

# Mixture subclass discriminant analysis link to restricted Gaussian model and other generalizations

Nikolaos Gkalelis, Vasileios Mezaris, *Member, IEEE*, Ioannis Kompatsiaris, *Senior Member, IEEE*, and Tania Stathaki

**Abstract**—In this paper, a theoretical link between mixture subclass discriminant analysis (MSDA) and restricted Gaussian model is first presented, and then two further discriminant analysis (DA) methods, fractional step MSDA (FSMSDA) and kernel MSDA (KMSDA) are proposed. Linking MSDA to an appropriate Gaussian model allows the derivation of a new DA method under the Expectation Maximization (EM) framework (EM-MSDA), that derives simultaneously the discriminant subspace as well as the maximum likelihood estimates. The two other proposed methods generalize MSDA in order to solve problems inherited from conventional discriminant analysis. FSMSDA solves the subclass separation problem, that is, the situation when the dimensionality of the discriminant subspace is strictly smaller than the rank of the inter-between-subclass scatter matrix. This is done by an appropriate weighting scheme and the utilization of an iterative algorithm for preserving useful discriminant directions. On the other hand, KMSDA uses the kernel trick to separate data with nonlinearly separable subclass structure. Extensive experimentation shows that the proposed methods outperform conventional MSDA and other LDA variants.

**Index Terms**—Feature extraction, discriminant analysis, mixture of Gaussians, probabilistic algorithms, clustering, pattern recognition, classification, machine learning.

## I. INTRODUCTION

In a natural environment, the high dimensional measurement signals, lying in the  $F$ -dimensional measurement space, usually represent patterns residing in a much lower,  $D$ -dimensional subspace embedded in the ambient measurement space [1]. Dimensionality reduction (DR) is an important component of statistical pattern classifiers that helps to overcome estimation problems in noisy high-dimensional environments, and thus, often results in improved classifier accuracy as well as lower storage and processing time requirements. A fundamental DR technique is linear discriminant analysis (LDA) [2]–[4]. Given a training set of  $C$  classes and  $N$  training observations represented with the block matrix  $\mathbf{X} = [\mathbf{X}_1, \dots, \mathbf{X}_C]$ , whose  $i$ -th block,  $\mathbf{X}_i = [\mathbf{x}_i^1, \dots, \mathbf{x}_i^{N_i}]$ , consists of the  $N_i$  observations  $\mathbf{x}_i^n \in \mathbb{R}^F$  of the  $i$ -th class, this method derives a discriminant subspace spanned by the column vectors of the transformation matrix  $\Psi \in \mathbb{R}^{F \times D}$  that maximizes the ratio

$$J_{LDA}(\Psi) = \frac{\text{tr}(\Psi^T \mathbf{S}_b \Psi)}{\text{tr}(\Psi^T \mathbf{S}_w \Psi)} \quad (1)$$

N. Gkalelis is with the Information Technologies Institute/Centre for Research and Technology Hellas (CERTH), Thessaloniki 57001, Greece, and also with the Department of Electrical and Electronic Engineering, Imperial College London, London SW7 2AZ, U.K. (email: gkalelis@iti.gr).

V. Mezaris and I. Kompatsiaris are with the Information Technologies Institute/Centre for Research and Technology Hellas (CERTH), Thessaloniki 57001, Greece (email: bmezaris@iti.gr; ikom@iti.gr).

T. Stathaki is with the Department of Electrical and Electronic Engineering, Imperial College London, London SW7 2AZ, U.K (email: t.stathaki@imperial.ac.uk).

of the between-class sum of squares  $\mathbf{S}_b = \sum_{i=1}^C \hat{p}_i (\hat{\boldsymbol{\mu}}_i - \hat{\boldsymbol{\mu}})(\hat{\boldsymbol{\mu}}_i - \hat{\boldsymbol{\mu}})^T$  to the within-class sum of squares  $\mathbf{S}_w = \sum_{i=1}^C \hat{p}_i \hat{\boldsymbol{\Sigma}}_i$ , where,  $\hat{p}_i = N_i/N$ ,  $\hat{\boldsymbol{\Sigma}}_i = \frac{1}{N_i} \sum_{n=1}^{N_i} (\mathbf{x}_i^n - \hat{\boldsymbol{\mu}}_i)(\mathbf{x}_i^n - \hat{\boldsymbol{\mu}}_i)^T$ ,  $\hat{\boldsymbol{\mu}}_i = \frac{1}{N_i} \sum_{n=1}^{N_i} \mathbf{x}_i^n$ ,  $\hat{\boldsymbol{\mu}} = \sum_{i=1}^C \hat{p}_i \hat{\boldsymbol{\mu}}_i$  are the estimated prior, the sample covariance matrix, the sample mean, and the total sample mean, respectively. This optimization problem turns out to be equivalent to the generalized eigenvalue decomposition  $\mathbf{S}_b \Psi = \mathbf{S}_w \Psi \Lambda$ , where the columns of  $\Psi$  are the generalized eigenvectors corresponding to the largest generalized eigenvalues in the diagonal matrix  $\Lambda$  [5].

Despite its elegant algebraic formulation, two important shortcomings of LDA restrict its use in real-world applications: a) The LDA criterion cannot be applied directly when the matrix  $\mathbf{S}_w$  is rank-deficient, a situation that occurs frequently in many applications involving *small sample size* (SSS) data. Several methods have been proposed to deal with this problem, including PCA+LDA [6], MMC LDA [7], dICA [8], and others. b) LDA faces difficulties in deriving a discriminant subspace when the classes are not linearly separable (a problem called hereafter *nonlinearity problem*). This problem has been mostly addressed by using kernel extensions of LDA, [9], [10] or methods that use local linear discriminant analyzers to learn the nonlinear data structure [2], [11]. However, the SSS problem remains, and to address it similar solutions to those discussed above are exploited for both the kernel-based [12], [13] and local-based [14] LDA variants.

Another strategy for solving the nonlinearity problem is to use a clustering procedure to derive a subclass division of the data, and then incorporate this information into the LDA criterion (again, the SSS problem is handled with techniques that overcome the rank-deficiency of  $\mathbf{S}_w$ , e.g. see [15]). The main advantage of this strategy over the methods described in the previous paragraph (especially over the kernel-based variants of LDA) is that it offers faster computation times during testing, because it only involves a single matrix multiplication. This is the underlying principle of mixture discriminant analysis (MDA) [16] that utilizes the following criterion

$$J_{MDA}(\Psi) = \frac{\text{tr}(\Psi^T \mathbf{S}_{bs} \Psi)}{\text{tr}(\Psi^T \mathbf{S}_{ws} \Psi)}, \quad (2)$$

where  $\mathbf{S}_{bs} = \sum_{i=1}^C \sum_{j=1}^{H_i} \hat{p}_{i,j} (\hat{\boldsymbol{\mu}}_{i,j} - \hat{\boldsymbol{\mu}})(\hat{\boldsymbol{\mu}}_{i,j} - \hat{\boldsymbol{\mu}})^T$  is the between-subclass scatter,  $\mathbf{S}_{ws} = \sum_{i=1}^C \sum_{j=1}^{H_i} \hat{p}_{i,j} \hat{\boldsymbol{\Sigma}}_{i,j}$  is the within-subclass scatter matrix,  $H_i$  denotes the number of subclasses of the  $i$ -th class, and  $\hat{p}_{i,j}$ ,  $\hat{\boldsymbol{\mu}}_{i,j}$ ,  $\hat{\boldsymbol{\Sigma}}_{i,j}$  are the estimated prior, sample mean and sample covariance matrix of the  $j$ -th subclass of class  $i$ .

As our target is to derive a subspace that best separates

observations of different classes, a better choice is to define a discriminant metric that favors the scatter of means between subclasses of different classes. This idea is exploited in subclass discriminant analysis (SDA) [17] that defines the following criterion

$$J_{SDA}(\Psi) = \text{tr}(\Psi^T \mathbf{S}_{bsb} \Psi) / \text{tr}(\Psi^T \Sigma_X \Psi), \quad (3)$$

where  $\mathbf{S}_{bsb} = \sum_{i=1}^{C-1} \sum_{j=1}^{H_i} \sum_{k=i+1}^C \sum_{l=1}^{H_k} \hat{p}_{i,j} \hat{p}_{k,l} (\hat{\boldsymbol{\mu}}_{i,j} - \hat{\boldsymbol{\mu}}_{k,l})(\hat{\boldsymbol{\mu}}_{i,j} - \hat{\boldsymbol{\mu}}_{k,l})^T$  is the inter-between-subclass scatter matrix, representing the scatter between the means of subclasses of different classes (inter-subclass scatter of means), and  $\Sigma_X = \frac{1}{N} \sum_{i=1}^C \sum_{j=1}^{H_i} \sum_{n=1}^{N_{i,j}} (\mathbf{x}_{i,j}^n - \hat{\boldsymbol{\mu}})(\mathbf{x}_{i,j}^n - \hat{\boldsymbol{\mu}})^T$  is the total covariance matrix.

Several extensions of MDA [18]–[21], and SDA [22]–[25] have been proposed, mainly seeking a more effective subclass partitioning procedure. In [26], mixture subclass discriminant analysis (MSDA) is presented, where it is explained that the use of the criterion

$$J_{MSDA}(\Psi) = \text{tr}(\Psi^T \mathbf{S}_{bsb} \Psi) / \text{tr}(\Psi^T \check{\Sigma}_X \Psi), \quad (4)$$

where,  $\check{\Sigma}_X = \mathbf{S}_{bsb} + \mathbf{S}_{ws}$ , is a better choice than the SDA criterion (3). Moreover, this algorithm assumes that the data have a Gaussian homoscedastic subclass structure and introduces an appropriate subclass partitioning procedure along with a nongaussianity criterion to derive the subclass division that optimizes the MSDA criterion. In [26], it was shown that in most cases MSDA outperforms SDA and other LDA variants. However, as we explain in the following, there is still room for further improving dimensionality reduction along the following directions:

1) *Link to Gaussian model*: In [16], [27], [28], it was shown that the LDA and MDA subspaces (defined by the column vectors of the respective projection matrix) coincide with the subspace that maximizes the log-likelihood function of Gaussian class densities or Gaussian mixture class densities, respectively, under the assumption that all class densities (or mixture component densities) are homoscedastic and that all class discriminant information is confined in a  $D$ -dimensional subspace of the  $F$ -dimensional measurement space. A respective link between MSDA (or SDA) and an appropriate Gaussian model has not yet been provided in the literature, and such a link could lead to a new DR approach.

2) *Subclass separation problem*: When the dimensionality of the LDA subspace is strictly lower than the rank of the between-class matrix, i.e.,  $D < C - 1$ , the projection of the class densities to the discriminant subspace may smear the neighboring classes in the measurement space, a situation described as the class-separation problem [29]–[31]. The same problem can equivalently occur to MSDA (and other subclass variants of LDA), i.e., neighboring subclasses in the original feature space may overlap in the projection subspace when the MSDA subspace dimensionality is strictly lower than the rank of the inter-between-subclass scatter matrix. We refer to this situation as the subclass separation problem.

3) *Subclass nonlinearity problem*: MSDA (and other subclass variants of LDA) can resolve the problem of nonlinearly separable classes as long as a subclass division that results in

linearly separable subclasses is identified. If this is not possible, a subclass-based approach that can deal with nonlinearly separable subclasses is desirable, often using an appropriate kernel to map the nonlinearly separable subclass divisions into a new space where they are linearly separable. For instance, in [32], [33] the kernel SDA (KSDA) method was shown to outperform a number of other approaches including kernel discriminant analysis (KDA) [9] and kernel support vector machines (KSVM) [34].

Inspired from the above discussion, in this paper we first provide an explicit link between MSDA and an appropriate Gaussian model, which allows the derivation of a new DA method under the Expectation Maximization framework (EM-MSDA). Furthermore, we present two additional methods, fractional-step MSDA (FSMSDA) and kernel MSDA (KMSDA), to alleviate the subclass separation problem of MSDA and to handle cases where MSDA subclasses are not linearly separable, respectively.

The rest of the paper is structured as follows: In Section II a link between MSDA and a Gaussian model is provided and EM-MSDA is derived, while in Sections III and IV, FSMSDA and KMSDA are presented. In Section V experimental results are reported and Section VI concludes the paper.

## II. LINK TO GAUSSIAN MODEL

In this section, we initially provide a Gaussian mixtures model formulation of the classification task, and then show how the Expectation Maximization (EM) algorithm [35]–[37] can be applied to estimate the unknown model parameters. Through this treatment we provide an explicit link between MSDA and the described Gaussian model, and consequently derive the EM-MSDA algorithm.

### A. Gaussian mixtures model

Let  $\omega_1, \dots, \omega_C$  be a finite set of  $C$  states of nature (classes) and  $(X, Y)$  be an  $\mathcal{X} \times I_C$ -valued random pair, where  $\mathcal{X} \subset \mathbb{R}^F$  is the space of observations and  $I_C = \{1, \dots, C\}$  is the class indicator variable [2], [3], [38]. Under this framework we model the  $i$ -th class-conditional probability density function  $p(\mathbf{x}|\omega_i)$  as a multivariate Gaussian mixture density of  $H_i$  component densities where the mixture components along all classes are homoscedastic [16], i.e.

$$p(\mathbf{x}|\omega_i) = \sum_{j=1}^{H_i} \pi_{i,j} \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_{i,j}), \quad (5)$$

where,  $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_{i,j}) = (\tau)^{-F/2} |\Sigma|^{-1/2} \exp((-1/2)\Delta(\mathbf{x}, \boldsymbol{\mu}_{i,j}))$  is the  $j$ -th component density (subclass) of the  $i$ -th mixture with constant  $\tau \approx 6.283185\dots$ , nonnegative mixing coefficient  $\pi_{i,j}$  (satisfying  $\sum_{j=1}^{H_i} \pi_{i,j} = 1$ ), mean vector  $\boldsymbol{\mu}_{i,j}$  and covariance matrix  $\Sigma$  shared along all mixture components. Moreover,  $\Delta(\mathbf{x}, \boldsymbol{\mu}_{i,j}) = (\mathbf{x} - \boldsymbol{\mu}_{i,j})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_{i,j})$  is the Mahalanobis distance between observation  $\mathbf{x}$  and the  $j$ -th component of class  $i$ .

We then wish to obtain a  $D < F$ -dimensionality reduction of the data which favors the separability of those subclasses that correspond to different classes. Consequently,

the parameter vector of the presented model is formed as  $\theta = [\pi_{1,1}, \boldsymbol{\mu}_{1,1}^T, \dots, \pi_{C,H_C}, \boldsymbol{\mu}_{C,H_C}^T, \text{vec}(\boldsymbol{\Sigma})^T, \text{vec}(\boldsymbol{\Psi})^T]^T$ , where  $\boldsymbol{\Psi} \in \mathbb{R}^{F \times D}$  is the required projection matrix for mapping the data into the reduced subspace,  $T$  is the vector transposition operator, and the  $\text{vec}()$  operator stacks the matrix columns to a vector.

### B. Log-likelihood function

For the estimation of the unknown parameters  $\theta$  we resort to the EM algorithm. The EM algorithm is based on the interpretation of the observed data set  $\mathbf{X}_i$  of  $i$ -th class as incomplete, where the missing part is a corresponding set  $\mathbf{Z}_i = [\mathbf{z}_i^1, \dots, \mathbf{z}_i^{N_i}]$  of categorical vectors  $\mathbf{z}_i^n = [z_{i,1}^n, \dots, z_{i,H_i}^n]^T$ , in which only a particular element  $z_{i,\kappa}^n$  equals to 1, indicating that  $\mathbf{x}_i^n$  was produced from the  $\kappa$ -th component (subclass) of the  $i$ -th mixture density. Under the above formulation and assuming that the  $C$  data matrices (blocks) of the block matrix  $\mathbf{X}$  (Section I) are independent as well as that the column vectors of the  $i$ -th block constitute a random sample from the population with density  $p(\mathbf{x}|\omega_i)$  (i.e., all observation vectors are independent and identically distributed (i.i.d.)), the log-likelihood function  $\mathcal{L}_1$  of the complete dataset would be (similar to [16], [27] – see Appendix A-A)

$$\begin{aligned} 2\mathcal{L}_1 = & \sum_{i=1}^C \sum_{j=1}^{H_i} 2\tilde{N}_{i,j} \ln \pi_{i,j} - NF \ln(2\pi) + N \ln(\det \boldsymbol{\Sigma}^{-1}) \\ & - \sum_{i=1}^C \sum_{j=1}^{H_i} \tilde{N}_{i,j} (\bar{\mathbf{x}}_{i,j} - \boldsymbol{\mu}_{i,j})^T \boldsymbol{\Sigma}^{-1} (\bar{\mathbf{x}}_{i,j} - \boldsymbol{\mu}_{i,j}) \\ & - \sum_{i=1}^C \sum_{n=1}^{N_i} \sum_{j=1}^{H_i} h_{i,j}^n (\mathbf{x}_i^n - \bar{\mathbf{x}}_{i,j})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i^n - \bar{\mathbf{x}}_{i,j}), \end{aligned} \quad (6)$$

where  $h_{i,j}^n$  are the responsibilities, i.e., the expected values of the categorical variables  $z_{i,j}^n$  for each data point, given by

$$h_{i,j}^n = \mathbb{E}[z_{i,j}^n] = \frac{\hat{\pi}_{i,j} \mathcal{N}(\mathbf{x}_i^n | \hat{\boldsymbol{\mu}}_{i,j})}{\sum_{j=1}^{H_i} \hat{\pi}_{i,j} \mathcal{N}(\mathbf{x}_i^n | \hat{\boldsymbol{\mu}}_{i,j})}, \quad (7)$$

and  $\bar{\mathbf{x}}_{i,j} = (1/\tilde{N}_{i,j}) \sum_{n=1}^{N_i} h_{i,j}^n \mathbf{x}_i^n$ ,  $\tilde{N}_{i,j} = \sum_{n=1}^{N_i} h_{i,j}^n$  are the weighted sample mean and the effective number of points of the  $j$ -th component of the  $i$ -th mixture respectively (note from (7) that  $\sum_{j=1}^{H_i} \tilde{N}_{i,j} = N_i$ ). Moreover,  $\ln \delta$  and  $\det \mathbf{A}$  denote the natural logarithm of number  $\delta$  and the determinant of matrix  $\mathbf{A}$ , respectively. We can rewrite (6) more compactly as

$$\begin{aligned} 2\mathcal{L}_1 = & \zeta - \sum_{i=1}^C \sum_{j=1}^{H_i} \tilde{N}_{i,j} (\bar{\mathbf{x}}_{i,j} - \boldsymbol{\mu}_{i,j})^T \boldsymbol{\Sigma}^{-1} (\bar{\mathbf{x}}_{i,j} - \boldsymbol{\mu}_{i,j}) \\ = & \zeta - \text{tr}\{\mathbf{N}(\bar{\mathbf{X}} - \mathbf{M})^T \boldsymbol{\Sigma}^{-1} (\bar{\mathbf{X}} - \mathbf{M})\} \end{aligned} \quad (8)$$

where  $\zeta$  is the part of the log-likelihood function that is independent of the true means ( $\zeta = \sum_{i=1}^C \sum_{j=1}^{H_i} 2\tilde{N}_{i,j} \ln \pi_{i,j} - NF \ln(2\pi) + N \ln(\det \boldsymbol{\Sigma}^{-1}) - \sum_{i=1}^C \sum_{n=1}^{N_i} \sum_{j=1}^{H_i} h_{i,j}^n (\mathbf{x}_i^n - \bar{\mathbf{x}}_{i,j})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i^n - \bar{\mathbf{x}}_{i,j})$ ), and  $\mathbf{M} = [\boldsymbol{\mu}_{1,1}, \dots, \boldsymbol{\mu}_{C,H_C}]$ ,  $\bar{\mathbf{X}} = [\bar{\mathbf{x}}_{1,1}, \dots, \bar{\mathbf{x}}_{C,H_C}]$  are the matrices of true means and weighted sample means respectively.

1) *Constrained M*: We wish to impose two constraints on the values of the true means as we explain in the following. Firstly, we require that the discriminant information is confined in a  $D$ -dimensional subspace of the original  $F$ -dimensional measurement space (e.g., see p. 339 [27], [16], [28]). Under this restriction the mean of the  $j$ -th mixture component of the  $i$ -th class density is expressed as

$$\boldsymbol{\mu}_{i,j} = \boldsymbol{\mu}_o + \boldsymbol{\Sigma} \boldsymbol{\Psi} \mathbf{v}_{i,j} \quad (9)$$

where,  $\boldsymbol{\Psi} \in \mathbb{R}^{F \times D}$  is a singular transformation matrix with uncorrelated column vectors that transforms  $\boldsymbol{\Sigma}$  into the unit matrix

$$\boldsymbol{\Psi}^T \boldsymbol{\Sigma} \boldsymbol{\Psi} = \mathbf{I}, \quad (10)$$

$\boldsymbol{\mu}_o$  is the total mean, and  $\mathbf{v}_{i,j} \in \mathbb{R}^D$  is the projection of  $\boldsymbol{\mu}_{i,j}$  into the lower-dimensional subspace. The latter is clear if we rearrange (9) to yield  $\mathbf{v}_{i,j} = \boldsymbol{\Psi}^T (\boldsymbol{\mu}_{i,j} - \boldsymbol{\mu}_o)$ . In matrix form (9) can be written as

$$\mathbf{M} = \mathbf{M}_o + \boldsymbol{\Sigma} \boldsymbol{\Psi} \boldsymbol{\Upsilon} \quad (11)$$

where,  $\mathbf{M}$  now is of column rank  $D$ ,  $\mathbf{M}_o = [\boldsymbol{\mu}_o, \dots, \boldsymbol{\mu}_o]$  is the  $F \times H$  matrix whose column vectors equal to the total mean  $\boldsymbol{\mu}_o$ , and  $\boldsymbol{\Upsilon} = [\mathbf{v}_{1,1}, \dots, \mathbf{v}_{C,H_C}]$  is the matrix with the projection coefficients of the mean vectors.

Secondly, we wish to penalize (6) such that in the lower dimensional subspace the between-subclass spread is emphasized relative to the within-subclass spread. We can impose this by penalizing (6) with the term  $\text{tr}\{\boldsymbol{\Upsilon} \mathbf{Q} \boldsymbol{\Upsilon}^T\}$ . The penalty matrix  $\mathbf{Q}$  is defined as

$$\mathbf{Q} = \mathbf{N} \mathbf{A}^{-1} \mathbf{N} - \mathbf{N} \quad (12)$$

where  $\mathbf{N} = \text{diag}(\tilde{N}_{1,1}, \dots, \tilde{N}_{C,H_C})$  is an  $H \times H$  diagonal matrix with diagonal elements the effective sample numbers of the respective mixture component,  $\mathbf{A}$  is a symmetric matrix that allows us to express the weighted inter-between-subclass scatter matrix

$$\mathbf{S}_{bsb}^w = \sum_{i=1}^{C-1} \sum_{j=1}^{H_i} \sum_{k=i+1}^C \sum_{l=1}^{H_k} \tilde{p}_{i,j} \tilde{p}_{k,l} (\bar{\mathbf{x}}_{i,j} - \bar{\mathbf{x}}_{k,l})(\bar{\mathbf{x}}_{i,j} - \bar{\mathbf{x}}_{k,l})^T, \quad (13)$$

in a matrix product form,

$$\mathbf{S}_{bsb}^w = \bar{\mathbf{X}} \mathbf{A} \bar{\mathbf{X}}^T, \quad (14)$$

and  $\bar{\mathbf{X}} = [\bar{\mathbf{x}}_{1,1}, \dots, \bar{\mathbf{x}}_{C,H_C}]$  is the matrix of the weighted means. That is, the matrix element  $\mathbf{A}_{i,j,k,l}$  that corresponds to  $\bar{\mathbf{x}}_{i,j}$  and  $\bar{\mathbf{x}}_{k,l}$  weighted means takes the value

$$\mathbf{A}_{i,j,k,l} = \begin{cases} \tilde{p}_{i,j}(1 - \tilde{p}_{i,j}), & \text{if } (i,j) = (k,l), \\ 0 & \text{if } i = k, j \neq l, \\ -\tilde{p}_{i,j} \tilde{p}_{k,l} & \text{else} \end{cases} \quad (15)$$

where  $\tilde{p}_{i,j} = \tilde{N}_{i,j}/N$ ,  $\tilde{p}_i = \sum_{j=1}^{H_i} \tilde{p}_{i,j} = N_i/N$ . Notice that the sum of the components of any row vector (or any column vector) of matrix  $\mathbf{A}$  equals to zero. Therefore, for any matrix with equal column vectors  $\mathbf{B} = [\mathbf{b}, \dots, \mathbf{b}]$  the matrix product  $\mathbf{A} \mathbf{B}^T$  will yield the zero matrix. We should also note that  $\mathbf{Q}$  is symmetric and that for  $\mathbf{A} = \mathbf{N}$ ,  $\mathbf{Q}$  and consequently the penalty term vanish, leading to the conventional MDA algorithm [16]. As we will explain in the

sequel, such a specialization of the penalty matrix  $\mathbf{Q}$  will lead to an interesting extension of the MDA algorithm that will provide a subspace equivalent to the MSDA subspace.

2) *Constrained  $\hat{\Sigma}$* : Similarly to MDA with centroid shrinking (p.171, [16]) we constrain the covariance matrix at the weighted within-subclass scatter matrix

$$\begin{aligned}\hat{\Sigma} &= \mathbf{S}_{ws}^w = \sum_{i=1}^C \sum_{j=1}^{H_i} \tilde{p}_{i,j} \hat{\Sigma}_{i,j}^w \\ &= \frac{1}{N} \sum_{i=1}^C \sum_{j=1}^{H_i} \sum_{n=1}^{N_i} h_{i,j}^n (\mathbf{x}_i^n - \bar{\mathbf{x}}_{i,j})(\mathbf{x}_i^n - \bar{\mathbf{x}}_{i,j})^T,\end{aligned}\quad (16)$$

where  $\hat{\Sigma}_{i,j}^w = (1/\tilde{N}_{i,j}) \sum_{n=1}^{N_i} h_{i,j}^n (\mathbf{x}_i^n - \bar{\mathbf{x}}_{i,j})(\mathbf{x}_i^n - \bar{\mathbf{x}}_{i,j})^T$  is the weighted sample covariance matrix of  $(i, j)$  component density.

Imposing the above constraints (10), (11), (16), and the penalty term (12) in (8), we finally arrive to the following penalized and restricted version of the log-likelihood function

$$2\mathcal{L}_2 = -\text{tr}\{\mathbf{N}(\bar{\mathbf{X}} - \mathbf{M}_o - \hat{\Sigma}\Psi\Upsilon)^T \hat{\Sigma}^{-1}(\bar{\mathbf{X}} - \mathbf{M}_o - \hat{\Sigma}\Psi\Upsilon)\} - \text{tr}\{\Upsilon\mathbf{Q}\Upsilon^T\} + \zeta, \quad (17)$$

where  $\Psi$  is constrained by (10).

### C. EM algorithm

The EM algorithm can be applied to obtain the maximum likelihood estimate (MLE) of the model parameters in (17). This algorithm alternates between two steps, the Expectation step (E-step) and the Maximization step (M-step), to produce a sequence of estimates until some convergence criterion is met.

1) *E-step*: During the E-step, the parameter values identified in the previous EM cycle are used to compute the responsibilities  $h_{i,j}^n$  using (7).

2) *M-step*: In this step, the unknown mixture parameters are estimated by maximizing (17). In particular, we need to estimate the mixing coefficients  $\pi_{i,j}$  and the true means  $\mu_{i,j}$  for each mixture component in (5).

*Estimation of  $\pi_{i,j}$* : The mixing coefficients are estimated by maximizing (17) subject to the constraint that  $\sum_{j=1}^{H_i} \pi_{i,j} = 1$ , giving (similarly to [37] – see Appendix A-B)

$$\hat{\pi}_{i,j} = \frac{\tilde{N}_{i,j}}{N_i}. \quad (18)$$

*Estimation of  $\mu_{i,j}$* : Now we proceed to estimating the true means in  $\mathbf{M}$ , or equivalently  $\mathbf{M}_o$ ,  $\Upsilon$ , and  $\Psi$ , that maximize (17) subject to  $\Psi^T \hat{\Sigma} \Psi = \mathbf{I}$ . In (17)  $\zeta$  is independent of  $\mathbf{M}$  and thus can be discarded from the optimization criterion. Moreover, the maximization of  $\mathcal{L}_2$  is equivalent to the minimization of  $-\mathcal{L}_2$  under the same conditions, leading us to the following optimization problem

$$\underset{\mathbf{M}_o, \Upsilon, \Psi}{\text{argmin}} \mathcal{L}_3 \quad \text{subject to} \quad \Psi^T \hat{\Sigma} \Psi = \mathbf{I}, \quad (19)$$

where,

$$\mathcal{L}_3 = \text{tr}\{\mathbf{N}(\bar{\mathbf{X}} - \mathbf{M}_o - \hat{\Sigma}\Psi\Upsilon)^T \hat{\Sigma}^{-1}(\bar{\mathbf{X}} - \mathbf{M}_o - \hat{\Sigma}\Psi\Upsilon)\} + \text{tr}\{\Upsilon\mathbf{Q}\Upsilon^T\}. \quad (20)$$

Setting  $\bar{\mathbf{y}}_{i,j} = \hat{\Sigma}^{-1/2} \bar{\mathbf{x}}_{i,j}$ ,  $\mathbf{v}_{i,j} = \hat{\Sigma}^{-1/2} \mu_{i,j}$ , and  $\mathbf{v}_o = \hat{\Sigma}^{-1/2} \mu_o$ , we can write  $\mathbf{v}_{i,j} = \mathbf{v}_o + \tilde{\Psi} \mathbf{v}_{i,j}$  or in matrix form

$$\mathbf{V} = \mathbf{V}_o + \tilde{\Psi} \Upsilon \quad (21)$$

where  $\tilde{\Psi} = \hat{\Sigma}^{1/2} \Psi$ ,  $\bar{\mathbf{Y}} = \hat{\Sigma}^{-1/2} \bar{\mathbf{X}} = [\bar{\mathbf{y}}_{1,1}, \dots, \bar{\mathbf{y}}_{C,H_C}]$ ,  $\mathbf{V} = \hat{\Sigma}^{-1/2} \mathbf{M} = [\mathbf{v}_{1,1}, \dots, \mathbf{v}_{C,H_C}]$ , and  $\mathbf{V}_o = \hat{\Sigma}^{-1/2} \mathbf{M}_o$ . Substituting this to (20) we arrive to

$$\begin{aligned}\mathcal{L}_3 &= \sum_{i=1}^C \sum_{j=1}^{H_i} \tilde{N}_{i,j} (\bar{\mathbf{y}}_{i,j} - \mathbf{v}_{i,j})^T (\bar{\mathbf{y}}_{i,j} - \mathbf{v}_{i,j}) \\ &= \text{tr}\{\mathbf{N}(\bar{\mathbf{Y}} - \mathbf{V})^T (\bar{\mathbf{Y}} - \mathbf{V})\} + \text{tr}\{\Upsilon\mathbf{Q}\Upsilon^T\} \\ &= \text{tr}\{\mathbf{N}(\bar{\mathbf{Y}} - \mathbf{V}_o - \tilde{\Psi}\Upsilon)^T (\bar{\mathbf{Y}} - \mathbf{V}_o - \tilde{\Psi}\Upsilon)\} \\ &\quad + \text{tr}\{\Upsilon\mathbf{Q}\Upsilon^T\}\end{aligned}\quad (22)$$

Setting the derivatives of  $\mathcal{L}_3$  in (22) with respect to the projection coefficients  $\Upsilon$  to zero we obtain

$$\frac{\partial \mathcal{L}_3}{\partial \Upsilon} = 0 \Rightarrow \Upsilon = \tilde{\Psi}^T (\bar{\mathbf{Y}} - \mathbf{V}_o) \mathbf{A} \mathbf{N}^{-1} \quad (23)$$

We can now expand (22) as

$$\begin{aligned}\mathcal{L}_3 &= \text{tr}\{\mathbf{N}(\bar{\mathbf{Y}} - \mathbf{Y}_o + \mathbf{Y}_o - \mathbf{V}_o - \tilde{\Psi}\Upsilon)^T \\ &\quad \times (\bar{\mathbf{Y}} - \mathbf{Y}_o + \mathbf{Y}_o - \mathbf{V}_o - \tilde{\Psi}\Upsilon)\} + \text{tr}\{\Upsilon\mathbf{Q}\Upsilon^T\} \\ &= \text{tr}\{\mathbf{N}(\bar{\mathbf{Y}} - \mathbf{Y}_o)^T (\bar{\mathbf{Y}} - \mathbf{Y}_o)\} \\ &\quad + \text{tr}\{\mathbf{N}(\mathbf{Y}_o - \mathbf{V}_o)^T (\mathbf{Y}_o - \mathbf{V}_o)\} \\ &\quad + \text{tr}\{\mathbf{N}\Upsilon^T \tilde{\Psi}^T \tilde{\Psi}\Upsilon\} + \text{tr}\{\Upsilon\mathbf{Q}\Upsilon^T\} \\ &\quad + 2 \text{tr}\{\mathbf{N}(\bar{\mathbf{Y}} - \mathbf{Y}_o)^T (\mathbf{Y}_o - \mathbf{V}_o)\} \\ &\quad - 2 \text{tr}\{\mathbf{N}(\bar{\mathbf{Y}} - \mathbf{V}_o)^T \tilde{\Psi}\Upsilon\}\end{aligned}\quad (24)$$

Reformulating the fifth term of (24) we see that it vanishes

$$\begin{aligned}\text{tr}\{\mathbf{N}(\bar{\mathbf{Y}} - \mathbf{Y}_o)^T (\mathbf{Y}_o - \mathbf{V}_o)\} \\ = \sum_{i=1}^C \sum_{j=1}^{H_i} \tilde{N}_{i,j} (\bar{\mathbf{y}}_{i,j} - \mathbf{y}_o)^T (\mathbf{y}_o - \mathbf{v}_o) = 0,\end{aligned}\quad (25)$$

Using (12), (23) and taking into account that  $\tilde{\Psi}^T \tilde{\Psi} = \Psi^T \hat{\Sigma} \Psi = \mathbf{I}$  the summand of the third and fourth term of (24) becomes

$$\begin{aligned}\text{tr}\{\mathbf{N}\Upsilon^T \tilde{\Psi}^T \tilde{\Psi}\Upsilon\} + \text{tr}\{\Upsilon\mathbf{Q}\Upsilon^T\} \\ = \text{tr}\{\Upsilon(\mathbf{N} + \mathbf{Q})\Upsilon^T\} = \text{tr}\{\Upsilon \mathbf{N} \mathbf{A}^{-1} \mathbf{N} \Upsilon^T\} \\ = \text{tr}\{\mathbf{A}(\bar{\mathbf{Y}} - \mathbf{V}_o)^T \tilde{\Psi} \tilde{\Psi}^T (\bar{\mathbf{Y}} - \mathbf{V}_o)\}\end{aligned}\quad (26)$$

and similarly using (23) the sixth term of (24) becomes

$$\begin{aligned}\text{tr}\{\mathbf{N}(\bar{\mathbf{Y}} - \mathbf{V}_o)^T \tilde{\Psi}\Upsilon\} \\ = \text{tr}\{\mathbf{A}(\bar{\mathbf{Y}} - \mathbf{V}_o)^T \tilde{\Psi} \tilde{\Psi}^T (\bar{\mathbf{Y}} - \mathbf{V}_o)\}\end{aligned}\quad (27)$$

Substituting (25), (26), (27) into (24) we arrive to

$$\begin{aligned}\mathcal{L}_3 &= \text{tr}\{\mathbf{N}(\bar{\mathbf{Y}} - \mathbf{Y}_o)^T (\bar{\mathbf{Y}} - \mathbf{Y}_o)\} \\ &\quad + \text{tr}\{\mathbf{N}(\mathbf{Y}_o - \mathbf{V}_o)^T (\mathbf{Y}_o - \mathbf{V}_o)\} \\ &\quad - \text{tr}\{\mathbf{A}(\bar{\mathbf{Y}} - \mathbf{V}_o)^T \tilde{\Psi} \tilde{\Psi}^T (\bar{\mathbf{Y}} - \mathbf{V}_o)\}\end{aligned}\quad (28)$$

Using the fact that  $\mathbf{A} \mathbf{V}_o^T = \mathbf{0}$ , the last term of (28) is simplified to

$$\begin{aligned}\text{tr}\{\mathbf{A}(\bar{\mathbf{Y}} - \mathbf{V}_o)^T \tilde{\Psi} \tilde{\Psi}^T (\bar{\mathbf{Y}} - \mathbf{V}_o)\} &= \text{tr}\{\mathbf{A} \bar{\mathbf{Y}}^T \tilde{\Psi} \tilde{\Psi}^T \bar{\mathbf{Y}}\} \\ &\quad + \text{tr}\{\mathbf{A} \mathbf{V}_o^T \tilde{\Psi} \tilde{\Psi}^T \mathbf{V}_o\} - 2 \text{tr}\{\mathbf{A} \mathbf{V}_o^T \tilde{\Psi} \tilde{\Psi}^T \bar{\mathbf{Y}}\} \\ &= \text{tr}\{\mathbf{A} \bar{\mathbf{Y}}^T \tilde{\Psi} \tilde{\Psi}^T \bar{\mathbf{Y}}\},\end{aligned}\quad (29)$$

and substituting this back to (28) we arrive to

$$\mathcal{L}_3 = \text{tr}\{\mathbf{N}(\bar{\mathbf{Y}} - \mathbf{Y}_o)^T(\bar{\mathbf{Y}} - \mathbf{Y}_o)\} + \text{tr}\{\mathbf{N}(\mathbf{Y}_o - \mathbf{V}_o)^T(\mathbf{Y}_o - \mathbf{V}_o)\} - \text{tr}\{\tilde{\Psi}^T \bar{\mathbf{Y}} \mathbf{A} \bar{\mathbf{Y}}^T \tilde{\Psi}\}, \quad (30)$$

where,  $\mathbf{Y}_o = \hat{\Sigma}^{-1/2} \mathbf{X}_o$ , and  $\mathbf{X}_o$  is the  $F \times H$  matrix whose column vectors equal to the weighted mean  $\mathbf{x}_o = (1/N) \sum_{i=1}^C \sum_{j=1}^{H_i} \sum_{n=1}^{N_i} h_{i,j}^n \mathbf{x}_i^n = \sum_{i=1}^C \sum_{j=1}^{H_i} \tilde{p}_{i,j} \bar{\mathbf{x}}_{i,j}$ . We now have to minimize (30) with respect to  $\mathbf{V}_o$ , or equivalently

$$\text{tr}\{\mathbf{N}(\mathbf{Y}_o - \mathbf{V}_o)^T(\mathbf{Y}_o - \mathbf{V}_o)\} = N(\mathbf{y}_o - \mathbf{v}_o)^T(\mathbf{y}_o - \mathbf{v}_o) \quad (31)$$

which is minimized for  $\mathbf{y}_o = \mathbf{v}_o$  and, thus, yielding  $\hat{\mu}_o = \mathbf{x}_o$  or in matrix form

$$\hat{\mathbf{M}}_o = \mathbf{X}_o. \quad (32)$$

Without loss of generality we can set  $\mathbf{x}_o = [0, \dots, 0]^T$  (e.g. setting  $\mathbf{X} \leftarrow \mathbf{X} - \mathbf{X}_o$ ). Substituting this back to (30) we arrive to

$$\mathcal{L}_3 = \text{tr}\{\hat{\Sigma}^{-1} \bar{\mathbf{X}} \mathbf{N} \bar{\mathbf{X}}^T\} - \text{tr}\{\Psi^T \bar{\mathbf{X}} \mathbf{A} \bar{\mathbf{X}}^T \Psi\} \quad (33)$$

where we have used the requirement that  $\Psi$  transforms the pooled covariance matrix  $\hat{\Sigma}$  into the unit matrix ( $\Psi^T \hat{\Sigma} \Psi = \mathbf{I}$ ). In (33) only the second term depends on the transformation matrix, and, thus, this matrix can be obtained by solving the following optimization problem

$$\underset{\Psi}{\text{argmax}} \text{tr}\{\Psi^T \mathbf{S}_{bsb}^w \Psi\} \quad \text{subject to} \quad \Psi^T \mathbf{S}_{ws}^w \Psi = \mathbf{I} \quad (34)$$

where we have used (14) and fixed  $\hat{\Sigma}$  according to (16). The solution to this problem is obtained by the set  $\{\psi_i | i = 1, \dots, D\}$  of the generalized eigenvectors of  $\mathbf{S}_{bsb}^w$  and  $\mathbf{S}_{ws}^w$  corresponding to the  $D$  largest eigenvalues  $\{\lambda_i | i = 1, \dots, D\}$  of the following generalized eigenvalue decomposition [2]

$$\mathbf{S}_{bsb}^w \Psi = \mathbf{S}_{ws}^w \Psi \Lambda \quad (35)$$

where  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_D)$ . Therefore, the subspace that maximizes the constrained log-likelihood function in (19) at each EM cycle coincides with the subspace that maximizes the MSDA criterion, where the scatter matrices in (4) are replaced by their weighted equivalent in each EM cycle. The MLE of the true means can now be computed by substituting (23), (32) into (11) and using the computed estimates of (16), (35) for  $\mathbf{S}_{ws}$  and  $\Psi$  respectively

$$\begin{aligned} \hat{\mathbf{M}} &= \mathbf{X}_o + \mathbf{S}_{ws}^w \Psi \Psi^T (\bar{\mathbf{X}} - \mathbf{X}_o) \mathbf{A} \mathbf{N}^{-1} \\ &= \mathbf{S}_{ws}^w \Psi \Psi^T \bar{\mathbf{X}} \mathbf{A} \mathbf{N}^{-1} \end{aligned} \quad (36)$$

where, we have assumed that  $\mathbf{X}_o = \mathbf{0}$ .

### D. Model selection

The Gaussian model described above as well as the derived EM algorithm assume that the number of mixing components in each Gaussian mixture density is provided. However, this information is rarely known. In order to estimate the optimum number of mixing components for each mixture density with respect to the given training set, we utilize an iterative procedure, where at each iteration a new Gaussian model is specified (with respect to the number of mixture components) and a nongaussianity measure  $\Phi$  is evaluated in order to assess the goodness of fit of the particular Gaussian model. This

iterative process is repeated until the nongaussianity measure  $\Phi$  converges to a small value as explained in the following.

Skewness and kurtosis can be used to provide an indication of how well a particular Gaussian mixture density fits the training data of a specific class [26], [39], [40]. Estimates of the weighted standardized skewness  $\hat{\beta}_{i,j,f}$  and kurtosis  $\hat{\gamma}_{i,j,f}$  along the  $f$ -th dimension regarding the  $j$ -th mixture component of the  $i$ -th class can be computed as follows

$$\hat{\beta}_{i,j,f} = \frac{\frac{1}{N_{i,j}} \sum_{n=1}^{N_i} h_{i,j}^n (x_{i,j,f}^n - \hat{\mu}_{i,j,f})^3}{\hat{\sigma}_{i,j,f}^3}, \quad (37)$$

$$\hat{\gamma}_{i,j,f} = \frac{\frac{1}{N_{i,j}} \sum_{n=1}^{N_i} h_{i,j}^n (x_{i,j,f}^n - \hat{\mu}_{i,j,f})^4}{\hat{\sigma}_{i,j,f}^4} - 3, \quad (38)$$

where  $x_{i,j,f}^n$  is the  $f$ -th element of  $\mathbf{x}_{i,j}^n$ , and  $\hat{\mu}_{i,j,f}, \hat{\sigma}_{i,j,f}$  are the sample mean and standard deviation of the  $j$ -th mixture of  $i$ -th class along the  $f$ -th dimension. The above estimates will be close to zero for Gaussian densities and deviate from zero the more the underlying density deviates from the Gaussian. We can thus obtain an estimate of the skewness  $\hat{\beta}_{i,j}$  and kurtosis  $\hat{\gamma}_{i,j}$  of the  $(i, j)$  component density by averaging along all dimensions, i.e.,  $\hat{\beta}_{i,j} = (1/F) \sum_{f=1}^F |\hat{\beta}_{i,j,f}|$ ,  $\hat{\gamma}_{i,j} = (1/F) \sum_{f=1}^F |\hat{\gamma}_{i,j,f}|$ , where  $|a|$  denotes absolute value of  $a$ . Similarly, we can define a nongaussianity measure regarding the Gaussian mixture density referring to the  $i$ -th class using

$$\Phi_i = \sum_{j=1}^{H_i} \hat{\pi}_{i,j} (\hat{\beta}_{i,j} + \hat{\gamma}_{i,j}). \quad (39)$$

A large value of  $\Phi_i$  will denote that the respective Gaussian mixture density does not fit well the underlying density function of the  $i$ -th class training data. Therefore, at each iteration this measure is used to select the mixture density that yielded the worst fit according to the following criterion

$$k = \underset{i=1, \dots, C}{\text{argmax}} (\Phi_i), \quad (40)$$

and the required number of mixture components referring to this mixture density is increased by one ( $H_k \leftarrow H_k + 1$ ). Similarly, at each iteration a total nongaussianity measure is defined for assessing the fitness of the current Gaussian model with respect to the overall training data set

$$\Phi = \sum_{i=1}^C \tilde{p}_i \Phi_i. \quad (41)$$

The value of  $\Phi$  is examined at each iteration, and the iterative procedure is completed upon the convergence of  $\Phi$  to a steady-state solution. The resulting EM-MSDA algorithm is outlined in Algorithm 1. Alternatively, a cross-validation criterion can be used to select the Gaussian model that provides the best empirical recognition rate.

## III. FRACTIONAL STEP MIXTURE SUBCLASS DISCRIMINANT ANALYSIS

In equivalence to the class separation problem of LDA [29]–[31], the subclass separation problem may occur when the dimensionality of the MSDA subspace  $D$  is strictly lower

**Algorithm 1** EM-MSDA**Input:** Annotated data set  $\mathbf{X}$ **Output:**  $\Psi$ 

- 1: Initialize:  $H_1 = \dots = H_C = 1$ ,  $H = C$ ,  $\Phi_i$  (39),  $\Phi$  (41)
- 2: **repeat**
- 3:   Compute class label  $k$  of class to repartition (40)
- 4:   Set:  $H_k \leftarrow H_k + 1$
- 5:   Repartition  $k$ -th class to  $H_k$  subclasses using k-means
- 6:   Initialize the MLE parameters  $\hat{\theta}$
- 7:   **repeat**
- 8:     *E-step:* Compute responsibilities  $h_{i,j}^n$  (7)
- 9:     *M-step:* Compute MLEs:  $\Psi$  (35),  $\theta$  (16), (36)
- 10:   **until** convergence of  $\hat{\theta}$
- 11:   Compute nongaussianity  $\Phi_i$  for each class (39)
- 12:   Compute total nongaussianity  $\Phi$  (41)
- 13: **until** convergence of  $\Phi$

than the rank of the inter-between-subclass scatter matrix ( $D < \text{rank}(\mathbf{S}_{bsb}) \leq \min(F, H - 1)$ ). When this happens, distinct subclasses in the measurement space may not separate well in the lower dimensional subspace. To demonstrate this problem we use an artificial dataset of two classes, where, the first class consists of two Gaussian subclasses  $\mathcal{N}_{1,1}$ ,  $\mathcal{N}_{1,2}$ , and the second class is a unimodal Gaussian  $\mathcal{N}_{2,1}$ . The means of the Gaussian distributions are  $\mu_{1,1} = [6 \ 22]^T$ ,  $\mu_{1,2} = [0 \ 0]^T$ ,  $\mu_{2,1} = [12 \ 22]^T$ , whereas a common covariance matrix is shared along all distributions  $\Sigma = [0.7 \ 0.3; 0.3 \ 0.7]$ , as depicted in Figure 1. Under these settings, we see that the one-dimensional projection transformation derived using MSDA ( $\psi_{MSDA}$ ) causes a large overlap between the subclasses  $\mathcal{N}_{1,1}$  and  $\mathcal{N}_{2,1}$ , which are close to each other, but well separated in the measurement space. This happens because the large subclass distance  $d_{1,2,2,1} = \|\mu_{1,2} - \mu_{2,1}\|^2$  dominates the MSDA criterion, and, thus, the derived projection transformation preserves well the separation of the subclasses  $\mathcal{N}_{1,2}$  and  $\mathcal{N}_{2,1}$ , while, on the other hand, merges the two subclasses that are close together in the measurement space,  $\mathcal{N}_{1,1}$  and  $\mathcal{N}_{2,1}$ .

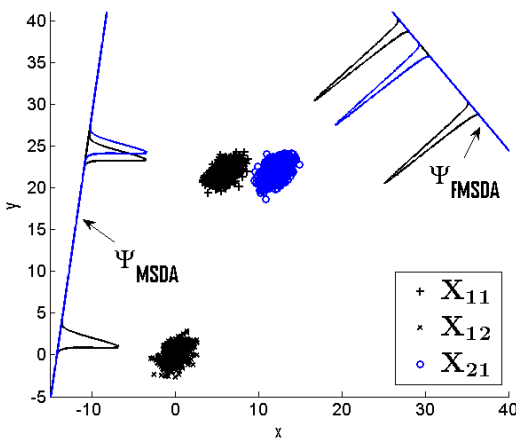


Fig. 1. Subclass separation problem.

To overcome the subclass separation problem, inspired from [31], we introduce the fractional-step MSDA (FMSDA) that

utilizes the following objective function

$$J_{FMSDA}(\Psi) = \frac{\text{tr}(\Psi^T \tilde{\mathbf{S}}_{bsb} \Psi)}{\text{tr}(\Psi^T \tilde{\Sigma}_{\mathbf{X}} \Psi)}, \quad (42)$$

where the inter-between-subclass scatter matrix is modified using an appropriate weighting function  $w_{i,j,k,l}$

$$\tilde{\mathbf{S}}_{bsb} = \sum_{i=1}^{C-1} \sum_{j=1}^{H_i} \sum_{k=i+1}^C \sum_{l=1}^{H_k} w_{i,j,k,l} (\mu_{i,j} - \mu_{k,l})(\mu_{i,j} - \mu_{k,l})^T, \quad (43)$$

and the modified covariance matrix accordingly becomes  $\tilde{\Sigma}_{\mathbf{X}} \equiv \tilde{\mathbf{S}}_{bsb} + \mathbf{S}_{ws}$ . The weighting function is a monotonically decreasing function defined as  $w_{i,j,k,l} = d_{i,j,k,l}^{-r}$ , where,  $d_{i,j,k,l} = \|\mu_{i,j} - \mu_{k,l}\|^2$  is the euclidian distance between the estimated means of subclasses  $(i, j)$  and  $(k, l)$ , and  $r$  is an integer number larger than two.

The FMSDA algorithm (Algorithm 2) starts with the application of the subclass partitioning procedure described in the previous section (Eqs. (37) to (41)) to derive a subclass division of the data. Then, the FMSDA criterion (42) is utilized to initialize the projection transformation matrix  $\Psi_D \in \mathbb{R}^{F \times D}$ , and an iterative algorithm is applied, where at each iteration  $\rho$  fractional steps are used for decreasing the dimensionality of the subspace by one. That is, at the  $t$ -th fractional step of the  $k$ -th iteration the data are projected in the  $k$ -th dimensional subspace using the transformation matrix  $\Psi_k \in \mathbb{R}^{F \times k}$ , scaled utilizing the following scaling transformation

$$\vartheta(\mathbf{y}, t) = \begin{cases} \alpha^t y_i, & i = k \\ y_i, & i = 1, \dots, k-1, \end{cases} \quad (44)$$

where  $\alpha = \exp(-\ln(\rho)/(\rho-1))$ , and the transformation matrix  $\Psi_k$  is recomputed using the projected and scaled data. At the end of this fractional procedure the last,  $k$ -th eigenvector of  $\Psi_k$  (i.e., the one that corresponds to the smallest eigenvalue of  $\Psi_k$ ) is discarded. The scaling transformation compresses the data along the direction of the last eigenvector of  $\Psi_k$ . This allows the subclass means that are along the direction of the  $k$ -th eigenvector to be increasingly weighted in the next fractional step, causing the  $k$ -dimensional subspace to reorient so that a useful projection direction is not discarded at the end of each iteration. A validation set is used to assess the performance of the derived projection matrix  $\Psi_k$  at each iteration, and the one that provided the best correct classification rate (CCR) is selected.

The main advantage of FMSDA (and also EM-MSDA) over kernel variants of LDA is that the projection matrix still constitutes a linear transformation, which can provide real time performance during the testing stage. On the other hand, in contrast to EM-MSDA that tends to optimize the fit of the subclasses and simultaneously seek the projection that maximizes the inter-subclass scatter of means, FMSDA derives an initial subclass structure of the data and gradually attempts to identify the subspace that provides the best empirical recognition rate.

#### IV. KERNEL MIXTURE SUBCLASS DISCRIMINANT ANALYSIS

The methods described in the previous sections will still not perform well when it is not possible to identify a sub-

**Algorithm 2** FMSDA

**Input:** Annotated set  $\mathbf{X}$ , validation set  $\mathbf{G}$ , parameters  $\rho, r$ 
**Output:**  $\Psi$ 


---

```

1: Initialize:  $H_1 = \dots = H_C = 1, H = C, \Phi_i$  (39),  $\Phi$  (41)
2: repeat
3:   Compute class label  $k$  of class to repartition (40)
4:   Set  $H_k \leftarrow H_k + 1$  and repartition  $k$ -th class
5:   Compute nongaussianity values  $\Phi_i$  (39) and  $\Phi$  (41)
6: until convergence of  $\Phi$ 
7: Compute  $\Psi_D$  (42), set  $D = \text{rank}(\tilde{\mathbf{S}}_{bsb})$  (43)
8: Set  $\text{CCR}_k = 0, k = 1, \dots, D$ 
9: for  $k = D$  to 1 do
10:  for  $t = 0$  to  $\rho - 1$  do
11:    Project training data:  $\mathbf{y} = \Psi_k^T \mathbf{x}$ 
12:    Apply scaling transformation:  $\tilde{\mathbf{y}} = \vartheta(\mathbf{y}, t)$ 
13:    Compute  $\tilde{\Psi}$  (42) using scaled data
14:    Set:  $\Psi_k \leftarrow \Psi_k \tilde{\Psi}$ 
15:  end for
16:  Discard the last ( $k$ -th) column of  $\Psi_k$ 
17:  Project and classify validation samples using  $\Psi_k$ 
18:  if sample  $\mathbf{g}_j$  is classified correctly then
19:     $\text{CCR}_k + +$ 
20:  end if
21: end for
22: Set:  $k_o = \text{argmax}_k(\text{CCR}_k); \Psi = \Psi_{k_o}$ 
    
```

---

class division that results in linearly separable classes [32], [33]. To deal with such cases, a nonlinear feature mapping  $\phi(\cdot) : \mathbb{R}^F \mapsto \mathcal{F}$  can be used to map the partitioned data into some high- or even infinite-dimensional feature space  $\mathcal{F}$ , where the data are expected to be linearly separable. Given a subclass partition of the data  $\mathbf{X} = [\mathbf{X}_{1,1}, \dots, \mathbf{X}_{C,H_C}]$ , where  $\mathbf{X}_{i,j} = [\mathbf{x}_{i,j}^1, \dots, \mathbf{x}_{i,j}^{N_{i,j}}]$  contains the observations of the  $(i, j)$  subclass, the transformation matrix  $\mathbf{W}$  that maximizes the MSDA criterion in  $\mathcal{F}$  can be computed from the following generalized eigenvalue problem

$$\mathbf{S}_{bsb}^\phi \mathbf{W} = \tilde{\Sigma}_{\mathbf{X}}^\phi \mathbf{W} \mathbf{\Lambda}^\phi \quad (45)$$

where,

$$\mathbf{S}_{bsb}^\phi = \sum_{i=1}^{C-1} \sum_{j=1}^{H_i} \sum_{k=i+1}^C \sum_{l=1}^{H_k} \hat{p}_{i,j} \hat{p}_{k,l} (\hat{\boldsymbol{\mu}}_{i,j}^\phi - \hat{\boldsymbol{\mu}}_{k,l}^\phi) (\hat{\boldsymbol{\mu}}_{i,j}^\phi - \hat{\boldsymbol{\mu}}_{k,l}^\phi)^T,$$

$$\tilde{\Sigma}_{\mathbf{X}}^\phi = \mathbf{S}_{bsb}^\phi + \mathbf{S}_{ws}^\phi, \quad \mathbf{S}_{ws}^\phi = \sum_{i=1}^C \sum_{j=1}^{H_i} \hat{p}_{i,j} \hat{\Sigma}_{i,j}^\phi,$$

are the inter-between-subclass scatter matrix, the within-subclass scatter matrix, the modified total sample covariance matrix, and  $\hat{\Sigma}_{i,j}^\phi = (1/N_{i,j}) \sum_{n=1}^{N_{i,j}} (\phi(\mathbf{x}_{i,j}^n) - \hat{\boldsymbol{\mu}}_{i,j}^\phi) (\phi(\mathbf{x}_{i,j}^n) - \hat{\boldsymbol{\mu}}_{i,j}^\phi)^T$ ,  $\hat{\boldsymbol{\mu}}_{i,j}^\phi = (1/N_{i,j}) \sum_{n=1}^{N_{i,j}} \phi(\mathbf{x}_{i,j}^n)$  are the sample covariance matrix and the sample mean of  $(i, j)$  subclass in  $\mathcal{F}$  respectively. To avoid working with the mapped data explicitly (which may be impossible in case of infinite dimensional feature space  $\mathcal{F}$ ) a kernel function formulated as an inner product in the feature space satisfying the Mercer's condition

is used [9]

$$k(\mathbf{x}_{i,j}^n, \mathbf{x}_{k,l}^m) = \phi(\mathbf{x}_{i,j}^n)^T \phi(\mathbf{x}_{k,l}^m). \quad (46)$$

Under mild conditions, any solution of  $\mathbf{W}$  must lie in the span of all the training samples [9], and, thus, it can be represented by a linear combination of the training samples as

$$\mathbf{W} = \Phi(\mathbf{X}) \mathbf{\Gamma} \quad (47)$$

where  $\Phi(\mathbf{X}) = [\phi(\mathbf{x}_{1,1}^1), \dots, \phi(\mathbf{x}_{C,H_C}^{N_{C,H_C}})]$  and  $\mathbf{\Gamma} \in \mathbb{R}^{N \times C-1}$  contains the expansion coefficients. Substituting (47) into (45) and multiplying from the left with  $\Phi(\mathbf{X})^T$  we get  $\Phi^T(\mathbf{X}) \mathbf{S}_{bsb}^\phi \Phi(\mathbf{X}) \mathbf{\Gamma} = \Phi^T(\mathbf{X}) \tilde{\Sigma}_{\mathbf{X}}^\phi \Phi(\mathbf{X}) \mathbf{\Gamma} \mathbf{\Lambda}^\phi$  or

$$\mathbf{S}_{bsb}^k \mathbf{\Gamma} = \tilde{\Sigma}_{\mathbf{X}}^k \mathbf{\Gamma} \mathbf{\Lambda}^\phi \quad (48)$$

where we set  $\mathbf{S}_{bsb}^k = \Phi^T(\mathbf{X}) \mathbf{S}_{bsb}^\phi \Phi(\mathbf{X})$ ,  $\mathbf{S}_{ws}^k = \Phi^T(\mathbf{X}) \mathbf{S}_{ws}^\phi \Phi(\mathbf{X})$ , and  $\tilde{\Sigma}_{\mathbf{X}}^k = \mathbf{S}_{bsb}^k + \mathbf{S}_{ws}^k$ . The mean and sample covariance matrix of the  $(i, j)$  subclass in  $\mathcal{F}$  can be written in matrix product form as  $\boldsymbol{\mu}_{i,j}^\phi = \Phi(\mathbf{X}_{i,j}) \mathbf{p}_{i,j}$  and  $\Sigma_{i,j}^\phi = (1/N_{i,j}) \Phi(\mathbf{X}_{i,j}) (\mathbf{I} - \mathbf{P}_{i,j}) \Phi^T(\mathbf{X}_{i,j})$  respectively, where,  $\mathbf{p}_{i,j}$  is a  $N_{i,j} \times 1$  vector and  $\mathbf{P}_{i,j}$  is a  $N_{i,j} \times N_{i,j}$  matrix with all elements equal to  $1/N_{i,j}$ . Using the above expressions, the scatter matrices in (48) can be entirely expressed by the kernel functions as follows

$$\begin{aligned} \mathbf{S}_{bsb}^k &= \sum_{i=1}^{C-1} \sum_{j=1}^{H_i} \sum_{k=i+1}^C \sum_{l=1}^{H_k} p_{i,j} p_{k,l} (\mathbf{K}_{i,j} \mathbf{p}_{i,j} - \mathbf{K}_{k,l} \mathbf{p}_{k,l}) \\ &\quad \times (\mathbf{K}_{i,j} \mathbf{p}_{i,j} - \mathbf{K}_{k,l} \mathbf{p}_{k,l})^T, \end{aligned} \quad (49)$$

$$\mathbf{S}_{ws}^k = \frac{1}{N} \sum_{i=1}^C \sum_{j=1}^{H_i} \mathbf{K}_{i,j} (\mathbf{I} - \mathbf{P}_{i,j}) \mathbf{K}_{i,j}^T \quad (50)$$

where,  $\mathbf{K}_{i,j} = \Phi^T(\mathbf{X}) \Phi(\mathbf{X}_{i,j})$ ,  $\mathbf{K}_{i,j} \in \mathbb{R}^{N \times N_{i,j}}$ , and, thus,  $\mathbf{\Gamma}$  can be easily computed from (48) using only kernel evaluations. The derived  $\mathbf{\Gamma}$  can then be used for the projection of a test sample  $\phi(\mathbf{x})$  in the discriminant subspace using

$$\mathbf{z} = \mathbf{W}^T \phi(\mathbf{x}) = \mathbf{\Gamma}^T \mathbf{k} \quad (51)$$

where  $\mathbf{k} = [k(\mathbf{x}_{1,1}^1, \mathbf{x}), \dots, k(\mathbf{x}_{C,H_C}^{N_{C,H_C}}, \mathbf{x})]^T$  and  $\mathbf{z}$  is the projection of  $\phi(\mathbf{x})$ .

The optimal subclass partition of the data is identified by exploiting the nongaussianity-based iterative algorithm described in Algorithms 1 and 2. Consequently, the KMSDA algorithm is presented in Algorithm 3. In certain cases, KMSDA may provide superior performance in comparison to EM-MSDA and FMSDA, however, at the cost of much higher computation time during both the training and testing stage, especially when large-scale training data sets are used (due to the large number of kernel evaluations for mapping the observations in the kernel space, and the associated computational burden of performing eigenanalysis in this space).

## V. EXPERIMENTS

In this section, we use 12 standard benchmarks (defining in total 19 classification tasks) to compare the proposed algorithms, EM-MSDA, FMSDA and KMSDA, with various linear and nonlinear methods, in particular with PCA [41], LDA [6], FDA [30], MDA [16], SMDA [21], SDA [17], MSDA [26], KDA [9] and KSDA [33].

**Algorithm 3** KMSDA**Input:** Annotated data set  $X$ **Output:**  $\Gamma$ 

- 1: Initialize:  $H_1 = \dots = H_C = 1$ ,  $H = C$ ,  $\Phi_i$  (39),  $\Phi$  (41)
- 2: **repeat**
- 3:   Compute class label  $k$  of class to repartition (40)
- 4:   Set:  $H_k \leftarrow H_k + 1$
- 5:   Repartition  $k$ -th class to  $H_k$  subclasses using k-means
- 6:   Compute  $\Phi_i$  (39) and total nongaussianity  $\Phi$  (41)
- 7: **until** convergence of  $\Phi$
- 8: Compute  $\Gamma$  (48)

*A. Datasets*

For the evaluation we use four datasets that belong to the UCI repository [42], two datasets from the Gunnar Ratsch’s Benchmark Datasets [43], and six datasets that have been widely used for face, object and video shot detection:

Dataset 1: The Monk problem [42] is based on an artificial dataset of 432 data points in  $\mathbb{N}_+^6$ . Three binary classification tasks have been defined, i.e., MONK1, MONK2 and MONK3. For each task, a portion of the data has been randomly selected for forming the training set, and all 432 samples are used as the test set. In addition, in the third task 5% of the training data have been annotated wrongly in order to simulate the effect of random noise contaminating the data.

Dataset 2: The Landsat data set (LSD) consists of 6 classes (red soil, cotton crop, grey soil, damp grey soil, soil with vegetation stubble, and very damp grey soil) and 6435 feature vectors in  $\mathbb{N}_+^{36}$ . A partition of the dataset to training set (4435 samples) and test set (2000 samples) is already provided in [42].

Dataset 3: The Wisconsin diagnostic breast cancer (WDBC) dataset [42] is used for the recognition of benign and malignant cells from diagnostic images. This database comprises 569 diagnostic images represented in  $\mathbb{R}^{30}$ .

Dataset 4: The multi-feature digit dataset (MDD) [42] consists of ten classes and 200 patterns per class, i.e. 2,000 patterns in total, where each class represents one handwritten numeral (“0”-“9”). Each pattern is represented in terms of 6 feature sets, extracted from a  $30 \times 48$  binary image, as follows: a) MDD-pix: 240 pixel averages in  $2 \times 3$  windows, b) MDD-four: 76 Fourier coefficients of the character shapes, c) MDD-fac: 216 profile correlations, d) MDD-kar: 64 Karhunen-Loeve coefficients, e) MDD-zer: 47 Zernike moments, f) MDD-mor: 6 morphological features. Each set of features defines a separate classification task.

Dataset 5: The ETH80 database [44] consists of 8 object classes, namely, apples, pears, cars, cows, horses, dogs, tomatoes, and cups. Each class contains color images of 10 different objects recorded from 41 different views spaced evenly over the upper viewing hemisphere, i.e., the database contains 3280 images in total. In our computations the classic COIL segmentation masks of  $128 \times 128$  pixels size provided in [44] are employed, resized to  $25 \times 30$  pixels size and scanned column-wise to form 750-dimensional feature vectors.

Dataset 6: A subset of the MediaMill Challenge dataset

is used for event recognition experiments. It consists of 492 shots belonging to one of five different sport events (baseball, basketball, football, golf, soccer). Each shot is represented by a 101-dimensional vector, where the  $\kappa$ -th component of this vector is in the range  $[0; 1]$ , expressing the degree of confidence that the  $\kappa$ -th concept (out of 101 concepts) is present in the shot [45]. These values are the output of SVM-based automatic concept detectors, thus represent highly-noisy data.

Datasets 7-10: Four face datasets were used in our experiments. The Sheffield face database [46] offers 575 gray-scale cropped facial images of 20 individuals, shown in a range of poses from profile to frontal views. The AT&T Database of Faces [47] contains 400 facial images of 40 individuals captured at different times, with varying lighting conditions, facial expressions, etc. The Extended Yale B (ExtYaleB) database [48] offers 2432 gray-scale cropped facial images of 38 individuals under 64 illumination conditions. The CMU Pose, Illumination, and Expression (PIE) database [49] is a collection of more than 40,000 facial images of 68 people captured across 13 different poses, under 43 different illumination conditions, and with four different expressions. For the Sheffield database, we downsampled the facial images to size  $32 \times 32$  pixels resolution using bicubic interpolation, and scanned them columnwise to retrieve a set of 575 feature vectors in  $\mathbb{R}^{1024}$ . For the rest of the face databases we used the preprocessed  $32 \times 32$  pixels resolution facial image sets of the Four Face database collection [50], [51].

Dataset 11: The Banana set [52] is a binary class dataset consisting of 5300 samples in  $\mathbb{R}^2$ . It is an artificial dataset created using a mixture of overlapping Gaussians.

Dataset 12: The Breast Cancer dataset [52] is a two-class dataset containing samples of 277 patients in  $\mathbb{R}^9$  (excluding the nine samples that contain unknown attribute values).

*B. Evaluation*

A division of the datasets described in the previous subsection to training and test sets is necessary in order to evaluate the proposed algorithms. Such a division is provided along with the data for Monk and LSD. For Banana and Breast Cancer, we used 50 random realizations for training/test sets for each dataset from the Gunnar Ratsch’s benchmark collection [52]. Similarly, for AT&T, ExtYaleB and PIE, 30 random realizations from the Four Face database collection [50], [51] were used, where the training set at each realization contains 10 images per subject for ExtYaleB and PIE, and 8 images per subject for AT&T. For each of the remaining datasets, we divided them following standard practices in similar works of the literature, e.g. [17], [33]. In particular, we have designed  $\varsigma$  cross-validation (CV) folds by selecting randomly  $\varpi\%$  of the samples of each class at each fold to form the test set, and used the rest of the samples as the training set. The number of folds  $\varsigma$  and the percentage of test samples  $\varpi\%$  for WDBC, MDD, ETH-80, Sheffield, and Mediamill dataset were set to  $(\varsigma, \varpi) = (1, 50), (5, 50), (10, 10), (30, 60)$  and  $(30, 20)$  respectively.

The optimal parameters of each method at each CV fold are selected using as primary metric the correct classification rate



(CCR). For this, the global-to-local search strategy is applied (e.g., see [9]), i.e., after globally searching using a coarse scale of the parameter space, a candidate interval where the optimal parameters might exist is retrieved, and then a finer inspection for identifying the optimal parameters within this interval is performed. For the subclass methods (SDA, MSDA, FMSDA, EM-MSDA, KMSDA) we optimize over the number of subclasses in each class, and consequently over the total number of subclasses. For the FMSDA method we additionally require the identification of the exponent  $r$  of the weighting functions in (43) and the number of fractional steps  $\rho \in \mathbb{N}_+$  for decreasing the subspace dimensionality by one. For the optimization of these parameters we search over the following values:  $r = 3, 4, \dots, 16$  and  $\rho = 3, 4, \dots, 20$ . Similarly, for the kernel-based methods (KDA, KSDA, KMSDA) we need to identify the optimal parameters of the kernel functions. In our experiments we used two types of base kernels: Gaussian radial basis function  $k(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|/2\sigma^2)$ ,  $\sigma \in \mathbb{R}_+$ , and the polynomial function  $k(\mathbf{x}_i, \mathbf{x}_j) = ((\mathbf{x}_i \mathbf{x}_j) + o)^\varrho$ ,  $o \in \mathbb{R}$ ,  $\varrho \in \mathbb{N}_+$ . For their parameters we search for the optimal values over the following ranges:  $o = 0, 1$ ,  $\varrho = 1, 2, \dots, 8$ ,  $\sigma = 0.1, 0.2, \dots, 4$ . We should also note that for the datasets whose number of training observations  $N$  is small compared to their dimensionality  $F$  (such as the Sheffield and ETH-80 datasets), the computation of the inverse of the MLE of the sample covariance matrix (16) by the EM-based methods, for instance SMDA and EM-MSDA, will be especially problematic (e.g. see [2], [53]). In these cases, we compute the inverse using the eigenvalue decomposition of the sample covariance matrix, keeping only the eigenvalue components whose eigenvalues are above a specific threshold [2].

The recognition performance of a method regarding a dataset is measured using the average CCR (ACCR) along all CV folds, i.e., at each CV fold the maximum correct classification rate (CCR) for the different set of parameters is retained, and the CCRs are averaged along all CV folds. Similarly, the ground truth labels and the predicted labels at each CV fold for each algorithm are retained, and the McNemar's hypothesis test [54], [55] with a significance level of 0.025 is used to evaluate the statistical significance of the difference in the performance between each of the proposed algorithms and any other algorithm used in our experiments. Moreover, to compare the computational complexity of the algorithms we recorded the testing and training times in minutes, on a Intel i7 2.8GHz machine, with respect to one CV fold for each method and each dataset. Except for MDA and SMDA, for which their R package implementation [21] is exploited, all the other algorithms are compared using an unoptimized Matlab implementation. The FMSDA algorithm was then used as the baseline algorithm to compute the speedup rate  $s_\kappa$  for the  $\kappa$ -th algorithm using  $s_\kappa = T_{fmsda}/T_\kappa$ , where  $T_{fmsda}$  and  $T_\kappa$  are the training (or testing) time concerning the FMSDA and the  $\kappa$ -th algorithm respectively.

The ACCRs of the methods along with the average dimensionality in the discriminant subspace are shown in Table I, while, the results of the statistical significance tests are shown in Table II. In the latter, a cell contains the symbol  $+$ ,  $*$  or  $\sim$  for FMSDA, EMMSDA or KMSDA respectively, in order to

denote that the improvement in performance achieved by the aforementioned methods in comparison to the method corresponding to the column of the table is statistically significant. Finally, the speedup rate for the training stage (left side of the comma) and testing stage (right side of the comma) of the algorithms on each dataset are depicted in Table III, where higher speedup values indicate faster computations. In every table we have divided the methods into three groups, namely, linear, subclass and kernel-based methods. With respect to this partitioning, for Tables I and III we have used bold digits and underlined-bold digits to denote the best performance rate within each group and along all methods respectively.

From Table I we can see that for the majority of the datasets the best ACCR among the linear subclass methods is provided by FMSDA (in 10 out of 19 classification tasks of Table I) or EM-MSDA (again in 10 out of 19 tasks). In overall the best ACCR among all methods is achieved by KMSDA (in 17 out of 19 tasks). We should also note that in many cases FMSDA and EM-MSDA outperform the kernel-based methods as well (including KMSDA in 2 classification tasks, while they match KMSDA's performance in another 2 tasks). Between FMSDA and EM-MSDA, we observe that the former tends to perform better when the data dimensionality is larger than the number of the samples, and at the same time many subclasses are necessary in order to capture the subclass structure of the data. In these cases, the training samples per subclass are limited and consequently the subclass covariance matrices are poorly estimated [53]. This adversely affects the performance of EM-based methods. For instance, the performance of SMDA and EM-MSDA on the ETH80 dataset (which contains 8 object classes and each object class 10 different objects) is considerably lower than that of FMSDA.

From the results in Table II we can also see that the performance improvements attained by the proposed methods are statistically significant for most of the datasets. From Table III, we additionally see that FMSDA and EM-MSDA provide systematically lower computation times during the testing stage among all linear subclass methods (and, as expected are also faster than the kernel-based methods). This quality of FMSDA and EM-MSDA is a critical advantage of them, especially for applications that require real-time or near real-time processing of large data volumes, such as event detection in video streams. Summarizing, we observe that the three proposed methods in most cases outperform the current state of the art as recently reported for KSDA in [33] and in the also very recent works [21], [26], at the same time offering competitive response times during the testing (recognition) stage.

## VI. CONCLUSIONS

Subclass DA methods are attractive alternatives to the kernel DA variants because they offer fast (often real-time) computations and comparable recognition performance. Furthermore, combining subclass partitioning and the kernel trick in a single DA method opens new possibilities for improved DA effectiveness. MSDA is a very recent subclass method, that utilizes an effective partitioning procedure to derive a



## ACKNOWLEDGMENT

This work was supported by the EC under contracts FP7-248984 GLOCAL and FP7-287911 LinkedTV.

## APPENDIX A

## DERIVATION OF EQUATIONS IN SECTION II

## A. Derivation of Eqs. (6) and (7)

The Gaussian mixture distribution concerning the  $i$ -th class in (5) can be derived in terms of latent variables [36], [37], as described in the following. Let  $Z_i \in \mathbb{R}^{H_i}$  be a categorical latent random vector concerning the  $i$ -th class, whose parameter space  $\mathcal{Z}_i$  is the standard base of  $\mathbb{R}^{H_i}$ , i.e.,  $\mathcal{Z}_i = \{\mathbf{e}_{i,1}, \dots, \mathbf{e}_{i,H_i}\}$ , where only the  $j$ -th element of the unit vector  $\mathbf{e}_{i,j}$  is equal to one and all other elements are equal to zero. Setting  $p(Z_i = \mathbf{e}_{i,j}) = \pi_{i,j}$  and  $p(\mathbf{x}|Z_i = \mathbf{e}_{i,j}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_{i,j})$  the marginal and conditional densities,  $p(\mathbf{z}_i)$  and  $p(\mathbf{x}|\mathbf{z}_i)$ , are expressed in terms of the mixing coefficients and mixture components respectively,  $p(\mathbf{z}_i) = \prod_{j=1}^{H_i} \pi_{i,j}^{z_{i,j}}$ ,  $p(\mathbf{x}|\mathbf{z}_i) = \prod_{j=1}^{H_i} \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_{i,j})^{z_{i,j}}$ . Thus, using the product rule of probability we can express the  $i$ -th class-conditional joint density as

$$\begin{aligned} p(\mathbf{x}, \mathbf{z}_i|\omega_i) &= p(\mathbf{z}_i|\omega_i)p(\mathbf{x}|\mathbf{z}_i, \omega_i) = p(\mathbf{z}_i)p(\mathbf{x}|\mathbf{z}_i) \\ &= \prod_{j=1}^{H_i} (\pi_{i,j} \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_{i,j}))^{z_{i,j}}, \end{aligned} \quad (52)$$

where we have used the fact that  $\mathbf{x}$  is conditionally independent of  $\omega_i$  given  $\mathbf{z}_i$ , and  $\mathbf{z}_i$  is independent of  $\omega_i$ . The  $i$ -th class-conditional marginal distribution of  $\mathbf{x}$  can then be written as

$$p(\mathbf{x}|\omega_i) = \sum_{\mathbf{z}_i} p(\mathbf{x}, \mathbf{z}_i|\omega_i) = \sum_{j=1}^{H_i} \pi_{i,j} \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_{i,j}), \quad (53)$$

which is a Gaussian mixture equivalent to (5), and, using the Bayes' rule the posterior distribution is also derived

$$p(\mathbf{z}_i|\mathbf{x}, \omega_i) = \frac{\prod_{j=1}^{H_i} (\pi_{i,j} \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_{i,j}))^{z_{i,j}}}{\sum_{j=1}^{H_i} \pi_{i,j} \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_{i,j})}. \quad (54)$$

Therefore, under the i.i.d. assumption, the likelihood of the complete data set is expressed as (p.108, [27])

$$\begin{aligned} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) &= \prod_{i=1}^C \prod_{n=1}^{N_i} p(\mathbf{x}_i^n, \mathbf{z}_i^n|\omega_i) \\ &= \prod_{i=1}^C \prod_{n=1}^{N_i} \prod_{j=1}^{H_i} (\pi_{i,j} \mathcal{N}(\mathbf{x}_i^n|\boldsymbol{\mu}_{i,j}))^{z_{i,j}^n}. \end{aligned} \quad (55)$$

while the posterior distribution takes the form

$$p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}) \propto \prod_{i=1}^C \prod_{n=1}^{N_i} \prod_{j=1}^{H_i} (\pi_{i,j} \mathcal{N}(\mathbf{x}_i^n|\boldsymbol{\mu}_{i,j}))^{z_{i,j}^n}, \quad (56)$$

where  $\mathbf{Z} = \{\mathbf{Z}_1, \dots, \mathbf{Z}_C\}$  is the set of all categorical vectors. Observing that the posterior distribution is independent over  $z_{i,j}^n$ , the expectation of the categorical variables can be derived

$$\mathbb{E}[z_{i,j}^n] = \frac{\sum_{j=1}^{H_i} z_{i,j}^n (\pi_{i,j} \mathcal{N}(\mathbf{x}_i^n|\boldsymbol{\mu}_{i,j}))^{z_{i,j}^n}}{\sum_{j=1}^{H_i} \pi_{i,j} \mathcal{N}(\mathbf{x}_i^n|\boldsymbol{\mu}_{i,j})}, \quad (57)$$

and simplifying the above, we arrive to the definition of the responsibilities in (7).

Moreover, from (56) the log likelihood of the complete data set is retrieved

$$\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) = \sum_{i=1}^C \sum_{n=1}^{N_i} \sum_{j=1}^{H_i} z_{i,j}^n (\ln \pi_{i,j} + \ln \mathcal{N}(\mathbf{x}_i^n|\boldsymbol{\mu}_{i,j})). \quad (58)$$

Applying the expectation operator to the above expression and substituting  $\mathbb{E}[z_{i,j,n}]$  from (7) the expectation of the complete data log-likelihood is expressed as

$$\begin{aligned} \mathbb{E}[\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})] &= \sum_{i=1}^C \sum_{n=1}^{N_i} \sum_{j=1}^{H_i} h_{i,j}^n (\ln \pi_{i,j} + \ln \mathcal{N}(\mathbf{x}_{i,n}|\boldsymbol{\mu}_{i,j}, \boldsymbol{\Sigma})) \\ &= \sum_{i=1}^C \sum_{j=1}^{H_i} \tilde{N}_{i,j} \ln \pi_{i,j} - \frac{NF}{2} \ln(2\pi) + \frac{N}{2} \ln |\boldsymbol{\Sigma}^{-1}| \\ &\quad - \frac{1}{2} \sum_{i=1}^C \sum_{n=1}^{N_i} \sum_{j=1}^{H_i} h_{i,j,n} (\mathbf{x}_{i,n} - \boldsymbol{\mu}_{i,j})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_{i,n} - \boldsymbol{\mu}_{i,j}). \end{aligned} \quad (59)$$

Using the identity  $(\mathbf{x}_i^n - \boldsymbol{\mu}_{i,j})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i^n - \boldsymbol{\mu}_{i,j}) = (\mathbf{x}_i^n - \bar{\mathbf{x}}_{i,j}^n)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i^n - \bar{\mathbf{x}}_{i,j}^n) + (\bar{\mathbf{x}}_{i,j} - \boldsymbol{\mu}_{i,j})^T \boldsymbol{\Sigma}^{-1} (\bar{\mathbf{x}}_{i,j} - \boldsymbol{\mu}_{i,j}) + 2(\mathbf{x}_i^n - \bar{\mathbf{x}}_{i,j}^n)^T \boldsymbol{\Sigma}^{-1} (\bar{\mathbf{x}}_{i,j} - \boldsymbol{\mu}_{i,j})$  along with the fact that  $\sum_{n=1}^{N_i} (\mathbf{x}_i^n - \bar{\mathbf{x}}_{i,j}^n)^T \boldsymbol{\Sigma}^{-1} (\bar{\mathbf{x}}_{i,j} - \boldsymbol{\mu}_{i,j}) = 0$ , and multiplying both sides by two, we arrive to (6).

## B. Derivation of Eq. (18)

The constraint that the mixing coefficients should sum to one can be incorporated in (17) using  $C$  lagrange multipliers  $\eta_i, i = 1, \dots, C$ . Therefore, we need to find the stationary point of

$$\begin{aligned} &\sum_{i=1}^C \sum_{n=1}^{N_i} \sum_{j=1}^{H_i} h_{i,j}^n (\ln \pi_{i,j} + \ln \mathcal{N}(\mathbf{x}_i^n|\boldsymbol{\mu}_{i,j})) \\ &\quad + \sum_{i=1}^C \eta_i (\sum_{j=1}^{H_i} \pi_{i,j} - 1) \end{aligned} \quad (60)$$

with respect to  $\pi_{i,j}$  and  $\eta_i$ . Optimizing over  $\pi_{i,j}$  we arrive to  $\tilde{N}_{i,j}/\pi_{i,j} + \eta_i = 0$ . If we multiply both sides with  $\pi_{i,j}$  and sum over all subclasses of the  $i$ -th class we get  $\eta_i = -N_i$ . Eliminating  $\eta_i$  we obtain (18).

## REFERENCES

- [1] Y. Ephraim and H. L. V. Trees, "A signal subspace approach for speech enhancement," *IEEE Trans. Speech Audio Process.*, vol. 3, no. 4, pp. 251–266, Jul. 1995.
- [2] K. Fukunaga, *Introduction to statistical pattern recognition (2nd ed.)*. San Diego, CA, USA: Academic Press Professional, Inc., 1990.
- [3] R. Duda, P. Hart, and D. Stork, *Pattern Classification, (2nd ed.)*. New York, USA: John Wiley & Sons, Inc., 2001.
- [4] S. Ji and J. Ye, "Generalized linear discriminant analysis: A unified framework and efficient model selection," *IEEE Trans. Neural Netw.*, vol. 19, no. 10, pp. 1768–1782, Oct. 2008.
- [5] C. B. Moler and G. W. Stewart, "An algorithm for generalized matrix eigenvalue problems," *SIAM Journal on Numerical Analysis*, vol. 10, no. 2, pp. 241–256, Apr. 1973.
- [6] P. N. Belhumeur, J. P. Hespanha, and D. J. Kriegman, "Eigenfaces vs. Fisherfaces: Recognition using class specific linear projection," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 19, no. 7, pp. 711–720, Jul. 1997.

- [7] H. Li, T. Jiang, and K. Zhang, "Efficient and robust feature extraction by maximum margin criterion," *IEEE Trans. Neural Netw.*, vol. 17, no. 1, pp. 157–165, Oct. 2006.
- [8] C. S. Dhir and S.-Y. Lee, "Discriminant independent component analysis," *IEEE Trans. Neural Netw.*, vol. 22, no. 6, pp. 845–857, Jun. 2011.
- [9] K. Muller, S. Mika, G. Ratsch, S. Tsuda, and B. Scholkopf, "An introduction to kernel-based learning algorithms," *IEEE Trans. Neural Netw.*, vol. 12, no. 2, pp. 181–202, Mar. 2001.
- [10] L. Wang, K. L. Chan, P. Xue, and L. Zhou, "A kernel-induced space selection approach to model selection in KLDA," *IEEE Trans. Neural Netw.*, vol. 19, no. 12, pp. 2116–2131, Dec. 2008.
- [11] Z. Li, D. Lin, and X. Tang, "Nonparametric discriminant analysis for face recognition," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 31, no. 4, pp. 755–761, Apr. 2009.
- [12] J. Lu, K. N. Plataniotis, and A. N. Venetsanopoulos, "Face recognition using kernel direct discriminant analysis algorithms," *IEEE Trans. Neural Netw.*, vol. 14, no. 1, pp. 117–126, Jan. 2003.
- [13] S. Zafeiriou, G. Tzimiropoulos, M. Petrou, and T. Stathaki, "Regularized kernel discriminant analysis with a robust kernel for face recognition and verification," *IEEE Trans. Neural Netw. and Learning Syst.*, vol. 23, no. 3, pp. 526–534, Mar. 2012.
- [14] Z. Fan, Y. Xu, and D. Zhang, "Local linear discriminant analysis framework using sample neighbors," *IEEE Trans. Neural Netw.*, vol. 22, no. 7, pp. 1119–1132, Jul. 2011.
- [15] B.-C. Kuo and K.-Y. Chang, "Feature extractions for small sample size classification problem," *IEEE Trans. Geosci. Remote Sens.*, vol. 45, no. 3, pp. 756–764, Mar. 2007.
- [16] T. Hastie and R. Tibshirani, "Discriminant analysis by Gaussian mixtures," *Journal of the Royal Statistical Society. Series B*, vol. 58, no. 1, pp. 155–176, Jul. 1996.
- [17] M. Zhu and A. Martinez, "Subclass discriminant analysis," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 28, no. 8, pp. 1274–1286, Aug. 2006.
- [18] S.-W. Kim and R. P. W. Duin, "On using a pre-clustering technique to optimize LDA-based classifiers for appearance-based face recognition," in *Proc. 12th Iberoamerican Congress on Pattern Recognition*, Vina del Mar-Valparaiso, Chile, Nov. 2007, pp. 466–476.
- [19] M. H. Yang, D. J. Kriegman, and N. Ahuja, "Face detection using multimodal density models," *Computer Vision and Image Understanding*, vol. 84, no. 2, pp. 264–284, Oct. 2001.
- [20] A. Pnevmatikakis and L. Polymenakos, "Subclass linear discriminant analysis for video-based face recognition," *Journal of Visual Communication and Image Representation*, vol. 20, no. 8, pp. 543–551, Nov. 2009.
- [21] L. Clemmensen, T. Hastie, D. Witten, and B. Ersboll, "Sparse discriminant analysis," *Technometrics*, vol. 53, no. 4, pp. 406–413, Nov. 2011.
- [22] D. Wu and K. L. Boyer, "Resilient subclass discriminant analysis," in *Proc. IEEE Int. Conf. on Computer Vision (ICCV 2009)*, Kyoto, Japan, Sep./Oct. 2009, pp. 389–396.
- [23] F. Oveis, "Subclass discriminant analysis using dynamic cluster formation for EEG-based brain-computer interface," in *Proc. IEEE/EMBS 4th Int. Conf. on Neural Engineering*, Antalya, Turkey, May 2009, pp. 303–306.
- [24] S.-W. Kim, "A pre-clustering technique for optimizing subclass discriminant analysis," *Pattern Recogn. Lett.*, vol. 31, no. 6, pp. 462–468, Apr. 2010.
- [25] N. Gkalelis, V. Mezaris, and I. Kompatsiaris, "High-level event detection in video exploiting discriminant concepts," in *Proc. 9th International Workshop on Content-Based Multimedia Indexing (CBMI 2011)*, Madrid, Spain, Jun. 2011, pp. 85–90.
- [26] —, "Mixture subclass discriminant analysis," *IEEE Signal Process. Lett.*, vol. 18, no. 5, pp. 319–332, May 2011.
- [27] K. V. Mardia, J. T. Kent, and J. M. Bibby, *Multivariate analysis*. Academic Press, 1979.
- [28] N. A. Campbell, "Canonical variate analysis - A general model formulation," *Australian & New Zealand Journal of Statistics*, vol. 26, no. 1, pp. 86–96, 1984.
- [29] D. Tao, X. Li, X. Wu, and S. J. Maybank, "Geometric mean for subspace selection," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 31, no. 2, pp. 260–274, Feb. 2009.
- [30] R. Lotlikar and R. Kothari, "Fractional-step dimensionality reduction," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 22, no. 6, pp. 623–627, Jun. 2000.
- [31] M. Loog, R. P. W. Duin, and R. Haeb-Umbach, "Multiclass linear dimension reduction by weighted pairwise Fisher criteria," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 23, no. 7, pp. 762–766, Jul. 2001.
- [32] B. Chen, L. Yuan, H. Liu, and Z. Bao, "Kernel subclass discriminant analysis," *Neurocomputing*, vol. 71, no. 1–3, pp. 455–458, Dec. 2007.
- [33] D. You, O. C. Hamsici, and A. M. Martinez, "Kernel optimization in discriminant analysis," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 33, no. 3, pp. 631–638, Mar. 2011.
- [34] V. Vapnik, *Statistical learning theory*. New York: Wiley, 1998.
- [35] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the EM algorithm," *Journal of the Royal Statistical Society. Series B (Methodological)*, vol. 39, no. 1, pp. 1–38, 1977.
- [36] G. McLachlan and D. Peel, *Finite Mixture Models*. New York: Wiley-Interscience, 2000.
- [37] C. M. Bishop, *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Secaucus, NJ, USA: Springer-Verlag New York, Inc., 2006.
- [38] L. Devroye, L. Györfi, and G. Lugosi, *A Probabilistic Theory of Pattern Recognition*. New York, USA: Springer, 1996.
- [39] N. Vlassis and A. Likas, "A kurtosis-based dynamic approach to gaussian mixture modeling," *IEEE Trans. Syst., Man, Cybern. A*, vol. 29, no. 4, pp. 393–399, Jul. 1999.
- [40] L. Wang and J. Ma, "A kurtosis and skewness based criterion for model selection on gaussian mixture," in *2nd Int. Conf. on BioMedical Engineering and Informatics*, Tianjin, China, Oct. 2009, pp. 1–5.
- [41] I. T. Jolliffe, *Principal Component Analysis*. New York, USA: Springer, Oct. 2002.
- [42] A. Frank and A. Asuncion, "UCI machine learning repository," 2010. [Online]. Available: <http://archive.ics.uci.edu/ml>
- [43] "Gunnar raetsch's benchmark datasets," <http://theoval.cmp.uea.ac.uk/~gcc/matlab/default.html#benchmarks>, accessed 2012-05-01.
- [44] B. Leibe and B. Schiele, "Analyzing appearance and contour based methods for object categorization," in *Proc. 2003 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, vol. 2, Madison, WI, USA, Jun. 2003, pp. II–409–15.
- [45] N. Gkalelis, V. Mezaris, and I. Kompatsiaris, "Automatic event-based indexing of multimedia content using a joint content-event model," in *ACM Multimedia 2010 (EiMM10)*, Firenze, Italy, Oct. 2010.
- [46] D. B. Graham and N. M. Allinson, "Characterizing virtual eigensignatures for general purpose face recognition," in *Face Recognition: From Theory to Applications, Computer and Systems Sciences*, H. Wechsler et al., Ed. NATO ASI Series F, 1998, vol. 163, pp. 446–456.
- [47] F. Samaria and A. Harter, "Parameterisation of a stochastic model for human face identification," in *Proc. 2nd IEEE Workshop on Applications of Computer Vision*, Sarasota FL, USA, Dec. 1994, pp. 138–142.
- [48] K. Lee, J. Ho, and D. Kriegman, "Acquiring linear subspaces for face recognition under variable lighting," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 27, no. 5, pp. 684–698, May 2005.
- [49] T. Sim, S. Baker, and M. Bsat, "The CMU pose, illumination, and expression database," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 25, no. 12, pp. 1615–1618, Dec. 2003.
- [50] D. Cai, X. He, Y. Hu, J. Han, and T. S. Huang, "Learning a spatially smooth subspace for face recognition," in *IEEE Conf. Comput. Vis. and Pattern Recognit.*, Minneapolis, Minnesota, USA, Jun. 2007, pp. 138–142.
- [51] "Four face databases in matlab format," <http://http://www.cad.zju.edu.cn/home/dengcai/Data/FaceData.html>, accessed 2012-05-01.
- [52] G. Rätsch, T. Onoda, and K.-R. Müller, "Soft margins for adaboost," *Mach. Learn.*, vol. 42, no. 3, pp. 287–320, Mar. 2001.
- [53] C. E. Thomaz, D. F. Gillies, and R. Q. Feitosa, "A new covariance estimate for bayesian classifiers in biometric recognition," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 14, no. 2, pp. 214–223, Feb. 2004.
- [54] Q. McNemar, "Note on the sampling error of the difference between correlated proportions or percentages," *Psychometrika*, vol. 12, no. 2, pp. 153–157, Jun. 1947.
- [55] S. Zafeiriou, G. Tzimiropoulos, M. Petrou, and T. Stathaki, "Regularized kernel discriminant analysis with a robust kernel for face recognition and verification," *IEEE Trans. Neural Netw.*, vol. 23, no. 3, pp. 526–534, Mar. 2012.
- [56] Y. Aksu, D. J. Miller, G. Kesidis, and Q. X. Yang, "Margin-maximizing feature elimination methods for linear and nonlinear kernel-based discriminant functions," *IEEE Trans. Neural Netw.*, vol. 21, no. 5, pp. 701–717, May 2010.
- [57] F. Song, D. Mei, and H. Li, "Feature selection based on linear discriminant analysis," in *IEEE Int. Conf. Intell. Syst. Design and Eng. Appl.*, vol. 1, Changsha, China, Oct. 2010, pp. 746–749.
- [58] S. Huh and D. Lee, "Linear discriminant analysis for signatures," *IEEE Trans. Neural Netw.*, vol. 21, no. 12, pp. 1990–1996, Dec. 2010.