

# COMPRESSIVE REGULARIZED DISCRIMINANT ANALYSIS OF HIGH-DIMENSIONAL DATA WITH APPLICATIONS TO MICROARRAY STUDIES

Muhammad Naveed Tabassum and Esa Ollila

Aalto University, Dept. of Signal Processing and Acoustics, P.O. Box 15400, FI-00076 Aalto, Finland

## ABSTRACT

We propose a modification of linear discriminant analysis, referred to as *compressive regularized discriminant analysis (CRDA)*, for analysis of high-dimensional datasets. CRDA is especially designed for feature elimination purpose and can be used as gene selection method in microarray studies. CRDA lends ideas from  $\ell_{q,1}$  norm minimization algorithms in the multiple measurement vectors (MMV) model and utilizes joint-sparsity promoting hard thresholding for feature elimination. A regularization of the sample covariance matrix is also needed as we consider the challenging scenario where the number of features (variables) is comparable or exceeding the sample size of the training dataset. A simulation study and four examples of real life microarray datasets evaluate the performances of CRDA based classifiers. Overall, the proposed method gives fewer misclassification errors than its competitors, while at the same time achieving accurate feature selection.

**Index Terms**— Classification, gene expression microarrays, joint-sparse recovery, regularized discriminant analysis.

## 1. INTRODUCTION

Sparse signal approximations are widely used in many applications such as regression or classification where variable-selection (i.e., ranking and selection of features) aims at reducing the number of variables (or features) without sacrificing accuracy measured by the test error. Reduction in the set of features facilitates interpretation as well as stabilizes estimation. This is often deemed necessary in the high-dimensional (HD) context where the number of features,  $p$ , is often several magnitudes larger than the number of observations,  $n$ , in the training dataset (i.e.,  $p \gg n$ ).

Many classification techniques assign a  $p$ -dimensional observation  $\mathbf{x}$  to one of the  $G$  classes (groups or populations) based on the following rule

$$\mathbf{x} \in \text{group} \left[ \tilde{g} = \arg \max_g d_g(\mathbf{x}) \right], \quad (1)$$

where  $d_g(\mathbf{x})$  is called the *discriminant function* for population  $g \in \{1, \dots, G\}$ . In linear discriminant analysis (LDA),

$d_g(\mathbf{x})$  is a linear function of  $\mathbf{x}$ ,  $d_g(\mathbf{x}) = \mathbf{x}^\top \beta_g + c_g$ , for some constant  $c_g \in \mathbb{R}$  and vector  $\beta_g \in \mathbb{R}^p$ . The vector  $\beta_g = \beta_g(\Sigma)$  depends on the unknown covariance matrix  $\Sigma$  of the populations (via its inverse matrix) which is commonly estimated by the pooled sample covariance matrix (SCM). In the HD setting, the SCM is no-longer invertible, and therefore regularized SCM (RSCM)  $\hat{\Sigma}$  is used for constructing an estimated discriminant function  $\hat{d}_g(\mathbf{x})$ . Such approaches are commonly referred to as regularized LDA methods, which we refer shortly as RDA. See e.g., [1–5].

Next note that if the  $i$ -th entry of  $\beta_g$  is zero, then the  $i$ -th feature does not contribute in the classification of  $g$ -th population. To eliminate unnecessary features, many authors have proposed to shrink  $\beta_g$  using element-wise soft-thresholding, e.g., as in shrunken centroids (SC)RDA method [1]. These methods are often difficult to tune because the shrinkage threshold parameter is the same across all groups, but different populations would often benefit from different shrinkage intensity. Consequently, they tend to yield rather higher false-positive (FP) rates.

Element-wise shrinkage does not achieve *simultaneous* feature selection as the eliminated feature from group  $i$  may still affect the discriminant function of group  $j$ . In this paper, we propose *compressive regularized discriminant analysis (CRDA)* that promotes simultaneous *joint-sparsity* to pick fewer and differentially expressed variables. CRDA lends ideas from mixed  $\ell_{q,1}$  norm minimization in the multiple measurement vectors (MMV) model [6], which is an extension of compressed sensing model to multivariate case. CRDA uses  $\ell_{q,1}$ -norm based hard-thresholding which has the advantage of having a shrinkage parameter that is much easier to tune: namely, joint-sparsity level  $K \in \{1, \dots, p\}$  instead of shrinkage threshold  $\Delta \in [0, \infty)$  as in SCRDA. Our approach also employs a different RSCM estimator compared to SCRDA. The used RSCM has the benefit of being able to attain the minimum mean squared error [7, 8] for an appropriate choice of the regularization parameter. The optimal pair of the tuning parameters can be found via cross validation (CV), but we also propose a computationally simpler approach that uses the RSCM proposed in [8]. This facilitates the computations considerably as only a single variable, the joint-sparsity level  $K$ , needs to be tuned.

The paper is organized as follows. Section 2 describes

The research was partially supported by the Academy of Finland grant no. 298118 which is gratefully acknowledged.

RDA and SVD based inversion of the used RSCM. In Section 3, the proposed CRDA as well as our tuning parameter selection criteria is introduced. Section 4 provides the results on simulation studies which explore both the feature-selection capability and misclassification errors of CRDA, and the competing methods. Classification results on four real microarray datasets are also provided. Section 5 concludes the paper.

## 2. REGULARIZED LDA

We are given a  $p$ -variate random vector  $\mathbf{x}$  which we need to classify into one of the  $G$  classes or populations. In LDA, one assumes that the class populations are  $p$ -variate multivariate normal (MVN) with a common positive definite symmetric covariance matrix  $\Sigma$  over each class but distinct class mean vectors  $\mu_g \in \mathbb{R}^p$ ,  $g = 1, \dots, G$ . The problem is then to classify  $\mathbf{x}$  to one of the MVN populations,  $\mathcal{N}_p(\mu_g, \Sigma)$ ,  $g = 1, \dots, G$ . Sometimes prior knowledge is available on proportions of each population and we denote by  $p_g$ ,  $g = 1, \dots, G$ , the prior probabilities of the classes ( $\sum_{g=1}^G p_g = 1$ ). LDA uses the rule (1) with discriminant function

$$d_g(\mathbf{x}) = \mathbf{x}^\top \beta_g - \frac{1}{2} \mu_g^\top \beta_g + \ln p_g,$$

where  $\beta_g = \Sigma^{-1} \mu_g$  for  $g = 1, \dots, G$ .

The LDA rule involves a set of unknown parameters, the class mean vectors  $\mu_g$  and the covariance matrix  $\Sigma$ . These are estimated from the *training dataset*  $\mathbf{X} = (\mathbf{x}_1 \cdots \mathbf{x}_n)$  that consists of  $n_g$  observations from each of the classes ( $g = 1, \dots, G$ ). Let  $c(i)$  denote the class label associated with the  $i$ -th observation, so  $c(i) \in \{1, \dots, G\}$ . Then  $n_g = \sum_{i=1}^n \mathbb{1}(c(i) = g)$  is the number of observations belonging to  $g$ -th population, and we denote by  $\pi_g = n_g/n$  the relative sample proportions. We assume observations in the training dataset are centered by the sample mean vectors of the classes,

$$\hat{\mu}_g = \bar{\mathbf{x}}_g = \frac{1}{n_g} \sum_{c(i)=g} \mathbf{x}_i. \quad (2)$$

Since  $\mathbf{X}$  is centered, the pooled (over groups) sample covariance matrix (SCM) can be written simply as

$$\mathbf{S} = \frac{1}{n} \mathbf{X} \mathbf{X}^\top.$$

In practice, an observation  $\mathbf{x}$  is classified using an estimated discriminant function,

$$\hat{d}_g(\mathbf{x}) = \mathbf{x}^\top \hat{\beta}_g - \frac{1}{2} \hat{\mu}_g^\top \hat{\beta}_g + \ln \pi_g, \quad (3)$$

where  $\hat{\beta}_g = \hat{\Sigma}^{-1} \hat{\mu}_g$ ,  $g = 1, \dots, G$  and  $\hat{\Sigma}$  is an estimator of  $\Sigma$ . Note that in (3) the prior probabilities  $p_g$ -s are replaced by their estimates,  $\pi_g$ -s. Commonly, the pooled SCM  $\mathbf{S}$  is used as an estimator  $\hat{\Sigma}$ . Since we are in the regime, where  $p \gg n$ , the pooled SCM is no longer invertible and hence

can not be used in (3). To avoid the singularity of the estimated covariance matrix, a commonly used approach in the literature (cf. [7, 8]) is to use a *regularized SCM (RSCM)*,

$$\hat{\Sigma} = \alpha \mathbf{S} + (1 - \alpha) \eta \mathbf{I} \quad (4)$$

where  $\eta = \text{Tr}(\mathbf{S})/p$ . SCRDA [1] uses an estimator  $\hat{\Sigma} = \alpha \mathbf{S} + (1 - \alpha) \mathbf{I}$ . However, (4) has some theoretical justification since with an appropriate (data dependent) choice  $\hat{\alpha}$ , the obtained RSCM in (4) will be a consistent minimum mean squared error (MMSE) estimator of  $\Sigma$ . Such choices of  $\alpha$  have been proposed, e.g., in [7] and in [8].

In the HD setup, the main computational burden is related with inverting the matrix  $\hat{\Sigma}$  in (4). The inversion can be done using the SVD-trick, as follows [1, 9]. The SVD of  $\mathbf{X}$  is

$$\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}^\top,$$

where  $\mathbf{U} \in \mathbb{R}^{p \times m}$ ,  $\mathbf{D} \in \mathbb{R}^{m \times m}$ ,  $\mathbf{V} \in \mathbb{R}^{n \times m}$  and  $m = \text{rank}(\mathbf{X})$ . Direct computation of SVD is time consuming and the trick is that  $\mathbf{V}$  and  $\mathbf{D}$  can be computed first from SVD of  $\mathbf{X}^\top \mathbf{X} = \tilde{\mathbf{V}} \tilde{\mathbf{D}}^2 \tilde{\mathbf{V}}^\top$ , which is only an  $n \times n$  matrix. Here  $\tilde{\mathbf{V}}$  is an orthogonal  $n \times n$  matrix whose first  $m$ -columns are  $\mathbf{V}$  and  $\tilde{\mathbf{D}}$  is an  $n \times n$  diagonal matrix whose upper left corner  $m \times m$  matrix is  $\mathbf{D}$ . After we compute  $\mathbf{V}$  and  $\mathbf{D}$  from SVD of  $\mathbf{X}^\top \mathbf{X}$ , we may compute  $\mathbf{U}$  from  $\mathbf{X}$  by  $\mathbf{U} = \mathbf{X} \mathbf{V} \mathbf{D}^{-1}$ . Then, using the SVD representation of the SCM,  $\mathbf{S} = (1/n) \mathbf{U} \mathbf{D}^2 \mathbf{U}^\top$ , and simple algebra, one obtains a simple formula for the inverse:

$$\hat{\Sigma}^{-1} = \mathbf{U} \left[ \left( \frac{\alpha}{n} \mathbf{D}^2 + (1 - \alpha) \eta \mathbf{I}_m \right)^{-1} - \frac{1}{(1 - \alpha) \eta} \mathbf{I}_m \right] \mathbf{U}^\top + \frac{1}{(1 - \alpha) \eta} \mathbf{I}_p, \quad (5)$$

where  $\eta = \text{Tr}(\mathbf{S})/p = \text{Tr}(\mathbf{D}^2)/np$ . This reduces the complexity from  $\mathcal{O}(p^3)$  to  $\mathcal{O}(pn^2)$  which is a significant saving in  $p \gg n$  case.

## 3. COMPRESSIVE RDA

### 3.1. Proposed CRDA Approach

In order to explain the proposed compressive RDA approach, we first write the discriminant rule in vector form as

$$\begin{aligned} \mathbf{d}(\mathbf{x}) &= (d_1(\mathbf{x}), \dots, d_G(\mathbf{x})) \\ &= \mathbf{x}^\top \mathcal{B} - \frac{1}{2} \text{diag}(\mathbf{M}^\top \mathcal{B}) + \ln \mathbf{p}, \end{aligned} \quad (6)$$

where  $\ln \mathbf{p} = (\ln p_1, \dots, \ln p_G)$ ,  $\mathbf{M} = (\mu_1 \ \dots \ \mu_G)$  and  $\mathcal{B} = (\beta_1 \ \dots \ \beta_G) = \Sigma^{-1} \mathbf{M}$ . Above notation  $\text{diag}(\cdot)$  extract the diagonal of the  $G \times G$  matrix  $\mathbf{A}$  into a vector, i.e.,  $\text{diag}(\mathbf{A}) = (a_{11}, \dots, a_{GG})$ . The discriminant function in (6) is linear in  $\mathbf{x}$  with coefficient matrix  $\mathcal{B} \in \mathbb{R}^{p \times G}$ . This means that if the  $i$ -th row of the coefficient matrix  $\mathcal{B}$  is a zero vector  $\mathbf{0}$ , then it implies that  $i$ -th predictor does not contribute

to the classification rule and hence can be eliminated. If the coefficient matrix  $\mathcal{B}$  is row-sparse, then the method can be potentially used as a simultaneous feature elimination procedure. In microarray data analysis, this means that gene  $i$  does not contribute in the classification procedure and thus the row-sparsity of the coefficient matrix allows, at the same time, identify differentially expressed genes.

In the MMV model [6], the goal is to achieve simultaneous sparse reconstruction (SSR) of the signal matrix. The task is to estimate the  $K$ -row-sparse signal matrix  $\mathcal{B}$ , given an observed measurement matrix  $\mathbf{Y}$  and an (over complete) basis matrix (or dictionary)  $\Phi$ .  $K$ -row-sparsity of  $\mathcal{B}$  means that only  $K$  rows of  $\mathcal{B}$  contain non-zero entries. Commonly, this goal is achieved by  $\ell_{q,1}$  mixed matrix norm minimization, where

$$\|\mathcal{B}\|_{q,1} = \sum_{i=1}^p \|\beta_{[i]}\|_q,$$

for some  $q \geq 1$ , where  $\beta_{[i]}$  denotes the  $i$ -th row of  $\mathcal{B}$ . Values  $q = 1, 2, \infty$  have been advocated in the literature. Many SSR algorithms use *hard-thresholding operator*  $H_K(\cdot, q)$ , defined as transform  $H_K(\mathcal{B}, q)$ , which retains the elements of the  $K$  rows of  $\mathcal{B}$  that possess largest  $\ell_q$  norm and set elements of the other rows to zero. This leads us to define our *compressive RDA* discriminant function as

$$\begin{aligned} \hat{\mathbf{d}}(\mathbf{x}) &= (\hat{d}_1(\mathbf{x}), \dots, \hat{d}_G(\mathbf{x})) \\ &= \mathbf{x}^\top \hat{\mathcal{B}} - \frac{1}{2} \text{diag}(\hat{\mathbf{M}}^\top \hat{\mathcal{B}}) + \ln \boldsymbol{\pi}, \end{aligned} \quad (7)$$

where  $\ln \boldsymbol{\pi} = (\ln \pi_1, \dots, \ln \pi_G)$ ,  $\hat{\mathbf{M}} = (\hat{\boldsymbol{\mu}}_1 \ \dots \ \hat{\boldsymbol{\mu}}_G)$  and

$$\hat{\mathcal{B}} = H_K(\hat{\boldsymbol{\Sigma}}^{-1} \hat{\mathbf{M}}, q)$$

where  $\hat{\boldsymbol{\Sigma}}$  has been defined in (4) and  $\hat{\boldsymbol{\mu}}_g$  are the sample mean vectors of the classes in (2). Fast formula to compute  $\hat{\boldsymbol{\Sigma}}^{-1}$  is given in (5).

Next let us draw attention to SCRDA [1] which uses  $\hat{\boldsymbol{\Sigma}} = \alpha \mathbf{S} + (1 - \alpha) \mathbf{I}$  instead of estimator in (4). Another difference is in its use of element-wise soft-shrinkage. Namely, SCRDA can also be written in the multivariate form (7), but using

$$\hat{\mathcal{B}} = \mathcal{S}_\Delta(\hat{\boldsymbol{\Sigma}}^{-1} \hat{\mathbf{M}}) \quad (8)$$

where  $\mathcal{S}_\Delta(\cdot)$  is the soft-thresholding function that is applied element-wise to its matrix-valued argument. That is, the  $(i, j)$ -th element  $\hat{b}_{ij}$  of  $\hat{\mathcal{B}}$  in (8) is

$$\hat{b}_{ij} = \mathcal{S}_\Delta(t_{ij}) = \text{sign}(t_{ij})(|t_{ij}| - \Delta)_+$$

where  $(t)_+ = \max(t, 0)$  for  $t \in \mathbb{R}$  and  $t_{ij}$  denotes the  $(i, j)$ -th element of  $\mathbf{T} = \hat{\boldsymbol{\Sigma}}^{-1} \hat{\mathbf{M}}$ . One disadvantage of SCRDA is the shrinkage thresholding parameter  $\Delta \in [0, \infty)$  which is the same across all groups, and different populations would often benefit from different shrinkage intensity. A sensible upper

**Table 1.** Classification results for the simulation setups I–III. Figures in bold-face indicate the best results in each column. For setup III, the false positive (FP) and detection rate (DR) are also reported. Results in parenthesis are obtained using  $(\hat{\alpha}_{ell}, \hat{K}_{cv})$  instead of  $(\hat{\alpha}_{cv}, \hat{K}_{cv})$ .

Methods	Setup I		Setup II	
	TE	NFS	TE	NFS
CRDA <sup><math>\ell_1</math></sup>	120 (116)	165 (163)	180 (174)	105 (101)
CRDA <sup><math>\ell_2</math></sup>	95 (94)	126 (120)	184 (182)	96 (105)
CRDA <sup><math>\ell_\infty</math></sup>	<b>84 (81)</b>	<b>112 (114)</b>	185 (177)	<b>94 (96)</b>
PLDA	117	301	<b>151</b>	148
SCRDA	97	227	291	349
NSC	89	290	277	440
Setup III				
Methods	TE	NFS	DR	FP
CRDA <sup><math>\ell_1</math></sup>	<b>46</b> (50)	<b>205</b> (259)	90 ( <b>94</b> )	<b>12</b> (27)
CRDA <sup><math>\ell_2</math></sup>	49 ( <b>46</b> )	240 ( <b>203</b> )	<b>92</b> (92)	23 ( <b>10</b> )
CRDA <sup><math>\ell_\infty</math></sup>	50 (52)	238 (252)	89 (92)	27 (27)
SCRDA	108	282	69	51

bound of  $\Delta$  is difficult to determine and is highly data dependent. The proposed CRDA on the other hand uses simple to tune joint-sparsity level  $K \in \{1, 2, \dots, p\}$  and has the benefit of offering simultaneous joint-sparse recovery, i.e., features are eliminated across all groups instead of group-wise.

### 3.2. Model (Parameters) Selection

We employ  $Q$ -fold CV to estimate the optimal pair  $(\hat{\alpha}_{cv}, \hat{K}_{cv})$  using a 2D grid of candidate values  $\{\alpha_i\}_{i=1}^I \times \{K_j\}_{j=1}^J$  of the tuning parameters, where  $\alpha \in [0, 1)$  and  $K \in [1, p] = \{1, 2, \dots, p\} \subset \mathbb{N}$ . Often there are several pairs that yield the minimal cross-validation error from the training dataset and each pair can exhibit varying degree of sparsity (number of features selected). Among them, we would prefer the pair that had minimal number of features. Since a pair with minimal CV error may not yield a classifier that is at the same time sparse, one may wish to set a lower bound for the number of features selected (NFS) in order to enhance the interpretability of the discriminant function.

Let  $\varepsilon_{cv}(\alpha, K)$  denote a CV error for a pair  $(\alpha, K)$ . To have a trade-off between a minimum (CV-based) training error  $\varepsilon_{cv} \in [1, n]$  and NFS, we use a threshold  $\varepsilon^{\text{th}} = \max(0.15 \cdot n, \varepsilon_{cv})$  and choose only the pairs which have CV error smaller than  $\varepsilon^{\text{th}}$ , i.e., pairs which verify  $\varepsilon_{cv}(\alpha, K) \leq \varepsilon^{\text{th}}$ . From these pairs, the final optimal pair  $(\hat{\alpha}_{cv}, \hat{K}_{cv})$  is chosen as the one that has the smallest NFS value. For finding the optimal pair, we utilize a uniform grid of 100  $K$ -values and a uniform grid of 25  $\alpha$ -values.

We compare the CV approach to computationally much lighter approach which uses the estimated parameter  $\hat{\alpha}_{ell}$

**Table 2.** Classification results for the four microarray datasets using 5-fold CV. Note that figures in bold-face indicate the best results. Results in parenthesis are obtained using  $(\hat{\alpha}_{ell}, \hat{K}_{cv})$  instead of  $(\hat{\alpha}_{cv}, \hat{K}_{cv})$ .

Methods	Ramaswamy <i>et al.</i> dataset		Yeoh <i>et al.</i> dataset		Sun <i>et al.</i> dataset		Nakayama <i>et al.</i> dataset	
	TE / 47	NFS	TE / 62	NFS	TE / 45	NFS	TE / 26	NFS
CRDA <sup><math>\ell_1</math></sup>	10.6 ( <b>9.9</b> )	2634 (4899)	9.6 (7.5)	2525 (4697)	12.5 ( <b>12.9</b> )	23320 (27416)	8.3 (7.9)	2941 (6952)
CRDA <sup><math>\ell_2</math></sup>	10.4 (10.3)	2683 (3968)	9.7 ( <b>6.0</b> )	2273 ( <b>4659</b> )	12.9 (13.3)	<b>20589</b> (23484)	7.9 (7.6)	3142 (7755)
CRDA <sup><math>\ell_\infty</math></sup>	<b>10.3</b> (10.3)	3405 (4530)	<b>9.3</b> (6.5)	<b>846</b> (4697)	<b>12.4</b> (13.5)	21354 ( <b>20207</b> )	7.6 (7.6)	<b>2719</b> ( <b>2340</b> )
PLDA	18.8	5023	NA	NA	15.2	21635	4.4	10479
SCRDA	24	14874	NA	NA	15.7	54183	<b>2.8</b>	22283
NSC	16.3	<b>2337</b>	NA	NA	15	30005	4.2	5908

given in [8]. We note that value of  $\hat{\alpha}_{ell}$  can be computed efficiently using the SVD trick. Given the optimal RSCM based on  $\hat{\alpha}_{ell}$  we then estimate the sparsity level  $K$  using CV estimate  $\hat{K}_{cv}$ . This reduces the computational cost significantly.

#### 4. NUMERICAL EXAMPLES

The simulation study investigates the performance of CRDA based classifiers using different simulation setups commonly used in the RDA literature (e.g. in [1, 3, 10, 11]) and draw a comparison with the available results, against the nearest shrunken centroids (NSC) [2], SCRDA [1] and PLDA [3]. For simulation setups I and II, we generate 1200 observations from MVN distribution,  $\mathcal{N}_p(\boldsymbol{\mu}_g, \mathbf{I}_p)$ , with equal probabilities for each of  $G = 4$  groups. The observations are divided into three sets: (i) the validation set with 100 observations finds the tuning parameters, (ii) then 100 observations in the training set estimate  $\hat{\Sigma}^{-1}$  and (iii) the rest 1000 form the test set for calculating misclassification test errors (TE). A total of  $T = 100$  out of  $p = 500$  features differ between the groups. In setup I,  $\boldsymbol{\mu}_g$  contains  $t = 25$  nonzeros for each group  $g$  and rest all zeros, i.e.,  $[\boldsymbol{\mu}_g]_i = 0.7$  for  $t(g-1) + 1 \leq i \leq t(g-1) + t$ . While,  $[\boldsymbol{\mu}_g]_i = \frac{2-1}{3}$  if  $i \leq 100$  and zero otherwise for setup II. Table 1 lists the average of the TE and NFS for each classifier using 25 MC trials and 5-fold CV.

The third simulation setup resembles real gene expression data. We generate  $n = 200$  training and 1000 test observations each having  $p = 10,000$  features. All groups have equal probabilities and follow MVN distribution  $\mathcal{N}_p(\boldsymbol{\mu}_g, \boldsymbol{\Sigma}_g)$  for  $g = 1, \dots, G = 3$ . We have  $\boldsymbol{\mu}_1 = \mathbf{0}_p$  and  $\boldsymbol{\mu}_2$  contains all zeros except first 200-entries (i.e., true positives) with value  $1/2$  and  $\boldsymbol{\mu}_3 = -\boldsymbol{\mu}_2$ . Each group employs following block-diagonal auto-regressive covariance-structure

$$\boldsymbol{\Sigma}_g = \boldsymbol{\Sigma}^{(\rho_g)} \oplus \boldsymbol{\Sigma}^{(-\rho_g)} \oplus \dots \oplus \boldsymbol{\Sigma}^{(\rho_g)} \oplus \boldsymbol{\Sigma}^{(-\rho_g)},$$

where  $\oplus$  indicates the direct sum (not the Kronecker sum) of 100 block matrices having the AR(1) covariance structure

$$[\boldsymbol{\Sigma}^{(\rho_g)}]_{1 \leq i, j \leq 100} = \rho_g^{|i-j|}$$

where  $\rho_g$  is the correlation which is different for each group, namely,  $\rho_1 = 0.5$ ,  $\rho_2 = 0.7$  and  $\rho_3 = 0.9$ . This setup mimics real microarray data as genes are correlated within a pathway and independent between the pathways. Table 1 reveals the higher accuracy of the proposed CRDA methods compared to SCRDA when measured by TE, NFS, detection rate (DR) and FP rates. The results are averaged over 10 Monte-Carlo trials using 10-fold CV.

Next we do a comparison based on real microarray datasets. A summary of the used datasets is given below:

Dataset	$N$	$p$	$G$	Disease
Ramaswamy <i>et al.</i> [12]	190	16,063	14	Cancer
Yeoh <i>et al.</i> [13]	248	12,625	6	Leukemia
Sun <i>et al.</i> [14]	180	54,613	4	Glioma
Nakayama <i>et al.</i> [15]	105	22,283	10	Sarcoma

We compute the results for each dataset over 10 training-test set splits, each with a random choice of training and test set containing 75% and 25% of the total  $N$  observations, respectively. The average results of classification and gene-selection by CRDA methods are given in Table 2 with available comparison results. The proposed CRDA based classifiers showcase better classification and feature-selection results for all simulation setups. Overall, it seems that  $\ell_2$  and  $\ell_\infty$ -norm based CRDA methods are doing better as compared to others. Moreover, the CRDA based on  $\ell_\infty$ -norm appears to have best overall performance. Note that the proposed CRDA classifiers outperform other methods with a significant margin in the case of Ramaswamy *et al.* (with 14 groups) and Sun *et al.* (of  $p = 54,613$  genes).

#### 5. DISCUSSIONS AND CONCLUSIONS

We proposed a modified version of LDA, called compressive regularized discriminant CRDA, for analysis of data sets in high dimension low sample size situations. CRDA was shown to outperform competing methods in most of the cases. It also had the best detection rate which illustrates that the method can be a useful tool for accurate selection of (differentially expressed) genes in microarray studies.

## 6. REFERENCES

- [1] Yaqian Guo, Trevor Hastie, and Robert Tibshirani, "Regularized linear discriminant analysis and its application in microarrays," *Biostatistics*, vol. 8, no. 1, pp. 86–100, 2006.
- [2] Robert Tibshirani, Trevor Hastie, Balasubramanian Narasimhan, and Gilbert Chu, "Class prediction by nearest shrunken centroids, with applications to dna microarrays," *Statistical Science*, pp. 104–117, 2003.
- [3] Daniela M Witten and Robert Tibshirani, "Penalized classification using fisher's linear discriminant," *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, vol. 73, no. 5, pp. 753–772, 2011.
- [4] Alok Sharma, Kuldip K Paliwal, Seiya Imoto, and Satoru Miyano, "A feature selection method using improved regularized linear discriminant analysis," *Machine vision and applications*, vol. 25, no. 3, pp. 775–786, 2014.
- [5] Emanuel Neto, Felix Biessmann, Harald Aurlien, Helge Nordby, and Tom Eichele, "Regularized linear discriminant analysis of eeg features in dementia patients," *Frontiers in aging neuroscience*, vol. 8, 2016.
- [6] M. F. Duarte and Y. C. Eldar, "Structured compressed sensing: From theory to applications," *IEEE Transactions on Signal Processing*, vol. 59, no. 9, pp. 4053–4085, 2011.
- [7] Olivier Ledoit and Michael Wolf, "A well-conditioned estimator for large-dimensional covariance matrices," *Journal of multivariate analysis*, vol. 88, no. 2, pp. 365–411, 2004.
- [8] Esa Ollila, "Optimal high-dimensional shrinkage covariance estimation for elliptical distributions," in *Proc. European Signal Processing Conference (EUSIPCO 2017)*, Kos, Greece, 2017, pp. 1689–1693.
- [9] Jerome Friedman, Trevor Hastie, and Robert Tibshirani, *The elements of statistical learning*, vol. 1, Springer series in statistics New York, 2001.
- [10] John A Ramey and Phil D Young, "A comparison of regularization methods applied to the linear discriminant function with high-dimensional microarray data," *Journal of Statistical Computation and Simulation*, vol. 83, no. 3, pp. 581–596, 2013.
- [11] Yan Zhou, Baoxue Zhang, Gaorong Li, Tiejun Tong, and Xiang Wan, "Gd-rda: A new regularized discriminant analysis for high-dimensional data," *Journal of Computational Biology*, 2017.
- [12] Sridhar Ramaswamy, Pablo Tamayo, Ryan Rifkin, Sayan Mukherjee, Chen-Hsiang Yeang, Michael Angelo, Christine Ladd, Michael Reich, Eva Latulippe, Jill P Mesirov, et al., "Multiclass cancer diagnosis using tumor gene expression signatures," *Proceedings of the National Academy of Sciences*, vol. 98, no. 26, pp. 15149–15154, 2001.
- [13] Eng-Juh Yeoh, Mary E Ross, Sheila A Shurtleff, W Kent Williams, Divyen Patel, Rami Mahfouz, Fred G Behm, Susana C Raimondi, Mary V Relling, Anami Patel, et al., "Classification, subtype discovery, and prediction of outcome in pediatric acute lymphoblastic leukemia by gene expression profiling," *Cancer cell*, vol. 1, no. 2, pp. 133–143, 2002.
- [14] Lixin Sun, Ai-Min Hui, Qin Su, Alexander Vortmeyer, Yuri Kotliarov, Sandra Pastorino, Antonino Passaniti, Jayant Menon, Jennifer Walling, Rolando Bailey, et al., "Neuronal and glioma-derived stem cell factor induces angiogenesis within the brain," *Cancer cell*, vol. 9, no. 4, pp. 287–300, 2006.
- [15] Robert Nakayama, Takeshi Nemoto, Hiro Takahashi, Tsutomu Ohta, Akira Kawai, Kunihiko Seki, Teruhiko Yoshida, Yoshiaki Toyama, Hitoshi Ichikawa, and Tadashi Hasegawa, "Gene expression analysis of soft tissue sarcomas: characterization and reclassification of malignant fibrous histiocytoma," *Modern pathology*, vol. 20, no. 7, pp. 749–759, 2007.