
This is an electronic reprint of the original article.
This reprint may differ from the original in pagination and typographic detail.

Raninen, Elias; Ollila, Esa

Optimal Pooling of Covariance Matrix Estimates Across Multiple Classes

Published in:

2018 IEEE International Conference on Acoustics, Speech, and Signal Processing, ICASSP 2018 - Proceedings

DOI:

[10.1109/ICASSP.2018.8461327](https://doi.org/10.1109/ICASSP.2018.8461327)

Published: 10/09/2018

Document Version

Peer-reviewed accepted author manuscript, also known as Final accepted manuscript or Post-print

Please cite the original version:

Raninen, E., & Ollila, E. (2018). Optimal Pooling of Covariance Matrix Estimates Across Multiple Classes. In *2018 IEEE International Conference on Acoustics, Speech, and Signal Processing, ICASSP 2018 - Proceedings* (Vol. 2018-April, pp. 4224-4228). Article 8461327 (Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing). IEEE. <https://doi.org/10.1109/ICASSP.2018.8461327>

This material is protected by copyright and other intellectual property rights, and duplication or sale of all or part of any of the repository collections is not permitted, except that material may be duplicated by you for your research use or educational purposes in electronic or print form. You must obtain permission for any other use. Electronic or print copies may not be offered, whether for sale or otherwise to anyone who is not an authorised user.

OPTIMAL POOLING OF COVARIANCE MATRIX ESTIMATES ACROSS MULTIPLE CLASSES

Elias Raninen and Esa Ollila*

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

ABSTRACT

The paper considers the problem of estimating the covariance matrices of multiple classes in a low sample support condition, where the data dimensionality is comparable to, or larger than, the sample sizes of the available data sets. In such conditions, a common approach is to shrink the class sample covariance matrices (SCMs) towards the pooled SCM. The success of this approach hinges upon the ability to choose the optimal regularization parameter. Typically, a common regularization level is shared among the classes and determined via a procedure based on cross-validation. We use class-specific regularization levels since this enables the derivation of the optimal regularization parameter for each class in terms of the minimum mean squared error (MMSE). The optimal parameters depend on the true unknown class population covariances. Consistent estimators of the parameters can, however, be easily constructed under the assumption that the class populations follow (unspecified) elliptically symmetric distributions. We demonstrate the performance of the proposed method via a simulation study as well as via an application to discriminant analysis using both synthetic and real data sets.

Index Terms— Covariance matrix estimation, regularization, elliptical distribution, classification.

1. INTRODUCTION

The problem of estimating the covariance matrices of K classes appears in classification problems as well as in graphical models, where in the latter the inverse of the covariance matrix gives information about the conditional dependence structure of the graph [1]. In the context of classification, the covariance matrices determine the shape of the data within each class in the feature space. In low sample support conditions, where the data dimensionality p is comparable to, or larger than, the sample sizes n_k available from each population, conventional SCMs are susceptible to high variance and may not even be positive definite and hence invertible. Therefore, regularized sample covariance matrix (RSCM) estimators are needed. Their usefulness has been validated in many challenging data analysis problems [2] [3] [4]. The

level of regularization in the RSCM is determined by the regularization parameter, which is typically chosen via cross-validation. However, procedures based on cross-validation can be computationally prohibitive, especially when p is large. Hence, simple methods that do not require excessive tuning of hyperparameters are more desirable, especially in applications where the aforementioned conditions are predominant, such as in remote sensing [5].

We present an analytical method for choosing the regularization levels of a particular form of the RSCM, where the individual class SCMs are shrunk towards the pooled SCM. This type of regularization often improves the estimation accuracy when there is low sample support and when the population covariance matrices share a common structure. In classification, this type of regularization is used in regularized discriminant analysis [2] (RDA), in which a common regularization level for the classes is determined via cross-validation. Class-specific regularization has been used in a slightly different RSCM formulation in [6], where the parameters were likewise chosen using a method based on cross-validation. Our approach is different in that we derive analytical expressions for the optimal regularization parameters of the classes. Even though the derived expressions depend on the true unknown population covariance matrices, consistent estimates can easily be constructed in a similar fashion as in [7], which studied the single covariance matrix estimation problem. We only need to assume that the class populations follow (possibly different unspecified) elliptical distributions with finite 4th order moments. In the conducted simulations, the proposed method performed not only better than the method based on cross-validation, but also enjoys the advantages of low computational complexity and ease of implementation.

The remainder of this paper is organized as follows. Section 2 introduces the proposed estimator and an analytical expression for the optimal class-specific regularization parameter is derived. In Section 3, we construct an estimator of the optimal regularization parameter under the assumption of (unspecified) elliptical populations. In Section 4, we assess the mean squared error (MSE) performance of the proposed method via a simulation study. In Section 5, the method is applied to discriminant analysis, where the classification performance is evaluated using both synthetic and real data examples. Section 6 concludes the paper.

*The research of E. Raninen was supported by the Academy of Finland grant No. 298118.

2. OPTIMAL REGULARIZATION PARAMETER

Consider observing independent and identically distributed (i.i.d.) p -dimensional samples from K different classes (or populations), each class having n_k number of samples. Let $\{\mathbf{x}_{k,i}\}_{i=1}^{n_k}$ denote the data set of the k th class. The scatter of the data in the p -dimensional feature space is characterized by the covariance matrices

$$\boldsymbol{\Sigma}_k = \mathbb{E}[(\mathbf{x}_k - \boldsymbol{\mu}_k)(\mathbf{x}_k - \boldsymbol{\mu}_k)^\top],$$

where $\boldsymbol{\mu}_k = \mathbb{E}[\mathbf{x}_k]$ and \mathbf{x}_k denotes a random vector from the k th class. Let $\bar{\mathbf{x}}_k = (1/n_k) \sum_{i=1}^{n_k} \mathbf{x}_{k,i}$ and

$$\mathbf{S}_k = \frac{1}{n_k - 1} \sum_{i=1}^{n_k} (\mathbf{x}_{k,i} - \bar{\mathbf{x}}_k)(\mathbf{x}_{k,i} - \bar{\mathbf{x}}_k)^\top$$

denote the sample mean vector and the (unbiased) SCM of the k th class. We are interested in an RSCM defined for class k as

$$\hat{\boldsymbol{\Sigma}}_k(\beta_k) = \beta_k \mathbf{S}_k + (1 - \beta_k) \mathbf{S}, \quad (1)$$

where $\beta_k \in [0, 1]$, and \mathbf{S} denotes the pooled SCM, i.e., $\mathbf{S} = \sum_{k=1}^K \pi_k \mathbf{S}_k$, where $\pi_k = n_k / (\sum_{j=1}^K n_j)$. For ease of notation, we will omit the subscript k from β_k hereafter. Our goal is to determine the optimal regularization parameter β^* for class k which minimizes the MSE between the regularized estimator (1) and the true population covariance matrix,

$$\beta^* = \arg \min_{\beta \in [0,1]} \mathbb{E}[\|\hat{\boldsymbol{\Sigma}}_k(\beta) - \boldsymbol{\Sigma}_k\|_{\mathbb{F}}^2], \quad (2)$$

where $\|\cdot\|_{\mathbb{F}}$ denotes the Frobenius matrix norm, i.e., $\|\mathbf{A}\|_{\mathbb{F}}^2 = \text{tr}(\mathbf{A}^\top \mathbf{A})$ for any square matrix \mathbf{A} .

Next, we derive the optimal regularization parameter. Write $L_k(\beta) = \mathbb{E}[\|\hat{\boldsymbol{\Sigma}}_k(\beta) - \boldsymbol{\Sigma}_k\|_{\mathbb{F}}^2]$. Then note that

$$\begin{aligned} L_k(\beta) &= \mathbb{E}[\|\beta \mathbf{S}_k + (1 - \beta) \mathbf{S} - \boldsymbol{\Sigma}_k\|_{\mathbb{F}}^2] \\ &= \beta^2 \mathbb{E}[\|\mathbf{S}_k - \mathbf{S}\|_{\mathbb{F}}^2] + \mathbb{E}[\|\boldsymbol{\Sigma}_k - \mathbf{S}\|_{\mathbb{F}}^2] \\ &\quad - 2\beta \mathbb{E}[\text{tr}((\mathbf{S}_k - \mathbf{S})(\boldsymbol{\Sigma}_k - \mathbf{S}))]. \end{aligned} \quad (3)$$

The second derivative $L_k''(\beta)$ is positive whenever $\mathbf{S}_k \neq \mathbf{S}$. Thus under this assumption, the loss function is strictly convex and the optimal regularization parameter value can be found by solving $L_k'(\beta^*) = 0$, which yields

$$\beta^* = \frac{\mathbb{E}[\text{tr}((\mathbf{S}_k - \mathbf{S})(\boldsymbol{\Sigma}_k - \mathbf{S}))]}{\mathbb{E}[\|\mathbf{S}_k - \mathbf{S}\|_{\mathbb{F}}^2]}. \quad (4)$$

It can be seen that if $\mathbf{S}_k \approx \boldsymbol{\Sigma}_k$, i.e., the SCM is close to the true covariance matrix, then $\beta^* \approx 1$, and the estimator gives all the weight to the SCM. By expanding the expressions in the numerator and denominator of (4), we get

$$\beta^* = \frac{\text{tr}(\boldsymbol{\Sigma}_k^2) - \mathbb{E}[\text{tr}(\mathbf{S}_k \mathbf{S})] - \text{tr}(\boldsymbol{\Sigma}_k \mathbf{S}) + \mathbb{E}[\text{tr}(\mathbf{S}^2)]}{\mathbb{E}[\text{tr}(\mathbf{S}_k^2)] - 2\mathbb{E}[\text{tr}(\mathbf{S}_k \mathbf{S})] + \mathbb{E}[\text{tr}(\mathbf{S}^2)]}, \quad (5)$$

where $\boldsymbol{\Sigma} = \mathbb{E}[\mathbf{S}] = \sum_{k=1}^K \pi_k \boldsymbol{\Sigma}_k$,

$$\begin{aligned} \mathbb{E}[\text{tr}(\mathbf{S}^2)] &= \sum_j \pi_j^2 \mathbb{E}[\text{tr}(\mathbf{S}_j^2)] + \sum_{i \neq j} \pi_i \pi_j \text{tr}(\boldsymbol{\Sigma}_i \boldsymbol{\Sigma}_j), \\ \mathbb{E}[\text{tr}(\mathbf{S}_k \mathbf{S})] &= \pi_k \mathbb{E}[\text{tr}(\mathbf{S}_k^2)] + \sum_{j=1, j \neq k}^K \pi_j \text{tr}(\boldsymbol{\Sigma}_k \boldsymbol{\Sigma}_j), \end{aligned}$$

and $\text{tr}(\boldsymbol{\Sigma}_k \boldsymbol{\Sigma}) = \pi_k \text{tr}(\boldsymbol{\Sigma}_k^2) + \sum_{j=1, j \neq k}^K \pi_j \text{tr}(\boldsymbol{\Sigma}_k \boldsymbol{\Sigma}_j)$. With some further algebra, β^* in (5) becomes

$$\beta^* = \frac{(1 - \pi_k) \text{tr}(\boldsymbol{\Sigma}_k^2) - \pi_k \mathbb{E}[\text{tr}(\mathbf{S}_k^2)] + \delta_k}{(1 - 2\pi_k) \mathbb{E}[\text{tr}(\mathbf{S}_k^2)] + \delta_k}, \quad (6)$$

where

$$\begin{aligned} \delta_k &= \sum_j \pi_j^2 \mathbb{E}[\text{tr}(\mathbf{S}_j^2)] - 2 \sum_{j=1, j \neq k}^K \pi_j \text{tr}(\boldsymbol{\Sigma}_k \boldsymbol{\Sigma}_j) \\ &\quad + \sum_{i \neq j} \pi_i \pi_j \text{tr}(\boldsymbol{\Sigma}_i \boldsymbol{\Sigma}_j). \end{aligned}$$

The value of the MSE at the optimum is then

$$\begin{aligned} L_k(\beta^*) &= (\beta^*)^2 \mathbb{E}[\text{tr}(\mathbf{S}_k^2)] + (1 - \beta^*)^2 \mathbb{E}[\text{tr}(\mathbf{S}^2)] \\ &\quad + (1 - 2\beta^*) \text{tr}(\boldsymbol{\Sigma}_k^2) + 2\beta^*(1 - \beta^*) \mathbb{E}[\text{tr}(\mathbf{S}_k \mathbf{S})] \\ &\quad - 2(1 - \beta^*) \text{tr}(\boldsymbol{\Sigma}_k \boldsymbol{\Sigma}), \end{aligned} \quad (7)$$

which follows from a straightforward calculation of (3).

The optimal regularization parameter depends on the unknown true covariance matrices. Hence, this parameter must be estimated, which forms the topic of the next section.

3. ESTIMATING THE OPTIMAL REGULARIZATION PARAMETER

An estimate $\hat{\beta}$ of the optimal regularization parameter β^* of class k in (6) is obtained by the estimates of $\text{tr}(\mathbf{S}_k^2)$, $\mathbb{E}[\text{tr}(\mathbf{S}_k^2)]$, and $\text{tr}(\boldsymbol{\Sigma}_i \boldsymbol{\Sigma}_j)$, where $i \neq j$. The last expression can simply be estimated with $\text{tr}(\mathbf{S}_i \mathbf{S}_j)$, which follows from independence, i.e., $\mathbb{E}[\text{tr}(\mathbf{S}_i \mathbf{S}_j)] = \text{tr}(\mathbb{E}[\mathbf{S}_i] \mathbb{E}[\mathbf{S}_j]) = \text{tr}(\boldsymbol{\Sigma}_i \boldsymbol{\Sigma}_j)$.

We now assume that the random sample of the k th class, $\{\mathbf{x}_{k,i}\}_{i=1}^{n_k}$, is from an elliptically symmetric distribution with mean vector $\boldsymbol{\mu}_k$ and covariance matrix $\boldsymbol{\Sigma}_k$, and that it possesses finite 4th order moments. Then the probability density function (p.d.f.) of $\mathbf{x}_{k,i}$ is of the form

$$f_k(\mathbf{x}) = C_k |\boldsymbol{\Sigma}_k|^{-1/2} g_k((\mathbf{x} - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k)),$$

where $g_k : [0, \infty) \rightarrow [0, \infty)$ is a fixed function, called the *density generator*, which is independent of \mathbf{x} , $\boldsymbol{\mu}_k$ and $\boldsymbol{\Sigma}_k$, and C_k is a normalizing constant ensuring that $f_k(\mathbf{x})$ integrates to 1. We denote this case as $\mathbf{x}_{k,i} \sim \mathcal{E}_p(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, g_k)$. The functional form of the density generator $g_k(\cdot)$ determines the elliptical distribution. For example, the multivariate normal (MVN) distribution is obtained when $g_k(t) = \exp(-t/2)$. In our derivations, we do not have to assume that the density generators are the same for each class or even specify the elliptical populations; e.g., $g_k(t)$ can be the generator of the

multivariate t -distribution, MVN, or any other elliptical distribution with finite 4th order moments. Under the assumption that $\mathbf{x}_{k,i} \sim \mathcal{E}_p(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, g_k)$ and by using the results from [8], one can show that

$$\mathbb{E} [\text{tr}(\mathbf{S}_k^2)] = p\eta_k^2 \{\tau_1(p + \gamma_k) + (\tau_2 + 1)\gamma_k\}, \quad (8)$$

where $\tau_1 = (n_k - 1)^{-1} + \kappa_k/n_k$, $\tau_2 = \kappa_k/n_k$, $\eta_k = \text{tr}(\boldsymbol{\Sigma}_k)/p$, $\gamma_k = p\text{tr}(\boldsymbol{\Sigma}_k^2)/\text{tr}(\boldsymbol{\Sigma}_k)^2$, and κ_k is the elliptical kurtosis parameter of the k th population,

$$\kappa_k = (1/3) \cdot \text{kurt}(x_j), \quad (9)$$

where x_j denotes any element of $\mathbf{x}_{k,i} \sim \mathcal{E}_p(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, g_k)$ and $\text{kurt}(x)$ denotes the kurtosis of a random variable x .

A consistent estimate $\hat{\beta}$ of β^* for the k th class is obtained given consistent estimates of η_k , γ_k and κ_k . An obvious estimate of η_k is $\hat{\eta}_k = \text{tr}(\mathbf{S}_k)/p$. For γ_k , we used the estimator

$$\hat{\gamma}_{\text{sgn},k} = p\text{tr}(\mathbf{S}_{\text{sgn},k}^2) - \frac{p}{n_k}, \quad (10)$$

which uses the *sample sign covariance matrix* [9],

$$\mathbf{S}_{\text{sgn},k} = \frac{1}{n_k} \sum_{i=1}^{n_k} \frac{(\mathbf{x}_{k,i} - \hat{\boldsymbol{\mu}}_k)(\mathbf{x}_{k,i} - \hat{\boldsymbol{\mu}}_k)^\top}{\|\mathbf{x}_{k,i} - \hat{\boldsymbol{\mu}}_k\|^2}, \quad (11)$$

where $\hat{\boldsymbol{\mu}}_k = \arg \min_{\boldsymbol{\mu}} \sum_{i=1}^{n_k} \|\mathbf{x}_{k,i} - \boldsymbol{\mu}\|$ is the spatial sample median [10] of the k th data set. It was shown in [11] that $\hat{\gamma}_{\text{sgn},k}$ is a consistent estimator of γ_k under the random matrix theory regime. Their proof considered the centered case ($\boldsymbol{\mu}_k = 0$), where the centering by the spatial sample median was not required in (11). We have not yet extended this result to the case of unknown center, however, the estimator (10) performed well in our simulations. As in [7], a consistent estimator of κ_k is

$$\hat{\kappa}_k = \max \left\{ -\frac{2}{p+2}, \frac{1}{3p} \sum_{j=1}^p \hat{q}_{k,j} \right\}, \quad (12)$$

where $\hat{q}_{k,j}$ is the sample kurtosis of the j th variable of the k th class, i.e., $\hat{q}_{k,j} = m_{k,j}^{(4)}/(m_{k,j}^{(2)})^2 - 3$, where $m_{k,j}^{(l)} = (1/n_k) \sum_{i=1}^{n_k} ((\mathbf{x}_{k,i})_j - (\bar{\mathbf{x}}_k)_j)^l$ denotes the l th order moment of the j th variable of the k th class.

As the final estimate of the optimal regularization parameter, we used $\max\{0, \min\{1, \hat{\beta}\}\}$ since β^* needs to be within $[0, 1]$.

4. SIMULATION STUDY

We now illustrate the performance of the proposed method via a simulation study. We generated training data of dimension $p = 20$ comprising $K = 4$ classes and $n = \sum_k n_k = 100$ samples. The data was generated from a Student's t_ν -distribution with $\nu = 10$ degrees of freedom. The first class, $k = 1$, had zero mean and the subsequent classes, $k = 2, 3$, and 4, had means in orthogonal directions such that $\|\boldsymbol{\mu}_k\| = 1 + k$. We simulated three different set-ups:

Table 1. The empirical NMSE \tilde{L}_k for the covariance estimates in the set-ups 1 to 3 (from top to down). The corresponding standard deviations are shown in parenthesis.

	\tilde{L}_1	\tilde{L}_2	\tilde{L}_3	\tilde{L}_4	Sum
Oracle	0.99 (0.41)	0.52 (0.16)	0.27 (0.06)	0.28 (0.03)	2.06 (0.48)
Prop 1	0.98 (0.38)	0.50 (0.15)	0.28 (0.06)	0.29 (0.03)	2.04 (0.45)
Pool	4.47 (1.15)	0.62 (0.21)	0.27 (0.06)	0.28 (0.03)	5.63 (1.42)
SCM	1.13 (0.48)	1.12 (0.41)	1.20 (0.67)	1.18 (0.53)	4.63 (1.07)
Oracle	2.13 (1.09)	0.70 (0.29)	0.32 (0.12)	0.24 (0.05)	3.40 (1.23)
Prop 1	2.07 (0.97)	0.67 (0.20)	0.31 (0.08)	0.24 (0.05)	3.29 (1.06)
Pool	6.80 (1.85)	0.96 (0.37)	0.32 (0.13)	0.24 (0.05)	8.32 (2.38)
SCM	2.89 (1.63)	1.42 (0.84)	0.92 (0.35)	0.73 (0.41)	5.95 (1.94)
Oracle	1.25 (0.71)	0.93 (0.40)	0.40 (0.17)	0.24 (0.09)	2.81 (0.90)
Prop 1	1.18 (0.60)	0.88 (0.30)	0.38 (0.12)	0.24 (0.07)	2.68 (0.72)
Pool	6.25 (1.77)	2.04 (0.73)	0.43 (0.23)	0.29 (0.05)	9.01 (2.71)
SCM	1.50 (1.05)	1.32 (0.75)	0.86 (0.34)	0.39 (0.31)	4.07 (1.36)

1. The true covariance matrices were $\boldsymbol{\Sigma}_k = k\mathbf{I}$, and the sample sizes were $n_k = 25$ for all k .
2. The true covariance matrices were $\boldsymbol{\Sigma}_k = k\mathbf{I}$, and the sample sizes were $n_k = 10 \cdot k$.
3. The true covariance matrices were generated such that the ij th entry of the covariance matrix was $(\boldsymbol{\Sigma}_k)_{ij} = k\rho_k^{|i-j|}$, where $\rho_1 = -0.6$, $\rho_2 = -0.2$, $\rho_3 = 0.2$, and $\rho_4 = 0.6$, and the sample sizes were $n_k = 10 \cdot k$.

We report the empirical normalized MSE (NMSE), $\tilde{L}_k(\hat{\boldsymbol{\Sigma}}_k) = \text{Ave} \|\hat{\boldsymbol{\Sigma}}_k - \boldsymbol{\Sigma}_k\|_F^2 / \|\boldsymbol{\Sigma}_k\|_F^2$, of our proposed estimator (Prop 1) and the oracle estimator (Oracle), which uses the true values of η_k , γ_k , κ_k , and $\boldsymbol{\Sigma}_k$ in (6) and (8). The NMSEs of the SCM and the pooled SCM (Pool) are also shown. The results were averaged over 300 Monte-Carlo (MC) trials and are given in Table 1. The proposed method (Prop 1) provided not only a significant improvement both over the SCM and the pooled SCM by yielding the smallest NMSEs and standard deviations, but was also somewhat robust to estimation errors in the sense that it was able to perform at the level of the oracle estimator.

5. APPLICATION TO DISCRIMINANT ANALYSIS

In linear and quadratic discriminant analysis, one uses a discriminant rule which assigns any new observation \mathbf{x} to class \hat{k} , such that it minimizes the quadratic discriminant function

$$\hat{k} = \arg \min_k (\mathbf{x} - \bar{\mathbf{x}}_k)^\top \hat{\boldsymbol{\Sigma}}_k^{-1} (\mathbf{x} - \bar{\mathbf{x}}_k) + \log |\hat{\boldsymbol{\Sigma}}_k|,$$

where $\hat{\boldsymbol{\Sigma}}_k$ denotes an estimator of $\boldsymbol{\Sigma}_k$. In quadratic discriminant analysis (QDA), one uses $\hat{\boldsymbol{\Sigma}}_k = \mathbf{S}_k$ for all k , whereas in linear discriminant analysis (LDA), one uses $\hat{\boldsymbol{\Sigma}}_k = \mathbf{S}$ for all k . In RDA [2], $\hat{\boldsymbol{\Sigma}}_k(\beta_k)$ of (1), using $\beta_k = \beta$ for all k , is further regularized towards a scaled identity matrix by

$$\hat{\boldsymbol{\Sigma}}_k(\alpha, \beta) = \alpha \hat{\boldsymbol{\Sigma}}_k(\beta) + (1 - \alpha) (\text{tr}(\hat{\boldsymbol{\Sigma}}_k(\beta))/p) \mathbf{I}, \quad (13)$$

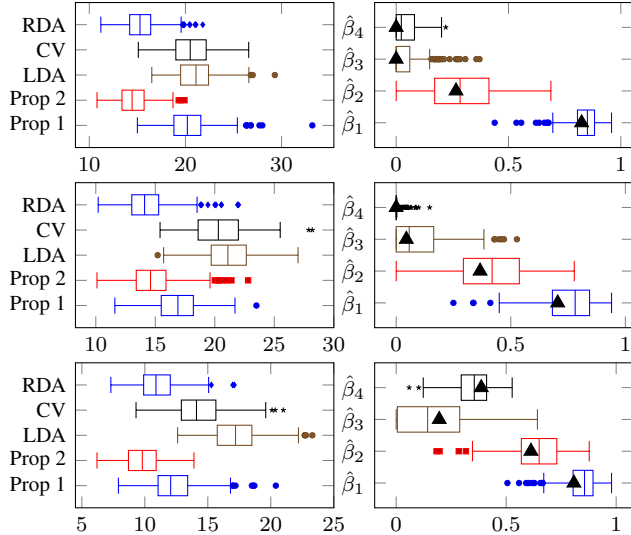


Fig. 1. The misclassification rate $\times 100$ and the corresponding boxplots of $\hat{\beta}_k$ for the set-ups 1, 2 and 3 shown from top to bottom. The black triangles denote the oracle values β_k^* .

and $\alpha, \beta \in [0, 1]$ are common over the classes and chosen via cross-validation. Additional regularization via parameter α helps to stabilize the estimates if the total sample size $n = \sum_k n_k$ is small compared to p . We also applied (13) to further shrink the covariance estimates obtained by our method using the approach in [7], thus acquiring the estimates of α_k by

$$\hat{\alpha}_k = \max \left\{ 0, \frac{T_k}{T_k + \frac{1}{n_k}(\hat{\kappa}_k(2\hat{\gamma}_k + p) + \hat{\gamma}_k + p)} \right\}, \quad (14)$$

where $T_k = \hat{\gamma}_k - 1$, yielding an estimator which we abbreviate Prop 2. Unlike in the single covariance matrix estimation problem of [7], Prop 2 does not have the interpretation of being the MMSE estimator.

In our next example, we used the same simulation set-ups as in Section 4. However, for each MC trial, an additional test data set comprising $10 \cdot n_k$ observations from each class was generated and used for computing the misclassification rate. In addition to our proposed estimators (Prop 1 and 2), the results are reported for LDA, RDA, and CV, of which the last uses only a common cross-validated β in (1) for shrinking the SCMs towards the pooled SCM. Using 10-fold cross-validation, the grid of the regularization parameter values for α and β ranged from 0 to 1 with a step of 0.05. In situations, where equal training errors were obtained with distinct values of α and β , we averaged the performance of all of them. The QDA estimate was omitted, since it does not exist if the SCM is not invertible. The average misclassification rates, over 300 MC trials, are given in Figure 1 along with the boxplots of the estimated regularization parameters $\hat{\beta}_k$. The black triangles in the boxplots denote the oracle values β_k^* computed using the

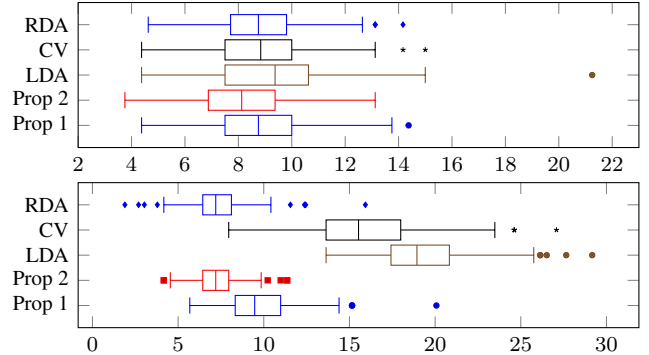


Fig. 2. Misclassification rate $\times 100$ for the glass and ionosphere data shown from top to bottom.

true values of $\eta_k, \gamma_k, \kappa_k$, and Σ_k in (6) and (8). As shown, the optimal regularization parameters were estimated reasonably well. Also, Prop 1 performed better than CV, whereas Prop 2 performed at the level of RDA, if not slightly better.

Next, we applied the proposed methods to the following real data examples obtained from the UCI Machine Learning Repository [12]. The **glass data set** had $p = 9$ variables and sample sizes of $n_1 = 51$ (window glass) and $n_2 = 163$ (non-window glass). The **ionosphere data set** had $p = 32$ variables¹ and sample sizes $n_1 = 126$ (bad radar return) and $n_2 = 225$ (good radar return). A fraction 1/4 of the samples from each class were chosen randomly as training data in each MC trial, and the remaining samples were used as test data. This made the training data sizes comparable to the dimension for at least one of the classes. The misclassification rates were averaged over 300 MC trials and are shown in Figure 2. In addition to the proposed estimators, the performances of LDA, CV, and RDA are given. As Figure 2 shows, Prop 1 performed better than CV, whereas Prop 2 and RDA performed equally well.

6. CONCLUSION

We considered joint covariance matrix estimation of multiple classes, where the SCMs of the classes are individually regularized towards the pooled SCM. We derived the optimal class-specific regularization parameters and showed how they could be estimated when the data is considered to be elliptically distributed. The conducted synthetic simulation study showed that the optimal regularization parameters could be estimated with an MSE performance close to the optimal oracle level. When applied to discriminant analysis classification, the conducted synthetic and real data simulations indicated that the performance of the proposed methods are on a par with, and often better than, the computationally more intensive methods based on cross-validation.

¹Two variables, which were zero for all samples in the classes, were removed from the original data.

7. REFERENCES

- [1] Patrick Danaher, Pei Wang, and Daniela M. Witten, “The joint graphical lasso for inverse covariance estimation across multiple classes,” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, vol. 76, no. 2, pp. 373–397, 2014.
- [2] Jerome H. Friedman, “Regularized discriminant analysis,” *Journal of the American Statistical Association*, vol. 84, no. 405, pp. 165–175, 1989.
- [3] Yaqian Guo, Trevor Hastie, and Robert Tibshirani, “Regularized linear discriminant analysis and its application in microarrays,” *Biostatistics*, vol. 8, no. 1, pp. 86–100, 2007.
- [4] Trevor Hastie, Robert Tibshirani, and Jerome Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, Springer, 2nd edition, 2009.
- [5] Tatyana V. Bandos, Lorenzo Bruzzone, and Gustavo Camps-Valls, “Classification of hyperspectral images with regularized linear discriminant analysis,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 47, no. 3, pp. 862–873, 2009.
- [6] Joseph P. Hoffbeck and David A. Landgrebe, “Covariance matrix estimation and classification with limited training data,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 18, no. 7, pp. 763–767, 1996.
- [7] Esa Ollila, “Optimal high-dimensional shrinkage covariance estimation for elliptical distributions,” in *Proc. 25th European Signal Processing Conference (EUSIPCO 2017)*, Kos, Greece, 2017, pp. 1639–1643.
- [8] David E. Tyler, “Radial estimates and the test for sphericity,” *Biometrika*, vol. 69, no. 2, pp. 429–436, 1982.
- [9] Christophe Croux, Esa Ollila, and Hannu Oja, “Sign and rank covariance matrices: statistical properties and application to principal components analysis,” in *Statistical data analysis based on the L1-norm and related methods*, pp. 257–269. Birkhäuser, Basel, 2002.
- [10] B. M. Brown, “Statistical uses of the spatial median,” *Journal of the Royal Statistical Society. Series B (Methodological)*, vol. 45, no. 1, pp. 25–30, 1983.
- [11] Teng Zhang and Ami Wiesel, “Automatic diagonal loading for Tyler’s robust covariance estimator,” in *IEEE Statistical Signal Processing Workshop (SSP 2016)*, 2016, pp. 1–5.
- [12] M. Lichman, “UCI machine learning repository,” <http://archive.ics.uci.edu/ml>, 2013.