on the AdaBoost.M2 idea of deriving a strong learner h_f by using the linear combination of weak (WSVM, in our case) learners h^1, h^2, \dots, h^T :

$$h_f(\mathbf{x}) = \arg \max_{y \in Y} \sum_{t=1}^T \tilde{\alpha}^t h^t(\mathbf{x}, y), \quad (7)$$

where $\tilde{\alpha}^t$ is computed in line 16. This expression offers the possibility of extending the DPCA methodology to multi-class problems using the Adaboost.M2 result. To see this, we shall remember that $h^t(x, y)$ in line 9 is computed through the function f , in expression (5), and rewrite expression (7) as:

$$h_f(\mathbf{x}) = \arg \max_{y \in Y} \left[\sum_{t=1}^T \tilde{\alpha}^t f(\langle \mathbf{x}, \phi_y^t \rangle + b_y^t) \right]. \quad (8)$$

But, from equation (5), we get:

$$f(\langle \mathbf{x}, \phi_y^t \rangle + b_y^t) = \frac{\langle \mathbf{x}, \phi_y^t \rangle + b_y^t - z_{min,y}^t}{z_{max,y}^t - z_{min,y}^t}. \quad (9)$$

Therefore, by substituting this expression into equation (8), and using the linearity of the inner product, we can show that:

$$h_f(\mathbf{x}) = \arg \max_{y \in Y} [\langle \mathbf{x}, \Phi_y \rangle + \psi_y], \quad (10)$$

where:

$$\Phi_y = \sum_{t=1}^T \tilde{\alpha}^t \frac{\phi_y^t}{z_{max,y}^t - z_{min,y}^t}, \psi_y = \sum_{t=1}^T \tilde{\alpha}^t \frac{(b_y^t - z_{min,y}^t)}{z_{max,y}^t - z_{min,y}^t},$$

with $\Phi_y \in \mathbb{R}^{m'}$ and $\psi_y \in \mathbb{R}$. The bias ψ_y can be incorporated in the inner product through a translation \bar{T}_y that satisfies $\langle \bar{T}_y, \Phi_y \rangle = \psi_y$, which renders:

$$h_f(x) = \arg \max_{y \in Y} \left[\sum_{i=1}^n (x_i + \bar{T}_{i,y}) \Phi_{i,y} \right]. \quad (11)$$

This expression is the key to generalize the DPCA technique for multi-class problems. Specifically, each feature i has a vector of weights $\Phi_{i,y}$ with size $|Y|$. So, for each feature i we need to seek for the most important weight, in absolute

Algorithm 1: Multi-Class.M2 DPCA Procedure

Input: Samples: $X = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2) \dots (\mathbf{x}_M, y_M)\}$; where $y_i \in Y$ and $Y = \{1, 2, 3, \dots, N\}$; Distribution D over the M examples; Percentage μ ;

- 1 Initialize the weight vector: $w_{i,y}^1 = \frac{D(i)}{|Y|-1}$, for $i = 1, \dots, M; y \in Y - \{y_i\}$
- 2 Calculate P_{pca} and the projected data $\bar{\mathbf{x}}_i = (P_{pca})^T \tilde{\mathbf{x}}_i$ where $\tilde{\mathbf{x}}_i = x_i - \hat{x}$, with, $\hat{x} = \frac{1}{M} \sum_{i=1}^M x_i$
- 3 Build the labeled projected data set $\bar{\Theta} = \{(\bar{\mathbf{x}}_1, y_1), (\bar{\mathbf{x}}_2, y_2) \dots (\bar{\mathbf{x}}_M, y_M)\}$
- 4 **for** $t = 1, \dots$ **to** T **do**
 - 5 for $y \neq y_i$: $q_{i,y}^t = \frac{w_{i,y}^t}{W_i^t}$; and set $D^t(i) = \frac{W_i^t}{\sum_{i=1}^N W_i^t}$
 - 6 **for** $y = 1, \dots$ **to** N **do**
 - 7 Build the subset $\bar{\Theta}^y$, given by expression (4);
 - 8 $(\phi_y^t, b_y^t) = WSVM(\bar{\Theta}^y, \mathcal{Y}, D^t, \mu)$ where $\mathcal{Y} = \{-1, 1\}$;
 - 9 Get hypothesis $h^t: X \times Y \rightarrow [0, 1]$, given by $h^t(\mathbf{x}, y) = f(\langle \mathbf{x}, \phi_y^t \rangle + b_y^t)$
 - 10 Compute:
$$e^t = \frac{1}{2} \sum_{i=1}^N D^t(i) \left(1 - h^t(\mathbf{x}_i, y_i) + \sum_{y \neq y_i} q_{i,y}^t h^t(\mathbf{x}_i, y) \right)$$
 - 11 **if** $e_t > 0.5$ **then**
 - 12 break;
 - 13 Calculate AdaBoost.M2 weights: $\alpha^t = \frac{1}{2} \ln \left(\frac{1-e^t}{e^t} \right)$;
 - 14 **for** $i = 1, \dots, N$ **and** $y \in Y - \{y_i\}$ **do**
 - 15 Update:
$$w_{i,y}^{t+1} = w_{i,y}^t \exp(-\alpha^t (1 - h^t(x_i, y_i) + h^t(x_i, y)));$$
- 16 Normalize $\tilde{\alpha}^t = \alpha^t / \sum_{j=1}^T \alpha^j, t = 1, 2, \dots, T$
- 17 **for** $i = 1, \dots$ **to** m' **do**
 - 18
$$|\Phi_{i,y}| = \left| \sum_{t=1}^T \tilde{\alpha}^t \frac{\phi_{i,y}^t}{z_{max,y}^t - z_{min,y}^t} \right|, y \in Y \quad (6)$$
 - 19 Compute $v(i) = \max_{y \in Y} \{|\Phi_{i,y}|\}$, $i = 1, 2, \dots, m'$
 - 20 Sort discriminant weights: $v(1) \geq v(2) \geq \dots \geq v(m')$
 - 21 Select the principal components following $v(i)$

Output: Discriminant principal components: $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_{m'}$

Algorithm 2: WSVM Procedure: Build a Weakened version of SVM.

Input: Labeled samples: $X = \{(\mathbf{x}_i, y_i), i = 1, 2, \dots, n'\}$ where $y_i \in Y$ is the label of the sample \mathbf{x}_i ;
Samples probability distribution $D(\mathbf{x}_i)$;
Percentage μ ;
Select \mathcal{J} so that, $\sum_{j \in \mathcal{J}} D(x_j) \leq (1 - \mu)$;
Select $(\mathbf{x}_i, y_i); i \in \mathcal{J}$, and define $D^* = D_{\mathcal{J}}$;
Compute the weighted data $X^* = \{(D_i^* \cdot \mathbf{x}_i, y_i), i \in \mathcal{J}\}$;
Compute the (weak) SVM hyperplane ϕ_{svm} using X^* ;
Output: WSVM hyperplane ϕ_{svm}, b .

value $|\Phi_{i,y_i}|$, which can be interpreted as a measure of the discriminant contribution of the corresponding feature. These values are used to generate the vector \mathbf{v} in line 19 of Algorithm 1. Next, we shall sort the obtained array in decreasing order, as performed in line 20 of the Algorithm 1, to get the global discriminant weights. The output of the Multi-Class.M2 DPCA procedure is the discriminant principal components $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_{m'}$ where \mathbf{q}_i is a PCA component selected according to its discriminant weight $v(i)$.

On the other hand, if we compute the LDA in the PCA space, we get $N - 1$ hyperplane directions $\phi_{lda}^i \in \mathbb{R}^{m'}, i = 1, 2, \dots, (N - 1)$. Consequently, we obtain in this case a LDA weight matrix $\phi_{lda}^{i,j}$, which can be processed according to lines 19-20 of Algorithm 1, by just replacing $\Phi_{i,y}$ by $\phi_{lda}^{i,j}$. The obtained global discriminant weights are named Multi-Class LDA-DPCA in the following sections.

We also aim to study the influence of the denominator $z_{max,y}^t - z_{min,y}^t$ in expression (6) of line 18 of Multi-Class.M2 DPCA algorithm. To perform this task, we test a version of the Algorithm 1 with equation (6) replaced by $|\Phi_{i,y}| = \left| \sum_{t=1}^T \tilde{\alpha}^t \phi_{i,y}^t \right|$. We call the obtained (non normalized) algorithm as the Multi-Class.M2 DPCA-NN.

IV. COMPUTATIONAL EXPERIMENTS

In this section we perform facial expression experiments using the Radboud (RaFD) [17] and the Japanese Female Facial Expression (JAFFE) image databases [18]. In order to save memory allocation along the Algorithm 1 execution, we convert each pose to gray scale and resize it to 50×50 before computation.

In the following, we consider the PCA as well as the Multi-Class DPCA discriminant technique, presented in [7], and the discriminant approaches explained in the section III: Multi-Class.M2 DPCA, Multi-Class.M2 DPCA-NN, and the Multi-Class LDA-DPCA. For evaluation of the discriminant principal components, the following separation tasks have been performed using frontal face images of the mentioned databases (see [14], section 6):

- **Three-Class** experiment: neutral, happiness, and sad samples;
- **Five-Class** experiment: neutral, happiness, sad, fear, and anger classes.

The recognition tasks experiments are carried out using the full rank PCA subspace with all non-zero eigenvalues. In these experiments we have assumed equal prior probabilities and misclassification costs for all the classes. On the PCA subspace, the mean of each class i has been calculated from the corresponding training images and the Mahalanobis distance from each class mean $\hat{\mathbf{x}}_i$ has been used to assign a test observation \mathbf{x}_r to either the different facial expressions. That is, we have assigned \mathbf{x}_r to class i that minimizes:

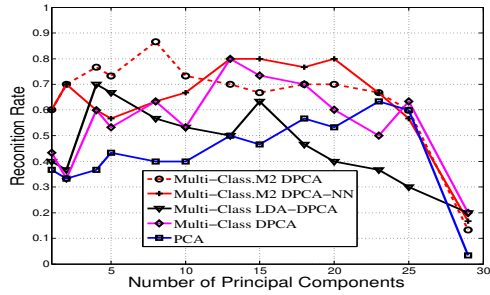
$$d_i(\mathbf{x}_r) = \sum_{j=1}^k \frac{1}{\lambda_j} (x_{rj} - \hat{x}_{ij})^2, \quad (12)$$

where λ_j is the corresponding eigenvalue, k is the number of principal components retained, x_{rj} and \hat{x}_{ij} are the projections of the sample \mathbf{x}_r and of the mean $\hat{\mathbf{x}}_i$, respectively, in the j th component considered.

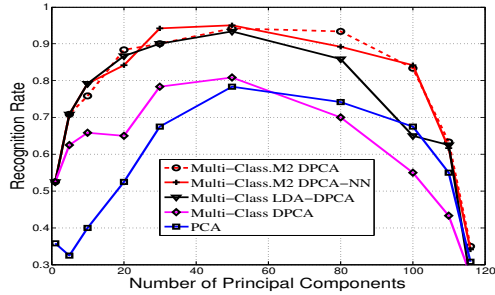
The Figure 3 shows the average recognition rates of the 10-fold cross validation experiments for PCA and the discriminant techniques for the three and five-class classification problems above mentioned. When analysing the Figures 3.(a)-(d) we notice that the Multi-Class.M2 DPCA and Multi-Class.M2 DPCA-NN achieve highest recognition rates or perform closer to the best one. Specifically, let us highlight the intervals where the highest recognition rates are not achieved by Multi-Class.M2 DPCA or Multi-Class.M2 DPCA-NN. This happens in the Figure 3.(a), for $25 < k < 29$, where the Multi-Class DPCA is the best technique. Also, in Figure 3.(b), the Multi-Class LDA-DPCA outperforms both Multi-Class.M2 DPCA and Multi-Class.M2 DPCA-NN in the range $10 \leq k \leq 17$. For $k > 110$ all the methods, except the Multi-Class DPCA, achieves the same accuracy. In Figure 3.(d) shows that in the range $8 \leq k \leq 20$ the Multi-Class LDA-DPCA technique is the best method. However, in all these cases, if we take the absolute value of the difference between the minimum classification rate obtained by Multi-Class.M2 DPCA and Multi-Class.M2 DPCA-NN and the maximum accuracy of the other ones, we get the values 1,4% for $k = 17$ and 3% for $k = 10$, for Figure 3.(b) and Figure 3.(d) respectively, which are not expressive values.

Moreover, although Multi-Class LDA-DPCA and PCA classification rates are equal to the highest accuracy in some intervals of Figures 3.(b)-(d), we must observe that the maxima of the recognition rates are obtained by Multi-Class.M2 DPCA (Figure 3.(a): 86% in $k = 8$; Figure 3.(d): 82% in $k = 120$) and Multi-Class.M2 DPCA-NN (Figure 3.(b): 94% in $k = 30$; Figure 3.(c): 66% in $k = 15$). However, in the case of Figure 3.(c) we shall notice also that Multi-Class.M2 DPCA gets an accuracy very close to the maximum (64%) but using just 10 components. These facts indicates a slight superiority of Multi-Class.M2 DPCA subspaces against the Multi-Class.M2 DPCA-NN ones for expression recognition tasks. Also, it is important to notice in Figures 3.(a)-(d) a degradation in the accuracy of all techniques for higher subspace dimensions, probably due to overfitting.

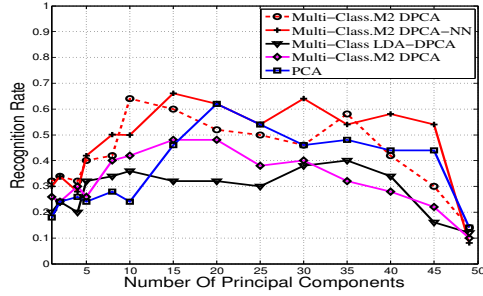
In the results to be presented bellow, we use only the Radboud database because it has gender and expression variations



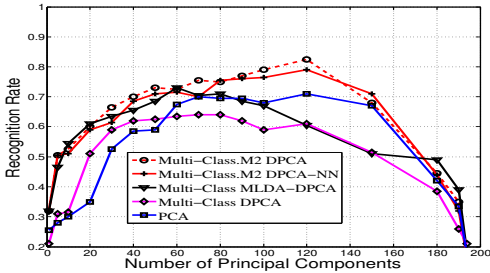
(a)



(b)



(c)



(d)

Fig. 3. Expression experiments using the JAFFE and Radboud databases. Average recognition rates of PCA components selected by the focused techniques:(a) Three-class tasks with JAFFE. (b) Three-class experiments using Radboud. (c) Five-class tasks with JAFFE. (d) Five-class experiments using Radboud.

which allow a more complete analysis of the discriminant components meaning (see [14] for more complete results). Table I, lists the 10 principal components with the highest discriminant weights given by the Multi-Class.M2 DPCA al-

gorithm and the counterparts, for discriminating the expression samples, when considering the three-class experiment with the Radboud database. We can observe that all the considered techniques have selected some distant PCA components among the first 10 most discriminant principal components. In the specific case of the of Multi-Class.M2 DPCA, it selected the 42th, 38th and 81th PCA components among its first 10 discriminant principal components. Since principal components with lower variances describe particular information related to few samples, these results confirm the ability of Multi-Class.M2 DPCA of zooming into the details of group differences. However, components with lower variances should count less for the global reconstruction than PCA components with higher variances. We expect some consequences of this fact in the reconstruction experiments, as we will see next.

Multi-Class.M2 DPCA	24	36	26	32	20	38	81	42	29	9
Multi-Class.M2 DPCA-NN	24	26	36	20	32	8	9	38	17	19
Multi-Class LDA-DPCA	24	26	20	71	81	36	56	19	59	21
Multi-Class DPCA	24	26	20	22	19	17	34	40	8	44

TABLE I

TOP 10 (FROM TOP TO BOTTOM AND LEFT TO RIGHT) DISCRIMINANT PRINCIPAL COMPONENTS, RANKED BY THE DISCRIMINANT TECHNIQUES, USING THE RADBOUD DATABASE FOR THREE-CLASS TASKS.

To understand the changes described by the principal components for the three-class separation tasks with the Radboud data set, we reconstruct one expressive feature by varying a discriminant principal component \mathbf{q}_i using the equation:

$$I = \hat{\mathbf{x}} + \delta \cdot \mathbf{q}_i, \quad (13)$$

where $\hat{\mathbf{x}}$ is the global mean, $\delta \in \{\pm j \cdot \bar{\lambda}^{0.5}, j = 0, \pm 3\}$, and $\bar{\lambda}$ is the average eigenvalue of the total covariance matrix of PCA. We choose $\bar{\lambda}$ instead of λ_i because some λ_i can be very small (or big) in this case, showing no changes (or color saturation) between the samples when we move along the corresponding principal components.

From Table I we notice that the first three columns do not show expressive differences between the selected principal components. Hence, in Figure 4 we illustrate the transformations on the forth PCA most expressive component contrasted with the forth discriminant principal component selected by the discriminant techniques to separate facial expressions.

In Figures 4.(m)-(o), it can be seen that the forth PCA most expressive direction captures essentially the changes in gender, which are the major variations of all the training samples. Owing to the fact that changes in facial expression are much less significant than the gender ones, the standard PCA is unable to capture such minor variations in its most expressive components. However, when we compare these results with the ones reconstructed by the forth most discriminant principal component selected by the other techniques, illustrated by Figures 4.(a)-(l), we can see that more distant principal component (see Table I) carry more information about expression variations than the first PCA ones. That is why the discriminant methods achieves, in general, higher

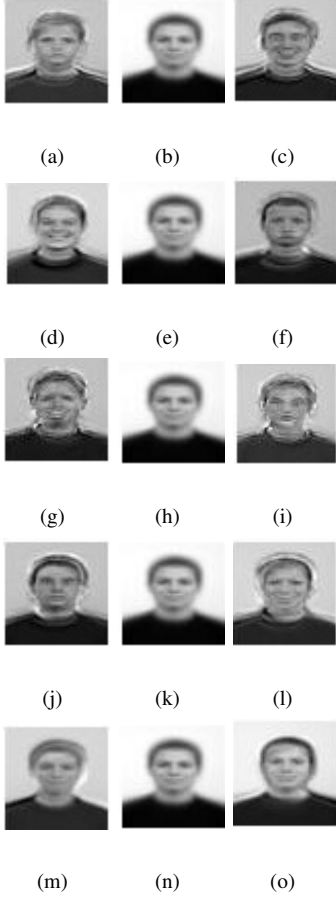


Fig. 4. Visualization of the changes described by the fourth principal direction, using the Radboud face database for three-class experiments, selected by: (a)-(c) Multi-Class.M2 DPCA. (d)-(f) Multi-Class.M2 DPCA-NN. (g)-(i) Multi-Class LDA-DPCA. (j)-(l) Multi-Class DPCA. (m)-(o) PCA.

classification rates than PCA. On the other hand, we can notice in Figure 4.(a)-(f) that the Multi-Class.M2 DPCA and Multi-Class.M2 DPCA-NN fourth discriminant principal components capture more clearly the facial expression with less artifacts in the reconstruction than the other discriminant techniques which agrees with the observed superiority of Multi-Class.M2 DPCA and Multi-Class.M2 DPCA-NN in the recognition experiments. An analogous result is observed for five-class classification tasks [14].

Now, it is worthwhile to consider the accumulated variance, for Radboud experiment, explained by the selected subspaces which is computed by:

$$Vacc^{l,i}(k) = \frac{\sum_{j=1}^k \lambda_j^l}{\sum_{j=1}^{m'} \lambda_j^l}, \quad (14)$$

where $i \in \{3, 5\}$, $l \in \{1, 2, 3, 4, 5\}$ with $l = 1$ corresponds to Multi-Class.M2 DPCA, $l = 2$ to Multi-Class.M2 DPCA-NN, $l = 3$ to Multi-Class LDA-DPCA, and $l = 4, 5$ correspond to Multi-Class DPCA and PCA, respectively. Also, λ_j^l is the variance associated to the j th component selected by the discriminant techniques l . The expression (14) is important for this discussions because we expect some correlation between the performance for reconstruction and the accumulated

variance in expression (14) due to the known fact that the components with larger variances keep global information related to features that most vary in the samples [6]. In Figure 5 we plot the result of expression (14).

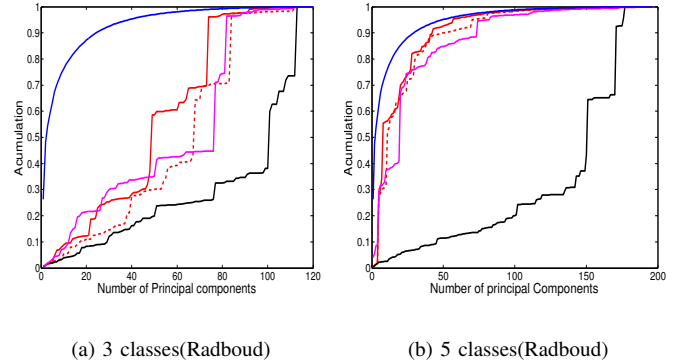


Fig. 5. Total variance computed by expression (14), explained by discriminant principal components selected by Multi-Class.M2 DPCA (dashed red line); Multi-Class.M2 DPCA-NN (solid red line); Multi-Class LDA-DPCA (black line); Multi-Class DPCA (magenta line); and PCA (blue line) using Radboud database.

We observe that the Multi-Class LDA-DPCA gives the lowest results for the $Vacc$ while PCA technique gives the largest accumulated variances for both three and five-class Radboud experiments. The $Vacc$ of the other considered techniques fall between Multi-Class LDA-DPCA and PCA in both Figures 5.(a),(b). Although PCA gives the largest values for $Vacc$ its recognition rates are, in general, outperformed by the discriminant principal components. Since PCA explains features that most vary in the samples the principal subspaces do not necessarily represent important discriminant directions to separate sample groups.

However, the reconstruction results are expected to give lower errors if we take components with higher variances, like PCA does. Besides, in the case of five classes, we expect a better reconstruction performance for Multi-Class.M2 DPCA and Multi-Class.M2 DPCA-NN once their accumulated variances are closer to the PCA ones in this case. To make clear these observations, let us quantify the reconstruction quality through the root mean squared error (RMSE), computed as follows:

$$RMSE^{l,i}(k) = \sqrt{\frac{\sum_{j=1}^M \|P \cdot I_k^l \cdot P^T \mathbf{x}_j - \mathbf{x}_j\|^2}{M}}, \quad (15)$$

where the index l and i follows the same maps used in expression (14), I_k^l is a truncated identity matrix that keeps the selected subspace with dimension k , $P = P_{pca}$, and $\|\cdot\|$ is the usual 2-norm.

In Figure 6 shows the RMSE for the reconstruction process for the subspaces given by the focused techniques. It is noticeable that for all experiments, PCA reconstruction performs equal or better than the multi-class discriminant components for all the simulated values of k . This observation agrees with the accumulated variance reported in Figure 5 which is such that $Vacc^{5,i}(k) \geq Vacc^{l,i}(k)$ for all values of k , l and $i = 3, 5$. Therefore, while in the classification tasks the PCA

method is, in general, outperformed by the Multi-Class DPCA method, in the reconstruction experiments the PCA subspaces become more efficient for almost all the simulated values.

Let us compare the reconstruction performance of the discriminant approaches. From Figure 6.(a), for $1 \leq k \leq 5$ the discriminant techniques performs equal to each other. For $5 \leq k \leq 116$ the RMSE for Multi-Class LDA-DPCA is the largest one. The other discriminant technique performs as: (a) For $5 \leq k \leq 42$, $RMSE^{1,3}(k) > RMSE^{2,3}(k) > RMSE^{4,3}(k)$; (b) For $42 \leq k \leq 116$, $RMSE^{1,3}(k) \geq RMSE^{4,3}(k) \geq RMSE^{2,3}(k)$. Therefore, the Multi-Class.M2 DPCA is equal or worst than the Multi-Class DPCA and Multi-Class.M2 DPCA-NN in terms of RMSE results. On the other hand, the Multi-Class.M2 DPCA-NN performs better than the Multi-Class DPCA and Multi-Class.M2 DPCA for $42 \leq k \leq 116$. This observations are in accordance with the corresponding accumulated variances in Figure 5.(a).

In the case of the RMSE for the five-class experiments with the Radboud, shown in Figure 6.(b), we observe an analogous behaviour for PCA and Multi-Class LDA-DPCA for $1 \leq k \leq 180$. All the other discriminant methods perform equal for $1 \leq k \leq 5$. For $5 \leq k \leq 21$ we notice that $RMSE^{2,5}(k) \leq RMSE^{4,5}(k) \leq RMSE^{1,5}(k)$. Next, in the range $21 \leq k \leq 194$ we have $RMSE^{2,5}(k) \leq RMSE^{1,5}(k) \leq RMSE^{4,5}(k)$. The reported behaviors agree with the accumulated variances in Figure 5.(b).

Therefore, in terms of reconstruction, the Multi-Class.M2 DPCA-NN is better or equal than the Multi-Class.M2 DPCA in both three and five-class experiments. The Multi-Class DPCA performs better than the Multi-Class.M2 DPCA-NN only for $5 \leq k \leq 42$ and three-class problems. The PCA is the best technique and the Multi-Class LDA-DPCA is the worse one.

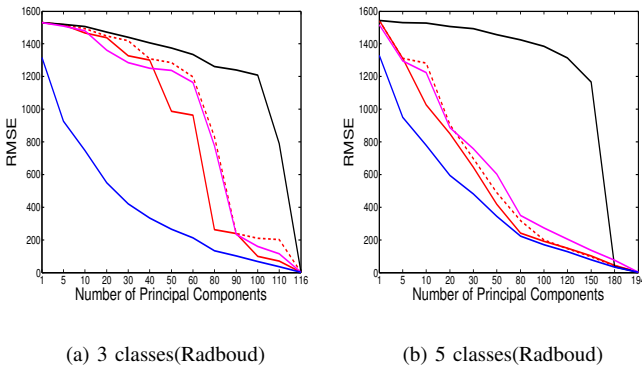


Fig. 6. RMSE computed by equation (15) for: Multi-Class.M2 DPCA (dashed red line); Multi-Class.M2 DPCA-NN (solid red line); Multi-Class LDA-DPCA (black line); Multi-Class DPCA (magenta line); and PCA (blue line).

V. CONCLUSION AND FUTURE WORKS

This paper introduces the Multi-Class.M2 DPCA and Multi-Class.M2 DPCA-NN algorithms for ranking PCA components computed from multi-class facial expression databases. The basic methodology has a computational complexity dominated by the AdaBoost.M2 algorithm plus PCA computation. The

facial expressions experiments show that, in general, the PCA components selected by Multi-Class.M2 DPCA and Multi-Class.M2 DPCA-NN allow higher recognition rates using less linear features than Multi-Class LDA-DPCA, Multi-Class DPCA and the standard PCA.

Further work is being undertaken to test the algorithm for more than 5 classes as well as with other applications. We shall replace the AdaBoost.M2 technique by the bagging one [11] as a direction to improve the classification performance when increasing the number of classes. Moreover, we need to improve its reconstruction in low dimensional subspaces.

REFERENCES

- [1] T. Hastie, R. Tibshirani, and J. Friedman, "The elements of statistical learning," *Springer*, 2001.
- [2] G. A. Giraldo, P. S. Rodrigues, E. C. Kitani, and C. E. Thomaz, "Dimensionality reduction, classification and reconstruction problems in statistical learning approaches," *Revista de Informatica Teorica e Aplicada (RITA)*, vol. 15, no. 1, pp. 141–173, 2008.
- [3] J. P. Cunningham and Z. Ghahramani, "Linear dimensionality reduction: Survey, insights, and generalizations," *J. of Mach. Learn. Research*, vol. 16, pp. 2859–2900, 2015.
- [4] H. Safavi and C.-I. Chang, "Projection pursuit-based dimensionality reduction," *Proc. SPIE*, vol. 6966, pp. 69 661H–69 661H–11, 2008.
- [5] C. E. Thomaz and G. A. Giraldo, "A new ranking method for principal components analysis and its application to face image analysis," *Image Vision Comput.*, vol. 28, no. 6, pp. 902–913, June 2010.
- [6] D. Swets and J. Weng, "Using discriminants eigenfeatures for image retrieval," *IEEE Trans. Patterns Anal. Mach. Intell.*, vol. 18(8), pp. 831–836, 1996.
- [7] T. Filisbino, D. Leite, G. Giraldo, and C. Thomaz, "Multi-class discriminant analysis based on svm ensembles for ranking principal components," in *36th Ibero-Latin Am. Cong. on Comp. Meth. in Eng. (CILAMCE)*, Nov 2015.
- [8] M. Zhu and A. M. Martinez, "Selecting principal components in a two-stage lda algorithm," in *CVPR'06*, June 2006, pp. 132–137.
- [9] V. N. Vapnik, *Statistical Learning Theory*. John Wiley & Sons, INC., 1998.
- [10] M. Loog, R. P. W. Duin, and R. Haeb-Umbach, "Multiclass linear dimension reduction by weighted pairwise fisher criteria," *IEEE Trans. Patterns Anal. Mach. Intell.*, vol. 23(7), pp. 762–766, 2001.
- [11] Z.-H. Zhou, *Ensemble Methods: Foundations and Algorithms*, 1st ed. Chapman & Hall/CRC, 2012.
- [12] E. Garcia and F. Lozano, "Boosting Support Vector Machines," in *Proceedings of International Conference of Machine Learning and Data Mining (MLDM'2007)*. Leipzig, Germany: Ibal publishing, Jul. 2007, pp. 153–167.
- [13] Y. Freund and R. E. Schapire, "A decision-theoretic generalization of on-line learning and an application to boosting," *Journal of computer and system sciences*, vol. 55, no. 1, pp. 119–139, 1997.
- [14] T. Filisbino, G. Giraldo, and C. Thomaz, "Multi-class discriminant analysis based on linear classifiers and adaboost.m2 for ranking principal components in face spaces," Tech. Rep., 2016. [Online]. Available: http://www.lncc.br/departamentos/producaocientificageral.php?vMenu=2&vTipo=13&vCabecalho=pesq&vTitulo=lncc&vDepto=&idt_responsavel=&vAno=2016&ano=2016&anof=2016&idt_linha_pesquisa=
- [15] M. Turk and A. Pentland, "Eigenfaces for recognition," *Journal of Cognitive Neuroscience*, vol. 3, pp. 71–86, 1991.
- [16] E. Yildizer, A. M. Balci, M. Hassan, and R. Alhaji, "Efficient content-based image retrieval using multiple support vector machines ensemble," *Expert Syst. Appl.*, vol. 39, pp. 2385–2396, Feb. 2012.
- [17] O. Langner, R. Dotsch, G. Bijlstra, D. H. J. Wigboldus, S. T. Hawk, and A. van Knippenberg, "Presentation and validation of the radboud faces database," *Cognition & Emotion*, vol. 24, no. 8, pp. 1377–1388, 2010.
- [18] M. J. Lyons, S. Akamatsu, M. Kamachi, J. Gyoba, and J. Budynek, "The japanese female facial expression (jaffe) database," 1998.