# Face Recognition Using Kernel Direct Discriminant Analysis Algorithms

Juwei Lu, K.N. Plataniotis, A.N. Venetsanopoulos

Bell Canada Multimedia Laboratory, The Edward S. Rogers Sr.

Department of Electrical and Computer Engineering

University of Toronto, Toronto, M5S 3G4, ONTARIO, CANADA

PUBLICATION FORMAT: REGULAR PAPER

AREA: IMAGE PROCESSING AND RECOGNITION

CORRESPONDENCE ADDRESS:

Prof. K.N. Plataniotis

Bell Canada Multimedia Laboratory

The Edward S. Rogers Sr. Department of Electrical and Computer Engineering

University of Toronto

10 King's College Road

Toronto, Ontario M5S 3G4, Canada

Tel: (416) 946-5605

Fax: (416) 978-4425

E-mail: kostas@dsp.toronto.edu

## Abstract

Techniques that can introduce low-dimensional feature representation with enhanced discriminatory power is of paramount importance in face recognition (FR) systems. It is well known that the distribution of face images, under a perceivable variation in viewpoint, illumination or facial expression, is highly nonlinear and complex. It is therefore not surprising that linear techniques, such as those based on Principle Component Analysis (PCA) or Linear Discriminant Analysis (LDA), cannot provide reliable and robust solutions to those FR problems with complex face variations. In this paper, we propose a kernel machine based discriminant analysis method, which deals with the nonlinearity of the face patterns' distribution. The proposed method also effectively solves the so-called "small sample size" (SSS) problem which exists in most FR tasks. The new algorithm has been tested, in terms of classification error rate performance, on the multi-view UMIST face database. Results indicate that the proposed methodology is able to achieve excellent performance with only a very small set of features being used, and its error rate is approximately 34% and 48% of those of two other commonly used kernel FR approaches, the Kernel-PCA (KPCA) and the Generalized Discriminant Analysis (GDA) respectively.

## Keywords

4

# I. INTRODUCTION

Within the last decade, face recognition (FR) has found a wide range of applications, from identity authentication, access control, and face-based video indexing/browsing, to human-computer interaction/communication. As a result, numerous FR algorithms have been proposed, and surveys in this area can be found in [1], [2], [3], [4], [5]. Two issues are central to all these algorithms: (i) feature selection for face representation, and (ii) classification of a new face image based on the chosen feature representation [6]. This work focuses on the issue of feature selection. The main objective is to find techniques that can introduce low-dimensional feature representation of face objects with enhanced discriminatory power. Among various solutions to the problem, the most successful are those appearance-based approaches, which generally operate directly on images or appearances of face objects and process the images as 2D holistic patterns, to avoid difficulties associated with 3D modeling, and shape or landmark detection [5].

Principle Component Analysis (PCA) and Linear Discriminant Analysis (LDA) are two classic tools widely used in the appearance-based approaches for data reduction and feature extraction. Many state-of-the-art FR methods, such as Eigenfaces [7] and Fisherfaces [8], are built on these two techniques or their variants. It is generally believed that when it comes to solving problems of pattern classification, LDA based algorithms outperform PCA based ones, since the former optimizes the low-dimensional representation of the objects with focus on the most discriminant feature extraction while the latter achieves simply object reconstruction. However, many LDA based algorithms suffer from the so-called *"small sample size problem"* (SSS) which exists in high-dimensional pattern recognition tasks where the number of available samples is smaller than the dimensionality of the samples. The traditional solution to the SSS problem is to utilize PCA concepts in conjunction with LDA (PCA+LDA) as it was done for example in Fisherfaces [8]. Recently, more effective solutions, called Direct LDA (D-LDA) methods, have been presented [9], [10]. Although successful in many cases, linear methods fail to deliver good performance when face patterns are subject to large variations in viewpoints, which results in a highly non-convex and complex distribution. The limited success of these methods should be attributed to their linear nature [11]. As a result, it is reasonable to assume that a better

solution to this inherent nonlinear problem could be achieved using nonlinear methods, such as the so-called kernel machine techniques [12], [13], [14], [15].

In this paper, motivated by the success that Support Vector Machines (SVM) [16], [17], [18], Kernel PCA (KPCA) [19] and Generalized Discriminant Analysis (GDA) [20] have in pattern regression and classification tasks, we propose a new kernel discriminant analysis algorithm for face recognition. The algorithm generalizes the strengths of the recently presented D-LDA and the kernel techniques while at the same time overcomes many of their shortcomings and limitations. Therefore, the proposed algorithm can be seen as an enhanced kernel D-LDA method (hereafter KDDA). Following the SVM paradigm, we first nonlinearly map the original input space to an implicit high-dimensional feature space, where the distribution of face patterns is hoped to be linearized and simplified. Then, a new variant of the D-LDA method is introduced to effectively solve the SSS problem and derive a set of optimal discriminant basis vectors in the feature space.

The rest of the paper is organized as follows. Since KDDA is built on D-LDA and GDA, in Section II, we start the analysis by briefly reviewing the two latter methods. Following that, KDDA is introduced and analyzed. The relationship of KDDA to D-LDA and GDA is also discussed. In Section III, two sets of experiments are presented to demonstrate the effectiveness of the KDDA algorithm on highly nonlinear, highly complex face pattern distributions. KDDA is compared, in terms of the classification error rate performance, to KPCA and GDA on the multi-view UMIST face database. Conclusions are summarized in Section IV.

## II. Methods

The problem to be solved is formally stated as follows: A set of $L$ training face images $\{\mathbf{z}_i\}_{i=1}^{L}$ is available. Each image is defined as a vector of length $N(= I_w \times I_h)$, *i.e.* $\mathbf{z}_i \in \mathbb{R}^N$, where $I_w \times I_h$ is the face image size and $\mathbb{R}^N$ denotes a $N$-dimensional real space. It is further assumed that each image belongs to one of $C$ classes $\{\mathbf{Z}_i\}_{i=1}^{C}$. The objective is to find a transformation $\varphi$, based on optimization of certain separability criteria, which produces a mapping $\mathbf{y}_i = \varphi(\mathbf{z}_i)$, with $\mathbf{y}_i \in \mathbb{R}^M$ that leads to an enhanced separability of different face objects.

## A. Generalized Discriminant Analysis (GDA)

For solving nonlinear problems, the classic LDA has been generalized to its kernel version, namely GDA [20]. Let $\phi : \mathbf{z} \in \mathbb{R}^N \rightarrow \phi(\mathbf{z}) \in \mathbb{F}$ be a nonlinear mapping from the input space to a high-dimensional feature space $\mathbb{F}$, where different classes of objects are supposed to be linearly separable. The idea behind GDA is to perform a classic LDA in the feature space $\mathbb{F}$ instead of the input space $\mathbb{R}^N$.

Let $\mathbf{S}_{BTW}$ and $\mathbf{S}_{WTH}$ be the between- and within-class scatter matrices in the feature space $\mathbb{F}$ respectively, expressed as follows:

$$\mathbf{S}_{BTW} = \frac{1}{L} \sum_{i=1}^{C} C_i(\bar{\phi}_i - \bar{\phi})(\bar{\phi}_i - \bar{\phi})^T \tag{1}$$

$$\mathbf{S}_{WTH} = \frac{1}{L} \sum_{i=1}^{C} \sum_{j=1}^{C_i} (\phi_{ij} - \bar{\phi}_i)(\phi_{ij} - \bar{\phi}_i)^T \tag{2}$$

where $\phi_{ij} = \phi(\mathbf{z}_{ij})$, $\bar{\phi}_i = \frac{1}{C_i} \sum_{j=1}^{C_i} \phi(\mathbf{z}_{ij})$ is the mean of class $\mathbf{Z}_i$, $\bar{\phi} = \frac{1}{L} \sum_{i=1}^{C} \sum_{j=1}^{C_i} \phi(\mathbf{z}_{ij})$ is the average of the ensemble, and $C_i$ is the element number in $\mathbf{Z}_i$, which leads to $L = \sum_{i=1}^{C} C_i$. LDA determines a set of optimal discriminant basis vectors, denoted by $\{\psi_k\}_{k=1}^{M}$, so that the ratio of the between- and within-class scatters is maximized [21]. Assuming $\Psi = [\psi_1, \ldots, \psi_M]$, the maximization can be achieved by solving the following eigenvalue problem,

$$\Psi = \arg\max_{\Psi} \frac{\left|(\Psi^T \mathbf{S}_{BTW} \Psi)\right|}{\left|(\Psi^T \mathbf{S}_{WTH} \Psi)\right|} \tag{3}$$

The feature space $\mathbb{F}$ could be considered as a "linearization space" [22], however, its dimensionality could be arbitrarily large, and possibly infinite. Fortunately, the exact $\phi(\mathbf{z})$ is not needed and the feature space can become implicit by using kernel methods, where dot products in $\mathbb{F}$ are replaced with a kernel function in the input space $\mathbb{R}^N$ so that the nonlinear mapping is performed implicitly in $\mathbb{R}^N$ [23], [24].

In FR tasks, the number of training samples, $L$, is in most cases much smaller than the dimensionality of $\mathbb{R}^N$ (for LDA) or $\mathbb{F}$ (for GDA) leading to a degenerated scatter matrix $\mathbf{S}_{WTH}$. Traditional methods, for example GDA and Fisherfaces [8], attempt to solve the so-called SSS problem by using techniques such as pseudo inverse or PCA to remove the null space of $\mathbf{S}_{WTH}$. However, it has been recently shown that the null space may contain the most significant discriminant information [9], [10].

*B. Direct LDA (D-LDA)*

Recently, Chen et al. [9] and Yang et al. [10] proposed the so-called direct LDA (D-LDA) algorithm that attempts to avoid the shortcomings existing in traditional solutions to the SSS problem. The basic idea behind the algorithm is that the null space of $\mathbf{S}_{WTH}$ may contain significant discriminant information if the projection of $\mathbf{S}_{BTW}$ is not zero in that direction, and that no significant information will be lost if the null space of $\mathbf{S}_{BTW}$ is discarded. Assuming, for example, that $\mathcal{A}$ and $\mathcal{B}$ represent the null spaces of $\mathbf{S}_{BTW}$ and $\mathbf{S}_{WTH}$ respectively, the complement spaces of $\mathcal{A}$ and $\mathcal{B}$ can be written as $\mathcal{A}' = \mathbb{R}^N - \mathcal{A}$ and $\mathcal{B}' = \mathbb{R}^N - \mathcal{B}$. Therefore, the optimal discriminant subspace sought by the D-LDA algorithm is the intersection space $(\mathcal{A}' \cap \mathcal{B})$.

The difference between Chen's method [9] and Yang's method [10] is that Yang's method firstly diagonalizes $\mathbf{S}_{BTW}$ to find $\mathcal{A}'$ when seek solution of Eq.3, while Chen's method firstly diagonalizes $\mathbf{S}_{WTH}$ to find $\mathcal{B}$. Although there is no significant difference between the two approaches, it may be intractable to calculate $\mathcal{B}$ when the size of $\mathbf{S}_{WTH}$ is large, which is the case in most FR applications. For example, the size of $\mathbf{S}_{WTH}$ and $\mathbf{S}_{BTW}$ amounts to $10304 \times 10304$ for face images of size $112 \times 92$ such as those used in our experiments. Fortunately, the rank of $\mathbf{S}_{BTW}$ is determined by $rank(\mathbf{S}_{BTW}) = min(N, C - 1)$, with $C$ the number of image classes, usually a small value in most of FR tasks, *e.g.* $C = 20$ in our experiments, resulting in $rank(\mathbf{S}_{BTW}) = 19$. $\mathcal{A}'$ can be easily found by solving eigenvectors of a $19 \times 19$ matrix rather than the original $10304 \times 10304$ matrix through the algebraic transformation proposed in [7]. The intersection space $(\mathcal{A}' \cap \mathcal{B})$ can be obtained by solving the null space of projection of $\mathbf{S}_{WTH}$ into $\mathcal{A}'$, with the projection being a small matrix of size $19 \times 19$. For the reasons explained above, we proceed by firstly diagonalizing the matrix $\mathbf{S}_{BTW}$ instead of $\mathbf{S}_{WTH}$ in the derivation of the proposed here algorithm.

*C. Kernel Direct Discriminant Analysis (KDDA).*

C.1 Eigen-analysis of $\mathbf{S}_{BTW}$ in the Feature Space

Following the general D-LDA framework, we start by solving the eigenvalue problem of $\mathbf{S}_{BTW}$, which can be rewritten here as follows,

$$\mathbf{S}_{BTW} = \sum_{i=1}^{C} \left( \sqrt{\frac{C_i}{L}} \left( \bar{\phi}_i - \bar{\phi} \right) \right) \left( \sqrt{\frac{C_i}{L}} \left( \bar{\phi}_i - \bar{\phi} \right) \right)^T = \sum_{i=1}^{C} \tilde{\bar{\phi}}_i \tilde{\bar{\phi}}_i^T = \Phi_b \Phi_b^T \qquad (4)$$

where $\tilde{\bar{\phi}}_i = \sqrt{\frac{C_i}{L}} \left( \bar{\phi}_i - \bar{\phi} \right)$, and $\Phi_b = \left[ \tilde{\bar{\phi}}_1 \cdots \tilde{\bar{\phi}}_c \right]$. Since the dimensionality of the feature space $\mathbb{F}$, denoted as $N'$, could be arbitrarily large or possibly infinite, it is intractable to directly compute the eigenvectors of the $(N' \times N')$ matrix $\mathbf{S}_{BTW}$. Fortunately, the first $m$ ($\le C-1$) most significant eigenvectors of $\mathbf{S}_{BTW}$, which correspond to non-zero eigenvalues, can be indirectly derived from the eigenvectors of the matrix $\Phi_b^T \Phi_b$ (with size $C \times C$) [7].

Computing $\Phi_b^T \Phi_b$, requires dot product evaluation in $\mathbb{F}$. This can be done in a manner similar to the one used in SVM, KPCA and GDA by utilizing kernel methods. For any $\phi(\mathbf{z}_i), \phi(\mathbf{z}_j) \in \mathbb{F}$, we assume that there exists a kernel function $k(\cdot)$ such that $k(\mathbf{z}_i, \mathbf{z}_j) = \phi(\mathbf{z}_i) \cdot \phi(\mathbf{z}_j)$. The introduction of the kernel function allows us to avoid the explicit evaluation of the mapping. Any function satisfying Mercer's condition can be used as a kernel, and typical kernel functions include polynomial function, radial basis function (RBF) and multi-layer perceptrons [17].

Using the kernel function, for two arbitrary classes $\mathbf{Z}_l$ and $\mathbf{Z}_h$, a $C_l \times C_h$ dot product matrix $K_{lh}$ can be defined as:

$$K_{lh} = (k_{ij})_{\substack{i=1,\cdots,C_l \\ j=1,\cdots,C_h}} , \quad where \quad k_{ij} = k(\mathbf{z}_{li}, \mathbf{z}_{hj}) = \phi_{li} \cdot \phi_{hj} \qquad (5)$$

For all of $C$ classes $\{\mathbf{Z}_i\}_{i=1}^{C}$, we then define a $L \times L$ kernel matrix $\mathbf{K}$,

$$\mathbf{K} = (K_{lh})_{\substack{l=1,\cdots,C \\ h=1,\cdots,C}} \qquad (6)$$

which allows us to express $\Phi_b^T \Phi_b$ as follows:

$$\Phi_b^T \Phi_b = \quad \frac{1}{L} \mathbf{B} \cdot (\mathbf{A}_{LC}^T \cdot \mathbf{K} \cdot \mathbf{A}_{LC} - \frac{1}{L}(\mathbf{A}_{LC}^T \cdot \mathbf{K} \cdot \mathbf{1}_{LC}) - \\ \frac{1}{L}(\mathbf{1}_{LC}^T \cdot \mathbf{K} \cdot \mathbf{A}_{LC}) + \frac{1}{L^2}(\mathbf{1}_{LC}^T \cdot \mathbf{K} \cdot \mathbf{1}_{LC})) \cdot \mathbf{B} \qquad (7)$$

where $\mathbf{B} = \mathbf{diag}\left[\sqrt{C_1} \cdots \sqrt{C_c}\right]$, $\mathbf{1}_{LC}$ is a $L \times C$ matrix with terms all equal to: one, $\mathbf{A}_{LC} = \mathbf{diag}\left[\mathbf{a}_{c_1} \cdots \mathbf{a}_{c_c}\right]$ is a $L \times C$ block diagonal matrix, and $\mathbf{a}_{c_i}$ is a $C_i \times 1$ vector with all terms equal to: $\frac{1}{C_i}$ (see Appendix I for a detailed derivation of Eq.7.).

Let $\lambda_i$ and $\mathbf{e}_i$ $(i = 1 \cdots C)$, be the $i$-th eigenvalue and corresponding eigenvector of $\Phi_b^T \Phi_b$, sorted in **decreasing** order of eigenvalues. Since $(\Phi_b \Phi_b^T)(\Phi_b \mathbf{e}_i) = \lambda_i (\Phi_b \mathbf{e}_i)$, $\mathbf{v}_i = \Phi_b \mathbf{e}_i$ is the eigenvector of $\mathbf{S}_{BTW}$. In order to remove the null space of $\mathbf{S}_{BTW}$, we only use its first $m$ $(\leq C - 1)$ eigenvectors: $\mathbf{V} = [\mathbf{v}_1 \cdots \mathbf{v}_m] = \Phi_b \mathbf{E}_m$ where $\mathbf{E}_m = [\mathbf{e}_1 \ldots \mathbf{e}_m]$, whose corresponding eigenvalues are greater than 0. It is not difficult to see that $\mathbf{V}^T \mathbf{S}_{BTW} \mathbf{V} = \Lambda_b$, with $\Lambda_b = \mathbf{diag}[\lambda_1^2 \cdots \lambda_m^2]$, a $m \times m$ diagonal matrix.

## C.2 Eigen-analysis of $\mathbf{S}_{WTH}$ in the Feature Space

Let $\mathbf{U} = \mathbf{V}\Lambda_b^{-1/2}$. Projecting $\mathbf{S}_{BTW}$ and $\mathbf{S}_{WTH}$ into the subspace spanned by $\mathbf{U}$, it can easily be seen that $\mathbf{U}^T \mathbf{S}_{BTW} \mathbf{U} = \mathbf{I}$ while $\mathbf{U}^T \mathbf{S}_{WTH} \mathbf{U}$ can be expanded as:

$$\mathbf{U}^T \mathbf{S}_{WTH} \mathbf{U} = (\mathbf{E}_m \Lambda_b^{-1/2})^T (\Phi_b^T \mathbf{S}_{WTH} \Phi_b)(\mathbf{E}_m \Lambda_b^{-1/2}) \tag{8}$$

Using the kernel matrix $\mathbf{K}$, a closed form expression of $\Phi_b^T \mathbf{S}_{WTH} \Phi_b$ can be obtained as follows,

$$\Phi_b^T \mathbf{S}_{WTH} \Phi_b = \frac{1}{L} \left(\mathbf{J}1 - \mathbf{J}2\right) \tag{9}$$

with $\mathbf{J}1$ and $\mathbf{J}2$ defined in Appendix II along with the detailed derivation of the expression in Eq.9.

We proceed by diagonalizing $\mathbf{U}^T \mathbf{S}_{WTH} \mathbf{U}$, a tractable matrix with size $m \times m$. Let $\mathbf{p}_i$ be the $i$-th eigenvector of $\mathbf{U}^T \mathbf{S}_{WTH} \mathbf{U}$, where $i = 1 \cdots m$, sorted in **increasing** order of the corresponding eigenvalue $\lambda_i'$. In the set of ordered eigenvectors, those that correspond to the smallest eigenvalues maximize the ratio in Eq.3, and should be considered the most discriminative features. Discarding the eigenvectors with the largest eigenvalues, the $M(\leq m)$ selected eigenvectors are denoted as $\mathbf{P} = [\mathbf{p}_1 \cdots \mathbf{p}_M]$. Defining a matrix $\mathbf{Q} = \mathbf{UP}$, we can obtain $\mathbf{Q}^T \mathbf{S}_{WTH} \mathbf{Q} = \Lambda_w$, with $\Lambda_w = \mathbf{diag}[\lambda_1' \cdots \lambda_M']$, a $M \times M$ diagonal matrix.

Based on the calculations presented above, a set of optimal discriminant feature vectors can be derived through $\Gamma = \mathbf{Q}\Lambda_w^{-1/2}$. The features form a low-dimensional subspace in $\mathbb{F}$, where the ratio in Eq.3 is maximized. Similar to the D-LDA framework, the subspace

obtained contains the intersection space $(\mathcal{A}' \cap \mathcal{B})$ shown in section II-B. However, it is possible that there exist eigenvalues with $\lambda_i' = 0$ in $\Lambda_w$. To alleviate the problem, threshold values were introduced in [10], where any value below the threshold $\epsilon$ is promoted to $\epsilon$ (a very small value). Obviously, performance heavily depends on the heuristic evaluation of the parameter $\epsilon$.

To robustify the approach, we propose a modified Fisher's criterion to be used instead of the conventional definition in Eq.3 when $\mathbf{U}^T \mathbf{S}_{WTH} \mathbf{U}$ is singular. The new criterion can be expressed as:

$$\Psi = \arg\max_{\Psi} \frac{\left|(\Psi^T \mathbf{S}_{BTW} \Psi)\right|}{\left|(\Psi^T \mathbf{S}_{BTW} \Psi) + (\Psi^T \mathbf{S}_{WTH} \Psi)\right|} \tag{10}$$

The modified Fisher's criterion of Eq.10 has been proved to be equivalent to the conventional one (Eq.3) in [25]. The expression $\mathbf{U}^T(\mathbf{S}_{BTW} + \mathbf{S}_{WTH})\mathbf{U}$ which is used in Eq.10 instead of the $\mathbf{U}^T \mathbf{S}_{WTH} \mathbf{U}$ can be shown to be non-singular by the following lemma.

*Lemma 1:* Suppose $\mathbf{D}$ is a real matrix of size $\mathcal{N} \times \mathcal{N}$, and can be represented by $\mathbf{D} = \Phi \Phi^T$ where $\Phi$ is a real matrix of size $\mathcal{N} \times \mathcal{M}$. Then, $(\mathbf{I} + \mathbf{D})$ is positive definite, *i.e.* $\mathbf{I} + \mathbf{D} > 0$, where $\mathbf{I}$ is a $\mathcal{N} \times \mathcal{N}$ identity matrix.

*Proof:* Since $\mathbf{D}^T = \mathbf{D}$, $(\mathbf{I} + \mathbf{D})$ is a real symmetric matrix. For any $\mathcal{N} \times 1$ non-zero real vector: $x$, $x^T(\mathbf{I} + \mathbf{D})x = x^T x + x^T \mathbf{D} x = x^T x + (\Phi^T x)^T(\Phi^T x) > 0$. According to [26], the matrix $(\mathbf{I} + \mathbf{D})$ that satisfies the above conditions is positive definite. ∎

Following a procedure similar to $\mathbf{S}_{BTW}$, $\mathbf{S}_{WTH}$ can be expressed as $\mathbf{S}_{WTH} = \Phi_w \Phi_w^T$, with $\mathbf{U}^T \mathbf{S}_{WTH} \mathbf{U} = (\mathbf{U}^T \Phi_w)(\mathbf{U}^T \Phi_w)^T$. Since $\mathbf{U}^T \mathbf{S}_{BTW} \mathbf{U} = \mathbf{I}$ and $(\mathbf{U}^T \mathbf{S}_{WTH} \mathbf{U})$ satisfies the conditions on $\mathbf{D}$ discussed in Lemma 1, $\mathbf{U}^T(\mathbf{S}_{BTW} + \mathbf{S}_{WTH})\mathbf{U}$ is positive definite. As a result, $\mathbf{Q}^T(\mathbf{S}_{BTW} + \mathbf{S}_{WTH})\mathbf{Q} = \Lambda_w$ is non-singular.

### C.3 Dimensionality Reduction and Feature Extraction

For any input pattern $\mathbf{z}$, its projection into the set of feature vectors, $\Gamma$, derived in Section II-C.2, can be calculated by

$$\mathbf{y} = \Gamma^T \phi(\mathbf{z}) = \left( \mathbf{E}_m \cdot \Lambda_b^{-1/2} \cdot \mathbf{P} \cdot \Lambda_w^{-1/2} \right)^T \left( \Phi_b^T \phi(\mathbf{z}) \right) \tag{11}$$

where $\Phi_b^T \phi(\mathbf{z}) = [\, \tilde{\bar{\phi}}_1 \quad \cdots \quad \tilde{\bar{\phi}}_c \,]^T \phi(\mathbf{z})$. Since

$$\tilde{\bar{\phi}}_i^T \phi(\mathbf{z}) = \left( \sqrt{\frac{C_i}{L}} \left( \bar{\phi}_i - \bar{\phi} \right) \right)^T \phi(\mathbf{z}) = \sqrt{\frac{C_i}{L}} \left( \frac{1}{C_i} \sum_{m=1}^{C_i} \phi_{im}^T \phi(\mathbf{z}) - \frac{1}{L} \sum_{p=1}^{C} \sum_{q=1}^{C_p} \phi_{pq}^T \phi(\mathbf{z}) \right) \tag{12}$$

we have

$$\Phi_b^T \phi(\mathbf{z}) = \frac{1}{\sqrt{L}} \mathbf{B} \cdot \left( \mathbf{A}_{LC}^T \cdot \gamma(\phi(\mathbf{z})) - \frac{1}{L} \mathbf{1}_{LC}^T \cdot \gamma(\phi(\mathbf{z})) \right) \tag{13}$$

where $\gamma(\phi(\mathbf{z})) = [\, \phi_{11}^T \phi(\mathbf{z}) \quad \phi_{12}^T \phi(\mathbf{z}) \quad \cdots \quad \phi_{cc_c}^T \phi(\mathbf{z}) \,]^T$ is a $L \times 1$ kernel vector.

Combining Eq.11 and Eq.13, we obtain

$$\mathbf{y} = \Theta \cdot \gamma(\phi(\mathbf{z})) \tag{14}$$

where $\Theta = \frac{1}{\sqrt{L}} \left( \mathbf{E}_m \cdot \Lambda_b^{-1/2} \cdot \mathbf{P} \cdot \Lambda_w^{-1/2} \right)^T \left( \mathbf{B} \cdot \left( \mathbf{A}_{LC}^T - \frac{1}{L} \mathbf{1}_{LC}^T \right) \right)$ is a $M \times L$ matrix which can be calculated offline. Thus, through Eq.14, a low-dimensional representation $\mathbf{y}$ on $\mathbf{z}$ with enhanced discriminant power, suitable for classification tasks, has been introduced.

## C.4 Comments

The KDDA method implements an improved D-LDA in a high-dimensional feature space using a kernel approach. Its main advantages can be summarized as follows:

1. KDDA introduces a nonlinear mapping from the input space to an implicit high-dimensional feature space, where the nonlinear and complex distribution of patterns in the input space is "linearized" and "simplified" so that conventional LDA can be applied. It is not difficult to see that KDDA reduces to D-LDA for $\phi(\mathbf{z}) = \mathbf{z}$. Thus, D-LDA can be viewed as a special case of the proposed KDDA framework.

2. KDDA effectively solves the SSS problem in the high-dimensional feature space by employing an improved D-LDA algorithm. Unlike the original D-LDA method of [10] zero eigenvalues of the within-class scatter matrix are never used as divisors in the improved one. In this way, the optimal discriminant features can be exactly extracted from both of inside and outside of $\mathbf{S}_{WTH}$'s null space.

3. In GDA, to remove the null space of $\mathbf{S}_{WTH}$, it is required to compute the pseudo inverse of the kernel matrix $\mathbf{K}$, which could be extremely ill-conditioned when certain kernels or kernel parameters are used. Pseudo inversion is based on inversion of the nonzero eigenvalues. Due to round-off errors, it is not easy to identify the true null eigenvalues. As a result, numerical stability problems often occur [14]. However, it can been seen from the derivation of KDDA that such problems are avoided in KDDA. The improvement can be observed also in experimental results reported in Fig.4:**A** and Fig.5:**A**.

The detailed steps for implementing the KDDA method are summarized in Fig.6.

## III. Experimental Results

Two sets of experiments are included in the paper to illustrate the effectiveness of the KDDA algorithm. In all experiments reported here, we utilize the UMIST face database [27], [28], a multi-view database, consisting of 575 gray-scale images of 20 subjects, each covering a wide range of poses from profile to frontal views as well as race, gender and appearance. All input images are resized into $112 \times 92$, a standardized image size commonly used in FR tasks. The resulting standardized input vectors are of dimensionality $N = 10304$. Fig.1 depicts some sample images of a typical subset in the UMIST database.

### A. Distribution of Multi-view Face Patterns

The distribution of face patterns is highly non-convex and complex, especially when the patterns are subject to large variations in viewpoints as is the case with the UMIST database. The first experiment aims to provide insights on how the KDDA algorithm linearizes and simplifies the face pattern distribution.

For the sake of simplicity in visualization, we only use a subset of the database, which contains 170 images of 5 randomly selected subjects (classes). Four types of feature bases are generalized from the subset by utilizing the PCA, KPCA, D-LDA and KDDA algorithms respectively. In the four subspaces produced, two are linear, produced by PCA and D-LDA, and two are nonlinear, produced by KPCA and KDDA. In the sequence, all of images are projected onto the four subspaces. For each image, its projections in the first two most significant feature bases of each subspace are visualized in Figs.2-3.

In Fig.2, the visualized projections are the first two most significant principal components extracted by PCA and KPCA, and they provide a low-dimensional representation for the samples, which can be used to capture the structure of data. Thus, we can roughly learn the original distribution of face samples from Fig.2:**A**, which is non convex and complex as we expected based on the analysis presented in the previous sections. In Fig.2:**B**, KPCA generalizes PCA to its nonlinear counterpart using a RBF kernel function : $k(\mathbf{z}_1, \mathbf{z}_2) = exp\left(\frac{-||\mathbf{z}_1 - \mathbf{z}_2||^2}{\sigma^2}\right)$ with $\sigma^2 = 5\mathbf{e}6$. However, it is hard to find any useful improvement for the purpose of pattern classification from Fig.2:**B**. It can be therefore concluded that the low-dimensional representation obtained by PCA like techniques, achieve

simply object reconstruction, and they are not necessarily useful for discrimination and classification tasks [8], [29].

Unlike PCA approaches, LDA optimizes the low-dimensional representation of the objects based on separability criteria. Fig.3 depicts the first two most discriminant features extracted by utilizing D-LDA and KDDA respectively. Simple inspection of Figs.2-3 indicates that these features outperform, in terms of discriminant power, those obtained using PCA like methods. However, subject to limitation of linearity, some classes are still non-separable in the D-LDA-based subspace as shown in Fig.3:**A**. In contrast to this, we can see the linearization property of the KDDA-based subspace, as depicted in Fig.3:**B**, where all of classes are well linearly separable when a RBF kernel with $\sigma^2 = 5\mathbf{e}6$ is used.

## B. Comparison with KPCA and GDA

The second experiment compares the classification error rate performance of the KDDA algorithm to two other commonly used kernel FR algorithms, KPCA and GDA. The FR procedure is completed in two stages:

1. Feature extraction. The overall database is randomly partitioned into two subsets: the training set and test set. The training set is composed of 120 images: 6 images per person are randomly chosen. The remaining 455 images are used to form the test set. There is no overlapping between the two. After training is over, both sets are projected into the feature spaces derived from the KPCA, GDA and KDDA methods.

2. Classification. This is implemented by feeding feature vectors obtained in step-1 into a nearest neighbor classifier. It should be noted at this point that, since the focus in this paper is on feature extraction, a simple classifier is always prefered so that the FR performance is not mainly contributed by the classifier but the feature selection algorithms. We anticipate that the classification accuracy of all the three methods compared here will improve if a more sophisticated classifier such as SVM is used instead of the nearest neighbor. However, such an experiment is beyond the scope of this paper. To enhance the accuracy of performance evaluation, the classification error rates reported in this work are averaged over 8 runs. Each run is based on a random partition of the database into the training and test sets. Following the framework introduced in [30], [6], [31], the average

error rate, denoted as $E_{ave}$, is given as follows,

$$E_{ave} = \left(\sum_{i=1}^{r} t_{mis}^i\right) / (r \cdot t) \tag{15}$$

where $r$ is the number of runs, $t_{mis}^i$ is the number of misclassifications for the $i$th run, and $t$ is the number of total test samples of each run.

To evaluate the overall performance of the three methods, two typical kernel functions: namely the RBF and the polynomial function, and a wide range of parameter values are tested. Sensitivity analysis is performed with respect to the kernel parameters and the number of used feature vectors, $M$. Figs.4-5 depict the average error rates ($E_{ave}$) of the three methods compared when the RBF and polynomial kernels are used.

The only kernel parameter for RBF is the scale value $\sigma^2$. Fig.4:**A** shows the error rates as functions of $\sigma^2$ within the range from 0.5**e**7 to 1.5**e**8, when the optimal number of feature vectors, $M = M_{opt}$, is used. The optimal feature number is a result of the existence of the peaking effect in the feature selection procedure. It is well known that the classification error initially declines with the addition of new features, attains a minimum, and then starts to increase [32]. The optimal number can be found by searching the number of used feature vectors that results in the minimal summation of the error rates over the variation range of $\sigma^2$. In Fig.4:**A**, $M_{opt} = 99$ is the value used for KPCA, while $M_{opt} = 19$ is used for GDA and KDDA. Fig.4:**B** depicts the error rates as functions of $M$ within the range from 5 to 19, when optimal $\sigma^2 = \sigma_{opt}^2$ is used. Similar to $M_{opt}$, $\sigma_{opt}^2$ is defined as the scale parameter that results in the minimal summation of the error rates over the variation range of $M$ for the experiment discussed here. In Fig.4:**B**, a value $\sigma_{opt}^2 = 1.5$**e**8 is found for KPCA, $\sigma_{opt}^2 = 5.3333$**e**7 for GDA and $\sigma_{opt}^2 = 1.3389$**e**7 for KDDA.

As such, the average error rates of the three methods with polynomial kernel ($k(\mathbf{z}_1, \mathbf{z}_2) = (a \cdot (\mathbf{z}_1 \cdot \mathbf{z}_2) + b)^d$) are shown in Fig.5. For the sake of simplicity, we only test the influence of $a$, while $b = 1$ and $d = 3$ are fixed. Fig.5:**A** depicts the error rates as functions of $a$ within the range from 1**e** − 9 to 5**e** − 8, where $M_{opt} = 100$ for KPCA, $M_{opt} = 19$ for GDA and KDDA. Fig.5:**B** shows the error rates as functions of $M$ within the range from 5 to 19 with $a_{opt} = 1$**e** − 9 for KPCA, $a_{opt} = 2.822$**e** − 8 for GDA and $a_{opt} = 1$**e** − 9 for KDDA, determined similarly to $\sigma_{opt}^2$ and $M_{opt}$.

Let $\alpha_M$ and $\beta_M$ be the average error rates of KDDA and any one of other two methods respectively, where $M = [5 \dots 19]$. From Fig.4:**B** and Fig.5:**B**, we can obtain an interesting quantity comparison: the average percentages of the error rate of KDDA over those of other methods by $\sum_{M=5}^{19} (\alpha_M/\beta_M)$. The results are tabulated in Table I. The average error rate of KDDA to KPCA and GDA are only about 34.375% and 47.765% respectively. It should be also noted that Fig.4:**A** and Fig.5:**A** reveal the numerical stability problems existing in practical implementations of GDA. Comparing the GDA performance to that of KDDA we can easily see that the later is more stable and predictable, resulting in a cost effective determination of parameter values during the training phase.

## IV. Conclusion

A new FR method has been introduced in this paper. The proposed method combines kernel-based methodologies with discriminant analysis techniques. The kernel function is utilized to map the original face patterns to a high-dimensional feature space, where the highly non-convex and complex distribution of face patterns is linearized and simplified, so that linear discriminant techniques can be used for feature extraction. The small sample size problem caused by high dimensionality of mapped patterns, is addressed by an improved D-LDA technique which exactly finds the optimal discriminant subspace of the feature space without any loss of significant discriminant information. Experimental results indicate that the performance of the KDDA algorithm is overall superior to those obtained by the KPCA or GDA approaches. In conclusion, the KDDA algorithm is a general pattern recognition method for nonlinearly feature extraction from high-dimensional input patterns without suffering from the SSS problem. We expect that in addition to face recognition, KDDA will provide excellent performance in applications where classification tasks are routinely performed, such as content-based image indexing and retrieval, video and audio classification.

## Acknowledgments

## APPENDIX I. COMPUTATION OF $(\Phi_b^T \Phi_b)$

Expanding $\Phi_b^T \Phi_b$, we have

$$\Phi_b^T \Phi_b = \left[ \tilde{\bar{\phi}}_1 \ \cdots \ \tilde{\bar{\phi}}_C \right]^T \left[ \tilde{\bar{\phi}}_1 \ \cdots \ \tilde{\bar{\phi}}_C \right] = \left( \tilde{\bar{\phi}}_i^{\ T} \tilde{\bar{\phi}}_j \right)_{\substack{i=1,\cdots,C \\ j=1,\cdots,C}} \tag{16}$$

where

$$\tilde{\bar{\phi}}_i^{\ T} \tilde{\bar{\phi}}_j = \frac{\sqrt{C_i C_j}}{N} \left( \bar{\phi}_i^T \bar{\phi}_j - \bar{\phi}_i^T \bar{\phi} - \bar{\phi}^T \bar{\phi}_j + \bar{\phi}^T \bar{\phi} \right) \tag{17}$$

We develop each term of Eq.17 according to the kernel matrix $\mathbf{K}$ as follows,

- $\bar{\phi}^T \bar{\phi} = \left( \frac{1}{L} \sum_{l=1}^{C} \sum_{k=1}^{C_l} \phi_{lk} \right)^T \left( \frac{1}{L} \sum_{h=1}^{C} \sum_{m=1}^{C_h} \phi_{hm} \right) = \frac{1}{L^2} \sum_{l=1}^{C} \sum_{k=1}^{C_l} \sum_{h=1}^{C} \sum_{m=1}^{C_h} (k_{km})_{lh}$
$\Rightarrow \left( \bar{\phi}^T \bar{\phi} \right)_{\substack{i=1,\cdots,C \\ j=1,\cdots,C}} = \frac{1}{L^2} \left( \mathbf{1}_{LC}^T \cdot \mathbf{K} \cdot \mathbf{1}_{LC} \right);$

- $\bar{\phi}^T \bar{\phi}_j = \left( \frac{1}{L} \sum_{l=1}^{C} \sum_{k=1}^{C_l} \phi_{lk} \right)^T \left( \frac{1}{C_j} \sum_{m=1}^{C_j} \phi_{jm} \right) = \frac{1}{L C_j} \sum_{l=1}^{C} \sum_{k=1}^{C_l} \sum_{m=1}^{C_j} (k_{km})_{lj}$
$\Rightarrow \left( \bar{\phi}^T \bar{\phi}_j \right)_{\substack{i=1,\cdots,C \\ j=1,\cdots,C}} = \frac{1}{L} \left( \mathbf{1}_{LC}^T \cdot \mathbf{K} \cdot \mathbf{A}_{LC} \right);$

- $\bar{\phi}_i^T \bar{\phi} = \left( \frac{1}{C_i} \sum_{m=1}^{C_i} \phi_{im} \right)^T \left( \frac{1}{L} \sum_{l=1}^{C} \sum_{k=1}^{C_l} \phi_{lk} \right) = \frac{1}{L C_i} \sum_{m=1}^{C_i} \sum_{l=1}^{C} \sum_{k=1}^{C_l} (k_{mk})_{il}$
$\Rightarrow \left( \bar{\phi}_i^T \bar{\phi} \right)_{\substack{i=1,\cdots,C \\ j=1,\cdots,C}} = \frac{1}{L} \left( \mathbf{A}_{LC}^T \cdot \mathbf{K} \cdot \mathbf{1}_{LC} \right);$

- $\bar{\phi}_i^T \bar{\phi}_j = \left( \frac{1}{C_i} \sum_{m=1}^{C_i} \phi_{im} \right)^T \left( \frac{1}{C_j} \sum_{n=1}^{C_j} \phi_{jn} \right) = \frac{1}{C_i C_j} \sum_{m=1}^{C_i} \sum_{n=1}^{C_j} (k_{mn})_{ij}$
$\Rightarrow \left( \bar{\phi}_i^T \bar{\phi}_j \right)_{\substack{i=1,\cdots,C \\ j=1,\cdots,C}} = \left( \mathbf{A}_{LC}^T \cdot \mathbf{K} \cdot \mathbf{A}_{LC} \right).$

Applying the above derivations into Eq.17, we obtain the Eq.7.

## APPENDIX II. COMPUTATION OF $\left( \Phi_b^T \mathbf{S}_{WTH} \Phi_b \right)$

Expanding $\Phi_b^T \mathbf{S}_{WTH} \Phi_b$, we have

$$\Phi_b^T \mathbf{S}_{WTH} \Phi_b = \left[ \tilde{\bar{\phi}}_1 \ \cdots \ \tilde{\bar{\phi}}_C \right]^T \mathbf{S}_{WTH} \left[ \tilde{\bar{\phi}}_1 \ \cdots \ \tilde{\bar{\phi}}_C \right] = \left( \tilde{\bar{\phi}}_i^{\ T} \mathbf{S}_{WTH} \tilde{\bar{\phi}}_j \right)_{\substack{i=1,\cdots,C \\ j=1,\cdots,C}} \tag{18}$$

where

$$\begin{aligned}
\tilde{\bar{\phi}}_i^{\ T} \mathbf{S}_{WTH} \tilde{\bar{\phi}}_j &= \frac{1}{L} \tilde{\bar{\phi}}_i^{\ T} \left( \sum_{l=1}^{C} \sum_{k=1}^{C_l} (\phi_{lk} - \bar{\phi}_l)(\phi_{lk} - \bar{\phi}_l)^T \right) \tilde{\bar{\phi}}_j \\
&= \frac{1}{L} \tilde{\bar{\phi}}_i^{\ T} \left( \sum_{l=1}^{C} \sum_{k=1}^{C_l} \phi_{lk} \phi_{lk}^T - \sum_{l=1}^{C} \bar{\phi}_l \left( \sum_{k=1}^{C_l} \phi_{lk}^T \right) - \sum_{l=1}^{C} \left( \sum_{k=1}^{C_l} \phi_{lk} \right) \bar{\phi}_l^T + \sum_{l=1}^{C} C_l \bar{\phi}_l \bar{\phi}_l^T \right) \tilde{\bar{\phi}}_j \\
&= \frac{1}{L} \left( \sum_{l=1}^{C} \sum_{k=1}^{C_l} \tilde{\bar{\phi}}_i^{\ T} \phi_{lk} \phi_{lk}^T \tilde{\bar{\phi}}_j - \sum_{l=1}^{C} C_l \tilde{\bar{\phi}}_i^{\ T} \bar{\phi}_l \bar{\phi}_l^T \tilde{\bar{\phi}}_j \right)
\end{aligned} \tag{19}$$

Firstly, expand the term $\sum_{l=1}^{C}\sum_{k=1}^{C_l}\tilde{\bar{\phi}}_i^T\phi_{lk}\phi_{lk}^T\tilde{\bar{\phi}}_j$ in Eq.19, and have

$$\sum_{l=1}^{C}\sum_{k=1}^{C_l}\tilde{\bar{\phi}}_i^T\phi_{lk}\phi_{lk}^T\tilde{\bar{\phi}}_j = \frac{\sqrt{C_iC_j}}{L}\sum_{l=1}^{C}\sum_{k=1}^{C_l}\left(\bar{\phi}_i^T\phi_{lk}\phi_{lk}^T\bar{\phi}_j - \bar{\phi}_i^T\phi_{lk}\phi_{lk}^T\bar{\phi} - \bar{\phi}^T\phi_{lk}\phi_{lk}^T\bar{\phi}_j + \bar{\phi}^T\phi_{lk}\phi_{lk}^T\bar{\phi}\right)$$

(20)

We develop each term of Eq.20 according to the kernel matrix $\mathbf{K}$ as follows,

$$\bullet \quad \sum_{l=1}^{C}\sum_{k=1}^{C_l}\bar{\phi}_i^T\phi_{lk}\phi_{lk}^T\bar{\phi}_j = \frac{1}{C_iC_j}\sum_{l=1}^{C}\sum_{k=1}^{C_l}\left(\sum_{m=1}^{C_i}\phi_{im}^T\phi_{lk}\right)\left(\sum_{n=1}^{C_j}\phi_{lk}^T\phi_{jn}\right)$$

$$= \frac{1}{C_iC_j}\sum_{l=1}^{C}\sum_{k=1}^{C_l}\sum_{m=1}^{C_i}\sum_{n=1}^{C_j}(k_{mk})_{il}(k_{kn})_{lj}$$

$$\Rightarrow \left(\sum_{l=1}^{C}\sum_{k=1}^{C_l}\bar{\phi}_i^T\phi_{lk}\phi_{lk}^T\bar{\phi}_j\right)_{\substack{i=1,\cdots,C\\j=1,\cdots,C}} = \left(\mathbf{A}_{LC}^T\cdot\mathbf{K}\cdot\mathbf{K}\cdot\mathbf{A}_{LC}\right);$$

$$\bullet \quad \sum_{l=1}^{C}\sum_{k=1}^{C_l}\bar{\phi}_i^T\phi_{lk}\phi_{lk}^T\bar{\phi} = \frac{1}{LC_i}\sum_{l=1}^{C}\sum_{k=1}^{C_l}\left(\sum_{n=1}^{C_i}\phi_{in}^T\phi_{lk}\right)\left(\sum_{h=1}^{C}\sum_{m=1}^{C_h}\phi_{lk}^T\phi_{hm}\right)$$

$$= \frac{1}{LC_i}\sum_{l=1}^{C}\sum_{k=1}^{C_l}\sum_{n=1}^{C_i}\sum_{h=1}^{C}\sum_{m=1}^{C_h}(k_{nk})_{il}(k_{km})_{lh}$$

$$\Rightarrow \left(\sum_{l=1}^{C}\sum_{k=1}^{C_l}\bar{\phi}_i^T\phi_{lk}\phi_{lk}^T\bar{\phi}\right)_{\substack{i=1,\cdots,C\\j=1,\cdots,C}} = \frac{1}{L}\left(\mathbf{A}_{LC}^T\cdot\mathbf{K}\cdot\mathbf{K}\cdot\mathbf{1}_{LC}\right);$$

$$\bullet \quad \sum_{l=1}^{C}\sum_{k=1}^{C_l}\bar{\phi}^T\phi_{lk}\phi_{lk}^T\bar{\phi}_j = \frac{1}{LC_j}\sum_{l=1}^{C}\sum_{k=1}^{C_l}\left(\sum_{h=1}^{C}\sum_{m=1}^{C_h}\phi_{hm}^T\phi_{lk}\right)\left(\sum_{n=1}^{C_j}\phi_{lk}^T\phi_{jn}\right)$$

$$= \frac{1}{LC_j}\sum_{l=1}^{C}\sum_{k=1}^{C_l}\sum_{h=1}^{C}\sum_{m=1}^{C_h}\sum_{n=1}^{C_j}(k_{mk})_{hl}(k_{kn})_{lj}$$

$$\Rightarrow \left(\sum_{l=1}^{C}\sum_{k=1}^{C_l}\bar{\phi}^T\phi_{lk}\phi_{lk}^T\bar{\phi}_j\right)_{\substack{i=1,\cdots,C\\j=1,\cdots,C}} = \frac{1}{L}(\mathbf{1}_{LC}^T\cdot\mathbf{K}\cdot\mathbf{K}\cdot\mathbf{A}_{LC});$$

$$\bullet \quad \sum_{l=1}^{C}\sum_{k=1}^{C_l}\bar{\phi}^T\phi_{lk}\phi_{lk}^T\bar{\phi} = \frac{1}{L^2}\sum_{l=1}^{C}\sum_{k=1}^{C_l}\left(\sum_{h=1}^{C}\sum_{m=1}^{C_h}\phi_{hm}^T\phi_{lk}\right)\left(\sum_{p=1}^{C}\sum_{q=1}^{C_p}\phi_{lk}^T\phi_{pq}\right)$$

$$= \frac{1}{L^2}\sum_{l=1}^{C}\sum_{k=1}^{C_l}\sum_{h=1}^{C}\sum_{m=1}^{C_h}\sum_{p=1}^{C}\sum_{q=1}^{C_p}(k_{mk})_{hl}(k_{kq})_{lp}$$

$$\Rightarrow \left(\sum_{l=1}^{C}\sum_{k=1}^{C_l}\bar{\phi}^T\phi_{lk}\phi_{lk}^T\bar{\phi}\right)_{\substack{i=1,\cdots,C\\j=1,\cdots,C}} = \frac{1}{L^2}(\mathbf{1}_{LC}^T\cdot\mathbf{K}\cdot\mathbf{K}\cdot\mathbf{1}_{LC}).$$

Defining $\mathbf{J}1 = \left(\sum_{l=1}^{C}\sum_{k=1}^{C_l}\tilde{\bar{\phi}}_i^T\phi_{lk}\phi_{lk}^T\tilde{\bar{\phi}}_j\right)_{\substack{i=1,\cdots,C\\j=1,\cdots,C}}$, we conclude:

$$\mathbf{J}1 = \frac{1}{L}\mathbf{B}\cdot(\mathbf{A}_{LC}^T\cdot\mathbf{K}\cdot\mathbf{K}\cdot\mathbf{A}_{LC} - \frac{1}{L}(A_{Nc}^T\cdot\mathbf{K}\cdot\mathbf{K}\cdot\mathbf{1}_{LC}) -$$
$$\frac{1}{L}(\mathbf{1}_{LC}^T\cdot\mathbf{K}\cdot\mathbf{K}\cdot\mathbf{A}_{LC}) + \frac{1}{L^2}(\mathbf{1}_{LC}^T\cdot\mathbf{K}\cdot\mathbf{K}\cdot\mathbf{1}_{LC}))\cdot\mathbf{B}$$

(21)

Expanding the term $\sum_{l=1}^{C}C_l\tilde{\bar{\phi}}_i^T\bar{\phi}_l\bar{\phi}_l^T\tilde{\bar{\phi}}_j$ in Eq.19, we obtain:

$$\sum_{l=1}^{C}\sum_{k=1}^{C_l}\tilde{\bar{\phi}}_i^T\bar{\phi}_l\bar{\phi}_l^T\tilde{\bar{\phi}}_j = \frac{\sqrt{C_iC_j}}{L}\sum_{l=1}^{C}C_l\left(\bar{\phi}_i^T\bar{\phi}_l\bar{\phi}_l^T\bar{\phi}_j - \bar{\phi}_i^T\bar{\phi}_l\bar{\phi}_l^T\bar{\phi} - \bar{\phi}^T\bar{\phi}_l\bar{\phi}_l^T\bar{\phi}_j + \bar{\phi}^T\bar{\phi}_l\bar{\phi}_l^T\bar{\phi}\right)$$

(22)

Using the kernel matrix $\mathbf{K}$, the terms in Eq.22 can be developed as follows,

- $\left( \sum\limits_{l=1}^{C} C_l \bar{\phi}_i^T \bar{\phi}_l \bar{\phi}_l^T \bar{\phi}_j \right)_{\substack{i=1,\cdots,C \\ j=1,\cdots,C}} = \mathbf{A}_{LC}^T \cdot \mathbf{K} \cdot \mathbf{W} \cdot \mathbf{K} \cdot \mathbf{A}_{LC},$

- $\left( \sum\limits_{l=1}^{C} C_l \bar{\phi}_i^T \bar{\phi}_l \bar{\phi}_l^T \bar{\phi} \right)_{\substack{i=1,\cdots,C \\ j=1,\cdots,C}} = \frac{1}{L} \left( \mathbf{A}_{LC}^T \cdot \mathbf{K} \cdot \mathbf{W} \cdot \mathbf{K} \cdot \mathbf{1}_{LC} \right),$

- $\left( \sum\limits_{l=1}^{C} C_l \bar{\phi}^T \bar{\phi}_l \bar{\phi}_l^T \bar{\phi}_j \right)_{\substack{i=1,\cdots,C \\ j=1,\cdots,C}} = \frac{1}{L} \left( \mathbf{1}_{LC}^T \cdot \mathbf{K} \cdot \mathbf{W} \cdot \mathbf{K} \cdot \mathbf{A}_{LC} \right),$

- $\left( \sum\limits_{l=1}^{C} C_l \bar{\phi}^T \bar{\phi}_l \bar{\phi}_l^T \bar{\phi} \right)_{\substack{i=1,\cdots,C \\ j=1,\cdots,C}} = \frac{1}{L^2} \left( \mathbf{1}_{LC}^T \cdot \mathbf{K} \cdot \mathbf{W} \cdot \mathbf{K} \cdot \mathbf{1}_{LC} \right),$

where $\mathbf{W} = \mathbf{diag} \left[ \mathbf{w}_1 \ \cdots \ \mathbf{w}_c \right]$ is a $L \times L$ block diagonal matrix, and $\mathbf{w}_i$ is a $C_i \times C_i$ matrix with terms all equal to: $\frac{1}{C_i}$.

Defining $\mathbf{J}2 = \left( \sum\limits_{l=1}^{C} \sum\limits_{k=1}^{C_l} \tilde{\bar{\phi}}_i^T \bar{\phi}_l \bar{\phi}_l^T \tilde{\bar{\phi}}_j \right)_{\substack{i=1,\cdots,C \\ j=1,\cdots,C}}$, and using the above derivations, we conclude that:

$$
\begin{aligned}
\mathbf{J}2 = \quad & \frac{1}{L}\mathbf{B} \cdot (\mathbf{A}_{LC}^T \cdot \mathbf{K} \cdot \mathbf{W} \cdot \mathbf{K} \cdot \mathbf{A}_{LC} - \frac{1}{L}(\mathbf{A}_{LC}^T \cdot \mathbf{K} \cdot \mathbf{W} \cdot \mathbf{K} \cdot \mathbf{1}_{LC}) - \\
& \frac{1}{L}(\mathbf{1}_{LC}^T \cdot \mathbf{K} \cdot \mathbf{W} \cdot \mathbf{K} \cdot \mathbf{A}_{LC}) + \frac{1}{L^2}(\mathbf{1}_{LC}^T \cdot \mathbf{K} \cdot \mathbf{W} \cdot \mathbf{K} \cdot \mathbf{1}_{LC})) \cdot \mathbf{B}
\end{aligned}
\tag{23}
$$

Thus,

$$
\Phi_b^T \mathbf{S}_{WTH} \Phi_b = \left( \tilde{\bar{\phi}}_i^T \mathbf{S}_{WTH} \tilde{\bar{\phi}}_j \right)_{\substack{i=1,\cdots,C \\ j=1,\cdots,C}} = \frac{1}{L}(\mathbf{J}1 - \mathbf{J}2) \tag{24}
$$

## References

[1] Ashok Samal and Prasana A.Iyengar, "Automatic recognition and analysis of human faces and facial expressions: A survey", *Pattern Recognition*, vol. 25, pp. 65–77, 1992.

[2] Dominique Valentin, J. O. Toole Herve Abdi Alice, and Garrison W. Cottrell, "Connectionist models of face processing: A survey", *Pattern Recognition*, vol. 27, no. 9, pp. 1209–1230, 1994.

[3] R. Chellappa, C.L. Wilson, and S. Sirohey, "Human and machine recognition of faces: A survey", *Proceedings of the IEEE*, vol. 83, pp. 705–740, 1995.

[4] Shaogang Gong, Stephen J McKenna, and Alexandra Psarrou, *"Dynamic Vision From Images to Face Recognition"*, Imperial College Press, World Scientific Publishing, May 2000.

[5] Matthew Turk, "A random walk through eigenspace", *IEICE Trans. Inf. & Syst.*, vol. E84-D, no. 12, pp. 1586–1695, December 2001.

[6] Stan Z. Li and Juwei Lu, "Face recognition using the nearest feature line method", *IEEE Transactions on Neural Networks*, vol. 10, pp. 439–443, 1999.

[7] Matthew A. Turk and Alex P. Pentland, "Eigenfaces for recognition", *Journal of Cognitive Neuroscience*, vol. 3, no. 1, pp. 71–86, 1991.

[8] P. N. Belhumeur, J. P. Hespanha, and D. J. Kriegman, "Eigenfaces vs. Fisherfaces: recognition using class specific linear projection", *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 19, no. 7, pp. 711–720, 1997.

[9] Li-Fen Chen, Hong-Yuan Mark Liao, Ming-Tat Ko, Ja-Chen Lin, and Gwo-Jong Yu, "A new LDA-based face recognition system which can solve the small sample size problem", *Pattern Recognition*, vol. 33, pp. 1713–1726, 2000.

[10] Hua Yu and Jie Yang, "A direct lda algorithm for high-dimensional data with application to face recognition", *Pattern Recognition*, vol. 34, pp. 2067–2070, 2001.

[11] M. Bichsel and A. P. Pentland, "Human face recognition and the face image set's topology", *CVGIP: Image Understanding*, vol. 59, pp. 254–261, 1994.

[12] Bernhard Schölkopf, Chris Burges, and Alex J. Smola, *"Advances in Kernel Methods - Support Vector Learning"*, MIT Press, Cambridge, MA, 1999.

[13] Kernel machines website: http://www.kernel-machines.org, 2000.

[14] A. Ruiz and P.E. López de Teruel, "Nonlinear kernel-based statistical pattern analysis", *IEEE Transactions on Neural Networks*, vol. 12, no. 1, pp. 16–32, January 2001.

[15] K.-R. Müller, S. Mika, G. Rätsch, K. Tsuda, and B. Schölkopf, "An introduction to kernel-based learning algorithms", *IEEE Transactions on Neural Networks*, vol. 12, no. 2, pp. 181–201, March 2001.

[16] C. Cortes and V. N. Vapnik, "Support vector networks", *Machine Learning*, vol. 20, pp. 273–297, 1995.

[17] V. N. Vapnik, *"The Nature of Statistical Learning Theory"*, Springer-Verlag, New York, 1995.

[18] B. Schölkopf, *"Support Vector Learning"*, Oldenbourg-Verlag, Munich, Germany, 1997.

[19] Schölkopf B., Smola A., and Müller K. R., "Nonlinear component analysis as a kernel eigenvalue problem", *Neural Computation*, vol. 10, pp. 1299–1319, 1999.

[20] G. Baudat and F. Anouar, "Generalized discriminant analysis using a kernel approach", *Neural Computation*, vol. 12, pp. 2385–2404, 2000.

[21] R.A. Fisher, "The use of multiple measures in taxonomic problems", *Ann. Eugenics*, vol. 7, pp. 179–188, 1936.

[22] M. A. Aizerman, E. M. Braverman, and L. I. Rozonoér, "Theoretical foundations of the potential function method in pattern recognition learning", *Automation and Remote Control*, vol. 25, pp. 821–837, 1964.

[23] V. N. Vapnik, *"Statistical Learning Theory"*, Wiley, New York, 1998.

[24] B. Scholkopf, S. Mika, C.J.C. Burges, P. Knirsch, K.-R. Muller, G. Ratsch, and A.J. Smola, "Input space versus feature space in kernel-based methods", *IEEE Transactions on Neural Networks*, vol. 10, no. 5, pp. 1000 –1017, September 1999.

[25] K. Liu, Y.Q. Cheng, J.Y. Yang, and X. Liu, "An efficient algorithm for foley-sammon optimal set of discriminant vectors by algebraic method", *Int. J. Pattern Recog. Artif. Intell.*, vol. 6, pp. 817–829, 1992.

[26] Roger A. Horn and Charles R. Johnson, *Matrix Analysis*, Cambridge University Press, 1992.

[27] Daniel Graham and Nigel Allinson, Web site of umist multi-view face database: "http://images.ee.umist.ac.uk/danny/database.html", *Image Engineering and Neural Computing Lab, UMIST, UK*, 1998.

[28] Daniel B Graham and Nigel M Allinson, "Characterizing virtual eigensignatures for general purpose face recognition", in *Face Recognition: From Theory to Applications, NATO ASI Series F, Computer and Systems Sciences*, H. Wechsler, P. J. Phillips, V. Bruce, F. Fogelman-Soulie, and T. S. Huang, Eds., vol. 163, pp. 446–456. 1998.

[29] D. L. Swets and J. Weng, "Using discriminant eigenfeatures for image retrieval", *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 18, pp. 831–836, 1996.

[30] Steve Lawrence, C. Lee Giles, A.C. Tsoi, and A.D. Back, "Face recognition: A convolutional neural network approach", *IEEE Transactions on Neural Networks*, vol. 8, no. 1, pp. 98–113, 1997.

[31] Meng Joo Er, Shiqian Wu, Juwei Lu, and Hock Lye Toh, "Face recognition with radial basis function (RBF) neural networks", *IEEE Transactions on Neural Networks*, vol. 13, no. 3, pp. 697–710, May 2002.

[32] Sarunas J. Raudys and Anil K. Jain, "Small sample size effects in statistical pattern recognition: Recommendations for practitioners", *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 13, no. 3, pp. 252–264, 1991.

## List of Tables

TABLE I

The average percentages of the error rate of KDDA over those of others.

| Kernel | RBF | Polynomial | (RBF+Polynomial)/2 |
|---|---|---|---|
| KDDA/KPCA | 33.669% | 35.081% | 34.375% |
| KDDA/GDA | 47.866% | 47.664% | 47.765% |

## LIST OF FIGURES
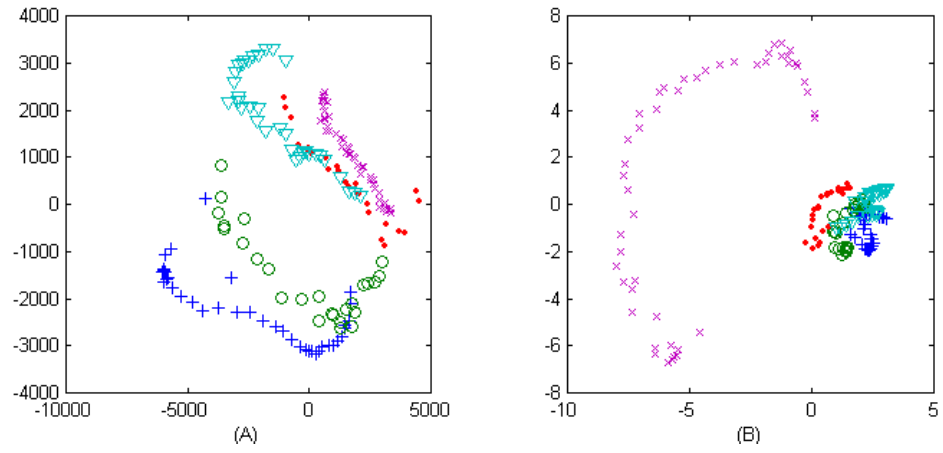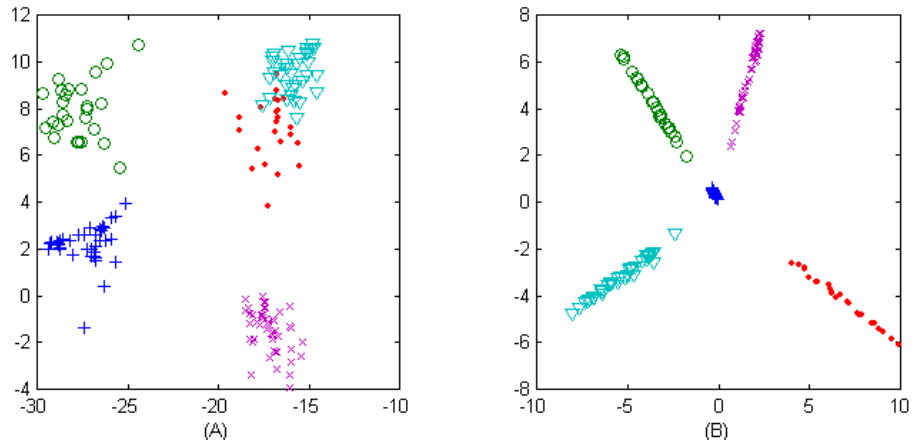
Fig. 1. Some face samples of one subject from the UMIST face database.

Fig. 2. Distribution of 170 samples of 5 subjects in PCA- and KPCA-based subspaces. **A**: PCA-based subspace ($\subset \mathbb{R}^N$). **B**: KPCA-based subspace ($\subset \mathbb{F}$).
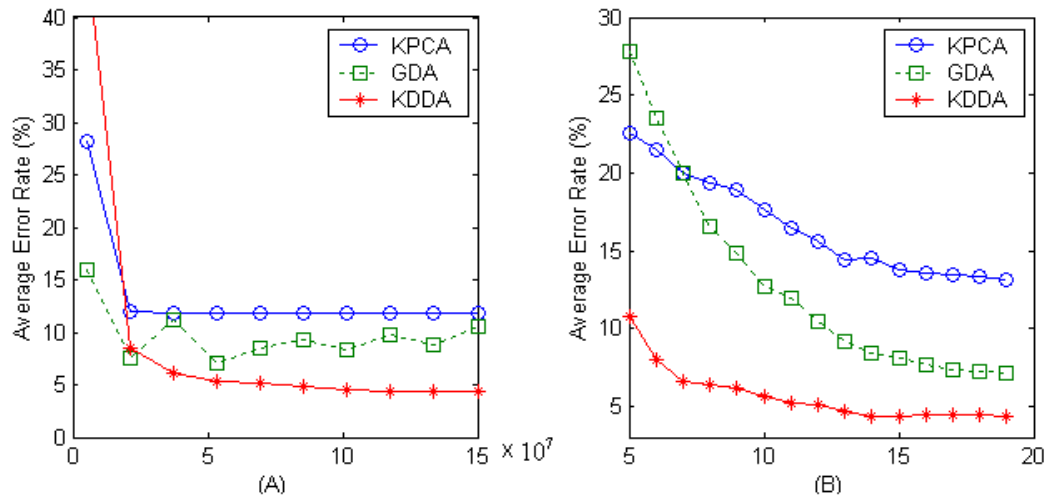
Fig. 3. Distribution of 170 samples of 5 subjects in D-LDA- and KDDA-based subspaces. **A**: D-LDA-based subspace ($\subset \mathbb{R}^N$). **B**: KDDA-based subspace ($\subset \mathbb{F}$).

Fig. 4. Comparison of error rates based on RBF kernel function. **A**: error rates as functions of $\sigma^2$. **B**: error rate as functions of $M$.
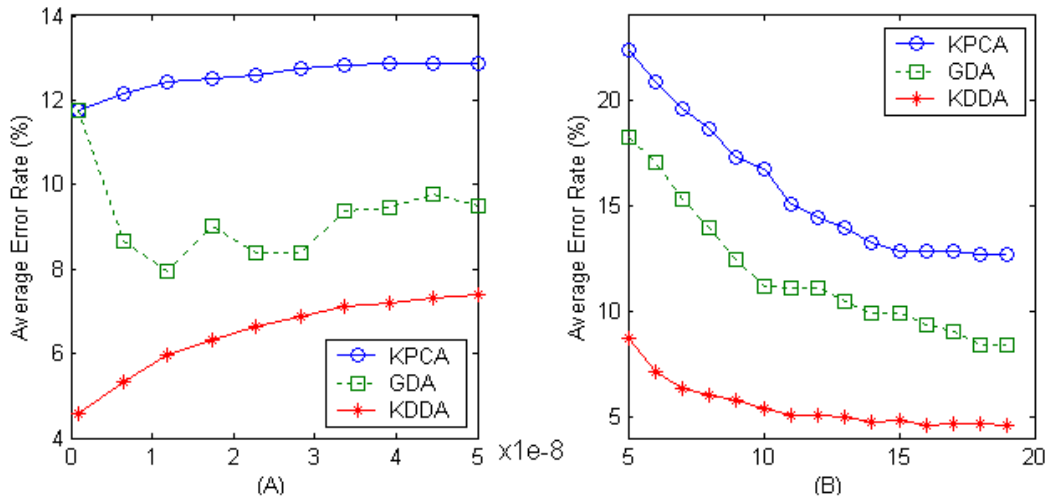
Fig. 5. Comparison of error rates based on Polynomial kernel function. **A**: error rates as functions of $a$. **B**: error rate as functions of $M$.

**Input:** A set of training face images $\{\mathbf{z}_i\}_{i=1}^L$, each of images is represented

   as a $L$-dimensional vector.

**Output:** A low-dimensional representation $\mathbf{y}$ of $\mathbf{z}$ with enhanced

   discriminatory power.

**Algorithm:**

   Step 1. Calculate kernel matrix $\mathbf{K}$ using Eq.6.

   Step 2. Calculate $\Phi_b^T\Phi_b$ using Eq.7, and find $\mathbf{E}_m$ and $\Lambda_b$ from $\Phi_b^T\Phi_b$

      in the way shown in section(II-C.1).

   Step 3. Calculate $\mathbf{U}^T\mathbf{S}_{WTH}\mathbf{U}$ using Eq.8 and Eq.9, and

      if $\left|\mathbf{U}^T\mathbf{S}_{WTH}\mathbf{U}\right| \neq 0$ then

         /* *using the conventional criterion in Eq.3 when* $\mathbf{U}^T\mathbf{S}_{WTH}\mathbf{U}$ *is nonsingular.* */

         Calculate $\mathbf{P}$ and $\Lambda_w$ from $\mathbf{U}^T\mathbf{S}_{WTH}\mathbf{U}$ as shown in section(II-C.2);

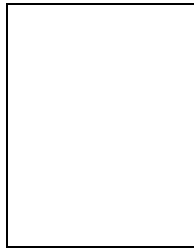      else /* *using the modified criterion in Eq.10 when* $\mathbf{U}^T\mathbf{S}_{WTH}\mathbf{U}$ *is singular.* */

         Calculate $\mathbf{P}$ and $\Lambda_w$ from $\mathbf{U}^T(\mathbf{S}_{BTW} + \mathbf{S}_{WTH})\mathbf{U}$ as shown in section(II-C.2);

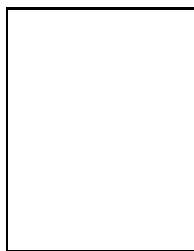   Step 4. Calculate $\Theta$ in Eq.14.

   Step 5. For input pattern $\mathbf{z}$, calculate its kernel matrix $\gamma(\phi(\mathbf{z}))$ in Eq.13.

   Step 6. The optimal discriminant feature representation of $\mathbf{z}$ can be obtained

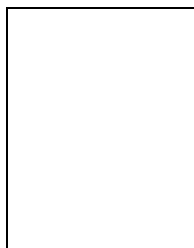      by $\mathbf{y} = \Theta \cdot \gamma(\phi(\mathbf{z}))$ based on Eq.14.

Fig. 6.   KDDA pseudo-code implementation

**Juwei Lu** received the B.Eng. degree in Electrical Engineering from Nanjing University of Aeronautics and Astronautics, China, in 1994, and the M.Eng. degree from the School of Electrical and Electronic Engineering, Nanyang Technological University, Singapore, in 1999. From July 1999 to January 2001, he was with the Center for Signal Processing, Singapore, as a Research Engineer. Currently, he is pursuing the Ph.D. degree in the Edward S. Rogers Sr. Department of Electrical and Computer Engineering, University of Toronto, Canada. His research interests include multimedia signal processing, face detection and recognition, kernel methods, support vector machines, neural networks, and Boosting technologies.

**K.N. (Kostas) Plataniotis** received his B. Eng. degree in Computer Engineering & Informatics from University of Patras, Greece in 1988 and his M.S and Ph.D degrees in Electrical Engineering from Florida Institute of Technology(Florida Tech) in Melbourne, Florida, in 1992 and 1994 respectively. He was with the Computer Technology Institute (C.T.I) in Patras, Greece from 1989 to 1991. He was a Postdoctoral Fellow at the Digital Signal & Image Processing Laboratory, Department of Electrical and Computer Engineering University of Toronto, Toronto, Canada from November 1994 to May 1997. From September 1997 to June 1999 he was an Assistant Professor with the School of Computer Science at Ryerson University, Toronto, Ontario. He is now an Assistant Professor with The Edward S. Rogers Sr. Department of Electrical and Computer Engineering at the University of Toronto in Toronto, Ontario, a Nortel Institute for Telecommunications Associate, and an Adjunct Professor in the Department of Mathematics, Physics and Computer Science at Ryerson University. His research interests include multimedia signal processing, intelligent and adaptive systems, and wireless communication systems. Dr. Plataniotis is a member of the Technical Chamber of Greece, and IEEE.

**Anastasios N. Venetsanopoulos** received the Diploma in Engineering degree from the National Technical University of Athens (NTU), Greece, in 1965, and the M.S., M.Phil., and Ph.D. degrees in Electrical Engineering from Yale University in 1966, 1968 and 1969 respectively. He joined the Department of Electrical and Computer Engineering of the University of Toronto in September 1968 as a Lecturer and he was promoted to Assistant Professor in 1970, Associate Professor in 1973, and Professor in 1981. Prof. A.N. Venetsanopoulos has served as Chair of the Communications Group and Associate Chair of the Department of Electrical Engineering. Between July 1997 - June 2001, he was Associate Chair: Graduate Studies of the Department of Electrical and Computer Engineering and was Acting Chair during the spring term of 1998-99. In 1999 a Chair in Multimedia was established in the ECE Department, made possible by a donation of $1.25M\ from\ Bell\ Canada, matched\ by\ 1.0M$ of university funds. Prof. A.N. Venetsanopoulos assumed the position as Inaugural Chairholder in July 1999 and two additional Assistant Professor positions became available in the same area. Since July 2001 he has served as the 12th Dean of the Faculty of Applied Science and Engineering of the University of Toronto.

Prof. A.N. Venetsanopoulos was on research leave at the Imperial College of Science and Technology, the National Technical University of Athens, the Swiss Federal Institute of Technology, the University of Florence and the Federal University of Rio de Janeiro, and has also served as Adjunct Professor at Concordia University. He

has served as lecturer in 138 short courses to industry and continuing education programs and as Consultant to numerous organizations; he is a contributor to twenty-nine (29) books, a co-author of Nonlinear Filters in Image Processing: Principles Applications (ISBN-0-7923-9049-0), and Artificial Neural Networks: Learning Algorithms, Performance Evaluation and Applications (ISBN-0-7923-9297-3), Fuzzy Reasoning in Information Decision and Control systems (ISBN-0-77293-2643-1) and Color Image Processing and Applications (ISBN-3-540-66953-1), and has published over 700 papers in refereed journals and conference proceedings on digital signal and image processing and digital communications.

Prof. Venetsanopoulos has served as Chair on numerous boards, councils and technical conference committees of the Institute of Electrical and Electronic Engineers (IEEE), such as the Toronto Section (1977-1979) and the IEEE Central Canada Council (1980-1982); he was President of the Canadian Society for Electrical Engineering and Vice President of the Engineering Institute of Canada (EIC) (1983-1986). He was a Guest Editor or Associate Editor for several IEEE journals and the Editor of the Canadian Electrical Engineering Journal (1981-1983). He is a member of the IEEE Communications, Circuits and Systems, Computer, and Signal Processing Societies of IEEE, as well as a member of Sigma Xi, the Technical Chamber of Greece, the European Association of Signal Processing, the Association of Professional Engineers of Ontario (APEO) and Greece. He was elected as a Fellow of the IEEE "for contributions to digital signal and image processing", he is also a Fellow of the EIC, and was awarded an Honorary Doctorate from the National Technical University of Athens, in October 1994. In October 1996 he was awarded the "Excellence in Innovation Award" of the Information Technology Research Centre of Ontario and Royal Bank of Canada, "for innovative work in color image processing and its industrial applications". In November 2000 he became Recipient of the "Millennium Medal of IEEE". In April 2001 he became a Fellow of the Canadian Academy of Engineering.