



**HAL**  
open science

## Optimal shrinkage for robust covariance matrix estimators in a small sample size setting

Karina Ashurbekova, Antoine Usseglio-Carleve, Florence Forbes, Sophie Achard

### ► To cite this version:

Karina Ashurbekova, Antoine Usseglio-Carleve, Florence Forbes, Sophie Achard. Optimal shrinkage for robust covariance matrix estimators in a small sample size setting. 2019. hal-02378034v1

**HAL Id: hal-02378034**

**<https://hal.science/hal-02378034v1>**

Preprint submitted on 24 Nov 2019 (v1), last revised 27 Mar 2021 (v3)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Optimal shrinkage for robust covariance matrix estimators in a small sample size setting

Karina Ashurbekova · Antoine  
Usseglio-Carleve · Florence Forbes ·  
Sophie Achard

Received: date / Accepted: date

**Abstract** When estimating covariance matrices, traditional sample covariance-based estimators are straightforward but suffer from two main issues: 1) a lack of robustness, which occurs as soon as the samples do not come from a Gaussian distribution or are contaminated with outliers and 2) a lack of data when the number of parameters to estimate is too large compared to the number of available observations, which occurs as soon as the covariance matrix dimension is greater than the sample size. The first issue can be handled by assuming samples are drawn from a heavy-tailed distribution, at the cost of more complex derivations, while the second issue can be addressed by shrinkage with the difficulty of choosing the appropriate level of regularization. In this work we offer both a tractable and optimal framework based on shrunk likelihood-based M-estimators. First, a closed-form expression is provided for a regularized covariance matrix estimator with an optimal shrinkage coefficient for any sample distribution in the elliptical family. Then, a complete inference procedure is proposed which can also handle both unknown mean and tail parameter, in contrast to most existing methods that focus on the covariance matrix parameter requiring pre-set values for the others. An illustration on synthetic and real data is provided in the case of the t-distribution with unknown mean and degrees-of-freedom parameters.

**Keywords** Covariance estimation · Small sample size · Shrinkage methods · Robust estimation · Elliptical distributions.

## 1 Introduction

Accurate estimation of covariance matrices (or more generally scale matrices) is fundamental in many areas of statistics. Examples include applications in finance [21], bioinformatics and classification [12]. Practitioners usually have to deal with two main difficulties. First, observations may deviate from the Gaussian distribution due to a particular data generating process or the presence of

---

K. Ashurbekova, A. Usseglio-Carleve, F. Forbes, S. Achard  
Univ. Grenoble Alpes, Inria, CNRS, Grenoble INP, LJK, 38000 Grenoble, France  
E-mail: firstname.lastname@inria.fr

outlying data. Ignoring this deviation may conduct to inadequate predictions and conclusions [9]. A widespread solution to design so-called robust estimators, is to consider heavy-tailed distributions which can better accommodate outliers. Among those, elliptical distributions have been studied as good candidates as they include tractable heavy-tailed distributions such as the  $t$ -distribution, whose tail is controlled by a single degrees-of-freedom (dof) parameter (see [33] or [4]). In addition, for elliptical distributions, robust estimators of the scale matrix  $\Sigma$  are provided by Maronna's M-estimators, defined as the solution  $\tilde{\Sigma}$  of a fixed-point relationship  $\tilde{\Sigma} = \mathbb{E} \left[ u \left( \mathbf{x}^\top \tilde{\Sigma}^{-1} \mathbf{x} \right) \mathbf{x} \mathbf{x}^\top \right]$  where  $u$  is function satisfying a set of general assumptions (see [26]). A second difficulty is then that the problem dimension may be too large compared to the number of available observations, which prevents accurate estimation when this feature is not explicitly taken into account. For example, if the dimension  $p$  of  $\Sigma$  is greater than the sample size  $n$ , previous Maronna's estimators are known to perform poorly [6]. As a consequence, many authors have proposed alternative estimators which can be divided into two main categories. A first set of approaches assumes structured matrices so as to reduce the number of parameters to estimate, while a second set of approaches aims at compensating the lack of samples with regularization or prior knowledge modelling. The first category includes attempts based on sparsity assumptions such as graphical Lasso, *e.g.* [3, 11, 44], and nodewise Lasso, *e.g.* [13, 27]. Besides not to be always satisfying in small sample size settings (see [19] for a recent review), these methods assume Gaussian observations and are therefore not suitable for elliptical distributions with heavy tails. Generalizations have been considered more recently that are more robust [2, 9], but they require a large number of  $\Sigma$  entries to be zero which may be too restrictive in some applications. In this work, we rather consider estimators in the second category based on shrinkage methods, introduced in [22]. In shrinkage methods, the considered estimators are convex combinations of an initial estimator and the identity matrix view as a regularization term. The construction of these estimators rely then on two main ingredients, the choice of the initial estimator to be regularized and the choice of the regularization parameter, or equivalently the weight of the identity matrix. As already mentioned, when aiming at robust inference, M-estimators are good initial basis. Following this line, the authors in [6] have proposed a shrinkage procedure, with an optimal shrinkage coefficient, for a particular case of M-estimators, called Tyler's estimator where the function  $u(t)$  is set to  $p/t$  [39]. This choice of  $u$  is motivated by the fact that if  $\mathbf{x}$  is elliptically distributed with mean  $\boldsymbol{\mu}$ , then the normalized vector  $\mathbf{z} = (\mathbf{x} - \boldsymbol{\mu}) / \|\mathbf{x} - \boldsymbol{\mu}\|$  follows an angular central Gaussian distribution. This approach has the advantage to be somewhat non-parametric and has shown a lot of merits in various settings [31, 34, 43]. Unfortunately, a serious limit is that it requires the mean  $\boldsymbol{\mu}$  to be known in advance so that the shape of the distribution cannot be taken into account when estimating the mean. This point has been highlighted in [36], which proposes to estimate  $\boldsymbol{\mu}$  assuming  $\mathbf{x}$  follows a Cauchy distribution (*i.e.* a  $t$ -distribution with dof parameter equal to 1), and as a follow-up more recently in [24] with a generalization to any  $t$ -distributions. However, in contrast to [6], none of these papers provide an optimal shrinkage coefficient. Although the effect of tuning this coefficient may be important, the issue is usually eliminated either by searching in a finite grid of values [24, 36]

or using cross-validation [38] in both cases at the cost of a higher computational complexity and time.

We aim at building on these previous approaches by providing both a flexible and optimal framework based on shrunk likelihood-based M-estimators. The distribution of  $\mathbf{x}$  is assumed to be elliptical so that the corresponding function  $u$  and the associated M-estimator can be derived straightforwardly from a maximum likelihood principle. We propose then a shrinkage version of this estimator with an explicit formula for the optimal shrinkage coefficient that depends on two moments of the radius of  $\mathbf{x}$ . Then, a complete inference procedure is proposed which does not require neither to pre-set the value of the mean nor that of the tail parameter. Explicit expressions of the optimal shrinkage coefficient are given for Gaussian and t-distributions and an algorithm for estimating both  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$  is proposed. Experiments on simulated and real data illustrate the good behavior of the proposed method in comparison to other existing methods such as Tyler's estimator, graphical Lasso, etc.

The paper is organized as follows. Section 2 recalls definitions and the main properties of elliptical distributions and M-estimators. The optimal shrinkage problem is addressed in Section 3 with a general formula for the optimal shrinkage coefficient. In the following Section 4, the optimal parameter value is given in the case of multivariate t-distributions together with a practical algorithm to estimate both the mean and covariance matrix in a potentially low sample size setting. The proposed estimator and algorithm are illustrated on simulated and real data respectively in Section 5 and 6. A conclusion ends the paper. At last, all proofs and supplementary results are provided in Appendix.

## 2 Preliminaries

### 2.1 Elliptical distributions

A continuous random vector  $\mathbf{x} \in \mathbb{R}^p$  follows a multivariate elliptical symmetric distribution if its probability density function (pdf) is of the form (see [5] or [20]):

$$p(\mathbf{x}) = C_{p,g} |\boldsymbol{\Sigma}|^{-1/2} g\left((\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right), \quad (1)$$

where  $\boldsymbol{\Sigma} \in \mathbb{R}^{p \times p}$  is the scale matrix with determinant  $|\boldsymbol{\Sigma}|$ ,  $\boldsymbol{\mu} \in \mathbb{R}^p$  is the location or mean vector,  $C_{p,g}$  is a normalizing constant so that  $p(\mathbf{x})$  integrates to one. The non-negative function  $g(\cdot)$  is called the density generator and determines the shape of the pdf. Also, it is important to note that elliptical distributions have the stochastic representation  $\mathbf{x} = \boldsymbol{\mu} + R\boldsymbol{\Lambda}\mathbf{U}$  [5], where  $R$  (called radius) is a non negative random variable,  $\boldsymbol{\Lambda}$  is a  $p \times p$  matrix so that  $\boldsymbol{\Lambda}\boldsymbol{\Lambda}^\top = \boldsymbol{\Sigma}$  and  $\mathbf{U}$  is a  $p$ -dimensional random vector following a uniform distribution on the unit sphere of dimension  $p$  ( $R$  and  $\mathbf{U}$  are independent). The radius  $R$  and the generator  $g$  are closely related. Indeed, according to Theorem 2.9 in [8], an elliptical distribution has a generator if and only if the random variable  $R$  has a density, and there exists a link between these two functions (Theorem 3 in [10] gives a similar result). Throughout this paper, we assume that our elliptical distribution has a generator, and the latter may be defined either by  $g$ , or by its radius  $R$ . In addition, to avoid

identifiability issues (for any scaling factor  $\gamma$ ,  $(\boldsymbol{\mu}, R/\gamma, \gamma\mathbf{A})$  and  $(\boldsymbol{\mu}, R, \mathbf{A})$  lead to the same distribution), we assume in the rest of the paper that the trace of  $\boldsymbol{\Sigma}$  is  $p$ , denoted by  $\text{tr}(\boldsymbol{\Sigma}) = p$ .

This family encompasses a lot of well known particular cases, like the Gaussian distribution (with  $g(t) = \exp(-t/2)$ ) and the Student distribution (also called t-distribution) with  $\nu > 0$  degrees of freedom (with  $g(t) = (1 + t/\nu)^{-(p+\nu)/2}$ ). Other examples include the Logistic [25], Kotz [29], Laplace [7] or Slash [1] distributions.

In this paper we consider the problem of the scale matrix estimation  $\boldsymbol{\Sigma}$  under the assumption that the data is elliptically distributed. It is an important task both in the case of known or unknown location or mean vector  $\boldsymbol{\mu}$ . A lot of methods have already been proposed. For instance, [46] focused on the widely used sample covariance matrix  $\hat{\mathbf{S}} = n^{-1} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$  as an estimator of  $\boldsymbol{\Sigma}$  (the mean vector  $\boldsymbol{\mu}$  is here considered as known, *i.e.* the data has previously been centered). However, being designed for the Gaussian distribution, this method is not suitable for the case of data with outliers. Moreover, it requires the existence of  $\mathbb{E}[\mathbf{x}\mathbf{x}^\top]$ , and this condition is not always fulfilled (see *e.g.* the Cauchy distribution). To overcome these difficulties, [39] proposed another estimator which is a particular case of Maronna's M-estimators [26] detailed in the next section.

## 2.2 M-estimators and Tyler's estimator

Let  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  be a set of  $n$  independent and identically distributed observations drawn from an elliptical distribution (1) with a known mean vector  $\boldsymbol{\mu}$ . Tyler [39] proposed a distribution-free estimator of the trace-normalized covariance matrix by working with the normalized observations  $\mathbf{z}_i = \frac{\mathbf{x}_i - \boldsymbol{\mu}}{\|\mathbf{x}_i - \boldsymbol{\mu}\|}$ . According to [40], each  $\mathbf{z}_i$  follows the angular central Gaussian distribution:

$$p(\mathbf{z}) = \frac{\Gamma(p/2)}{2\pi^{p/2}} |\boldsymbol{\Sigma}|^{-1/2} \left( \mathbf{z}^\top \boldsymbol{\Sigma}^{-1} \mathbf{z} \right)^{-p/2}. \quad (2)$$

Then maximum likelihood principle leads to an implicit estimator  $\tilde{\boldsymbol{\Sigma}}$ , solution of

$$\tilde{\boldsymbol{\Sigma}} = \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{z}_i \mathbf{z}_i^\top}{\mathbf{z}_i^\top \tilde{\boldsymbol{\Sigma}}^{-1} \mathbf{z}_i}. \quad (3)$$

A fixed point algorithm is thus usually used to estimate  $\boldsymbol{\Sigma}$  with a final normalization step to ensure  $\text{tr}(\tilde{\boldsymbol{\Sigma}}) = p$ . Tyler's estimator may then also be seen as a particular case of Maronna's M-estimator. Existence and uniqueness of  $\tilde{\boldsymbol{\Sigma}}$  are discussed in [39]. In particular, it is mentioned that condition  $n > p$  is required. Otherwise, according to [28], matrix  $\tilde{\boldsymbol{\Sigma}}$  is singular and this estimator is no longer suitable. In the case  $p \geq n$ , a regularized Tyler's estimator has been proposed, based on shrinkage methods [6] as specified in the following section.

## 2.3 Regularized Tyler's estimator

Inspired by the shrinkage method of Ledoit and Wolf [22], the authors in [6] extended Tyler's method to the high dimensional setting introducing the following

regularized fixed point equations. The  $t^{\text{th}}$  iteration is indicated with index ( $t$ ):

$$\tilde{\Sigma}^{(t+1)} = (1 - \rho) \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{z}_i \mathbf{z}_i^\top}{\mathbf{z}_i^\top \Sigma^{(t)-1} \mathbf{z}_i} + \rho \mathbf{I}, \quad (4)$$

$$\Sigma^{(t+1)} = \frac{\tilde{\Sigma}^{(t+1)}}{\text{tr}(\tilde{\Sigma}^{(t+1)})/p}. \quad (5)$$

Here  $0 \leq \rho \leq 1$  is a constant which is called shrinkage coefficient. The case of  $\rho = 0$  corresponds to the standard non regularized Tyler's estimator while  $\rho = 1$  reduces the estimator to the identity matrix. The term  $\rho \mathbf{I}$  ensures that the estimator is well-conditioned at each iteration. Both existence and uniqueness of the limit of the procedure (4)-(5) are proved in [6]. The choice of  $\rho$  is also discussed. As in [22], the authors in [6] propose to find parameter  $\rho$  by minimizing the mean-squared error (MSE) between the true matrix  $\Sigma$  and the so-called "clairvoyant estimator":

$$\tilde{\Sigma}_\rho = (1 - \rho) \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{z}_i \mathbf{z}_i^\top}{\mathbf{z}_i^\top \Sigma^{-1} \mathbf{z}_i} + \rho \mathbf{I}. \quad (6)$$

Thus,  $\rho$  is chosen as the solution  $\rho_T^*$  of:

$$\rho_T^* = \arg \min_{\rho} \mathbb{E} \left[ \left\| (1 - \rho) \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{z}_i \mathbf{z}_i^\top}{\mathbf{z}_i^\top \Sigma^{-1} \mathbf{z}_i} + \rho \mathbf{I} - \Sigma \right\|_F^2 \right], \quad (7)$$

where  $\|\cdot\|_F$  is the Frobenius norm. The solution can be seen as the value of  $\rho$  which minimizes the distance between the true  $\Sigma$  and its shrunked deformation. Following the above criteria, an explicit formula for  $\rho_T^*$  is obtained under the assumption  $\text{tr}(\Sigma) = p$ :

$$\rho_T^* = \frac{p^2 + (1 - 2/p)\text{tr}(\Sigma^2)}{(p^2 - np - 2n) + (n + 1 + 2(n - 1)/p)\text{tr}(\Sigma^2)}. \quad (8)$$

In the next section, we propose to generalize this last result, to the case when  $\boldsymbol{\mu}$  is not known and for all M-estimators when the data is sampled from a specified elliptically symmetric distribution (1). Under a criteria similar to (7), we provide a closed-form expression for the optimal shrinkage coefficient.

### 3 Optimal shrinkage for M-estimators

Let  $\mathbf{x}_1, \dots, \mathbf{x}_n$  denote *i.i.d* random realizations of a  $p$ -variate elliptical random vector  $\mathbf{x} = \boldsymbol{\mu} + R\mathbf{A}\mathbf{U}$ . In this section, we suppose that  $\boldsymbol{\mu}$  is known, and consider the class of Maronna's estimators satisfying :

$$\tilde{\Sigma} = m(\tilde{\Sigma}), \quad (9)$$

$$\text{with } m(\Sigma) = \frac{1}{n} \sum_{i=1}^n u\left(\frac{(\mathbf{x}_i - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x}_i - \boldsymbol{\mu})}{p}\right) (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top. \quad (10)$$

By taking  $u(t) = p/t$  and  $\boldsymbol{\mu} = \mathbf{0}$ , we recover Tyler's estimator (3). Some other examples of functions  $u$  are  $u(t) = 1$  [32], the Huber's function [18], or the Student maximum likelihood-based function  $(p + \nu)/(t + \nu)$  [26]. As proposed in [32], in this paper we consider a regularized estimator:

$$\tilde{\boldsymbol{\Sigma}}_\rho = \beta m(\boldsymbol{\Sigma}) + \alpha \mathbf{I}, \quad (11)$$

where  $\alpha = \rho$  and  $\beta = 1 - \rho$ . This choice of  $\alpha$  and  $\beta$  will be discussed and justified later (see proposition 2 below). We define the following criteria, similar to (7), for the choice of  $\rho$ . The optimal  $\rho^*$  will be chosen such as to minimize the MSE between the "clairvoyant estimator"  $\tilde{\boldsymbol{\Sigma}}_\rho$  and  $\boldsymbol{\Sigma}$ :

$$\mathbb{E} \left[ \left\| \frac{1 - \rho}{n} \sum_{i=1}^n u \left( (\mathbf{x}_i - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right) (\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{x}_i - \boldsymbol{\mu})^\top + \rho \mathbf{I} - \boldsymbol{\Sigma} \right\|_F^2 \right]. \quad (12)$$

The next theorem presents the closed-form solution to problem (12) for elliptical distributions. Note that the result holds for any function  $u$ . Alternatively, another MSE criterion has been considered in [32], for which it is also possible to give the optimal  $\rho$  value. This criterion is not further considered in our work but we provide the corresponding optimal  $\rho$  formula in Appendix 8.6.

**Theorem 1 (Optimal shrinkage coefficient)** *Let  $\mathbf{x} = \boldsymbol{\mu} + R\mathbf{A}\mathbf{U}$  be a stochastic representation of the elliptically distributed variable  $\mathbf{x}$ . The oracle coefficient  $\rho^*$  which minimizes (12) is, under the condition  $\text{tr}(\boldsymbol{\Sigma}) = p$ ,*

$$\rho^* = \frac{\text{tr}(\boldsymbol{\Sigma}^2) \left( f_1^2 \frac{n-1}{np^2} + f_2 \frac{2}{np(p+2)} - f_1 \frac{1}{p} \right) + f_2 \frac{p}{n(p+2)} - f_1 + p}{\text{tr}(\boldsymbol{\Sigma}^2) \left( f_1^2 \frac{n-1}{np^2} + f_2 \frac{2}{np(p+2)} \right) + f_2 \frac{p}{n(p+2)} - 2f_1 + p}, \quad (13)$$

where  $f_1 = \mathbb{E} [u(R^2)R^2]$  and  $f_2 = \mathbb{E} [u(R^2)^2R^4]$ .

The proof is provided in Appendix 8.1. When considering for  $u$  the Tyler's function  $u(t) = p/t$ , the result in [6] is recovered. Indeed, it follows  $f_1 = p$  and  $f_2 = p^2$ , which leads in turns to the optimal shrinkage coefficient formula (8).

We can now discuss the choice of the function  $u$ . For this purpose, we propose to take a function motivated by the likelihood. Indeed, the maximum likelihood estimator (MLE) of the covariance matrix minimizes the negative log-likelihood function:

$$\mathcal{L}(\boldsymbol{\Sigma}) = -\frac{2}{n} \sum_{i=1}^n \ln \left( g \left( (\mathbf{x}_i - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right) \right) - \ln |\boldsymbol{\Sigma}^{-1}|. \quad (14)$$

The previous equation leads to an implicit estimator of  $\boldsymbol{\Sigma}$ , obtained through a fixed point algorithm. However, we recall that this approach is no longer suitable if  $p > n$ . In that case, similarly to the approach of [32], we consider the penalized cost function:

$$\mathcal{L}_\rho(\boldsymbol{\Sigma}) = - (1 - \rho) \frac{2}{n} \sum_{i=1}^n \ln \left( g \left( (\mathbf{x}_i - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right) \right) - \ln |\boldsymbol{\Sigma}^{-1}| + \rho \text{tr}(\boldsymbol{\Sigma}^{-1}). \quad (15)$$

In what follows, we suppose that the generator  $g$ , or equivalently the density of  $R$ , is differentiable. The solution  $\tilde{\Sigma}_\rho$  which minimizes the penalized cost function  $\mathcal{L}_\rho(\Sigma)$  can be expressed as:

$$\tilde{\Sigma}_\rho = (1 - \rho) \frac{1}{n} \sum_{i=1}^n u \left( (\mathbf{x}_i - \boldsymbol{\mu})^\top \tilde{\Sigma}_\rho^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right) (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top + \rho \mathbf{I}, \quad (16)$$

with  $u(t) = -2g'(t)/g(t)$ .

It is interesting to note that penalizations are linked to prior choice for  $\Sigma$  in a Bayesian framework. In (15) above, the  $\text{tr}(\Sigma^{-1})$  penalization corresponds to an inverse Wishart prior where the scale matrix hyperparameter is the identity matrix. For a more general scale matrix hyperparameter  $\mathbf{T}$ , the penalty would be  $\text{tr}(\Sigma^{-1}\mathbf{T})$  and we would get a regularized estimator similar to (16) with a penalty term replaced by  $\rho\mathbf{T}$ . Theorem 1 can then be generalized to this case. This result and its proof are given in Appendix 8.2.

We thus defined a regularized M-estimator with function  $u(t) = -2g'(t)/g(t)$ . The main difficulty is now to calculate the terms  $f_1$  and  $f_2$  in Theorem 1. The next result provides a general formula for  $f_1$  while  $f_2$  has to be computed differently for each distribution.

### 3.1 Theoretical value of $\rho^*$

Let  $u(t) = -2g'(t)/g(t)$  (the radius  $R$  is thus supposed to have a differentiable pdf). It is remarkable that for most elliptical distributions,  $f_1 = p$  in that case.

**Proposition 1 (Value of  $f_1$ )** *Let  $\mathbf{x} = \boldsymbol{\mu} + R\mathbf{A}\mathbf{U}$ , where  $R$  is a positive random variable with a differentiable pdf  $f_R$  such that  $f_R(r) > 0$  for all  $r > 0$ ,  $rf_R(r) \rightarrow 0$  as  $r$  goes to either 0 or  $+\infty$ . If  $u(t) = -2g'(t)/g(t)$ , then  $f_1 = p$ .*

The proof is provided in Appendix 8.3. The following corollary follows straightforwardly.

**Corollary 1** *Under the assumptions of Theorem 1 and Proposition 1:*

$$\rho^* = \frac{\text{tr}(\Sigma^2) \left( f_2 \frac{2}{np(p+2)} - \frac{1}{n} \right) + f_2 \frac{p}{n(p+2)}}{\text{tr}(\Sigma^2) \left( 1 - \frac{1}{n} + f_2 \frac{2}{np(p+2)} \right) + f_2 \frac{p}{n(p+2)} - p}. \quad (17)$$

Note that the conditions of Proposition 1 encompass all the well known elliptical distribution (Gaussian, Student and all the examples mentioned in the paper). If  $f_1$  has a general expression, this is unfortunately not the case for  $f_2$ . In the following section, we focus on the multivariate t-distribution, for which these two terms can be explicitly calculated. We also provide for illustration the case of the Kotz-type distribution. Beforehand, we first justify the choice of  $\alpha = \rho$  and  $\beta = 1 - \rho$  in (11).



### 3.2 Optimal regularized estimators

Consider the regularized covariance matrix in the general form:

$$\tilde{\Sigma}_{\alpha\beta} = \beta m(\Sigma) + \alpha \mathbf{I}, \quad (18)$$

where  $\alpha \geq 0, \beta \geq 0$ . Next proposition provides the relationship between the oracle shrinkage coefficients  $\alpha$  and  $\beta$  for functions  $u$  which are derived from maximum likelihood estimation of the covariance matrix of an elliptical distribution. The coefficients are found under the same criteria as in (12). Thus we are looking for the optimal  $\alpha$  and  $\beta$  such that:

$$(\alpha^*, \beta^*) = \arg \min_{\alpha, \beta \geq 0} \mathbb{E} \left[ \left\| \frac{\beta}{n} \sum_{i=1}^n u \left( (\mathbf{x}_i - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right) (\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{x}_i - \boldsymbol{\mu})^\top + \alpha \mathbf{I} - \Sigma \right\|_F^2 \right]. \quad (19)$$

**Proposition 2** *Let  $\mathbf{x} = \boldsymbol{\mu} + R\mathbf{A}\mathbf{U}$  be a stochastic representation of an elliptically distributed variable  $\mathbf{x}$ . The coefficients  $\alpha^*$  and  $\beta^*$  are related through the following formula:  $\alpha = \frac{\text{tr}(\Sigma)}{p} (1 - \frac{f_1}{p} \beta)$ , where  $f_1 = \mathbb{E} [u(R^2)R^2]$ .*

The proof is provided in Appendix 8.4. From the proposition we directly get that, under the assumption  $\text{tr}(\Sigma) = p$ , coefficient  $\alpha = 1 - \frac{f_1}{p} \beta$  and under Proposition 1 for functions  $u$  coming from maximum likelihood, it comes moreover that  $\alpha = 1 - \beta$ . However, general formulas for  $\alpha$  and  $\beta$  are also provided in the proof of Proposition 2 in the Appendix.

## 4 Regularized covariance matrix estimator for the multivariate t-distribution

Theorem 1 provides a general result for all elliptical distributions. However, equation (13) depends on the moments  $f_1 = \mathbb{E} [u(R^2)R^2]$  and  $f_2 = \mathbb{E} [(u(R^2))^2 R^4]$ , so that to be useful in practice, expressions for  $\rho^*$  require further developments. In the following section, we provide further results allowing the specification of  $\rho^*$  in important cases such as the Gaussian and t-distributions. For both distributions, Proposition 1 applies and  $f_1 = p$ . The following result provides the values of  $f_2$ . In the sequel,  $\chi_k^2$  denotes the Chi-squared distribution with  $k$  degrees of freedom and  $F_{p,\nu}$  denotes the Fisher distribution with dof parameters  $p$  and  $\nu$ .

**Proposition 3 (Some values of  $f_2$ )** *Let  $\mathbf{x} = \boldsymbol{\mu} + R\mathbf{A}\mathbf{U}$ , where  $R$  is a positive random variable with a differentiable pdf, and  $u(t) = -2g'(t)/g(t)$ .*

1. *If  $R^2$  is distributed as  $\chi_p^2$ , then  $\mathbf{x}$  follows a Gaussian distribution,  $u(t) = 1$  and therefore:*

$$\mathbb{E} \left[ u \left( R^2 \right)^2 R^4 \right] = p(p+2). \quad (20)$$

2. If  $R^2$  is distributed as  $\left(\frac{1}{2\lambda}\chi_{2q+p-2}^2\right)^{1/s}$ , then  $\mathbf{x}$  follows a Kotz-type distribution,  $u(t) = 2(1-q)/t + 2s\lambda t^{s-1}$  and therefore:

$$\mathbb{E}\left[u\left(R^2\right)^2 R^4\right] = p(p+2s) + 4s(q-1). \quad (21)$$

3. If  $R^2$  is distributed as  $pF_{p,\nu}$ , then  $\mathbf{x}$  follows a  $t$ -distribution with  $\nu > 0$  degrees of freedom,  $u(t) = (p+\nu)/(t+\nu)$  and therefore :

$$\mathbb{E}\left[u\left(R^2\right)^2 R^4\right] = \frac{(\nu+p)(p+2)p}{p+\nu+2}. \quad (22)$$

The proof is provided in Appendix 8.5. The above formulas for  $f_2$  are consistent with the fact that when  $\nu$  goes to  $+\infty$ , the  $t$ -distribution tends to the Gaussian distribution and expression (22) tends to (20). A similar check can be done using that the Gaussian distribution is a particular case of Kotz-type distributions with  $s = 1$ ,  $q = 1$  and  $\lambda = 1/2$ .

Combining Theorem 1 and Propositions 1 and 3, the optimal shrinkage coefficient can be specified for the above distributions. In the following result, we restrict to the Gaussian and  $t$ -distributions.

**Corollary 2 (Optimal shrinkage coefficient for multivariate Gaussian and  $t$ -distributions)** *Assume  $\text{tr}(\boldsymbol{\Sigma}) = p$ , the optimal shrinkage coefficient is given by,*

1. For the Gaussian distribution:

$$\rho^* = \frac{\text{tr}(\boldsymbol{\Sigma}^2) + p^2}{\text{tr}(\boldsymbol{\Sigma}^2)(n+1) + p^2 - pn}. \quad (23)$$

2. For the  $t$ -distribution with  $\nu > 0$  degrees of freedom:

$$\rho^* = \frac{\text{tr}(\boldsymbol{\Sigma}^2) \left(1 + \frac{\nu}{p} - \frac{2}{p}\right) + p(\nu+p)}{\text{tr}(\boldsymbol{\Sigma}^2) \left((n+1) \left(\frac{\nu}{p} + 1\right) + \frac{2}{p}(n-1)\right) + (p+\nu)(p-n) - 2n}. \quad (24)$$

We now have an explicit formula for our optimal shrinkage coefficient in the  $t$ -distribution case including (for  $\nu = 0$ ) the Tyler's coefficient  $\rho_T^*$  specified in (8). In practice, it still remains to compute estimations for  $\text{tr}(\boldsymbol{\Sigma}^2)$  and  $\nu$  to get values for  $\rho^*$ . This is specified in the next section.

#### 4.1 Estimation of the oracle coefficients

The oracle value  $\rho^*$  cannot be implemented as such because it is a function of unknown quantities  $\text{tr}(\boldsymbol{\Sigma}^2)$  and  $\nu$ . Following [6, 30], a good candidate for plug-in for  $\boldsymbol{\Sigma}$  is the trace-normalised normalised sample covariance matrix  $\mathbf{S}$  defined by:

$$\mathbf{S} = \frac{p}{n} \sum_{i=1}^n \frac{(\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top}{\|\mathbf{x}_i - \boldsymbol{\mu}\|^2}. \quad (25)$$

As observed in [6], the only requirement on such an estimator is that it provides a good approximation to  $\text{tr}(\boldsymbol{\Sigma}^2)$ . It does not have to be well-conditioned nor does it have to be an accurate estimator of the true  $\boldsymbol{\Sigma}$ . As a matter of fact,  $\text{tr}(\mathbf{S}^2)$  has been shown in [45] to be highly robust for the estimation of  $\text{tr}(\boldsymbol{\Sigma}^2)$  when  $\text{tr}(\boldsymbol{\Sigma}) = p$ . Regarding the dof parameter  $\nu$ , for t-distributions, the norm  $\|\mathbf{x}\|_2$  is regularly varying with tail index  $1/\nu$  (see [16]), *i.e.* :

$$\forall t > 0, \lim_{z \rightarrow +\infty} \frac{\mathbb{P}(\|\mathbf{x}\|_2 > tz)}{\mathbb{P}(\|\mathbf{x}\|_2 > z)} = t^{-1/\nu}. \quad (26)$$

Different estimators of the tail index are available in the literature, the most popular and widespread being the Hill estimator, introduced in [17]. By taking the inverse of the latter, we define an estimator  $\hat{\nu}_{k_n}$  of  $\nu$  :

$$\hat{\nu}_{k_n} = \left( \frac{1}{k_n} \sum_{i=1}^{k_n} \ln \left( \frac{\|\mathbf{x}_{[i]}\|_2}{\|\mathbf{x}_{[k_n+1]}\|_2} \right) \right)^{-1}, \quad (27)$$

where  $\mathbf{x}_{[i]}$  denoted the ordered observations such that  $\|\mathbf{x}_{[1]}\|_2 \geq \dots \geq \|\mathbf{x}_{[k_n+1]}\|_2 \geq \dots \geq \|\mathbf{x}_{[n]}\|_2$ . The Hill estimator, and therefore  $\hat{\nu}_{k_n}$ , are related to a number  $k_n$ . On a theoretical point of view, the latter has to fulfill  $k_n \rightarrow +\infty$  and  $k_n/n \rightarrow 0$  as  $n \rightarrow +\infty$ , and leads to a compromise between a variable and biased estimation of  $\nu$ . Indeed, a small  $k_n$  leads to a low biased and high variable estimation, while a large  $k_n$  increases the bias and reduces the variance (see Section 3.2 in [15] for details). Therefore, the choice of  $k_n$ , usually chosen as  $\lfloor n^b \rfloor$ ,  $0 < b < 1$ , is an important point, and is discussed in [42] in the t-distribution case. According to [42], a choice of  $b$  around  $4/(\nu + 4)$  is suitable.

#### 4.2 Joint covariance matrix and mean estimation

In the previous sections we assumed that the mean vector  $\boldsymbol{\mu}$  was known. When the mean vector is unknown one can use the sample mean  $\hat{\boldsymbol{\mu}} = 1/n \sum_{i=1}^n \mathbf{x}_i$  as an estimator. However, it is well-known that the sample mean may have a bad performance in the presence of outliers or simply in high-dimensional case. Thus, the estimation of the covariance matrix can be severely degraded by an inaccurate estimation of the mean vector  $\boldsymbol{\mu}$ . To overcome this difficulty, we focus on the joint mean - covariance matrix estimation as the solution of a system of equations of the following form defining Maronna's M-estimators [26]:

$$0 = \frac{1}{n} \sum_{i=1}^n h \left( (\mathbf{x}_i - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right) (\mathbf{x}_i - \boldsymbol{\mu}), \quad (28)$$

$$\boldsymbol{\Sigma} = \frac{1}{n} \sum_{i=1}^n u \left( (\mathbf{x}_i - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right) (\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{x}_i - \boldsymbol{\mu})^\top, \quad (29)$$

where the functions  $h$  and  $u$  satisfy a set of general assumptions stated in [26]. We propose to extend this approach to the high-dimensional/small sample size case by replacing the covariance matrix by its regularized estimator. Moreover, as

before, the trace normalized covariance matrix is considered. It follows the iterative algorithm below for  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ . Denoting the  $t^{\text{th}}$  iteration with index ( $t$ ):

$$\boldsymbol{\mu}^{(t+1)} = \frac{\sum_{i=1}^n h \left( (\mathbf{x}_i - \boldsymbol{\mu}^{(t)})^\top \boldsymbol{\Sigma}^{(t)-1} (\mathbf{x}_i - \boldsymbol{\mu}^{(t)}) \right) \mathbf{x}_i}{\sum_{i=1}^n h \left( (\mathbf{x}_i - \boldsymbol{\mu}^{(t)})^\top \boldsymbol{\Sigma}^{(t)-1} (\mathbf{x}_i - \boldsymbol{\mu}^{(t)}) \right)}, \quad (30)$$

$$\begin{aligned} \tilde{\boldsymbol{\Sigma}}^{(t+1)} &= (1 - \rho) \frac{1}{n} \sum_{i=1}^n u \left( (\mathbf{x}_i - \boldsymbol{\mu}^{(t+1)})^\top \boldsymbol{\Sigma}^{(t)-1} (\mathbf{x}_i - \boldsymbol{\mu}^{(t+1)}) \right) \\ &\quad \times (\mathbf{x}_i - \boldsymbol{\mu}^{(t+1)}) (\mathbf{x}_i - \boldsymbol{\mu}^{(t+1)})^\top + \rho \mathbf{I}, \end{aligned} \quad (31)$$

$$\boldsymbol{\Sigma}^{(t+1)} = \frac{\tilde{\boldsymbol{\Sigma}}^{(t+1)}}{\text{tr}(\tilde{\boldsymbol{\Sigma}}^{(t+1)})/p}. \quad (32)$$

The optimal coefficient  $\rho$ , in the general case for any distribution in the elliptical family, can be found using the obtained result (13). This value is a function of  $\text{tr}(\boldsymbol{\Sigma}^2)$  and of the parameters defining the distribution of the radius  $R$ . While the term  $\text{tr}(\boldsymbol{\Sigma}^2)$  can be estimated using the trace-normalized normalized sample matrix  $\mathbf{S}$ , the estimation of the  $R$  distribution parameters may not be obvious and requires additional examination. One possible solution is to set these parameters to fixed values representing the prior knowledge about the distribution of  $R$ . Another solution is to choose the parameters values corresponding to the heaviest tail case within the chosen distribution subclass. For example, for the multivariate  $t$ -distribution, the authors in [36] propose to set  $\nu = 1$ , in other words, to focus on the Cauchy distribution, which corresponds to a heavy-tail representative among the subclass of  $t$ -distributions. To keep more flexibility on the choice of these parameters, they could also be estimated using a maximum likelihood approach, possibly using the EM-algorithm. The EM algorithm is tractable for a subclass of the elliptical family referred to as Gaussian scale mixtures (GSM) [14]. GSM distributions include the generalized Gaussian distribution, the multivariate  $t$ -distribution, the Pearson type VII distribution, etc.

Hereafter, Algorithm 1 provides a simple algorithm for the joint mean - covariance matrix estimation for the multivariate  $t$ -distribution. This algorithm does not require an additional step for the estimation of the dof  $\nu$  in contrast to an EM-algorithm implementation that would iteratively update  $\nu$  as well as  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ . Instead, this parameter is evaluated once through the Hill estimator (27) and employed then for both the estimation of  $\rho^*$  and that of  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ .

## 5 Results on simulated data

In this section we conduct a simulation study to illustrate the performance of the proposed shrinkage approach through Algorithm 1. In our experiments, an autoregressive (AR) covariance structure which satisfies  $\text{tr}(\boldsymbol{\Sigma}) = p$  is considered :

$$(\boldsymbol{\Sigma})_{ij} = r^{|i-j|}, \quad r \in (0, 1). \quad (35)$$

When  $r$  tends to 0, then  $\boldsymbol{\Sigma}$  is close to an identity matrix; when  $r$  tends to 1, then  $\boldsymbol{\Sigma}$  is close to a singular matrix of rank 1. As pointed out earlier, the choice of  $b$  for

---

**Algorithm 1** Small sample size covariance matrix and mean estimation for a multivariate t-distribution
 

---

- 1: Initialize  $\Sigma$  to  $\Sigma^{(0)}$  an arbitrary positive definite matrix and  $\mu$  to  $\mu^{(0)}$  an arbitrary vector
- 2: Estimate dof  $\hat{\nu}$  using (27) and  $\hat{\rho}^*$  using (24) and (25) with  $\mu$  set to  $\mu^{(0)}$ :

$$\hat{\rho}^* = \frac{\text{tr}(\mathbf{S}^2) \left(1 + \frac{\hat{\nu}}{p} - \frac{2}{p}\right) + p(\hat{\nu} + p)}{\text{tr}(\mathbf{S}^2) \left((n+1) \left(\frac{\hat{\nu}}{p} + 1\right) + \frac{2}{p}(n-1)\right) + (p + \hat{\nu})(p - n) - 2n}.$$

- 3: Iterate the following steps until convergence. On step  $(t+1)$ :
  - 3.1: Update the mean vector  $\mu$  as :

$$\mu^{(t+1)} = \frac{\sum_{i=1}^n \bar{w}_i^{(t+1)} \mathbf{x}_i}{\sum_{i=1}^n \bar{w}_i^{(t+1)}}, \text{ where } \bar{w}_i^{(t+1)} = \frac{\hat{\nu} + p}{\hat{\nu} + (\mathbf{x}_i - \mu^{(t)})^\top (\Sigma^{(t)})^{-1} (\mathbf{x}_i - \mu^{(t)})}. \quad (33)$$

- 3.2: Compute matrix  $\tilde{\Sigma}^{(t+1)}$  :

$$\tilde{\Sigma}^{(t+1)} = (1 - \hat{\rho}^*) \frac{p + \hat{\nu}}{n} \sum_{i=1}^n \frac{(\mathbf{x}_i - \mu^{(t+1)})(\mathbf{x}_i - \mu^{(t+1)})^\top}{(\mathbf{x}_i - \mu^{(t+1)})^\top (\Sigma^{(t)})^{-1} (\mathbf{x}_i - \mu^{(t+1)}) + \hat{\nu}} + \hat{\rho}^* \mathbf{I}, \quad (34)$$

- 3.3: Update  $\Sigma$  as the trace-normalized  $\tilde{\Sigma}^{(t+1)}$  :

$$\Sigma^{(t+1)} = p \frac{\tilde{\Sigma}^{(t+1)}}{\text{tr}(\tilde{\Sigma}^{(t+1)})}.$$


---

computing the degrees of freedom parameter  $\hat{\nu}_{k_n}$  in (27) is an important issue. The optimal theoretical value of this parameter depends on the true  $\nu$  through formula  $b \approx 4/(\nu+4)$ . Here we choose  $b = 0.25$  which corresponds to a suitable value for  $\nu \leq 12$ . This allows a robust estimation of  $\nu$  both for heavy-tailed distributions (*e.g.* with  $\nu = 1$  corresponding to a Cauchy distribution) and light-tailed distributions (*e.g.*  $\nu = 12$ ).

In the first experiment, we simulate data from a multivariate t-distribution in dimension  $p = 50$ , with the following different choices of dof parameter  $\nu = \{1, 2, 3, 6, 10\}$  and various AR schemes with  $r = \{0.1, 0.5, 0.9\}$ . The mean  $\mu$  is set to the vector with all components equal to 5. For each pair of parameters  $(\nu, r)$  and for a sample size  $n$  varying from 5 to 50, 100 data sets are generated leading to estimations  $\hat{\mu}_s, \hat{\Sigma}_s$  for  $s = 1$  to 100. The performance of a method is then assessed using the normalized mean square-error (NMSE) for both  $\Sigma$  and  $\mu$ :

$$NMSE(\Sigma) = \frac{\mathbb{E} \left\{ \|\hat{\Sigma} - \Sigma\|_F^2 \right\}}{\|\Sigma\|_F^2} \approx \frac{1}{100} \sum_{s=1}^{100} \frac{\|\hat{\Sigma}_s - \Sigma\|_F^2}{\|\Sigma\|_F^2}, \quad (36)$$

$$NMSE(\mu) = \frac{\mathbb{E} \left\{ \|\hat{\mu} - \mu\|_F^2 \right\}}{\|\mu\|_F^2} \approx \frac{1}{100} \sum_{s=1}^{100} \frac{\|\hat{\mu}_s - \mu\|_F^2}{\|\mu\|_F^2}. \quad (37)$$

For each data set, four different algorithms are then used leading to four different estimators of  $\Sigma$  and three different estimators of  $\mu$ :

- Algorithm 1 referred to as "t-dist" where both  $\Sigma$  and  $\mu$  are computed iteratively;

- Algorithm 1 restricted to steps 3.2 and 3.3, setting  $\boldsymbol{\mu}$  to the sample mean (sm) and estimating  $\boldsymbol{\Sigma}$  on the data centered by the sample mean (sm). We denote this algorithm by "t-dist sm";
- Algorithm 1 with the theoretical oracle value  $\rho^*(\boldsymbol{\Sigma}, \nu)$  in (24) obtained with the true values of  $\boldsymbol{\Sigma}$  and  $\nu$  denoted as "Oracle".
- The algorithm proposed in [6] that estimates  $\boldsymbol{\Sigma}$  on the data centered by the sample mean, referred to as "Tyler sm".

To illustrate the impact of the estimation of  $\boldsymbol{\mu}$ , a comparison of the  $NMSE(\boldsymbol{\mu})$  values is also provided for the three estimators resulting from the above algorithms, namely :

- the sample mean denoted by "sm" ;
- the estimation of  $\boldsymbol{\mu}$  from Algorithm 1 denoted by "t-dist";
- the  $\boldsymbol{\mu}$  estimation from Algorithm 1 with value  $\rho^*(\boldsymbol{\Sigma}, \nu)$  obtained with the true  $\boldsymbol{\Sigma}$  and  $\nu$  in (24) denoted as "Oracle".

The results for the covariance matrix and the mean vector  $\boldsymbol{\mu}$  are illustrated in Figure 1 while more complete results are shown in Figures 7 and 8 in Appendix 8.7. In all scenarios, the proposed "t-dist" algorithm 1 is consistently the closest to the "Oracle" procedure, which as expected always provides the best results with the minimal  $NMSE(\boldsymbol{\Sigma})$  and  $NMSE(\boldsymbol{\mu})$ . More specifically, the "t-dist" performance is very close to the ideal oracle estimator with increasing  $n$ . For both "Oracle" and "t-dist" methods,  $NMSE(\boldsymbol{\Sigma})$  steadily decreases as  $n$  increases. In contrast, the "Tyler sm" and "t-dist sm" algorithms show in some cases irregular NMSE curves due to a bad estimation of  $\boldsymbol{\mu}$  by the sample mean. As illustrated in Figure 1(b) and Figure 8(a,b,c), this is particularly so for small dof like  $\nu = 1$ . This confirms the potential limits of methods that do not estimate  $\boldsymbol{\mu}$  accurately as pointed out in [36]. As regards the impact of  $\nu$ , it is not easy to illustrate separately the effect of  $\nu$  from that of  $\boldsymbol{\mu}$ . To do so we consider the comparison of "Tyler sm" and "t-dist sm" algorithms. In the first step, both procedures employ the sample mean vector to center the original data. In the second step, the covariance matrix is computed using the iteration given in (4) for Tyler's estimator with  $\hat{\rho}_{T^*}$  defined in (8) and iteration given in (34) for "t-dist sm" with  $\hat{\rho}^*$  in (24). Tyler's procedure can be viewed as an extreme case of the multivariate t-distribution with  $\nu = 0$ . Accordingly, we expect that the difference between the two estimators becomes more significant for larger values of  $\nu$ . The better results provided by "t-dist sm" over "Tyler sm" is shown more specifically on Figure 3 for  $\nu = 10$  where the main gains appear for small sample sizes (Figure 3(b)). Algorithm "t-dist sm" also outperforms "Tyler sm" for smaller dof like  $\nu = 1$  but with an increasing gain as  $n$  becomes larger (see Figure 3(c,d)). However, for small  $\nu$  the differences are somewhat less visible due to larger NMSEs coming mainly from a bad estimation of  $\boldsymbol{\mu}$ . When  $\nu$  is large, the resulting curve for "t-dist sm" coincides with the one for "t-dist". This is not surprising as the t-distribution tends to the Gaussian distribution when  $\nu$  tends to  $\infty$  so that the mean vector in (33) becomes closer to the sample mean. This result is confirmed by Figure 1(d,f) and Figure 8(j,k,l) for  $\nu = 6$  and Figure 8(m), (n), (o) for  $\nu = 10$ . Overall estimating the dof parameter  $\nu$  in (34) is important, especially in the small sample size regime, and more generally because it allows a better estimation of the mean vector  $\boldsymbol{\mu}$  which appears to have

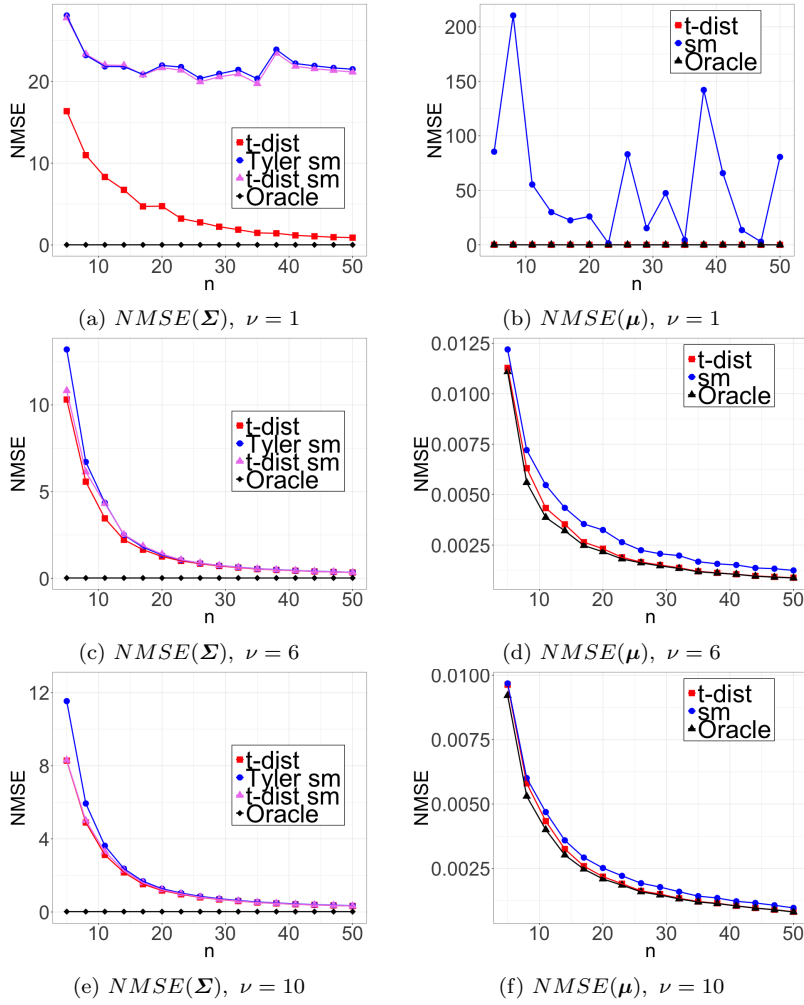


Fig. 1: Multivariate  $t$ -distribution with AR( $r$ ) covariance structure ( $p = 50$ ,  $r = 0.1$ ,  $\nu \in \{1, 6, 10\}$  and  $\mu$  is set to a vector of 5). Normalized mean squared-errors for  $\Sigma$  (first column) and  $\mu$  (second column) are computed over 100 simulated samples of  $n$  observations each with  $n$  varying from 5 to 50.

a critical impact of the covariance structure estimation. For both aspects, our proposed procedure improves over the regularized Tyler's algorithm. Regarding the impact of the choice of  $r$ , it does not seem to lead to significantly different conclusions (see Figures 7 and 8 in the Appendix).

## 6 Application to real data.

Robust estimation of covariance matrices is especially needed for real data where we know that Gaussian hypotheses are generally not true. This is the case for

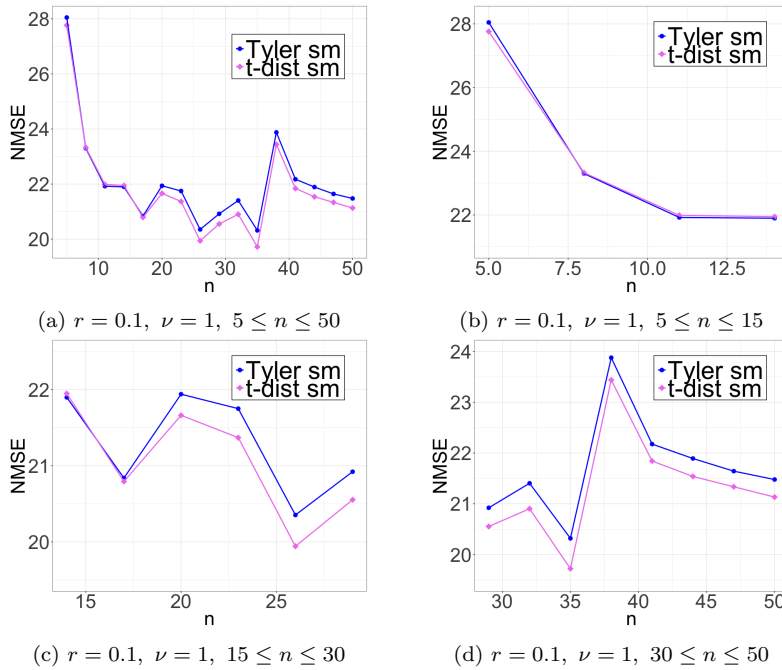


Fig. 2: Multivariate  $t$ -distribution with AR( $r$ ) covariance structure ( $p = 50$ ,  $r = 0.1$ ,  $\nu = 10$  and  $\boldsymbol{\mu}$  is set to a vector of 5): comparison of "Tyler sm" and "t-dist sm" algorithms when the mean is fixed to the sample mean. Normalized mean squared-errors are computed over 100 simulated samples of  $n$  observations each. Varying values of  $n$  from 5 to 50 (a) are also plotted separately for better layout in (b,c,d).

the inference of brain connectivity. Thanks to noninvasive neuroimaging, brain recordings are now available to follow the activity of the brain during a task or at rest. Using functional magnetic resonance imaging (fMRI), one volume of the brain is acquired every one second or less for several minutes. Usually each volume is composed of thousands of voxels that are gathered into a set of hundreds of parcels or brain regions. Each brain region is then associated to a time series. These data are still too complex to provide an easy visualisation and interpretation. Brain connectivity graphs or networks are then constructed by defining nodes as brain regions and edges as connections between time series associated to these brain regions. This allows a spatio-temporal modeling of the brain while functioning.

In this application, to quantify the links between time series, edges are associated to partial correlations read on inverse covariance matrices. Inference of brain connectivity graphs depends then on accurate estimation of covariance or precision matrices. In order to compare objectively different methods, we use a test-retest dataset selected from a larger dataset publicly released as part of the Human Connectome Project (HCP), WU-Minn Consortium (<https://www.humanconnectome.org/>). More details can be found in [37]. We select 100 subjects who have been scanned twice in two different sessions of about 15 minutes each. These two ses-



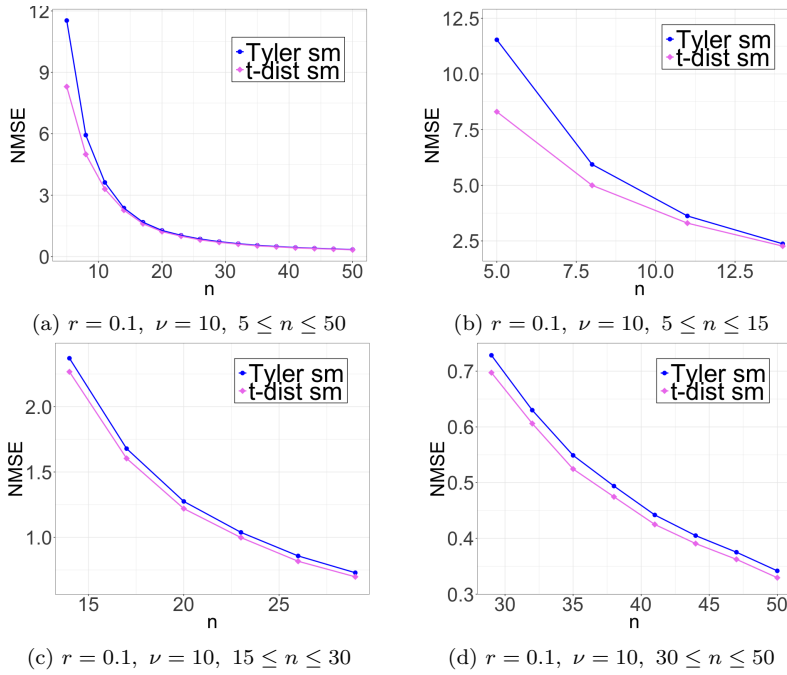


Fig. 3: Moderate-tailed  $t$ -distribution ( $\nu = 10$ ) with AR( $r$ ) covariance structure ( $p = 50$ ,  $r = 0.1$  and  $\boldsymbol{\mu}$  is set to a vector of 5): comparison of "Tyler sm" and "t-dist sm" algorithms when the mean is fixed to the sample mean. Normalized mean squared-errors are computed over 100 simulated samples of  $n$  observations each. Varying values of  $n$  from 5 to 50 (a) are also plotted separately in (b,c,d).

sions are divided into two sub-sessions of half duration and denoted respectively by  $S_{11}, S_{12}$  for session 1 and  $S_{21}, S_{22}$  for session 2. Shorter sessions are more common in practice and they represent a higher challenge for the tested approaches. Following [37], fMRI time series are analyzed through their wavelets decompositions providing vectors of wavelets coefficients resulting in datasets of size  $n = 547$ . The number of brain regions is set to  $p = 90$  based on a commonly used parcellation of the brain into 90 regions [41].

The reliability of each structure learning approach is then evaluated by measuring for each subject, the graph properties differences for the four pairs of sub-sessions coming from different sessions, namely  $(S_{11}, S_{21})$ ,  $(S_{11}, S_{22})$ ,  $(S_{12}, S_{21})$ ,  $(S_{12}, S_{22})$ . For each datasets, five different approaches are considered: sample partial correlation *sample pc*, shrinkage using the Ledoit-Wolf's estimator *lw* in [22], Maximum-likelihood for a  $t$ -distribution using EM *t-dist em* [23], graphical Lasso *glasso* using [11] (the tuning parameter  $\lambda$  in *glasso* algorithm for each subject in each sub-session was obtained by cross-validation with  $k = 3$  folds) and our shrinkage method for a  $t$ -distribution *t-dist shrink*. In order to produce graphs with a fixed comparable number of edges, we apply soft-thresholding to each obtained matrix. For each subject in each sub-session, we then obtain an adjacency matrix that defines an unweighted graph for which a graph metric called *global efficiency*

is computed. This metric is related to the communication efficiency of a node  $i$  with all other nodes (detailed information can be found in [35]). If  $G = (V, E)$  denotes a graph with  $V$  as its set of  $p$  vertices and  $E$  as its edge set, the global efficiency  $Eglob_i$  is defined as the inverse of the harmonic mean of the set of the minimum path lengths  $L_{ij}$  between node  $i \in V$  and all other nodes  $j \in V$  in the graph:

$$Eglob_i = \frac{1}{p-1} \sum_{j \in V} \frac{1}{L_{ij}}. \quad (38)$$

Here  $p$  is the number of brain regions. Then by averaging these global efficiency values over all nodes, one value of this parameter is derived for a given graph. Consequently, for a given pair of sub-sessions, 200 global efficiency values, one per each subject in each sub-session, are computed for a given pre-set percentage of edges in the graphs.

### 6.1 Brain connectivity graphs

As an illustration, Figure 4 displays, for five subjects and sub-sessions  $S_{11}$  and  $S_{22}$ , the global efficiency computed for each region of the brain with either sample partial correlations or partial correlations using our shrinkage method. The global efficiency values are on average between 0.35 for the right Precentral region and 0.53 for the left Putamen region. Regions with high global efficiency greater than 0.5 include the post Cingulum, Amygdala, Frontal Middle Orbital, Occipital Inferior and Thalamus regions. Whereas regions with low efficiency less than 0.4 include the Precentral, Postcentral, Parietal Superior and Frontal Superior regions. The qualitative comparison between the two sub-sessions highlights a higher similarity and reproducibility between sessions with our shrinkage method *t-dist shrink*, Figure 4 (b), than with sample partial correlation *sample pc*, Figure 4 (a). Similar results are observed for other subjects and other pairs of sub-sessions.

### 6.2 Test-retest reliability

To quantify more specifically the differences between the five tested methods, we evaluate their ability to provide similar results between two sessions via the so-called intraclass correlation coefficient (ICC) between the sessions. Using the global efficiency values for each subject in each session, we compute their within-subject ( $s_w$ ) and between-subject ( $s_b$ ) mean square differences, as detailed in the Appendix of [37]. In our case, with two sessions, the ICC is then given by

$$ICC = \frac{s_b - s_w}{s_b + s_w}. \quad (39)$$

When for each subject, similar global efficiency values are found in the two sessions, then the ICC is close to 1 and the reliability is high. In contrast, ICC is close to 0 when the reliability is low. The ICC may take negative values when the variance within subjects is larger than between subjects. This is due to statistical errors given a particular dataset and should be considered as a non reliable estimation.

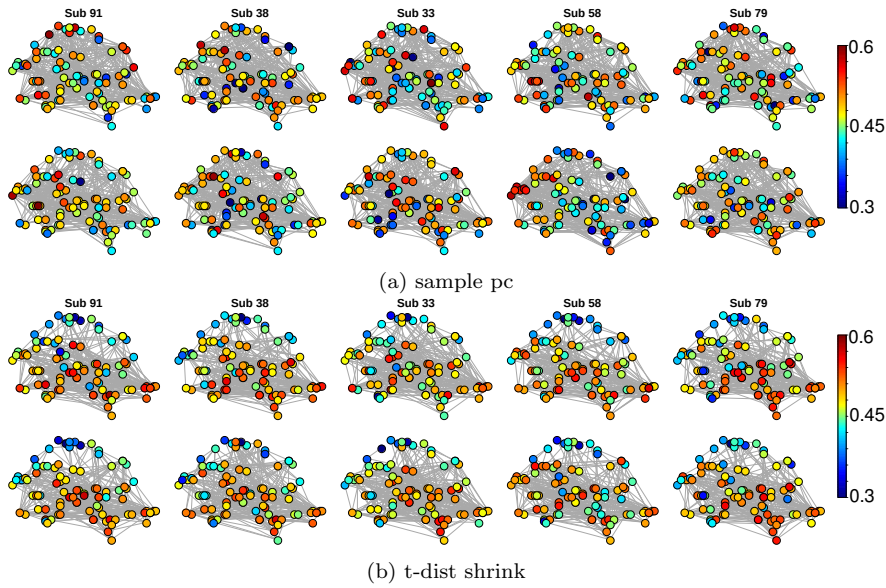


Fig. 4: Global efficiency per brain region (node) for brain graphs with 10% of edges. Five subjects are displayed (columns). Two sessions  $S_{11}$  (first and third rows) and  $S_{22}$  (second and fourth rows) are compared using (a) sample partial correlations and (b) our proposed shrinkage approach. A high global efficiency means that the node/region is well connected to other nodes. Hubs generally show a high global efficiency. In contrast, a low global efficiency means that the shortest path length is large, and is typical of less connected nodes.

Figure 5 represents, for sessions  $S_{11}$  and  $S_{22}$ , ICC values with respect to the pre-set percentage of edges in the graphs, referred to as the cost. Shrinkage methods ( $lw$  and  $t\text{-dist shrink}$ ) have the largest range of costs where ICC is above 0.4, Figure 5(a). This is confirmed by the computation of  $p$ -values, Figure 5 (b), where  $p$ -values allow to check whether the ICC is significantly larger than zero. The ICC computed using empirical partial correlation and EM for a  $t$ -distribution are nearly equal to 0 for any cost showing a poor reliability of global efficiency. Glasso is also showing poor performance because of the difficulties to choose the regularisation parameters. For additional confirmation that shrinkage schemes provide better results than sample partial correlation, we further check  $s_b$  values as an increase in  $s_b$  may artificially increase ICC values. Figure 6 displays  $s_b$  and  $s_w$  values obtained with the two shrinkage methods against the sample partial correlation values. Figure 6 (b) shows clearly that  $s_b$  behaves similarly in all three methods which allows then a fair comparison between them. In contrast, Figure 6 (a) shows that there is a clear decrease of  $s_w$  using methods based on shrinkage. This confirms that shrinkage estimators such as  $lw$  and  $t\text{-dist shrink}$  show very good similar performances and improve the reliability of global efficiency on this test-retest dataset.

Similar conclusions hold for the other pairs of sessions. The corresponding ICC and  $p$ -values plots can be found in the Appendix, Figures 9 and 10.

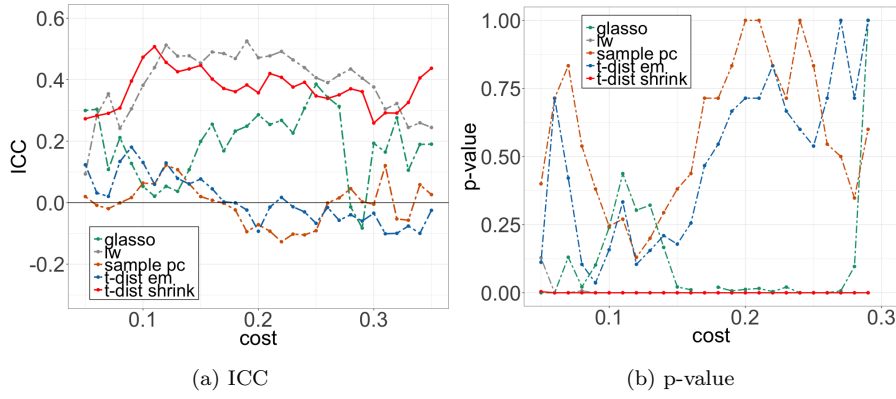


Fig. 5: Intra-class correlation coefficient (ICC) between two fMRI sessions ( $S_{11}$  and  $S_{22}$ ) (a) and associated p-values (b) with respect to the pre-set percentage of edges in the graphs (*cost*). The ICC values are shown for the various estimators considered in this study. The larger the ICC, the higher the consistency between the two sessions and the higher the estimator reliability.

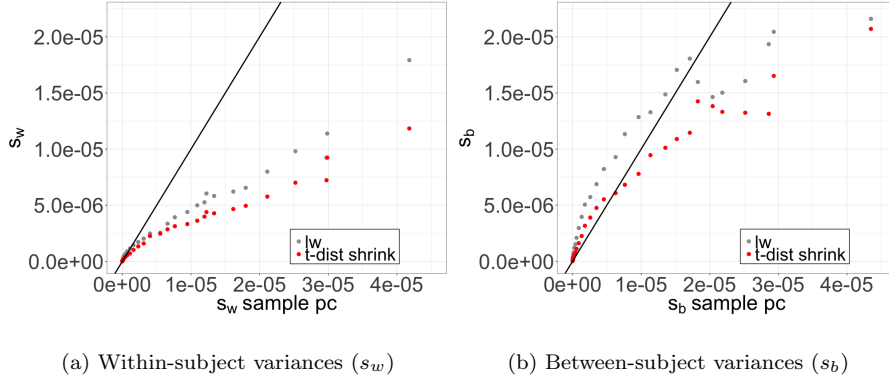


Fig. 6: Within-subject  $s_w$  (a) and between-subject  $s_b$  (b) variances, for sessions  $S_{11}$  and  $S_{22}$ , using different estimators. Values of  $s_w$  and  $s_b$  for two shrinkage methods (*lw* and *t-dist shrink*) are plotted against values obtained using sample partial correlation. The black line indicates the line of equal values. Red and grey dots correspond to varying percentages of edges in the graphs.

## 7 Conclusion

In this paper, we address the issue of robust covariance matrix estimation in settings where the sample size is small compared to the number of parameters and the sample mean is not known a priori. Elliptical distributions are considered to improve the robustness of the approaches. In particular, we focus on Student's *t*-distributions for their ability to model heavy tails and to handle outliers. A regularisation approach based on shrinkage is then used to face the relative lack of

data. These two aspects are combined and lead to a penalized maximum likelihood based estimator assuming the observations follow a multivariate Student's t-distribution. The proposed approach is showed to fulfill a number of theoretical results for more general elliptical distributions, and it has the advantage to be implemented easily in practice.

Among regularized robust estimators, the proposed estimator has several desirable properties: 1) The penalization level or regularizing coefficient is not tuned manually but estimated via a closed-form formula deriving from a minimum mean squared-error principle, 2) prior knowledge on the mean and degree-of-freedom parameter values are not needed and both these parameters can be estimated in a data driven way, at last 3) the efficient algorithm that is derived shows good estimation accuracy when compared to the standard Tyler's estimator on simulated data and to additional standard methods on real data. In particular our experiments confirm the importance of a good estimation of the mean and the potential advantage of methods that aim at estimating both the mean and covariance matrix.

### **Acknowledgment**

Authors acknowledge LabEx PERSYVAL-Lab (ANR-11-LABX-0025-01). Authors also acknowledge the Grenoble Alpes Data Institute that is supported by the French National Research Agency under the "Investissements d'avenir" program (ANR-15-IDEX-02).

## 8 Appendix

### 8.1 Proof of Theorem 1

Following the notations of [6], we define matrix  $\tilde{\mathbf{C}}$  as:

$$\tilde{\mathbf{C}} = m(\boldsymbol{\Sigma}) = \frac{1}{n} \sum_{i=1}^n u \left( (\mathbf{x}_i - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right) (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top.$$

Then  $\tilde{\boldsymbol{\Sigma}}_\rho = (1 - \rho)\tilde{\mathbf{C}} + \rho\mathbf{I}$  and deriving (12) w.r.t.  $\rho$  leads to:

$$\begin{aligned} \rho^* &= \frac{\mathbb{E} \left[ \text{tr} \left\{ (\mathbf{I} - \tilde{\mathbf{C}}) (\boldsymbol{\Sigma} - \tilde{\mathbf{C}}) \right\} \right]}{\mathbb{E} \left[ \left\| \mathbf{I} - \tilde{\mathbf{C}} \right\|_F^2 \right]} \\ &= \frac{\mathbb{E} \left[ \text{tr} (\tilde{\mathbf{C}}^2) \right] - \mathbb{E} \left[ \text{tr} (\tilde{\mathbf{C}}) \right] - \mathbb{E} \left[ \text{tr} (\tilde{\mathbf{C}} \boldsymbol{\Sigma}) \right] + \text{tr} (\boldsymbol{\Sigma})}{\mathbb{E} \left[ \text{tr} (\tilde{\mathbf{C}}^2) \right] - 2\mathbb{E} \left[ \text{tr} (\tilde{\mathbf{C}}) \right] + p}. \end{aligned} \quad (40)$$

Since the vectors  $\mathbf{x}_i = \boldsymbol{\mu} + R_i \boldsymbol{\Lambda} \mathbf{U}_i$  for  $1 \leq i \leq n$  are elliptically distributed and  $\boldsymbol{\Sigma} = \boldsymbol{\Lambda} \boldsymbol{\Lambda}^\top$ , then (see *e.g.* [5]):

$$(\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top = R_i^2 \boldsymbol{\Lambda} \mathbf{U}_i \mathbf{U}_i^\top \boldsymbol{\Lambda}^\top \quad \text{and} \quad (\mathbf{x}_i - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) = R_i^2.$$

Thus,

$$\begin{aligned} \tilde{\mathbf{C}} &= \frac{1}{n} \sum_{i=1}^n u(R_i^2) R_i^2 \boldsymbol{\Lambda} \mathbf{U}_i \mathbf{U}_i^\top \boldsymbol{\Lambda}^\top, \\ \text{and } \text{tr} (\tilde{\mathbf{C}}) &= \frac{1}{n} \sum_{i=1}^n u(R_i^2) R_i^2 \text{tr} (\boldsymbol{\Lambda} \mathbf{U}_i \mathbf{U}_i^\top \boldsymbol{\Lambda}^\top). \end{aligned}$$

Using that  $R_i$  and  $\mathbf{U}_i$  are independent and that  $\mathbb{E} [\mathbf{U}_i \mathbf{U}_i^\top] = \frac{1}{p} \mathbf{I}$ , it comes,

$$\begin{aligned} \mathbb{E} \left[ \text{tr} (\tilde{\mathbf{C}}) \right] &= \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[ u(R_i^2) R_i^2 \right] \text{tr} (\boldsymbol{\Lambda} \mathbb{E} [\mathbf{U}_i \mathbf{U}_i^\top] \boldsymbol{\Lambda}^\top) \\ &= f_1 \text{tr} (\boldsymbol{\Sigma}) / p. \end{aligned} \quad (41)$$

And similarly,  $\mathbb{E} \left[ \text{tr} (\tilde{\mathbf{C}} \boldsymbol{\Sigma}) \right] = f_1 \text{tr} (\boldsymbol{\Sigma}^2) / p$ .

Let us now deal with the term  $\mathbb{E} \left[ \text{tr} (\tilde{\mathbf{C}}^2) \right]$ :

$$\tilde{\mathbf{C}}^2 = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n u(R_i^2) R_i^2 u(R_j^2) R_j^2 \boldsymbol{\Lambda} \mathbf{U}_i \mathbf{U}_i^\top \boldsymbol{\Lambda}^\top \boldsymbol{\Lambda} \mathbf{U}_j \mathbf{U}_j^\top \boldsymbol{\Lambda}^\top,$$

and using the mutual independence of all  $\mathbf{U}_i$ 's and  $R_i$ 's and their common distribution, it comes,

$$\begin{aligned}\mathbb{E} \left[ \text{tr}(\tilde{\mathbf{C}}^2) \right] &= \frac{1}{n^2} \left( \sum_{i \neq j} \mathbb{E} \left[ u(R^2) R^2 \right]^2 \text{tr} \left( (\mathbf{A} \mathbb{E} \left[ \mathbf{U}^\top \mathbf{U} \right] \mathbf{A}^\top)^2 \right) \right. \\ &\quad \left. + \sum_{i=1}^n \mathbb{E} \left[ (u(R^2) R^2)^2 \right] \mathbb{E} \left[ \text{tr} \left( (\mathbf{U}^\top \mathbf{A}^\top \mathbf{A} \mathbf{U})^2 \right) \right] \right) \\ &= \frac{n-1}{n} \frac{f_1^2 \text{tr}(\boldsymbol{\Sigma}^2)}{p^2} + \frac{1}{n} f_2 \mathbb{E} \left[ (\mathbf{U}^\top \mathbf{A}^\top \mathbf{A} \mathbf{U})^2 \right].\end{aligned}$$

Using that  $\mathbf{A}^\top \mathbf{A} = \mathbf{A}$  where  $\mathbf{A}$  is the diagonal matrix containing the eigenvalues of  $\boldsymbol{\Sigma}$  denoted by  $a_1, \dots, a_p$  and denoting the elements of  $\mathbf{U}$  as  $\mathbf{U} = (u_1, \dots, u_p)^\top$ , the last expectation can be computed as,

$$\begin{aligned}\mathbb{E} \left[ (\mathbf{U}^\top \mathbf{A} \mathbf{U})^2 \right] &= \mathbb{E} \left[ \sum_{k=1}^p a_k^2 u_k^4 + \sum_{k \neq m} a_k a_m u_k^2 u_m^2 \right] = \sum_{k=1}^p a_k^2 \mathbb{E} \left[ u_k^4 \right] \\ &\quad + \sum_{k \neq m} a_k a_m \mathbb{E} \left[ u_k^2 u_m^2 \right].\end{aligned}$$

At last, it comes from Theorem 5 in [10] that,  $\mathbb{E} \left[ u_k^4 \right] = \frac{3}{p(p+2)}$  and  $\mathbb{E} \left[ u_k^2 u_m^2 \right] = \frac{1}{p(p+2)}$  for  $k \neq m$ . Then,

$$\begin{aligned}\mathbb{E} \left[ (\mathbf{U}^\top \mathbf{A} \mathbf{U})^2 \right] &= \frac{1}{p(p+2)} \left( 3 \sum_{k=1}^p a_k^2 + \sum_{k \neq m} a_k a_m \right) \\ &= \frac{1}{p(p+2)} \left( 2\text{tr}(\mathbf{A}^2) + (\text{tr}(\mathbf{A}))^2 \right) = \frac{1}{p(p+2)} \left( 2\text{tr}(\boldsymbol{\Sigma}^2) + \text{tr}(\boldsymbol{\Sigma})^2 \right).\end{aligned}$$

We finally get:

$$\begin{aligned}\mathbb{E} \left[ \text{tr}(\tilde{\mathbf{C}}^2) \right] &= \frac{1}{n^2} \left[ n \left( f_2 \frac{1}{p(p+2)} \left( 2\text{tr}(\boldsymbol{\Sigma}^2) + \text{tr}(\boldsymbol{\Sigma})^2 \right) \right) + (n^2 - n) \left( f_1^2 \frac{1}{p^2} \text{tr}(\boldsymbol{\Sigma}^2) \right) \right] \\ &= \text{tr}(\boldsymbol{\Sigma}^2) \left( f_1^2 \frac{n-1}{np^2} + f_2 \frac{2}{np(p+2)} \right) + \text{tr}(\boldsymbol{\Sigma})^2 f_2 \frac{1}{np(p+2)}.\end{aligned}$$

The shrinkage coefficient  $\rho^*$  thus takes the form:

$$\rho^* = \frac{\text{tr}(\boldsymbol{\Sigma}^2) \left( f_1^2 \frac{n-1}{np^2} + f_2 \frac{2}{np(p+2)} - f_1 \frac{1}{p} \right) + \text{tr}(\boldsymbol{\Sigma})^2 f_2 \frac{1}{np(p+2)} + \text{tr}(\boldsymbol{\Sigma}) \left( 1 - \frac{1}{p} f_1 \right)}{\text{tr}(\boldsymbol{\Sigma}^2) \left( f_1^2 \frac{n-1}{np^2} + f_2 \frac{2}{np(p+2)} \right) + \text{tr}(\boldsymbol{\Sigma})^2 f_2 \frac{1}{np(p+2)} - 2 \frac{1}{p} f_1 \text{tr}(\boldsymbol{\Sigma}) + p},$$

and the result is proved. Note also that this expression of  $\rho^*$  is by construction always in  $[0, 1]$ . Indeed looking at expression (40) and using that  $\text{tr}(\boldsymbol{\Sigma}) = p$ , it comes that the difference between the numerator and denominator is equal to  $\frac{f_1}{p} (p - \text{tr}(\boldsymbol{\Sigma}^2))$ , which is always negative because

$$\text{tr}((\boldsymbol{\Sigma} - \mathbf{I})^2) = \text{tr}(\boldsymbol{\Sigma}^2) - 2\text{tr}(\boldsymbol{\Sigma}) + \text{tr}(\mathbf{I}^2) = \text{tr}(\boldsymbol{\Sigma}^2) - 2p + p = \text{tr}(\boldsymbol{\Sigma}^2) - p \geq 0.$$

Similarly, we can check that  $\rho^*$  is always positive.  $\square$

## 8.2 Generalized regularized estimators

Theorem 1 is generalized to a penalty term equal to  $\rho \mathbf{T}$  where  $\mathbf{T}$  is a positive definite matrix. The regularized estimators we consider are generalized to:

$$\boldsymbol{\Sigma}_\rho = (1 - \rho)m(\boldsymbol{\Sigma}) + \rho \mathbf{T}.$$

As previously, parameter  $\rho$  is chosen to minimize the MSE between the "clairvoyant estimator"  $\boldsymbol{\Sigma}_\rho$  and  $\boldsymbol{\Sigma}$ :

$$\mathbb{E} \left[ \left\| \frac{1 - \rho}{n} \sum_{i=1}^n u \left( (\mathbf{x}_i - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right) (\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{x}_i - \boldsymbol{\mu})^\top + \rho \mathbf{T} - \boldsymbol{\Sigma} \right\|_F^2 \right]. \quad (42)$$

**Theorem 2** *If  $\mathbf{x} = \boldsymbol{\mu} + R\Lambda\mathbf{U}$ , the oracle coefficient  $\rho^*$  which minimizes (42) is, under the conditions of Proposition 1 and under  $\text{tr}(\boldsymbol{\Sigma}) = p$ ,*

$$\rho^* = \frac{\text{tr}(\boldsymbol{\Sigma}^2) \left( f_2 \frac{2}{np(p+2)} - \frac{1}{n} \right) + f_2 \frac{p}{n(p+2)}}{\text{tr}(\boldsymbol{\Sigma}^2) \left( 1 - \frac{1}{n} + f_2 \frac{2}{np(p+2)} \right) + f_2 \frac{p}{n(p+2)} - 2\text{tr}(\boldsymbol{\Sigma}\mathbf{T}) + \text{tr}(\mathbf{T}^2)}. \quad (43)$$

*Proof* With  $\tilde{\mathbf{C}}$  defined as in (45), deriving w.r.t.  $\rho$  in (42) leads to:

$$\begin{aligned} \rho^* &= \frac{\mathbb{E} \left[ \text{tr} \left\{ (\boldsymbol{\Sigma} - \tilde{\mathbf{C}}) (\mathbf{T} - \tilde{\mathbf{C}}) \right\} \right]}{\mathbb{E} \left[ \left\| \mathbf{T} - \tilde{\mathbf{C}} \right\|_F^2 \right]} \\ &= \frac{\mathbb{E} \left[ \text{tr}(\tilde{\mathbf{C}}^2) \right] - \mathbb{E} \left[ \text{tr}(\tilde{\mathbf{C}}\mathbf{T}) \right] - \mathbb{E} \left[ \text{tr}(\tilde{\mathbf{C}}\boldsymbol{\Sigma}) \right] + \text{tr}(\boldsymbol{\Sigma}\mathbf{T})}{\mathbb{E} \left[ \text{tr}(\tilde{\mathbf{C}}^2) \right] - 2\mathbb{E} \left[ \text{tr}(\tilde{\mathbf{C}}\mathbf{T}) \right] + \text{tr}(\mathbf{T}^2)}. \end{aligned}$$

The only difference with the case  $\mathbf{T} = \mathbf{I}$  of Theorem 1 is that now:

$$\mathbb{E} \left[ \text{tr}(\tilde{\mathbf{C}}\mathbf{T}) \right] = \text{tr} \left( \mathbb{E} \left[ \tilde{\mathbf{C}}\mathbf{T} \right] \right) = \text{tr} \left( \mathbb{E} \left[ \tilde{\mathbf{C}} \right] \mathbf{T} \right) = \text{tr} \left( \frac{1}{p} f_1 \boldsymbol{\Sigma}\mathbf{T} \right) = \frac{1}{p} f_1 \text{tr}(\boldsymbol{\Sigma}\mathbf{T}).$$

The other quantities have been already computed in Appendix 8.1. It follows,

$$\rho^* = \frac{\text{tr}(\boldsymbol{\Sigma}^2) \left( f_1^2 \frac{n-1}{np^2} + f_2 \frac{2}{np(p+2)} - f_1 \frac{1}{p} \right) + \text{tr}(\boldsymbol{\Sigma})^2 f_2 \frac{1}{np(p+2)} + \text{tr}(\boldsymbol{\Sigma}\mathbf{T}) \left( 1 - \frac{1}{p} f_1 \right)}{\text{tr}(\boldsymbol{\Sigma}^2) \left( f_1^2 \frac{n-1}{np^2} + f_2 \frac{2}{np(p+2)} \right) + \text{tr}(\boldsymbol{\Sigma})^2 f_2 \frac{1}{np(p+2)} - 2\frac{1}{p} f_1 \text{tr}(\boldsymbol{\Sigma}\mathbf{T}) + \text{tr}(\mathbf{T}^2)}.$$

Finally, with the condition  $\text{tr}(\boldsymbol{\Sigma}) = p$  and  $f_1 = p$  from Proposition 1, we obtain (43).



### 8.3 Proof of Proposition 1

Let us denote  $f_R$  the pdf of  $R$ . According to Theorem 2.9 in [8], we have:

$$f_R(r) = \frac{2\pi^{\frac{p}{2}}}{\Gamma(\frac{p}{2})} r^{p-1} g(r^2) \Leftrightarrow g(r) = \frac{\Gamma(\frac{p}{2})}{2\pi^{\frac{p}{2}}} r^{(1-p)/2} f_R(\sqrt{r}).$$

Recalling that  $u(t) = -2g'(t)/g(t)$  we can easily rewrite  $\mathbb{E}[u(R^2)R^2]$  as:

$$\int_0^\infty u(r^2)r^2 f_R(r)dr = -4 \frac{\pi^{\frac{p}{2}}}{\Gamma(\frac{p}{2})} \int_0^\infty r^{p+1} g'(r^2)dr.$$

Using the previous equation, it is straightforward to prove:

$$g'(r^2) = \frac{\Gamma(\frac{p}{2})}{2\pi^{\frac{p}{2}}} \left[ \left( \frac{1-p}{2} \right) r^{-1-p} f_R(r) + \frac{1}{2} r^{-p} f'_R(r) \right],$$

hence, using an integration by part, we get:

$$\begin{aligned} \mathbb{E}[u(R^2)R^2] &= (p-1) \int_0^\infty f_R(r)dr - \int_0^\infty r f'_R(r)dr \\ &= p - \lim_{r \rightarrow +\infty} r f_R(r) + \lim_{r \rightarrow 0} r f_R(r) = p. \quad \square \end{aligned} \quad (44)$$

### 8.4 Proof of Proposition 2

As before we use the following notation:

$$\tilde{\mathbf{C}} = m(\boldsymbol{\Sigma}) = \frac{1}{n} \sum_{i=1}^n u \left( (\mathbf{x}_i - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right) (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top. \quad (45)$$

Taking the derivative in (19) with respect to  $\alpha$  we get

$$\alpha^* = \frac{\text{tr}(\boldsymbol{\Sigma}) - \beta^* \mathbb{E}[\text{tr}(\tilde{\mathbf{C}})]}{p} = \frac{\text{tr}(\boldsymbol{\Sigma})}{p} \left( 1 - \frac{f_1}{p} \beta^* \right) \quad (46)$$

from (41) and Proposition 1. So that, under  $\text{tr}(\boldsymbol{\Sigma}) = p$

$$\alpha^* = \left( 1 - \frac{f_1}{p} \beta^* \right),$$

and when  $f_1 = p$ ,  $\alpha^* = (1 - \beta^*)$ .

Similarly, after derivation of (19), the expression for  $\beta^*$  takes the form:

$$\beta^* = \frac{\mathbb{E}[\text{tr}(\tilde{\mathbf{C}}\boldsymbol{\Sigma})] - \alpha^* \mathbb{E}[\text{tr}(\tilde{\mathbf{C}})]}{\mathbb{E}[\text{tr}(\tilde{\mathbf{C}}^2)]}. \quad (47)$$

□

## 8.5 Proof of Proposition 3

We prove each case separately.

1. If  $R^2$  is a Chi-squared distribution with  $p$  degrees of freedom, then:

$$g(t) = \frac{\Gamma\left(\frac{p}{2}\right)}{2\pi^{\frac{p}{2}}} t^{(1-p)/2} f_R(\sqrt{t}) = (2\pi)^{-p/2} \exp\left(-\frac{t}{2}\right).$$

It thus comes  $u(t) = 1$ , and  $f_2 = \mathbb{E}[\chi_p^4]$ . The moments of the Chi-squared distribution lead to  $f_2 = p(p+2)$ .

2. Some straightforward calculations lead to:

$$g(t) = \frac{\Gamma\left(\frac{p}{2}\right)}{\Gamma\left(\frac{2q+p-2}{2s}\right)} \frac{\lambda^{\frac{2q+p-2}{2s}}}{\pi^{\frac{p}{2}}} t^{q-1} \exp(-\lambda t^s)$$

and therefore

$$u(t) = 2\frac{1-q}{t} + 2\lambda s t^{s-1}.$$

By noticing that  $u(t)^2 t^2 = 4(1-q)^2 + 8\lambda s(1-q)t^s + 4\lambda^2 s^2 t^{2s}$ , the moment  $\mathbb{E}[u(R^2)^2 R^4]$  may thus be rewritten as follows:

$$4(1-q)^2 + 4s(1-q)\mathbb{E}\left[\chi_{\frac{2q+p-2}{s}}^2\right] + s^2\mathbb{E}\left[\left(\chi_{\frac{2q+p-2}{s}}^2\right)^2\right].$$

Using the moments of the Chi-squared distribution concludes the proof.

3. If  $R^2/p$  is a Fisher distribution with  $p$  and  $\nu$  degrees of freedom, it is known that  $\mathbf{x}$  follows a  $p$ -variate t-distribution with  $\nu$  degrees of freedom. In this case, the generator  $g(t)$  is given by  $(\nu\pi)^{-p/2} (1+t/\nu)^{-(p+\nu)/2}$ , and therefore

$$u(t) = (p+\nu)/(t+\nu).$$

. In addition,

$$\begin{aligned} \mathbb{E}\left[u(R^2)^2 R^4\right] &= \int_0^\infty (u(r^2))^2 r^4 f_R(r) dr \\ &= \int_0^\infty 4 \frac{(g'(r^2))^2}{(g(r^2))^2} r^4 \frac{2\pi^{\frac{p}{2}}}{\Gamma\left(\frac{p}{2}\right)} r^{p-1} g(r^2) dr \\ &= \int_0^\infty 8 \frac{(g'(r^2))^2}{g(r^2)} \frac{\pi^{\frac{p}{2}}}{\Gamma\left(\frac{p}{2}\right)} r^{p+3} dr \\ &= 8 \frac{\pi^{\frac{p}{2}}}{\Gamma\left(\frac{p}{2}\right)} \int_0^\infty \frac{(g'(r^2))^2}{g(r^2)} r^{p+3} dr. \end{aligned}$$

The ratio  $(g'(r^2))^2/g(r^2)$  may be rewritten here:

$$\frac{(g'_{r^2}(r^2))^2}{g(r^2)} = \frac{1}{4} \frac{\Gamma\left(\frac{p+\nu}{2}\right)}{\Gamma\left(\frac{\nu}{2}\right) \pi^{\frac{p}{2}}} \nu^{-\frac{p}{2}-2} (\nu+p)^2 \left(1 + \frac{r^2}{\nu}\right)^{-\frac{\nu+p}{2}-2}.$$

Combining the previous relationships obtained, it comes:

$$\mathbb{E} \left[ u(R^2)^2 R^4 \right] = \frac{\Gamma\left(\frac{p+\nu}{2}\right)}{\Gamma\left(\frac{\nu}{2}\right)\Gamma\left(\frac{p}{2}\right)} (\nu+p)^2 \int_0^\infty \left(\frac{r^2}{\nu}\right)^{\frac{p}{2}+1} \left(1 + \frac{r^2}{\nu}\right)^{-\frac{\nu+p}{2}-2} d\left(\frac{r^2}{\nu}\right).$$

Using the known moments of the t-distribution leads to:

$$\mathbb{E} \left[ u(R^2)^2 R^4 \right] = \frac{\Gamma\left(\frac{p+\nu}{2}\right)}{\Gamma\left(\frac{\nu}{2}\right)\Gamma\left(\frac{p}{2}\right)} (\nu+p)^2 \frac{\Gamma\left(\frac{p}{2}+2\right)\Gamma\left(\frac{\nu}{2}\right)}{\Gamma\left(\frac{p+\nu}{2}+2\right)} = \frac{(\nu+p)(p+2)p}{p+\nu+2}.$$

□

### 8.6 Alternative MSE criterion

The authors in [32] propose to estimate the regularization parameter  $\rho$  by minimizing the following alternative MSE criterion:

$$\rho^* = \arg \min_{\rho} E \left[ \left\| \boldsymbol{\Sigma}^{-1} \tilde{\boldsymbol{\Sigma}}_{\rho} - \frac{1}{p} \text{tr} \left( \boldsymbol{\Sigma}^{-1} \tilde{\boldsymbol{\Sigma}}_{\rho} \right) \mathbf{I} \right\|_F^2 \right].$$

Similarly, we provide below the optimal  $\rho$  in this case.

**Theorem 3** *Under the assumption  $\text{tr}(\boldsymbol{\Sigma}^{-1}) = p$ , the oracle estimate  $\rho^*$  is given by:*

$$\rho^* = \frac{f_2(p-2 + p \text{tr}(\boldsymbol{\Sigma}))}{f_2(p-2 + p \text{tr}(\boldsymbol{\Sigma})) + np^2(p+2)(p^{-1} \text{tr}(\boldsymbol{\Sigma}^{-2}) - 1)}.$$

*Proof*

$$\begin{aligned} \rho^* &= \arg \min_{\rho} E \left[ \left\| \boldsymbol{\Sigma}^{-1} \tilde{\boldsymbol{\Sigma}}_{\rho} - \frac{1}{p} \text{tr} \left( \boldsymbol{\Sigma}^{-1} \tilde{\boldsymbol{\Sigma}}_{\rho} \right) \mathbf{I} \right\|_F^2 \right] \\ &= \arg \min_{\rho} \left[ \text{tr} \left( \boldsymbol{\Sigma}^{-2} E \left[ \tilde{\boldsymbol{\Sigma}}_{\rho}^2 \right] \right) - \frac{1}{p} E \left[ \left( \text{tr} \left( \boldsymbol{\Sigma}^{-1} \tilde{\boldsymbol{\Sigma}}_{\rho} \right) \right)^2 \right] \right]. \end{aligned} \quad (48)$$

We start with term  $\text{tr} \left( \boldsymbol{\Sigma}^{-1} \tilde{\boldsymbol{\Sigma}}_{\rho} \right)$ :

$$\begin{aligned} \text{tr} \left( \boldsymbol{\Sigma}^{-1} \tilde{\boldsymbol{\Sigma}}_{\rho} \right) &= \text{tr} \left( \boldsymbol{\Sigma}^{-1} (1-\rho) \tilde{\mathbf{C}} + \rho \boldsymbol{\Sigma}^{-1} \right) \\ &= \text{tr} \left( (1-\rho) \frac{1}{n} \sum_{i=1}^n u \left( (\mathbf{x}_i - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right) \right. \\ &\quad \left. \times \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{x}_i - \boldsymbol{\mu})^\top + \rho \boldsymbol{\Sigma}^{-1} \right) \\ &= (1-\rho) \frac{1}{n} \sum_{i=1}^n u \left( (\mathbf{x}_i - \hat{\boldsymbol{\mu}})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right) \\ &\quad \times \text{tr} \left( (\mathbf{x}_i - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right) + \rho \text{tr}(\boldsymbol{\Sigma}^{-1}) \\ &= (1-\rho) \frac{1}{n} \sum_{i=1}^n u(R_i^2) R_i^2 + \rho \text{tr}(\boldsymbol{\Sigma}^{-1}). \end{aligned}$$

Thus:

$$\begin{aligned}
E \left[ \left( \text{tr} \left( \boldsymbol{\Sigma}^{-1} \tilde{\boldsymbol{\Sigma}}_{\rho} \right) \right)^2 \right] &= (1 - \rho)^2 \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n u(R_i^2) R_i^2 u(R_j^2) R_j^2 \\
&\quad + 2(1 - \rho) \rho \text{tr}(\boldsymbol{\Sigma}^{-1}) \frac{1}{n} \sum_{i=1}^n u(R_i^2) R_i^2 + \rho^2 \left( \text{tr}(\boldsymbol{\Sigma}^{-1}) \right)^2 \\
&= (1 - \rho)^2 \frac{1}{n^2} \left( n f_2 + (n^2 - n) f_1^2 \right) \\
&\quad + 2(1 - \rho) \rho \text{tr}(\boldsymbol{\Sigma}^{-1}) f_1 + \rho^2 \left( \text{tr}(\boldsymbol{\Sigma}^{-1}) \right)^2.
\end{aligned}$$

Then we compute  $\text{tr} \left( \boldsymbol{\Sigma}^{-2} E \left[ \tilde{\boldsymbol{\Sigma}}_{\rho}^2 \right] \right)$ :

$$\begin{aligned}
E \left[ \tilde{\boldsymbol{\Sigma}}_{\rho}^2 \right] &= (1 - \rho)^2 E \left[ \tilde{\mathbf{C}}^2 \right] + 2\rho(1 - \rho) E \left[ \tilde{\mathbf{C}} \right] + \rho^2 \mathbf{I} \\
\text{tr} \left( \boldsymbol{\Sigma}^{-2} E \left[ \tilde{\boldsymbol{\Sigma}}_{\rho}^2 \right] \right) &= (1 - \rho)^2 \text{tr} \left( \boldsymbol{\Sigma}^{-2} E \left[ \tilde{\mathbf{C}}^2 \right] \right) + 2\rho(1 - \rho) \text{tr} \left( \boldsymbol{\Sigma}^{-2} E \left[ \tilde{\mathbf{C}} \right] \right) \\
&\quad + \rho^2 \text{tr} \left( \boldsymbol{\Sigma}^{-2} \right) \\
\text{tr} \left( \boldsymbol{\Sigma}^{-2} E \left[ \tilde{\mathbf{C}} \right] \right) &= \text{tr} \left( \boldsymbol{\Sigma}^{-2} \frac{1}{p} f_1 \boldsymbol{\Sigma} \right) = \text{tr} \left( \frac{1}{p} f_1 \boldsymbol{\Sigma}^{-1} \right) = \frac{1}{p} f_1 \text{tr} \left( \boldsymbol{\Sigma}^{-1} \right) \\
\text{tr} \left( \boldsymbol{\Sigma}^{-2} E \left[ \tilde{\mathbf{C}}^2 \right] \right) &= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n E \left\{ u(R_i^2) R_i^2 u(R_j^2) R_j^2 \right\} \\
&\quad \times \text{tr} \left( \boldsymbol{\Sigma}^{-2} E \left\{ \boldsymbol{\Lambda} \mathbf{U}_i \mathbf{U}_i^{\top} \boldsymbol{\Lambda}^{\top} \boldsymbol{\Lambda} \mathbf{U}_j \mathbf{U}_j^{\top} \boldsymbol{\Lambda}^{\top} \right\} \right).
\end{aligned}$$

We first compute  $\text{tr} \left( \boldsymbol{\Sigma}^{-2} E \left[ \boldsymbol{\Lambda} \mathbf{U}_i \mathbf{U}_i^{\top} \boldsymbol{\Lambda}^{\top} \boldsymbol{\Lambda} \mathbf{U}_j \mathbf{U}_j^{\top} \boldsymbol{\Lambda}^{\top} \right] \right)$ :

$$\begin{aligned}
&\text{tr} \left( \boldsymbol{\Sigma}^{-2} E \left[ \boldsymbol{\Lambda} \mathbf{U}_i \mathbf{U}_i^{\top} \boldsymbol{\Lambda}^{\top} \boldsymbol{\Lambda} \mathbf{U}_j \mathbf{U}_j^{\top} \boldsymbol{\Lambda}^{\top} \right] \right) \\
&= \text{tr} \left( E \left[ \left( \boldsymbol{\Lambda}^{\top} \right)^{-1} \boldsymbol{\Lambda}^{-1} \left( \boldsymbol{\Lambda}^{\top} \right)^{-1} \boldsymbol{\Lambda}^{-1} \boldsymbol{\Lambda} \mathbf{U}_i \mathbf{U}_i^{\top} \boldsymbol{\Lambda}^{\top} \boldsymbol{\Lambda} \mathbf{U}_j \mathbf{U}_j^{\top} \boldsymbol{\Lambda}^{\top} \right] \right) \\
&= E \left[ \text{tr} \left\{ \mathbf{A}^{-1} \mathbf{U}_i \mathbf{U}_i^{\top} \mathbf{A} \mathbf{U}_j \mathbf{U}_j^{\top} \right\} \right].
\end{aligned}$$

Now for the term  $E \left[ \text{tr} \left\{ \mathbf{A}^{-1} \mathbf{U}_i \mathbf{U}_i^{\top} \mathbf{A} \mathbf{U}_j \mathbf{U}_j^{\top} \right\} \right]$ :

$$E \left[ \text{tr} \left\{ \mathbf{A}^{-1} \mathbf{U}_i \mathbf{U}_i^{\top} \mathbf{A} \mathbf{U}_j \mathbf{U}_j^{\top} \right\} \right] = \begin{cases} \frac{2p + \text{tr}(\boldsymbol{\Sigma}) \text{tr}(\boldsymbol{\Sigma}^{-1})}{p(p+2)}, & i = j; \\ \text{tr} \left( \mathbf{A}^{-1} \frac{1}{p} \mathbf{A} \frac{1}{p} \right) = \text{tr} \left( \frac{1}{p^2} \mathbf{I} \right) = \frac{1}{p}, & i \neq j. \end{cases}$$

Thus:

$$\text{tr} \left( \boldsymbol{\Sigma}^{-2} E \left[ \tilde{\mathbf{C}}^2 \right] \right) = \frac{1}{n^2} \left[ n f_2 \frac{2p + \text{tr}(\boldsymbol{\Sigma}) \text{tr}(\boldsymbol{\Sigma}^{-1})}{p(p+2)} + (n^2 - n) f_1^2 \frac{1}{p} \right].$$

So that the term to minimize in (48) is:

$$\begin{aligned}
& tr \left( \boldsymbol{\Sigma}^{-2} E \left[ \tilde{\boldsymbol{\Sigma}}_{\rho}^2 \right] \right) - \frac{1}{p} E \left[ \left( tr \left( \boldsymbol{\Sigma}^{-1} \tilde{\boldsymbol{\Sigma}}_{\rho} \right) \right)^2 \right] \\
&= (1 - \rho)^2 tr \left( \boldsymbol{\Sigma}^{-2} E \left[ \tilde{\mathbf{C}}^2 \right] \right) + 2\rho(1 - \rho) tr \left( \boldsymbol{\Sigma}^{-2} E \left[ \tilde{\mathbf{C}} \right] \right) + \rho^2 tr \left( \boldsymbol{\Sigma}^{-2} \right) \\
&- (1 - \rho)^2 \frac{1}{n^2} \frac{1}{p} \left( n f_2 + (n^2 - n) f_1^2 \right) - 2(1 - \rho) \rho \frac{1}{p} tr \left( \boldsymbol{\Sigma}^{-1} \right) f_1 - \rho^2 \frac{1}{p} \left( tr \left( \boldsymbol{\Sigma}^{-1} \right) \right)^2 \\
&= (1 - \rho)^2 \left[ tr \left( \boldsymbol{\Sigma}^{-2} E \left[ \tilde{\mathbf{C}}^2 \right] \right) - \frac{1}{n^2} \frac{1}{p} \left( n f_2 + (n^2 - n) f_1^2 \right) \right] \\
&+ 2(1 - \rho) \rho \left[ tr \left( \boldsymbol{\Sigma}^{-2} E \left[ \tilde{\mathbf{C}} \right] \right) - \frac{1}{p} tr \left( \boldsymbol{\Sigma}^{-1} \right) f_1 \right] \\
&+ \rho^2 \left[ tr \left( \boldsymbol{\Sigma}^{-2} \right) - \frac{1}{p} \left( tr \left( \boldsymbol{\Sigma}^{-1} \right) \right)^2 \right]. \tag{49}
\end{aligned}$$

Let us denote:

$$\begin{aligned}
m_1 &= tr \left( \boldsymbol{\Sigma}^{-2} E \left[ \tilde{\mathbf{C}}^2 \right] \right) - \frac{1}{n^2} \frac{1}{p} \left( n f_2 + (n^2 - n) f_1^2 \right), \\
m_2 &= tr \left( \boldsymbol{\Sigma}^{-2} E \left[ \tilde{\mathbf{C}} \right] \right) - \frac{1}{p} tr \left( \boldsymbol{\Sigma}^{-1} \right) f_1, \\
m_3 &= tr \left( \boldsymbol{\Sigma}^{-2} \right) - \frac{1}{p} \left( tr \left( \boldsymbol{\Sigma}^{-1} \right) \right)^2. \tag{50}
\end{aligned}$$

Taking the derivative w.r.t.  $\rho$  in (49) we get:

$$\rho^* = \frac{m_1 - m_2}{m_1 - 2m_2 + m_3}.$$

With  $m_1 = \frac{1}{n} \frac{1}{p} f_2 \frac{p-2+tr(\boldsymbol{\Sigma})tr(\boldsymbol{\Sigma}^{-1})}{p+2}$  and  $m_2 = 0$ , we get the final formula:

$$\rho^* = \frac{f_2 (p - 2 + tr(\boldsymbol{\Sigma})tr(\boldsymbol{\Sigma}^{-1}))}{f_2 (p - 2 + tr(\boldsymbol{\Sigma})tr(\boldsymbol{\Sigma}^{-1})) + np(p + 2)tr(\boldsymbol{\Sigma}^{-2}) - n(p + 2)(tr(\boldsymbol{\Sigma}^{-1}))^2},$$

which under the assumption  $tr(\boldsymbol{\Sigma}^{-1}) = p$  gives:

$$\rho^* = \frac{f_2 (p - 2 + ptr(\boldsymbol{\Sigma}))}{f_2 (p - 2 + ptr(\boldsymbol{\Sigma})) + np^2(p + 2)(p^{-1}tr(\boldsymbol{\Sigma}^{-2}) - 1)}.$$

## 8.7 Supplementary plots for simulated data

The following Figures show the results of the simulated data study for all pairs of tested parameters  $(r, \nu)$ .

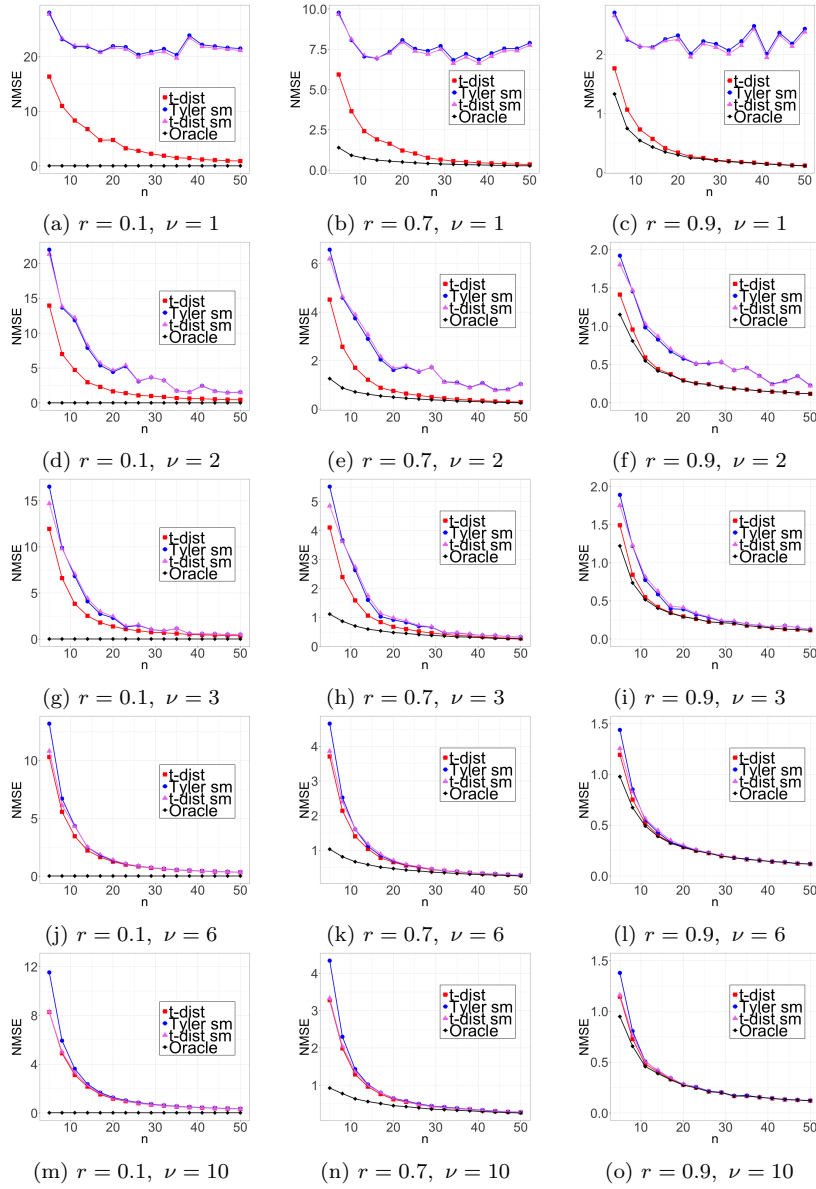


Fig. 7: Multivariate  $t$ -distribution with AR( $r$ ) covariance structure ( $p = 50$ ,  $r \in \{0.1, 0.5, 0.9\}$ ,  $\nu \in \{1, 2, 3, 6, 10\}$  and  $\boldsymbol{\mu}$  is set to a vector of 5): Normalized mean squared-errors for  $\boldsymbol{\Sigma}$  computed over 100 simulated samples of  $n$  observations each with  $n$  varying from 5 to 50.

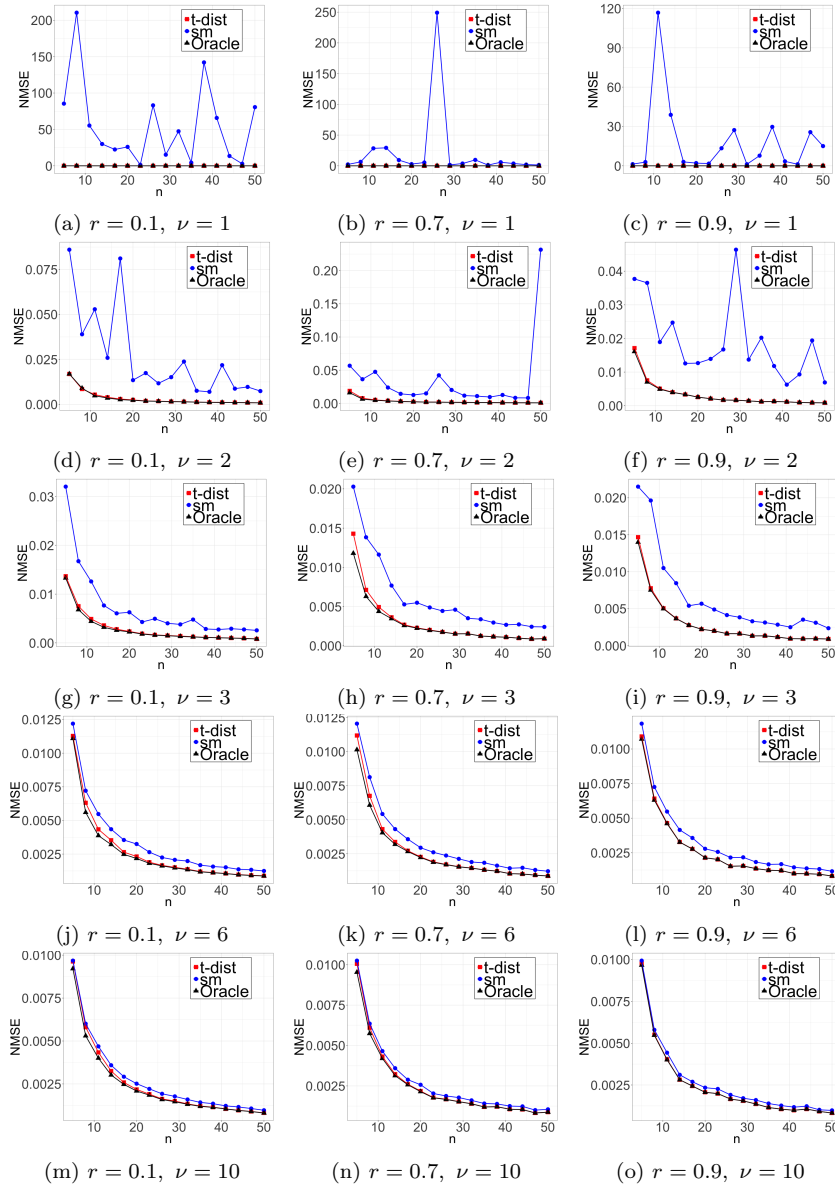


Fig. 8: Multivariate  $t$ -distribution with AR( $r$ ) covariance structure ( $p = 50$ ,  $r \in \{0.1, 0.5, 0.9\}$ ,  $\nu \in \{1, 2, 3, 6, 10\}$  and  $\boldsymbol{\mu}$  is set to a vector of 5): Normalized mean squared-errors for  $\boldsymbol{\mu}$  computed over 100 simulated samples of  $n$  observations each with  $n$  varying from 5 to 50.

## 8.8 Supplementary plots for real data

The following Figures 9 and 10 show the ICC values for the tested methods on different pairs of sessions and their associated  $p$ -values.

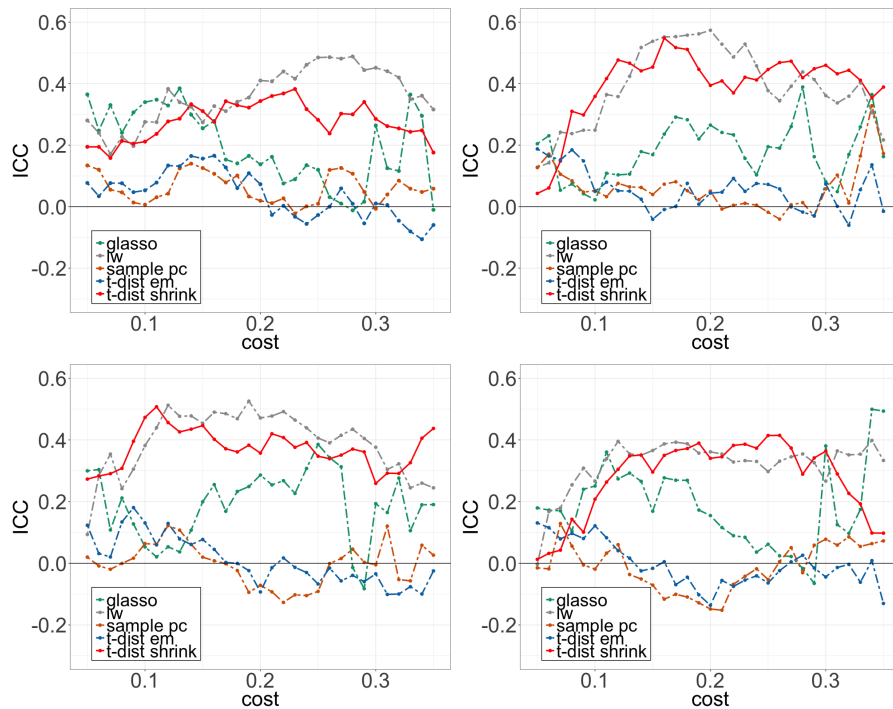


Fig. 9: Intraclass correlation coefficient (ICC) between pairs of fMRI sessions with respect to the pre-set percentage of edges in the graphs (cost). The ICC values are shown for the various estimators considered in this study. Pairs of sessions are from upper left to bottom right  $(S_{11}, S_{21})$ ,  $(S_{12}, S_{22})$ ,  $(S_{11}, S_{22})$  and  $(S_{12}, S_{21})$ .

## References

1. Arslan, O., Genç, A.I.: A generalization of the multivariate slash distribution. *Journal of Statistical Planning and Inference* **139**(3), 1164–1170 (2009)
2. Ashurbekova, K., Achard, S., Forbes, F.: Robust structure learning using multivariate T-distributions. In: 50e Journées de la Statistique (JdS2018), Saclay, France (2018)
3. Banerjee, O., El Ghaoui, L., dAspremont, A.: Model selection through sparse maximum likelihood estimation for multivariate Gaussian or binary data. *Journal of Machine learning research* **9**(Mar), 485–516 (2008)
4. Bodnar, T., Schmid, W.: A test for the weights of the global minimum variance portfolio in an elliptical model. *Metrika* **67**(2), 127 (2008)
5. Cambanis, S., Huang, S., Simons, G.: On the theory of elliptically contoured distributions. *Journal of Multivariate Analysis* **11**, 368–385 (1981)



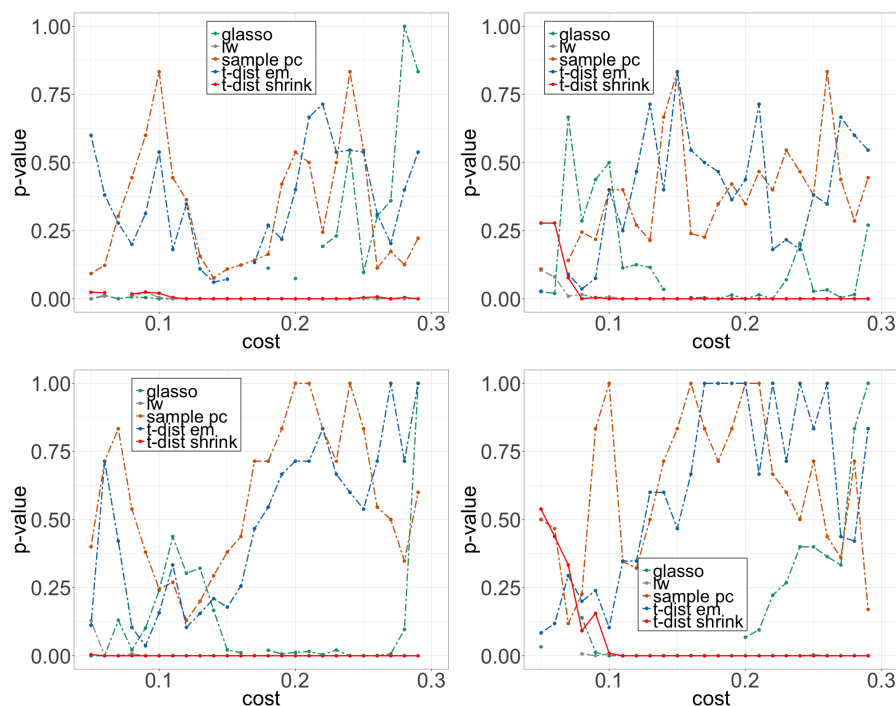


Fig. 10: Associated p-values for ICC with respect to the pre-set percentage of edges in the graphs (cost), for the various estimators considered in this study and different pairs of sessions. Pairs of sessions are from upper left to bottom right  $(S_{11}, S_{21})$ ,  $(S_{12}, S_{22})$ ,  $(S_{11}, S_{22})$  and  $(S_{12}, S_{21})$ .

6. Chen, Y., Wiesel, A., Hero, A.O.: Robust Shrinkage Estimation of High-Dimensional Covariance Matrices. *IEEE Transactions on Signal Processing* **59**(9), 4097–4107 (2011)
7. Eltoft, T., Kim, T., Lee, T.W.: On the multivariate Laplace distribution. *IEEE Signal Processing Letters* **13**(5), 300–303 (2006)
8. Fang, K.T., Kotz, S., Ng, K.W.: Symmetric multivariate and related distributions. Chapman and Hall (1990)
9. Finegold, M., Drton, M.: Robust graphical modeling of gene networks using classical and alternative t-distributions. *The Annals of Applied Statistics* **5**(2A), 1057–1080 (2011)
10. Frahm, G.: Generalized elliptical distributions: Theory and applications. Ph.D. thesis, Universität zu Köln (2004)
11. Friedman, J., Hastie, T., Tibshirani, R.: Sparse inverse covariance estimation with the graphical lasso. *Biostatistics* **9**(3), 432–441 (2008)
12. Friedman, J.H.: Regularized discriminant analysis. *Journal of the American statistical association* **84**(405), 165–175 (1989)
13. Giraud, C.: Estimation of Gaussian graphs by model selection. *Electron. J. Statist.* **2**, 542–563 (2008)
14. Gómez-Sánchez-Manzano, E., Gómez-Villegas, M., Marín, J.: Sequences of elliptical distributions and mixtures of normal distributions. *Journal of multivariate analysis* **97**(2), 295–310 (2006)
15. de Haan, L., Ferreira, A.: Extreme value theory: an introduction. Springer Science & Business Media (2006)
16. Hashorva, E.: On the regular variation of elliptical random vectors. *Statistics & probability letters* **76**(14), 1427–1434 (2006)

17. Hill, B.M.: A simple general approach to inference about the tail of a distribution. *The Annals of Statistics* **3** (1975)
18. Huber, P.J.: Robust estimation of a location parameter. *Ann. Math. Statist.* **35**(1), 73–101 (1964)
19. Jankova, J., van de Geer, S.: Inference in high-dimensional graphical models (2018). URL <https://arxiv.org/abs/1801.08512>
20. Kelker, D.: Distribution theory of spherical distributions and a location-scale parameter generalization. *Sankhya: The Indian Journal of Statistics, Series A* **32**(4), 419–430 (1970)
21. Ledoit, O., Wolf, M.: Improved estimation of the covariance matrix of stock returns with an application to portfolio selection. *Journal of empirical finance* **10**(5), 603–621 (2003)
22. Ledoit, O., Wolf, M.: A well-conditioned estimator for large-dimensional covariance matrices. *Journal of Multivariate Analysis* **88**(2), 365–411 (2004)
23. Liu, C., Rubin, D.B.: ML estimation of the t distribution using EM and its extensions, ECM and ECME. *Statistica Sinica* pp. 19–39 (1995)
24. Liu, J., Palomar, D.: Regularized robust estimation of mean and covariance matrix for incomplete data. *Signal Processing* **165**, 278 – 291 (2019)
25. Malik, H.J., Abraham, B., et al.: Multivariate logistic distributions. *The Annals of Statistics* **1**(3), 588–590 (1973)
26. Maronna, R.A.: Robust M-Estimators of Multivariate Location and Scatter. *The Annals of Statistics* **4**(1), 51–67 (1976)
27. Meinshausen, N., Bühlmann, P.: High-dimensional graphs and variable selection with the lasso. *The Annals of Statistics* **34**(3), 1436–1462 (2006)
28. Muirhead, R.J.: *Aspects of multivariate statistical theory*, vol. 197. John Wiley & Sons (2009)
29. Nadarajah, S.: The Kotz-type distribution with applications. *Statistics: A Journal of Theoretical and Applied Statistics* **37**(4), 341–358 (2003)
30. Ollila, E.: Optimal high-dimensional shrinkage covariance estimation for elliptical distributions. 2017 25th European Signal Processing Conference (EUSIPCO) pp. 1639–1643 (2017)
31. Ollila, E., Oja, H., Koivunen, V.: Complex-valued ICA based on a pair of generalized covariance matrices. *Computational Statistics & Data Analysis* **52**, 3789–3805 (2008)
32. Ollila, E., Tyler, D.E.: Regularized  $M$ -estimators of scatter matrix. *IEEE Transactions on Signal Processing* **62**(22), 6059–6070 (2014)
33. Owen, J., Rabinovitch, R.: On the class of elliptical distributions and their applications to the theory of portfolio choice. *The Journal of Finance* **38**(3), 745–752 (1983)
34. Pascal, F., Chitour, Y., Ovarlez, J., Forster, P., Larzabal, P.: Covariance Structure Maximum-Likelihood Estimates in Compound Gaussian Noise: Existence and Algorithm Analysis. *IEEE Transactions on Signal Processing* **56**(1), 34–48 (2008)
35. Rubinov, M., Sporns, O.: Complex network measures of brain connectivity: uses and interpretations. *NeuroImage* **52**(3), 1059–1069 (2010)
36. Sun, Y., Babu, P., Palomar, D.P.: Regularized robust estimation of mean and covariance matrix under heavy-tailed distributions. *IEEE Transactions on Signal Processing* **63**(12), 3096–3109 (2015)
37. Termenon, M., Delon-Martin, C., Jaillard, A., Achard, S.: Reliability of graph analysis of resting state fMRI using test-retest dataset from the human connectome project. *Neuroimage* **142**(15), 172–187 (2016)
38. Tong, J., Hu, R., Xi, J., Xiao, Z., Guo, Q., Yu, Y.: Linear shrinkage estimation of covariance matrices using low-complexity cross-validation. *Signal Process.* **148**(C), 223–233 (2018)
39. Tyler, D.E.: A Distribution-free  $M$ -estimator of Multivariate Scatter. *The Annals of Statistics* **15**(1), 234–251 (1987)
40. Tyler, D.E.: Statistical analysis for the angular central Gaussian distribution on the sphere. *Biometrika* **74**(3), 579–589 (1987)
41. Tzourio-Mazoyer, N., Landeau, B., Papathanassiou, D., Crivello, F., Etard, O., Delcroix, N., Mazoyer, B., Joliot, M.: Automated anatomical labeling of activations in SPM using a macroscopic anatomical parcellation of the MNI MRI single-subject brain. *NeuroImage* **15**(1), 273–289 (2002)
42. Usseglio-Carleve, A.: Estimation of conditional extreme risk measures from heavy-tailed elliptical random vectors. *Electronic Journal of Statistics* **12**, 4057–4093 (2018)
43. Wiesel, A.: Geodesic convexity and covariance estimation. *IEEE Transactions on Signal Processing* **60**, 6182–6189 (2012)

- 
44. Yuan, M., Lin, Y.: Model selection and estimation in the Gaussian graphical model. *Biometrika* **94**, 19–35 (2007)
  45. Zhang, T., Wiesel, A.: Automatic diagonal loading for Tyler’s robust covariance estimator. In: 2016 IEEE Statistical Signal Processing Workshop (SSP), pp. 1–5. IEEE (2016)
  46. Zhao, Y.: Covariance matrices of quadratic forms in elliptical distributions. *Statistics & Probability Letters* **21**(2), 131–140 (1994)