

L1-Norm Heteroscedastic Discriminant Analysis under Mixture of Gaussian Distributions

Wenming Zheng, *Member, IEEE*, Cheng Lu, Zhouchen Lin, *Fellow, IEEE*, Tong Zhang, Zhen Cui, Wankou Yang

Abstract—Fisher’s criterion is one of the most popular discriminant criteria for feature extraction. It is defined as the generalized Rayleigh quotient of the between-class scatter distance to the within-class scatter distance. Consequently, Fisher’s criterion does not take advantage of the discriminant information in the class covariance differences and hence its discriminant ability largely depends on the class mean differences. If the class mean distances are relatively large compared with the within-class scatter distance, Fisher’s criterion based discriminant analysis methods may achieve a good discriminant performance. Otherwise, it may not deliver good results. Moreover, we observe that the between-class distance of Fisher’s criterion is based on the ℓ_2 norm, which would be disadvantageous to separate the classes with smaller class mean distances. To overcome the drawback of Fisher’s criterion, in this paper we firstly derive a new discriminant criterion, expressed as a *mixture of absolute generalized Rayleigh quotients* (MAGRQ), based on a Bayes error upper bound estimation, where mixture of Gaussians is adopted to approximate the real distribution of data samples. Then, the criterion is further modified by replacing ℓ_2 norm with ℓ_1 one to better describe the between-class scatter distance, such that it would be more effective to separate the different classes. Moreover, we propose a novel ℓ_1 -norm heteroscedastic discriminant analysis method based on the new discriminant analysis (L1-HDA/GM) for heteroscedastic feature extraction, in which the optimization problem of L1-HDA/GM can be efficiently solved by using eigenvalue decomposition approach. Finally, we conduct extensive experiments on four real data sets and demonstrate that the proposed method achieves much competitive results compared to the state-of-the-art methods.

Index Terms—L1-norm heteroscedastic discriminant analysis, Heteroscedastic discriminant criterion, Fisher’s discriminant cri-

terion, Rayleigh quotient, Feature extraction

I. INTRODUCTION

Linear feature extraction plays a crucial role in statistical pattern recognition [1][2][3]. The goal of linear feature extraction can be regarded as seeking a transformation matrix that transforms the input data from the original high dimensional space to a low dimensional space while preserving some useful information. Fisher’s linear discriminant analysis (FLDA) [4] is one of the most popular linear feature extraction methods, which aims to find a set of optimal discriminant vectors such that the projections of the training samples onto these vectors have maximal between-class scatter distance and minimal within-class scatter distance. This is realized by solving a series of discriminant vectors that maximize Fisher’s discriminant criterion, defined as the generalized Rayleigh quotient of the between-class scatter distance to the within-class scatter distance. Over the past several decades, Fisher’s criterion based discriminative feature extraction methods had been successfully applied to face recognition [5], image retrieval [6], and speech recognition [7]. More recently, Yang et al. [8] adopted the Fisher’s criterion to enhance the discriminative ability of the sparse coefficient matrix in the sparse representation model [9]. Although the Fisher’s criterion had been shown to be very effective in practical applications, it should be noted that this criterion was developed under the homoscedastic distributions of the class data samples. Since the Fisher’s criterion is characterized by the ratio of the between-class scatter distance to the within-class scatter distance, it may not deliver a good discriminant performance when the class mean distances are relatively small compared with the within-class scatter distance. Take the electroencephalogram (EEG) feature extraction as an example, we cannot determine what features are the most discriminative ones according to Fisher’s criterion because the EEG signal conditioned on each class is often assumed to have a zero mean [10] and hence Fisher’s ratio will always be zero. In such a case, only the class covariance matrices can be utilized to extract the discriminant features [11]. Consequently, how to define a good discriminant criterion for extracting the useful discriminant features from the class covariance matrices is the major goal of this study.

In order to utilize the discriminant information from both class means and class covariance matrices, many heteroscedastic discriminant criteria have been proposed during the past several years [12], [13], [14], [15], [17], [18], which result in various heteroscedastic discriminant analysis (HDA) methods. Here we divide them into the following three categories. The

Manuscript received July 30, 2017; revised February 7, 2018; accepted July 25, 2018. This work was supported by the National Basic Research Program of China under Grants 2015CB351704 and 2015CB352502, the National Natural Science Foundation of China under Grants 61572009, 61625301, 61731018, and 61772276, the Jiangsu Provincial Key Research and Development Program under Grant BE2016616, and the support of Qualcomm, and Microsoft Research Asia. (*Corresponding author: Wenming Zheng.*)

Wenming Zheng is with the Key Laboratory of Child Development and Learning Science, Ministry of Education, School of Biological Sciences and Medical Engineering, Southeast University, Nanjing 210096, China (E-mail: wenming_zheng@seu.edu.cn).

Cheng Lu and Tong Zhang are with the Key Laboratory of Child Development and Learning Science, Ministry of Education, School of Information Science and Engineering, Southeast University, Nanjing 210096, China (E-mail: cheng.lu@seu.edu.cn; tongzhang@seu.edu.cn).

Zhouchen Lin is affiliated with the Key Laboratory of Machine Perception, Ministry of Education, School of Electronics Engineering and Computer Science, Peking University, Beijing 100871, China, and also affiliated with the Cooperative Medianet Innovation Center, Shanghai Jiao Tong University, Shanghai 200240, China (E-mail: zlin@pku.edu.cn).

Zhen Cui is with the Key Laboratory of Intelligent Perception and Systems for High-Dimensional Information, Ministry of Education, School of Computer Science and Engineering, Nanjing University of Science and Technology, Nanjing 210094, China (E-mail: zhen.cui@njust.edu.cn).

Wankou Yang is with the School of Automation, Southeast University, Nanjing 210096, China (E-mail: youngwankou@yeah.net).

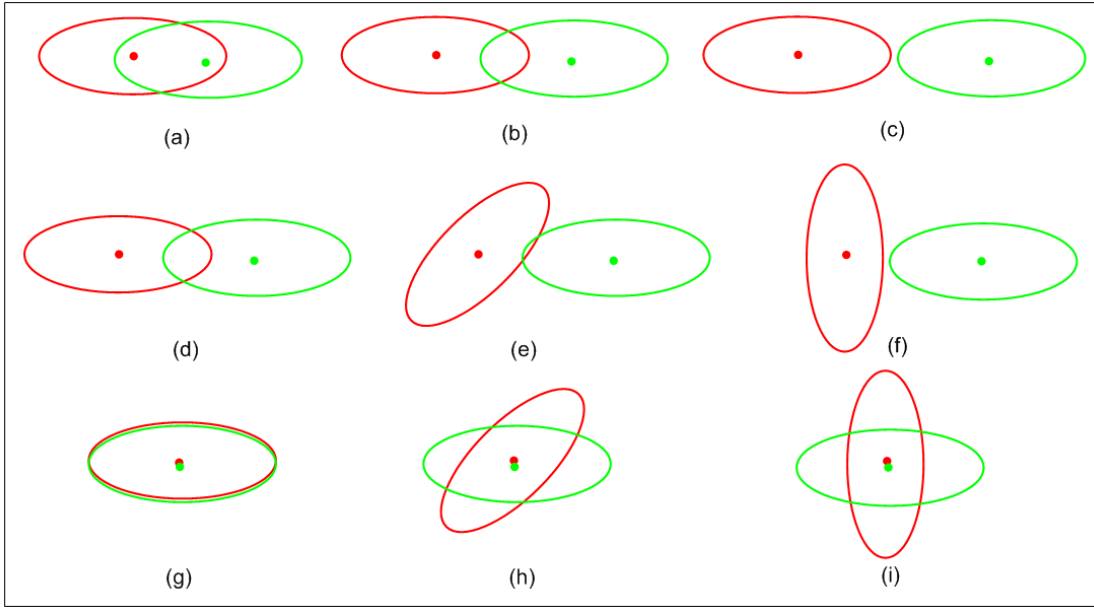


Fig. 1. The separability between two classes of data sets depicted by red color and green color, respectively. Figures (a)-(c) illustrate the examples that the separability of two classes improves with the increases of the class mean distances, whereas figures (d) - (f) and (g) - (i) illustrate the examples that the separability of two classes improves with the increase of the difference of class covariance matrices, respectively.

first one is derived under the maximum-likelihood framework [12], [13], [14]. A representative method of this category was proposed by Kumar and Andreou [12] for speech recognition. The second category is derived based on Chernoff distance or Bhattacharyya distance [1][16], where a representative method, denoted by HDA/Chernoff, is based on the Chernoff criterion [15]. A similar method to HDA/Chernoff is the approximate information discriminant analysis (AIDA) method [14], which uses the so-called μ -measure [17] as the discriminant criterion. The third category investigates hybrid linear feature extraction scheme for the heteroscedastic discriminant analysis (HDA/HLFE)[19], [18], in which the discriminative information are respectively extracted from the class means and class covariance matrices. Since HDA/HLFE is derived under the assumption of single Gaussian distribution of each class data samples, it should be noted that HDA/HLFE may not be well suitable for the cases where the class data samples abide by Gaussian mixture distribution. In addition, it is also notable that the discriminant vectors of HDA/HLFE are learned in two separated subspaces, such that the learned discriminant vectors may not be optimal in terms of Bayes error since the Bayes error is characterized by both class means and class covariance matrices simultaneously. For all of the aforementioned HDA methods, a common limitation of these methods is the suffering of the so-called small sample size problem [20], i.e., these methods all require that the number of samples in each class be larger than the dimension of the data space in order to guarantee the non-singularity of the class covariance matrices.

In addition to the HDA methods, there are other discriminant analysis approaches that have been proposed in recent years to overcome the drawbacks of FLDA, e.g., multi-view learning (or multi-modal learning) methods [50][51],

subclass methods [21][22], kernel-based methods [23], [24], or deep neural network method [52]. The multi-view learning (or multi-modal learning) methods are mainly related with the feature extraction problems of learning from the data represented by multiple distinct feature sets [53]. The subclass methods, e.g., the subclass discriminant analysis (SDA)[21], deal with the discriminative feature extraction by dividing each class samples into several subclasses, which enables this method more powerful than the FLDA method in extracting discriminative features. The kernel-based discriminant analysis methods (KDA) [23] are the nonlinear extension of FLDA via kernel trick [25] to solve the drawbacks of FLDA. In KDA, the input data samples are mapped by a nonlinear mapping from the input data space to a high-dimensional reproducing kernel Hilbert space (RKHS), such that the non-separable data samples of the input data space become separable in RKHS. As a consequence, performing feature extraction in RKHS using FLDA results in the nonlinear feature extraction in the original input data space. Similar to the kernel-based learning methods, the deep neural network methods can also extract the nonlinear features via nonlinear neural network learning.

Although the aforementioned methods are proposed to overcome the drawbacks of FLDA, most of them are developed under the Fisher's criterion, i.e., minimizing the within-class scatter distance and maximizing the between-class scatter distance. Hence, some of the limitations of Fisher's criterion, such as the difficulty of extracting the discriminant information lying in the class covariance differences, may still exist to some extent for these methods.

In this paper, we develop a new discriminant criterion under the distributions of Gaussian and mixture of Gaussians, respectively, for heteroscedastic discriminant problems, which can be expressed as a *mixture of absolute generalized Rayleigh*

quotients (MAGRQ). Preliminary applications of this work to non-frontal facial expression recognition and EEG classification had investigated in [26], [28], [27].

To show the physical meaning of our MAGRQ criterion, let's firstly consider a special two-class heteroscedastic case as shown in Fig.1, in which the first row illustrates three examples of two-class homoscedastic data sets (denoted by red and green colors), whereas the second and third rows illustrate another six examples of two-class data sets with same class means but different covariance matrices. From Fig.1, we can see that the separability of the two class data sets is closely related with both class means and class covariance matrices. In particular, it is notable that even the between-class distances are the same, e.g., the figures (d) - (f) and (g) - (i), the two-class data sets associated with the largest covariance matrices difference could be best separated. Especially, in figures (g) - (i), the class means are almost overlapped and hence the traditional FLDA would not be applicable, whereas the HDA method could largely separate the two class data sets. Fig.2 illustrates a

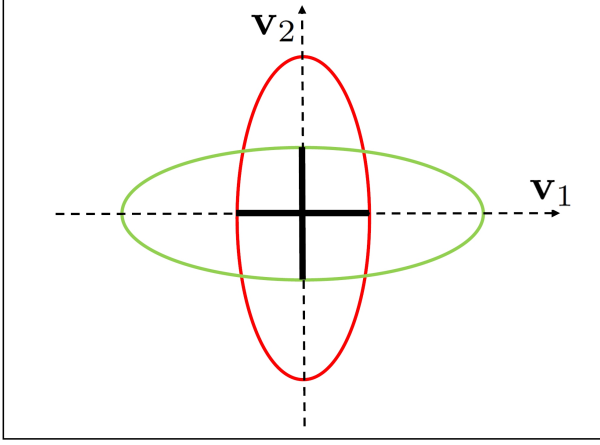


Fig. 2. An example where the class means are equal but the class covariance matrices are different. In this case, the FLDA method is not applicable because the between-class scatter matrix becomes a zero matrix.

special case of two-class heteroscedastic discriminant problem, where the class means are equal (zero) but the class covariance matrices are different. It is obvious that Fisher's criterion cannot be used in this scenario because of the zero class means. Now let \mathbf{v} denote a projection vector such that the projections of two data samples \mathbf{x} and \mathbf{y} onto this projection vector are $\mathbf{v}^T \mathbf{x}$ and $\mathbf{v}^T \mathbf{y}$, where we suppose that \mathbf{x} is from class 1 and \mathbf{y} is from class 2. Intuitively, to best distinguish the data samples between the two classes, we should minimize their overlapping parts as much as possible. To this end, we may expect that one class has smaller scatter distances whereas the other one has larger scatter distances, which means that we have to seek a projection vector \mathbf{v} such that the projection of one class has a smaller variance whereas the other one has a larger variance. This can be modeled as the following maximization problem:

$$\begin{aligned} \max_{\mathbf{v}^T \mathbf{v} = 1} |\text{var}(\mathbf{v}^T \mathbf{x}) - \text{var}(\mathbf{v}^T \mathbf{y})| &= |\mathbf{v}^T \Sigma_{\mathbf{x}} \mathbf{v} - \mathbf{v}^T \Sigma_{\mathbf{y}} \mathbf{v}| \\ &= |\mathbf{v}^T (\Sigma_{\mathbf{x}} - \Sigma_{\mathbf{y}}) \mathbf{v}|, \end{aligned} \quad (1)$$

where $\Sigma_{\mathbf{x}}$ and $\Sigma_{\mathbf{y}}$ denote the covariance matrix of classes 1 and 2, respectively. According to (1), we can obtain that the two most discriminative vectors for distinguishing between the two classes in Fig.2 should be \mathbf{v}_1 and \mathbf{v}_2 . This is because the projections of the data samples onto these two projection vectors will have a minimal overlapping parts (indicated by thick lines).

In [29], Malina proposed an extended Fisher's criterion that is expressed as the similar form of MAGRQ. Unfortunately, Malina's criterion is limited to two-class feature extraction problems and it is proposed empirically and hence lacks a rigorous theoretical justification. In contrast to Malina's criterion, our MAGRQ criterion is obtained based on a rigorous theoretical derivation. More specifically, we develop an upper bound of Bayes error under single Gaussian distribution assumption of each class data set and then extend to the case of mixture of Gaussian distributions. We also show that minimizing the upper bounds of Bayes error in both cases will result in the similar MAGRQ discriminant criterion, where a larger value of the MAGRQ criterion would lead to a smaller bound of the Bayes error. Additionally, it seems that our MAGRQ criterion is also related with the multi-view or multi-modal learning problems such as the works addressed in [50] and [51], there are significantly different between them. Specifically, the multi-view learning or multi-modal learning are mainly targeted at the problems of learning from data of multiple distinct feature sets. In contrast to these methods, the proposed MAGRQ criterion is developed under a single feature set. Moreover, the multi-view learning or multi-modal learning methods of [50] and [51] are developed without considering the discriminant information lying in the class covariance differences, whereas the proposed MAGRQ criterion aims to extract this kind of discriminant information.

In dealing with the discriminant analysis problem, it is well-known that both between-class scatter distance and within-class scatter distance can be formulated as the ℓ_2 norm operation [30]. Since the ℓ_2 norm is more sensitive to the influence of outliers, discriminant analysis based ℓ_1 norm had received increasing interests of researchers [30][31][32][33][34][35] in order to boost the robustness of the discriminant analysis methods. In [30], Wang et al. firstly introduced the ℓ_1 -norm distance metric for learning robust common spatial filters from EEG data samples contaminated by noises. The basic idea was further adopted to deal with the robust discriminative feature extraction of FLDA by Zhong et al. [31], Wang et al. [33], and Zheng et al. [32], respectively, which was referred to as the L1-FLDA method here.

Despite of the success of L1-FLDA in robust discriminative feature extraction, it is interesting to see that replacing ℓ_2 norm with ℓ_1 norm in the between-class scatter distance would be advantageous to increase the discrimination ability of FLDA, whereas it would not be a good choice to replace the ℓ_2 norm with ℓ_1 norm for the within-class scatter distance. This is because the use of ℓ_1 norm tends to suppressing the contribution of the well-separated classes (with longer between-class scatter distance) and hence can emphasize more on the classes (with smaller between-class scatter distance) that are difficult to be separated. For the within-class scatter

distance, we expect to minimize the within-scatter distance so as to achieving better discrimination, which means that we should focus more on the classes with longer within-scatter scatter distances rather than on those classes with smaller within-class scatter distances. In this sense, using the ℓ_2 norm would be more advantageous than the ℓ_1 one to describe the within-class scatter distance. Fig.3 shows an example of a set of data samples with three classes to illustrate the scenarios that more attentions should be focused on in order to achieve better discrimination, in which the distances indicted by thicker lines imply that they are more important than those indicated by thinner lines in order to separate the different classes. Consequently, to emphasize more on the classes with smaller between-class scatter distance, the ℓ_1 norm could be adopted to describe the between-class scatter distance. On the contrary, to emphasize more on the classes with larger within-class scatter distance, the ℓ_2 norm could be adopted to describe the within-class scatter distance. According to the above analysis, we extend the MAGRQ criterion by replacing the ℓ_2 norm with the ℓ_1 one in the between-class scatter distance, and hereafter propose the ℓ_1 norm based MAGRQ (L1-MAGRQ) criterion.

Based on the aforementioned L1-MAGRQ criterion, in this paper we propose a novel heteroscedastic discriminant analysis method under the mixture of Gaussian distribution (L1-HDA/GM) of each class data samples. Moreover, we also propose an efficient algorithm to solve the optimal discriminant vector sets of L1-HDA/GM, in which only the principal eigenvalue decomposition problems are involved, which can be efficiently solved by using the power iteration approach and the rank-one-update (ROU) technique [36]. Additionally, although L1-HDA/GM can be seen as the extension of our preliminary works in [26], [28], [27], it improves the previous works by using ℓ_1 norm to replace the ℓ_2 one for describing the between-class scatter distance. Specifically, under the ℓ_1 norm distance metric, the feature extraction of L1-HDA/GM

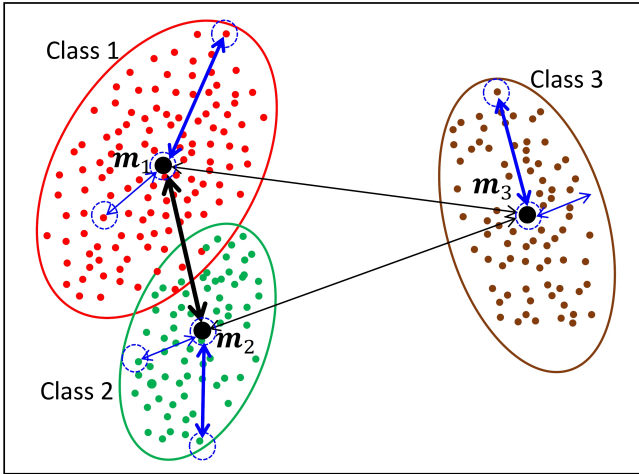


Fig. 3. An example of a set of data set with three classes to illustrate the scenarios that more attentions should be focused on, in which the distances indicted by thicker lines imply that they are more important than those indicated by thinner lines in order to separate the different classes.

will pay more attention to the non-separated pairwise classes, which makes it more powerful in extracting the discriminative features.

The remainder of this paper is organized as follows. In section II, we briefly introduce the Bayes error upper bound under both Gaussian and mixture of Gaussian distributions. In sections III, we develop the MAGRQ criterion based on Bayes error upper bound estimation, and then propose the simplified version of the L1-HDA/GM method as for the case when the number of Gaussian components is fixed at 1. In IV, we propose the complete L1-HDA/GM method. The experiments are presented in section V and section VI concludes the paper.

II. BAYES ERROR UPPER BOUND UNDER GAUSSIAN AND MIXTURE OF GAUSSIAN DISTRIBUTIONS

In this section, we briefly introduce an upper bound of the Bayes error under the single Gaussian distribution assumption and then extend it to the case of mixture of Gaussian distributions, which are the basis of deriving our MAGRQ criterion in sections III and IV, respectively.

A. Bayes Error Upper Bound Under Single Gaussian Distribution

Suppose that we are given a set of d -dimensional vector set $\mathbf{X} = \{\mathbf{x}_i^j | i = 1, \dots, c; j = 1, \dots, N_i\}$, where $\mathbf{x}_i^j \in \mathbb{R}^d$ be a sample vector, c and N_i denote the number of classes and the number of data samples in the i -th class, respectively. Let $p_i(\mathbf{x})$ and P_i denote the distribution and the prior probability of the i -th class, respectively. Assume that the distribution of the i -th class is Gaussian, i.e., $p_i(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \mathbf{m}_i, \Sigma_i)$, where $\mathcal{N}(\mathbf{x} | \mathbf{m}_i, \Sigma_i)$ is expressed by

$$\mathcal{N}(\mathbf{x} | \mathbf{m}_i, \Sigma_i) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma_i|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \mathbf{m}_i)^T \Sigma_i^{-1} (\mathbf{x} - \mathbf{m}_i) \right\},$$

\mathbf{m}_i and Σ_i denote the class mean and the class covariance matrix, respectively. Then, the Bayes error between class i and j can be expressed as [1]:

$$\varepsilon = \int \min(P_i p_i(\mathbf{x}), P_j p_j(\mathbf{x})) d\mathbf{x}, \quad (2)$$

By applying the following inequality to (2):

$$\min(a, b) \leq \sqrt{ab}, \quad \forall a, b \geq 0, \quad 0 \leq s \leq 1, \quad (3)$$

we obtain that the Bayes error can be bounded as the following form [1]:

$$\varepsilon \leq \int \sqrt{P_i P_j p_i(\mathbf{x}) p_j(\mathbf{x})} d\mathbf{x} \triangleq \sqrt{P_i P_j} \varepsilon_{ij}, \quad (4)$$

where

$$\varepsilon_{ij} = \int \sqrt{p_i(\mathbf{x}) p_j(\mathbf{x})} d\mathbf{x}. \quad (5)$$

Substituting the expressions of $p_i(\mathbf{x})$ into (5), we obtain that ε_{ij} can be calculated by

$$\varepsilon_{ij} = \exp \left\{ -\frac{1}{8} \Delta \mathbf{m}_{ij}^T \bar{\Sigma}_{ij}^{-1} \Delta \mathbf{m}_{ij} \right\} \left(\frac{\sqrt{|\Sigma_i| |\Sigma_j|}}{|\bar{\Sigma}_{ij}|} \right)^{\frac{1}{2}}, \quad (6)$$

where $\bar{\Sigma}_{ij} = \frac{1}{2} (\Sigma_i + \Sigma_j)$ and $\Delta \mathbf{m}_{ij} = \mathbf{m}_i - \mathbf{m}_j$.

B. Bayes Error Upper Bound Under Mixture of Gaussian Distributions

The aforementioned Bayes error upper bound is obtained under the assumption of single Gaussian distribution of $p_i(\mathbf{x})$. Assume that the class probability density function $p_i(\mathbf{x})$ is a mixture of Gaussians, i.e., $p_i(\mathbf{x})$ can be expressed as the form:

$$p_i(\mathbf{x}) = \sum_{r=1}^{K_i} \pi_{ir} \mathcal{N}(\mathbf{x} | \mathbf{m}_{ir}, \Sigma_{ir}), \quad (7)$$

where $0 \leq \pi_{ir} \leq 1$ ($\sum_{r=1}^{K_i} \pi_{ir} = 1$) are called the mixing coefficients, and K_i is the number of Gaussian mixture components.

Let $\mathcal{N}(\mathbf{x} | \mathbf{m}_{ir}, \Sigma_{ir}) \triangleq \mathcal{N}_{ir}$. Then, from (2), we obtain that the Bayes error between class i and j can be bounded by

$$\begin{aligned} \varepsilon &= \int \min(P_i p_i(\mathbf{x}), P_j p_j(\mathbf{x})) d\mathbf{x} \\ &\leq \sum_{r=1}^{K_i} \sum_{l=1}^{K_j} \int \min\{P_i \pi_{ir} \mathcal{N}_{ir}, P_j \pi_{jl} \mathcal{N}_{jl}\} d\mathbf{x} \\ &\leq \sum_{r=1}^{K_i} \sum_{l=1}^{K_j} \sqrt{P_i \pi_{ir} P_j \pi_{jl}} \varepsilon_{ij}^{rl}, \end{aligned} \quad (8)$$

where

$$\varepsilon_{ij}^{rl} = \exp \left\{ -\frac{1}{8} \Delta \mathbf{m}_{ij}^{rlT} (\bar{\Sigma}_{ij}^{rl})^{-1} \Delta \mathbf{m}_{ij}^{rl} \right\} \left(\frac{\sqrt{|\Sigma_{ir}| |\Sigma_{jl}|}}{|\bar{\Sigma}_{ij}^{rl}|} \right)^{\frac{1}{2}}, \quad (9)$$

where $\bar{\Sigma}_{ij}^{rl} = \frac{1}{2}(\Sigma_{ir} + \Sigma_{jl})$ and $\Delta \mathbf{m}_{ij}^{rl} = \mathbf{m}_{ir} - \mathbf{m}_{jl}$.

In what follows, we will limit our attention to deriving the MAGRQ criterion based on the Bayes error bound in (4) and (9), respectively. We first derive the MAGRQ criterion under single Gaussian distribution and then extend to the case of mixture of Gaussian distributions.

III. MAGRQ CRITERION FOR HDA UNDER SINGLE GAUSSIAN DISTRIBUTION

In this section, we develop the MAGRQ criterion based on the Bayes error upper bound estimation and then propose a novel HDA method based on this criterion.

A. MAGRQ Criterion Under Single Gaussian Distribution

Assume that the data samples of each class abide by the single Gaussian distribution, then we obtain that, when project the samples to 1D by a vector $\omega \in \mathbb{R}^d$, the distribution of the projected samples in i th class data set will become $\tilde{p}_i(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \omega^T \mathbf{m}_i, \omega^T \Sigma_i \omega)$, and the upper bound ε_{ij} becomes:

$$\begin{aligned} \varepsilon_{ij}(\omega) &= \exp \left\{ -\frac{1}{8} \frac{(\omega^T \Delta \mathbf{m}_{ij})^2}{\omega^T \bar{\Sigma}_{ij} \omega} \right\} \left(\frac{\omega^T \Sigma_i \omega \omega^T \Sigma_j \omega}{(\omega^T \bar{\Sigma}_{ij} \omega)^2} \right)^{\frac{1}{4}} \\ &= \exp \left\{ -\frac{1}{8} \frac{(\omega^T \Delta \mathbf{m}_{ij})^2}{\omega^T \bar{\Sigma}_{ij} \omega} \right\} \left(1 - \left(\frac{\omega^T \Delta \Sigma_{ij} \omega}{\omega^T \bar{\Sigma}_{ij} \omega} \right)^2 \right)^{\frac{1}{4}}, \end{aligned} \quad (10)$$

where $\Delta \mathbf{m}_{ij} = \mathbf{m}_i - \mathbf{m}_j$ and $\Delta \Sigma_{ij} = \frac{1}{2}(\Sigma_i - \Sigma_j)$.

To minimize the Bayes error, we should minimize its upper bound. Hence, based on (4) and (10), we should maximize both $\frac{(\omega^T \Delta \mathbf{m}_{ij})^2}{\omega^T \bar{\Sigma}_{ij} \omega}$ and $\frac{|\omega^T \Delta \Sigma_{ij} \omega|}{\omega^T \bar{\Sigma}_{ij} \omega}$, which results in the following two-class heteroscedastic discriminant criterion:

$$J_{ij}(\omega) = \frac{(\omega^T \Delta \mathbf{m}_{ij})^2 + |\omega^T \Delta \Sigma_{ij} \omega|}{\omega^T \bar{\Sigma}_{ij} \omega}, \quad (11)$$

We call the criterion of (11) as the *pairwise* mixture of absolute generalized Rayleigh quotient (MAGRQ) criterion. This criterion is the Malina's discriminant criterion [29] for two-class feature extraction. From the definition of $J_{ij}(\omega)$ in (11), we can see that the two-class MAGRQ criterion can be seen as the mixture of Fisher's criterion and Fukunaga-Koontz criterion [37], in which the first part corresponds to Fisher's criterion whereas the later one corresponds to the Fukunaga-Koontz criterion.

On the other hand, from the expression of (11), we obtain that the physical meaning of the MAGRQ criterion can be explained as the simultaneous optimization of the following two parts:

$$\begin{cases} \max(\omega^T \Delta \mathbf{m}_{ij})^2 + |\omega^T \Delta \Sigma_{ij} \omega|, \\ \min \omega^T \bar{\Sigma}_{ij} \omega, \end{cases} \quad (12)$$

For multi-class cases, the maximization problem of (12) can be extended by maximizing the pairwise summation of the two parts of (12), i.e.,

$$\begin{cases} \max \sum_{i,j} P_i P_j [(\omega^T \Delta \mathbf{m}_{ij})^2 + |\omega^T \Delta \Sigma_{ij} \omega|], \\ \min \sum_{i,j} P_i P_j \omega^T \bar{\Sigma}_{ij} \omega = \min 2\omega^T \bar{\Sigma} \omega, \end{cases} \quad (13)$$

where $\bar{\Sigma} = \sum_{i=1}^c P_i \Sigma_i$.

From (13), we define the multiclass MAGRQ criterion as the following form:

$$J(\omega) = \frac{\|\omega^T \mathbf{B}\|_2^2 + \sum_{i < j} P_i P_j |\omega^T \Delta \Sigma_{ij} \omega|}{\omega^T \bar{\Sigma} \omega}, \quad (14)$$

where $P_{ij} = P_i P_j$, and

$$\begin{aligned} \mathbf{B} &= [\sqrt{P_{12}} \Delta \mathbf{m}_{12}, \dots, \sqrt{P_{1c}} \Delta \mathbf{m}_{1c}, \sqrt{P_{23}} \Delta \mathbf{m}_{23}, \dots \\ &\quad \dots, \sqrt{P_{(c-1)c}} \Delta \mathbf{m}_{(c-1)c}]. \end{aligned} \quad (15)$$

The between-class scatter distance of the multiclass MAGRQ criterion in (14) is based on ℓ_2 norm, and hence it is referred to as the ℓ_2 norm based MAGRQ (L2-MAGRQ) criterion.

On the other hand, as what we have pointed out in section II, using ℓ_1 norm would be more advantageous than ℓ_2 norm in separating the different classes. Hence, we use ℓ_1 norm to replace the ℓ_2 one for describing the between-class scatter distance of $J(\omega)$, which means that we use $\|\omega^T \mathbf{B}\|_1^2$ to replace $\|\omega^T \mathbf{B}\|_2^2$ in the nominator part of $J(\omega)$, resulting in the following ℓ_1 norm based MAGRQ (L1-MAGRQ) criterion:

$$J_1(\omega) = \frac{\|\omega^T \mathbf{B}\|_1^2 + \sum_{i < j} P_i P_j |\omega^T \Delta \Sigma_{ij} \omega|}{\omega^T \bar{\Sigma} \omega}. \quad (16)$$

Based on the L2-MAGRQ criterion defined in (14) and L1-MAGRQ criterion defined in (16), we can develop two HDA methods under the Gaussian distribution, which are respectively denoted by L2-HDA/G and L1-HDA/G. In what follows, we will firstly provide the detailed algorithm description of

the L1-HDA/G method in Section III-B and then address the L2-HDA/G algorithm based on the L1-HDA/G algorithm in Section III-C.

B. L1-HDA/G Algorithm

Suppose that we want to obtain k discriminant vectors, denoted by $\omega_1, \dots, \omega_k$, of L1-HDA/G. Then we sequentially define the k discriminant vectors as follows:

Let $\omega_1, \dots, \omega_r$ be the first r discriminant vectors. Then the $(r+1)$ th discriminant vector is defined by

$$\omega_{r+1} = \arg \max_{\omega} J_1(\omega), \text{ s.t. } \omega^T \mathbf{S}_t \omega_j = 0, \forall j \leq r, \quad (17)$$

where \mathbf{S}_t is the covariance matrix of all data samples such that the discriminant vectors are statistically uncorrelated [38].

Let $\omega = \bar{\Sigma}^{-\frac{1}{2}} \alpha$ and

$$\begin{cases} \Delta \hat{\Sigma}_{ij} = P_{ij} \bar{\Sigma}^{-\frac{1}{2}} \Delta \Sigma_{ij} \bar{\Sigma}^{-\frac{1}{2}}, \\ \hat{\mathbf{B}} = \bar{\Sigma}^{-\frac{1}{2}} \mathbf{B}, \end{cases} \quad (18)$$

Then, we obtain that solving the optimization problem (17) is equivalent to solving the following optimization problem:

$$\alpha_{r+1} = \arg \max_{\alpha} \hat{J}_1(\alpha), \text{ s.t. } \alpha^T \mathbf{U}_r = \mathbf{0}^T, \quad (19)$$

where $\mathbf{U}_r = [\hat{\mathbf{S}}_t \alpha_1, \dots, \hat{\mathbf{S}}_t \alpha_r]$, $\hat{\mathbf{S}}_t = \bar{\Sigma}^{-\frac{1}{2}} \mathbf{S}_t \bar{\Sigma}^{-\frac{1}{2}}$, and

$$\hat{J}_1(\alpha) = \frac{\|\alpha^T \hat{\mathbf{B}}\|_1^2 + \sum_{i < j} |\alpha^T \Delta \hat{\Sigma}_{ij} \alpha|}{\alpha^T \alpha}. \quad (20)$$

The absolute value signs in the expression of $\hat{J}_1(\alpha)$ make the optimization of (20) difficult. So we introduce two $c \times c$ skew symmetric sign matrices $\mathbf{U} = ((\mathbf{U})_{ij})_{c \times c}$ and $\mathbf{V} = ((\mathbf{V})_{ij})_{c \times c}$, where $(\mathbf{U})_{ij}, (\mathbf{V})_{ij} \in \{+1, -1\}$, where $(\mathbf{U})_{ij}$ and $(\mathbf{V})_{ij}$ denote the i -th row and j -th column element of \mathbf{U} and \mathbf{V} , respectively. Let \mathbf{u} denote the vector by concatenating entries of \mathbf{U} according to the following form:

$$\mathbf{u} = [(\mathbf{U})_{12}, \dots, (\mathbf{U})_{1c}, (\mathbf{U})_{23}, \dots, (\mathbf{U})_{(c-1)c}]^T.$$

Denote Ω as the set of all sign matrices and define

$$\mathbf{T}(\mathbf{U}, \mathbf{V}) = \mathbf{B} \mathbf{u} \mathbf{u}^T \mathbf{B}^T + \sum_{i < j} (\mathbf{V})_{ij} \Delta \hat{\Sigma}_{ij}. \quad (21)$$

Then we obtain that

$$\begin{aligned} \alpha^T \mathbf{T}(\mathbf{U}, \mathbf{V}) \alpha &= (\alpha^T \mathbf{B} \mathbf{u})^2 + \sum_{i < j} (\mathbf{V})_{ij} \alpha^T \Delta \hat{\Sigma}_{ij} \alpha \\ &\leq \|\alpha^T \mathbf{B}\|_1^2 + \sum_{i < j} |\alpha^T \Delta \hat{\Sigma}_{ij} \alpha|. \end{aligned} \quad (22)$$

From (22), we obtain that the optimization problem of (20) can be formulated as the following one:

$$\hat{J}_1(\alpha) = \max_{\mathbf{U}, \mathbf{V} \in \Omega, \|\alpha\|=1} \alpha^T \mathbf{T}(\mathbf{U}, \mathbf{V}) \alpha. \quad (23)$$

From (23), we obtain that

$$\begin{aligned} \max_{\alpha} \hat{J}_1(\alpha) &= \max_{\|\alpha\|=1} \max_{\mathbf{U}, \mathbf{V} \in \Omega} \alpha^T \mathbf{T}(\mathbf{U}, \mathbf{V}) \alpha \\ &= \max_{\mathbf{U}, \mathbf{V} \in \Omega} \max_{\|\alpha\|=1} \alpha^T \mathbf{T}(\mathbf{U}, \mathbf{V}) \alpha. \end{aligned} \quad (24)$$

By observing (24), we can see that: If the sign matrix \mathbf{U} and \mathbf{V} are fixed, then the optimal discriminant vector is the *normalized* (we will not re-emphasize this in the sequel) eigenvector associated with the largest eigenvalue of the matrix $\mathbf{T}(\mathbf{U}, \mathbf{V})$. Solving the principal eigenvector of $\mathbf{T}(\mathbf{U}, \mathbf{V})$ can be easily realized via the power iteration method. So our problem of maximizing $\hat{J}_1(\alpha)$ is changed to finding the optimal sign matrices \mathbf{U} and \mathbf{V} such that the largest eigenvalue of $\mathbf{T}(\mathbf{U}, \mathbf{V})$ is maximized.

In what follows, we will propose a greedy algorithm to find the suboptimal sign matrices \mathbf{U} and \mathbf{V} . To begin with, we introduce the following theorem:

Theorem 1. Let $\alpha^{(1)}$ be the principal eigenvector of $\mathbf{T}(\mathbf{U}_1, \mathbf{V}_1)$. Define \mathbf{U}_2 and \mathbf{V}_2 as

$$\begin{cases} (\mathbf{U}_2)_{ij} = \text{sign}(\alpha^{(1)T} \Delta \hat{\mathbf{m}}_{ij}), \\ (\mathbf{V}_2)_{ij} = \text{sign}(\alpha^{(1)T} \Delta \hat{\Sigma}_{ij} \alpha^{(1)}), \end{cases} \quad (25)$$

where

$$\text{sign}(a) = \begin{cases} 1, & \text{if } a \geq 0 \\ -1, & \text{Others.} \end{cases}$$

Suppose that $\alpha^{(2)}$ is the principal eigenvector of $\mathbf{T}(\mathbf{U}_2, \mathbf{V}_2)$. Then, we have

$$\alpha^{(2)T} \mathbf{T}(\mathbf{U}_2, \mathbf{V}_2) \alpha^{(2)} \geq \alpha^{(1)T} \mathbf{T}(\mathbf{U}_1, \mathbf{V}_1) \alpha^{(1)}. \quad (26)$$

Proof: See appendix A. \square

Thanks to Theorem 1, we are able to improve the sign matrices step by step.

To solve the discriminant vector α_{r+1} , we introduce Propositions 1 and 2 below. Their proofs can be easily obtained from [39]:

Proposition 1. Let $\mathbf{Q}_r \mathbf{R}_r$ be the QR decomposition of \mathbf{U}_r , where \mathbf{R}_r is an $r \times r$ upper triangular matrix. Then α_{r+1} defined in (19) is the principal eigenvector corresponding to the largest eigenvalue of the following matrix $(\mathbf{I}_d - \mathbf{Q}_r \mathbf{Q}_r^T) \mathbf{T}(\mathbf{U}, \mathbf{V}) (\mathbf{I}_d - \mathbf{Q}_r \mathbf{Q}_r^T)$.

Proposition 2. Suppose that $\mathbf{Q}_r \mathbf{R}_r$ is the QR decomposition of \mathbf{U}_r . Let $\mathbf{U}_{r+1} = (\mathbf{U}_r \hat{\mathbf{S}}_t \alpha_{r+1})$, $\mathbf{q} = \hat{\mathbf{S}}_t \alpha_{r+1} - \mathbf{Q}_r (\mathbf{Q}_r^T \hat{\mathbf{S}}_t \alpha_{r+1})$, and $\mathbf{Q}_{r+1} = \begin{pmatrix} \mathbf{Q}_r & \frac{\mathbf{q}}{\|\mathbf{q}\|} \end{pmatrix}$. Then $\mathbf{Q}_{r+1} \begin{pmatrix} \mathbf{R}_r & \mathbf{Q}_r^T \hat{\mathbf{S}}_t \alpha_{r+1} \\ \mathbf{0} & \|\mathbf{q}\| \end{pmatrix}$ is the QR decomposition of \mathbf{U}_{r+1} .

The above two propositions provide an efficient approach for solving (19): Proposition 1 makes it possible to use the power method to solve (19), while Proposition 2 makes it possible to update \mathbf{Q}_{r+1} from \mathbf{Q}_r by adding a single column. Moreover, it should be noted that

$$\begin{aligned} \mathbf{I}_d - \mathbf{Q}_{r+1} \mathbf{Q}_{r+1}^T &= \prod_{i=1}^{r+1} (\mathbf{I}_d - \mathbf{q}_i \mathbf{q}_i^T) \\ &= (\mathbf{I}_d - \mathbf{Q}_r \mathbf{Q}_r^T) (\mathbf{I}_d - \mathbf{q}_{r+1} \mathbf{q}_{r+1}^T), \end{aligned} \quad (27)$$

where \mathbf{q}_i is the i th column of \mathbf{Q}_r . Eqn. (27) makes it possible to update $(\mathbf{I}_d - \mathbf{Q}_{r+1} \mathbf{Q}_{r+1}^T) \mathbf{T}(\mathbf{U}_r, \mathbf{V}_r) (\mathbf{I}_d - \mathbf{Q}_{r+1} \mathbf{Q}_{r+1}^T)$ from $(\mathbf{I}_d - \mathbf{Q}_r \mathbf{Q}_r^T) \mathbf{T}(\mathbf{U}_r, \mathbf{V}_r) (\mathbf{I}_d - \mathbf{Q}_r \mathbf{Q}_r^T)$ by the ROU technique.

Here it should be noted that the initial setting of the sign matrices \mathbf{U} and \mathbf{V} in $\mathbf{T}(\mathbf{U}, \mathbf{V})$ may influence the optimality

Algorithm 1: Solving optimal vectors ω_i ($i = 1, \dots, k$) of L1-HDA/G.

Input:

- Input data set $\{\mathbf{x}_i^j | i = 1, \dots, c; j = 1, \dots, N_i\}$, class label vector \mathbf{L} , where $N_1 + \dots + N_c = N$.

Initialization:

- Compute $\Sigma_i, \bar{\Sigma}_{ij}, \bar{\Sigma}, \mathbf{B}, \mathbf{S}^t, \hat{\Sigma}_i = \bar{\Sigma}^{-\frac{1}{2}} \Sigma_i \bar{\Sigma}^{-\frac{1}{2}}, \hat{\mathbf{B}} = \bar{\Sigma}^{-\frac{1}{2}} \mathbf{B}, \hat{\mathbf{S}}_t = \bar{\Sigma}^{-\frac{1}{2}} \mathbf{S}_t \bar{\Sigma}^{-\frac{1}{2}}, \mathbf{m}_i, \hat{\mathbf{m}}_i = \bar{\Sigma}^{-\frac{1}{2}} \mathbf{m}_i$;
- Initialize the discriminant vectors α_i ;

For $i = 1, 2, \dots, k$, **Do**

- 1) Compute $\Delta \hat{\Sigma}_{pq} \leftarrow \frac{\hat{\Sigma}_p - \hat{\Sigma}_q}{2}$, ($p < q$);
- 2) Set $\mathbf{U}, \mathbf{V} \leftarrow$ zero matrix, and compute

$$\begin{cases} (\mathbf{U}_1)_{pq} \leftarrow \text{sign}(\alpha_i^T \Delta \hat{\mathbf{m}}_{pq}); \\ (\mathbf{V}_1)_{pq} \leftarrow \text{sign}(\alpha_i^T \Delta \hat{\Sigma}_{pq} \alpha_i). \end{cases}$$
- 3) **While** $\mathbf{U} \neq \mathbf{U}_1$ **and** $\mathbf{V} \neq \mathbf{V}_1$, **Do**
 - a) Set $\mathbf{U} \leftarrow \mathbf{U}_1$ and $\mathbf{V} \leftarrow \mathbf{V}_1$, compute $\mathbf{T}(\mathbf{U}, \mathbf{V})$ and the principal eigenvector, α_i , of $\mathbf{T}(\mathbf{U}, \mathbf{V})$;
 - b) Compute

$$\begin{cases} (\mathbf{U}_1)_{pq} \leftarrow \text{sign}(\alpha_i^T \Delta \hat{\mathbf{m}}_{pq}); \\ (\mathbf{V}_1)_{pq} \leftarrow \text{sign}(\alpha_i^T \Delta \hat{\Sigma}_{pq} \alpha_i). \end{cases}$$
- 4) Update \mathbf{Q}_i : $\mathbf{Q}_i \leftarrow (\mathbf{Q}_{i-1} \frac{\mathbf{q}_i}{\|\mathbf{q}_i\|_2})$, where

$$\begin{cases} \mathbf{q}_1 \leftarrow \hat{\mathbf{S}}_t \alpha_1 \text{ and } \mathbf{Q}_1 \leftarrow \frac{\mathbf{q}_1}{\|\mathbf{q}_1\|_2}, & \text{if } i=1; \\ \mathbf{q}_i \leftarrow \hat{\mathbf{S}}_t \alpha_i - \mathbf{Q}_{i-1}(\mathbf{Q}_{i-1}^T \hat{\mathbf{S}}_t \alpha_i), & \text{otherwise.} \end{cases}$$
- 5) Update $\hat{\Sigma}_p$ and $\hat{\mathbf{B}}$:

$$\hat{\Sigma}_p \leftarrow (\mathbf{I} - \mathbf{q}_i \mathbf{q}_i^T) \hat{\Sigma}_p (\mathbf{I} - \mathbf{q}_i \mathbf{q}_i^T), \hat{\mathbf{B}} \leftarrow (\mathbf{I} - \mathbf{q}_i \mathbf{q}_i^T) \hat{\mathbf{B}};$$
- 6) Compute $\omega_i = \bar{\Sigma}^{-\frac{1}{2}} \alpha_i$, and set $\omega_i \leftarrow \omega_i / \|\omega_i\|$;

Output: $\omega_1, \dots, \omega_k$.

of the solution. Consequently, to obtain a better solution, we may initialize the sign matrices \mathbf{U} and \mathbf{V} based on the optimal discriminant vectors solved by other discriminant analysis algorithms. For example, suppose that α is the optimal discriminant vector solved by the HDA/Chernoff algorithm [15], then the initial sign matrix of \mathbf{U} and \mathbf{V} can be obtained as follows:

$$(\mathbf{U})_{pq} \leftarrow \text{sign}(\alpha^T \Delta \hat{\mathbf{m}}_{pq}), \quad (28)$$

$$(\mathbf{V})_{pq} \leftarrow \text{sign}(\alpha^T \Delta \hat{\Sigma}_{pq} \alpha). \quad (29)$$

We summarize the algorithm for solving the first k discriminant vectors of L1-HDA/G in Algorithm 1.

C. L2-HDA/G Algorithm

In the aforementioned section, we had developed a set of optimal discriminative vectors of L1-HDA/G based on the L1-MAGRQ criterion. Similarly, if the L2-MAGRQ criterion is adopted, then we can obtain an optimal set of L2-HDA/G discriminative vectors. Specifically, the optimal discriminative vectors of L2-HDA/G can be sequentially obtained: suppose that we have obtained the first r optimal discriminant vectors of L2-HDA/G, denoted by $\omega_1, \dots, \omega_r$, then the $(r+1)$ th

discriminant vector is defined by

$$\omega_{r+1} = \arg \max_{\omega} J(\omega), \text{ s.t. } \omega^T \mathbf{S}_t \omega_j = 0, \forall j \leq r. \quad (30)$$

The optimization problem of (30) is equivalent of the following one:

$$\alpha_{r+1} = \arg \max_{\alpha^T \mathbf{U} = 0^T} \hat{J}(\alpha), \quad (31)$$

where

$$\hat{J}(\alpha) = \frac{\|\alpha^T \hat{\mathbf{B}}\|_2^2 + \sum_{i < j} |\alpha^T \Delta \hat{\Sigma}_{ij} \alpha|}{\alpha^T \alpha}, \quad (32)$$

in which \mathbf{U}_r and $\hat{\mathbf{B}}$ are defined in Section III-B.

Noting that the ℓ_2 norm metric can be easily computed without the absolute value operation, the expression of $\mathbf{T}(\mathbf{U}, \mathbf{V})$ in (21) can be replaced by $\mathbf{T}(\mathbf{V})$ defined as follows:

$$\mathbf{T}(\mathbf{V}) = \mathbf{B} \mathbf{B}^T + \sum_{i < j} (\mathbf{V})_{ij} \Delta \hat{\Sigma}_{ij}. \quad (33)$$

As a result, the L2-HDA/G algorithm can be obtained with a simple modification of the L1-HDA/G algorithm shown in Algorithm 1, which can be summarized in the following Algorithm 2.

Algorithm 2: Solving optimal vectors ω_i ($i = 1, \dots, k$) of L2-HDA/G.

Input:

- Input data set $\{\mathbf{x}_i^j | i = 1, \dots, c; j = 1, \dots, N_i\}$, class label vector \mathbf{L} , where $N_1 + \dots + N_c = N$.

Initialization:

- Compute $\Sigma_i, \bar{\Sigma}_{ij}, \bar{\Sigma}, \mathbf{B}, \mathbf{S}^t, \hat{\Sigma}_i = \bar{\Sigma}^{-\frac{1}{2}} \Sigma_i \bar{\Sigma}^{-\frac{1}{2}}, \hat{\mathbf{B}} = \bar{\Sigma}^{-\frac{1}{2}} \mathbf{B}, \hat{\mathbf{S}}_t = \bar{\Sigma}^{-\frac{1}{2}} \mathbf{S}_t \bar{\Sigma}^{-\frac{1}{2}}, \mathbf{m}_i, \hat{\mathbf{m}}_i = \bar{\Sigma}^{-\frac{1}{2}} \mathbf{m}_i$;
- Initialize the discriminant vectors α_i ;

For $i = 1, 2, \dots, k$, **Do**

- 1) Compute $\Delta \hat{\Sigma}_{pq} \leftarrow \frac{\hat{\Sigma}_p - \hat{\Sigma}_q}{2}$, ($p < q$);
- 2) Set $\mathbf{V} \leftarrow$ zero matrix, and compute

$$(\mathbf{V}_1)_{pq} \leftarrow \text{sign}(\alpha_i^T \Delta \hat{\Sigma}_{pq} \alpha_i).$$

3) **While** $\mathbf{V} \neq \mathbf{V}_1$, **Do**

- a) Set $\mathbf{V} \leftarrow \mathbf{V}_1$ and compute $\mathbf{T}(\mathbf{V})$ and the principal eigenvector, α_i , of $\mathbf{T}(\mathbf{V})$;
- b) Compute $(\mathbf{V}_1)_{pq} \leftarrow \text{sign}(\alpha_i^T \Delta \hat{\Sigma}_{pq} \alpha_i)$.

4) Update \mathbf{Q}_i : $\mathbf{Q}_i \leftarrow (\mathbf{Q}_{i-1} \frac{\mathbf{q}_i}{\|\mathbf{q}_i\|_2})$, where

$$\begin{cases} \mathbf{q}_1 \leftarrow \hat{\mathbf{S}}_t \alpha_1 \text{ and } \mathbf{Q}_1 \leftarrow \frac{\mathbf{q}_1}{\|\mathbf{q}_1\|_2}, & \text{if } i=1; \\ \mathbf{q}_i \leftarrow \hat{\mathbf{S}}_t \alpha_i - \mathbf{Q}_{i-1}(\mathbf{Q}_{i-1}^T \hat{\mathbf{S}}_t \alpha_i), & \text{otherwise.} \end{cases}$$

5) Update $\hat{\Sigma}_p$ and $\hat{\mathbf{B}}$:

$$\hat{\Sigma}_p \leftarrow (\mathbf{I} - \mathbf{q}_i \mathbf{q}_i^T) \hat{\Sigma}_p (\mathbf{I} - \mathbf{q}_i \mathbf{q}_i^T), \hat{\mathbf{B}} \leftarrow (\mathbf{I} - \mathbf{q}_i \mathbf{q}_i^T) \hat{\mathbf{B}};$$

6) Compute $\omega_i = \bar{\Sigma}^{-\frac{1}{2}} \alpha_i$, and set $\omega_i \leftarrow \omega_i / \|\omega_i\|$;

Output: $\omega_1, \dots, \omega_k$.

TABLE I
COMPARISON OF COMPUTATIONAL COMPLEXITY BETWEEN L1-HDA/G ALGORITHM AND L2-HDA/G ALGORITHM.

Algorithm	Computational complexity of each step in the algorithm						
	Initialization	1)	2)	3)	4)	5)	6)
Algorithm 1	$O(cd^2N)+O(cd^3)$	$O(c^2)$	$O(c^2d^2)$	$O(c^2d)+O(c^2d^2)$	$O(d^2)$	$O(d^2)$	$O(d^2)$
Algorithm 2	$O(cd^2N)+O(cd^3)$	$O(c^2)$	$O(c^2d^2)$	$O(d^3)+O(c^2d^2)$	$O(d^2)$	$O(d^2)$	$O(d^2)$

D. Computational Analysis of L1-HDA/G and L2-HDA/G

According to the detailed algorithm description of L1-HDA/G shown in Algorithm 1, we can obtain the computational complexity of L1-HDA/G algorithm. Specifically, in the initialization part, the computational complexity of calculating the class covariance matrices is $O(cd^2N)$ and the computational complexity of calculating the transformation matrices $(\hat{\Sigma}_i, \hat{\mathbf{B}}, \text{ and } \hat{\mathbf{S}}_i)$ is $O(cd^3)$. In calculating each discriminative vector ω_i , the computational complexity of step 1) and 2) are $O(d^2)$ and $O(c^2d^2)$, respectively. The computational complexity of calculating $\mathbf{T}(\mathbf{U}, \mathbf{V})$ in step 3) is $O(c^2d)+O(c^2d^2)$, and the complexity of solving the principal eigenvector of $\mathbf{T}(\mathbf{U}, \mathbf{V})$ is only $O(d^2)$ (e.g., using power method). The complexity of updating \mathbf{U} and \mathbf{V} in step 3) is $O(c^2d)+O(c^2d^2)$. In addition, it is easy to check that the computational complexity of step 4), 5), and 6) are $O(d^2)$, $O(d^2)$, and $O(d^2)$, respectively. According to the aforementioned analysis, we summarize the computational complexity of Algorithm 1 in Table I.

In contrast to the L1-HDA/G algorithm, the major difference of L2-HDA/G lies in the calculation of $\mathbf{T}(\mathbf{U}, \mathbf{V})$ (replaced by $\mathbf{T}(\mathbf{V})$ in L2-HDA/G). The computational complexity of calculating $\mathbf{T}(\mathbf{V})$ is $O(d^3)+O(c^2d^2)$, which would be a bit more than calculating the value of $\mathbf{T}(\mathbf{U}, \mathbf{V})$ in the L1-HDA/G algorithm. The detailed computational complexity of L2-HDA/G is also summarized in Table I.

IV. L1-MAGRQ CRITERION UNDER MIXTURE OF GAUSSIAN DISTRIBUTIONS AND L1-HDA/GM

In this section, we generalize the L2-MAGRQ criterion and the L1-MAGRQ criterion from the single Gaussian distribution to mixture of Gaussian distributions. Then, we propose the L2-HDA/GM method and the L1-HDA/GM method. If the data samples of each class abide by mixture of Gaussian distributions, then we have the following theorem with respect to the projected samples:

Theorem 2. Suppose that the distribution function of the i th class is a mixture of Gaussians, i.e.,

$$p_i(\mathbf{x}) = \sum_{r=1}^{K_i} \pi_{ir} \mathcal{N}(\mathbf{x} | \mathbf{m}_{ir}, \Sigma_{ir}).$$

Then, the class distribution function $\tilde{p}_i(\omega^T \mathbf{x})$ of the projected samples $\omega^T \mathbf{x}$ is also a mixture of Gaussians, i.e.,

$$\tilde{p}_i(\omega^T \mathbf{x}) = \sum_{r=1}^{K_i} \pi_{ir} \mathcal{N}(\omega^T \mathbf{x} | \omega^T \mathbf{m}_{ir}, \omega^T \Sigma_{ir} \omega). \quad (34)$$

Proof: See appendix B. \square

Thanks to Theorem 2, we obtain that the two-class Bayes error bound expressed in (8) can be replaced by

$$\varepsilon \leq \sum_{r=1}^{K_1} \sum_{l=1}^{K_j} \sqrt{P_i \pi_{ir} P_j \pi_{jl}} \varepsilon_{ij}^{rl}(\omega), \quad (35)$$

where $\varepsilon_{ij}^{rl}(\omega)$ is formulated as:

$$\varepsilon_{ij}^{rl}(\omega) = \exp \left\{ -\frac{1}{8} \frac{(\omega^T \Delta \mathbf{m}_{ij}^{rl})^2}{\omega^T \bar{\Sigma}_{ij}^{rl} \omega} \right\} \left(1 - \left(\frac{\omega^T \Delta \Sigma_{ij}^{rl} \omega}{\omega^T \bar{\Sigma}_{ij}^{rl} \omega} \right)^2 \right)^{\frac{1}{4}}, \quad (36)$$

Similar to the derivation in section III, from the Bayes error upper bound shown in (35) and (36), we obtain the following two-class MAGRQ criterion under the mixture of Gaussian distributions:

$$J_{ij}(\omega) = \sum_{r,l} \pi_{ir} \pi_{jl} \frac{(\omega^T \Delta \mathbf{m}_{ij}^{rl})^2 + |\omega^T \Delta \Sigma_{ij}^{rl} \omega|}{\omega^T \bar{\Sigma}_{ij}^{rl} \omega}. \quad (37)$$

Note that in real applications, the number of samples are often insufficient for estimating a mixture of Gaussians with different Σ_{ir} 's. To remedy this issue, we may assume that the matrices $\{\Sigma_{ir}\}$ are identical, say equal to Σ_i . In this case, the two-class MAGRQ criterion in (37) can be expressed as the following form as for the case of the mixture of Gaussian distribution:

$$J_{ij}(\omega) = \frac{\sum_{r,l} \pi_{ir} \pi_{jl} (\omega^T \Delta \mathbf{m}_{ij}^{rl})^2}{\omega^T \bar{\Sigma}_{ij} \omega} + \frac{|\omega^T \Delta \Sigma_{ij} \omega|}{\omega^T \bar{\Sigma}_{ij} \omega}. \quad (38)$$

where $\Delta \Sigma_{ij}$ and $\bar{\Sigma}_{ij}$ are the same as those defined in section III-A.

For multiclass case, we define the following L2-MAGRQ criterion based on the two-class MAGRQ criterion of (38):

$$J(\omega) \triangleq \frac{\|\omega^T \tilde{\mathbf{B}}\|_2^2}{\omega^T \bar{\Sigma} \omega} + \frac{\sum_{i < j} P_{ij} |\omega^T \Delta \Sigma_{ij} \omega|}{\omega^T \bar{\Sigma} \omega}, \quad (39)$$

where

$$\tilde{\mathbf{B}} = [\sqrt{P_{12}} \mathbf{B}_{12}, \dots, \sqrt{P_{1c}} \mathbf{B}_{1c}, \sqrt{P_{23}} \mathbf{B}_{23}, \dots, \dots, \sqrt{P_{(c-1)c}} \mathbf{B}_{(c-1)c}], \quad (40)$$

$$\mathbf{B}_{ij} = [\sqrt{\pi_{i1} \pi_{j1}} \Delta \mathbf{m}_{ij}^{11}, \sqrt{\pi_{i1} \pi_{j2}} \Delta \mathbf{m}_{ij}^{12}, \dots, \dots, \sqrt{\pi_{iK_i} \pi_{jK_j}} \Delta \mathbf{m}_{ij}^{N_i N_j}]. \quad (41)$$

In this case, we can obtain the following L1-MAGRQ criterion under the mixture of Gaussian distribution:

$$J_1(\omega) \triangleq \frac{\|\omega^T \tilde{\mathbf{B}}\|_1^2}{\omega^T \bar{\Sigma} \omega} + \frac{\sum_{i < j} P_{ij} |\omega^T \Delta \Sigma_{ij} \omega|}{\omega^T \bar{\Sigma} \omega}. \quad (42)$$

Finally, based on the L1-MAGRQ criterion and the L2-MAGRQ criterion defined in (42) and (39), we can define the optimal discriminant vector set of L1-HDA/GM and L2-HDA/GM, respectively, which are similar as those defined in (17) and (30), respectively, and the solution methods are also similar as those shown in Algorithms 1 and Algorithm 2, respectively.

V. EXPERIMENTS

In this section, we will evaluate the discriminant performance of the proposed methods on four real data databases, i.e., the Multi-PIE and the BU-3DFE facial expression databases [42][45], the EEG database in “BCI competition 2005” - data set IIIa [47], and the UCI database [56]. Since L1-HDA/G is only a special case of L1-HDA/GM as for the case when the number of Gaussian components equals to 1, in the following experiments we only adopt the L1-HDA/GM method to conduct the experiments. Moreover, we also use the L2-HDA/GM method to conduct the same experiments. For comparison purpose, several state-of-the-art discriminant analysis methods are adopted to conduct the same experiments, which include the FLDA method [5], the AIDA method [14], the HDA/Chernoff method [15], the HDA/HLFE method [18], the SDA method [21], and the MvDN method [52]. In addition, we also conduct the experiment without any feature extraction and refer it as the Baseline method. Noting that the experiments aim to evaluate the discriminative feature extraction performance of the various methods, in the following experiments we only adopt simple classifiers, such as K-nearest neighbor (KNN) and the linear classifier, to produce the classification results in order to compare the discriminative power of the extracted features¹. In real applications, one may use more complex classifiers, such as SVM [40] and Adaboost [41], to further enhance the classification performance.

A. Experiment on Multi-PIE Facial Expression Database

In this experiment, we will use the famous Multi-PIE database [42] to evaluate the performance of the various methods. This is a multi-view facial expression database consisting of 755,370 facial images of 337 subjects. Similar to [43], we choose 4200 facial images from 100 subjects as those previously used in [43] to conduct the experiment, in which each subject contains 42 facial images that cover 7 facial views (0° , 15° , 30° , 45° , 60° , 75° , and 90°) and 6 facial expressions (Disgust, Neutral, Scream, Smile, Squint, and Surprise). Fig. 4 shows 42 facial images covering the 6 facial expressions and 7 facial views of one subject in the Multi-PIE database.

We explore two kinds of facial feature extraction schemes to evaluate the proposed methods. The first one is to adopt local binary patterns (LBP) [44] to extract 5015 features from each facial image, and the second one is to extract deep learning features learned via deep neural networks as used in face recognition [55]. The details of extracting both kinds of features are summarized as follows:

¹Otherwise, one may not be able to separate the contribution from feature extraction and that from classifier.

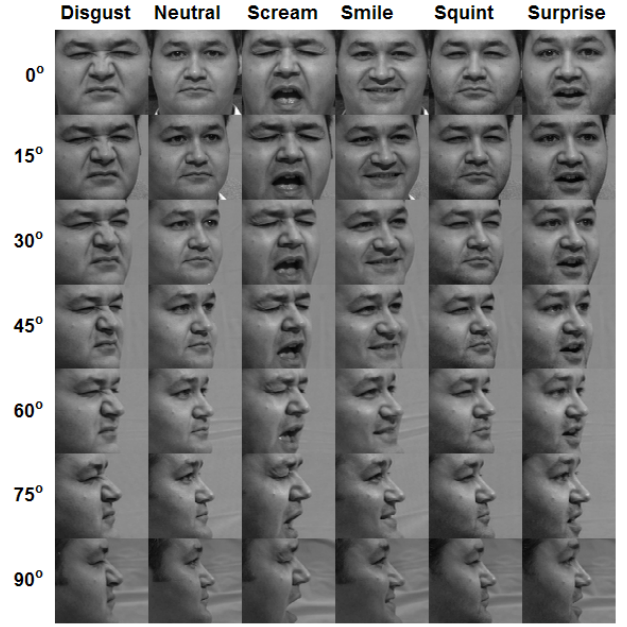


Fig. 4. Examples of the 42 facial images covering 6 facial expressions and 7 facial views of one subject in the Multi-PIE database.

- To extract the LBP facial features, we use the multi-scale face region division scheme proposed in [43] to obtain 85 facial regions and then extract a 59-dimensional LBP feature vectors from each region, which results in 85 LBP feature vectors for each facial image. Finally, we concatenate all the 85 LBP feature vectors into a 5015-dimensional feature vector.
- To extract the deep learning features, we focus our attentions to the existing deep learning neural network model that had been successfully used in extracting facial features. For this purpose, we firstly utilize the real-world affective faces (RAF) database [54] to fine-tune the deep-face VGG model (VGG-Face) [55]. Then, based on the fine-tuned VGG-Face model, we further extract the deep facial expression features by fitting the facial images of Multi-PIE database into the model. In this way, we finally obtain a set of facial features with dimensionality of 4096 taken from fc7 layer of the VGG-Face model.

In order to visualize the data distribution associated with each facial expression, we project the LBP feature points associated with the same facial expression onto the first two principal components obtained by principal component analysis (PCA) and then depict the distribution of the projection points. Fig. 5 shows the distributions of projection points with respect to all the 6 facial expression classes. From Fig. 5 we can see that, due to the multiview property of the facial images, the distribution of the data points associated with each facial expression demonstrates a multimodal form with 7 clusters. In this case, it would be very possible that using single Gaussian function could not be enough to characterize the distribution of the multiview facial feature points. For this reason, for both L1-HDA/GM and L2-HDA/GM methods, we use the mixture of Gaussians to describe the distribution of the facial feature points, in which the number of Gaussian components is set

TABLE II
AVERAGE CLASSIFICATION RATES (%) OF VARIOUS METHODS WITH RESPECT TO EACH FACIAL EXPRESSION ON THE MULTI-PIE DATABASE.

		Disgust	Neutral	Scream	Smile	Squint	Surprise	Overall
# samples		700	700	700	700	700	700	4200
# subjects		100	100	100	100	100	100	100
dimensionality of LBP feature		5015	5015	5015	5015	5015	5015	5015
dimensionality of deep learning feature		4096	4096	4096	4096	4096	4096	4096
LBP Feature	Baseline	58.93	71.14	76.93	68.29	41.07	80.64	66.17
	FLDA	73.57	74.21	90.71	79.29	66.21	89.86	78.98
	AIDA [14]	63.43	74.43	85.64	76.14	65.86	91.36	76.14
	HDA/HLFE [18]	69.29	78.64	90.00	78.36	66.43	87.00	78.29
	HDA/Chernoff [15]	64.93	76.21	86.07	76.07	66.57	90.29	76.69
	SDA [21]	73.14	77.79	90.79	79.79	73.86	91.71	81.18
	MvDN [52]	66.00	77.50	91.00	75.00	80.50	95.50	80.92
	L2-HDA/GM	74.00	78.50	90.64	80.57	73.07	92.07	81.48
	L1-HDA/GM	73.86	78.64	89.93	80.86	73.86	92.79	81.65
Deep Learning Feature	Baseline	29.36	79.29	78.00	65.71	49.21	85.14	64.45
	FLDA	68.07	68.00	92.14	72.29	60.29	88.14	74.82
	AIDA [14]	32.29	58.21	71.00	48.14	35.79	69.64	52.51
	HDA/HLFE [18]	63.43	66.71	91.50	73.86	59.93	88.50	73.99
	HDA/Chernoff [15]	55.86	62.43	87.57	67.00	48.36	82.14	67.23
	SDA [21]	66.5	70.14	91.86	71.00	64.64	89.50	75.61
	MvDN [52]	62.00	68.00	91.50	79.00	76.50	91.50	78.42
	L2-HDA/GM	67.43	70.79	92.07	73.29	65.29	90.07	76.49
	L1-HDA/GM	68.79	70.00	92.14	73.93	64.50	89.57	76.49

TABLE III
THE COMPUTATIONAL EFFICIENCY OF THE VARIOUS METHODS IN TERMS OF THE CPU RUNNING TIME (SECOND) IN THE TRAINING STAGE ON THE MULTI-PIE DATABASE.

Feature	Baseline	FLDA	AIDA	HDA/HLFE	HDA/Chernoff	SDA	MvDN	L2-HDA/GM	L1-HDA/GM
LBP Feature	0.01	0.45	0.52	1.38	2.43	0.65	1209	11.16	9.62
DL Feature	0.02	3.71	14.29	61.95	132.90	5.19	2566	77.05	65.92

to be the number of facial views (= 7) and each Gaussian component corresponds to the data samples of one facial view.

To evaluate the recognition performance of the various methods, we adopted the experimental protocol used in [43] to carry out this experiment. According to this protocol, cross-validation strategy is used to design the experiment, in which we randomly partition the 100 subjects into two subsets, in which the first subset contains the facial images of 80 subjects and the second subset contains the facial images of 20 subjects. Then, we choose the first subset as the training data set and the second subset as the testing data set. Consequently, the training data set contains a total of 3360 facial images whereas the testing one contains 840 facial images. Then we train the discriminant vectors of the various discriminant algorithms on the training data set and evaluate the performance using the testing data set. Moreover, considering that the dimensionality of the feature space is relatively larger than the number of each class samples, a PCA operation on the training data set is used to reduce the dimensionality of the feature vectors, such that the covariance matrix of each class data samples is non-singular. We conduct 10 trials of experiments in this database, and in each of the 10 trials, new training and testing data sets are chosen to evaluate the recognition performance

of the various methods. Finally, we average the results of all the experimental trials to obtain the final recognition rate. Table II shows the average recognition rates with respect to each facial expression and the overall recognition rates of the various methods on the Multi-PIE facial expression database.

From Tables II, we observe the following three major points:

- The average recognition accuracies of using LBP features are higher than those of using deep learning features. This is most likely due to the fact that VGG-face model is fine-tuned using other facial expression database, i.e., the RAF database, instead of the Multi-PIE database. Since RAF database is irrelevant to the facial expression images to be tested, the extracted facial features may not well capture the discriminative information of the facial images of Multi-PIE.
- The L1-HDA/GM method and the L2-HDA/GM method achieve much competitive recognition rates for most of the linear discriminant analysis methods, where the highest overall recognition accuracy (=81.65%) is achieved by L1-HDA/GM. The better recognition accuracies may attribute to the use the mixture of Gaussians to approximate the multimodal distribution of the feature vectors.
- The SDA method and the MvDN method also achieve

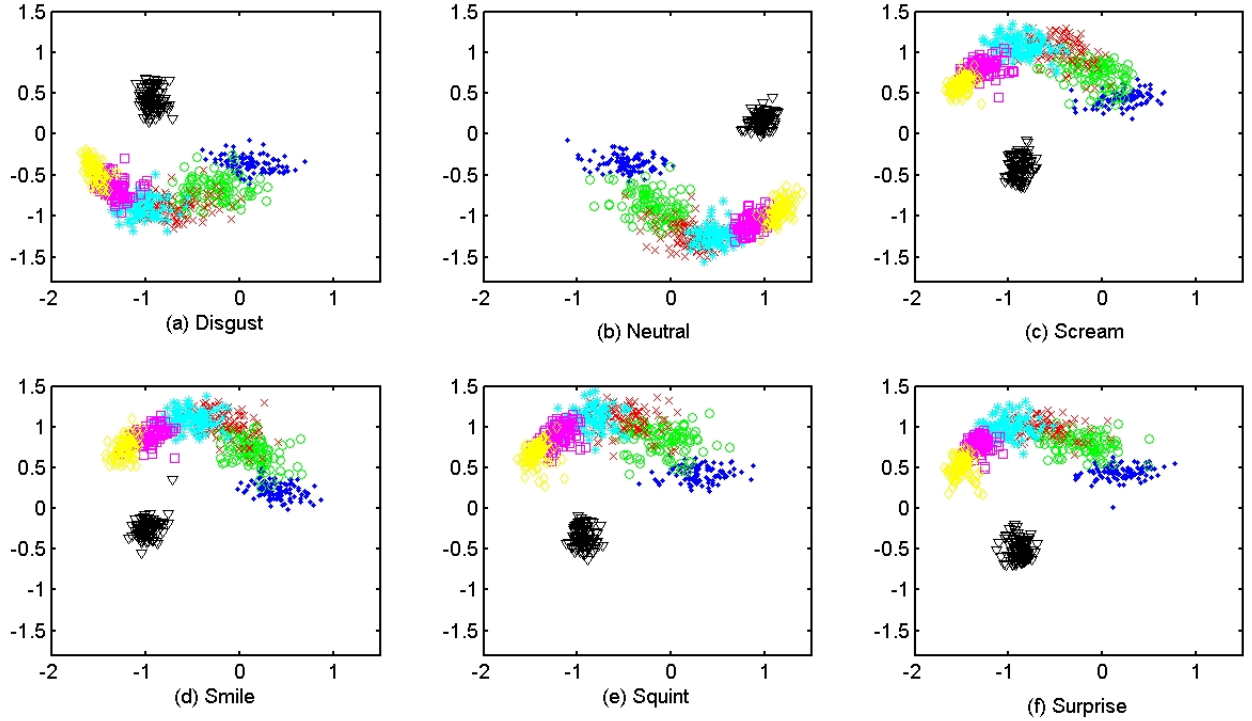


Fig. 5. The distributions of LBP facial features projected by the first two principal components of PCA with respect to all the 6 facial expression classes of Multi-PIE database, where the distribution of the data points of each facial expression demonstrates a multimodal form.

competitive recognition performance compared with the other methods. This is most likely due to the fact that S-DA is actually a special case of our L2-HDA/GM method when the class covariance matrices of data samples are equal. In addition, it is interesting to see from Fig.5 see that the scatter of the data points are similar and hence the corresponding class covariance matrices would be similar. As a result, the difference between our L2-HDA/GM method and SDA in this experiment is trivial and hence they achieve similar better recognition results. As for the MvDN method, we note that it is a nonlinear feature extraction method and hence its better recognition performance may largely attribute to the nonlinear feature learning trick.

Moreover, to evaluate the computational efficiency of the various methods in the training stage, we also compare the CPU running time (second) among the various methods, where the calculation of CPU time is based on the computation platform of Matlab2017B software with CPU i7-7700k and 16GB memory. Table III summarizes the results of the various methods under the two kinds of facial features, i.e. LBP feature versus deep learning (DL) feature, where the parameter scale of the MvDN method is as high as 6×10^5 . From Table III, we can see that the computational complexity of the proposed methods are very competitive to the state-of-the-art HDA methods, such as HDA/HLFE and HDA/Chernoff, which are much less than the MvDN method. Additionally, from Table III, we can also see that the CPU running time of L1-HDA/G is a bit less than L2-HDA/G, which coincides with the computational analysis in Section III-D.

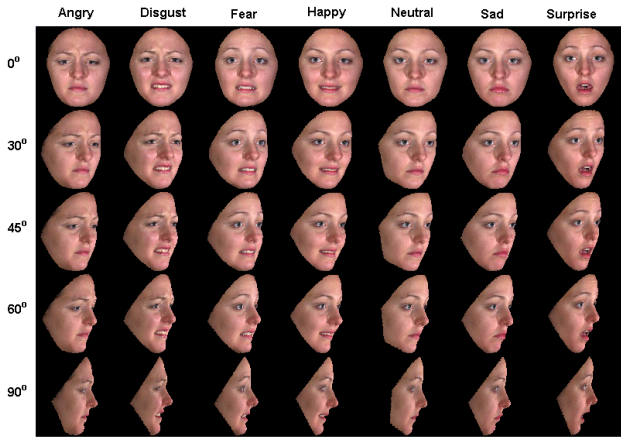
B. Experiment on BU-3DFE Facial Expression Database

In this experiment, we will evaluate the discriminative performance of the proposed methods on the BU-3DFE database, which was developed by Yin et. al. [45] at Binghamton University. The BU-3DFE database contains 2400 3D facial expression models of 100 subjects, which cover 6 basic facial expressions (Anger, Disgust, Fear, Happy, Sad, and Surprise) with 4 levels of intensities. Based on the 3D facial expression model, a set of 12000 multiview 2D facial images are generated [26], which covers 5 yaw facial views (0° , 30° , 45° , 60° , and 90°). Fig. 6 shows the examples of 30 facial images of one subject corresponding to 6 basic facial expressions and 5 facial views in the BU-3DFE database.

In this database, we also explore two kinds of facial feature extraction schemes to evaluate the proposed methods. One is to extract the scale invariant feature transform (SIFT) [46] feature and the other one is to extract the deep learning feature, in which the deep learning feature extraction procedure is similar to the one used in the Multi-PIE database. To extract the SIFT features, we utilized the 83 landmark points obtained by projecting 83 landmark 3D points located on each 3D facial expression model onto 2D image and extract a set of 83 SIFT feature vectors with 128 dimensionality from each facial image. Then, we concatenate all 83 SIFT feature vectors into a 10624-dimensional feature vector to describe the facial image. In addition, to visualize the distribution of the feature points associated with each facial expression, we reduce the dimensionality of the facial feature vectors of the same facial expression from the 10624-dimensional space to 2-dimensional subspace using PCA, and then depict the

TABLE IV
AVERAGE CLASSIFICATION RATES (%) OF VARIOUS METHODS WITH RESPECT TO EACH FACIAL EXPRESSION ON THE BU-3DFE DATABASE.

		Happy	Sad	Angry	Fear	Surprise	Disgust	Overall
# samples		2000	2000	2000	2000	2000	2000	12000
# subjects		100	100	100	100	100	100	100
dimensionality		10624	10624	10624	10624	10624	10624	10624
dimensionality of deep learning feature		4096	4096	4096	4096	4096	4096	4096
LBP Feature	Baseline	88.25	62.28	75.78	30.75	82.53	54.83	65.73
	FLDA	84.55	75.50	75.15	60.60	89.40	73.20	76.40
	AIDA [14]	85.95	75.35	75.63	56.00	88.85	72.40	75.70
	HDA/HLFE [18]	87.75	77.83	77.35	55.15	89.25	71.35	76.45
	HDA/Chernoff [15]	86.30	74.65	73.50	56.23	87.38	66.93	74.16
	SDA [21]	86.70	78.28	77.68	62.28	90.83	74.23	78.33
	MvDN [52]	84.00	73.75	77.00	66.25	88.13	71.88	77.10
	L2-HDA/GM	86.53	78.38	77.63	62.55	90.78	74.53	78.40
	L1-HDA/GM	85.88	78.23	77.88	62.98	90.98	74.25	78.36
Deep Learning Feature	Baseline	65.90	44.38	57.40	26.05	67.70	38.05	49.91
	FLDA	71.25	52.55	58.23	4.88	78.83	56.65	60.73
	AIDA [14]	67.83	38.60	50.43	34.25	68.45	35.50	49.18
	HDA/HLFE [18]	65.33	51.80	50.73	44.68	74.88	54.75	57.03
	HDA/Chernoff [15]	71.83	40.78	53.40	31.30	73.25	43.05	52.27
	SDA [21]	76.95	56.78	61.40	49.93	83.15	63.25	65.24
	MvDN [52]	60.87	67.00	58.50	77.75	55.125	78.00	66.21
	L2-HDA/GM	77.13	56.58	62.53	49.50	82.95	62.50	65.20
	L1-HDA/GM	76.63	56.65	61.88	50.05	82.33	62.48	65.00



Multi-view face images generated from BU-3DFE database

Fig. 6. Examples of the 2D facial images of one subject in the BU-3DFE database with respect to the 6 facial expressions and 5 facial views.

distribution of the data points. Fig.7 shows the distributions of projection points with respect to all the 6 facial expressions. From Fig.7 we observe that the distribution of the data points of each facial expression demonstrates a multimodal form with 5 clusters. Consequently, for both L1-HDA/GM and L2-HDA/GM methods, the mixture of Gaussians with 5 components is used to describe the distribution of the facial feature points for the proposed methods and each Gaussian component also corresponds to the data samples of one facial view.

In the experiments, we use the same cross-validation experimental setting as what we have done in section V-A, i.e., there are 10 trials of experiments are conducted and in each

trial of experiments we partition the whole data set into a training set and a testing set, where the training set contains a total of 9600 facial images of 80 subjects whereas the testing one contains 2400 facial images of 20 subjects. In addition, PCA is also used to reduce the dimensionality of the feature vectors such that the class covariance matrices become non-singular. Finally, the experimental results of all the 10 trials are averaged as the overall recognition rate. Table IV shows the recognition results of the various methods on the BU-3DFE facial expression database.

From Tables IV, we observe the similar experimental results as those obtained on the Multi-PIE database. That is, the average recognition accuracies of using SIFT features are higher than those of using deep learning features and both L1-HDA/GM and L2-HDA/GM achieve higher recognition accuracies than most of the other discriminant analysis methods. Similar to the experiments on Multi-PIE, the major reason of the better recognition results of the proposed methods may attribute to the use the mixture of Gaussians to approximate the multimodal distribution of the feature vectors. Again, we see that the SDA method and the MvDN method also achieve better recognition results compared to the other state-of-the-art linear discriminant analysis methods because of the similar reason of the proposed methods.

C. Experiment on UCI Data Sets

In the above two experiments, we note that the L1-HDA/GM method does not achieve significant improvement in contrast to the L2-HDA/GM method. This is probably due to the fact that the separabilities between the pairwise classes in both BU-3DFE and Multi-PIE databases are similar, such that

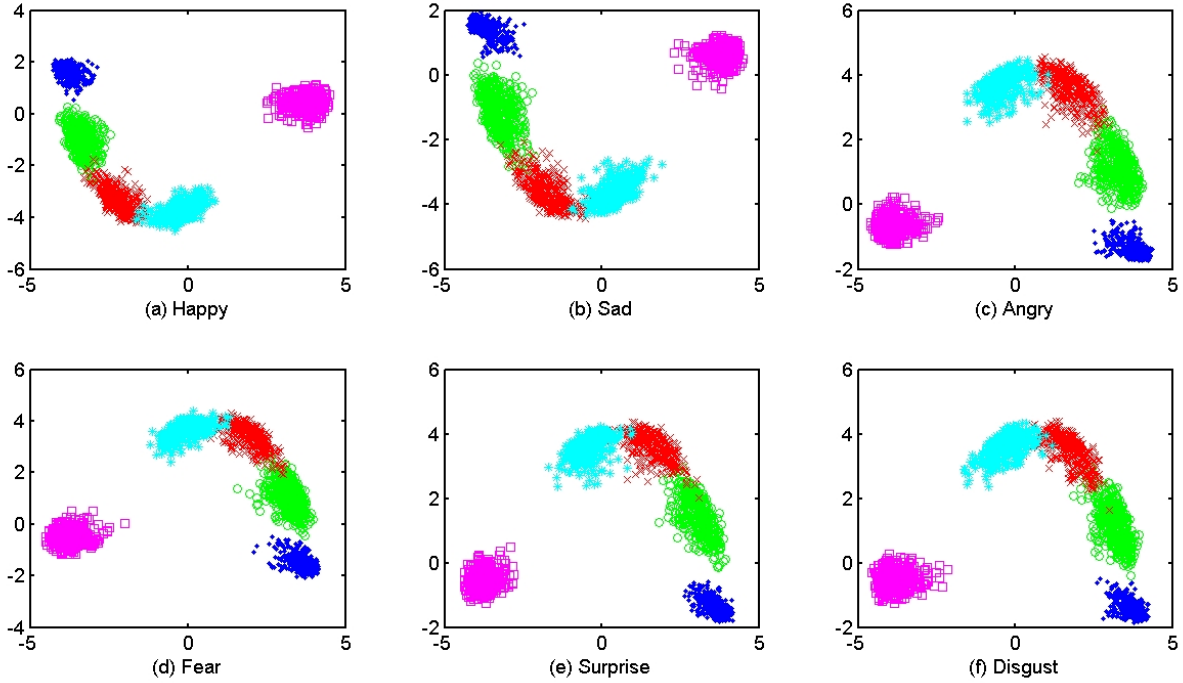


Fig. 7. The distributions of SIFT facial features projected by the first two principal components of PCA with respect to all the 6 facial expression classes of BU-3DFE database.

the advantages of L1-HDA/GM cannot be well reflected. To further compare the recognition performances between L1-HDA/GM and L2-HDA/GM, in this section we will conduct more experiments on more databases, in which the UCI database [56] previously used in [15] is adopted for evaluating the discriminant performance. There are totally 9 data sets are explored in the experiments, which are listed as follows:

- 1) Wisconsin breast cancer (WBC);
- 2) BUPA liver disorder (BUPA);
- 3) Pima indians diabetes (PID);
- 4) Wisconsin diagnostic breast cancer (WDBC);
- 5) Cleveland heart-disease (CHD);
- 6) Thyroid gland (TG);
- 7) Landsat satellite (LS);
- 8) Multifeature digit (Zernike moments) (MD);
- 9) Vowel context (VC);

For each one of the 9 data sets, we randomly divide it into two subsets, with an approximately equal size of samples. Then, we choose one subset for training the algorithms, and use the other one for testing the discriminant performance. We swap the training subset and the testing subset such that each subset is used as the training data set once. For each training data set, we utilize the nearest neighbor clustering to divide the data samples belonging to the same class into K subclasses. In this case, we can use the mixture of Gaussian model with K components to describe the distribution of the data samples of each class.

Similar to [15], before the experiments, a PCA is performed on the training set to reduce the dimensionality of the data samples such that the covariance matrix of each data set is non-singular. Throughout the experiments, we use the quadratic

classifier for classifying the testing data. For each one of the 9 data sets, the average error rate is used to evaluate the various discriminant methods. Table V summarizes the main properties of the 9 UCI data sets and the average error rates of various feature extraction methods, where “#PC” in the fourth row shows how many principal components we use after the PCA processing. From Table V, we can observe that L1-HDA/GM outperforms L2-HDA/GM in all these 9 data sets. Another observation from Table V is that, for both L1-HDA/GM and L2-HDA/GM, using the mixture of Gaussian model to describe the data samples of each class could achieve better than using single Gaussian model.

D. Experiment on EEG Data Sets

In this experiment, we will evaluate the effectiveness of the proposed heteroscedastic discriminant analysis in dealing with the feature extraction problem as for the case when the class means are the same. To this end, we focus on an EEG classification problem whose target is to recognize the motor imagery tasks based on the EEG signal, i.e., recognizing what kind of motor imagery task that the EEG signal corresponds to.

The data set used in this experiment is from the “BCI competition 2005” - data set IIIa [47]. This data set consists of recordings from three subjects (k3b, k6b, and l1b), which performed four different motor imagery tasks (left/right hand, one foot, or tongue) according to a cue. During the experiments, the EEG signal is recorded in 60 channels, using the left mastoid as reference and the right mastoid as ground. The EEG was sampled at 250 Hz and was filtered between 1 and 50 Hz with the notch filter on. Each trial lasted for 7 seconds,

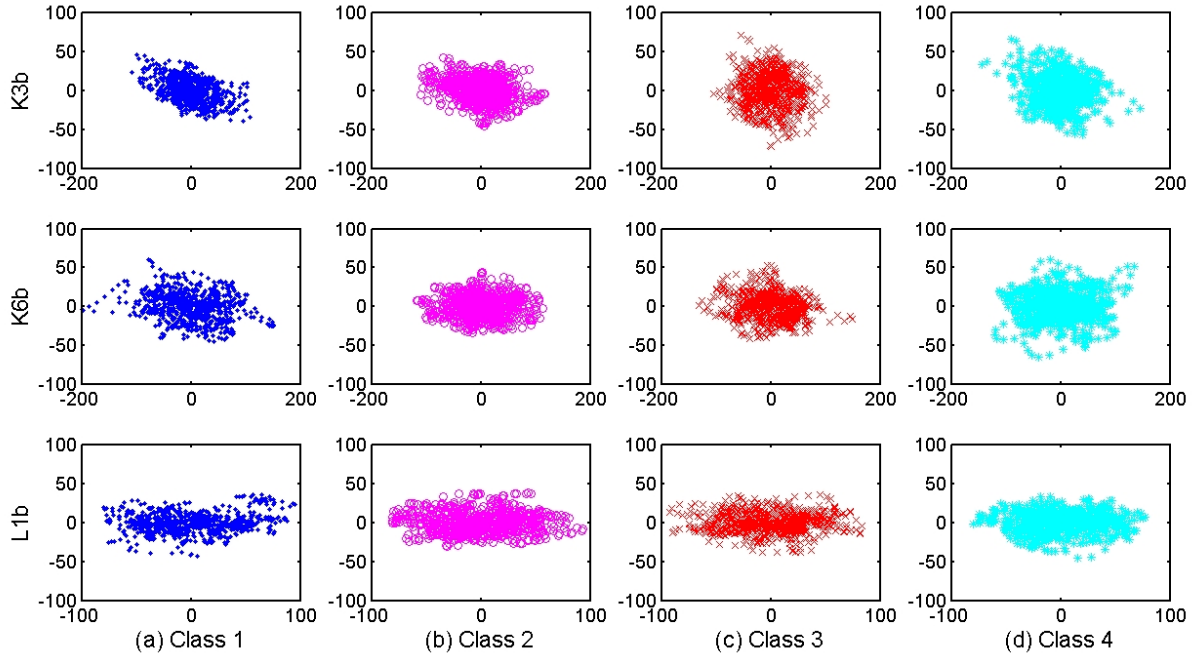


Fig. 8. Examples of the EEG data samples of one trial with respect to four EEG classes and three subjects. The four figures of each row show the data point distributions corresponding to the four EEG classes of one subject, whereas the three figures of each column show the data point distribution corresponding to the same class of the three subjects (K3b, K6b, and L1b).

with the motor imagery performed during the last 4 seconds of each trial. For subjects k6b and l1b, a total of 60 trials per condition were recorded. For subject k3b, a total of 90 trials per condition were recorded. In addition, similar to the method in [10], we discard the four trials of subject k6b with missing data. The EEG data samples associated with the same class are preprocessed such that their class means equal to zero vector [27]. Fig.8 shows examples of the preprocessed EEG data samples of one trial with respect to different EEG classes and subjects, in which the four figures of each row show the data point distributions corresponding to four EEG classes of one subject while the three figures of each column show the data point distribution corresponding to the same class of three subjects, respectively. From Fig.8 we can see that the sample mean in each class is a zero point and hence only the class covariance differences can be utilized to extract the discriminative features.

Similar to [27], for each trial of the EEG data, we only use parts of the sample points, i.e., from No.1001 to No.1750, as the experiment data. Consequently, each trial contains 750 data points. In the experiment, we adopt two-fold cross validation strategy to perform the experiment, i.e., we divide all the EEG trials into two groups and select one as training data set and the other one as testing data set, and then we swap the training and testing data set to repeat the experiment. Since the class means of the EEG data samples equal to the zero vector, the traditional Fisher's criterion based approach, such as FLDA and SDA, cannot be applied in this experiment. Consequently, in this experiment we only evaluate the discriminative feature extraction performance of the following four heteroscedastic discriminant analysis methods, i.e., the AIDA method, the

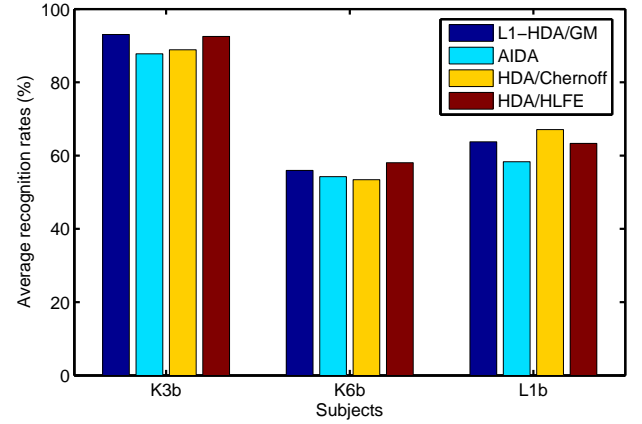


Fig. 9. The average recognition rates of various methods on the EEG data sets of three subjects.

HDA/Chernoff method, the HDA/HLFE method, and the proposed L1-HDA/GM method.

As for the EEG classification, we extract the same log-transformation variance used in [11] to represent the final EEG feature of each trial. Then, we use linear classifier to perform the EEG classification. Fig.9 shows the average recognition rates of the four methods on the three subjects, from which we can clearly see that the proposed L1-HDA/GM method achieves much competitive experimental results compared to the best experimental results obtained by the other state-of-the-art methods.

TABLE V
UCI BENCHMARK DATASETS AND THE AVERAGE ERROR RATES (%) OF SEVERAL METHODS

Data set		WBC	BUPA	PID	WDBC	CHD	TG	LS	MD	VC
# samples		682	345	768	569	297	215	6435	2000	990
# subjects		2	2	2	2	2	3	6	10	11
dimensionality		9	6	8	30	13	5	36	47	10
# PC		9	6	8	7	13	5	36	33	10
baseline	/	4.48	40.12	31.69	9.15	41.62	5.57	12.15	23.98	9.58
FLDA	/	4.56	43.14	29.35	6.14	24.22	4.54	11.71	19.90	1.16
AIDA	/	4.4	36.28	32.71	6.74	24.89	6.18	18.04	23.86	1.33
HDA/HLFE	/	4.56	43.14	29.35	6.14	22.00	4.54	13.71	19.90	1.16
HDA/Chernoff	/	4.42	39.28	31.13	5.88	23.05	3.63	13.22	23.03	1.38
MvDN	/	3.89	35.00	29.01	6.01	21.13	4.15	8.99	16.24	1.44
SDA [21]	$K = 1$	3.77	39.71	28.31	7.19	24.33	3.63	13.51	20.30	1.52
	$K = 2$	4.20	36.57	29.61	6.49	24.67	4.55	13.12	19.85	1.52
	$K = 3$	4.49	43.14	31.82	6.84	24.67	4.55	13.18	19.50	1.52
	$K = 4$	4.78	37.43	33.38	6.84	26.00	4.55	13.51	18.70	1.52
	$K = 5$	4.78	37.43	32.08	6.84	27.00	4.55	14.24	18.70	1.52
L2-HDA/GM	$K = 1$	5.21	36.57	41.10	6.32	27.00	6.59	13.65	19.40	1.11
	$K = 2$	4.35	36.57	34.68	6.32	23.00	5.45	11.74	18.35	1.11
	$K = 3$	4.35	36.57	31.75	6.32	23.00	5.45	11.43	17.90	1.11
	$K = 4$	4.35	36.57	31.43	6.32	22.17	5.45	11.13	17.90	1.11
	$K = 5$	4.35	36.57	31.17	6.32	22.17	5.45	11.13	17.90	1.11
L1-HDA/GM	$K = 1$	4.86	43.71	30.13	6.93	22.67	4.09	14.40	19.75	1.06
	$K = 2$	3.77	36.14	30.13	6.93	20.67	3.64	11.98	18.98	1.06
	$K = 3$	3.77	34.14	29.74	5.79	20.67	3.64	10.83	17.80	1.06
	$K = 4$	3.77	34.14	29.22	5.79	20.67	3.64	10.43	17.55	1.06
	$K = 5$	3.77	34.14	28.44	5.79	20.67	3.64	10.43	17.55	1.06

VI. CONCLUSIONS AND DISCUSSIONS

In this paper, we have proposed a novel L2-MAGRQ criterion based on the Bayes error upper bound estimation. This criterion can be seen as a generalization of the traditional Fisher's criterion aiming to overcome the limitations of Fisher's criterion in the case of heteroscedastic distribution of the data samples in each class. The L2-MAGRQ criterion is further modified by replacing the ℓ_2 norm operation in the between-class scatter distance with ℓ_1 norm, resulting in the L1-MAGRQ criterion. Two kinds of heteroscedastic discriminant analysis methods, L1-HDA/G (L2-HDA/G) and L1-HDA/GM (L2-HDA/GM), based on the L1-MAGRQ (L2-MAGRQ) criterion are respectively proposed for discriminative feature extraction, in which L1-HDA/G (L2-HDA/G) corresponds to the case where the distribution of each class data set is Gaussian whereas L1-HDA/GM (L2-HDA/GM) corresponds to the mixture of Gaussian distributions (it is notable that L1-HDA/G (L2-HDA/G) is actually a special case of L1-HDA/GM (L2-HDA/GM) when the mixture of Gaussian distributions reduces to the Gaussian distribution). Moreover, we also propose an efficient algorithm to compute the optimal discriminant vectors of L1-HDA/GM (L2-HDA/GM) by solving a series of principal eigenvectors, which can both methods presented in this paper for solving the optimal discriminant vectors of L1-HDA/GM (L2-HDA/GM) is a greedy algorithm, it is easier to develop a non-greedy algorithm for L1-HDA/GM (L2-HDA/GM) by referring the method proposed in [48]

and [49]. The experiments on four real databases had been conducted to evaluate the discriminative performance of the proposed methods. The experimental results demonstrate that the proposed L1-HDA/GM method achieves better recognition performance than most of the state-of-the-art methods, which may attribute to the use of mixture of Gaussian distributions and the Bayes error estimation. In addition, in the multi-view facial expression recognition experiments, we see that the SDA method also achieves the similar better experimental results as ours. This is most likely due to the fact that SDA is a special case of L2-HDA/GM as for the case of similar class covariance matrices.

Additionally, the proposed L1-HDA/GM (L2-HDA/GM) methods can also be used to improve the current graph-based subspace learning performance. For example, in [57], Peng et al. constructed an ℓ_2 norm based sparse similarity graph for robust subspace learning under the sparse representation framework, in which the sparse relationships among the data points are preserved in the low-dimensional subspace. It is notable that the feature extraction part of the method proposed by Peng et al. [57] is actually an unsupervised subspace learning approach, which could not fully utilize the class label information of the data points to improve the discriminative ability of the extracted features. By adopting the proposed HDA methods, however, we may be able to learn a more discriminative subspace for the feature extraction purpose.

In addition, in our experiments, we can see that the MvDN method proposed in [52] achieved very competitive results with our L1-HDA/GM (L2-HDA/GM) methods. This is very

likely due to the fact that MvDN is actually a nonlinear feature extraction method implemented by using deep neural networks approach (hence the computational complexity of MvDN would be larger than the other methods, see Table III for more details). In contrast to MvDN, both proposed L1-HDA/GM and L2-HDA/GM are linear feature extraction methods. Nevertheless, it is notable that we are able to adopt the similar non-linear learning trick of using deep neural network to realize the proposed L1-HDA/GM (L2-HDA/GM) algorithm to further improve discriminative feature extraction performance, which would be our future work.

APPENDIX

APPENDIX A: PROOF OF THEOREM 1.

Proof: Suppose that $\alpha^{(2)}$ is the principal eigenvector of $\mathbf{T}(\mathbf{U}_2, \mathbf{V}_2)$, i.e.,

$$\alpha^{(2)} = \arg \max_{\|\alpha\|=1} \alpha^T \mathbf{T}(\mathbf{U}_2, \mathbf{V}_2) \alpha. \quad (43)$$

Then by the definition of $\alpha^{(2)}$, we have

$$\alpha^{(2)T} \mathbf{T}(\mathbf{U}_2, \mathbf{V}_2) \alpha^{(2)} \geq \alpha^{(1)T} \mathbf{T}(\mathbf{U}_2, \mathbf{V}_2) \alpha^{(1)}. \quad (44)$$

On the other hand, we have

$$\begin{aligned} \alpha^{(1)T} \mathbf{T}(\mathbf{U}_2, \mathbf{V}_2) \alpha^{(1)} &= \left(\sum_{i < j} (\mathbf{U}_2)_{ij} \alpha^{(1)T} \Delta \hat{\mathbf{m}}_{ij} \right)^2 \\ &\quad + \sum_{i < j} (\mathbf{V}_2)_{ij} \alpha^{(1)T} \Delta \hat{\mathbf{S}}_{ij} \alpha^{(1)}. \end{aligned} \quad (45)$$

From (25) and (45), we have

$$\begin{aligned} \alpha^{(1)T} \mathbf{T}(\mathbf{U}_2, \mathbf{V}_2) \alpha^{(1)} &= \left(\sum_{i < j} |\alpha^{(1)T} \Delta \hat{\mathbf{m}}_{ij}| \right)^2 \\ &\quad + \sum_{i < j} |\alpha^{(1)T} \Delta \hat{\mathbf{S}}_{ij} \alpha^{(1)}| \\ &\geq \left(\sum_{i < j} (\mathbf{U}_1)_{ij} \alpha^{(1)T} \Delta \hat{\mathbf{m}}_{ij} \right)^2 + \sum_{i < j} (\mathbf{V}_1)_{ij} \alpha^{(1)T} \Delta \hat{\mathbf{S}}_{ij} \alpha^{(1)} \\ &= \alpha^{(1)T} \mathbf{T}(\mathbf{U}_1, \mathbf{V}_1) \alpha^{(1)}. \end{aligned} \quad (46)$$

Combine (44) and (46), we have

$$\alpha^{(2)T} \mathbf{T}(\mathbf{U}_2, \mathbf{V}_2) \alpha^{(2)} \geq \alpha^{(1)T} \mathbf{T}(\mathbf{U}_1, \mathbf{V}_1) \alpha^{(1)}. \quad (47)$$

APPENDIX B: PROOF OF THEOREM 2.

For simplicity of the derivation, we denote the class distribution function of $\mathbf{x} \in \mathbf{X}_i$ by $p_x(\mathbf{x})$ and the distribution function of $y = \omega^T \mathbf{x}$ by $p_y(y)$, i.e., $p_x(\mathbf{x}) = p_i(\mathbf{x} | \mathbf{x} \in \mathbf{X}_i)$ and $p_y(y = \omega^T \mathbf{x}) = \tilde{p}_i(\omega^T \mathbf{x} | \mathbf{x} \in \mathbf{X}_i)$. Then the characteristic function of \mathbf{x} is

$$\begin{aligned} \phi_x(\mathbf{t}) &= E(e^{j\mathbf{t}^T \mathbf{x}}) = \sum_{r=1}^{K_i} \pi_{ir} \int_{\mathbf{x}} e^{j\mathbf{t}^T \mathbf{x}} \mathcal{N}(\mathbf{m}_{ir}, \Sigma_{ir}) d\mathbf{x} \\ &= \sum_{r=1}^{K_i} \pi_{ir} \exp \left(j\mathbf{t}^T \mathbf{m}_{ir} - \frac{1}{2} \mathbf{t}^T \Sigma_{ir} \mathbf{t} \right). \end{aligned} \quad (48)$$

where $j^2 = -1$.

Then the characteristic function of y is:

$$\begin{aligned} \phi_y(\xi) &= E\{e^{j\xi y}\} = E\{e^{j\xi \omega^T \mathbf{x}}\} = \phi_x(\xi \omega) \\ &= \sum_{r=1}^{K_i} \pi_{ir} \exp \left(j\xi \omega^T \mathbf{m}_{ir} - \frac{1}{2} \xi^2 \omega^T \Sigma_{ir} \omega \right). \end{aligned} \quad (49)$$

So the density distribution function of y is

$$\begin{aligned} p_y(\eta) &= \frac{1}{2\pi} \int_{\xi} \phi_y(\xi) e^{-j\xi \eta} d\xi \\ &= \sum_{r=1}^{K_i} \pi_{ir} \frac{1}{2\pi} \int_{\xi} e^{-j\xi \eta} \exp \left(j\xi \omega^T \mathbf{m}_{ir} - \frac{1}{2} \xi^2 \omega^T \Sigma_{ir} \omega \right) d\xi \\ &= \sum_{r=1}^{K_i} \pi_{ir} \frac{1}{2\pi} \int_{\xi} \exp \left\{ -\frac{1}{2} [\xi^2 \omega^T \Sigma_{ir} \omega + 2j\xi(\eta - \omega^T \mathbf{m}_{ir})] \right\} d\xi \\ &= \sum_{r=1}^{K_i} \pi_{ir} \frac{1}{2\pi} \int_{\xi} \exp \left\{ -\frac{1}{2} [(\omega^T \Sigma_{ir} \omega)(\xi + \frac{j(\eta - \omega^T \mathbf{m}_{ir})}{\omega^T \Sigma_{ir} \omega})^2 + \frac{(\eta - \omega^T \mathbf{m}_{ir})^2}{\omega^T \Sigma_{ir} \omega}] \right\} d\xi \\ &= \sum_{r=1}^{K_i} \pi_{ir} \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} \frac{(\eta - \omega^T \mathbf{m}_{ir})^2}{\omega^T \Sigma_{ir} \omega} \right\} \\ &\quad \times \frac{1}{\sqrt{2\pi}} \int_{\xi} \exp \left\{ -\frac{1}{2} [(\omega^T \Sigma_{ir} \omega)(\xi + \frac{j(\eta - \omega^T \mathbf{m}_{ir})}{\omega^T \Sigma_{ir} \omega})^2] \right\} d\xi \\ &= \sum_{r=1}^{K_i} \pi_{ir} \frac{1}{\sqrt{2\pi(\omega^T \Sigma_{ir} \omega)}} \exp \left\{ -\frac{1}{2} \frac{(\eta - \omega^T \mathbf{m}_{ir})^2}{\omega^T \Sigma_{ir} \omega} \right\} \\ &= \sum_{r=1}^{K_i} \pi_{ir} \mathcal{N}(\eta | \omega^T \mathbf{m}_{ir}, \omega^T \Sigma_{ir} \omega). \end{aligned} \quad (50)$$

This completes the proof of Theorem 2.

REFERENCES

- [1] K. Fukunaga, "Introduction to statistical pattern recognition (second edition)", Academic Press, New York, 1990.
- [2] H. Zhu, F. Meng, J. Cai, S. Lu, "Beyond pixels: A comprehensive survey from bottom-up to semantic image segmentation and cosegmentation," *Journal of Visual Communication & Image Representation*, Vol.34, pp.12-27, 2016.
- [3] H. Zhu, R. Vial, S. Lu, X. Peng, H. Fu, Y. Tian, X. Cao, "YoTube: searching action proposal via recurrent and static regression networks," *IEEE Transactions on Image Processing*, Vol.27, No.6, pp.2609-2622, 2018.
- [4] R. O. Duda and P. E. Hart, "Pattern classification and scene analysis," John Wiley & Sons, Inc., New York, 1973.
- [5] 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.
- [6] D. Swets, J. Weng, "Using discriminant eigenfeatures for image retrieval," *IEEE Transactions on pattern analysis and machine intelligence*, Vol.18, No.8, pp.831-836, 1996.
- [7] R. Haeb-Umbach, H. Ney, "Linear discriminant analysis for improved large vocabulary continuous speech recognition," *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, pp.13-16, 1992.

- [8] M. Yang, L. Zhang, X. Feng, D. Zhang, "Sparse representation based Fisher discrimination dictionary learning for image classification," *International Journal of Computer Vision*, Vol.109, pp.209-232, 2014.
- [9] X. Peng, H. Tang, L. Zhang, Z. Yi, S. Xiao, "A unified framework for representation-based subspace clustering of out-of-sample and large-scale data," *IEEE Transactions on Neural Networks and Learning Systems*, Vol.27, No.12, pp.2499-2512, 2016.
- [10] M. Grosse-Wentrup and M. Buss, "Multiclass common spatial patterns and information theoretic feature extraction," *IEEE Transactions on Biomedical Engineering*, Vol.55, No.8, pp.1991-2000, 2008.
- [11] H. Ramoser, J. Mueller-Gerking, & G. Pfurtscheller, "Optimal spatial filtering of single trial EEG during imaged hand movement," *IEEE Transactions on Rehabilitation Engineering*, Vol.8, No.4, pp.441-446, 2000.
- [12] N. Kumar and A. G. Andreou, "Heteroscedastic discriminant analysis and reduced rank hmms for improved speech recognition," *Speech Communication*, Vol.26, pp.462-467, 1998.
- [13] G. Saon, M. Padmanabhan, R. Gopinath, and S. Chen, "Maximum likelihood discriminant feature spaces," *Proceedings of International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, pp.129-132, 2000.
- [14] K. Das and Z. Nenadic, "Approximate information discriminant analysis: A computationally simple heteroscedastic feature extraction technique," *Pattern Recognition*, Vol.41, pp.1548-1557, 2008.
- [15] M. Loog and R. P. W. Duin, "Linear dimensionality reduction via a heteroscedastic extension of LDA: The Chernoff criterion," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Vol.26, No.6, pp.732-739, 2004.
- [16] Y. K. Noh, J. Hamm, F. Park, et al., "Fluid Dynamic Models for Bhattacharyya-based Discriminant Analysis," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017.
- [17] Z. Nenadic, "Information discriminant analysis: feature extraction with an information-theoretic objective," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Vol.29, No.8, pp.1394-1407, 2007.
- [18] P. F. Hsieh, D. S. Wang, and C. W. Hsu, "A linear feature extraction for multiclass classification problems based on class mean and covariance discriminant information," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Vol.28, No.2, pp.223-235, 2006.
- [19] P.F. Hsieh, and D.A. Landgrebe, "Linear feature extraction for multiclass problems", *Proc. IEEE Int'l Geoscience and Remote Sensing Symp.*, Vol.4, pp.2050-2052, 1998.
- [20] W. Zheng, L. Zhao, and C. Zou, "An efficient algorithm to solve the small sample size problem for LDA," *Pattern Recognition*, Vol.37, No.5, pp.1077-1079, 2004.
- [21] M. Zhu, and A. M. Martinez, "Subclass discriminant analysis," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Vol.28, No.8, pp.1274-1286, 2006.
- [22] H. Wan, H. Wang, G. Guo, X. Wei, "Separability-oriented subclass discriminant analysis," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017.
- [23] S. Mika, G. Rätsch, J. Weston, B. Schölkopf, K.-R. & Müller, "Fisher discriminant analysis with kernels," in *Proceedings of IEEE int'l Workshop neural networks for Signal Processing IX*, pp.41-48, 1999.
- [24] M. H. Yang, "Kernel eigenfaces vs. kernel fisherfaces: face recognition using kernel methods," In *Proceedings of the Fifth IEEE International Conference on Automatic Face and Gesture Recognition*, 2002.
- [25] J. Shawe-Taylor and N. Cristianini, "Kernel methods for pattern analysis," Cambridge University Press, 2004.
- [26] W. Zheng, H. Tang, Z. Lin, & T.S. Huang, "A novel approach to expression recognition from non-frontal face images," *Proceedings of IEEE International Conference on Computer Vision (ICCV2009)*, pp.1901-1908, 2009.
- [27] W. Zheng and Z. Lin, "Optimizing multi-class spatio-spectral filters via Bayes error estimation for EEG classification," *Neural Information Processing Systems (NIPS)*, 2009.
- [28] W. Zheng, H. Tang, Z. Lin, & T.S. Huang, "Emotion recognition from arbitrary view facial images," *Proceedings of European Conference on Computer Vision (ECCV2010)*, pp.490-503, 2010.
- [29] W. Malina, "On an extended fisher criterion for feature selection", *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Vol.3, No.5, pp.611-614, 1981.
- [30] H. Wang, Q. Tang, W. Zheng, "L1-norm-based common spatial patterns," *IEEE Transactions on Biomedical Engineering*, Vol.59, No.3, pp.653-662, 2012.
- [31] F. Zhong, J. Zhang, "Linear discriminant analysis based on L1-norm maximization," *IEEE Transactions on Image Processing*, Vol.22, No.8, pp.3018-3027, 2013.
- [32] W. Zheng, Z. Lin, H. Wang, "L1-norm kernel discriminant analysis via Bayes error bound optimization for robust feature extraction", *IEEE transactions on neural networks and learning systems*, Vol.25, No.4, pp.793-805, 2014.
- [33] H. Wang, X. Lu, Z. Hu, W. Zheng, "Fisher discriminant analysis with L1-norm," *IEEE transactions on cybernetics*, Vol.44, No.6, pp.828-842, 2014.
- [34] Q. Ye, J. Yang, F. Liu, "L1-norm Distance Linear Discriminant Analysis Based on An Effective Iterative Algorithm," *IEEE Transactions on Circuits and Systems for Video Technology*, 2016.
- [35] Q. Ye, J. Yang, F. Liu, C. Zhao, N. Ye, and T. Yin, "L1-norm Distance Linear Discriminant Analysis Based on An Effective Iterative Algorithm," *IEEE Transactions on Circuits and Systems for Video Technology*, 2017. DOI 10.1109/TCSVT.2016.2596158
- [36] W. Zheng, Z. Lin, X. Tang, "A rank-one update algorithm for fast solving kernel FoleyCSammon optimal discriminant vectors," *IEEE Transactions on Neural Networks*, Vol.21, No.3, pp.393-403, 2010.
- [37] K. Fukunaga, W. Koontz, "Application of the Karhunen-Loeve expansion to feature selection and ordering," *IEEE Transactions on Computer*, Vol.C-19, No.4, pp.311-318, 1970.
- [38] Z. Jin, J. Y. Yang, Z. S. Hu, and Z. Lou, "Face recognition based on the uncorrelated discriminant transformation," *Pattern Recognition*, Vol.34, pp.1405-1416, 2001.
- [39] W. Zheng, "Heteroscedastic feature extraction for texture classification," *IEEE Signal Processing Letters*, Vol.16, No.9, pp.766-769, 2009.
- [40] C. Cortes, V. Vapnik, "Support-vector network," *Machine Learning*, Vol.20, pp.273-297, 1995.
- [41] P. Viola, M. Jones, "Rapid object detection using a boosted cascade of simple features," *Proceedings of CVPR*, pp.511-518, 2001.
- [42] R. Gross, I. Matthews, J. Cohn, T. Kanade, S. Baker, "Multi-PIE," *Image and Vision Computing*, Vol.28, pp.807-813, 2010.
- [43] W. Zheng, "Multi-view facial expression recognition based on group sparse reduced-rank regression," *IEEE Transactions on Affective Computing*, Vol.5, No.1, pp.71-85, 2014.
- [44] T. Ojala, M. Pietikainen, and I. Maenpaa, "Multi-resolution gray-scale and rotation invariant texture classification with local binary patterns," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, Vol.24, pp.971-987, 2002.
- [45] L. Yin, X. Wei, Y. Sun, J. Wang, M.J. Rosato, "A 3D facial expression database for facial behavior research," *Proceedings of 7th Int. Conf. on Automatic Face and Gesture Recognition*, pp.211-216, 2006.
- [46] D. Lowe, "Distinctive image features from scale-invariant keypoints," *International Journal of Computer Vision*, Vol.60, No.2, pp.91-110, 2004.
- [47] G. Blankertz, K. R. Müller, D. Krusienski, G. Schalk, J. R. Wolpaw, A. Schloegl, G. Pfurtscheller, J. R. Millan, M. Schroeder, and N. Birbaumer, "The BCI competition III: Validating alternative approaches to actual BCI problems", *IEEE Transactions on Rehabilitation Engineering*, Vol.14, No.2, pp.153-159, 2006.
- [48] F. Nie, H. Huang, C. Ding, D. Luo, H. Wang, "Robust principal component analysis with non-greedy L1-norm maximization," *Proceedings of the Twenty-second International Joint Conference on Artificial Intelligence*, pp.1433-1438, 2011.
- [49] Y. Liu, Q. Gao, S. Miao, X. Gao, F. Nie, and Y. Li, "A non-greedy algorithm for L1-norm LDA," *IEEE Transactions on Image Processing*, Vol.26, No.2, pp.684-695, 2017.
- [50] T. Diethe, D. R. Hardoon, J. Shawe-Taylor, "Constructing nonlinear discriminants from multiple data views," *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pp.328-343, 2010.
- [51] H. Wang, C. Ding, H. Huang, "Multi-label linear discriminant analysis," *European Conference on Computer Vision (ECCV2010)*, pp.126-139, 2010.
- [52] M. Kan, S. Shan, & X. Chen, "Multi-view deep network for cross-view classification," In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp.4847-4855, 2016.
- [53] S. Sun, "A survey of multi-view machine learning," *Neural Computing and Applications*, Vol.23, No.7-8, pp.2031-2038, 2013.
- [54] Shan Li, Weihong Deng, and JunPing Du, "Reliable crowdsourcing and deep locality-preserving learning for expression recognition in the wild," *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp.2584-2593, 2017.
- [55] O. Parkhi, A. Vedaldi, A. Zisserman, "Deep face recognition," in *Proceedings of British Machine Vision*, pp.1-12, 2015.
- [56] "UCI repository of machine learning database," <http://www.ics.uci.edu/~mllearn/MLRepository>, 2004.

- [57] X. Peng, Z. Yu, H. Tang, "Constructing the L2-graph for robust subspace learning and subspace clustering," IEEE Transactions on Cybernetics, Vol.47, No.4, pp.1053-1066, 2017.



Wenming Zheng (M'08) received the B.S. degree in computer science from Fuzhou University, Fuzhou, China, in 1997, the M.S. degree in computer science from Huaqiao University, Quanzhou, China, in 2001, and the Ph.D. degree in signal processing from Southeast University, Nanjing, China, in 2004. Since 2004, he has been with the Research Center for Learning Science, Southeast University. He is currently a Professor with the Key Laboratory of Child Development and Learning Science, Ministry of Education, Southeast University. His research

interests include affective computing, pattern recognition, machine learning, and computer vision. He is an associated editor of IEEE Transactions on Affective Computing, an associated editor of Neurocomputing and also an associate editors-in-chief of The Visual Computer.



Cheng Lu received the B.S. and M.S. degree from the School of Computer Science and Technology, Anhui University, China, in 2013 and 2017, respectively. Currently, he is a Ph.D. candidate in the School of Information Science and Engineering, Southeast University, under the supervision of professor Wenming Zheng. His research interests include speech emotion recognition, machine learning and pattern recognition.



Zhouchen Lin (M00-SM08-F18) is currently a professor with the Key Laboratory of Machine Perception, School of Electronics Engineering and Computer Science, Peking University. His research interests include computer vision, image processing, machine learning, pattern recognition, and numerical optimization. He is an area chair of CVPR 2014/2016/2019, ICCV 2015, NIPS 2015/2018 and AAAI 2019, and a senior program committee member of AAAI 2016/2017/2018 and IJCAI 2016/2018. He is an associate editor of the IEEE Transactions

on Pattern Analysis and Machine Intelligence and the International Journal of Computer Vision. He is an IAPR Fellow and IEEE Fellow.



Tong Zhang received the B.S. degree in Department of Information Science and Technology, Southeast University, China, in 2011, the M.S. degree in Research Center for Learning Science, Southeast University, China, in 2014. Currently, he is pursuing the Ph.D. degree in information and communication engineering in Southeast University, China. His interests include pattern recognition, machine learning and computer vision.



Zhen Cui Received the Ph.D. degree in computer science from Institute of Computing Technology (ICT), Chinese Academy of Science (CAS), Beijing, in Jun. 2014. He was a Research Fellow in the Department of Electrical and Computer Engineering at National University of Singapore (NUS) from Sep. 2014 to Nov. 2015. He also spent half a year as a Research Assistant on Nanyang Technological University (NTU) from Jun. 2012 to Dec. 2012. Now he is a professor of Nanjing University of Science and Technology, China. His research interests

cover computer vision, pattern recognition and machine learning, especially focusing on deep learning, manifold learning, sparse coding, face detection/alignment/recognition, object tracking, image super resolution, emotion analysis, etc.



Wankou Yang received the B.S., M.S., and Ph.D. degrees from the School of Computer Science and Technology, Nanjing University of Science and Technology, China, in 2002, 2004, and 2009, respectively. From 2009 to 2011, he was a Post-Doctoral Fellow with the School of Automation, Southeast University, China. From 2011 to 2016, he was an Assistant Professor with the School of Automation, Southeast University, where he is currently an Associate Professor. His research interests include pattern recognition and computer vision.