Robust estimation of precision matrices under cellwise contamination

G. Tarr*, S. Müller, N. C. Weber

School of Mathematics and Statistics, University of Sydney, NSW 2006, Australia

Abstract

There is a great need for robust techniques in data mining and machine learning contexts where many standard techniques such as principal component analysis and linear discriminant analysis are inherently susceptible to outliers. Furthermore, standard robust procedures assume that less than half the observation rows of a data matrix are contaminated, which may not be a realistic assumption when the number of observed features is large. This work looks at the problem of estimating covariance and precision matrices under cellwise contamination. We consider using a robust pairwise covariance matrix as an input to various regularisation routines, such as the graphical lasso, QUIC and CLIME. To ensure the input covariance matrix is positive semidefinite, we use a method that transforms a symmetric matrix of pairwise covariances to the nearest covariance matrix. The result is a potentially sparse precision matrix that is resilient to moderate levels of cellwise contamination. Since this procedure is not based on subsampling it scales well as the number of variables increases.

Keywords: Precision matrix, Covariance matrix, Robust estimation, Data mining 2010 MSC: 62G35, 62H20, 62H30

1. Introduction

Often the aim of data mining and statistics is to extract information about the relationships between the variables and identify any features or structure in the data. The covariance matrix, $\Sigma = var(\mathbf{y})$, where $\mathbf{y} \sim \mathbf{F}$, the distribution of the true data generating process, and its inverse, the precision matrix $\mathbf{\Theta} = \Sigma^{-1}$ are fundamental components of many statistical procedures, such as principal component analysis (PCA) and linear discriminant analysis. However, it is well known that the classical covariance matrix is inherently non-robust to outliers and suffers from distortion in its eigenstructure in high dimensions (Johnstone, 2001). This paper combines pairwise covariance matrix estimation with recent regularisation routines currently used in bioinformatics and machine learning to produce an estimated precision matrix that is robust to moderate levels of cellwise contamination.

The need for robust statistics in data mining and associated fields is well known, see Barnett and Lewis (1994) for a general overview. In particular, it is desirable for learning algorithms to be stable with respect to noisy features and unusual fluctuations in the inputs. For example Li (2004) considers robust incremental PCA applied to multi-view face modelling and Mavroeidis and Marchiori (2014) consider the stability of sparse PCA in the context of feature selection in microarray gene expression data. Other situations where

^{*}Corresponding author

Email address: gtar4178@uni.sydney.edu.au (G. Tarr)

robust techniques are important include speech recognition and neural networks, see Gales and van Dalen (2007) and Bieroza et al. (2011), respectively.

In the statistics literature, robust estimation of covariance matrices has received much attention in the past, notably the minimum volume ellipsoid and minimum covariance determinant (MCD) estimators, projection type estimators and M-estimators, see Hubert et al. (2008) for a survey. Furthermore, research into covariance matrix estimation and its applications is ongoing, see for example Filzmoser et al. (2014) who use the MCD estimator to construct robust Mahalanobis distances to identify local multivariate outliers; Hubert et al. (2014) who study the shape bias of a range of existing robust covariance matrix estimators; or Cator and Lopuhaä (2010, 2012) who consider asymptotic expansions and establish asymptotic normality for general MCD estimators.

An alternative approach is to estimate the covariance matrix in a component-wise manner based on a robust estimator of scale as outlined by Ma and Genton (2001). It is well known that the resulting symmetric matrix is not guaranteed to be positive definite (PD). Methods to ensure the resulting estimator is PD have previously been explored by Rousseeuw and Molenberghs (1993) with notable updates in the robustness literature by Maronna and Zamar (2002) and quite separately in the finance literature by Higham (2002). Alqallaf et al. (2002) also proposed a pairwise approach to covariance matrix estimation by means of first Winsorising the data. The resulting Quadrant Covariance estimate does not necessarily require a transformation to ensure the result is positive definite.

In practice, it is often the precision matrix, the inverse of the covariance matrix, that is primarily of interest. This is the case, for example, in Gaussian graphical model selection. As such, this paper is primarily concerned with robustly estimating the precision matrix. While there is an obvious link between covariance matrices and precision matrices, it is not obvious that a good (robust) estimator for one results in a good estimator for the other. We will employ robust pairwise covariance matrices as a starting point for various regularisation techniques to facilitate the estimation of robust, potentially sparse, precision matrices.

Classical robust estimators assume that contamination occurs within a restricted subset of the observation vectors, however, in recent years there has been interest in developing robust estimators that perform well under cellwise contamination. The cellwise contamination model was initially explored in a data mining context by Alqallaf et al. (2002) and later defined comprehensively by Alqallaf et al. (2009). This form of contamination is prevalent in large, automatically generated data sets, found in data mining and bioinformatics, where there is often little quality control over the inputs. Cellwise contamination is common in the context of missing data, however, it represents a philosophical divergence from the traditional approach to robustness. Recent examples where the problem of cellwise contamination have arisen include, Farcomeni (2014), Van Aelst et al. (2012) and Agostinelli et al. (2014).

We perform a detailed simulation study to assess the performance of a variety of precision matrix estimators in the presence of cellwise contamination over a number of scenarios and levels of p while keeping the sample size fixed. Our results are distilled from a comprehensive range of performance indices. We outline these indices and consider their applicability to the various scenarios in the supplementary material accompanying this article.

We show that the pairwise nature of the covariance estimates enables the resulting precision matrix to have a higher level of robustness than when using standard robust covariance matrix estimation procedures in the presence of cellwise contamination. This is a novel result and a significant first step towards dealing with cellwise contamination in this context.

The remainder of this paper is structured as follows. Section 2 outlines the cellwise contamination model and highlights why standard robust techniques fail in this setting. Sections 3 and 4 outline the theory for existing pairwise covariance matrix estimation techniques and regularisation routines and we propose a

new procedure which combines robust pairwise covariance matrix estimation with regularisation. Sections 5 and 6 present the results of an extensive simulation study and Section 7 summarises the important findings.

2. Cellwise contamination

Consider a data set $\mathbf{X} \in \mathbb{R}^{n \times p}$ consisting of *n* observations on *p* variables. Classically, even the most robust procedures are designed such that they only work when at most half of the rows in \mathbf{X} have contamination present.

Alqallaf et al. (2009) formally outline the cellwise contamination model as an extension of the standard Tukey-Huber contamination model which was first introduced in the univariate location-scale setup (Tukey, 1962; Huber, 1964). Consider the data generating process for the *n* rows in **X**, $\mathbf{x}_i = (\mathbf{I} - \mathbf{B}_i)\mathbf{y}_i + \mathbf{B}_i\mathbf{z}_i$, where $\mathbf{y}_i \sim \mathbf{F}$, the distribution of well-behaved data, $\mathbf{z}_i \sim \mathbf{G}$, some outlier generating distribution and $\mathbf{B}_i = \text{diag}(B_1, \ldots, B_p)$ is a diagonal matrix, where B_1, \ldots, B_p are Bernoulli random variables, $B_j \sim \mathcal{B}(1, \varepsilon_j)$. When **y**, **B** and **z** are independent we have a situation that is similar to the missing completely at random model, where the missingness does not depend on the values of **y**, see, for example, Little and Rubin (2002).

The structure of \mathbf{B}_i determines the contamination model. If B_1, \ldots, B_p are fully dependent, then $\mathbf{B}_i = U_i \mathbf{I}$, where $U_i \sim \mathcal{B}(1, \varepsilon)$, and we recover the fully dependent contamination model, the standard model on which classical robust procedures are based. In this setting, the probability that an observation is uncontaminated, $1 - \varepsilon$, is independent of the dimensionality. Furthermore, the proportion of contaminated observations is preserved under affine equivariant transformations.

In contrast, if B_1, \ldots, B_p are mutually independent we have the fully independent contamination model, where each element of \mathbf{x}_i is drawn from \mathbf{F} or \mathbf{G} independently of the other p - 1 elements in \mathbf{x}_i . That is, contaminating observations occur independently at the univariate level. In this setting, it may be be unreasonable to assume that less than half the rows have contamination. Furthermore, if p is large and there is only one outlier in an observation vector, then down-weighting the entire observation may be wasteful.

If the data matrix is randomly contaminated in this elementwise manner, as the number of variables increases, the chance that more than half the rows are contaminated increases exponentially. Formally, let ε be the probability that any particular element in a data matrix is contaminated. Assuming the contamination is randomly scattered throughout the data matrix, the probability that any particular row has no contamination is $(1 - \varepsilon)^p$, which quickly decays towards zero even for small values of ε . For example, if p = 30 and $\varepsilon = 0.1$, then the probability that any particular row remains uncontaminated is only 4%. This is demonstrated graphically in Figure 1. The plot on the left shows a 100×30 data matrix where 10% of the cells have been contaminated, the white cells. While virtually all the rows of the data matrix have at least one contaminated element, the majority of cells remain uncontaminated in the sense that they are still real measurements from the underlying data generating process. Even if $\varepsilon = 0.03$, the probability that any particular row is uncontaminated is 40%, however with a sample size of 100, this translates to a 98% chance that at least half the rows are contaminated, in which case standard robust methods fail.

It is important to note that the fully independent contamination model lacks affine equivariance, in the sense that linear combinations of columns of a contaminated data set result in "outlier propagation" (Alqallaf et al., 2009). As such, affine equivariance is not an achievable outcome for any estimator in this setting.

Existing research into the problem of cellwise contamination has focussed on coordinatewise procedures, that only operate on one column at a time. Croux et al. (2003) consider an approach based on "alternating regressions" using weighted L_1 regression, Maronna and Yohai (2008) use a coordinatewise procedure for principal component analysis. Liu et al. (2003) have an application involving the singular



Figure 1: On the left, a heat map of a data matrix with 30 variables and 100 observations. 10% of the cells have been contaminated and are shown as white cells, while the uncontaminated cells are in various shades of grey. On the right, the probability that any particular row (observations) in the data matrix will be contaminated, $1 - (1 - \varepsilon)^p$, over a range of ε , the proportion of cells affected by cellwise contamination.

value decomposition of microarray data and De la Torre and Black (2001) consider cellwise contamination in the context of computer vision.

We show that a pairwise approach is able to cope with much higher levels of cellwise contamination than existing classical robust estimators. In the simulations in Section 5 we do not use the fully independent contamination model, rather, we impose restrictions on the amount of contamination in each variable. As such the contamination is no longer strictly independent, however, the advantage is that we are able to assess the impact over various known levels of contamination in each variable.

3. Pairwise covariance matrix estimation

A pairwise approach to estimating covariance matrices in the presence of cellwise contamination has previously been explored by Alqallaf et al. (2002) where the classical correlation coefficient was applied to a Winsorised data set. Instead of transforming the underlying data, our approach is to take the p(p-1)/2 pairs of variables and robustly estimate the covariance between each pair. The primary advantage of this approach is robustness to cellwise contamination in the data set. The main disadvantage is that the resulting symmetric matrix is not guaranteed to be positive semidefinite or affine equivariant. However, as noted

earlier, in the cellwise contamination model, affine equivariance is unachievable as there is the potential for all rows to have a contaminated cell, hence linear combinations of the rows propagate the contamination.

A simple method for turning scale estimators into covariance estimators was introduced by Gnanadesikan and Kettenring (1972) and brought to prominence in the context of robust estimation by Ma and Genton (2001). The idea is based on the identity,

$$\operatorname{cov}(X,Y) = \frac{1}{4\alpha\beta} \left[\operatorname{var}(\alpha X + \beta Y) - \operatorname{var}(\alpha X - \beta Y) \right],\tag{1}$$

where X and Y are random variables. In general, X and Y may have different scales, hence it is standard to let $\alpha = 1/\sqrt{\operatorname{var}(X)}$ and $\beta = 1/\sqrt{\operatorname{var}(Y)}$. A robust covariance estimator is found by replacing the variance in (1) with (squared) robust scale estimators. We will focus on the estimators Q_n (Rousseeuw and Croux, 1993), the τ -scale as described in Maronna and Zamar (2002) and an estimator that is somewhat less robust but highly efficient, P_n , the interquartile range of the pairwise means, and its adaptively trimmed variant \tilde{P}_n with adaptive trimming parameter d = 3 (Tarr et al., 2012). We also consider the interquartile range (IQR) and the median absolute deviation from the median (MAD). A recent discussion on the efficiency of various robust scale estimators can be found in Tarr et al. (2012).

Using identity (1), a symmetric matrix full of pairwise covariances can easily be constructed, however, there is no guarantee that the result will be positive semidefinite. The two methods outlined below overcome this limitation by appropriately adjusting the eigenvalues of the symmetric matrix to ensure that they are all positive, and hence ensuring a positive definite result.

3.1. Orthogonalised Gnanadesikan Kettenring procedure

To overcome the possible lack of positive semidefiniteness in a matrix of pairwise covariances, Maronna and Zamar (2002) propose a modification based on the observation that the eigenvalues of a covariance matrix are the variances along the directions given by the respective eigenvectors. Essentially a principal components decomposition is performed and the covariance matrix is reconstructed using robust variance estimates of the principal component vectors in place of the original eigenvalues. This procedure is known as the Orthogonalised Gnanadesikan Kettenring (OGK) estimator. Note that even if the original covariance matrix was already positive definite, applying the OGK procedure will not necessarily return the same matrix.

Maronna and Zamar (2002) and Maronna et al. (2006, p. 207) suggest that the OGK estimator can be improved by iterating the procedure and then using this estimate to find robust Mahalanobis distances for each observation vector. These are then used to screen for outliers before applying the classical covariance estimator to the cleaned data, resulting in a procedure known as the reweighted OGK. This is done in an effort to increase efficiency and to make the result "more equivariant". In terms of the impact of not being affine equivariant, Maronna and Zamar (2002) note that "although the worst case may differ from the original data, for most transformations the results are very similar" and "the lack of equivariance is not a serious concern in our estimates".

Regardless, neither the OGK method nor the reweighted OGK method is able to cope with cellwise contamination. The issue of outlier propagation means that the number of contaminated principal components could easily be greater than 50% even for small levels of cellwise contamination. Hence, the robust variance estimates that are used in place of the eigenvalues will no longer be valid estimates – they will be in breakdown. Furthermore, the reweighting step will often needlessly exclude many observation vectors where there is only one contaminated cell.

3.2. Nearest positive definite matrix procedure

Higham (2002) considers the problem of computing the nearest positive definite (NPD) matrix to a given symmetric matrix. The motivation stems from finance, where sample covariance matrices are constructed from vectors of stock returns, however, the problem arises when not all stocks are observed every day. In this setting, classical covariances may be computed on a pairwise basis using data drawn only from days where both stocks have data available. The resulting covariance matrix is not guaranteed to be PD because it has been built from inconsistent data sets. Motivated by the same problem, Løland et al. (2013) propose both a pseudo-likelihood and a Bayesian approach to find PD estimates of pairwise covariances. However, their approach relies on expert knowledge to formulate priors for the pairwise covariances.

The NPD procedure is similar to the OGK procedure in that it performs a spectral decomposition and then updates the eigenvalues to ensure that the result is PD. However, it does not rely on linear transformations of the original dataset and hence is not affected by the "outlier propagation" issue associated with cellwise contamination. Formally, for an arbitrary symmetric $p \times p$ matrix **A**, the aim is to find the distance

$$\gamma(\mathbf{A}) = \min\{\|\mathbf{A} - \mathbf{W}\|_F : \mathbf{W} \text{ is a symmetric PD matrix}\},\tag{2}$$

and the resulting matrix that achieves this minimum distance. Higham (2002) uses the Frobenius norm, $\|\mathbf{B}\|_F = \sqrt{\operatorname{tr}(\mathbf{B'B})}$, as it is "the easiest norm to work with for this problem and also being the natural choice from the statistical point of view".

While Higham (2002) considers a variety of weighting mechanisms, in the simplest case without specifying any weights, the procedure is quite straightforward. The final estimate is $\hat{\mathbf{W}} = \mathbf{E}\hat{\mathbf{A}}\mathbf{E}'$, where $\mathbf{E}\mathbf{A}\mathbf{E}'$ is the spectral decomposition of \mathbf{A} , with $\mathbf{A} = \text{diag}(\lambda_1, \dots, \lambda_p)$ and $\hat{\mathbf{A}} = \text{diag}(\max\{\lambda_i, \delta\})$, where δ is a small positive constant. In contrast to the OGK procedure, if the initial symmetric matrix is already PD, then the NPD method simply returns the original pairwise covariance matrix.

In the presence of cellwise contamination the NPD method outperforms the OGK method. However, the NPD method often results in estimated matrices with a number of extremely small eigenvalues which give poorly conditioned estimates, i.e. the condition number of these estimators is very high as is the entropy loss, which involves the log of the eigenvalues. In general, it is not recommended to use either the OGK nor the NPD in isolation when there is cellwise contamination present. Even in the presence of standard row-wise contamination, the NPD method is not recommended due to its propensity to return poorly conditioned estimates.

4. Precision matrix estimation

Many statistical procedures are primarily concerned with the precision matrix, the inverse of a covariance matrix, rather than the covariance matrix itself. For example, finding Mahalanobis distances and performing linear discriminant analysis both require an estimate of $\Theta = \Sigma^{-1}$. Finding good precision matrix estimates has been a focus of many investigators over a long period of time, the first major contribution being Dempster (1972).

The following routines take as an input an estimated covariance matrix and output a regularised precision matrix. In Section 5 we demonstrate the advantages of using a robust pairwise covariance matrix estimate as the input to these regularisation routines.

4.1. GLASSO

A natural way to estimate Θ is by maximising the log-likelihood of the data. With Gaussian observations, the log-likelihood takes the form,

$$\log |\Theta| - tr(S\Theta), \tag{3}$$

where **S** is an estimate of the covariance matrix of the data. Maximising (3) with respect to Θ leads to the MLE, **S**⁻¹. In general, **S**⁻¹ will not be sparse, in the sense that it will contain no elements exactly equal to zero. Furthermore in p > n situations **S** will be singular so the MLE cannot be computed. Yuan and Lin (2007) consider minimising the penalised negative log-likelihood,

$$\operatorname{tr}(\mathbf{S}\boldsymbol{\Theta}) - \log |\boldsymbol{\Theta}| + \lambda \sum_{i,j} |\theta_{ij}|, \tag{4}$$

over the set of PD matrices where λ is a tuning parameter to control the amount of shrinkage. Friedman et al. (2008) refer to this estimator as the graphical lasso (GLASSO) and note that it has two major advantages over (3): the solution is PD for all $\lambda > 0$ even if **S** is singular, and for large values of λ the resulting estimate, $\hat{\Theta}$, will be sparse.

4.2. QUIC

The QUIC method solves the same minimisation problem as the GLASSO. The improvement in speed comes from noticing that the Gaussian log-likelihood component of (4) is twice differentiable and strictly convex which lends itself to a quadratic approximation and hence faster convergence (Hsieh et al., 2011). On the other hand, the penalty term is convex but not differentiable and so is treated separately.

The QUIC routine, as implemented in the R package QUIC, explicitly includes a step that ensures positive definiteness of the precision matrix for each iteration. Work has recently been undertaken to scale the QUIC estimator to scenarios with a million variables, see Hsieh et al. (2013) for details.

4.3. CLIME

An alternative to maximising the penalised log-likelihood is to use the constrained ℓ_1 minimisation approach to sparse precision matrix estimation (CLIME), implemented in the R package clime (Cai et al., 2011, 2012). The CLIME routine uses linear programming to solve the following (convex) optimisation problem,

$$\Theta^{\star} = \min |\Theta|_1$$
 subject to: $|\mathbf{S}\Theta - \mathbf{I}|_{\infty} \leq \lambda$,

where **S** is the sample covariance matrix and $|\mathbf{A}|_1 = \sum_{i,j} |a_{ij}|$ is the elementwise ℓ_1 norm of a matrix, **A**, and $|\mathbf{A}|_{\infty} = \max_{i,j} |a_{ij}|$ is the elementwise infinity norm. No symmetry requirements are placed on Θ^* so a symmetrising step is applied to obtain the final solution, $\hat{\mathbf{\Theta}}$,

$$\hat{\theta}_{ij} = \hat{\theta}_{ji} = \theta_{ij}^{\star} \mathbb{I}\{|\theta_{ij}^{\star}| \le |\theta_{ji}^{\star}|\} + \theta_{ji}^{\star} \mathbb{I}\{|\theta_{ij}^{\star}| > |\theta_{ji}^{\star}|\}.$$

Theorem 1 of Cai et al. (2011) shows that the resulting $\hat{\Theta}$ is PD with high probability.

Our simulations show that there is little difference between using CLIME and QUIC – the key point is that both appear to perform well in the presence of cellwise contamination when the input matrix is based on pairwise robust covariance estimates and it has been made PD using the NPD routine.

5. Simulation study for p < n

This section presents the results of an extensive simulation study to assess how well various robust covariance estimation techniques perform when used as an input to the regularisation routines outlined previously.

The proposed estimator begins by finding the covariances between all p(p-1)/2 pairs of variables. For the scale estimator underlying the robust covariance estimator, we consider Q_n , the τ -scale, the MAD and the IQR. We also consider the P_n estimator and two adaptively trimmed variants \tilde{P}_n , with trimming



Figure 2: Heat maps of the three kinds of precision matrices used to generate the data when p = 30.

parameters d = 3 and d = 5, see Tarr et al. (2012) for further details about these estimators. The pairwise covariances are arranged in a symmetric, though not necessarily PD, matrix. The symmetric matrix is transformed to a PD matrix using either the OGK method or the NPD method before being input into the GLASSO, QUIC or CLIME regularisation routines. For comparison purposes we also include the classical covariance estimator and the MCD as initial covariance matrix estimates.

5.1. Design

The simulated data follows a multivariate Gaussian distribution with n = 100 observations, $\mathcal{N}(\mathbf{0}, \mathbf{\Theta}^{-1})$. The precision matrices we select as the basis for the data generating process represent a broad range of scenarios that occur in practice and are similar to those used in Cai et al. (2011) and Hsieh et al. (2011). In particular, we consider three types of precision matrices, $\mathbf{\Theta}$, as shown in Figure 2 and outlined below.

- 1. Banded precision matrices, with elements $\theta_{ij} = 0.6^{|i-j|}$, such that the values of the entries decay the further they are from the main diagonal.
- 2. Sparse precision matrices, with randomly allocated non-zero entries, where $\Theta = \mathbf{B} + \delta \mathbf{I}$ with each off diagonal entry in **B** generated independently, where $P(b_{ij} = 0.5) = 0.1$ and $P(b_{ij} = 0) = 0.9$ and δ is chosen such that the condition number of the matrix equals *p*. The matrix is then standardised to have diagonal components equal to one. This scenario will be referred to as scattered sparsity.
- 3. Dense precision matrices, where Θ has all off diagonal elements equal to 0.5 and diagonal elements equal to 1.

The outliers were generated independently for each variable. In our simulations we allow the number of contaminated observations within each variable to increase up to a maximum of 25 observations (out of n = 100). In this way we have complete control over the total number of contaminated cells. The distribution of the outliers is a t_{10} distribution scaled by either a factor of 10 for extreme outliers or $\sqrt{10}$ for moderate outliers. The moderate outliers are perhaps closer to what one might expect in a real data set. However, the focus here is primarily on the extreme outliers where the overwhelming majority of the unusual observations lie well outside the cloud of standard observations. The extreme nature of the outliers serves to demark clearly estimators that have effectively broken down from those that are still capable of giving ballpark correct results. In both cases, the outliers are symmetrically distributed.

Each of the regularisation routines require a tuning parameter. At each replication, the tuning parameter was obtained by training on a separate (uncontaminated) randomly generated data set drawn from the true data generating process. For the training data, a sequence of precision matrices was obtained and the value

of the tuning parameter corresponding to the smallest entropy loss was then used for that replication. In practice, there was a small amount of variability in the choice of tuning parameter within each scenario and dimensionality. Furthermore, the QUIC and GLASSO routines almost always picked the same tuning parameter and the CLIME routine was free to choose slightly smaller tuning parameters. For example, in the scattered sparsity scenario with n = 100 and p = 60 the training set for CLIME resulted in a tuning parameter of around 0.09 whereas for QUIC and GLASSO it was closer to 0.12. In practice the tuning parameter may be chosen through cross validation, a BIC-type criterion such as in Yuan and Lin (2007), or it can be adjusted in an ad-hoc way until a desired level of sparsity is achieved.

We performed an extensive investigation into the most appropriate way to compare estimated precision matrices in the presence of cellwise contamination. We considered a number of matrix norms, the Frobenius norm, one norm, infinity norm and spectral norm, as well as the log determinant, the condition number and the entropy loss. The definition and behaviour of these performance indices under cellwise contamination are outlined in the supplementary material. We found that the most appropriate measure was the entropy loss defined as, $L_E(\Theta, \hat{\Theta}) = tr(\Theta^{-1}\hat{\Theta}) - \log det(\Theta^{-1}\hat{\Theta}) - p$.

As in Lin and Perlman (1985), we report the results in terms of the percentage relative improvement in average loss (PRIAL),

$$PRIAL(\hat{\boldsymbol{\Theta}}) = \frac{L_E(\boldsymbol{\Theta}, \hat{\boldsymbol{\Theta}}_0) - L_E(\boldsymbol{\Theta}, \hat{\boldsymbol{\Theta}})}{L_E(\boldsymbol{\Theta}, \hat{\boldsymbol{\Theta}}_0)} \times 100,$$

over N = 100 replication of each design, where $\hat{\Theta}_0$ is the estimated precision matrix after a regularisation technique has been applied to the classical sample covariance matrix for uncontaminated data. It is important to note that this is an extremely harsh benchmark to set.

5.2. Results

5.2.1. No contamination

Any good robust method should give comparable results to the classical non-robust method it is replacing when presented with a clean dataset. Table 1 presents the PRIAL results for the no contamination scenario. As the PRIAL results are relative to the base case for each routine, Table 1 cannot be used to compare the performance of the CLIME routine to the QUIC routine.

In the uncontaminated case, the OGK method substantially outperforms the NPD method. Overall, the methods appear to improve as the dimensionality increases, however, this is more a reflection of the deteriorating absolute performance of the baseline classical covariance matrix estimate.

For p = 30, the pairwise methods outperform the MCD, however the MCD method uses $\lfloor n + p + 1 \rfloor/2$ observations so when p = 90 the resulting estimator is the classical covariance estimate applied to 95 out of a total n = 100 observations. Hence, it is not surprising that the PRIAL for the MCD method is so close to zero. In fact, the MCD is not recommended for use when n < 2p (Rousseeuw et al., 2013).

The reweighted OGK (OGKw) methods essentially perform outlier detection and deletion before returning a classical covariance estimate of the cleaned data set. The performance of these methods is broadly similar over all the various initial scale estimates. Though not shown in Table 1, the MAD performs particularly poorly under both the OGK and the NPD corrections and would not be recommended for use.

As would be expected, given the solid Gaussian performance of P_n (see Tarr et al. (2012)), the methods based on P_n outperform those based on the τ -scale and Q_n . The relative deterioration in performance for the robust methods compared to the classical method is comparable to that in the simple univariate scale case. For example, the univariate scale estimator P_n has an asymptotic Gaussian relative efficiency of 86%.

5.2.2. Cellwise contamination

There are a number of ways to compare and contrast the various estimators. We consider data with n = 100 observations from three different data generating processes, across four dimensions, p = 15, 30, 60 or 90, contaminated with either moderate or extreme outliers, as explained in Section 5.1. We present figures for the extreme case only but comment on both based on the results of N = 100 replications of each design. We implement an array of initial covariance estimation techniques and process these through the GLASSO, QUIC and CLIME regularisation routines. Finally, as outlined in the supplementary material to this article, there are a number of performance indices that are considered. This section extracts and synthesises the key results.

We first consider the effect of dimensionality on the performance of the various estimators. A typical example is shown in Figure 3 where we plot the PRIAL results for the precision matrix resulting from the CLIME procedure for various input covariance matrices across different amounts of extreme contamination in each variable. The original data was generated assuming a banded precision matrix, however the trend holds true for scattered sparsity and dense precision matrices as well as for the QUIC and GLASSO procedures.

For relatively low dimensions, such as in the bottom panel of Figure 3 where p = 15, there is clearly an advantage to using the NPD method over the OGK method once there is more than a few percent of observations in each variable being contaminated. To avoid clutter, only the OGK method with P_n has been included in the plots, however, it is representative of the performance of the other scale estimators when used in conjunction with the OGK method.

As the dimensionality increases, the OGK and the MCD methods deteriorate faster. When p = 90, as outlined in the previous section, the MCD method behaves like the classical method. The OGK method performs similarly poorly as outlier propagation can lead to more than half of the elements in each principle component vector being contaminated. Hence, the eigenvalues in the spectral decomposition are replaced with robust estimates of scale that may no longer be valid.

Remarkably, the NPD methods perform consistently well. Their performance, relative to the classical method with no contamination improves as the number of variables increases. The P_n based method performs well for low levels of contamination, however once the proportion of contaminated cells is greater

		CLIME			QUIC			
		<i>p</i> = 30	p = 60	p = 90	<i>p</i> = 30	p = 60	p = 90	
au-scale	OGK	-17.1	-15.1	-12.7	-13.0	-9.5	-8.1	
	OGKw	-34.8	-29.3	-20.3	-26.9	-15.9	-15.6	
	NPD	-31.5	-31.3	-27.6	-25.3	-16.7	-15.4	
Qn	OGK	-16.5	-13.3	-11.7	-13.6	-10.5	-9.3	
	OGKw	-34.7	-28.5	-12.3	-27.0	-15.6	-14.6	
	NPD	-40.6	-36.9	-31.4	-32.5	-21.9	-20.3	
P _n	OGK	-13.6	-12.7	-10.9	-11.8	-9.6	-8.8	
	OGKw	-33.7	-27.1	-17.5	-26.2	-14.9	-14.5	
	NPD	-19.9	-18.4	-18.2	-15.7	-11.4	-11.2	
MCD		-53.1	-19.5	-3.7	-56.2	-19.5	-3.7	

Table 1: PRIAL results for the various estimators when there is no contamination present.



Figure 3: PRIAL results for a selection of estimators applied to data generated with a banded precision matrix with extreme outliers for p = 90 (top), p = 30 (middle) and p = 15 (bottom) using the CLIME regularisation procedure.

than 10% it does not perform as well as the other pairwise methods due to its lower breakdown value.

It is interesting to note that the adaptively trimmed P_n with adaptive trimming parameter d = 3, \tilde{P}_n , follows a somewhat different trajectory to the rest of the NPD type estimators. It maintains a relatively high level of performance even for quite high levels of contamination. This is due to the extreme nature of the contamination making the adaptive trimming extremely effective in identifying and excluding the errant observations. The advantage of \tilde{P}_n is lost when the contaminating distribution has only moderately sized outliers, in which case all the NPD pairwise methods perform comparably because \tilde{P}_n does almost no trimming.

To summarise, for p = 30, p = 60 and p = 90, using a pairwise method in conjunction with the NPD procedure as an input into the CLIME regularisation routine, the increase in entropy loss can be contained to less than double that of the classical method without contamination if the proportion of cellwise contamination is less than 10%.

The same pattern holds true when using the QUIC or the GLASSO regularisation routines. To demonstrate this consider Figure 4 where the PRIAL results are shown for CLIME, QUIC and the GLASSO under the banded precision matrix scenario with extreme outliers and p = 60. As would be expected the QUIC and GLASSO results are essentially identical, and largely consistent with the CLIME results in the top panel.

Consulting the raw entropy loss numbers reveals that the CLIME method gives slightly lower average entropy loss measurements, particularly for very high levels of contamination. In practice it does not matter what regularisation routine is used, the benefits of taking a pairwise approach to covariance estimation in the presence of cellwise contamination will still hold.

The NPD pairwise approach is a major improvement over standard robust estimators. An example of this is given in Figure 5 where we present the average PRIAL results for the QUIC estimator with p = 30 for the scenarios illustrated in Figure 2. Across all scenarios the same general pattern holds, the classical method and the OGK and MCD methods fail quite rapidly whereas the NPD approach offers much greater resilience to the cellwise contamination.

For the banded precision matrix scenario, top panel of Figure 5, the NPD based methods under the various robust scale estimators give similar results with P_n having a slight advantage over the others for low levels of contamination whereas Q_n has an advantage for higher contamination proportions.

For the scattered precision matrix and the dense precision matrix scenarios, \overline{P}_n gives the best results. The advantage of the adaptive trimming procedure is lost when the outliers are not so extreme, however, in such scenarios the adaptive trimming approach performs no worse than the other NPD methods.

We previously established that matrix norms are not a good performance measure for precision matrices. In terms of the other performance indicators, for all scenarios considered the log condition number remained bounded, suggesting that all three regularisation routines return well conditioned precision matrix estimates regardless of the level of contamination or the data generating process.

The NPD also performed well in terms of the log determinant performance index. As with the entropy loss, there appears to be an advantage to using \tilde{P}_n over the other scale estimators in each of the scenarios. Unlike with the entropy loss, the advantage of the adaptive trimming procedure is still evident even when the contamination is less extreme.

It is also constructive to see how $\hat{\Sigma} = \hat{\Theta}^{-1}$, the inverse of the estimated regularised precision matrix, compares with the true covariance matrix Σ . Figure 6 presents the average entropy loss and Frobenius norm results for the resulting estimated covariance matrices after regularisation using the CLIME procedure. We see similar trends to those outlined earlier. While using P_n alone does not perform well when the amount of contamination in each variable is large, the adaptive trimming procedure gives excellent results. The other pairwise methods also perform quite well. However, as we would expect, the classical method and



Figure 4: PRIAL results for a selection of estimators applied to data generated with a banded precision matrix with extreme outliers for p = 60 using CLIME (top), QUIC (middle) and GLASSO (bottom).



Figure 5: PRIAL results for a selection of estimators applied to data generated with a banded precision matrix (top), scattered precision matrix (middle) and dense precision matrix (bottom) with extreme outliers for p = 30 using the QUIC routine.

standard robust techniques, MCD and OGK fail quite rapidly. In general, the patterns for the matrix norms applied to $\hat{\Sigma}$ are very similar to those of the entropy loss for $\hat{\Theta}$. That is, the robust covariance matrices that are obtained at the end of the proposed procedure perform similarly well to the robust regularised precision matrices.

There are also important differences between $\hat{\Sigma}$ and the initial pairwise covariance matrix obtained after applying the NPD procedure. While the matrix norms for the the initial pairwise matrices were comparable to those for $\hat{\Sigma}$, the entropy loss and log determinant results for the initial pairwise covariance matrices were much worse due to small eigenvalues resulting from the NPD method. As such, we would not recommend simply applying the NPD procedure to a pairwise covariance matrix without further regularisation.

5.3. Gaussian graphical discovery rates

Another way to analyse the performance of a precision matrix estimator is through the lens of a Gaussian graphical model. When the data follow a multivariate Gaussian distribution, pairwise conditional independence between variables X_j and X_k holds if and only if $\theta_{jk} = 0$, therefore inferring linkages between variables corresponds to identifying the nonzero elements of $\Theta = (\theta_{jk})$, see Lauritzen (1996) for further details. Hence, rather than focussing on overall measures of similarity between the estimated precision and the true precision matrix, it can be informative to see how often the estimated precision matrix identifies the correct non-zero elements from the true precision matrix.

Figure 7 shows visually how well the QUIC estimator performs in the presence of cellwise contamination. When there is no contamination, all methods appear to perform similarly well in terms of their ability to correctly identify the true non-zero elements in the precision matrix. However, in the presence of 10% extreme contamination, the classical covariance and the MCD approach both fail to identify any structure as they tend to return overly dense precision matrices. On the other hand, the pairwise robust methods are, on average, still able to identify the underlying structure.

In the machine learning literature, the Matthews correlation coefficient (MCC), also known as the ϕ coefficient in the statistics literature, is often used to assess the ability of an estimator to identify the true non-zero elements in a precision matrix (Matthews, 1975). It takes into account the number of true positives (TP), false positives (FP), true negatives (TN) and false negatives (FN),

$$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}.$$

Typical MCC results are given in Figure 8 for the QUIC procedure. When cellwise contamination is introduced, the MCD, OGK and classical covariance approaches lose their ability to identify the true structure in the precision matrix quite quickly. The pairwise methods are much more resilient. As the number of contaminated observations in each variable increases, the ability of the pairwise methods to identify the true structure decreases gradually. When there is no contamination and p = 60, the classical method has an MCC of 0.39 compared with an MCC of 0.35 for the pairwise method based on P_n . When there is 5% extreme contamination in each variable, the MCC for the P_n based method is still at 0.29, while the classical approach is at 0.05. This pattern of results is virtually unchanged for the pairwise methods when the contamination is less extreme.

6. Simulation study for p > n

Of particular interest in a data mining context is the case when the number of variables p, is larger than the number of observations n. In order to perform simulations in a reasonable amount of time we reformulated the simulation settings in Section 5 such that samples of size n = 50 were drawn with dimension



Figure 6: Average entropy loss results for the precision matrices resulting from the CLIME procedure (top) and average entropy loss and Frobenius norm results for the resulting covariance matrix estimates (middle and bottom) for p = 60 with scattered sparsity and extreme outliers.



Figure 7: Heat maps showing how often each element in the precision matrix is identified as being non-zero using the QUIC routine over 100 replications. The top half have no contamination and the bottom half have 10% extreme contamination.



Figure 8: Matthews correlation coefficient results for the QUIC procedure with extreme outliers, p = 30 (top), p = 60 (middle) and p = 90 (bottom) in the scattered sparsity precision matrix scenario.

p = 15, 30, 60 and 90. The same types of precision matrices were considered as in the previous section, though we also looked at the case where the condition number of the sparse precision matrix was much larger than the dimension. Given the similar performance of the various regularisation routines found in the previous section, we restricted attention to the QUIC routine.

The results are similar to those found in the p < n setting. Figure 9 presents the MCC results for the case when p = 60 and n = 50. We can draw a direct comparison between Figure 9 and the middle panel of Figure 8 where p = 60 but n = 100. An important difference is that when p > n, the MCC values are lower across all levels of contamination, indicating that it is more difficult to recover the support of a Gaussian graphical model. For example, the classical approach had a MCC of 0.39 when n = 100, but only 0.29 when n = 50. Figure 10 allow us to compare the relative performance of the pairwise techniques as the condition number of the true precision matrix increases from 60 to 1000. The baseline value for the entropy loss is 8.8 for the classical approach when the condition number is 60, which increases slightly to 10.0 when the condition number is 1000. We observe that the performance of the various pairwise estimators decreases slightly as the condition number increases. For example when there is 8% contamination in each variable, the adaptively trimmed P_n estimator, has a PRIAL of -46% when the condition number is 60 (top panel), which decreases to -72% when the condition number is 1000 (bottom panel).



Figure 9: Matthews correlation coefficients results for the QUIC procedure for a scattered precision matrix with p = 60, n = 50, and extreme outliers.

Figure 11 demonstrates the impact of changing the extremity of the outliers when p = 90 and n = 50. As described in Section 5, the outliers are generated from a multivariate t_{10} distribution with scale matrix $k \mathbf{I}_p$, for k = 10, 50 and 100. With moderate outliers, as shown in the top panel of Figure 11, all the robust procedures perform comparably and the classical approach is not affected too badly. As the extremity of the outliers increases, shown in the bottom two panels of Figure 11 the performance of the classical approach deteriorates quickly. Between the middle and bottom panels, there is little difference in the performance of the high-breakdown value robust estimators, indicating that they have stabilised and are likely to continue giving the same result even if the existing outliers were moved further away. It is clear that the method based on P_n , with its lower breakdown value, continues to be affected by the size of contamination when there is a large proportion of contamination in each variable and we would expect its performance to continue to deteriorate if the outlier generating distribution was even more extreme.

Note that P_n , the adaptively trimmed P_n , remains the most stable as the extremity of the outliers in-



Figure 10: Entropy loss PRIAL results for the QUIC procedure with a scattered precision matrix, p = 60, n = 50 and extreme outliers. The condition number of the underlying precision matrix is 60 in the top panel and 1000 in the bottom panel.

creases. In fact, when there is 10% cellwise contamination the PRIAL remains at -59% whether k = 50 or 100, which is down from -47% when k = 10. This can be attributed to the adaptive trimming correctly identifying the vast majority of the contaminated bivariate observations when the outliers are so extreme.

7. Conclusion

A pairwise approach to covariance estimation has a natural resilience to the type of cellwise contamination seen in high dimensional scenarios where classical robust procedures, such as the MCD, *M*-estimators, Quadrant Covariance and OGK, tend to fail.

We considered a broad range of scenarios: from dense precision matrices, as typically found in standard analyses with $n \gg p$; to banded precision matrices that often occur in time series settings and may also be representative of scenarios with block diagonal precision matrices; as well as scattered sparsity, where the linkages between variables are not known beforehand and can show up anywhere within the precision matrix, as is often found in settings where p > n.

After careful consideration of the various performance indices available in the multivariate setting, outlined in the supplementary material, our primary choice of performance measure was the entropy loss. When appropriate, we showed that the entropy loss returned similar conclusions to other performance indices, such as the Frobenius norm and log determinant.

We have shown that combining robust pairwise covariance estimation with the NPD method and regularisation techniques such as the CLIME, QUIC or GLASSO yield precision matrices that are robust to cellwise contamination. The additional advantages of the regularisation techniques, such as the promotion of sparsity also carry through. While it was expected, given that they are solving the same minimisation problem, it is reassuring to find that the QUIC estimates are virtually indistinguishable from the standard GLASSO approach in all scenarios considered here. Furthermore, it did not appear to matter which of the three considered regularisation routines was applied, as all gave broadly similar results in the various scenarios considered. This is comforting given the current pace of research in this area, with new procedures being suggested frequently.

We demonstrated that the proposed approach maintains its ability to identify the true precision matrix structure, as measured by the Matthews correlation coefficient, under moderate levels of contamination.

We also investigated what happens when the resulting precision matrix is inverted to find the corresponding covariance matrix estimate. Applying the same performance indices to the resulting covariance matrices, we found they perform similarly well to the underlying precision matrix.

The simulation study allowed for quite high levels of arbitrary contamination in multivariate data sets. As such, the pairwise techniques based on the standard P_n estimator unsurprisingly did not perform as well as Q_n and τ -scale estimators, however, the adaptively trimmed P_n , \tilde{P}_n with trimming parameter d = 3 typically performed extremely well, due to its ability to detect and trim extreme outliers in bivariate space. Finally, we showed that the performance of the proposed technique continues to perform well even when p is moderately larger than n.

Acknowledgements

The authors would like to thank two anonymous referees for their helpful comments on an earlier draft.

References

Agostinelli, C., Leung, A., Yohai, V., Zamar, R., 2014. Robust estimation of multivariate location and scatter in the presence of cellwise and casewise contamination. Technical Report. arXiv:1406.6031.



Figure 11: Results for the QUIC procedure for a banded precision matrix with p = 90 and n = 50. The outlier generating distribution is a multivariate t_{10} distribution with scale matrix 10 I_{90} (top), 50 I_{90} (middle) and 100 I_{90} (bottom).

- Alqallaf, F., Van Aelst, S., Yohai, V.J., Zamar, R.H., 2009. Propagation of outliers in multivariate data. The Annals of Statistics 37, 311–331. doi:10.1214/07-A0S588.
- Alqallaf, F.A., Konis, K.P., Martin, R.D., Zamar, R.H., 2002. Scalable robust covariance and correlation estimates for data mining, in: Proceedings of the Eighth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, ACM, New York, NY, USA. pp. 14–23. doi:10.1145/775047.775050.
- Barnett, V., Lewis, T., 1994. Outliers in statistical data. Wiley Series in Probability and Mathematical Statistics: Applied Probability and Statistics. 3rd ed., Wiley, New York.
- Bieroza, M., Baker, A., Bridgeman, J., 2011. Classification and calibration of organic matter fluorescence data with multiway analysis methods and artificial neural networks: an operational tool for improved drinking water treatment. Environmetrics 22, 256–270. doi:10.1002/env.1045.
- Cai, T., Liu, W., Luo, X., 2011. A constrained ℓ_1 minimization approach to sparse precision matrix estimation. Journal of the American Statistical Association 106, 594–607. doi:10.1198/jasa.2011.tm10155.
- Cai, T., Liu, W., Luo, X., 2012. clime: Constrained ℓ_1 -minimization for inverse (covariance) matrix estimation. URL: http://CRAN.R-project.org/package=clime. R package version 0.4.1.
- Cator, E.A., Lopuhaä, H.P., 2010. Asymptotic expansion of the minimum covariance determinant estimators. Journal of Multivariate Analysis 101, 2372–2388. doi:10.1016/j.jmva.2010.06.009.
- Cator, E.A., Lopuhaä, H.P., 2012. Central limit theorem and influence function for the MCD estimators at general multivariate distributions. Bernoulli 18, 520–551. doi:10.3150/11-BEJ353.
- Croux, C., Filzmoser, P., Pison, G., Rousseeuw, P.J., 2003. Fitting multiplicative models by robust alternating regressions. Statistics and Computing 13, 23–36. doi:10.1023/A:1021979409012.
- Dempster, A.P., 1972. Covariance selection. Biometrics 28, 157–175. doi:10.2307/2528966.

Dey, D.K., Srinivasan, C., 1985. Estimation of a covariance matrix under Stein's loss. The Annals of Statistics 13, 1581–1591. doi:10.1214/aos/1176349756.

- Farcomeni, A., 2014. Robust constrained clustering in presence of entry-wise outliers. Technometrics 56, 102–111. doi:10.1080/00401706.2013.826148.
- Filzmoser, P., Ruiz-Gazen, A., Thomas-Agnan, C., 2014. Identification of local multivariate outliers. Statistical Papers 55, 29–47. doi:10.1007/s00362-013-0524-z.
- Friedman, J.H., Hastie, T., Tibshirani, R., 2008. Sparse inverse covariance estimation with the graphical lasso. Biostatistics 9, 432–441. doi:10.1093/biostatistics/kxm045.
- Gales, M.J.F., van Dalen, R.C., 2007. Predictive linear transforms for noise robust speech recognition, in: IEEE Workshop on Automatic Speech Recognition & Understanding, IEEE. pp. 59–64.
- Gentle, J., 2007. Matrix Algebra Theory, Computations and Applications in Statistics. Springer, New York.
- Gnanadesikan, R., Kettenring, J.R., 1972. Robust estimates, residuals, and outlier detection with multiresponse data. Biometrics 28, 81–124. doi:10.2307/2528963.
- Gupta, M., Srivastava, S., 2010. Parametric Bayesian estimation of differential entropy and relative entropy. Entropy 12, 818–843. doi:10.3390/e12040818.
- Higham, N.J., 2002. Computing the nearest correlation matrix a problem from finance. IMA Journal of Numerical Analysis 22, 329–343. doi:10.1093/imanum/22.3.329.
- Hsieh, C.J., Dhillon, I.S., Ravikumar, P.K., Sustik, M.A., 2011. Sparse inverse covariance matrix estimation using quadratic approximation, in: Shawe-Taylor, J., Zemel, R., Bartlett, P., Pereira, F., Weinberger, K. (Eds.), Advances in Neural Information Processing Systems, pp. 2330–2338.
- Hsieh, C.J., Sustik, M.A., Dhillon, I.S., Ravikumar, P.K., Poldrack, R., 2013. BIG & QUIC: Sparse inverse covariance estimation for a million variables, in: Burges, C., Bottou, L., Welling, M., Ghahramani, Z., Weinberger, K. (Eds.), Advances in Neural Information Processing Systems, pp. 3165–3173.
- Huber, P.J., 1964. Robust estimation of a location parameter. Annals of Mathematical Statistics 35, 73–101. doi:10.1214/aoms/1177703732.
- Hubert, M., Rousseeuw, P.J., Vakili, K., 2014. Shape bias of robust covariance estimators: an empirical study. Statistical Papers 55, 15–28. doi:10.1007/s00362-013-0544-8.
- Hubert, M., Rousseeuw, P.J., Van Aelst, S., 2008. High-breakdown robust multivariate methods. Statistical Science 23, 92–119. doi:10.1214/08834230700000087.
- James, W., Stein, C., 1961. Estimation with quadratic loss, in: Neyman, J. (Ed.), Proceedings of the Fourth Berkeley Symposium on Mathematical Statististics and Probability, University of California Press, Berkeley. pp. 361–379.
- Johnstone, I.M., 2001. On the distribution of the largest eigenvalue in principal components analysis. The Annals of Statistics 29, 295–327. doi:10.1214/aos/1009210544.
- Lauritzen, S., 1996. Graphical Models. Oxford University Press, New York.

Ledoit, O., Wolf, M., 2004. A well-conditioned estimator for large-dimensional covariance matrices. Journal of Multivariate

Analysis 88, 365-411. doi:10.1016/S0047-259X(03)00096-4.

- Li, Y., 2004. On incremental and robust subspace learning. Pattern Recognition 37, 1509 1518. doi:10.1016/j.patcog.2003. 11.010.
- Lin, S.P., Perlman, M.D., 1985. A Monte Carlo comparison of four estimators of a covariance matrix, in: Krishnaiah, P.R. (Ed.), Proceedings of the Sixth International Symposium on Multivariate Analysis, North-Holland, Amsterdam. pp. 411–429.
- Little, R., Rubin, D., 2002. Statistical Analysis with Missing Data. 2nd ed., Wiley, Hoboken.
- Liu, L., Hawkins, D.M., Ghosh, S., Young, S.S., 2003. Robust singular value decomposition analysis of microarray data. Proceedings of the National Academy of Sciences of the United States of America 100, 13167–13172. doi:10.1073/pnas. 1733249100.
- Løland, A., Huseby, R.B., Hjort, N.L., Frigessi, A., 2013. Statistical corrections of invalid correlation matrices. Scandinavian Journal of Statistics 40, 807–824. doi:10.1111/sjos.12035.
- Ma, Y., Genton, M.G., 2001. Highly robust estimation of dispersion matrices. Journal of Multivariate Analysis 78, 11–36. doi:10.1006/jmva.2000.1942.
- Mardia, K.V., Kent, J.T., Bibby, J.M., 1979. Multivariate Analysis. Probability and Mathematical Statistics, Academic Press, London.
- Maronna, R.A., Martin, R.D., Yohai, V.J., 2006. Robust Statistics. Wiley, London.
- Maronna, R.A., Yohai, V.J., 2008. Robust low-rank approximation of data matrices with elementwise contamination. Technometrics 50, 295–304. doi:10.1198/004017008000000190.
- Maronna, R.A., Zamar, R.H., 2002. Robust estimates of location and dispersion for high-dimensional datasets. Technometrics 44, 307–317. doi:10.1198/004017002188618509.
- Matthews, B.W., 1975. Comparison of the predicted and observed secondary structure of T4 phage lysozyme. Biochimica et Biophysica Acta 405, 442–451. doi:10.1016/0005-2795(75)90109-9.
- Mavroeidis, D., Marchiori, E., 2014. Feature selection for k-means clustering stability: theoretical analysis and an algorithm. Data Mining and Knowledge Discovery 28, 918–960. doi:10.1007/s10618-013-0320-3.
- Rousseeuw, P.J., Croux, C., 1993. Alternatives to the median absolute deviation. Journal of the American Statistical Association 88, 1273–1283. doi:10.1080/01621459.1993.10476408.
- Rousseeuw, P.J., Croux, C., Todorov, V., Ruckstuhl, A., Salibian-Barrera, M., Verbeke, T., Koller, M., Maechler, M., 2013. robustbase: Basic Robust Statistics. URL: http://CRAN.R-project.org/package=robustbase. R package version 0.9-10.
- Rousseeuw, P.J., Molenberghs, G., 1993. Transformation of non positive semidefinite correlation matrices. Communications in Statistics Theory and Methods 22, 965–984. doi:10.1080/03610928308831068.

Stein, C., 1956. Some Problems in Multivariate Analysis, Part I. Technical Report 6. Stanford University. Stanford.

- Tarr, G., Müller, S., Weber, N.C., 2012. A robust scale estimator based on pairwise means. Journal of Nonparametric Statistics 24, 187–199. doi:10.1080/10485252.2011.621424.
- Tarr, G., Müller, S., Weber, N.C., 2014. Robust estimation of precision matrices under cellwise contamination. Under review .
- De la Torre, F., Black, M.J., 2001. Robust principal component analysis for computer vision, in: International Conference on Computer Vision, IEEE, Vancouver. pp. 362–369. doi:10.1109/ICCV.2001.937541.

Tukey, J.W., 1962. The future of data analysis. Annals of Mathematical Statistics 33, 1–67. doi:doi:10.1214/aoms/1177704711.

- Van Aelst, S., Vandervieren, E., Willems, G., 2012. A Stahel-Donoho estimator based on Huberized outlyingness. Computational Statistics & Data Analysis 56, 531 542. doi:10.1016/j.csda.2011.08.014.
- Wilks, S.S., 1932. Certain generalizations in the analysis of variance. Biometrika 24, 471–494. doi:10.2307/2331979.
- Won, J.H., Lim, J., Kim, S.J., Rajaratnam, B., 2013. Condition-number-regularized covariance estimation. Journal of the Royal Statistical Society: Series B (Statistical Methodology) 75, 427–450. doi:10.1111/j.1467-9868.2012.01049.x.
- Yang, R., Berger, J.O., 1994. Estimation of a covariance matrix using the reference prior. The Annals of Statistics 22, 1195–1211. doi:10.1214/aos/1176325625.
- Yuan, M., Lin, Y., 2007. Model selection and estimation in the Gaussian graphical model. Biometrika 94, 19–35. doi:10.1093/ biomet/asm018.

Supplementary material for: Robust estimation of precision matrices under cellwise contamination

G. Tarr*, S. Müller, N. C. Weber

School of Mathematics and Statistics, University of Sydney, NSW 2006, Australia

Abstract

This supplementary material gives a summary of the performance indices used in the article "Robust estimation of covariance and precision matrices under cellwise contamination" by Tarr et al. (2014).

1. Performance indices

We require a way to assess the performance of various covariance and precision matrix estimators under cellwise contamination. There are a range of possible ways to measure how close an estimated matrix is to the true value. In order to assess the performance of our proposed estimators, we first need to identify which performance indices are appropriate. One class of performance indices considered are matrix norms which measure the size of a matrix. The second class looks at how closely the estimated precision (or covariance) matrix reflects the nature of the theoretical precision (or covariance) matrix, through either the determinant, condition number or an overall entropy loss index. Let Σ denote the true covariance matrix and $\Theta = \Sigma^{-1}$ denote the true precision matrix. In this section we define the performance measures and consider their appropriateness in the context of cellwise contamination.

1.1. Measures of performance

1.1.1. Matrix norms

The Frobenius norm is perhaps the most common matrix norm, it is an element-wise norm, the Euclidean norm of **A** treated as if it were a vector of length p^2 , $\|\mathbf{A}\|_F = \sqrt{\operatorname{tr}(\mathbf{A'A})}$. An alternative way of constructing a matrix norm is to take a vector norm and use it to generate a matrix norm of the form $\|\mathbf{A}\| = \sup_{\|\mathbf{x}\|=1} \|\mathbf{Ax}\|$, where $\|\cdot\|$ on the left is the induced (or operator) norm and $\|\cdot\|$ on the right is a vector norm. Examples of induced norms are the L_p norms. For example, the one norm, $\|\mathbf{A}\|_1 = \max_j \sum_{i=1}^n |a_{ij}|$; the infinity norm, $\|\mathbf{A}\|_{\infty} = \max_i \sum_{j=1}^n |a_{ij}|$; and the spectral norm, $\|\mathbf{A}\|_2 = \sigma_{\max}(\mathbf{A})$, where $\sigma_{\max}(\mathbf{A})$ is the largest singular value of \mathbf{A} . When \mathbf{A} is nonsingular $\|\mathbf{A}^{-1}\|_2 = 1/\sigma_{\min}(\mathbf{A})$ where $\sigma_{\min}(\mathbf{A})$ is the smallest singular value of \mathbf{A} .

In our experiments, we apply the matrix norms to $\mathbf{A} = \mathbf{\Theta}_0 - \mathbf{I}$ where $\mathbf{\Theta}_0 = \mathbf{\Theta}^{-1} \hat{\mathbf{\Theta}}$. While $\mathbf{\Theta}$ and $\hat{\mathbf{\Theta}}$ are symmetric, it is not the case that the product of two symmetric matrices yields a symmetric matrix, hence in general $\mathbf{\Theta}_0 \neq \mathbf{\Theta}'_0$, so in practice the one norm and the infinity norm may yield different results.

Preprint submitted to Elsevier

^{*}Corresponding author

Email address: gtar4178@uni.sydney.edu.au (G. Tarr)

For each norm, one may naïvely assume that the closer to zero the better, however, in Section 1.2 we demonstrate that this is not always the case, particularly when estimating precision matrices in the presence of outliers.

Aside from matrix norms, there are a few other commonly employed performance indices.

1.1.2. Entropy loss

The entropy loss, as suggested by Stein (1956) and featured in James and Stein (1961) and also Dey and Srinivasan (1985), when applied to precision matrices is defined as,

$$L_E(\mathbf{\Theta}, \hat{\mathbf{\Theta}}) = \operatorname{tr}(\mathbf{\Theta}^{-1} \hat{\mathbf{\Theta}}) - \log \operatorname{det}(\mathbf{\Theta}^{-1} \hat{\mathbf{\Theta}}) - p$$
$$= \sum_{i=1}^{p} (\lambda_i - \log \lambda_i) - p,$$

where λ_i , i = 1, ..., p, are the eigenvalues of Θ_0 . Stein (1956) notes that this function is "somewhat arbitrary" but it is convex in $\hat{\Theta}$ and assuming that Θ and $\hat{\Theta}$ are positive semidefinite, $L_E(\Theta, \hat{\Theta}) \ge 0$ with equality if and only if $\Theta = \hat{\Theta}$.

However, the entropy loss is not as "arbitrary" as it may seem at first. Note that the Kullback-Leibler divergence from $N_1(\mu, \Sigma_1)$ to $N_2(\mu, \Sigma_2)$ is, $D_{\text{KL}}(N_1, N_2) = \frac{1}{2}L_E(\Sigma_1, \Sigma_2)$. The close link between the entropy loss and the Kullback-Leibler or Bregman divergence loss is shown in a Bayesian context by Gupta and Srivastava (2010).

There is also a clear link between the entropy loss and the likelihood ratio test for H_0 : $\Sigma = \Sigma^*$ with unknown mean assuming the data come from a Gaussian distribution, $-2(l_1 - l_0) = nL_E(\Sigma^*, S)$, where l_0 and l_1 are the log-likelihoods under the null and alternative hypotheses, see, for example Mardia et al. (1979, p. 126).

The entropy loss is used extensively as a basis for developing and assessing improved precision and covariance matrix estimators, for example in Lin and Perlman (1985), Yang and Berger (1994) and more recently, Won et al. (2013).

1.1.3. Log determinant

In the multivariate Gaussian setting, Wilks (1932) names the determinant of the covariance matrix, det(Σ), the generalised variance. The generalised precision is similarly defined as the determinant of the precision matrix, det(Θ). This idea can be used as the basis for a performance index. Consider the log of the determinant of the standardised covariance or precision matrix, $L_D(\Theta_0) = \log \det(\Theta_0) = \sum_{i=1}^p \log \lambda_i$. The determinant of an identity matrix is 1, so the optimal value of $L_D(\Theta_0)$ is 0. Positive (negative) log determinant results indicate that the generalised variance or precision is being over (under) estimated. Note that, $L_D(\Theta_0) = -L_D(\Sigma_0)$. Thus, methods that underestimate the generalised variance will overestimate the generalised precision.

The log determinant is a very crude performance index which can be dominated by one eigenvalue that is very close to zero. Furthermore, it is incorporated as part of the entropy loss so there is little need to focus on it in the results of the simulation studies.

1.1.4. Log condition number

Formally, the condition number of a square matrix is the product of the norm of the matrix and the norm of its inverse, $\kappa(\Theta_0) = ||\Theta_0^{-1}|| \cdot ||\Theta_0||$, and hence depends on the kind of matrix-norm. It is common to use the spectral norm, in which case the condition number is the ratio of the largest to the smallest non-zero singular value of the matrix. The condition number associated with the systems of equations, Ax = b, gives



Figure 1: A series of boxplots showing the distribution of the largest through to the smallest eigenvalues of an estimated covariance matrix, **S**, over N = 100 samples from $\mathcal{N}(\mathbf{0}, \mathbf{I})$ with n = 100 and p = 30, 60 and 90.

a bound on how inaccurate the solution may be. A system is said to be ill-conditioned if small changes in the inputs, A and b, result in large changes in the solution, x.

Consider the performance index defined as the log of the condition number of Θ_0 ,

$$L_{\kappa}(\boldsymbol{\Theta}_{0}) = \log \kappa(\boldsymbol{\Theta}_{0}) = \log(\sigma_{\max}(\boldsymbol{\Theta}_{0})) - \log(\sigma_{\min}(\boldsymbol{\Theta}_{0})).$$

The condition number for an identity matrix is 1 and the condition number for a singular matrix is infinity, so the $0 \le L_{\kappa}(\Theta_0) \le \infty$.

Gentle (2007) notes that while the condition number of a matrix provides a useful indication of its ability to solve linear equations accurately, it can be misleading at times when the rows (or columns) of the matrix have very different scales. That is, the condition number can be changed by simply scaling the rows or columns which does not actually make a linear system of equations any better or worse conditioned. This is known as artificial ill-conditioning.

In the context of the sample covariance matrix, **S**, Ledoit and Wolf (2004) note that "when the ratio p/n is less than one but not negligible, the sample covariance matrix is invertible but numerically ill-conditioned, which means that inverting it amplifies estimation error dramatically." Won et al. (2013) go further stating that "the eigenstructure [of **S**] tends to be systematically distorted unless p/n is extremely small, resulting in numerically ill-conditioned estimators for Σ ." Figure 1 demonstrates the systematic deterioration in the eigenstructure as $p/n \rightarrow 1$. The eigenvalues of the true covariance matrix are all identically 1, however this is not reflected in the eigenvalues of the estimated sample covariance matrices.

As with the log determinant, the log condition number is not a particularly discerning performance index. To assess whether a robust estimator provides reasonable estimates, the most that it can contribute is whether or not the log condition number remains bounded.

1.1.5. Quadratic loss

Another index that is frequently used in the literature to assess the performance of covariance matrix estimators is the quadratic loss. The exact specification varies from paper to paper, for example Ledoit and Wolf (2004) define it as, $\|\hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}\|_F^2$. An alternative specification of the quadratic loss, more in line with the entropy loss, is used in Won et al. (2013), $\|\hat{\boldsymbol{\Sigma}}\boldsymbol{\Sigma}^{-1} - \mathbf{I}\|_F^2$. It is obvious that the quadratic loss is intrinsically linked to the Frobenius norm.

1.2. Behaviour of performance indices

The performance indices outlined in this section are typically used to compare competing estimators in uncontaminated data sets. Contaminated data will have potentially severe implications for structure and size of the estimated precision and covariance matrices, and it is not clear how these indices will behave in such settings. As such, we begin our investigation by exploring how these indices react to the presence of gross outliers in a data set.

The model used to assess the behaviour of the various performance indices is typical of that used in the simulation study, n = 100 observations drawn from the standard multivariate Gaussian distribution, $\mathcal{N}(\mathbf{0}, \mathbf{I})$, with p = 30.

1.2.1. Inflated variances

This section explores how the various performance indices react if we artificially inflate the variance of the first variable, i.e. increase the value of s_{11} in the sample covariance matrix, $\mathbf{S} = (s_{ij})$, based on a single sample of uncontaminated data. In this simple case, where $\Sigma = \Theta = \mathbf{I}$, the matrix norms are applied to $\mathbf{S} - \mathbf{I}$ or $\mathbf{S}^{-1} - \mathbf{I}$ and the entropy loss, log determinant and log condition number are simply applied to \mathbf{S} or \mathbf{S}^{-1} . Note that in this setting the one norm and the infinity norm will give identical results and so only the results for the one norm are shown.

Figure 2 shows the behaviour of the various indices when applied to these adjusted covariance matrices and Figure 3 presents the same for the resulting precision matrices, $\hat{\Theta} = S^{-1}$. The horizontal axis shows the size of s_{11} , the artificially inflated variance of the first variable in the sample covariance matrix.

In Figure 2, the majority of the performance indices behave similarly in the covariance case – there is an overall positive trend as the variance of the first variable increases. The spectral norm and Frobenius norm both increase uniformly with s_{11} . The one norm, and correspondingly the infinity norm (not shown), remains flat as the sum of the absolute value of the elements in another column (or row) remains larger than the first column (row) up until the point where $s_{11} \approx 2$, at which point the one norm (and infinity norm) increase linearly with s_{11} . For large s_{11} , it is clear that all considered performance indices register that the adjusted **S** matrix is no longer close to the true value, **I**.

In stark contrast, Figure 3 shows the indices when applied to the resulting precision matrix (after the first entry in the covariance matrix has been artificially inflated). As expected, the condition number of the resulting inverse is identical to that of the original covariance matrix and $L_D(\mathbf{S}^{-1}) = -L_D(\mathbf{S})$, however the other indices exhibit somewhat different behaviour. In particular the matrix norms tend to decrease, rather than increase. This is explained by noting that as the first element in the covariance matrix is artificially inflated, the first row and column of the precision matrix decay towards zero while the other elements are more or less constant. This is demonstrated in Figure 4 where $\hat{\theta}_{11}$ and $\hat{\theta}_{21}$ both tend towards zero whereas the other main diagonal and off diagonal elements remain quite stable as the level of contamination increases.

In Figure 3, the Frobenius norm, an elementwise norm, exhibits a minimum turning point before levelling off. This is due to the first row and column of the precision matrix converging rapidly to zero, often from relatively large starting points, whereas the convergence of the other elements is not as drastic and not necessarily shrinking towards zero, hence the upward trend.

The entropy loss broadly exhibits similar behaviour in both Figures 2 and 3. The minimum turning point in Figure 3 is somewhat similar to that of the Frobenius norm and is explained by noting that this is due to the sum of the eigenvalues decaying quite quickly before levelling off as s_{11} increases, whereas the sum of the logs of the eigenvalues decays much more slowly. Regardless, it is clear that the entropy loss tends to reflect the impact of the inflated variance in both the covariance matrix and the resulting precision matrix.



Figure 2: The impact of artificially inflating the size of the top left element of the sample covariance matrix, s_{11} , on each of the performance indices when applied to the resulting covariance matrix.



Figure 3: The impact of artificially inflating the size of the top left element of the sample covariance matrix, s_{11} , on each of the performance indices when applied to the resulting precision matrix.



Figure 4: The change in few elements of $\hat{\Theta} = \mathbf{S}^{-1}$ as s_{11} is artificially inflated. Main diagonal elements are in the top row and the off diagonal elements are in the bottom row. Note that $\mathbf{S} = (s_{ij})$ and $\hat{\Theta} = (\hat{\theta}_{ij})$.

1.2.2. Contamination in the data

Instead of directly manipulating the estimated covariance matrix, consider introducing contamination into the original data set and observing what effect that has on the performance metrics applied to the covariance and resulting precision matrix. For each level of contamination we take N = 1000 samples from $\mathcal{N}(0, \mathbf{I})$. For each sample we estimate the classical sample covariance matrix, \mathbf{S} , and take the inverse to obtain $\hat{\mathbf{\Theta}} = \mathbf{S}^{-1}$.

Figures 5 and 6 show the behaviour of the various loss indices over N = 1000 replications. The horizontal axis represents the number of contaminated observations within each of the p = 30 variables. Starting with an uncontaminated multivariate Gaussian distribution, we then progressively add one contaminated observation to each variable until there are 24 contaminated observations within each variable. The contamination is performed by assigning to each randomly selected cell the value of 10.

In general, cellwise outlying contamination will destroy any existing dependence structure and inflate the main diagonal of the covariance matrix, resulting in increases for the entropy loss and matrix norms applied to the covariance as seen in Figure 5. The log determinant of the estimated covariance matrix trends upwards, demonstrating the over estimated generalised variance with increasing levels of contamination. Apart from a relatively minor spike when there is only one contaminated observation in each variable, the condition number is not adversely affected by increasing levels of contamination, reflecting the stabilised eigenvalues of the resulting covariance matrix. This demonstrates that in this setting, the condition number is not an appropriate index against which to compare the performance of competing robust estimators.

The interpretation of the performance indices when they are applied to the resulting precision matrix is more complicated. We see in Figure 6 that there is still structure present in the precision matrix, in the sense that there is a main diagonal behaving distinctly from the off diagonal elements. However, all the elements



Figure 5: The impact of randomly contaminating a certain number of cells in each variable of a 100×30 data matrix on the various performance indices when applied to the resulting covariance matrix.



Figure 6: The impact of randomly contaminating a certain number of cells in each variable of a 100×30 data matrix on the various performance indices when applied to the resulting precision matrix.

tend to shrink towards zero. Hence, for large amounts of contamination, $\hat{\Theta} - I \approx -I$ and so the matrix norms tend to converge to $\| - I \|$.

As in the previous scenario, the Frobenius norm exhibits a minimal turning point before plateauing. This is explained by noting that while the introduction of contamination has an immediate shrinkage effect on the main diagonal of the precision matrix, including one or two influential observations in each variable induces an artificially high level of correlation between some variables. Hence, it can take some time for the off diagonal elements to stabilise. When more than a few outlying cells are present in each variable, the artificial correlation structure wanes and hence the Frobenius norm applied to the resulting precision matrix trends towards $||\mathbf{I}||_F = \sqrt{p}$. Hence, in this contamination setting the matrix norms appear to be useful only when applied to the covariance matrix, not the precision matrix.

As in the previous scenario, the entropy loss behaves consistently for both the covariance and precision matrix. Similarly to Figure 3, it exhibits a slight drop when only one cell in each variable is contaminated after which it increases as the proportion of contaminated cells grows. As such, the entropy loss is the preferred performance index when looking across both covariance and precision matrix estimators.