

# Linear Dimension Reduction for Multiple Heteroscedastic Multivariate Normal Populations

Songthip T. Ounpraseuth<sup>1</sup>, Phil D. Young<sup>2</sup>, Johanna S. van Zyl<sup>2</sup>,  
Tyler W. Nelson<sup>2</sup>, Dean M. Young<sup>2</sup>

<sup>1</sup>Department of Biostatistics, University of Arkansas for Medical Sciences, Little Rock, AK, USA

<sup>2</sup>Department of Statistical Science, Baylor University, Waco, TX, USA

Email: [stounpraseuth@uams.edu](mailto:stounpraseuth@uams.edu), [philip\\_young@baylor.edu](mailto:philip_young@baylor.edu), [tyler\\_nelson@baylor.edu](mailto:tyler_nelson@baylor.edu),  
[johanna\\_vanzyl@baylor.edu](mailto:johanna_vanzyl@baylor.edu), [dean\\_young@baylor.edu](mailto:dean_young@baylor.edu)

Received 16 April 2015; accepted 19 June 2015; published 24 June 2015

Copyright © 2015 by authors and Scientific Research Publishing Inc.

This work is licensed under the Creative Commons Attribution International License (CC BY).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

---

## Abstract

For the case where all multivariate normal parameters are known, we derive a new linear dimension reduction (*LDR*) method to determine a low-dimensional subspace that preserves or nearly preserves the original feature-space separation of the individual populations and the Bayes probability of misclassification. We also give necessary and sufficient conditions which provide the smallest reduced dimension that essentially retains the Bayes probability of misclassification from the original full-dimensional space in the reduced space. Moreover, our new *LDR* procedure requires no computationally expensive optimization procedure. Finally, for the case where parameters are unknown, we devise a *LDR* method based on our new theorem and compare our *LDR* method with three competing *LDR* methods using Monte Carlo simulations and a parametric bootstrap based on real data.

## Keywords

Linear Transformation, Bayes Classification, Feature Extraction, Probability of Misclassification

---

## 1. Introduction

The fact that the Bayes probability of misclassification (*BPMC*) of a statistical classification rule does not increase as the dimension or feature space increases, provided the class-conditional probability densities are known, is well-known. However, in practice when parameters are estimated and the feature-space dimension is

large relative to the training-sample sizes, the performance or efficacy of a sample discriminant rule may be considerably degraded. This phenomenon gives rise to a paradoxical behavior that [1] has called the *curse of dimensionality*.

An exact relationship between the expected probability of misclassification (*EPMC*), training-sample sizes, feature-space dimension, and actual parameters of the class-conditional densities is challenging to obtain. In general, as the classifier becomes more complex, the ratio of sample size to dimensionality must increase at an exponential rate to avoid the curse of dimensionality. The authors [2] have suggested a ratio of at least ten times as many training samples per class as the feature dimension increases. Hence, as the number of feature variables  $p$  becomes large relative to the training-sample sizes  $n_i$ ,  $i=1, \dots, m$ , where  $m$  is the number of classes, one might wish to use a smaller number of the feature variables to improve the classifier performance or computational efficiency. This approach is called feature subset selection.

Another effective approach to obtain a reduced dimension to avoid the curse of dimensionality is linear dimension reduction (*LDR*). Perhaps the most well-known *LDR* procedure for the  $m$ -class problem is linear discriminant analysis (*LDA*) from [3], which is a generalization of the linear discriminant function (*LDF*) derived in [4] for the case  $m=2$ . The *LDA LDR* method determines a vector  $\mathbf{a}$  that maximizes the ratio of between-class scatter to average within-class scatter in the lower-dimensional space. The sample within-class scatter matrix is

$$\mathbf{S}_W \equiv \sum_{i=1}^m \alpha_i \mathbf{S}_i, \quad (1)$$

where  $\mathbf{S}_i$  is the estimated sample covariance matrix for the  $i$ th class, and the sample between-class scatter matrix is

$$\mathbf{S}_B \equiv \sum_{i=1}^m \alpha_i (\bar{\mathbf{x}}_i - \bar{\mathbf{x}})(\bar{\mathbf{x}}_i - \bar{\mathbf{x}})', \quad (2)$$

where  $\bar{\mathbf{x}}_i$  is the sample mean vector for class  $\Pi_i$ ,  $\alpha_i$  is its *a priori* probability of class membership,  $\bar{\mathbf{x}} = \sum_{i=1}^m \alpha_i \bar{\mathbf{x}}_i$  is the estimated overall mean, and  $i=1, \dots, m$ . For  $m=2$  with  $\alpha_1 = \alpha_2$ , [4] determined the vector  $\mathbf{a}$  that maximizes the criterion function

$$\mathbf{J}_F(\mathbf{a}) = \frac{\mathbf{a}' \mathbf{S}_B \mathbf{a}}{(\mathbf{a}' \mathbf{S}_W \mathbf{a})},$$

which is achieved by an eigenvalue decomposition of  $\mathbf{S}_W^{-1} \mathbf{S}_B$ . An attractive feature of *LDA* as a *LDR* method is that it is computationally appealing; however, it attempts to maximally separate class means and does not incorporate the discriminatory information contained in the differences of the class covariance matrices. Many alternative approaches to the *LDA LDR* method have been proposed. For example, canonical variables can be viewed as an extension to the *LDF* when  $m > 2$  and  $q > 1$  (see [5]). Other straightforward extensions of the *LDF* include those in [6] [7].

Extensions of *LDA* that incorporate information on the differences in covariance matrices are known as *heteroscedastic linear dimension reduction (HLDR)* methods. The authors [8] have proposed an eigenvalue-based *HLDR* approach for  $m > 2$ , utilizing the so-called *Chernoff criterion*, and have extended the well-known *LDA* method using directed distance matrices that can be considered a generalization of (2). Additional *HLDR* methods have been proposed by authors such as [9]-[10].

Using results by [13] that characterize linear sufficient statistics for multivariate normal distributions, we develop an explicit *LDR* matrix  $\mathbf{B} \in \mathbb{R}_{q \times p}$  such that

$$\mathbf{x} \rightarrow \mathbf{y} = \mathbf{B}\mathbf{x},$$

where  $\mathbf{x} \in \mathbb{R}_{p \times 1}$ ,  $\mathbf{y} \in \mathbb{R}_{q \times 1}$ , and  $\mathbb{R}_{m \times n}$  denotes the space of all  $m \times n$  real matrices. Using the Bayes classification procedure in which we assume equal costs of misclassification and that all class parameters are known, we determine the reduced dimension  $q < p$  that is the smallest reduced dimension for which there exists a *LDR* matrix  $\mathbf{B} \in \mathbb{R}_{q \times p}$  that preserves all of the classification information originally contained in the  $p$ -dimensional feature space. We then derive a linear transformation that assigns  $\mathbf{x}$  to  $\Pi_k$  if and only if the corres-

ponding Bayes classification rule assigns  $\mathbf{B}\mathbf{x}$  to  $\Pi_k$ , where  $k \in \{1, \dots, m\}$ . We refer to this method as the *SY LDR* procedure.

Moreover, we use Monte Carlo simulations to compare the classification efficacy of the *BE* method, sliced inverse regression (*SIR*), and sliced average variance estimation (*SAVE*) found in [14]-[16], respectively, with the *SY* method.

The remainder of this paper is organized as follows. We begin with a brief introduction to the Bayes quadratic classifier in Section 2 and introduce some preliminary results that we use to prove our new *LDR* method in Section 3. In Section 4, we provide conditions under which the Bayes quadratic classification rule is preserved in the low-dimensional space and derive a new *LDR* matrix. We establish a *SVD*-based approximation to our *LDR* procedure along with an example of low-dimensional graphical representations in Section 5. We describe the four *LDR* methods that we compare using Monte Carlo simulations in Section 6. We present five Monte Carlo simulations in which we compare the competing *LDR* procedures for various population parameter configurations in Section 7. In Section 8, we compare the four methods using bootstrap simulations for a real data example, and, finally, we offer a few concluding remarks in Section 9.

## 2. The Bayes Quadratic Classifier

The Bayesian statistical classifier discriminates based on the probability density functions  $p(\Pi_i|\mathbf{x})$ ,  $i=1, \dots, m$ , of each class. The Bayes classifier is optimal in the sense that it maximizes the class *a posteriori* probability provided all class distributions and corresponding parameters are known. That is, suppose we have  $m$  classes,  $\Pi_1, \dots, \Pi_m$ , with assumed known *a priori* probabilities  $\alpha_1, \dots, \alpha_m$ , respectively. Also, let  $p(\cdot|\Pi_i)$  denote the  $p$ -dimensional multivariate normal density corresponding to population  $\Pi_i$ ,  $i=1, \dots, m$ . The goal of statistical decision theory is to obtain a decision rule that assigns an unlabeled observation  $\mathbf{x}$  to  $\Pi_k$  if  $p(\Pi_k|\mathbf{x})$  is the maximum overall *a posteriori*  $p(\Pi_i|\mathbf{x})$ ,  $i=1, \dots, m$ . Then,

$$\lambda(\Pi_i|\Pi_j) = \begin{cases} 0, & i = j \\ 1, & i \neq j \end{cases}$$

the Bayes classifier assigns  $\mathbf{x}$  to class  $\Pi_k$  if

$$p(\Pi_k|\mathbf{x}) > p(\Pi_j|\mathbf{x}), \quad j=1, \dots, m; j \neq k.$$

This decision rule partitions the measurement or feature space into  $m$  disjoint regions  $R_{\Pi_i}$ , where  $i=1, \dots, m$ , such that  $\mathbf{x}$  is assigned to class  $\Pi_k$  if  $\mathbf{x} \in R_{\Pi_k}$ . Using Bayes' rule, the *a posteriori* probabilities of class membership  $p(\Pi_k|\mathbf{x})$  can be defined as

$$p(\Pi_k|\mathbf{x}) = \frac{\alpha_k p(\mathbf{x}|\Pi_k)}{p(\mathbf{x})}.$$

One can re-express the Bayes classification as the following:

Assign  $\mathbf{x}$  to  $\Pi_k$  if

$$\alpha_k p(\mathbf{x}|\Pi_k) > \alpha_j p(\mathbf{x}|\Pi_j), \quad j=1, \dots, m; j \neq k. \quad (3)$$

This decision rule is known as the *Bayes' classification rule*. Let  $\Pi_i$  be modeled as a  $p$ -dimensional multivariate normal distribution, and let

$$d_i(\mathbf{x}) \equiv \ln|\Sigma_i| - 2[\ln(\alpha_i)] + (\mathbf{x} - \boldsymbol{\mu}_i)' \Sigma_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i), \quad i=1, \dots, m. \quad (4)$$

The Bayes decision rule (4) is to classify the unlabeled observation  $\mathbf{x}$  into the class  $\Pi_k$  such that  $d_k(\mathbf{x}) = \min\{d_i(\mathbf{x}); i=1, \dots, m\}$ . The classification rule defined by (4) is known as the *quadratic discriminant function (QDF)*, or the *quadratic classifier*.

## 3. Preliminary Results

The following notation will be used throughout the remainder of the paper. We let  $\mathbb{R}_p^>$  denote the set of  $p \times p$

positive definite matrices and  $\mathbb{R}_{p \times p}^S$  denote the set of  $p \times p$  symmetric matrices. Also, we let  $\mathbf{A}^+ \in \mathbb{R}_{m \times n}$  represent the Moore-Penrose pseudo-inverse of  $\mathbf{A} \in \mathbb{R}_{n \times m}$ .

The proof of the main result for the derivation of our new *LDR* method requires the following notation and lemmas. Let  $\mathbf{M} \in \mathbb{R}_{p \times (m-1)(p+1)}$  be

$$\mathbf{T} \equiv [\mathbf{d}_2 - \mathbf{d}_1 \mid \cdots \mid \mathbf{d}_m - \mathbf{d}_1 \mid \mathbf{E}_2 - \mathbf{E}_1 \mid \cdots \mid \mathbf{E}_m - \mathbf{E}_1], \quad (5)$$

where  $\mathbf{d}_i \in \mathbb{R}_{p \times 1}$ ,  $\mathbf{E}_i \in \mathbb{R}_{p \times p}^S$  such that  $\text{rank}(\mathbf{E}_i) = p$ , and  $\mathbf{d}_1 \neq \mathbf{d}_k$  and  $\mathbf{E}_1 \neq \mathbf{E}_k$  for at least one value of  $k$ , where  $2 \leq k \leq m$  and  $i = 1, \dots, m$ . Also, let  $\text{rank}(\mathbf{T}) = 1 \leq q < p$ , and let  $\mathbf{F} \in \mathbb{R}_{p \times q}$  and  $\mathbf{G} \in \mathbb{R}_{q \times (m-1)(p+1)}$  be matrix components of a full-rank decomposition of  $\mathbf{T}$  so that  $\mathbf{T} = \mathbf{F}\mathbf{G}$  with  $\text{rank}(\mathbf{F}) = \text{rank}(\mathbf{G}) = q$ . Then, the Moore-Penrose pseudoinverse of  $\mathbf{T}$  is  $\mathbf{T}^+ = \mathbf{G}^+\mathbf{F}^+$ , and, also,  $\mathbf{T}\mathbf{T}^+ = \mathbf{F}\mathbf{F}^+$  and  $\mathbf{T}\mathbf{T}^+\mathbf{T} = \mathbf{F}\mathbf{F}^+\mathbf{T} = \mathbf{T}$ . This result implies that for  $i = 1, \dots, m$ ,

(i)  $(\mathbf{I} - \mathbf{F}\mathbf{F}^+)(\mathbf{d}_i - \mathbf{d}_1) = 0$  and

(ii)  $(\mathbf{I} - \mathbf{F}\mathbf{F}^+)(\mathbf{E}_i - \mathbf{E}_1) = 0$ .

We now state and prove three lemmas that we use in the proof of our main result.

**Lemma 1** For  $\mathbf{T} = \mathbf{F}\mathbf{G}$  in (5), where  $\mathbf{E}_i \in \mathbb{R}_{p \times p}^S$ ,  $\mathbf{d}_i \in \mathbb{R}_{p \times 1}$ ,  $\mathbf{F} \in \mathbb{R}_{p \times q}$ , and  $\mathbf{G} \in \mathbb{R}_{q \times (m-1)(p+1)}$  such that  $\text{rank}(\mathbf{F}) = \text{rank}(\mathbf{G}) = q$ , and  $i = 1, \dots, m$ , we have that

(a)  $\mathbf{F}\mathbf{F}^+(\mathbf{E}_i - \mathbf{E}_1) = (\mathbf{E}_i - \mathbf{E}_1)\mathbf{F}\mathbf{F}^+$ ,

(b)  $\mathbf{F}\mathbf{F}^+\mathbf{E}_i = \mathbf{E}_i\mathbf{F}\mathbf{F}^+$ , and

(c)  $(\mathbf{I} - \mathbf{F}\mathbf{F}^+)\mathbf{E}_i = \mathbf{E}_1(\mathbf{I} - \mathbf{F}\mathbf{F}^+)$ .

*Proof.* Part (a) follows from the fact that  $\mathbf{E}_i = \mathbf{E}_i'$ ,  $i = 1, \dots, m$ , and from (ii) above. Parts (b) and (c) follow directly from (a).

**Lemma 2** For  $\mathbf{T} = \mathbf{F}\mathbf{G}$  in (5), where  $\mathbf{E}_i \in \mathbb{R}_{p \times p}^S$ ,  $\mathbf{d}_i \in \mathbb{R}_{p \times 1}$ ,  $\mathbf{F} \in \mathbb{R}_{p \times q}$ , and  $\mathbf{G} \in \mathbb{R}_{q \times (m-1)(p+1)}$  such that  $\text{rank}(\mathbf{F}) = \text{rank}(\mathbf{G}) = q$ , and  $i = 1, \dots, m$ , we have that  $[\mathbf{F}^+\mathbf{E}_i\mathbf{F}^{+'}]^{-1} = \mathbf{F}'\mathbf{E}_i^{-1}\mathbf{F}$ .

*Proof.* Because  $\mathbf{F}^+\mathbf{E}_i\mathbf{F}^{+'} \in \mathbb{R}_{q \times q}$  and  $\text{rank}(\mathbf{F}^+\mathbf{E}_i\mathbf{F}^{+'}) = q$ , the result follows because

$$(\mathbf{F}^+\mathbf{E}_i\mathbf{F}^{+'})(\mathbf{F}'\mathbf{E}_i^{-1}\mathbf{F}) = \mathbf{F}^+\mathbf{E}_i[\mathbf{F}\mathbf{F}^+]^{-1}\mathbf{E}_i^{-1}\mathbf{F} = \mathbf{F}^+\mathbf{F}\mathbf{F}^+\mathbf{E}_i\mathbf{E}_i^{-1}\mathbf{F} = \mathbf{I}_q.$$

**Lemma 3** Let  $\mathbf{T} = \mathbf{F}\mathbf{G}$  in (5), where  $\mathbf{E}_i \in \mathbb{R}_{p \times p}^S$ ,  $\mathbf{d}_i \in \mathbb{R}_{p \times 1}$ ,  $\mathbf{F} \in \mathbb{R}_{p \times q}$  and  $\mathbf{G} \in \mathbb{R}_{q \times (m-1)(p+1)}$  such that  $\text{rank}(\mathbf{F}) = \text{rank}(\mathbf{G}) = q$ , and  $i = 1, \dots, m$ . Also, let  $\mathbf{C} = \mathbf{R}[\mathbf{I} - \mathbf{F}\mathbf{F}^+] \in \mathbb{R}_{(p-q) \times p}$ , where  $\mathbf{R} \in \mathbb{R}_{(p-q) \times p}$  such that  $\text{rank}(\mathbf{C}) = p - q$ . Then,

(a)  $\mathbf{C}\mathbf{d}_i = \mathbf{C}\mathbf{d}_1$ ,

(b)  $\mathbf{C}\mathbf{d}_i + \mathbf{C}\mathbf{E}_i\mathbf{F}^{+'}(\mathbf{F}^+\mathbf{E}_i\mathbf{F}^{+'})^{-1}(\mathbf{y} - \mathbf{F}^+\mathbf{d}_i) = \mathbf{C}\mathbf{d}_1$ , where  $\mathbf{y} \in \mathbb{R}_{p \times 1}$ , and

(c)  $\mathbf{C}\mathbf{E}_i\mathbf{C}' - \mathbf{C}\mathbf{E}_i\mathbf{F}^{+'}(\mathbf{F}^+\mathbf{E}_i\mathbf{F}^{+'})^{-1}\mathbf{F}^+\mathbf{E}_i\mathbf{C}' = \mathbf{C}\mathbf{E}_i\mathbf{C}'$ .

*Proof.* The proof of part (a) of Lemma 3 follows from (i), and we have that

$$\mathbf{C}\mathbf{E}_i\mathbf{F}^{+'}(\mathbf{F}^+\mathbf{E}_i\mathbf{F}^{+'})^{-1} = \mathbf{C}\mathbf{E}_i\mathbf{F}^{+'}\mathbf{F}'\mathbf{E}_i^{-1}\mathbf{F} = \mathbf{C}\mathbf{F} = 0, \quad (6)$$

for each  $i \in \{1, \dots, m\}$ . Hence, (b) and (c) follow from (6).

#### 4. Linear Dimension Reduction for Multiple Heteroscedastic Nonsingular Normal Populations

We now derive a new *LDR* method that is motivated by results on linear sufficient statistics derived by [13] and by a linear feature selection theorem given in [17]. The theorem provides necessary and sufficient conditions for which a low-dimensional linear transformation of the original data will preserve the *BPMC* in the original feature space. Also, the theorem provides a representation of the *LDR* matrix.

**Theorem 1** Let  $\Pi_i$  be a  $p$ -dimensional multivariate normal population with a priori probability  $\alpha_i > 0$ , mean  $\boldsymbol{\mu}_i$ , and covariance matrix  $\boldsymbol{\Sigma}_i$ ,  $i = 1, \dots, m$ , such that  $\text{rank}(\boldsymbol{\Sigma}_i) = p$  and  $\boldsymbol{\Sigma}_1 \neq \boldsymbol{\Sigma}_j$  for some  $j \in \{2, \dots, m\}$ .

Next, let

$$\mathbf{M} \equiv \left[ \Sigma_2^{-1} \boldsymbol{\mu}_2 - \Sigma_1^{-1} \boldsymbol{\mu}_1 \mid \cdots \mid \Sigma_m^{-1} \boldsymbol{\mu}_m - \Sigma_1^{-1} \boldsymbol{\mu}_1 \mid \Sigma_2 - \Sigma_1 \mid \cdots \mid \Sigma_m - \Sigma_1 \right]. \quad (7)$$

Finally, let  $\mathbf{M} = \mathbf{F}\mathbf{G}$  be a full-rank decomposition of  $\mathbf{M}$ , where  $\text{rank}(\mathbf{M}) = q < p$ . Then, the  $p$ -dimensional Bayes procedure assigns  $\mathbf{x}$  to  $\Pi_k$  if and only if the  $q$ -dimensional Bayes procedure assigns  $\mathbf{F}^+ \mathbf{x}$  to  $\Pi_k$  for  $k \in \{1, \dots, m\}$ .

*Proof.* Let

$$\mathbf{w}_i = \begin{bmatrix} \mathbf{y}_i \\ \mathbf{u}_i \end{bmatrix} = \begin{bmatrix} \mathbf{F}^+_{q \times p} \\ \mathbf{C}_{(p-q) \times p} \end{bmatrix} \mathbf{x}_i,$$

where  $\mathbf{x}_i \sim N(\boldsymbol{\mu}_i, \Sigma_i)$ ,  $i = 1, \dots, m$ . Let  $\mathbf{H} \equiv \begin{bmatrix} \mathbf{F}^+ & \mathbf{C}' \end{bmatrix} \in \mathbb{R}_{p \times p}$  be full-rank with  $\mathbf{C} \equiv \mathbf{R}(\mathbf{I} - \mathbf{F}\mathbf{F}^+)$ , where  $\mathbf{R} \in \mathbb{R}_{(p-q) \times p}$  such that  $\text{rank}(\mathbf{R}) = p - q$ . Then  $\mathbf{w}_i \sim N(\mathbf{H}\boldsymbol{\mu}_i, \mathbf{H}\Sigma_i\mathbf{H}')$ , where

$$\mathbf{H}\boldsymbol{\mu}_i = \begin{bmatrix} \mathbf{F}^+ \boldsymbol{\mu}_i \\ \mathbf{C}\boldsymbol{\mu}_i \end{bmatrix} \text{ and } \mathbf{H}\Sigma_i\mathbf{H}' = \begin{bmatrix} \mathbf{F}^+ \Sigma_i \mathbf{F}^+ & \mathbf{F}^+ \Sigma_i \mathbf{C}' \\ \mathbf{C} \Sigma_i \mathbf{F}^+ & \mathbf{C} \Sigma_i \mathbf{C}' \end{bmatrix},$$

$i = 1, \dots, m$ . By Lemma 3, we conclude that  $E(\mathbf{u}_i | \mathbf{y}_i) = \mathbf{C}\boldsymbol{\mu}_i$  and  $\text{Var}(\mathbf{u}_i | \mathbf{y}_i) = \mathbf{C}\Sigma_i\mathbf{C}'$  for  $i = 1, \dots, m$ . Thus,  $\mathbf{w}_i \sim N(\mathbf{F}^+ \boldsymbol{\mu}_i, \mathbf{F}^+ \Sigma_i \mathbf{F}^+) \cdot N(\mathbf{C}\boldsymbol{\mu}_i, \mathbf{C}\Sigma_i\mathbf{C}')$ . That is,  $p(\mathbf{u} | \mathbf{y}, \Pi_k)$  does not depend on class  $k \in \{2, \dots, m\}$ .

Recall that for  $i, j = 1, \dots, m$ ,  $i \neq j$ , the  $p$ -variate Bayes procedure defined by (3) assigns  $\mathbf{x}$  to  $\Pi_j$  if and only if

$$\begin{aligned} \alpha_j p(\mathbf{x} | \Pi_j) > \alpha_i p(\mathbf{x} | \Pi_i) &\Leftrightarrow \alpha_j p(\mathbf{w} | \Pi_j) > \alpha_i p(\mathbf{w} | \Pi_i) \\ &\Leftrightarrow \alpha_j p(\mathbf{u} | \mathbf{y}, \Pi_j) p(\mathbf{y} | \Pi_j) > \alpha_i p(\mathbf{u} | \mathbf{y}, \Pi_i) p(\mathbf{y} | \Pi_i) \\ &\Leftrightarrow \alpha_j p(\mathbf{y} | \Pi_j) > \alpha_i p(\mathbf{y} | \Pi_i). \end{aligned}$$

Therefore, the original  $p$ -variate Bayes classification assignment is preserved by the linear transformation  $\mathbf{y} = \mathbf{F}^+ \mathbf{x}$ .

Theorem 1 is important in that if its conditions hold, we obtain a *LDR* matrix for the reduced  $q$ -dimensional subspace such that the *BPMC* in the  $q$ -dimensional space is equal to the *BPMC* for the original  $p$ -dimensional feature space. In other words, provided the conditions in Theorem 1 hold, we have that the *LDR* matrix  $\mathbf{F}^+ \in \mathbb{R}_{q \times p}$  exists and that  $\text{BPMC}_p = \text{BPMC}_q$ , where  $\text{rank}(\mathbf{M}) = q < p$ .

With the following corollary, we demonstrate that for two multivariate normal populations such that  $\Sigma_1 = \Sigma_2 = \Sigma$  and  $\boldsymbol{\mu}_2 \neq \boldsymbol{\mu}_1$ , our *LDR* matrix derived in Theorem 1 reduces to the *LDF* of [4].

**Corollary 1** Assuming we have two multivariate normal populations  $N(\boldsymbol{\mu}_1, \Sigma)$  and  $N(\boldsymbol{\mu}_2, \Sigma)$ , the proposed *LDR* matrix in Theorem 1 reduces to  $\mathbf{M} = \Sigma^{-1}(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)$ , which is the well-known Fisher's *LDF*.

*Proof.* The proof is immediate from (7).

## 5. Low-Dimensional Graphical Representations for Heteroscedastic Multivariate Normal Populations

### 5.1. Low-Dimensional LDR Using the SVD

If  $\text{rank}(\mathbf{M}) = p$ , one cannot use Theorem 1 to directly obtain a  $q \times p$  *LDR* matrix that preserves the full-feature *BPMC*. Also, in some situations when Theorem 1 holds, we may desire to determine a low-dimensional representation with dimension less than  $q$ , say  $r$ , where  $1 \leq r < q < p$ . Even if the required conditions of Theorem 1 hold, we may desire to determine a low-dimensional representation with dimension less than  $q$ , say  $r$ , where  $1 \leq r < q < p$ . Thus, we seek to construct an  $r$ -dimensional representation which preserves as much of the original  $p$ -dimensional *BPMC* as possible.

One method of obtaining an  $r$ -dimensional *LDR* matrix,  $1 \leq r \leq q < p$ , is the *SVD* approximation to  $\mathbf{M}$  in (7). We use the following theorem from [17] to determine such an  $r$ -dimensional *LDR* matrix.

**Theorem 2** Let  $\mathbf{C}_{s,t}^{(p)}$  denote the class of all  $s \times t$  real matrices of rank  $p$ , and let  $\mathbf{C}_{s,t}^{(r)}$  denote the class of all  $s \times t$  real matrices of rank  $r$ , where  $1 \leq r < q < p$ . If  $\mathbf{A}_p \in \mathbf{C}_{s,t}^{(p)}$  and  $\mathbf{A}_r \in \mathbf{C}_{s,t}^{(r)}$ , given by  $\mathbf{A}_r = \mathbf{U}\mathbf{D}_r\mathbf{V}'$ , then

$$\|\mathbf{A}_p - \mathbf{A}_r\| < \|\mathbf{A}_p - \mathbf{X}\| \quad \text{for all } \mathbf{X} \in \mathbf{C}_{s,t}^{(r)},$$

where  $\mathbf{A}_p = \mathbf{U}\mathbf{D}_p\mathbf{V}'$ ,  $\mathbf{D}_p = \text{diag}(d_1, \dots, d_p)$ ,  $\mathbf{D}_r = \text{diag}(d_1, \dots, d_r, 0_{r+1}, \dots, 0_p)$ , and  $\|\mathbf{A}_p\|$  is the usual Euclidean or Frobenius norm of a matrix  $\mathbf{A}_p$ , given by

$$\|\mathbf{A}_p\| = \left( \sum_{i=1}^s \sum_{j=1}^t |a_{ij}|^2 \right)^{1/2} = \left( \sum_{i=1}^p d_i^2 \right)^{1/2}.$$

Furthermore,  $\|\mathbf{A}_p - \mathbf{A}_r\| = \left( \sum_{i=r+1}^p d_i^2 \right)^{1/2}$ .

Using Theorems 1 and 2, we now construct a linear transformation for projecting high-dimensional data onto a low-dimensional subspace when all class distribution parameters are known. Let  $\mathbf{M} = \mathbf{U}\mathbf{D}_p\mathbf{V}$  be the SVD of the matrix  $\mathbf{M}$ , where  $\mathbf{D}_p \equiv \text{diag}(\lambda_1, \dots, \lambda_p)$  for  $\lambda_i$ ,  $i = 1, \dots, p$ , which are the singular values of  $\mathbf{M}$  with  $\lambda_i \geq \lambda_j$  for  $1 \leq i < j \leq p$ ,  $\lambda_q > 0$ ,  $\lambda_{q+1} \geq 0$  for  $1 \leq q \leq p$ . Let  $\mathbf{F} = \mathbf{U}\mathbf{D}_p$  and define  $\mathbf{D}_r \equiv \text{diag}(\lambda_1, \dots, \lambda_r, 0_{r+1}, \dots, 0_p)$  with  $\lambda_i \geq \lambda_j$  for  $1 \leq i < j \leq q$ . From Theorem 2, we have that  $\mathbf{M}_r = \mathbf{U}\mathbf{D}_r\mathbf{V}'$  is a rank- $r$  approximation of  $\mathbf{M}$ , and, therefore, a rank- $r$  approximation to  $\mathbf{F}$  is  $\mathbf{F}_r = \mathbf{U}\mathbf{D}_r$ . Thus,  $\mathbf{F}_r'$  is an  $r \times p$  LDR matrix that yields an  $r$ -dimensional representation of the original  $p$ -dimensional class models. One can also use  $\mathbf{F}_r^+$  to construct low-dimensional representations of high-dimensional class densities.

We next provide an example to demonstrate the efficacy of Theorems 1 and 2 to determine low-dimensional representations for multiple multivariate normal populations with known mean vectors and covariance matrices. In the example, we display the simplicity of Theorem 1 to formulate a low-dimensional representation for three populations ( $m = 3$ ) with unequal covariance matrices and original dimension  $p = 6$ . Note that unlike the low-dimensional representation of [18], our Theorem 1 does not restrict the reduced dimension to be  $r = 1$ .

### 5.2. Example

Consider the configuration  $N(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$ ,  $N(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$ , and  $N(\boldsymbol{\mu}_3, \boldsymbol{\Sigma}_3)$ , where

$$\boldsymbol{\mu}_1 = [0, 0, 0, 0, 0, 0],$$

$$\boldsymbol{\mu}_2 = [1, 1, 1, 1, 1, 1],$$

$$\boldsymbol{\mu}_3 = [2, 2, 2, 2, 2, 2],$$

$$\boldsymbol{\Sigma}_1 = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}, \quad \boldsymbol{\Sigma}_2 = \begin{bmatrix} 2 & 1 & 1 & 1 & 1 & 1 \\ 1 & 2 & 1 & 1 & 1 & 1 \\ 1 & 1 & 2 & 1 & 1 & 1 \\ 1 & 1 & 1 & 2 & 1 & 1 \\ 1 & 1 & 1 & 1 & 2 & 1 \\ 1 & 1 & 1 & 1 & 1 & 2 \end{bmatrix}, \quad \text{and} \quad \boldsymbol{\Sigma}_3 = \begin{bmatrix} 2 & 1 & 0 & 1 & 1 & 1 \\ 1 & 2 & 0 & 1 & 1 & 1 \\ 0 & 0 & 7 & 0 & 0 & 0 \\ 1 & 1 & 0 & 2 & 1 & 1 \\ 1 & 1 & 0 & 1 & 2 & 1 \\ 1 & 1 & 0 & 1 & 1 & 2 \end{bmatrix}.$$

We have  $\text{rank}(\mathbf{M}) = 2$ , and, thus, by Theorem 1, the six-dimensional multivariate normal densities can be compressed to the dimension  $q = 2$  without increasing the BPMC.

Using Theorem 1, we have that an optimal two-dimensional representation space is  $\text{Col}(\mathbf{F})$ , where

$$\mathbf{F} = \begin{bmatrix} 0.3800 & 0.3800 & 0.5272 & 0.3800 & 0.3800 & 0.3800 \\ 0.2358 & 0.2358 & -0.8497 & 0.2358 & 0.2358 & 0.2358 \end{bmatrix}'. \tag{8}$$

The optimal two-dimensional ellipsoidal representation is shown in **Figure 1**.

We can also determine a one-dimensional representation of the three multivariate normal populations through application of the *SVD* described in Theorem 2 applied to the matrix  $\mathbf{M}$  given in (7). A one-dimensional representation space is column one of the matrix  $\mathbf{F}$  in (8), and the graphical representation of this configuration of univariate normal densities is depicted in Figure 2.

### 6. Four LDR Methods for Statistical Discrimination

In this section, we present and describe the four *LDR* methods that we wish to compare and contrast in Sections 7 and 8.

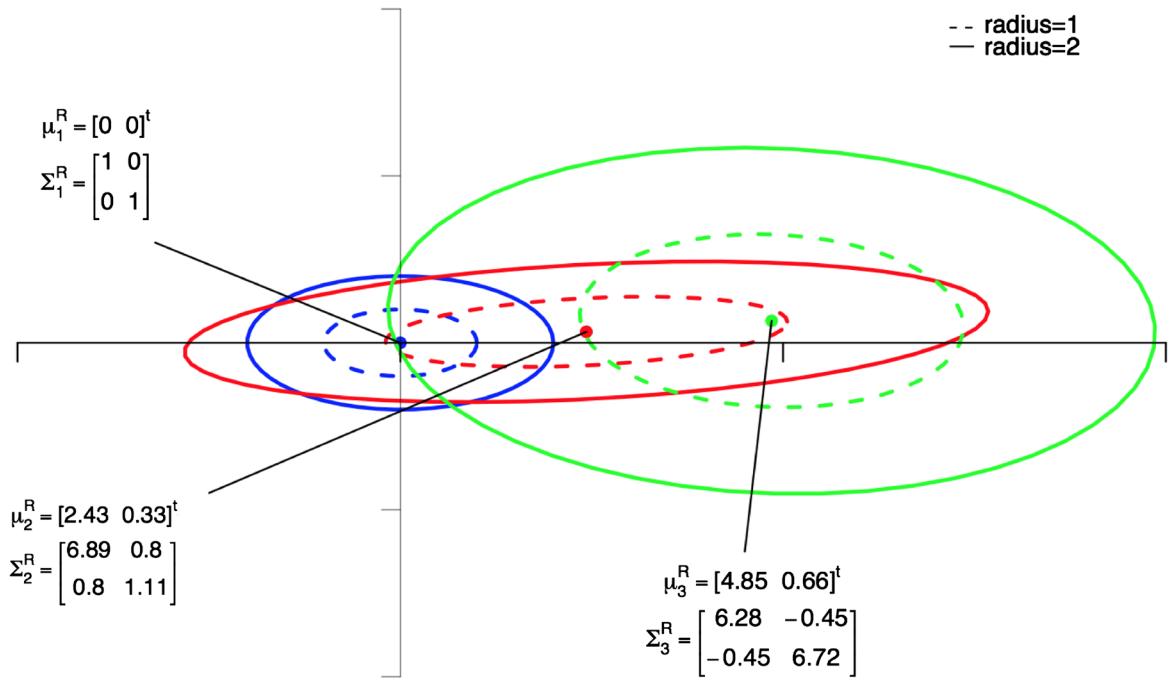


Figure 1. The optimal two-dimensional representation for normal densities from the example in Section 5.2.

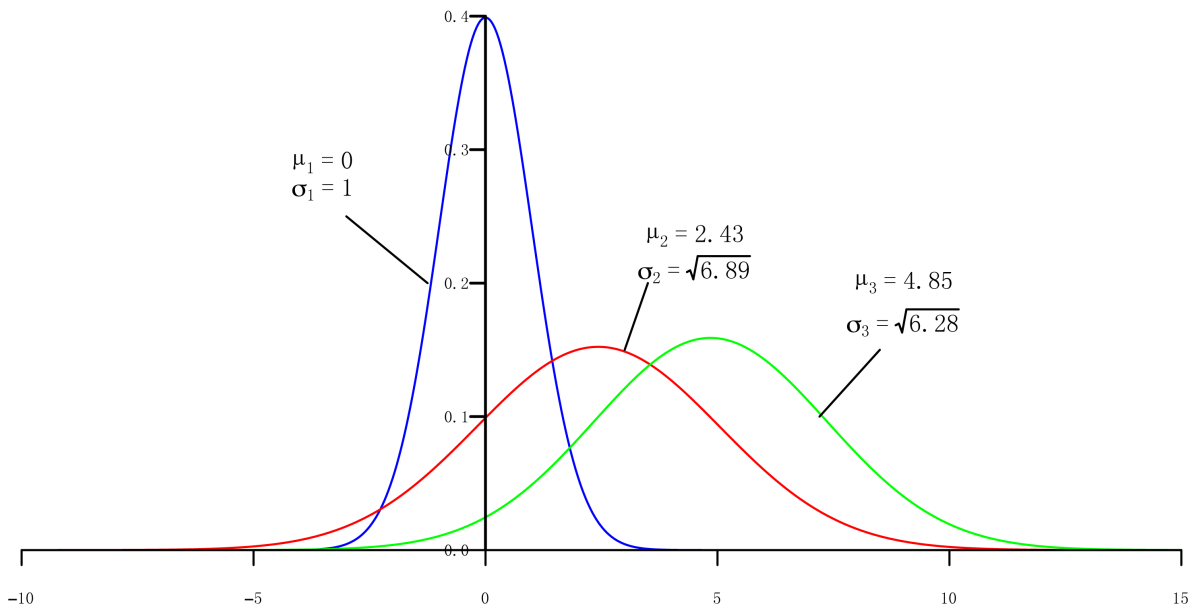


Figure 2. The optimal one-dimensional representation for normal densities from the example in Section 5.2.

### 6.1. The SY Method

In Theorem 1, we assume the parameters  $\boldsymbol{\mu}_i$  and  $\boldsymbol{\Sigma}_i$ ,  $i=1, \dots, m$ , are known, but in reality this assumption is rarely the case. In a sampling situation, we can replace the columns of the  $\mathbf{M}$  matrix in (7) by their sample estimators yielding

$$\hat{\mathbf{M}} \equiv \left[ S_2^{-1}\bar{\mathbf{x}}_2 - S_1^{-1}\bar{\mathbf{x}}_1 \mid S_3^{-1}\bar{\mathbf{x}}_3 - S_1^{-1}\bar{\mathbf{x}}_1 \mid \dots \mid S_m^{-1}\bar{\mathbf{x}}_m - S_1^{-1}\bar{\mathbf{x}}_1 \mid S_2 - S_1 \mid \dots \mid S_m - S_1 \right],$$

provided  $n_i > p$ ,  $i=1, \dots, m$ . Our estimator  $\hat{\mathbf{M}}$ , along with Theorems 1 and 2, yields a LDR technique based on the selection of an  $r$ -dimensional hyperplane determined from a rank- $r$  approximation to the full-rank matrix  $\hat{\mathbf{M}}$ .

Thus, using the SVD, we let  $\hat{\mathbf{M}} = \mathbf{U}\mathbf{D}_p\mathbf{V}$ , where  $\mathbf{D}_p \equiv \text{diag}(\lambda_1, \dots, \lambda_p)$ ,  $\lambda_j > 0$  are the singular values of  $\hat{\mathbf{M}}$  for  $j=1, \dots, p$ , and let  $\hat{\mathbf{F}} = \mathbf{U}\mathbf{D}_p$ . Also, let

$$\mathbf{D}_r \equiv \text{diag}(\lambda_1, \dots, \lambda_r, 0_{r+1}, \dots, 0_p) \quad (9)$$

with  $\lambda_j \geq \lambda_l$  for  $1 \leq j < l \leq p$  and  $1 \leq r < p$ . From Theorem 2, we have that  $\tilde{\mathbf{M}}_r = \mathbf{U}\mathbf{D}_r\mathbf{V}'$  is a rank- $r$  approximation of  $\hat{\mathbf{M}}$ , and a rank- $r$  approximation to  $\mathbf{F}$  is  $\tilde{\mathbf{F}}_r = \mathbf{U}\mathbf{D}_r$ . Thus,  $\tilde{\mathbf{F}}_r'$  is our new  $r \times p$  LDR matrix.

Provided  $\left[ \left( \sum_{j=r+1}^p \lambda_j \right) / \left( \sum_{i=1}^p \lambda_i \right) \right]$  is relatively small,  $\tilde{\mathbf{F}}_r'$  will yield an  $EPMC(r)$  such that

$EPMC(p) \approx EPMC(r)$ , and in certain population parameter configurations, we may have that  $EPMC(r) < EPMC(p)$ . We refer to our new LDR matrix as the SY LDR method.

### 6.2. The BE Method

A second LDR method presented by [14] is

$$\mathbf{U} \equiv \left[ (\boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2)^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \mid \dots \mid (\boldsymbol{\Sigma}_i + \boldsymbol{\Sigma}_j)^{-1}(\boldsymbol{\mu}_i - \boldsymbol{\mu}_j) \mid \dots \mid (\boldsymbol{\Sigma}_{m-1} + \boldsymbol{\Sigma}_m)^{-1}(\boldsymbol{\mu}_{m-1} - \boldsymbol{\mu}_m) \right]$$

for  $1 \leq i < j \leq m$ , provided all multivariate normal population parameters are known. For the unknown parameter case, an estimator of  $\mathbf{U}$  is then

$$\hat{\mathbf{U}} \equiv \left[ (S_1 + S_2)^{-1}(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2) \mid \dots \mid (S_i + S_j)^{-1}(\bar{\mathbf{x}}_i - \bar{\mathbf{x}}_j) \mid \dots \mid (S_{m-1} + S_m)^{-1}(\bar{\mathbf{x}}_{m-1} - \bar{\mathbf{x}}_m) \right]. \quad (10)$$

The LDR matrix derived in [14], which we refer to as the BE matrix, was also obtained when we determined low-rank approximation to  $\hat{\mathbf{U}}$  by using the SVD. That is, let  $\hat{\mathbf{U}} = \mathbf{R}\mathbf{D}_p\mathbf{S}'$  be the SVD of  $\hat{\mathbf{U}}$ , where

$$\mathbf{D}_p \equiv \text{diag}(\lambda_1, \dots, \lambda_p) \quad (11)$$

with  $\lambda_j \geq \lambda_l$  for  $1 \leq j < l \leq p$ , and let  $\hat{\mathbf{H}} = \mathbf{R}\mathbf{D}_p$ . Define  $\mathbf{D}_r$  as in (9) with  $\lambda_j \geq \lambda_l$  for  $1 \leq j \leq r$ . Then,  $\tilde{\mathbf{U}} = \mathbf{R}\mathbf{D}_r\mathbf{S}'$  is a rank- $r$  approximation of  $\hat{\mathbf{U}}$ , and  $\tilde{\mathbf{H}} = \mathbf{R}\mathbf{D}_r = \mathbf{R}_r \in \mathbb{R}^{p \times r}$  is a rank- $r$  approximation of  $\hat{\mathbf{H}}$ . Thus, the BE LDR matrix is  $\tilde{\mathbf{H}}_r'$ , or, equivalently,  $\mathbf{R}_r'$ , where the  $j^{\text{th}}$  eigenvector corresponds to the  $j^{\text{th}}$  largest singular value and  $j=1, \dots, r$ .

The BE LDR approach is based on the rotated differences in the means. The LDR matrix (10) uses a type of pooled covariance matrix estimator for the precision matrices. However, the BE method does not incorporate all of the information contained in the different individual covariance matrices. Another disadvantage of the BE LDR approach is that it is limited to a reduced dimension that depends on the number of classes,  $m$ . For  $m=2$ , BE allows one to reduce the data to only one dimension, regardless of the full-feature vector dimension. Therefore, one may lose some discriminatory information with the application of the BE LDR method when the covariance matrices are considerably different.

### 6.3. Sliced Inverse Regression (SIR)

The next LDR method we consider is sliced inverse regression (SIR), which was proposed in [16]. Assuming all population parameters are known, we first define  $\boldsymbol{\Sigma}_w \equiv \sum_{i=1}^m \boldsymbol{\Sigma}_i / m$  as the within-group covariance matrix and  $\boldsymbol{\Sigma}_B \equiv \sum_{i=1}^m (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})' / m$  as the between-group covariance matrix, where  $\boldsymbol{\mu} \equiv \sum_{i=1}^m \boldsymbol{\mu}_i / m$  is the overall



population mean. As its criterion matrix, the *SIR LDR* method uses

$$\mathbf{M}_{SIR} \equiv \Gamma^{-1/2} \Sigma_B \Gamma^{-1/2}, \quad (12)$$

where  $\Gamma \equiv \Sigma_B + \Sigma_W$  is the marginal covariance matrix of  $\mathbf{x}$ . When population parameters must be estimated, an estimator of (12) is

$$\hat{\mathbf{M}}_{SIR} \equiv \hat{\Gamma}^{-1/2} \mathbf{S}_B \hat{\Gamma}^{-1/2}, \quad (13)$$

where  $\hat{\Gamma} \equiv \mathbf{S}_B + \mathbf{S}_W$  with  $\mathbf{S}_W$  and  $\mathbf{S}_B$  given in (1) and (2), respectively. Let  $\hat{\mathbf{M}}_{SIR} = \mathbf{G} \mathbf{D}_p \mathbf{P}'$  be the SVD of (13), and let  $\hat{\mathbf{K}} = \hat{\Gamma}^{-1/2} \mathbf{G}_p$ , where  $\mathbf{G}_p$  is composed of eigenvectors of (13) such that the  $j^{\text{th}}$  eigenvector of (13) corresponds to the  $j^{\text{th}}$  largest singular value of (13),  $j = 1, \dots, p$ . Then,  $\hat{\mathbf{K}}_r = \hat{\Gamma}^{-1/2} \mathbf{G}_r$  is a rank- $r$  approximation of  $\hat{\mathbf{K}}$ , and the  $r \times p$  *SIR LDR* matrix is  $\hat{\mathbf{K}}'_r$ , which is composed of the eigenvectors corresponding to the  $r$  largest singular values of  $\mathbf{K}$ .

#### 6.4. Sliced Average Variance Estimation (SAVE)

The last *LDR* method we consider is sliced average variance estimation (*SAVE*), which has been proposed in [19]-[20]. The *SAVE* method uses the  $p \times p$  criterion matrix  $\mathbf{M}_{SAVE}^* \equiv \sum_{i=1}^m (\mathbf{I}_p - \Gamma^{-1/2} \Sigma_i \Gamma^{-1/2})^2 / m$ , where  $\Gamma \equiv \Sigma_B + \Sigma_W$ . We use a form of *SAVE* given in [21], which is

$$\mathbf{M}_{SAVE} \equiv (\Gamma^{-1/2} \Sigma_B \Gamma^{-1/2})^2 + \Gamma^{-1/2} \Sigma_{\mathbf{r}} \Gamma^{-1/2}, \quad (14)$$

where  $\Sigma_{\mathbf{r}} \equiv \sum_{i=1}^m (\Sigma_i - \Sigma_W) \Gamma^{-1} (\Sigma_i - \Sigma_W) / m$ . An estimator of  $\mathbf{M}_{SAVE}$  is

$$\hat{\mathbf{M}}_{SAVE} \equiv (\hat{\Gamma}^{-1/2} \mathbf{S}_B \hat{\Gamma}^{-1/2})^2 + \hat{\Gamma}^{-1/2} \mathbf{S}_{\mathbf{r}} \hat{\Gamma}^{-1/2}, \quad (15)$$

where,  $\mathbf{S}_{\mathbf{r}} \equiv \sum_{i=1}^m (\mathbf{S}_i - \mathbf{S}_W) \hat{\Gamma}^{-1} (\mathbf{S}_i - \mathbf{S}_W) / m$  with  $\mathbf{S}_W$  and  $\mathbf{S}_B$  given in (1) and (2), respectively. Next, let  $\hat{\mathbf{M}}_{SAVE} = \mathbf{B} \mathbf{D}_p \mathbf{Q}'$  be the SVD of (15), and let  $\hat{\mathbf{L}} = \hat{\Gamma}^{-1/2} \mathbf{B}_p$ , where  $\mathbf{B}_p$  is composed of the eigenvectors of (15) such that the  $j^{\text{th}}$  column of  $\mathbf{B}_p$  corresponds to the  $j^{\text{th}}$  of (15) with the  $j^{\text{th}}$  largest singular value,  $j = 1, \dots, p$ . Then,  $\hat{\mathbf{L}}_r = \hat{\Gamma}^{-1/2} \mathbf{B}_r$  is a rank- $r$  approximation of  $\hat{\mathbf{L}}$ , and, thus, the  $r \times p$  *SAVE LDR* matrix is  $\hat{\mathbf{L}}'_r$ . An alternative representation of the  $r$ -dimensional *SAVE LDR* matrix is  $\mathbf{B}_r$ , which is composed of the eigenvectors corresponding to the  $r$  largest singular values of  $\mathbf{B}$ .

### 7. A Monte Carlo Comparison of Four LDR Methods for Statistical Classification

Here, we compare our new *SY LDR* method derived above to the *BE*, *SIR*, and *SAVE LDR* methods. Specifically, we evaluate the classification efficacy in terms of the *EPMC* for the *SY*, *BE*, *SIR*, and *SAVE LDR* methods using Monte Carlo simulations for five different configurations of multivariate normal populations with  $p = 10$ . We have generated 10,000 training-sample and test datasets from the appropriate multivariate normal distributions for each parameter configuration. The test data were assigned to either population class  $\Pi_1$  or  $\Pi_2$  when  $m = 2$  or  $\Pi_1$ ,  $\Pi_2$ , or  $\Pi_3$  when  $m = 3$  using the sample *QDF* corresponding to (4). We have applied the four competing *LDR* matrices  $\tilde{\mathbf{F}}'_r$ ,  $\tilde{\mathbf{H}}'_r$ ,  $\tilde{\mathbf{K}}'_r$ , and  $\tilde{\mathbf{L}}'_r$  and calculated the *EPMCs* for the full-dimensional *QDF* and for the four reduced-dimensional *QDFs* by averaging the estimated conditional error rate over all training samples. We examined the effect of the training-sample sizes on the four *LDR* methods using sample sizes  $n_i = (2.5)p$  and  $n_i = 5p$ ,  $i = 1, 2$  or  $i = 1, 2, 3$ .

For the *SY* and *SAVE LDR* approaches, we reduce the dimension to  $r = 1, 2, 3$ . In general, for  $m$  populations, we remark that the *BE LDR* method can reduce the feature vector to at most the dimension  $r = m(m-1)/2$  because that is the column dimension of the matrix  $\hat{\mathbf{U}} \in \mathbb{R}_{p \times r}$ . This limitation is potentially a major drawback when the ratio  $m/p$  is small and especially when  $m = 2$ . In particular, for the *BE* and *SIR LDR* approaches, we can reduce only to the dimension  $r = 1$  when  $m = 2$ . When  $m = 3$ , the *BE LDR* method can be applied to reduce the dimensions to at most  $r = 2, 3$ , and for the *SIR LDR* method, we can reduce the number of original features to at most  $r = m - 1$ . However, the *SY LDR* method avoids this shortcoming and allows one to reduce the original feature vector to our choice of reduced dimension  $r$ , where  $1 \leq r < p$ , for any finite number of

populations  $m$ .

The five Monte Carlo simulations were generated using the programming language R. **Table 1** gives a description of the number of populations and the theoretically optimal rank of the four  $LDR$  indices for each configuration. In **Figures 3-7**, we display the corresponding  $EPMCs$  of the four competing  $LDR$  methods for the various population configurations and values of  $n_i$  and  $r$ , where  $i = 1, \dots, m$ . The estimated standard error of all  $EPMCs$  in the following tables was less than 0.001. For each configuration, we also calculated  $\hat{M}$ ,  $\hat{U}$ ,  $\hat{K}$ , and  $\hat{L}$  along with their respective singular values using the  $SVD$ . These singular values contain information concerning the amount of discriminatory information available in each reduced dimension. In the subsequent subsections, we use the following notation:  $EPMC_r(SY)$ ,  $EPMC_r(BE)$ ,  $EPMC_r(SIR)$ , and  $EPMC_r(SAVE)$  denote the estimated  $EPMCs$  for the  $SY$ ,  $BE$ ,  $SIR$ , and  $SAVE$   $LDR$  methods, respectively, for each of the appropriate reduced dimensions.

### 7.1. Configuration 1: $m = 2$ with Moderately Different Covariance Matrices

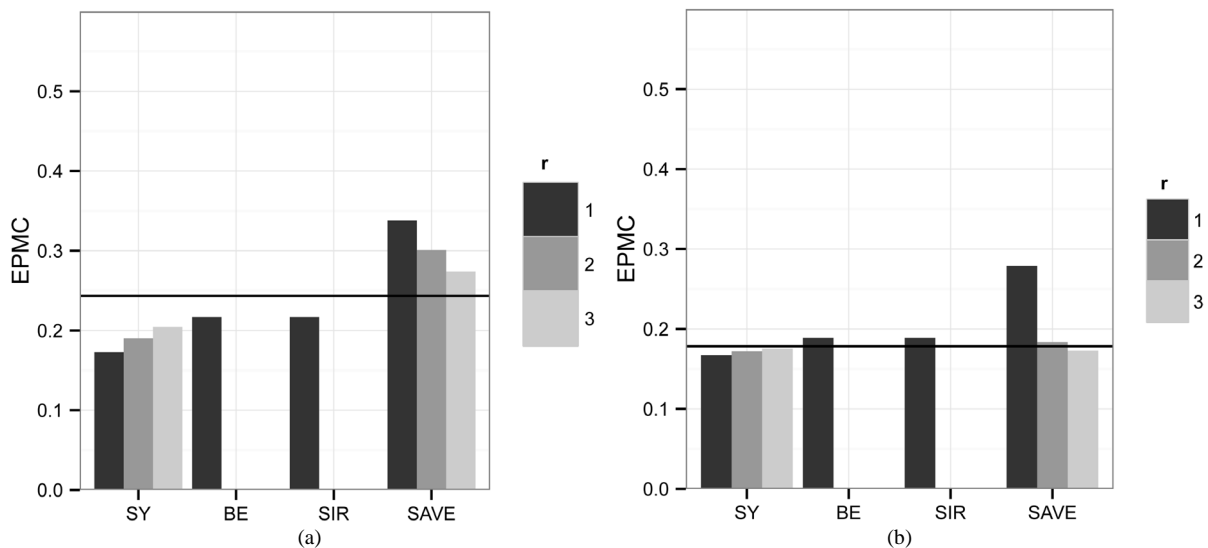
The first population configuration we examined was composed of two multivariate normal populations  $N_p(\mathbf{0}, \Sigma_1)$  and  $N_p(\mu_2, \Sigma_2)$ , where  $p = 10$ ,

$$\mu_1 = [0, 0, 0, 0, 0, 0, 0, 0, 0, 0]',$$

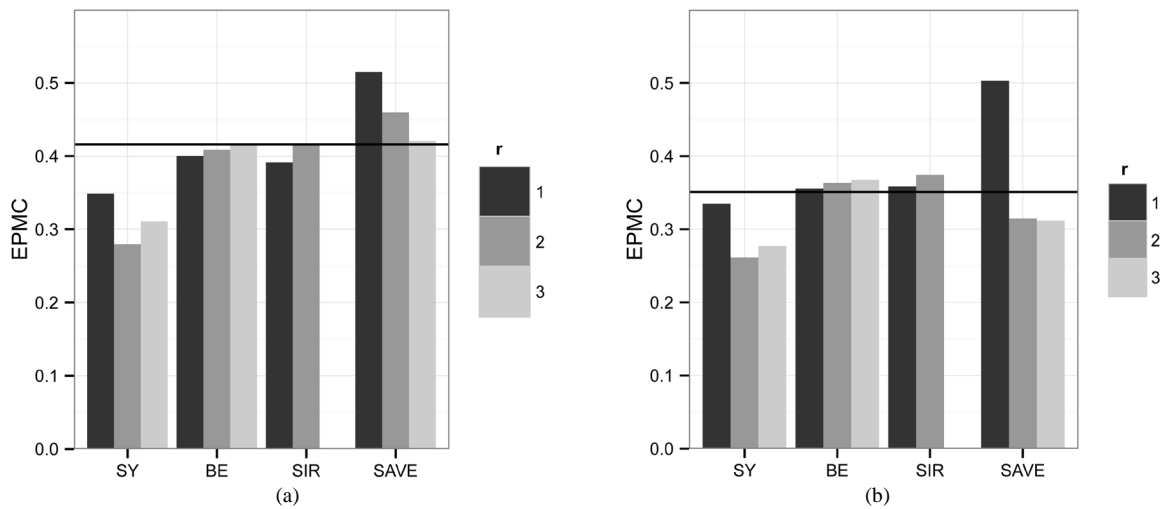
$$\mu_2 = [5, 5, 5, 5, 5, 5, 5, 5, 5, 5]',$$

**Table 1.** A description of Monte Carlo simulation parametric configurations and singular values in Section 2 with unequal covariance matrices and  $p = 10$ .

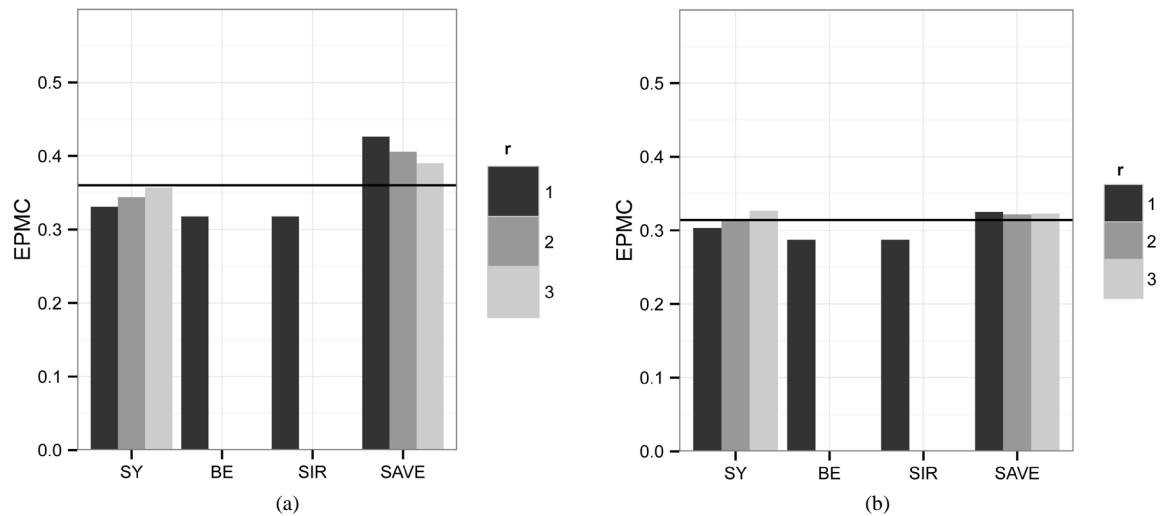
Configuration	Means	Rank				
		$m$	$M$	$U$	$K$	$L$
1	Moderately separated unequal means	2	3	1	1	2
2	Relatively close unequal means	3	2	1	1	2
3	Relatively close unequal means	2	2	1	1	2
4	Relatively close unequal means	3	2	1	1	2
5	Close unequal means	2	4	3	2	4



**Figure 3.** Simulation results from Section 7.1. The horizontal bar in each graph represents the  $EPMC_{10}$ . (a)  $n_i = 25$ ; (b)  $n_i = 50$ .



**Figure 4.** Simulation results from Configuration 2. The horizontal bar in each graph represents the  $EPMC_{10}$  with no dimension reduction. (a)  $n_i = 25$  ; (b)  $n_i = 50$  .



**Figure 5.** Simulation results from Section 7.3. The horizontal bar in each graph represents the  $EPMC_{10}$  . (a)  $n_i = 25$  ; (b)  $n_i = 50$  .

$$\Sigma_1 = \begin{bmatrix} 10.00 & 4.00 & 5.00 & 4.00 & 3.00 & 4.00 & 4.00 & 5.00 & 4.00 & 3.00 \\ 4.00 & 10.00 & 5.00 & 2.00 & 4.00 & 3.00 & 3.00 & 5.00 & 4.00 & 3.00 \\ 5.00 & 5.00 & 10.00 & 5.00 & 5.00 & 3.00 & 4.00 & 4.00 & 4.00 & 4.00 \\ 4.00 & 2.00 & 5.00 & 10.00 & 3.00 & 4.00 & 2.00 & 3.00 & 4.00 & 3.00 \\ 3.00 & 4.00 & 5.00 & 3.00 & 10.00 & 3.00 & 4.00 & 5.00 & 3.00 & 3.00 \\ 4.00 & 3.00 & 3.00 & 4.00 & 3.00 & 12.00 & 3.00 & 4.00 & 4.00 & 4.00 \\ 4.00 & 3.00 & 4.00 & 2.00 & 4.00 & 3.00 & 14.00 & 2.00 & 2.00 & 2.00 \\ 5.00 & 5.00 & 4.00 & 3.00 & 5.00 & 4.00 & 2.00 & 12.00 & -0.50 & -0.50 \\ 4.00 & 4.00 & 4.00 & 4.00 & 3.00 & 4.00 & 2.00 & -0.50 & 14.00 & -1.00 \\ 3.00 & 3.00 & 4.00 & 3.00 & 3.00 & 4.00 & 2.00 & -0.50 & -1.00 & 11.00 \end{bmatrix},$$

and

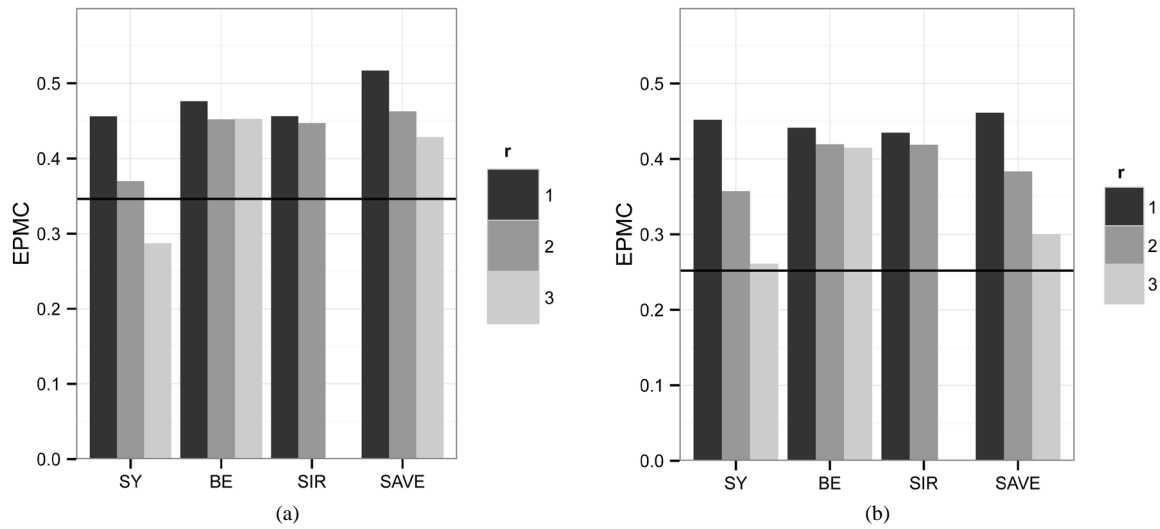


Figure 6. Simulation results from Configuration 4. The horizontal bar in each graph represents the  $EPMC_{10}$  with no dimension reduction. (a)  $n_i = 25$  ; (b)  $n_i = 50$  .

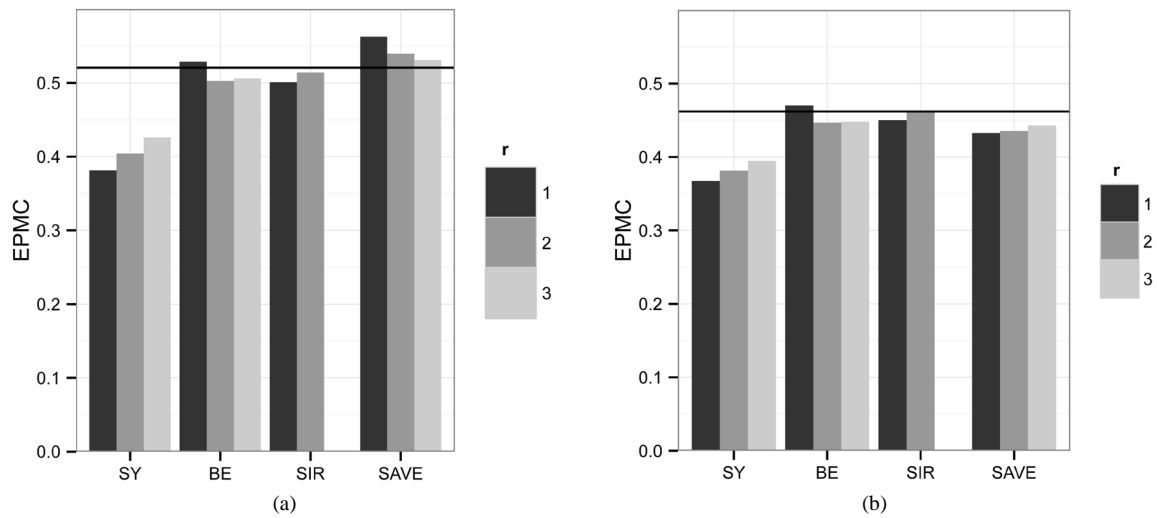


Figure 7. Simulation results from Section 7.5. The horizontal bar in each graph represents the  $EPMC_{10}$ . (a)  $n_i = 25$  ; (b)  $n_i = 50$  .

$$\Sigma_2 = \begin{bmatrix} 16.00 & 10.00 & 11.00 & 10.00 & 9.00 & 10.00 & 10.00 & 11.00 & 10.00 & 9.00 \\ 10.00 & 16.00 & 11.00 & 8.00 & 10.00 & 9.00 & 9.00 & 11.00 & 10.00 & 9.00 \\ 11.00 & 11.00 & 16.00 & 11.00 & 11.00 & 9.00 & 10.00 & 10.00 & 10.00 & 10.00 \\ 10.00 & 8.00 & 11.00 & 16.00 & 9.00 & 10.00 & 8.00 & 9.00 & 10.00 & 9.00 \\ 9.00 & 10.00 & 11.00 & 9.00 & 16.00 & 9.00 & 10.00 & 11.00 & 9.00 & 9.00 \\ 10.00 & 9.00 & 9.00 & 10.00 & 9.00 & 18.00 & 9.00 & 10.00 & 10.00 & 10.00 \\ 10.00 & 9.00 & 10.00 & 8.00 & 10.00 & 9.00 & 20.00 & 8.00 & 8.00 & 8.00 \\ 11.00 & 11.00 & 10.00 & 9.00 & 11.00 & 10.00 & 8.00 & 22.00 & 9.50 & 9.50 \\ 10.00 & 10.00 & 10.00 & 10.00 & 9.00 & 10.00 & 8.00 & 9.50 & 24.00 & 9.00 \\ 9.00 & 9.00 & 10.00 & 9.00 & 9.00 & 10.00 & 8.00 & 9.50 & 9.00 & 21.00 \end{bmatrix}$$

Here,  $\text{rank}(\Sigma_2 - \Sigma_1) = 1$ , which implies  $\text{rank}(\mathbf{M}) = 2$  because  $(\Sigma_2^{-1}\mu_2 - \Sigma_1^{-1}\mu_1) \notin \text{span}(\Sigma_2 - \Sigma_1)$ . The singular values of  $\hat{\mathbf{M}}$  in Table 2 indicate that most of the classificatory information can be captured when  $r = 1$ , because the subsequent singular values are small relative to the first. Also, the BE LDR technique loses classificatory information from pooling the acutely dissimilar pair of covariance matrices.

When  $n_i = 25$ , the EPMC was reduced by the SY, BE, and SIR LDR methods but not for the SAVE LDR method. This effect occurred because the training-sample size  $n_i = 25$ ,  $i = 1, 2$ , are small relative to the full-feature dimensionality  $p = 10$ , and, therefore, insufficient data were available to accurately estimate the  $p(p+3)$  total population parameters. Thus, by reducing the full-feature dimension  $p = 10$  to dimension  $r \ll 10$ , we considerably increased the ratio of the training-sample size relative to the original new dimension so that  $n_i/r \gg n_i/p$  to  $n_i/r$ , where  $r \ll p$ . Thus, we achieved improved parameter estimates in the  $r$ -dimensional subspaces. Not surprisingly, for  $n_i = 25$ ,  $i = 1, 2$ , we found that  $[EPMC_{10} - EPMC_1(SY)] \approx 0.06$  and for  $n_i = 50$ ,  $[EPMC_{10} - EPMC_1(SY)] \approx 0.01$ , which demonstrated the advantage of employing LDR in the classification process, and, more specifically, demonstrated the value of the SY LDR method. Additionally, as  $n_i$  increased, the  $EPMC_r(\cdot)$  approached  $EPMC_p(\cdot)$  for all four LDR methods as  $r$  increased.

In addition, the SIR LDR method did not utilize discriminatory information contained in the differences of the covariance matrices when  $n_i = 25$ ,  $i = 1, 2$ . However, for  $n_i = 50$ ,  $i = 1, 2$ , all four LDR methods yielded essentially the same EPMC, except for SAVE when  $r = 1$ .

### 7.2. Configuration 2: $m = 3$ with Two Similar Covariance Matrices and One Spherical Covariance Matrix

The second Monte Carlo simulation used a configuration with the three multivariate normal populations  $N_p(\mu_1, \Sigma_1)$ ,  $N_p(\mu_2, \Sigma_2)$ , and  $N_p(\mu_3, \Sigma_3)$  where  $p = 10$ ,

$$\mu_1 = [0, 0, 0, 0, 0, 0, 0, 0, 0, 0]',$$

$$\mu_2 = [1, 1, 1, 1, 1, 1, 1, 1, 1, 1]',$$

$$\mu_3 = [2, 2, 2, 2, 2, 2, 2, 2, 2, 2]',$$

Table 2. Summary of singular values for the four competing LDR methods for Configuration 1.

Singular value	$n_i = 25$				$n_i = 50$			
	$\hat{\mathbf{M}}$	$\hat{\mathbf{U}}$	$\hat{\mathbf{K}}$	$\hat{\mathbf{L}}$	$\hat{\mathbf{M}}$	$\hat{\mathbf{U}}$	$\hat{\mathbf{K}}$	$\hat{\mathbf{L}}$
$e_1$	72.16	0.38	0.57	0.74	67.98	0.19	0.52	0.66
$e_2$	16.27			0.54	11.67			0.36
$e_3$	12.21			0.43	8.93			0.25
$e_4$	9.31			0.34	6.88			0.18
$e_5$	7.22			0.25	5.34			0.13
$e_6$	5.46			0.18	4.06			0.09
$e_7$	3.99			0.11	2.98			0.05
$e_8$	2.73			0.07	2.05			0.03
$e_9$	1.67			0.03	1.23			0.01
$e_{10}$	0.75			0.01	0.53			<0.01

$$\Sigma_1 = \begin{bmatrix} 1.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 \\ 0.00 & 1.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 \\ 0.00 & 0.00 & 1.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 \\ 0.00 & 0.00 & 0.00 & 1.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 \\ 0.00 & 0.00 & 0.00 & 0.00 & 1.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 \\ 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 1.00 & 0.00 & 0.00 & 0.00 & 0.00 \\ 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 1.00 & 0.00 & 0.00 & 0.00 \\ 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 1.00 & 0.00 & 0.00 \\ 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 1.00 & 0.00 \\ 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 1.00 \end{bmatrix},$$

$$\Sigma_2 = \begin{bmatrix} 2.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 \\ 1.00 & 2.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 \\ 1.00 & 1.00 & 2.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 \\ 1.00 & 1.00 & 1.00 & 2.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 \\ 1.00 & 1.00 & 1.00 & 1.00 & 2.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 \\ 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 2.00 & 1.00 & 1.00 & 1.00 & 1.00 \\ 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 2.00 & 1.00 & 1.00 & 1.00 \\ 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 2.00 & 1.00 & 1.00 \\ 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 2.00 & 1.00 \\ 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 2.00 \end{bmatrix},$$

and

$$\Sigma_3 = \begin{bmatrix} 2.00 & 1.00 & 0.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 \\ 1.00 & 2.00 & 0.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 \\ 0.00 & 0.00 & 10.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 \\ 1.00 & 1.00 & 0.00 & 2.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 \\ 1.00 & 1.00 & 0.00 & 1.00 & 2.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 \\ 1.00 & 1.00 & 0.00 & 1.00 & 1.00 & 2.00 & 1.00 & 1.00 & 1.00 & 1.00 \\ 1.00 & 1.00 & 0.00 & 1.00 & 1.00 & 1.00 & 2.00 & 1.00 & 1.00 & 1.00 \\ 1.00 & 1.00 & 0.00 & 1.00 & 1.00 & 1.00 & 1.00 & 2.00 & 1.00 & 1.00 \\ 1.00 & 1.00 & 0.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 2.00 & 1.00 \\ 1.00 & 1.00 & 0.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 1.00 & 2.00 \end{bmatrix}.$$

In this configuration, the population means are unequal but relatively close. Moreover,  $\Sigma_2$  and  $\Sigma_3$  are unequal but notably more similar to one another than to  $\Sigma_1$ . However, the variance of the third feature in  $\Pi_3$  is significantly greater than the population variances of the third feature for either  $\Pi_1$  or  $\Pi_2$ .

As a result of markedly different covariance matrices, the *BE* and *SIR LDR* methods are not ideal because both methods aggregated the sample covariance matrices. The *SY LDR* method, however, attempts to estimate each individual covariance matrix for all three populations, and uses this information that yielded *SY* as the superior *LDR* procedure. **Table 3** gives the singular values of each of the *LDR* methods considered here.

For the larger sample-size scenario,  $n_i = 50$ ,  $i = 1, 2$ , the *SAVE LDR* was more competitive with *SY* than either *BE* or *SIR* when  $r > 1$ . However, the *SY LDR* remained the preferred *LDR* method. In **Figure 4**, we added noise to the *SY* method when we used  $r > 2$ , and, we eliminated essential discriminatory information for the *SY* procedure when we chose  $r = 1$ . Thus, the optimal choice for the reduced dimension for the *SY LDR* method was  $r = 2$ , regardless of the training-sample size of  $n_i$ ,  $i = 1, 2, 3$ . Furthermore, in our simulation we determined

**Table 3.** Summary of singular values for the four competing *LDR* methods applied to Configuration 2.

Singular value	$n_i = 25$				$n_i = 50$			
	$\hat{M}$	$\hat{U}$	$\hat{K}$	$\hat{L}$	$\hat{M}$	$\hat{U}$	$\hat{K}$	$\hat{L}$
$e_1$	15.18	1.31	0.55	1.25	14.27	1.15	0.51	1.18
$e_2$	8.92	0.90	0.11	0.61	8.89	0.47	0.06	0.43
$e_3$	4.07	0.43		0.49	2.09	0.29		0.28
$e_4$	2.51			0.40	1.50			0.21
$e_5$	1.84			0.32	1.21			0.16
$e_6$	1.49			0.25	1.01			0.13
$e_7$	1.22			0.19	0.84			0.10
$e_8$	1.00			0.14	0.70			0.07
$e_9$	0.79			0.10	0.56			0.05
$e_{10}$	0.58			0.06	0.42			0.03

that  $[EPMC_{10}(SY) - EPMC_2(SY)] \approx 0.10$  for  $n_i = 50$  and  $[EPMC_{10}(SY) - EPMC_2(SY)] \approx 0.15$  for  $n_i = 25$ ,  $i = 1, 2, 3$ . This significant reduction from  $EPMC_{10}$  demonstrated a significant benefit of dimension reduction in general and the *SY LDR* in particular.

### 7.3. Configuration 3: $m = 2$ with Relatively Close Means and Different But Similar Covariance Matrices

In this configuration, we have  $N_p(\mu_1, \Sigma_1)$  and  $N_p(\mu_2, \Sigma_2)$  with  $p = 10$ ,

$$\mu_1 = [0, 0, 0, 0, 0, 0, 0, 0, 0, 0]',$$

$$\mu_2 = [2, 2, 2, 2, 2, 2, 2, 2, 2, 2]',$$

$$\Sigma_1 = \begin{bmatrix} 10.00 & 4.00 & 5.00 & 4.00 & 3.00 & 4.00 & 4.00 & 5.00 & 4.00 & 3.00 \\ 4.00 & 10.00 & 5.00 & 2.00 & 4.00 & 3.00 & 3.00 & 5.00 & 4.00 & 3.00 \\ 5.00 & 5.00 & 10.00 & 5.00 & 5.00 & 3.00 & 4.00 & 4.00 & 4.00 & 4.00 \\ 4.00 & 2.00 & 5.00 & 10.00 & 3.00 & 4.00 & 2.00 & 3.00 & 4.00 & 3.00 \\ 3.00 & 4.00 & 5.00 & 3.00 & 12.00 & 3.00 & 4.00 & 5.00 & 3.00 & 3.00 \\ 4.00 & 3.00 & 3.00 & 4.00 & 3.00 & 9.00 & 3.00 & 4.00 & 4.00 & 4.00 \\ 4.00 & 3.00 & 4.00 & 2.00 & 4.00 & 3.00 & 14.00 & 2.00 & 2.00 & 2.00 \\ 5.00 & 5.00 & 4.00 & 3.00 & 5.00 & 4.00 & 2.00 & 12.00 & 1.00 & -0.50 \\ 4.00 & 4.00 & 4.00 & 4.00 & 3.00 & 4.00 & 2.00 & 1.00 & 14.00 & -1.00 \\ 3.00 & 3.00 & 4.00 & 3.00 & 3.00 & 4.00 & 2.00 & -0.50 & -1.00 & 11.00 \end{bmatrix},$$

and

$$\Sigma_2 = \begin{bmatrix} 8.00 & 2.00 & 3.00 & 4.00 & 1.00 & 2.00 & 2.00 & 3.00 & 2.00 & 1.00 \\ 2.00 & 8.00 & 3.00 & 2.00 & 2.00 & 1.00 & 1.00 & 3.00 & 2.00 & 1.00 \\ 3.00 & 3.00 & 8.00 & 5.00 & 3.00 & 1.00 & 2.00 & 2.00 & 2.00 & 2.00 \\ 4.00 & 2.00 & 5.00 & 10.00 & 3.00 & 4.00 & 2.00 & 3.00 & 4.00 & 3.00 \\ 1.00 & 2.00 & 3.00 & 3.00 & 10.00 & 1.00 & 2.00 & 3.00 & 1.00 & 1.00 \\ 2.00 & 1.00 & 1.00 & 4.00 & 1.00 & 7.00 & 1.00 & 2.00 & 2.00 & 2.00 \\ 2.00 & 1.00 & 2.00 & 2.00 & 2.00 & 1.00 & 12.00 & 0.00 & 0.00 & 0.00 \\ 3.00 & 3.00 & 2.00 & 3.00 & 3.00 & 2.00 & 0.00 & 10.00 & -1.00 & -2.50 \\ 2.00 & 2.00 & 2.00 & 4.00 & 1.00 & 2.00 & 0.00 & -1.00 & 12.00 & -3.00 \\ 1.00 & 1.00 & 2.00 & 3.00 & 1.00 & 2.00 & 0.00 & -2.50 & -3.00 & 9.00 \end{bmatrix}.$$

As in Configuration 1, we have that  $rank(\Sigma_2 - \Sigma_1) = 1$ , and, hence,  $rank(\mathbf{M}) = 2$ . In Configuration 3, the *BE* and *SIR LDR* methods outperformed the *SY* and *SAVE LDR* methods because of the similarity in the covariance matrices. This phenomenon occurred because both *LDR* methods aggregated the sample covariance matrices, which resulted in less-variable covariance matrix estimators and, therefore, smaller values of  $EPMC_1(BE)$  and  $EPMC_1(SIR)$ . From Table 4, we see that the first singular value for  $\hat{\mathbf{M}}$  was considerably less predominant than the first singular value for  $\hat{\mathbf{M}}$  in Configuration 1. This result explained the inferiority of the *SY* method for this configuration, regardless of the chosen  $n_i$ ,  $i = 1, 2$ .

For  $r = 1$ , we have that  $[EPMC_{10} - EPMC_1(SIR)] \approx [EPMC_{10} - EPMC_1(BE)] \approx 0.04$  for  $n_i = 25$  and  $[EPMC_{10} - EPMC_1(SIR)] \approx [EPMC_{10} - EPMC_1(BE)] \approx 0.06$  for  $n_i = 50$ . Again, we exemplify the value of *LDR* as a classification tool. In this configuration, not only was  $EPMC_{10}$  considerably reduced as  $n_i$  increased, but the difference between  $[EPMC_{10} - EPMC_1(BE)]$  and  $[EPMC_{10} - EPMC_1(SIR)]$  decreased as well. This example demonstrated the fact that the *SY* method is not a uniformly superior *LDR* approach even when covariance matrices are unequal. However,  $EPMC_1(SY)$  was not much greater than either  $EPMC_1(BE)$  or  $EPMC_1(SIR)$ . Here, *SAVE* is not as competitive as the three other *LDR* methods because it does not use all the information in the difference of the covariance matrices and means to obtain a better information-preserving subspace.

#### 7.4. Configuration 4: $m = 3$ with Two Similar Covariance Matrices Except for the First Two Dimensions

In this situation, we have three multivariate normal populations:  $N_p(\boldsymbol{\mu}_1, \Sigma_1)$ ,  $N_p(\boldsymbol{\mu}_2, \Sigma_2)$ , and  $N_p(\boldsymbol{\mu}_3, \Sigma_3)$ ,

**Table 4.** Summary of singular values for the four competing *LDR* methods for Configuration 3.

Singular value	$n_i = 25$				$n_i = 50$			
	$\hat{\mathbf{M}}$	$\hat{\mathbf{U}}$	$\hat{\mathbf{K}}$	$\hat{\mathbf{L}}$	$\hat{\mathbf{M}}$	$\hat{\mathbf{U}}$	$\hat{\mathbf{K}}$	$\hat{\mathbf{L}}$
$e_1$	27.91	0.45	0.42	0.61	22.98	0.370	0.35	0.44
$e_2$	14.84			0.49	10.48			0.28
$e_3$	10.73			0.39	7.69			0.21
$e_4$	8.12			0.30	5.86			0.15
$e_5$	6.18			0.23	4.50			0.11
$e_6$	4.65			0.16	3.39			0.08
$e_7$	3.41			0.10	2.50			0.50
$e_8$	2.41			0.06	1.78			0.03
$e_9$	1.55			0.03	1.16			0.01
$e_{10}$	0.76			0.01	0.60			<0.01



where  $p = 10$ ,

$$\boldsymbol{\mu}_1 = [0, 0, 0, 0, 0, 0, 0, 0, 0, 0]^T,$$

$$\boldsymbol{\mu}_2 = [1, 0, 1, 0, 1, 0, 1, 0, 1, 0]^T,$$

$$\boldsymbol{\mu}_3 = [-1, -1, -1, -1, -1, -1, -1, -1, -1, -1]^T,$$

$$\boldsymbol{\Sigma}_1 = \begin{bmatrix} 2.00 & 0.80 & 1.00 & 0.80 & 0.60 & 0.80 & 0.80 & 1.00 & 0.80 & 0.60 \\ 0.80 & 2.00 & 1.00 & 0.40 & 0.80 & 0.60 & 0.60 & 1.00 & 0.80 & 0.60 \\ 1.00 & 1.00 & 2.00 & 1.00 & 1.00 & 0.60 & 0.80 & 0.80 & 0.80 & 0.80 \\ 0.80 & 0.40 & 1.00 & 2.00 & 0.60 & 0.80 & 0.40 & 0.60 & 0.80 & 0.60 \\ 0.60 & 0.80 & 1.00 & 0.60 & 2.00 & 0.60 & 0.80 & 1.00 & 0.60 & 0.60 \\ 0.80 & 0.60 & 0.60 & 0.80 & 0.60 & 2.40 & 0.60 & 0.80 & 0.80 & 0.80 \\ 0.80 & 0.60 & 0.80 & 0.40 & 0.80 & 0.60 & 2.80 & 0.40 & 0.40 & 0.40 \\ 1.00 & 1.00 & 0.80 & 0.60 & 1.00 & 0.80 & 0.40 & 2.40 & -0.10 & -0.10 \\ 0.80 & 0.80 & 0.80 & 0.80 & 0.60 & 0.80 & 0.40 & -0.10 & 2.80 & -0.20 \\ 0.60 & 0.60 & 0.80 & 0.60 & 0.60 & 0.80 & 0.40 & -0.10 & -0.20 & 2.20 \end{bmatrix},$$

$$\boldsymbol{\Sigma}_2 = \begin{bmatrix} 20.00 & 0.80 & 1.00 & 0.80 & 0.60 & 0.80 & 0.80 & 1.00 & 0.80 & 0.60 \\ 0.80 & 40.00 & 1.00 & 0.40 & 0.80 & 0.60 & 0.60 & 1.00 & 0.80 & 0.60 \\ 1.00 & 1.00 & 2.00 & 1.00 & 1.00 & 0.60 & 0.80 & 0.80 & 0.80 & 0.80 \\ 0.80 & 0.40 & 1.00 & 2.00 & 0.60 & 0.80 & 0.40 & 0.60 & 0.80 & 0.60 \\ 0.60 & 0.80 & 1.00 & 0.60 & 2.00 & 0.60 & 0.80 & 1.00 & 0.60 & 0.60 \\ 0.80 & 0.60 & 0.60 & 0.80 & 0.60 & 2.40 & 0.60 & 0.80 & 0.80 & 0.80 \\ 0.80 & 0.60 & 0.80 & 0.40 & 0.80 & 0.60 & 2.80 & 0.40 & 0.40 & 0.40 \\ 1.00 & 1.00 & 0.80 & 0.60 & 1.00 & 0.80 & 0.40 & 2.40 & -0.10 & -0.10 \\ 0.80 & 0.80 & 0.80 & 0.80 & 0.60 & 0.80 & 0.40 & -0.10 & 2.80 & -0.20 \\ 0.60 & 0.60 & 0.80 & 0.60 & 0.60 & 0.80 & 0.40 & -0.10 & -0.20 & 2.20 \end{bmatrix},$$

and

$$\boldsymbol{\Sigma}_3 = \begin{bmatrix} 35.00 & 0.80 & 1.00 & 0.80 & 0.60 & 0.80 & 0.80 & 1.00 & 0.80 & 0.60 \\ 0.80 & 40.00 & 1.00 & 0.40 & 0.80 & 0.60 & 0.60 & 1.00 & 0.80 & 0.60 \\ 1.00 & 1.00 & 2.00 & 1.00 & 1.00 & 0.60 & 0.80 & 0.80 & 0.80 & 0.80 \\ 0.80 & 0.40 & 1.00 & 2.00 & 0.60 & 0.80 & 0.40 & 0.60 & 0.80 & 0.60 \\ 0.60 & 0.80 & 1.00 & 0.60 & 2.00 & 0.60 & 0.80 & 1.00 & 0.60 & 0.60 \\ 0.80 & 0.60 & 0.60 & 0.80 & 0.60 & 2.40 & 0.60 & 0.80 & 0.80 & 0.80 \\ 0.80 & 0.60 & 0.80 & 0.40 & 0.80 & 0.60 & 2.80 & 0.40 & 0.40 & 0.40 \\ 1.00 & 1.00 & 0.80 & 0.60 & 1.00 & 0.80 & 0.40 & 2.40 & -0.10 & -0.10 \\ 0.80 & 0.80 & 0.80 & 0.80 & 0.60 & 0.80 & 0.40 & -0.10 & 2.80 & -0.20 \\ 0.60 & 0.60 & 0.80 & 0.60 & 0.60 & 0.80 & 0.40 & -0.10 & -0.20 & 2.20 \end{bmatrix}.$$

For the fourth configuration, the covariance matrices are considerably different from one another, which benefits both the *SY* and *SAVE LDR* methods. Specifically, the *SY LDR* procedure uses information contained in the unequal covariance matrices to determine classificatory information contained in  $S_j - S_1$  and in  $S_j^{-1}\bar{x}_j - S_1^{-1}\bar{x}_1$ ,  $j = 2, 3$ . The largest  $EPMC_r$  for all four *LDR* methods occurred when  $r = 1$ . Generally, the  $EPMC_r$  de-

creased as  $r$  increased. For the  $SY$  method, we have  $\text{rank}(\mathbf{M}) = 4$ , which clarifies the reason that  $EPMC_r(SY)$  and  $r$  were inversely related. As one can see in **Table 5**, the first three singular values of  $\hat{\mathbf{M}}$  were relatively large. However, the pooling of the sample covariance matrices used in  $BE$  and  $SIR$  tended to obscure classificatory information in the sample covariance matrices. The only improvement from the full dimension we found for the values of  $r$  and  $n_i$  considered here was the  $SY$  method when  $r = 3$  for  $n_i = 25$ ,  $i = 1, 2, 3$ . Specifically, we have that  $[EPMC_{10} - EPMS_3(SY)] \approx 0.05$  when  $n_i = 25$ , although  $[EPMC_{10} - EPMS_3(SY)] \approx -0.01$  when  $n_i = 50$ ,  $i = 1, 2, 3$ .

This population configuration illustrated the fact that we cannot always choose  $r = 1$  and expect to see a reduction in  $EPMS_p$ . However, we can often reduce the  $EPMS_p$  if we use both a judicious choice of  $r$  and an appropriate  $LDR$  method.

### 7.5. Configuration 5: $m = 3$ with Diverse Population Covariance Matrices

In Configuration 5, we have three multivariate normal populations:  $N_p(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$ ,  $N_p(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$ , and  $N_p(\boldsymbol{\mu}_3, \boldsymbol{\Sigma}_3)$ , where  $p = 10$ ,

$$\boldsymbol{\mu}_1 = [0, 0, 0, 0, 0, 0, 0, 0, 0, 0]',$$

$$\boldsymbol{\mu}_2 = [4.43, 4.43, 4.43, 4.43, 4.43, 4.43, 4.43, 4.43, 4.43, 4.43]',$$

$$\boldsymbol{\mu}_3 = [8, 8, 8, 8, 8, 8, 8, 8, 8, 8]',$$

with

$$\boldsymbol{\Sigma}_1 = \begin{bmatrix} 15.01 & 0.81 & 1.25 & 1.13 & -2.10 & -2.30 & 2.43 & -3.30 & -0.87 & -1.11 \\ 0.81 & 26.10 & 1.51 & -0.74 & 0.89 & 4.35 & 1.25 & 1.85 & -0.50 & -0.39 \\ 1.25 & 1.51 & 24.55 & -5.57 & 3.96 & 1.62 & -0.27 & 0.49 & -5.97 & 0.87 \\ 1.13 & -0.74 & -5.57 & 29.36 & -3.58 & -0.89 & 2.21 & -3.71 & -0.52 & -2.19 \\ -2.10 & 0.89 & 3.96 & -3.58 & 20.17 & 6.05 & -5.20 & 2.22 & -0.80 & -2.69 \\ -2.30 & 4.35 & 1.62 & -0.89 & 6.05 & 40.18 & -5.18 & 3.83 & 1.94 & 0.51 \\ 2.43 & 1.25 & -0.27 & 2.21 & -5.20 & -5.18 & 17.93 & -0.17 & -3.09 & 0.49 \\ -3.30 & 1.85 & 0.49 & -3.71 & 2.22 & 3.83 & -0.17 & 26.05 & -0.54 & 3.34 \\ -0.87 & -0.50 & -5.97 & -0.52 & -0.80 & 1.94 & -3.09 & -0.54 & 16.30 & -1.04 \\ -1.11 & -0.39 & 0.87 & -2.19 & -2.69 & 0.51 & 0.49 & 3.34 & -1.04 & 26.84 \end{bmatrix},$$

**Table 5.** Summary of singular values for the four competing  $LDR$  methods for Configuration 4.

Singular value	$n_i = 25$				$n_i = 50$			
	$\hat{\mathbf{M}}$	$\hat{\mathbf{U}}$	$\hat{\mathbf{K}}$	$\hat{\mathbf{L}}$	$\hat{\mathbf{M}}$	$\hat{\mathbf{U}}$	$\hat{\mathbf{K}}$	$\hat{\mathbf{L}}$
$e_1$	59.76	1.93	0.48	0.90	56.56	0.86	0.44	0.72
$e_2$	36.37	0.66	0.16	0.69	37.02	0.60	0.11	0.56
$e_3$	7.72	0.31		0.52	5.36	0.16		0.31
$e_4$	5.21			0.42	3.53			0.24
$e_5$	3.97			0.34	2.65			0.19
$e_6$	3.12			0.27	2.09			0.14
$e_7$	2.46			0.21	1.66			0.11
$e_8$	1.91			0.16	1.31			0.08
$e_9$	1.44			0.11	1.00			0.06
$e_{10}$	0.98			0.07	0.70			0.03

$$\Sigma_2 = \begin{bmatrix} 45.01 & 30.81 & 31.25 & 31.13 & 27.90 & 27.70 & 32.43 & 26.70 & 29.13 & 28.89 \\ 30.81 & 56.10 & 31.51 & 29.26 & 30.89 & 34.35 & 31.25 & 31.85 & 29.50 & 29.61 \\ 31.25 & 31.51 & 54.55 & 24.43 & 33.96 & 31.62 & 29.73 & 30.49 & 24.03 & 30.87 \\ 31.13 & 29.26 & 24.43 & 59.36 & 26.42 & 29.11 & 32.21 & 26.29 & 29.48 & 27.81 \\ 27.90 & 30.89 & 33.96 & 26.42 & 50.17 & 36.05 & 24.80 & 32.22 & 29.20 & 27.31 \\ 27.70 & 34.35 & 31.62 & 29.11 & 36.05 & 70.18 & 24.82 & 33.83 & 31.94 & 30.51 \\ 32.43 & 31.25 & 29.73 & 32.21 & 24.80 & 24.82 & 47.93 & 29.83 & 26.91 & 30.49 \\ 26.70 & 31.85 & 30.49 & 26.29 & 32.22 & 33.83 & 29.83 & 56.05 & 29.46 & 33.34 \\ 29.13 & 29.50 & 24.03 & 29.48 & 29.20 & 31.94 & 26.91 & 29.46 & 46.30 & 28.96 \\ 28.89 & 29.61 & 30.87 & 27.81 & 27.31 & 30.51 & 30.49 & 33.34 & 28.96 & 56.84 \end{bmatrix},$$

and

$$\Sigma_3 = \begin{bmatrix} 75.01 & 60.81 & 61.25 & 61.13 & 57.90 & 57.70 & 62.43 & 56.70 & 59.13 & 58.89 \\ 60.81 & 86.10 & 61.51 & 59.26 & 60.89 & 64.35 & 61.25 & 61.85 & 59.50 & 59.61 \\ 61.25 & 61.51 & 84.55 & 54.43 & 63.96 & 61.62 & 59.73 & 60.49 & 54.03 & 60.87 \\ 61.13 & 59.26 & 54.43 & 89.36 & 56.42 & 59.11 & 62.21 & 56.29 & 59.48 & 57.81 \\ 57.90 & 60.89 & 63.96 & 56.42 & 80.17 & 66.05 & 54.80 & 62.22 & 59.20 & 57.31 \\ 57.70 & 64.35 & 61.62 & 59.11 & 66.05 & 100.18 & 54.82 & 63.83 & 61.94 & 60.51 \\ 62.43 & 61.25 & 59.73 & 62.21 & 54.80 & 54.82 & 77.93 & 59.83 & 56.91 & 60.49 \\ 56.70 & 61.85 & 60.49 & 56.29 & 62.22 & 63.83 & 59.83 & 86.05 & 59.46 & 63.34 \\ 59.13 & 59.50 & 54.03 & 59.48 & 59.20 & 61.94 & 56.91 & 59.46 & 76.30 & 58.96 \\ 58.89 & 59.61 & 60.87 & 57.81 & 57.31 & 60.51 & 60.49 & 63.34 & 58.96 & 86.84 \end{bmatrix}.$$

Here, we have three considerably different covariance matrices. For this population configuration, the *SY* method outperformed the three other *LDR* methods for the reduced dimensions  $r=1,2,3$ . Examining the singular values of  $\hat{M}$ , we see that most of the discriminatory information was contained in the first transformed dimension. This fact was illustrated with how  $EPMC_r(SY)$  and  $r$  were directly related, regardless of the training-sample size  $n_i$ ,  $i=1,2,3$ .

The *BE* and *SIR LDR* methods did not perform as well here because of the pooling of highly diverse estimated covariance matrices. While the *SAVE LDR* procedure was the least effective *LDR* method for  $n_i=25$ ,  $EPMC(SAVE)$  and  $n_i$  were inversely related. While each of the four *LDR* methods decreased from  $EPMC_{10}$  for some combination of  $n_i$  and  $r$ , the largest reduction in  $EMPC$  was for the *SY LDR* method where  $[EPMC_{10} - EPMC_1(SY)] \approx 0.09$  for  $n_i=50$  and  $[EPMC_{10} - EPMC_1(SY)] \approx 0.14$  for  $n_i=25$ ,  $i=1,2,3$ . This result implied that the *SY LDR* method was especially useful when  $n_i/p$  was relatively small and the covariance matrices were considerably different.

## 8. A Parametric Bootstrap Simulation

In the following parametric bootstrap simulation, we use a real dataset to obtain the population means and covariance matrices for three multivariate normal populations. The chosen dataset comes from the University of California at Irvine Machine Learning Repository, which describes the diagnoses of cardiac Single Proton Emission Computed Tomography (*SPECT*) images. Each patient in the study is classified into two categories: normal or abnormal. The dataset contains 267 *SPECT* image sets of patients. Each observation consists of 44 continuous features for each patient. However, for our simulation, we chose only ten of the 44 features. The ten selected features were F2R, F6R, F7S, F9S, F11R, F11S, F14S, F16R, F17R, and F19S. Hence, we performed the parametric bootstrap Monte Carlo simulation with two populations:  $N(\mu_1, \Sigma_1)$  and  $N(\mu_2, \Sigma_2)$ , where

$$\mu_1 = [61.68, 65.85, 70.37, 70.48, 66.01, 68.21, 60.18, 73.78, 71.68, 63.66]^T,$$

**Table 6.** Summary of singular values for the four competing *LDR* methods for Configuration 6.

Singular value	$n_i = 25$				$n_i = 50$			
	$\hat{M}$	$\hat{U}$	$\hat{K}$	$\hat{L}$	$\hat{M}$	$\hat{U}$	$\hat{K}$	$\hat{L}$
$e_1$	691.33	0.37	0.36	0.72	683.07	0.13	0.30	0.53
$e_2$	71.16	0.14	0.11	0.56	50.37	0.07	0.06	0.31
$e_3$	50.13	0.08		0.46	35.78	0.04		0.24
$e_4$	39.44			0.37	28.39			0.19
$e_5$	32.07			0.30	23.19			0.15
$e_6$	26.27			0.24	19.10			0.12
$e_7$	21.43			0.18	15.67			0.09
$e_8$	17.22			0.14	12.13			0.07
$e_9$	13.41			0.09	9.89			0.05
$e_{10}$	9.59			0.06	7.08			0.03

$$\mu_2 = [62.60, 67.90, 71.60, 71.60, 71.00, 69.40, 62.30, 75.90, 72.90, 66.00]'$$

$$\Sigma_1 = \begin{bmatrix} 131.90 & 62.83 & 6.01 & -1.82 & -6.75 & 0.79 & -16.74 & 3.69 & 30.97 & 9.27 \\ 62.83 & 127.44 & 0.44 & 21.72 & -6.27 & 10.91 & -12.34 & 4.41 & 4.48 & 26.22 \\ 6.01 & 0.44 & 42.82 & 20.46 & 27.11 & 21.93 & 26.23 & 16.45 & 16.68 & 14.73 \\ -1.82 & 21.72 & 20.46 & 43.16 & 13.24 & 22.12 & 12.11 & 15.18 & 10.83 & 17.75 \\ -6.75 & -6.27 & 27.11 & 13.24 & 71.60 & 15.62 & 20.40 & 13.33 & 26.02 & 3.17 \\ 0.79 & 10.91 & 21.93 & 22.12 & 15.62 & 48.53 & 15.54 & 8.80 & 6.24 & 13.60 \\ -16.74 & -12.34 & 26.23 & 12.11 & 20.40 & 15.54 & 49.80 & 12.69 & 3.44 & -0.59 \\ 3.69 & 4.41 & 16.45 & 15.18 & 13.33 & 8.80 & 12.69 & 47.52 & 25.45 & 5.78 \\ 30.97 & 4.48 & 16.68 & 10.83 & 26.02 & 6.24 & 3.44 & 25.45 & 62.92 & 21.23 \\ 9.27 & 26.22 & 14.73 & 17.75 & 3.17 & 13.60 & -0.59 & 5.78 & 21.23 & 87.24 \end{bmatrix},$$

and

$$\Sigma_2 = \begin{bmatrix} 33.54 & 18.26 & -1.10 & 0.69 & 3.29 & -10.69 & 10.57 & -2.99 & 9.97 & 13.07 \\ 18.26 & 27.35 & -5.60 & -1.60 & -9.43 & -7.19 & 2.95 & -0.94 & 12.64 & 24.29 \\ -1.10 & -5.60 & 9.11 & 3.83 & 0.79 & 0.46 & 7.71 & -3.77 & -3.31 & -8.71 \\ 0.69 & -1.60 & 3.83 & 23.26 & 8.36 & 10.46 & -4.50 & 4.59 & 10.97 & 11.93 \\ 3.29 & -9.43 & 0.79 & 8.36 & 18.43 & 10.29 & -0.93 & 6.64 & 10.79 & -8.50 \\ -10.69 & -7.19 & 0.46 & 10.46 & 10.29 & 29.11 & -12.57 & 4.27 & 4.03 & -5.50 \\ 10.57 & 2.95 & 7.71 & -4.50 & -0.93 & -12.57 & 28.52 & -3.88 & -6.05 & -7.57 \\ -2.99 & -0.94 & -3.77 & 4.59 & 6.64 & 4.27 & -3.88 & 19.41 & 17.56 & 9.50 \\ 9.97 & 12.64 & -3.31 & 10.97 & 10.79 & 4.03 & -6.05 & 17.56 & 46.07 & 20.93 \\ 13.07 & 24.29 & -8.71 & 11.93 & -8.50 & -5.50 & -7.57 & 9.50 & 20.93 & 46.57 \end{bmatrix}.$$

For this dataset,  $rank(\mathbf{M}) = rank(\hat{\mathbf{M}}) = 10$ ,  $rank(\mathbf{U}) = rank(\hat{\mathbf{U}}) = 1$ ,  $rank(\mathbf{K}) = rank(\hat{\mathbf{K}}) = 2$ , and  $rank(\mathbf{L}) = rank(\hat{\mathbf{L}}) = 10$ . **Table 6** gives the singular values for each of the *LDR* methods considered here.

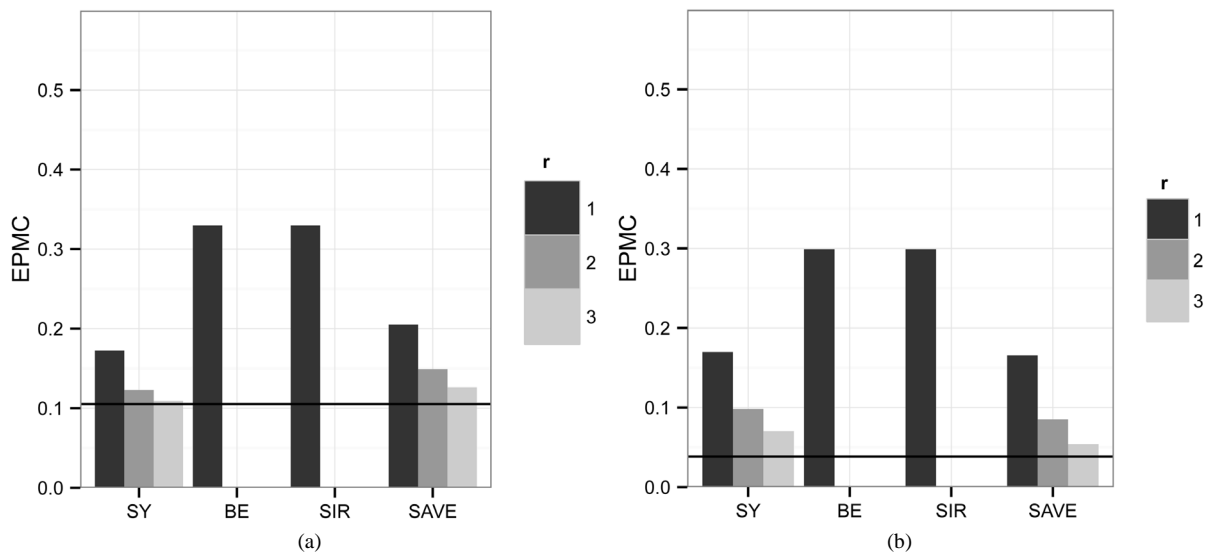
**Figure 8** gives the *EPMC*s for the reduced dimensions 1, 2, and 3 for each of the *LDR* methods. In this example, a considerable amount of discriminatory information is contained in the covariance matrices, which are ex-

tremely different. Hence, not surprisingly, neither the *BE* nor the *SIR LDR* methods performed well for  $r=1,2,3$ . Furthermore,  $EPMC_r(SY)$ ,  $EPMC_r(SIR)$ , and  $EPMC_r(SAVE)$  were inversely related to  $r$ , regardless of the training-sample size  $n_i$ ,  $i=1,2$ . From **Table 7**, we see that one reason that  $EPMC_r(SY)$  increased as  $r$  decreased is that the singular values are relatively large up to  $r=7$ . Thus, when we chose  $r=1,2,3$ , we discarded some necessary discriminatory information.

For  $n_i=50$ ,  $i=1,2$ , *SAVE* slightly outperformed the other three *LDR* methods, though the *SY LDR* approach was very competitive. For  $n_i=25$ ,  $i=1,2$ , we found that  $SY EPMC_3(SY)$  was the smallest error rate and that  $EPMC_{10} \approx EPMC_3(SY)$ . This particular example verified that one can implement the *SY* method and obtain excellent results even though the conditions of Theorem 1 do not essentially hold. This example also illustrated the prospect that, with a judicious choice of  $r$  combined with the appropriate *LDR* method, we can significantly reduce the feature dimension while still preserving the *EMPC*.

**Table 7.** Summary of singular values for the four competing *LDR* methods for the parametric bootstrap example.

Singular value	$n_i = 25$				$n_i = 50$			
	$\hat{M}$	$\hat{U}$	$\hat{K}$	$\hat{L}$	$\hat{M}$	$\hat{U}$	$\hat{K}$	$\hat{L}$
$e_1$	390.43	0.23	0.36	0.94	275.15	0.15	0.28	0.92
$e_2$	198.20			0.86	175.62			0.83
$e_3$	118.13			0.74	111.28			0.71
$e_4$	74.01			0.56	69.53			0.50
$e_5$	52.16			0.42	47.01			0.31
$e_6$	37.65			0.30	35.02			0.21
$e_7$	25.72			0.20	24.92			0.13
$e_8$	16.73			0.12	15.91			0.07
$e_9$	9.50			0.05	8.60			0.03
$e_{10}$	4.27			0.02	3.81			0.01



**Figure 8.** Simulation results from Section 8. The horizontal bar in each graph represents the  $EPMC_{10}$ . (a)  $n_i = 25$ ; (b)  $n_i = 50$ .

## 9. Discussion

In this paper, while all population parameters are known, we have presented a simple and flexible algorithm for a low-dimensional representation of data from multiple multivariate normal populations with different parametric configurations. Also, we have given necessary and sufficient conditions for attaining the subspace of smallest dimension  $q < p$ , which preserves the original Bayes classification assignments. We have provided a constructive proof for obtaining a low-dimensional representation space when certain population parameter conditions are satisfied. Under a special case for the two-class multivariate normal problem with equal nonsingular covariance structures, our proposed *LDR* transformation is the *LDF* in [4]. Moreover, we have also extended our concept proposed in Theorem 1 to cases where the conditions in Theorem 1 are not satisfied through the application of the *SVD* given in Theorem 2.

We have presented several advantages of our proposed low-dimensional representation method. First, our method is not restricted to a one-dimensional representation regardless of the number of populations, unlike the transformation introduced by [14]. Second, our method allows for equal and unequal covariance structures. Third, the original feature dimension  $p$  does not significantly impact the computational complexity. Also, under certain conditions, one-dimensional representations of populations with unequal covariance structures can be accomplished without an appreciable increase in the *EMPC*.

Furthermore, we have derived a *LDR* method for realistic cases where the class population parameters are unknown and the linear sufficient matrix  $M$  in (7) must be estimated using training data. Using Monte Carlo simulation studies, we have compared the performance of the *SY LDR* method with three other *LDR* procedures derived by [14] [16] [19]. In the Monte Carlo simulation studies, we have demonstrated that the full-dimension *EMPC* can sometimes be decreased by implementing the *SY LDR* method when the training-sample size is small relative to the total number of estimated parameters.

Finally, we have extended our concept proposed in Theorem 1 through the application of the *SVD* given in Theorem 2 to cases where one might wish to considerably reduce the original feature dimension. Our new *LDR* approach can yield excellent results provided the population covariance matrices are sufficiently different.

## References

- [1] Bellman, R. (1961) Adaptive Control Processes: A Guide Tour. John Wiley, Princeton University, Princeton.
- [2] Jain, A.K. and Chandrasekaran, B. (1982) Dimensionality and Sample Size Considerations in Pattern Recognition Practice. In: Krishnaiah, P.R. and Kanal, L.N., Eds., *Handbook of Statistics*, Volume 2, Elsevier, Amsterdam, 201-213. [http://dx.doi.org/10.1016/s0169-7161\(82\)02042-2](http://dx.doi.org/10.1016/s0169-7161(82)02042-2)
- [3] Rao, C.R. (1948) The Utilization of Multiple Measurements in Problems of Biological Classification. *Journal of the Royal Statistical Society: Series B*, **10**, 159-203.
- [4] Fisher, R.A. (1936) The Use of Multiple Measurements in Taxonomic Problems. *Annals of Eugenics*, **7**, 179-188. <http://dx.doi.org/10.1111/j.1469-1809.1936.tb02137.x>
- [5] McLachlan, G.J. (1992) Discriminant Analysis and Statistical Pattern Recognition. John Wiley, New York. <http://dx.doi.org/10.1002/0471725293>
- [6] Decell, H.P. and Mayekar, S.M. (1977) Maximum Likelihood from Incomplete Data via the EM Algorithm. *Journal of the Royal Statistical Society, Series B*, **39**, 1-38.
- [7] Okada, T. and Tomita, S. (1984) An Extended Fisher Criterion for Feature Extraction—Malina's Method and Its Problems. *Electronics and Communications in Japan*, **67**, 10-16. <http://dx.doi.org/10.1002/ecja.4400670603>
- [8] Loog, M. and Duin, P.W. (2004) Linear Dimensionality Reduction via a Heteroscedastic Extension of LDA: The Chernoff Criterion. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **26**, 732-739. <http://dx.doi.org/10.1109/TPAMI.2004.13>
- [9] Fukunaga, K. (1990) Introduction to Statistical Pattern Recognition. 2nd Edition, Academic Press, Boston.
- [10] Hennig, C. (2004) Asymmetric Linear Dimension Reduction for Classification. *Journal of Computational and Graphical Statistics*, **13**, 930-945. <http://dx.doi.org/10.1198/106186004X12740>
- [11] Kumar, N. and Andreou, A.G. (1996) A Generalization of Linear Discriminant Analysis in a Maximum Likelihood Framework. *Proceedings of the Joint Statistical Meeting*.
- [12] Young, D.M., Marco, V.R. and Odell, P.L. (1987) Quadratic Discrimination: Some Results on Optimal Low-Dimensional Representation. *The Journal of Statistical Planning and Inference*, **17**, 307-319.

- [http://dx.doi.org/10.1016/0378-3758\(87\)90122-4](http://dx.doi.org/10.1016/0378-3758(87)90122-4)
- [13] Peters, B.C., Redner, R. and Decell, H.P. (1978) Characterization of Linear Sufficient Statistics. *Sankhya*, **40**, 303-309.
- [14] Brunzell, H. and Eriksson, J. (2000) Feature Reduction for Classification of Multidimensional Data. *Pattern Recognition*, **33**, 1741-1748. [http://dx.doi.org/10.1016/S0031-3203\(99\)00142-9](http://dx.doi.org/10.1016/S0031-3203(99)00142-9)
- [15] Cook, R.D. and Weisberg, S. (1991) Discussion of “Sliced Inverse Regression for Dimension Reduction” by K.-C. Li. *Journal of the American Statistical Association*, **86**, 328-332.
- [16] Li, K.-C. (1991) Sliced Inverse Regression for Dimension Reduction. *Journal of the American Statistical Association*, **86**, 316-327. <http://dx.doi.org/10.1080/01621459.1991.10475035>
- [17] Tubbs, J.D., Coberly, W.A. and Young, D.M. (1982) Linear Dimension Reduction and Bayes Classification with Unknown Population Parameters. *Pattern Recognition*, **15**, 167-172. [http://dx.doi.org/10.1016/0031-3203\(82\)90068-1](http://dx.doi.org/10.1016/0031-3203(82)90068-1)
- [18] Schervish, M.J. (19884) Linear Discrimination for Three Known Normal Populations. *Journal of Statistical Planning and Inference*, **10**, 167-175. [http://dx.doi.org/10.1016/0378-3758\(84\)90068-5](http://dx.doi.org/10.1016/0378-3758(84)90068-5)
- [19] Cook, R.D. (2000) SAVE: A Method for Dimension Reduction and Graphics in Regression. *Communications in Statistics—Theory and Methods*, **29**, 2109-2121. <http://dx.doi.org/10.1080/03610920008832598>
- [20] Cook, R.D. and Yin, X. (1991) Dimension Reduction and Visualization in Discriminant Analysis. *Australian and New Zealand Journal of Statistics*, **43**, 147-199. <http://dx.doi.org/10.1111/1467-842X.00164>
- [21] Vellila, S. (2012) A Note on the Structure of the Quadratic Subspace in Discriminant Analysis. *Statistics and Probability Letters*, **82**, 739-747. <http://dx.doi.org/10.1016/j.spl.2011.12.020>