

Cellwise robust regularized discriminant analysis

Aerts S, Wilms I.



Cellwise robust regularized discriminant analysis

Stéphanie Aerts*¹ and Ines Wilms²

Abstract

Quadratic and Linear Discriminant Analysis (QDA/LDA) are the most often applied classification rules under normality. In QDA, a separate covariance matrix is estimated for each group. If there are more variables than observations in the groups, the usual estimates are singular and cannot be used anymore. Assuming homoscedasticity, as in LDA, reduces the number of parameters to estimate. This rather strong assumption is however rarely verified in practice. Regularized discriminant techniques that are computable in high-dimension and cover the path between the two extremes QDA and LDA have been proposed in the literature. However, these procedures rely on sample covariance matrices. As such, they become inappropriate in presence of cellwise outliers, a type of outliers that is very likely to occur in high-dimensional datasets. In this paper, we propose cellwise robust counterparts of these regularized discriminant techniques by inserting cellwise robust covariance matrices. Our methodology results in a family of discriminant methods that (i) are robust against outlying cells, (ii) cover the gap between LDA and QDA and (iii) are computable in high-dimension. The good performance of the new methods is illustrated through simulated and real data examples. As a by-product, visual tools are provided for the detection of outliers.

Keywords

Cellwise robust precision matrix; Classification; Discriminant analysis; Penalized estimation.

1 Introduction

Consider a training set $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ of N observations of dimension p , each belonging to one of K groups G_1, \dots, G_K , with n_k observations in the k -th group and $N = \sum_{k=1}^K n_k$. Discriminant analysis methods aim to construct a decision rule based on \mathbf{X} that automatically assigns a new observation \mathbf{x} to one of the K groups. If the group conditional densities

¹*HEC- Liège, University of Liège, stephanie.aerts@ulg.ac.be, Phone: +324327272*

²*Leuven Statistics Research Centre (LStat), KU Leuven, ines.wilms@kuleuven.be*

$f_k(\mathbf{x})$ are known, the Bayes classifier $\delta(\cdot)$ assigns \mathbf{x} to the group with maximum posterior probability. In *quadratic discriminant analysis* (QDA), where the conditional distributions are assumed Gaussian $N_p(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$, this yields the rule

$$\delta(\mathbf{x}) = \arg \min_k ((\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Theta}_k (\mathbf{x} - \boldsymbol{\mu}_k) - \log(\det \boldsymbol{\Theta}_k) - 2 \log \pi_k), \quad (1)$$

where $\boldsymbol{\Theta}_k := \boldsymbol{\Sigma}_k^{-1}$ is the k -th group precision matrix. This rule splits the measurement space into K disjoint regions with quadratic boundaries. In the special case of *linear discriminant analysis* (LDA), homoscedasticity is further assumed, yielding linear boundaries. Even when the population group precision matrices substantially differ, LDA is often used because it might improve estimation accuracy by reducing the number of parameters to estimate.

In practice, the group parameters $\boldsymbol{\mu}_k$ and $\boldsymbol{\Sigma}_k$ are commonly estimated by the arithmetic mean $\bar{\mathbf{x}}_k$ and the sample covariance matrix $\hat{\boldsymbol{\Sigma}}_k$ in QDA, or the sample pooled covariance matrix $\hat{\boldsymbol{\Sigma}}_{\text{pool}}$ in LDA. These procedures will be denoted by s-QDA and s-LDA from now on. However, when $p \approx n_k$, these estimators become highly inaccurate and for $p > n_k$, their inverse cannot be computed anymore. In the sequel, *high-dimension* refers to such settings. Regularized estimators have been successful in obtaining accurate estimates of $\boldsymbol{\Theta}_k$ in high-dimension. They do this by biasing away the estimates from the sample covariance matrices. Particular focus is given to *sparse* precision matrices where many elements are estimated as zero, see e.g. [25]. Xu et al. [24] propose to plug the popular *Graphical Lasso* [11] sparse precision matrices in quadratic rule (1). Similarly, one can obtain a sparse pooled precision matrix for LDA. However, QDA and LDA are two extreme cases and their underlying assumptions, i.e. all distinct covariance matrices in QDA, and homoscedasticity in LDA, are rather strong. Therefore, several regularized discriminant methods have been proposed that cover the path between LDA and QDA, see e.g [11] or [17].

All these procedures take the sample covariance matrices $\hat{\boldsymbol{\Sigma}}_k$ as input. These estimators are however not robust to outliers, i.e. atypical observations. Therefore, the proposed methods inherit their lack of robustness. Since outliers frequently occur in high-dimensional datasets, their possible presence should be accounted for. Several high-dimensional procedures have been proposed to *detect* outliers (see e.g. [9] or [27] for a review). Our focus is on how to *deal* with outliers in regularized discriminant analysis.

To robustify the discriminant rule (1), one could think of replacing the group means and covariance matrices by standard robust estimates. Croux and Dehon [4] use the S-estimator, while Hubert and Van Driessen [16] and Filzmoser et al. [10] insert the MCD estimator

into rule (1). Nevertheless, these estimators are not computable anymore in high-dimension. For high-dimensional datasets, robust discriminant methods have been investigated in [15] and [23]. However, these methods circumvent the high-dimensionality problem by applying a two-step procedure. First, a robust dimension reduction technique is applied. Then a robust discriminant rule, using standard robust location and covariance matrix estimates, is computed in this low-dimensional subspace.

Another issue with the standard robust estimators is that they usually downweight an observation even if only one of its components is contaminated. In high-dimensional datasets where many variables are measured on a small number of observations, this may result in a huge loss of information. For such high-dimensional datasets, the *cellwise contamination* model (see [3]), where each observation may contain at least one contaminated component, is more appropriate. The development of cellwise robust procedures only appeared recently, see e.g. [22] and [1] for low-dimensional datasets, or [20] and [5] for high-dimensional ones.

In this paper, we use cellwise robust covariance matrix estimates as an input for regularized discriminant methods. As a result, we obtain discriminant methods that deal with two important topics in applied statistics: regularized estimation and the presence of outliers in high-dimensional datasets. The resulting family of discriminant methods has clear advantages: (i) the methods are robust against cellwise outliers, (ii) as a by-product, they provide a way to detect both rowwise and cellwise outliers, (iii) they cover the path between LDA and QDA, and (iv) they are computable in high-dimension without requiring an initial dimension reduction technique.

The remainder of this article is structured as follows. In Section 2, we review several non robust regularized discriminant methods. We propose cellwise robust counterparts in Section 3. Simulation studies in Section 4 compare the proposed methods and show their good performance in contaminated and uncontaminated settings. Finally, we analyze two real data sets in Section 5. We find that the proposed cellwise robust discriminant methods improve the classification performance. Furthermore, two visual tools for outlier detection are provided. The conclusions are outlined in Section 6.

2 Regularized discriminant methods

To classify a new observation \mathbf{x} in one of the K groups on the basis of the discriminant rule (1), we need estimators of the group means $\boldsymbol{\mu}_k$ and precision matrices $\boldsymbol{\Theta}_k$. The usual

estimators for the group means are the average means $\bar{\mathbf{x}}_k$. In this section, we review several procedures to obtain high-dimensional precision matrix estimates that can then simply be plugged into (1). Their cellwise robust version will be discussed in Section 3.

GL-LDA and GL-QDA. Starting from the sample covariance matrix in the k -th group $\hat{\Sigma}_k$, the *Graphical Lasso* [11] maximizes the L_1 penalized log-likelihood

$$\hat{\Theta}_{k,\text{GL}} := \underset{\Theta_k}{\operatorname{argmax}} \quad n_k \log \det(\Theta_k) - n_k \operatorname{tr}(\Theta_k \hat{\Sigma}_k) - \lambda_1 \sum_{i \neq j} |\theta_{k,ij}|, \quad (2)$$

subject to the constraint that Θ_k is positive definite. Here, $\theta_{k,ij}$ is the element (i, j) of Θ_k and $\lambda_1 \geq 0$ is a regularization parameter. The L_1 -norm of the off-diagonal elements ensures that problem (2) can be solved even when the dimension p exceeds the group sample size n_k . For large values of λ_1 , many off-diagonal elements in $\hat{\Theta}_k$ will be equal to zero. Under normality, this can be interpreted as conditional independence between the corresponding variables in the specific group. Problem (2) can be solved using the R-package `huge` [26].

We denote by GL-QDA, the quadratic classifier obtained by computing the Graphical Lasso in each group and by plugging $\hat{\Theta}_{1,\text{GL}}, \dots, \hat{\Theta}_{K,\text{GL}}$ into (1), see [24]. A regularized estimate of the pooled precision matrix can be obtained in a similar way by using

$$\hat{\Sigma}_{\text{pool}} := \sum_{k=1}^K \frac{n_k}{n - K} \hat{\Sigma}_k \quad (3)$$

as input in (2). We denote by GL-LDA the resulting linear classifier.

JGL-DA. The GL-QDA discriminant rule does not exploit the potential similarities between the groups since it estimates the K precision matrices independently. On the other hand, the homoscedasticity assumption behind GL-LDA ignores the group specificities that may be of particular interest in the classification context. To encourage similar sparsity patterns across the groups, Price et al. [17] (and Gao et al. [13] in model based clustering) use the *Joint Graphical Lasso* (JGL) [7]. The JGL estimates are

$$\begin{aligned} \left(\hat{\Theta}_{k,\text{JGL}} \right)_{k=1}^K &:= \underset{\Theta_1, \dots, \Theta_K}{\operatorname{argmax}} \quad \sum_{k=1}^K n_k \log \det(\Theta_k) - n_k \operatorname{tr}(\Theta_k \hat{\Sigma}_k) \\ &\quad - \lambda_1 \sum_{k=1}^K \sum_{i \neq j} |\theta_{k,ij}| - \lambda_2 \sum_{k < k'} \sum_{i,j} |\theta_{k,ij} - \theta_{k',ij}|, \quad (4) \end{aligned}$$

subject to the constraint that $\Theta_1, \dots, \Theta_K$ are positive definite. The first penalty in (4), with regularization parameter $\lambda_1 \geq 0$, is the same as in (2). The second penalty in (4), with regularization parameter $\lambda_2 \geq 0$, encourages similar sparsity patterns across the groups and similar signs and values for the non-zero elements. The GL-QDA estimator corresponds to the particular case with $\lambda_2 = 0$. Large values of the similarity parameter λ_2 yield precision matrices with many similar elements across the groups. Varying the parameter λ_2 provides a variety of classifiers that lie in between LDA and QDA. Problem (4) can be solved using the R-package JGL [6]. We denote by JGL-DA the discriminant rule obtained by plugging the precision matrices of equation (4) into rule (1).

RDA. Friedman [12] proposes another regularized discriminant method, denoted by RDA from now on. Like JGL-DA, it gives a path from LDA to QDA and it is computable in high-dimension. Unlike JGL-DA, it does not provide sparse precision matrix estimates.

RDA starts by computing a convex combination of the group specific and pooled sample covariance matrices

$$\widehat{\Sigma}_k^{\rho_1} = (1 - \rho_1)\widehat{\Sigma}_k + \rho_1\widehat{\Sigma}_{\text{pool}}, \quad (5)$$

where $0 \leq \rho_1 \leq 1$ is a regularization parameter. Then, the resulting estimator (5) is shrunken towards a multiple of the identity matrix

$$\widehat{\Sigma}_{k,\text{RDA}} = (1 - \rho_2)\widehat{\Sigma}_k^{\rho_1} + \frac{\rho_2}{p}\text{tr}(\widehat{\Sigma}_k^{\rho_1})\mathbf{I}_p, \quad (6)$$

with a second regularization parameter $0 \leq \rho_2 \leq 1$. The s-QDA solution corresponds to $\rho_1 = \rho_2 = 0$ while the s-LDA solution is obtained with $\rho_2 = 0, \rho_1 = 1$. The resulting estimators $\widehat{\Sigma}_{k,\text{RDA}}$ can then be inverted to obtain regularized precision matrix estimates to be plugged into (1).

3 Cellwise robust discriminant methods

The estimators from Section 2 all use the sample covariance matrices and/or the pooled covariance as input and are therefore not robust against cellwise outliers. To obtain cellwise robust discriminant methods, we start with computing initial cellwise robust covariance matrices \mathbf{S}_k and the corresponding pooled covariance \mathbf{S}_{pool} . These cellwise robust covariance matrices are used to replace the sample covariance matrices $\widehat{\Sigma}_k$ and $\widehat{\Sigma}_{\text{pool}}$ as input of the

Graphical Lasso in (2), the Joint Graphical Lasso in (4) or RDA in (5) and (6). Then, we use these precision matrices along with robust mean estimates, namely the vector of marginal medians, in the discriminant rule (1). As a result, we obtain discriminant methods that are both robust against cellwise outliers and easily computable in high-dimension.

These cellwise robust counterparts of the regularized discriminant methods of Section 2 are denoted by rGL-LDA, rGL-QDA, rJGL-DA and rRDA from now on. The code to obtain these estimators is made available on <http://feb.kuleuven.be/ines.wilms/software>.

3.1 Cellwise robust covariance matrix estimates

We estimate each bivariate covariance between variables X^i and X^j by

$$s_{ij} = \widehat{\text{scale}}(X^i)\widehat{\text{scale}}(X^j)\widehat{\text{corr}}(X^i, X^j), \quad (7)$$

as in Croux and Öllerer [5]. We use the robust Q_n -estimator [18] as $\widehat{\text{scale}}(\cdot)$ and the Kendall's correlation estimator as $\widehat{\text{corr}}(\cdot)$. For a bivariate sample $(x_1^i, x_1^j), \dots, (x_n^i, x_n^j)$, this correlation estimator is defined as

$$\widehat{\text{corr}}_K(X^i, X^j) = \frac{2}{n(n-1)} \sum_{l < m} \text{sign} \left((x_l^i - x_m^i)(x_l^j - x_m^j) \right).$$

By using signs rather than numerical values, Kendall correlation is a robust correlation measure and, hence, can cope with outliers.

A $\mathcal{O}(n \log n)$ algorithm to compute it is available in the `pcaPP` package [8]. The covariance matrices obtained by estimating each pairwise covariance as in (7) are denoted by \mathbf{S}_k and the corresponding pooled covariance by \mathbf{S}_{pool} . Croux and Öllerer [5] show that replacing the sample covariances by these robust estimates in the Graphical Lasso estimator (2) results in robust precision matrices with 50% breakdown point against cellwise contamination. Note that few proposals of cellwise robust covariance estimators that are computable in high-dimension are available in the literature. Alternatives can be found in [2] and [20].

3.2 Selection of the regularization parameters

The regularized methods rGL-LDA, rGL-QDA, rJGL-DA and rRDA all depend on one or two regularization parameters. To select the regularization parameters of a given method, we consider a grid of values, i.e. a one-dimensional grid for rGL-QDA and rGL-LDA, and a two-dimensional grid for rJGL-DA and rRDA. For each (combination of) regularization

parameter(s), we compute the corresponding precision matrix estimates $\widehat{\Theta}_1, \dots, \widehat{\Theta}_K$ (or $\widehat{\Theta}_{\text{pool}}$ for LDA) and we search for the optimal ones minimizing the Bayesian Information Criterion

$$\text{BIC} = \sum_{k=1}^K \left[n_k \text{tr} \left(\mathbf{S}_k \widehat{\Theta}_k \right) - n_k \log(\det(\widehat{\Theta}_k)) \right] + \log(N) \text{df}, \quad (8)$$

where df denotes the degrees of freedom of the model. The cellwise robust covariance matrices \mathbf{S}_k are replaced by \mathbf{S}_{pool} in LDA. Note that $\widehat{\Theta}_1, \dots, \widehat{\Theta}_K$ and df depend on the regularization parameters.

We take the degrees of freedom df to be the total number of *distinct* non-zero elements in $\widehat{\Theta}_1, \dots, \widehat{\Theta}_K$. Danaher et al. [7] replace df by the total number of non-zero elements in the estimated precision matrices. As such, model complexity is however overestimated since it does not take into account the fact that some elements may be identical across the $\widehat{\Theta}_k$.

Grid bounds. For each regularization parameter, we consider a logarithmic spaced grid of five values between chosen upper and lower bounds, i.e. the grid consists of the exponential of five equally spaced values between the logarithms of the bounds. The lower bound is a tenth of the upper bound. Below, we discuss the choice of the upper bound for each of the proposed methods. Up to our knowledge, no upper or lower bound for such a grid is proposed in the literature for the multiple group setting.

For rGL-QDA, we take as upper bound for λ_1

$$\lambda_{1,\max}^{\text{rGL-QDA}} = \max_k \max_{i,j} n_k |(\mathbf{S}_k)_{ij} - \mathbf{I}_{ij}|, \quad (9)$$

where \mathbf{I} is the identity matrix. This value is the maximum over the K upper bounds considered by default in the `huge` package when performing the Graphical Lasso in each group. The upper bound for rGL-LDA can be obtained by replacing \mathbf{S}_k by \mathbf{S}_{pool} in (9).

For rJGL-DA, we consider

$$\lambda_{1,\max}^{\text{rJGL-DA}} = \max_k \max_{i,j;i \neq j} n_k |(\mathbf{S}_k)_{ij}|,$$

as upper bound for λ_1 since $\lambda_1 \geq \lambda_{1,\max}$ is a sufficient condition for all the off-diagonal elements of the solution to be zero [7]. For λ_2 , we propose the heuristic upper bound

$$\lambda_{2,\max}^{\text{rJGL-DA}} = \max_k \max_{i,j} n_k |(\mathbf{S}_{\text{pool}})_{ij} - (\mathbf{S}_k)_{ij}|.$$

We take this bound since $\lambda_2 > \lambda_{2,\max}^{\text{rJGL-DA}}$ is a necessary condition for \mathbf{S}_{pool} to satisfy the KKT conditions of problem (4) (for $\lambda_1 = 0$, if \mathbf{S}_{pool} has full rank and $K = 2$).

For rRDA, the regularization parameters ρ_1 and ρ_2 in equations (5) and (6) are the coefficients of convex combinations. Hence, we set the upper bounds to one.

4 Simulation study

In this section, we compare the performance of the regularized discriminant methods (cfr. Section 2) and their cellwise robust counterparts (cfr. Section 3) through a simulation study. The considered non robust methods are s-LDA, s-QDA, GL-LDA, GL-QDA, JGL-DA and RDA. The cellwise robust counterparts are respectively r-LDA, r-QDA, rGL-LDA, rGL-QDA, rJGL-DA and rRDA. The former two are obtained by plugging the cellwise robust covariance matrices \mathbf{S}_k (or the pooled matrix \mathbf{S}_{pool}) directly into rule (1). Regularization parameters for the regularized methods are selected based on the BIC equation (8). For the non robust versions, we use $\widehat{\Sigma}_k$ instead of \mathbf{S}_k in the BIC equation (8).

Up to our knowledge, no comparison of all the considered discriminant methods has been made in the literature. Hence, it is worth analyzing their relative performances in depth. To this end, we consider both uncontaminated and contaminated settings.

Performance measures. For all the settings described below, we simulate 1000 training datasets consisting of K groups for each of which we estimate the mean and precision matrix. These estimates are then used to construct a discriminant rule. Additionally, for each training dataset, we generate a test dataset that consists of $N_{\text{test}} = 1000$ observations. For each observation of the test set, we randomly select (with equal probability) one of the K populations, then randomly draw a value from this population. We evaluate the discriminant methods in terms of *classification performance* and *estimation accuracy*.

To evaluate classification performance, we use the test datasets to compute the *average percentage of correct classification* over the 1000 simulation runs. The higher the average percentage, the better the classification performance.

To measure estimation accuracy, we report the *Kullback Leibler (KL) distance*. Under the normal model, the KL-distance from the model with estimated precision matrices $\widehat{\Theta}_1, \dots, \widehat{\Theta}_K$ to the model with true precision matrices $\Theta_1, \dots, \Theta_K$ is

$$\text{KL}(\widehat{\Theta}_1, \dots, \widehat{\Theta}_K; \Theta_1, \dots, \Theta_K) = \left(\sum_{k=1}^K -\log \det(\widehat{\Theta}_k \Theta_k^{-1}) + \text{tr}(\widehat{\Theta}_k \Theta_k^{-1}) \right) - Kp.$$

The lower the KL-distance, the more accurate the estimates. If all the estimates are equal to the true precision matrices, the KL-distance is equal to zero.

4.1 Uncontaminated scheme

We consider two different scenarios. In the first one, the number of groups is set to $K = 10$, with group sample sizes $n_k = 30$ and varying dimension $p = 5, 10, 30$. The precision matrices of the first five groups are equal. They have diagonal elements equal to one and zero off-diagonal elements except in the upper left block of size 2×2 , where the off-diagonal elements are set to 0.9. For the precision matrices of the next five groups, the 2×2 block with 0.9 off-diagonal elements is located in the lower right corner. The mean vectors of the first five groups have elements all equal to zero except element k in the k -th group that is equal to 3. The mean vectors of the remaining five groups have the opposite sign.

In the second scenario, the number of groups is set to $K = 6$, with group sample size $n_k = 30$ and dimension $p = 50$. The covariance matrices are chosen as in [12]: they are all diagonal with elements $(9(i-1)/(p-1)+1)^2$ for $i = 1, \dots, p$ in the first three groups, and $(9(p-i)/(p-1)+1)^2$ for $i = 1, \dots, p$ in the three other groups. The condition number is thus the same for all the groups but the low and high variance subspaces of groups 1 to 3 and 4 to 6 are complementary. All the elements of the mean vector of group k are set to zero except the k -th element in the first three groups, and the $p-k$ -th element in the three other groups that are equal to $\log(p)$.

Results. Table 1 gives the correct classification percentages and KL-distances for the non robust (top) and robust methods (bottom) in the two scenarios. Robust methods are generally expected to perform well also in uncontaminated settings. In both scenarios, the correct classification percentages and KL-distances of the robust techniques are, overall, quite similar to those of the non robust ones, regardless of the dimension p . The use of the cellwise robust discriminant methods (instead of the non robust ones) thus only results in a small statistical efficiency loss.

Next, we compare the standard discriminant methods s-LDA and s-QDA to the regularized ones GL-LDA, GL-QDA. For their cellwise robust counterparts, the same conclusions can be made. In a low dimension (Scenario 1, $p = 5$), all the methods show similar classification performance. As soon as the dimension increases (Scenario 1, $p = 10$ and $p = 30$), GL-LDA and GL-QDA perform much better. Regularization is necessary to improve both the

classification performance and estimation accuracy of the standard methods. For instance, in dimension $p = 10$, the standard quadratic classifier s-QDA suffers from low estimation accuracy, i.e. its KL-distance is 3 times that of its regularized version GL-QDA. This, in turn, negatively impacts its classification performance. In dimension $p = 30$, s-QDA is not computable anymore as indicated by NA in Table 1. Its regularized version, on the contrary, is always computable and yields good classification performance and high estimation accuracy. Also for the robust methods, although still computable in this setting³, r-QDA yields poor classification performance and estimation accuracy compared to rGL-QDA.

Further improvement over GL-LDA and GL-QDA can be obtained by using JGL-DA. The latter not only yields a high correct classification percentage but also the most accurate precision matrix estimates in each setting. When several groups share similar sparsity patterns (as in the considered simulation settings), estimation accuracy can be considerably improved by using a method that covers the path between LDA and QDA, like JGL-DA. In other (unreported) simulation studies expected to favour either GL-LDA or GL-QDA, we also find JGL-DA to be a tough competitor with respect to both classification performance and estimation accuracy. The results are available from the authors upon request. The same overall conclusions hold when comparing the cellwise robust discriminant methods.

Differences between the regularized methods are even more marked in the second scenario (see Table 1). JGL-DA still attains the best correct classification performance and estimation accuracy, followed by GL-QDA. The standard s-QDA is not computable since p is too large. The methods GL-LDA and RDA show poor classification performance and low estimation accuracy. In this scenario, the groups are more difficult to distinguish since they have complementary low and high variance subspaces, and the mean differences lie in the direction with the lowest variance. The precision matrix estimates used in GL-LDA, RDA (and s-LDA) artificially inflate the lowest variances by pooling and enlarging the smallest eigenvalues. This results in poor performance of these techniques. JGL-DA, on the contrary, lies in between LDA and QDA without shrinking the precision matrices towards each other as a whole. Only the coefficients that are similar across groups will be estimated identically. Hence, JGL-DA shows good performance in terms of both correct classification and estimation accuracy.

³Kendall's coefficient uses all the possible pairs of observations. As such, the corresponding covariance matrices may be invertible for $p > n_k$ but they become singular as soon $p > n_k(n_k - 1)/2$.

Table 1: Percentages of correct classification (CC) and KL-distance for the non robust discriminant methods (top) and their cellwise robust counterparts (bottom), averaged over 1000 simulation runs.

			s-LDA	s-QDA	GL-LDA	GL-QDA	JGL-DA	RDA
Scenario 1 $K = 10$	$p = 5$	CC	78.4	79.0	78.5	81.4	81.9	78.4
		KL	12.65	7.53	12.73	3.60	3.01	12.64
	$p = 10$	CC	82.7	76.6	83.4	85.6	86.1	82.8
		KL	14.06	39.36	13.46	13.60	3.68	14.00
	$p = 30$	CC	77.7	NA	80.5	83.0	83.5	75.4
		KL	30.29	NA	21.87	40.41	5.03	58.67
Scenario 2 $K = 6$	$p = 50$	CC	23.5	NA	25.7	54.5	71.4	25.7
		KL	223.90	NA	158.12	85.43	78.67	155.72

			r-LDA	r-QDA	rGL-LDA	rGL-QDA	rJGL-DA	rRDA
Scenario 1 $K = 10$	$p = 5$	CC	77.0	78.2	76.9	79.3	79.6	77.0
		KL	14.96	8.16	15.97	8.90	8.15	15.07
	$p = 10$	CC	81.4	69.3	82.2	84.5	85.1	81.4
		KL	14.57	104.11	14.00	13.88	4.44	14.51
	$p = 30$	CC	76.1	59.7	77.4	79.7	80.1	73.5
		KL	22.86	139.18	22.98	44.57	11.02	59.01
Scenario 2 $K = 6$	$p = 50$	CC	24.0	70.2	24.6	51.9	61.4	25.5
		KL	177.28	238.08	159.64	93.01	104.73	156.07

4.2 Contaminated scheme

We compare the performance of the non robust and robust discriminant methods in the presence of cellwise outliers. To this end, we add contamination to the settings from Section 4.1. In each training set, we randomly replace a given proportion of the cells in each group. The considered cellwise contamination percentages are $\varepsilon = 5\%$ and 10% in the first scenario and 1% in the second one. The test datasets are generated as in Section 4.1 without contamination.

In the first scenario, each contaminated cell is drawn from a normal distribution $N(-10, 0.2)$ in the first five groups and from $N(10, 0.2)$ in the five other groups. This shift contamination may drive the estimated means of the first five groups in the direction of the means of the remaining groups (and vice-versa) and inflate the sample covariance estimates. In the second scenario, each contaminated cell is drawn from a normal distribution with a large variance

$N(0, 50)$.

Results. We compare the non robust discriminant methods to their robust counterparts in terms of correct classification. Figure 1 shows the results for scenario 1 (left panel: 5% of contaminated cells, right panel: 10% of contaminated cells), Figure 2 for scenario 2. For each method, the boxplot of correct classification percentages of the 1000 simulation runs for the non robust version is displayed on the left while the boxplot on the right corresponds to its cellwise robust counterpart.

In all the considered contaminated settings, the cellwise robust methods maintain their good classification performance. On the contrary, the outlying cells mislead all the considered non robust methods. As a result, their correct classification percentages considerably decrease and their KL-distances (unreported) are high.

The higher the dimension p and/or the higher the contamination proportion ε , the better the performance of the cellwise robust estimators relative to their non robust version. For instance, keeping ε fixed to 5%, rJGL-DA leads to an increase in correct classification performance of (on average) 19 percentage points over JGL-DA in dimension $p = 5$, while this gain increases to 32 percentage points in dimension $p = 30$. Likewise, keeping the dimension $p = 5$ fixed but varying the proportion of contaminated cells, rJGL-DA improves classification performance by 19 percentage points over JGL-DA when $\varepsilon = 5\%$, and this gain doubles when $\varepsilon = 10\%$. The deteriorating performance of the non robust methods when the dimension increases is expected in this cellwise contamination scheme. Indeed, the probability that an observation has at least one contaminated cell is $1 - (1 - \varepsilon)^p$. In the first scenario, for 5% of contamination and $p = 5$, already 22.6% of the observations are expected to be contaminated, and more than 78% if $p = 30$. For 10% of cellwise contamination, the presence of contaminated cells is expected for nearly a half of the observations if $p = 5$, and almost all of them if $p = 30$.

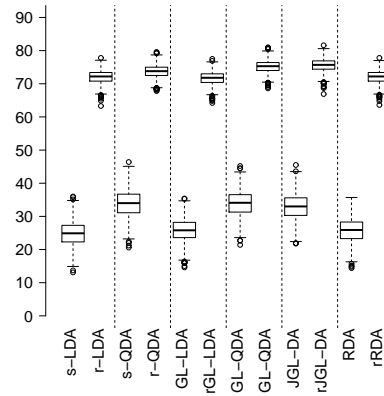
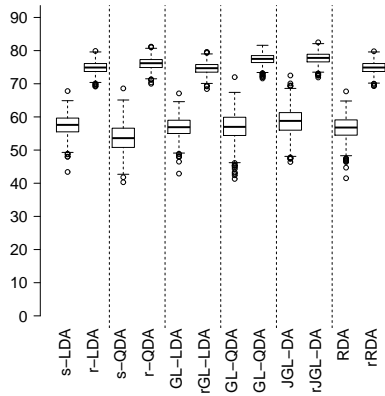
Results for the cellwise robust estimators are summarized in Table 2. The main findings are similar to those detailed in the uncontaminated case. As the dimension increases in Scenario 1, the regularized techniques rGL-QDA and rJGL-DA are among the best in terms of correct classification. rJGL-DA achieves considerably lower KL-distances than rGL-QDA. In the second scenario, r-QDA, rJGL-DA and rGL-QDA yield the best correct classification rates but the estimation accuracy of r-QDA is much worse than rJGL-DA and rGL-QDA.

5% cellwise contamination

10% cellwise contamination

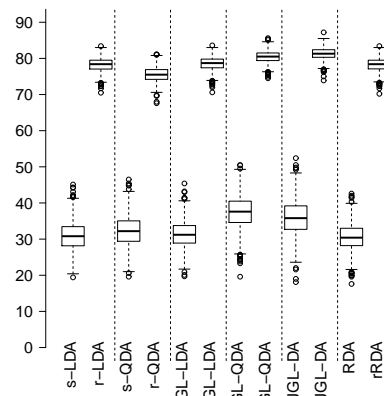
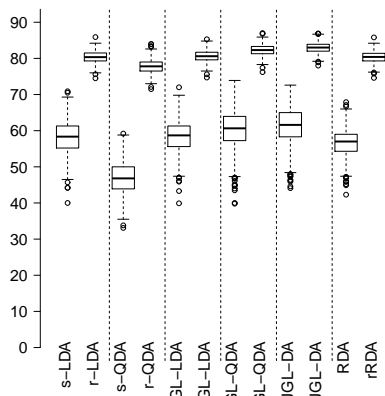
$p = 5$

$p = 5$



$p = 10$

$p = 10$



$p = 30$

$p = 30$

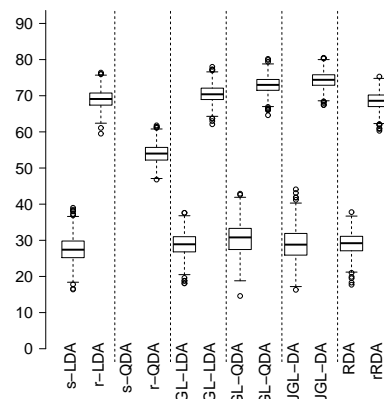
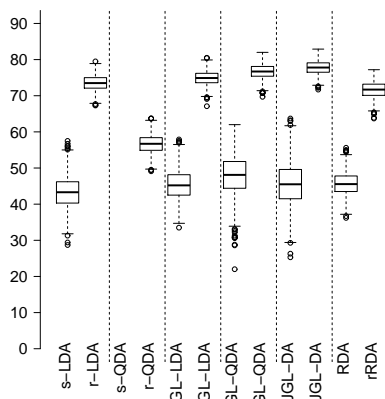


Figure 1: Scenario 1 : Boxplots of correct classification percentages. Left panel : 5% cellwise contamination, right panel: 10% cellwise contamination.

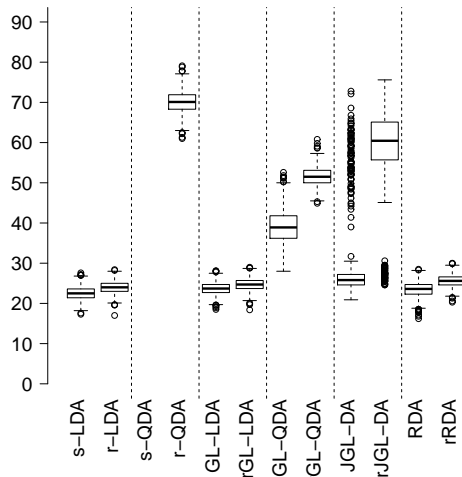


Figure 2: Scenario 2 : Boxplots of correct classification percentages. 1% cellwise contamination.

To summarize, the proposed cellwise robust discriminant methods have better classification performance and estimation accuracy than their non-robust counterparts in presence of cellwise outliers. Among the robust techniques, the regularized methods outperform the standard ones in high dimension. Among the cellwise robust regularized methods, rJGL-DA attains the best overall performance. This method relies on a similarity parameter that, when varied, covers the path from QDA to LDA. Choosing this parameter in a data-driven way makes rJGL-DA a tough competitor not only in presence of many similar precision matrices but also under LDA/QDA assumptions.

5 Examples

In this section, we illustrate the performance of the proposed methods on two real datasets. The first example, the *forest soil* dataset, is known in the robust discriminant analysis literature. The second example, the *phoneme* dataset, includes a large number of variables and has been used in papers on regularized discriminant analysis. We find most cellwise robust discriminant methods to attain better classification performance than their non robust counterparts. We also provide two visual tools showing the outlying observations and outlying cells in the analysed datasets.

The robust mean and precision matrix estimates proposed in this article can be used for

Table 2: Percentages of correct classification (CC) and KL-distance for the cellwise robust discriminant methods, averaged over 1000 simulation runs.

			r-LDA	r-QDA	rGL-LDA	rGL-QDA	rJGL-DA	rRDA
Scenario 1 - $\varepsilon = 5\%$ $K = 10$	$p = 5$	CC	74.9	76.1	74.6	77.4	77.7	74.9
		KL	20.83	12.530	22.32	15.31	14.07	20.95
	$p = 10$	CC	80.4	77.8	80.6	82.3	82.9	80.4
		KL	22.58	21.25	25.13	23.83	15.83	22.68
	$p = 30$	CC	73.5	56.7	74.9	76.7	77.8	71.6
		KL	33.90	116.03	40.26	70.25	23.21	68.95
Scenario 1 - $\varepsilon = 10\%$ $K = 10$	$p = 5$	CC	72.0	73.7	71.7	75.2	75.6	72.0
		KL	26.13	16.41	27.84	20.47	18.73	26.24
	$p = 10$	CC	78.3	75.5	78.6	80.4	81.3	78.3
		KL	30.78	25.54	34.20	33.40	23.38	30.93
	$p = 30$	CC	69.1	54.0	70.5	73.0	74.3	68.6
		KL	51.52	106.68	62.86	98.31	41.18	85.74
Scenario 2 - $\varepsilon = 1\%$ $K = 6$	$p = 50$	CC	24.0	70.0	24.7	51.6	58.6	25.6
		KL	176.60	227.57	161.81	98.79	117.88	156.82

outlier detection. Usual outlier detection methods are based on the computation of distances. The sample mean and covariance matrix are, however, heavily influenced by outliers. Therefore, distances computed from them may be large for the clean observations/cells and small for the outlying ones. As a result, the actual outliers are not detected. This effect is known as the *masking effect* [19]. It can be avoided by computing robust distances from robust location and precision matrix estimates like the ones proposed in Section 3. They allow to pinpoint both outlying observations, i.e. rowwise outliers, and outlying cells, i.e. cellwise outliers.

To find rowwise outliers in each group, robust Mahalanobis distances are computed for each observation \mathbf{x}_i of the group,

$$D_i = \sqrt{(\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k)^T \hat{\boldsymbol{\Theta}}_k (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k)},$$

where $\hat{\boldsymbol{\mu}}_k$ is the vector of marginal medians in group k and $\hat{\boldsymbol{\Theta}}_k$ is one of the precision matrix estimates defined in Section 3. Observations with a Mahalanobis distance above the square root of the 0.99 quantile of the chi-square distribution with p degrees of freedom are flagged as rowwise outliers.

To find cellwise outliers in each group, we follow the approach of [1]. For each cell x_{ij} in

group k , a cellwise standardized distance is computed

$$d_{ij} = \frac{x_{ij} - m_{k,j}}{t_{k,j}}, \quad (10)$$

where $m_{k,j}$ and $t_{k,j}$ estimate respectively the marginal location and scale of the j th variable in group k . We replace $m_{k,j}$ in (10) by the median of the j th variable in group k . For the scale, we replace $t_{k,j}$ by the square root of element (j, j) of $\hat{\Theta}_k^{-1}$. Cells with standardized distance exceeding the square root of the $0.99^{1/(n_k p)}$ quantile of a chi-square distribution with one degree of freedom are flagged as cellwise outliers.

5.1 Forest soil data

The *forest soil* dataset, available in the R-package `rrcovHD` [21], contains measurements on $N = 58$ soil pits in the Hubbard Brook Experimental Forest in north-central New Hampshire of 1983. For each soil sample, the exchangeable cations of calcium, magnesium, potassium and sodium ($p = 4$) are reported. The pit location can be classified in $K = 3$ types of forest: spruce-fir ($n_1 = 11$), high elevation hardwood ($n_2 = 23$) and low elevation hardwood ($n_3 = 24$). Some unusual soil samples are present in the dataset, as already noticed in [23]. Note that the group sample sizes are low compared to the dimension p . We compare the classification performance of the proposed cellwise robust and non robust discriminant methods. The robust location and precision matrix estimates are then used to construct an outlier detection map.

Classification performance. As the sample size is low, the same dataset is used to construct and evaluate the discriminant rule. Table 3 summarizes the percentages of correct classification for the non robust and robust methods. This dataset is characterized by strong overlapping groups, which causes overall low correct classification rates. We observe that the robust discriminant methods reduce the influence of the unusual soil samples and yield better correct classification rates. Although the dimension is not high, the regularized techniques rJGL-DA and rGL-QDA are to be preferred since the sample size is low.

Outlier detection. Since the cellwise robust methods perform better, outliers might be present. To characterize the rowwise outliers, we compute robust Mahalanobis distances. To characterize the cellwise outliers, we compute robust cellwise standardized distances. Figure 3 visualizes the detected outliers using the rGL-LDA (left panel), rGL-QDA (middle) and

Table 3: Percentage of correct classification. Forest soil data, $K = 3$, $p = 4$, $N = 58$.

s-LDA	s-QDA	GL-LDA	GL-QDA	JGL-DA	RDA
56.9	56.9	56.9	55.1	56.9	56.9
r-LDA	r-QDA	rGL-LDA	rGL-QDA	rJGL-DA	rRDA
60.3	63.8	60.3	65.6	67.2	60.3

rJGL-DA (right panel) estimates. Observations flagged as rowwise outliers are colored in blue, cellwise outliers in red. The black vertical lines split the observations according to their group membership: spruce-fir (top), high elevation hardwood (middle) and low elevation hardwood (bottom).

The same four rowwise outliers are highlighted by the three estimation methods. They all have a sodium measurement (last column) that is much higher than expected, while their other components are in line with the majority of the data. Unlike standard robust methods, the proposed cellwise discriminant methods result in less loss of information since they do not drop the entire row because of the presence of only one outlying component. The methods also detect, overall, the same outlying cells, mainly among the measurements for sodium. Used together, these rowwise and cellwise outlier detection techniques allow a deeper comprehension of the dataset.

5.2 Phoneme data

The *phoneme* dataset [14] contains $N = 1717$ observations corresponding to the record of a male voice pronouncing one of $K = 2$ similar sounds, either *aa* or *ao*. The aim is to build a classifier of these sounds on the basis of the $p = 256$ log-periodograms representing the log intensity of the sound across 256 frequencies.

Classification performance. Since there are enough observations in each group, we split the dataset into a training and a test set (with 60%/ 40% of the observations respectively). We evaluate the performance of the methods on the test set. To diminish the influence of the split, this procedure is repeated ten times. The average percentage of correctly classified observations is computed and reported in Table 4. For the majority of the considered methods, the robust ones perform better than their non robust counterparts. Among the cellwise robust techniques, the regularized methods rGL-LDA and rGL-QDA improve classification performance compared to respectively r-LDA and r-QDA.

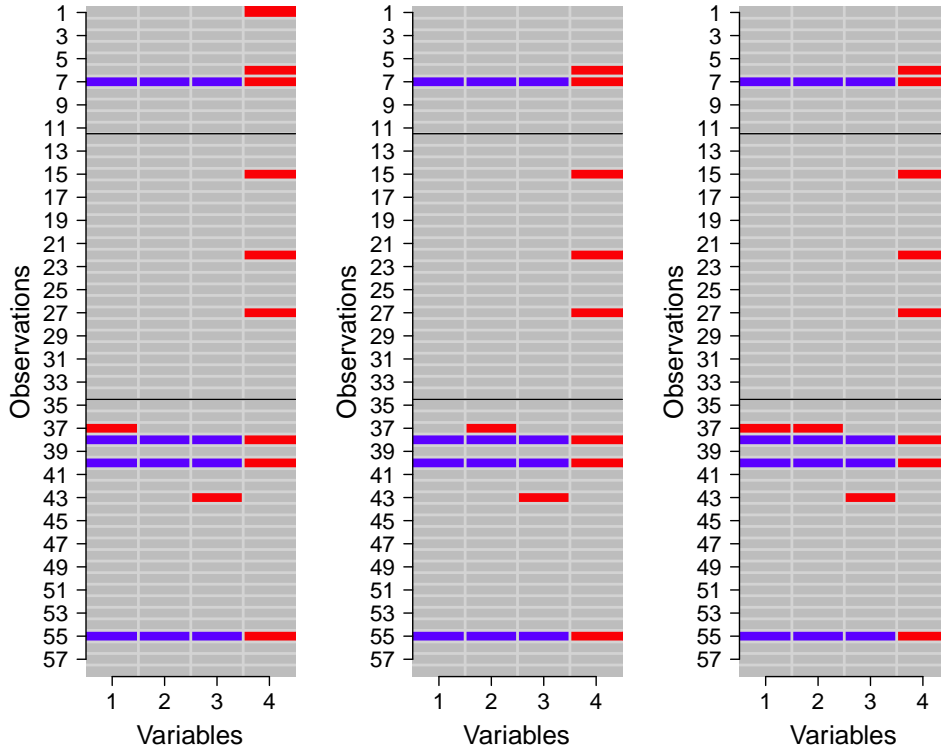


Figure 3: Outlier detection map of Forest soil data: detected rowwise outliers (blue) and cellwise outliers (red) by rGL-LDA (left), rGL-QDA (middle) and rJGL-DA (right).

Outlier detection. As in the first data example, the robust mean and precision matrix estimates can be used to detect outliers. Here, the dataset is however too large for a clear visual inspection via the outlier detection map (cfr. Figure 3). Another useful visualisation tool to detect rowwise outliers is the plot of the robust squared Mahalanobis distances versus the observation numbers, as represented in Figure 4. The distances are computed with the rGL-LDA estimates in the left panel, and with the rGL-QDA and rJGL-DA estimates in the middle and right panels. The vertical line corresponds to the 0.99 chi-square quantile with $p = 256$ degrees of freedom. Observations beyond this threshold are considered as rowwise outliers. For all the methods, outliers are detected in both groups. The same extreme outlying rows are highlighted by all the estimation methods.

The different methods also detect overall the same outlying cells (not shown), although JGL-DA highlights many more cellwise outliers. Again, the multivariate outlying behaviour

Table 4: Average percentage of correct classification. Phoneme dataset $K = 2$, $p = 256$, $N_{\text{train}} = 1030$, $N_{\text{test}} = 687$.

s-LDA	s-QDA	GL-LDA	GL-QDA	JGL-DA	RDA
77.7	62.4	81.4	74.9	78.4	78.2
r-LDA	r-QDA	rGL-LDA	rGL-QDA	rJGL-DA	rRDA
81.1	74.7	81.7	76.0	76.7	73.3

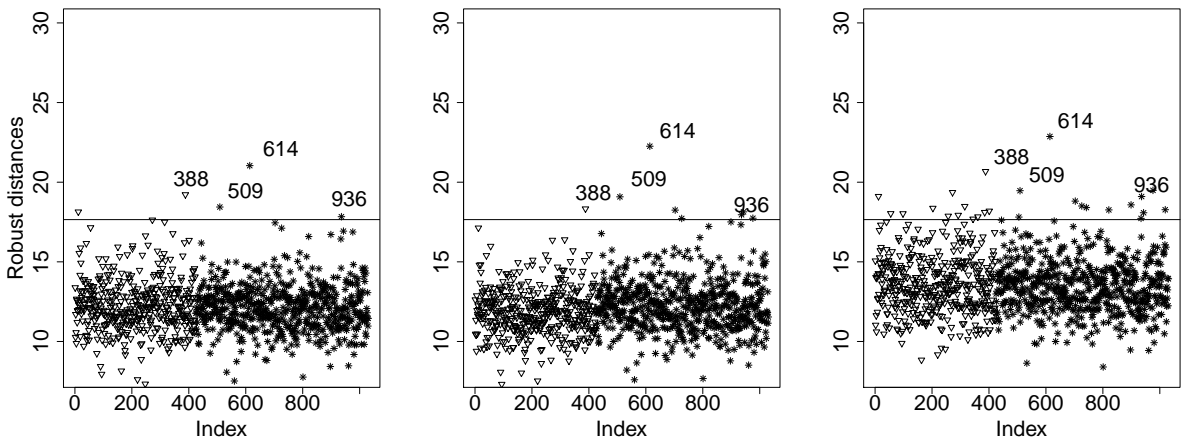


Figure 4: Robust Mahalanobis distances computed with the rGL-LDA (left), the rGL-QDA (second panel) and the rJGL-DA (right panel) estimates. Observations in the first group are represented by a ∇ and those in the second by a $*$.

of some entire rows may be explained by only one abnormal component, as is the case for observations 388, 509 and 614. On the contrary, the extreme rowwise outlier 936 has a high Mahalanobis distance while none of its components is flagged by the cellwise outlier detection procedure. Its outlying behaviour is thus caused by a particular relation between the components rather than by one of its components. Hence, it is important to consider both the rowwise and cellwise outlier detection procedures in combination.

6 Discussion

In this paper, cellwise robust discriminant analysis methods are proposed. We discuss all the implementation issues and we make the code publicly available. The proposed discriminant methods enjoy several important advantages. They are robust against cellwise outliers, a type

of outliers that is very likely to occur in high-dimensional datasets. We provide visual tools to detect both rowwise and cellwise outliers (cfr. Section 4). Furthermore, contrary to LDA that makes the strong assumption of homoscedasticity, we consider methods that emphasizes the true difference between the group covariance matrices while making use of their similarities. As such, the proposed methods lie in between LDA and QDA. Our simulations show that these approaches result in better classification performance as well as higher estimation accuracy. Finally, many discriminant methods require dimension reduction techniques to be computable in high-dimension (see [15], [23]). Our methods, in contrast, are computable even when the dimension exceeds the group sample size without requiring any preprocessing step. By using sparse precision matrix estimates, we reduce the effect of uninformative variables.

The proposed regularized methods depend on regularization parameters that are tuned in a data-driven way. In this paper, we focus on parameter selection via minimization of the BIC. In the classification context, one may also be interested in selecting parameters so as to maximize the expected out-of-sample correct classification rate, which can be obtained by L -fold cross validation. This method is, however, much more time consuming. Furthermore, in our simulations, it, overall, does not outperform BIC model selection in terms of correct classification. Even more, parameter selection by L -fold cross validation sometimes achieve higher KL-distances.

Acknowledgements. We thank Prof. C. Croux who provided insights and critical comments that helped improving the manuscript. We gratefully acknowledge support from the FWO (Research Foundation Flanders, contract number 12M8217N).

References

- [1] C. Agostinelli, A. Leung, C.J. Yohai, and R. H. Zamar. Robust estimation of multivariate location and scatter in the presence of cellwise and casewise contamination. *TEST*, 24(3):441–461, 2015.
- [2] F. A. Alqallaf, Martin R.D. Komis, K. P., and R.H. Zamar. Scalable robust covariance and correlation estimates for data mining. *Proceedings of the Eighth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 14–23, 2002.

- [3] F.A. Alqallaf, S. Van Aelst, V.J. Yohai, and R.H. Zamar. Propagation of outliers in multivariate data. *The Annals of Statistics*, 37:311–331, 2009.
- [4] C. Croux and C. Dehon. Robust linear discriminant analysis using S-estimators. *The Canadian Journal of Statistics*, 29(3):473–493, 2001.
- [5] C. Croux and V. Öllerer. *Modern Multivariate and Robust Methods*, chapter Robust high-dimensional precision matrix estimation. Springer, 2015.
- [6] P. Danaher. *JGL: Performs the Joint Graphical Lasso for sparse inverse covariance estimation on multiple classes*, 2013. URL <https://CRAN.R-project.org/package=JGL>. R package version 2.3.
- [7] P. Danaher, P. Wang, and D. Witten. The joint graphical lasso for inverse covariance estimation across multiple classes. *Journal of the Royal Statistical Society, Series B*, 76: 373–397, 2014.
- [8] P. Filzmoser and H. Fritz. *pcaPP : Robust PCA by Projection Pursuit*, 2006. URL <https://CRAN.R-project.org/package=pcaPP>. R package version 1.0.
- [9] P. Filzmoser, R. Maronna, and M. Werner. Outlier identification in high dimension. *Computational Statistics and Data Analysis*, 52:1694–1711, 2008.
- [10] P. Filzmoser, K. Hron, and M. Templ. Discriminant analysis for compositional data and robust parameter estimation. *Computational Statistics*, 27:585–604, 2012.
- [11] J. Friedman, T. Hastie, and R. Tibshirani. Sparse inverse covariance estimation with the graphical lasso. *Biostatistics*, 9(3):432–441, 2008.
- [12] J. H. Friedman. Regularized discriminant analysis. *Journal of the American Statistical Association*, 84:165–175, 1989.
- [13] C. Gao, Y. Zhu, X. She, and W. Pan. Estimation of multiple networks in gaussian mixture models. *Electronic Journal of Statistics*, 10:1133–1154, 2016.
- [14] T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning, Data Mining, Inference and Prediction, Second Edition*. Springer Verlag, New York, 2009.
- [15] M. Hubert and S. Engelen. Robust PCA and classification in biosciences. *Bioinformatics*, 20(11):1728–1736, 2004.

- [16] M. Hubert and K. Van Driessen. Fast and robust discriminant analysis. *Computational Statistics and Data Analysis*, 45(2):301–320, 2004.
- [17] B. Price, C. Geyer, and A. Rothman. Ridge fusion in statistical learning. *Journal of Computational and Graphical Statistics*, 24(2):439–454, 2015.
- [18] P. Rousseeuw and C. Croux. Alternatives to the median absolute deviation. *Journal of the American Statistical Association*, 88(424):1273–1283, 1993.
- [19] P. J. Rousseeuw and A.M. Leroy. *Robust regression and outlier detection*. John Wiley and Sons, New-York, 1987.
- [20] G. Tarr, S. Müller, and N.C. Weber. Robust estimation of precision matrices under cellwise contamination. *Computational Statistics and Data Analysis*, 93:404–420, 2015.
- [21] V. Todorov. *rrcovHD : Robust Multivariate Methods for High dimensional data*, 2016. URL <https://CRAN.R-project.org/package=rrcovHD>. R package version 0.2-5.
- [22] S. Van Aelst. Stahel-Donoho estimation for high dimensional data. *International Journal of Computer Mathematics*, 93:628–639, 2016.
- [23] K. Vanden Branden and M. Hubert. Robust classification in high dimension based on the SIMCA method. *Chemometrics and Intelligent Laboratory Systems*, 79:10–21, 2005.
- [24] B. Xu, K. Huang, King I., C. Liu, J. Sun, and N. Satoshi. Graphical lasso quadratic discriminant function and its application to character recognition. *Neurocomputing*, 129:33–40, 2014.
- [25] T. Yuan and J. Wang. A coordinate descent algorithm for sparse positive definite matrix estimation. *Statistical Analysis and Data Mining*, 6(5):431–442, 2013.
- [26] T. Zhao, H. Liu, K. Roeder, J. Lafferty, and L. Wasserman. The huge package for high-dimensional undirected graph estimation in R. *Journal of Machine Learning Research*, 13:1059–1062, 2012.
- [27] A. Zimek, E. Schubert, and H-P. Kriegel. A survey on unsupervised outlier detection in high dimensional numerical data. *Statistical Analysis and Data Mining*, 5:363–476, 2012.

FACULTY OF ECONOMICS AND BUSINESS
Naamsestraat 69 bus 3500
3000 LEUVEN, BELGIË
tel. + 32 16 32 66 12
fax + 32 16 32 67 91
info@econ.kuleuven.be
www.econ.kuleuven.be

