

Sparsity-based Cholesky Factorization and Its Application to Hyperspectral Anomaly Detection

Ahmad W. Bitar*, Jean-Philippe Ovarlez*† and Loong-Fah Cheong‡

*SONDRA/CentraleSupélec, Plateau du Moulon, 3 rue Joliot-Curie, F-91190 Gif-sur-Yvette, France

†ONERA, DEMR/TSI, Chemin de la Hunière, 91120 Palaiseau, France

‡National University of Singapore (NUS), Singapore, Singapore

Abstract—Estimating large covariance matrices has been a longstanding important problem in many applications and has attracted increased attention over several decades. This paper deals with two methods based on pre-existing works to impose sparsity on the covariance matrix via its unit lower triangular matrix (aka Cholesky factor) \mathbf{T} . The first method serves to estimate the entries of \mathbf{T} using the Ordinary Least Squares (OLS), then imposes sparsity by exploiting some generalized thresholding techniques such as Soft and Smoothly Clipped Absolute Deviation (SCAD). The second method directly estimates a sparse version of \mathbf{T} by penalizing the negative normal log-likelihood with L_1 and SCAD penalty functions. The resulting covariance estimators are always guaranteed to be positive definite. Some Monte-Carlo simulations as well as experimental data demonstrate the effectiveness of our estimators for hyperspectral anomaly detection using the Kelly anomaly detector.

Keywords—Hyperspectral anomaly detection, covariance matrix, sparsity, Cholesky factor.

I. INTRODUCTION

An airborne hyperspectral imaging sensor is capable of simultaneously acquiring the same spatial scene in a contiguous and multiple narrow spectral wavelength (color) bands [1], [2], [3]. When all the spectral bands are stacked together, the resulting hyperspectral image (HSI) is a three dimensional data cube; each pixel in the HSI is a p -dimensional vector, $\mathbf{x} = [x_1, \dots, x_p]^T \in \mathbb{R}^p$, where p designates the total number of spectral bands. With the rich information afforded by the high spectral dimensionality, hyperspectral imagery has found many applications in various fields such as agriculture [4], [5], mineralogy [6], military [7], [8], [9], and in particular, target detection [1], [2], [10], [11], [7], [12], [13]. In many situations of practical interest, we do not have sufficient a priori information to specify the statistics of the target class. More precisely, the target's spectra is not provided to the user. This unknown target is referred as « anomaly » [14] having a very different spectra from the background (e.g., a ship at sea).

Different Gaussian-based anomaly detectors have been proposed in the literature [15], [16], [17], [18], [19]. The detection performance of these detectors mainly depend on the true unknown covariance matrix (of the background surrounding the test pixel) whose entries have to be carefully estimated specially in large dimensions. Due to the fact that in hyperspectral imagery, the number of covariance matrix parameters to estimate grows with the square of the spectral dimension, it becomes impractical to use traditional covariance estimators where the target detection performance can deteriorate significantly. Many a time, the researchers assume that

compounding the large dimensionality problem can be alleviated by leveraging on the assumption that the true unknown covariance matrix is sparse, namely, many entries are zero.

This paper outlines two simple methods based on pre-existing works in order to impose sparsity on the covariance matrix via its Cholesky factor \mathbf{T} . The first method imposes sparsity by exploiting thresholding operators such as Soft and SCAD on the OLS estimate of \mathbf{T} . The second method directly estimates a sparse version of \mathbf{T} by penalizing the negative normal log-likelihood with L_1 and SCAD penalty functions.

Summary of Main Notations: Throughout this paper, we depict vectors in lowercase boldface letters and matrices in uppercase boldface letters. The notation $(\cdot)^T$ stands for the transpose, while $|\cdot|$, $(\cdot)^{-1}$, $(\cdot)'$, $\det(\cdot)$, and $\mathbb{1}$ are the absolute value, the inverse, the derivative, the determinant, and indicator function, respectively. For any $z \in \mathbb{R}$, we define $\text{sign}(z) = 1$ if $z > 0$, $\text{sign}(z) = 0$ if $z = 0$ and $\text{sign}(z) = -1$ if $z < 0$.

II. BACKGROUND AND SYSTEM OVERVIEW

Suppose that we observe a sample of n independent and identically distributed p -random vectors, $\{\mathbf{x}_i\}_{i \in [1, n]}$, each follows a multivariate Gaussian distribution with zero mean and unknown covariance matrix $\Sigma = [\sigma_{g,l}]_{p \times p}$. The first traditional estimator we consider in this paper is the Sample Covariance Matrix (SCM), defined as $\hat{\Sigma}_{SCM} = [\hat{\sigma}_{g,l}]_{p \times p} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T$. In order to address the positivity definiteness constraint problem of $\hat{\Sigma}_{SCM}$, Pourahmadi [20] has modeled the covariance matrices via linear regressions. This is done by letting $\hat{\mathbf{x}} = [\hat{x}_1, \dots, \hat{x}_p]^T \in \mathbb{R}^p$, and consider each element \hat{x}_t , $t \in [1, p]$, as the linear least squares predictor of x_t based on its $t-1$ predecessors $\{x_j\}_{j \in [1, t-1]}$. In particular, for $t \in [1, p]$, let

$$\hat{x}_t = \sum_{j=1}^{t-1} C_{t,j} x_j, \quad \mathbf{T} \Sigma \mathbf{T}^T = \mathbf{D}. \quad (1)$$

where \mathbf{T} is a unit lower triangular matrix with $-C_{t,j}$ in the (t, j) th position for $t \in [2, p]$ and $j \in [1, t-1]$, and \mathbf{D} is a diagonal matrix with $\theta_t^2 = \text{var}(\epsilon_t)$ as its diagonal entries, where $\epsilon_t = x_t - \hat{x}_t$ is the prediction error for $t \in [1, p]$. Note that for $t = 1$, let $\hat{x}_1 = E(x_1) = 0$, and hence, $\text{var}(\epsilon_1) = \theta_1^2 = E[(x_1)^2]$. Given a sample $\{\mathbf{x}_i\}_{i \in [1, n]}$, with $n > p$, a natural estimate of \mathbf{T} and \mathbf{D} , denoted as $\hat{\mathbf{T}}_{OLS}$ and $\hat{\mathbf{D}}_{OLS}$ in this paper, is simply done by plugging in the OLS estimates of the regression coefficients and residual variances

in (1), respectively. In this paper, we shall denote the second traditional estimator by $\hat{\Sigma}_{OLS} = \hat{\mathbf{T}}_{OLS}^{-1} \hat{\mathbf{D}}_{OLS} \hat{\mathbf{T}}_{OLS}^{-T}$.

Obviously, when the spectral dimension p is considered large compared to the number of observed data n , both $\hat{\Sigma}_{SCM}$ and $\hat{\Sigma}_{OLS}$ face difficulties in estimating Σ without an extreme amount of errors. Realizing the challenges brought by the high dimensionality, researchers have thus circumvent these challenges by proposing various regularization techniques to consistently estimate Σ based on the assumption that the covariance matrix is sparse. Recently, Bickel et al. [21] proposed a banded version of $\hat{\Sigma}_{SCM}$, denoted as $B_m(\hat{\Sigma}_{SCM})$ in this paper, with $B_m(\hat{\Sigma}_{SCM}) = [\hat{\sigma}_{g,l} \mathbb{1}(|g - l| \leq m)]$, where $0 \leq m < p$ is the banding parameter. However, this kind of regularization does not always guarantee positive definiteness. In [22], a class of generalized thresholding operators applied on the off-diagonal entries of $\hat{\Sigma}_{SCM}$ have been discussed. These operators combine shrinkage with thresholding and have the advantage to estimate the true zeros as zeros with high probability. These operators (e.g., Soft and SCAD), though simple, do not always guarantee positive definiteness of the thresholded version of $\hat{\Sigma}_{SCM}$. In [23], the covariance matrix is constrained to have an eigen decomposition which can be represented as a sparse matrix transform (SMT) that decomposes the eigen-decomposition into a product of very sparse transformations. The resulting estimator, denoted as $\hat{\Sigma}_{SMT}$ in this paper, is always guaranteed to be positive definite.

In addition to the above review, some other works have attempted to enforce sparsity on the covariance matrix via its Cholesky factor \mathbf{T} . Hence, yielding sparse covariance estimators that are always guaranteed to be positive definite. For example, in [24], Pourahmadi et al. proposed to smooth the first few subdiagonals of $\hat{\mathbf{T}}_{OLS}$ and set to zero the remaining subdiagonals. In [25], Huang et al. proposed to directly estimate a sparse version of \mathbf{T} by penalizing the negative normal log-likelihood with a L_1 -norm penalty function. Hence, allowing the zeros to be irregularly placed in the Cholesky factor. This seems to be an advantage over the work in [24].

We put forth two simple methods for imposing sparsity on the covariance matrix via its Cholesky factor \mathbf{T} . The first method is based on the work in [22], but attempts to render $\hat{\Sigma}_{OLS}$ sparse by thresholding its Cholesky factor $\hat{\mathbf{T}}_{OLS}$ using operators such as Soft and SCAD. The second method aims to generalize the work in [25] in order to be used for various penalty functions. The two methods allow the zeros to be irregularly placed in the Cholesky factor.

Clearly, in real world hyperspectral imagery, the true covariance model is not known, and hence, there is no prior information on its degree of sparsity. However, enforcing sparsity on the covariance matrix seems to be a strong assumption, but can be critically important if the true covariance model for a given HSI is indeed sparse. That is, taking advantage of the possible sparsity in the estimation can potentially improve the target detection performance, as can be seen from the experimental results later. On the other hand, while the true covariance model may not be sparse (or not highly sparse), there should be no worse detection results than to those of the traditional estimators $\hat{\Sigma}_{SCM}$ and $\hat{\Sigma}_{OLS}$.

We evaluate our estimators for hyperspectral anomaly detection using the Kelly anomaly detector [26]. More precisely, we first perform a thorough evaluation of our estimators

on some Monte-Carlo simulations for three true covariance models of different sparsity levels. From our experiments in Subsection IV-A, the detection results show that in truly sparse models, our estimators improve the detection significantly with respect to the traditional ones, and have competitive results with state-of-the-art [21], [22], [23]. When the true model is not sparse, we find that empirically our estimators still have no worse detection results than to those of $\hat{\Sigma}_{SCM}$ and $\hat{\Sigma}_{OLS}$. Next, in Subsection IV-B, our estimators are evaluated on experimental data where a good target detection performances are obtained comparing to the traditional estimators and state-of-the-art. In all the experiments later, we observe that $\hat{\Sigma}_{OLS}$ achieves higher target detection results than to those of $\hat{\Sigma}_{SCM}$.

III. MAIN CONTRIBUTIONS

Before describing the two methods, we want to recall the definition for $\hat{\Sigma}_{OLS}$. Given a sample $\{\mathbf{x}_i\}_{i \in [1, n]}$, we have:

$$x_{i,t} = \sum_{j=1}^{t-1} C_{t,j} x_{i,j} + \epsilon_{i,t}, \quad t \in [2, p], \quad i \in [1, n]. \quad (2)$$

By writing (2) in vector-matrix form for any $t \in [2, p]$, one obtains the simple linear regression model:

$$\mathbf{y}_t = \mathbf{A}_{n,t} \boldsymbol{\beta}_t + \mathbf{e}_t, \quad (3)$$

where $\mathbf{y}_t = [x_{1,t}, \dots, x_{n,t}]^T \in \mathbb{R}^n$, $\mathbf{A}_{n,t} = [x_{i,j}]_{n \times (t-1)}$, $\boldsymbol{\beta}_t = [C_{t,1}, \dots, C_{t,t-1}]^T \in \mathbb{R}^{(t-1)}$, and $\mathbf{e}_t = [\epsilon_{1,t}, \dots, \epsilon_{n,t}]^T \in \mathbb{R}^n$.

When $n > p$, the OLS estimate of $\boldsymbol{\beta}_t$ and the corresponding residual variance are plugged in \mathbf{T} and \mathbf{D} for each $t \in [2, p]$, respectively. At the end, one obtains the estimator $\hat{\Sigma}_{OLS} = \hat{\mathbf{T}}_{OLS}^{-1} \hat{\mathbf{D}}_{OLS} \hat{\mathbf{T}}_{OLS}^{-T}$. Note that $\hat{\mathbf{T}}_{OLS}$ has $-\hat{C}_{t,j}^{OLS}$ in the (t, j) th position for $t \in [2, p]$ and $j \in [1, t-1]$.

A. Generalized thresholding based Cholesky Factor

For any $0 \leq \lambda \leq 1$, we define a matrix thresholding operator $Th(\cdot)$ and denote by $Th(\hat{\mathbf{T}}_{OLS}) = [Th(-\hat{C}_{t,j}^{OLS})]_{p \times p}$ the matrix resulting from applying a specific thresholding operator $Th(\cdot) \in \{\text{Soft, SCAD}\}$ to each element of the matrix $\hat{\mathbf{T}}_{OLS}$ for $t \in [2, p]$ and $j \in [1, t-1]$.

We consider the following minimization problem:

$$Th(\hat{\mathbf{T}}_{OLS}) = \underset{\mathbf{T}}{\operatorname{argmin}} \sum_{t=2}^p \sum_{j=1}^{t-1} \left\{ \frac{1}{2} (\hat{C}_{t,j}^{OLS} - C_{t,j})^2 + p_\lambda \{|C_{t,j}|\} \right\} \quad (4)$$

where $p_\lambda \in \{p_\lambda^{L_1}, p_{\lambda,a>2}^{SCAD}\}$. We have $p_\lambda^{L_1}(|C_{t,j}|) = \lambda |C_{t,j}|$, and $p_{\lambda,a>2}^{SCAD}(|C_{t,j}|) = \begin{cases} \lambda |C_{t,j}| & \text{if } |C_{t,j}| \leq \lambda \\ \frac{-|C_{t,j}|^2 - 2a\lambda|C_{t,j}| + \lambda^2}{2(a-1)} & \text{if } \lambda < |C_{t,j}| \leq a\lambda \\ \frac{(a+1)\lambda^2}{2} & \text{if } |C_{t,j}| > a\lambda \end{cases}$.

Solving (4) with $p_\lambda^{L_1}$ and $p_{\lambda,a>2}^{SCAD}$, yields a closed-form Soft and SCAD thresholding rules, respectively [22], [27]. The value $a = 3.7$ was recommended by Fan and Li [27]. Despite the application here is different than in [27], for simplicity, we use the same value throughout the paper.

We shall designate the two thresholded matrices by $\hat{\mathbf{T}}_{Soft}$ and $\hat{\mathbf{T}}_{SCAD}$, that apply Soft and SCAD on $\hat{\mathbf{T}}_{OLS}$, respectively. We denote our first two estimators as:

$$\boxed{\hat{\Sigma}_{OLS}^{Soft} = \hat{\mathbf{T}}_{Soft}^{-1} \hat{\mathbf{D}}_{OLS} \hat{\mathbf{T}}_{Soft}^{-T}}$$

$$\boxed{\hat{\Sigma}_{OLS}^{SCAD} = \hat{\mathbf{T}}_{SCAD}^{-1} \hat{\mathbf{D}}_{OLS} \hat{\mathbf{T}}_{SCAD}^{-T}}$$

Note that in [22], the authors have demonstrated that for a non sparse true covariance model, the covariance matrix does not suffer any degradation when thresholding is applied to the off-diagonal entries of $\hat{\Sigma}_{SCM}$. However, this is not the case for the target detection problem where the inverse covariance is used; we found that, and in contrast to our estimators, the scheme in [22] has a deleterious effect on the detection performance when compared to those of $\hat{\Sigma}_{SCM}$ and $\hat{\Sigma}_{OLS}$.

B. A generalization of the estimator in [25]

We present the same concept in [25], but by modifying the procedure by which the entries of \mathbf{T} have been estimated. Note that $\det(\mathbf{T}) = 1$ and $\Sigma = \mathbf{T}^{-1} \mathbf{D} \mathbf{T}^{-T}$. It follows that $\det(\Sigma) = \det(\mathbf{D}) = \prod_{t=1}^p \theta_t^2$. Hence, the negative normal log-likelihood of $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{p \times n}$, ignoring an irrelevant constant, satisfies:

$$\Lambda = -2 \log(L(\Sigma, \mathbf{x}_1, \dots, \mathbf{x}_n)) = n \log(\det(\mathbf{D})) + \mathbf{X}^T (\mathbf{T}^T \mathbf{D}^{-1} \mathbf{T}) \mathbf{X} = n \log(\det(\mathbf{D})) + (\mathbf{T} \mathbf{X})^T \mathbf{D}^{-1} (\mathbf{T} \mathbf{X}) = n \sum_{t=1}^p \log \theta_t^2 + \sum_{t=1}^p \sum_{i=1}^n \epsilon_{i,t}^2 / \theta_t^2.$$

By adding a penalty function $\sum_{t=2}^p \sum_{j=1}^{t-1} p_\alpha \{|C_{t,j}|\}$ to Λ , where $p_\alpha \in \{p_\alpha^{L_1}, p_{\alpha,a>2}^{SCAD}\}$ (see subsection III. A) with $\alpha \in [0, \infty)$, we have:

$$n \log \theta_1^2 + \sum_{i=1}^n \frac{\epsilon_{i,1}^2}{\theta_1^2} + \sum_{t=2}^p \left(n \log \theta_t^2 + \sum_{i=1}^n \frac{\epsilon_{i,t}^2}{\theta_t^2} + \sum_{j=1}^{t-1} p_\alpha \{|C_{t,j}|\} \right) \quad (5)$$

Obviously, minimizing (5) with respect to θ_1^2 and θ_t^2 gives the solutions $\hat{\theta}_1^2 = \frac{1}{n} \sum_{i=1}^n \epsilon_{i,1}^2 = \frac{1}{n} \sum_{i=1}^n x_{i,1}^2$ and $\hat{\theta}_t^2 = \frac{1}{n} \sum_{i=1}^n \epsilon_{i,t}^2 = \frac{1}{n} \sum_{i=1}^n (x_{i,t} - \sum_{j=1}^{t-1} C_{t,j} x_{i,j})^2$, respectively.

It remains to estimate the entries of \mathbf{T} by minimizing (5) with respect to β_t . From equation (2) and (3), the minimization problem to solve for each $t \in [2, p]$ is:

$$\begin{aligned} \hat{\beta}_t &= \underset{\beta_t}{\operatorname{argmin}} \sum_{i=1}^n \frac{\epsilon_{i,t}^2}{\theta_t^2} + \sum_{j=1}^{t-1} p_\alpha \{|C_{t,j}|\} \\ &= \underset{\beta_t}{\operatorname{argmin}} \frac{1}{\theta_t^2} \sum_{i=1}^n \left(x_{i,t} - \sum_{j=1}^{t-1} C_{t,j} x_{i,j} \right)^2 + \sum_{j=1}^{t-1} p_\alpha \{|C_{t,j}|\} \quad (6) \\ &= \underset{\beta_t}{\operatorname{argmin}} \frac{1}{\theta_t^2} \|\mathbf{y}_t - \mathbf{A}_{n,t} \beta_t\|_F^2 + \sum_{j=1}^{t-1} p_\alpha \{|C_{t,j}|\} \end{aligned}$$

By denoting $l(\beta_t) = \frac{1}{\theta_t^2} \|\mathbf{y}_t - \mathbf{A}_{n,t} \beta_t\|_F^2$ and $r(\beta_t) = \sum_{j=1}^{t-1} p_\alpha \{|C_{t,j}|\} = \sum_{j=1}^{t-1} r_j(C_{t,j})$, we solve β_t iteratively using the General Iterative Shrinkage and Thresholding (GIST) algorithm [28]:

$$\begin{aligned} \hat{\beta}_t^{(k+1)} &= \underset{\beta_t}{\operatorname{argmin}} l(\beta_t^{(k)}) + r(\beta_t) + \nabla l(\beta_t^{(k)})^T (\beta_t - \beta_t^{(k)}) \\ &\quad + \frac{w^{(k)}}{2} \|\beta_t - \beta_t^{(k)}\|^2 \quad (7) \\ &= \underset{\beta_t}{\operatorname{argmin}} 0.5 \|\beta_t - \mathbf{u}_t^{(k)}\|^2 + \frac{1}{w^{(k)}} r(\beta_t) \end{aligned}$$

where $\mathbf{u}_t^{(k)} = \beta_t^{(k)} - \nabla l(\beta_t^{(k)}) / w^{(k)}$, and $w^{(k)}$ is the step size initialized using the Barzilai-Browein rule [29].

By decomposing (7) into independent (t-1) univariate optimization problems, we have for $j = 1, \dots, t-1$:

$$C_{t,j}^{(k+1)} = \underset{C_{t,j}}{\operatorname{argmin}} 0.5 \|C_{t,j} - u_{t,j}^{(k)}\|^2 + \frac{1}{w^{(k)}} r_j(C_{t,j}) \quad (8)$$

where $\mathbf{u}_t^{(k)} = [u_{t,1}^{(k)}, \dots, u_{t,t-1}^{(k)}]^T \in \mathbb{R}^{(t-1)}$.

By solving (8) with the L_1 -norm penalty, $p_\alpha^{L_1}$, we have the following closed form solution:

$$C_{t,j,(L_1)}^{(k+1)} = \operatorname{sign}(u_{t,j}^{(k)}) \max(0, u_{t,j}^{(k)} - \alpha / w^{(k)}) \quad (9)$$

For the SCAD penalty function, $p_{\alpha,a>2}^{SCAD}$, we can observe that it contains three parts for three different conditions (see Subsection III-A). In this case, by recasting problem (8) into three minimization sub-problems for each condition, and after solving them, one can obtain the following three sub-solutions $h_{t,j}^1$, $h_{t,j}^2$, and $h_{t,j}^3$, where:

$$\begin{aligned} h_{t,j}^1 &= \operatorname{sign}(u_{t,j}^{(k)}) \min(\alpha, \max(0, |u_{t,j}^{(k)}| - \alpha / w^{(k)})), \\ h_{t,j}^2 &= \operatorname{sign}(u_{t,j}^{(k)}) \min(a\alpha, \max(\alpha, \frac{w^{(k)}|u_{t,j}^{(k)}|(a-1)-a\alpha}{w^{(k)}(a-2)})), \\ h_{t,j}^3 &= \operatorname{sign}(u_{t,j}^{(k)}) \max(a\lambda, |u_{t,j}^{(k)}|). \end{aligned}$$

Hence, we have the following closed form solution:

$$\begin{aligned} C_{t,j,(SCAD)}^{(k+1)} &= \underset{q_{t,j}}{\operatorname{argmin}} 0.5 \|q_{t,j} - u_{t,j}^{(k)}\|^2 + \frac{1}{w^{(k)}} r_j(q_{t,j}) \quad (10) \\ \text{s.t. } q_{t,j} &\in \{h_{t,j}^1, h_{t,j}^2, h_{t,j}^3\} \end{aligned}$$

At the end, we denote our last two estimators as:

$$\boxed{\hat{\Sigma}_{L_1} = \hat{\mathbf{T}}_{L_1}^{-1} \hat{\mathbf{D}} \hat{\mathbf{T}}_{L_1}^{-T}}$$

$$\boxed{\hat{\Sigma}_{SCAD} = \hat{\mathbf{T}}_{SCAD}^{-1} \hat{\mathbf{D}} \hat{\mathbf{T}}_{SCAD}^{-T}}$$

where $\hat{\mathbf{T}}_{L_1}$ and $\hat{\mathbf{T}}_{SCAD}$ have respectively $-\hat{C}_{t,j,(L_1)}$ and $-\hat{C}_{t,j,(SCAD)}$ in the (t, j) th position for $t \in [2, p]$ and $j \in [1, t-1]$, whereas $\hat{\mathbf{D}}$ has the entries $(\hat{\theta}_1^2, \hat{\theta}_t^2)$ on its diagonal. Note that in [25], the authors have used the local quadratic approximation (LQA) of the L_1 -norm in order to get a closed form solution for β_t in equation (6). Our algorithm is now more general since after exploiting the GIST algorithm to solve (6), it can be easily extended to some other penalties such as SCAD [27], Capped-L1 penalty [30], [31], [32], Log Sum Penalty[33], Minimax Concave Penalty [34] etc. and they all have closed-form solutions [28]. In this paper, we are only interested to the L_1 and SCAD penalty functions.

IV. HYPERSPECTRAL ANOMALY DETECTION

We first describe the Kelly anomaly detector [26] used for the detection evaluation. Next, we present two subsections to gauge the detection performances of our estimators $\{\hat{\Sigma}_{OLS}^{Soft}, \hat{\Sigma}_{OLS}^{SCAD}, \hat{\Sigma}_{L_1}, \hat{\Sigma}_{SCAD}\}$ when compared to the traditional ones $\{\hat{\Sigma}_{SCM}, \hat{\Sigma}_{OLS}\}$ and state-of-the-art: $\hat{\Sigma}_{SMT}$ [23], $B_k(\hat{\Sigma}_{SCM})$ [21], and the two estimators that apply Soft and SCAD thresholding on the off-diagonal entries of $\hat{\Sigma}_{SCM}$ in [22], and which will be denoted in the following experiments as $\hat{\Sigma}_{SCM}^{Soft}$ and $\hat{\Sigma}_{SCM}^{SCAD}$, respectively. Note that the tuning

Models	Σ	$\hat{\Sigma}_{SCM}$	$\hat{\Sigma}_{OLS}$	$\hat{\Sigma}_{OLS}^{Soft}$	$\hat{\Sigma}_{OLS}^{SCAD}$	$\hat{\Sigma}_{L_1}$	$\hat{\Sigma}_{SCAD}$	$\hat{\Sigma}_{SMT}$	$B_k(\hat{\Sigma}_{SCM})$	$\hat{\Sigma}_{SCM}^{Soft}$	$\hat{\Sigma}_{SCM}^{SCAD}$
Model 1	0.9541	0.7976	0.8331	0.9480	0.9480	0.9509	0.9509	0.9503	0.9509	0.9509	0.9509
Model 2	0.9540	0.7977	0.8361	0.9124	0.9124	0.9264	0.9264	0.9184	0.9478	0.9274	0.9270
Model 3	0.9541	0.7978	0.8259	0.8169	0.8257	0.8236	0.8261	0.7798	0.5321	0.5969	0.5781
MUSE	Not known	0.6277	0.6575	0.9620	0.9643	0.8844	0.8844	0.7879	0.9277	0.7180	0.7180

Table 1. A List of Area Under Curve (AUC) values of our estimators $\hat{\Sigma}_{OLS}^{Soft}$, $\hat{\Sigma}_{OLS}^{SCAD}$, $\hat{\Sigma}_{L_1}$, $\hat{\Sigma}_{SCAD}$ when compared to some others.

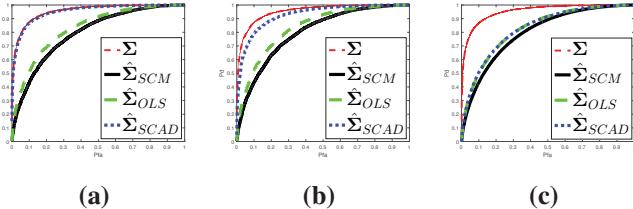


Fig. 1. ROC curves three Models. (a) Model 1. (b) Model 2. (c) Model 3.

parameter λ (in subsection III-A) and α (in subsection III-B) are chosen automatically using a 5-fold crossvalidated loglikelihood procedure (see Subsection 4.2 in [25] for details).

Suppose the following signal model:

$$\begin{cases} H_0 : \mathbf{x} = \mathbf{n}, & \mathbf{x}_i = \mathbf{n}_i, \quad i = 1, \dots, n \\ H_1 : \mathbf{x} = \gamma \mathbf{d} + \mathbf{n}, & \mathbf{x}_i = \mathbf{n}_i, \quad i = 1, \dots, n \end{cases} \quad (11)$$

where $\mathbf{n}_1, \dots, \mathbf{n}_n$ are n i.i.d p -vectors, each follows a multivariate Normal distribution $\mathcal{N}(\mathbf{0}, \Sigma)$. \mathbf{d} is an unknown steering vector and which denotes the presence of an anomalous signal with unknown amplitude $\gamma > 0$. After some calculation (refer to [26] and both Subsection II. B and Remark II. 1 in [35] for details), the Kelly anomaly detector is described as follows:

$$D_{KellyAD\hat{\Sigma}}(\mathbf{x}) = \mathbf{x}^T \hat{\Sigma}_{SCM}^{-1} \mathbf{x} \stackrel{H_1}{\gtrless} \delta, \quad (12)$$

where δ is a prescribed threshold value. In the following two subsections, the detection performances of the estimators, when are plugged in $D_{KellyAD,\hat{\Sigma}}$ are evaluated by the Receiver Operating Characteristics (ROC) curves and their corresponding Area Under Curves (AUC) values.

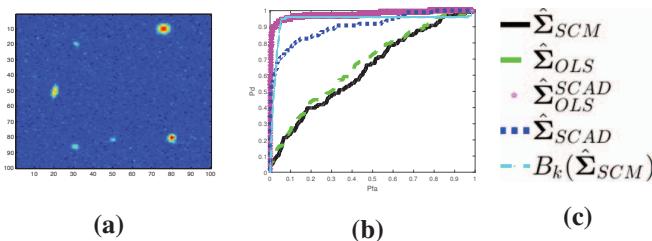


Fig. 2. (a) MUSE HSI (average). (b) ROC curves for MUSE. (c) Legend.

A. Monte-Carlo simulations

The experiments are conducted on three covariance models:

- Model 1: $\Sigma = \mathbf{I}$, the identity matrix,
- Model 2: the Autoregressive model order 1, AR(1), $\Sigma = [\sigma_{gl}]_{p \times p}$, where $\sigma_{gl} = c^{|g-l|}$, for $c = 0.3$,
- Model 3: $\Sigma = [\sigma_{gl}]_{p \times p}$, where $\sigma_{gl} = (1 - ((|g-l|)/r))_+$, for $r = p/2$: the triangular matrix.

Model 1 is very sparse and model 2 is approximately sparse. Model 3 with $r = p/2$ is considered the least sparse [22] among the three models we consider.

The computations have been made through 10^5 Monte-Carlo trials and the ROC curves are drawn for a signal to noise ratio equal to 15dB. We choose $n = 80$ for covariance estimation under Gaussian assumption, and set $p = 60$. The artificial anomaly we consider is a vector containing normally distributed pseudorandom numbers (to have fair results, the same vector is used for the three models). The ROC curves for Model 1, 2 and 3 are shown in Fig. 1(a), 1(b) and 1(c), respectively, and their corresponding AUC values are presented in Table 1. For a clear presentation of the figures, we only exhibit the ROC curves for Σ , $\hat{\Sigma}_{SCM}$, $\hat{\Sigma}_{OLS}$.

The highest AUC values are shown in bold in Table 1. For both Model 1 and 2, our estimators significantly improve the detection performances comparing to those of the traditional estimators ($\hat{\Sigma}_{SCM}$, $\hat{\Sigma}_{OLS}$), and have competitive detection results with state-of-the-art. An important finding is that even for a non sparse covariance model (that is, Model 3), our estimators do not show a harm on the detection when compared to those of $\hat{\Sigma}_{SCM}$, $\hat{\Sigma}_{OLS}$. Despite $\hat{\Sigma}_{OLS}^{Soft}$, $\hat{\Sigma}_{OLS}^{SCAD}$ and $\hat{\Sigma}_{L_1}$ have slightly lower AUC values than for $\hat{\Sigma}_{OLS}$, this is still a negligible degradation on the detection. Thus, considering that $\hat{\Sigma}_{OLS}^{Soft}$, $\hat{\Sigma}_{OLS}^{SCAD}$ and $\hat{\Sigma}_{L_1}$ have no worse detection results than to that of $\hat{\Sigma}_{OLS}$ is still acceptable.

B. Application on experimental data

Our estimators are now evaluated for galaxy detection on the Multi Unit Spectroscopic Explorer (MUSE) data cube (see [36]). It is a 100×100 image and consists of 3600 bands in wavelengths ranging from 465-930 nm. We used one band of each 60, so that 60 bands in total. Figure 2(a) exhibits the mean power in dB over the 60 bands. The covariance matrix is estimated using a sliding window of size 9×9 , having $n = 80$ secondary data (after excluding only the test pixel). The mean has been removed from the given HSI. Figure 2(b) exhibits the ROC curves [37] of our estimators when compared to some others, and their AUC values are shown in Table 1. Note that the curves for $\hat{\Sigma}_{SMT}$, $\hat{\Sigma}_{SCM}^{Soft}$ and $\hat{\Sigma}_{SCM}^{SCAD}$ are not drawn but their AUC values are shown in Table 1.

The estimators $\hat{\Sigma}_{OLS}^{Soft}$, $\hat{\Sigma}_{OLS}^{SCAD}$ achieve higher detection results than for all the others, whereas both $\hat{\Sigma}_{L_1}$ and $\hat{\Sigma}_{SCAD}$ achieve only a lower AUC values than for $B_k(\hat{\Sigma}_{SCM})$.

V. CONCLUSION

Two methods are outlined to impose sparsity on the covariance matrix via its Cholesky factor \mathbf{T} . Some Monte-Carlo simulations as well as experimental data demonstrate the effectiveness (in terms of anomaly detection) of the two proposed methods using the Kelly anomaly detector.

REFERENCES

- [1] G. Shaw and D. Manolakis, "Signal processing for hyperspectral image exploitation," *IEEE Signal Processing Magazine*, vol. 19, no. 1, pp. 12–16, Jan 2002.
- [2] D. Manolakis, D. Marden, and G. Shaw, "Hyperspectral image processing for automatic target detection applications," *Lincoln Laboratory Journal*, vol. 14, no. 1, pp. 79–116, 2003.
- [3] D. G. Manolakis, R. B. Lockwood, and T. W. Cooley, *Hyperspectral Imaging Remote Sensing: Physics, Sensors, and Algorithms*. Cambridge University Press, 2016.
- [4] N. K. Patel, C. Patnaik, S. Dutta, A. M. Shekh, and A. J. Dave, "Study of crop growth parameters using airborne imaging spectrometer data," *International Journal of Remote Sensing*, vol. 22, no. 12, pp. 2401–2411, 2001. [Online]. Available: <http://www.tandfonline.com/doi/abs/10.1080/01431160117383>
- [5] B. Datt, T. R. McVicar, T. G. van Niel, D. L. B. Jupp, and J. S. Pearlman, "Preprocessing eo-1 hyperion hyperspectral data to support the application of agricultural indexes," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 41, pp. 1246–1259, Jun. 2003.
- [6] B. Hörig, F. Kühn, F. Oschütz, and F. Lehmann, "HyMap hyperspectral remote sensing to detect hydrocarbons," *International Journal of Remote Sensing*, vol. 22, pp. 1413–1422, May 2001.
- [7] D. Manolakis and G. Shaw, "Detection algorithms for hyperspectral imaging applications," *Signal Processing Magazine, IEEE*, vol. 19, no. 1, pp. 29–43, 2002.
- [8] D. W. J. Stein, S. G. Beaven, L. E. Hoff, E. M. Winter, A. P. Schaum, and A. D. Stocker, "Anomaly detection from hyperspectral imagery," *IEEE Signal Processing Magazine*, vol. 19, no. 1, pp. 58–69, Jan 2002.
- [9] M. T. Eismann, A. D. Stocker, and N. M. Nasrabadi, "Automated hyperspectral cueing for civilian search and rescue," *Proceedings of the IEEE*, vol. 97, no. 6, pp. 1031–1055, June 2009.
- [10] D. Manolakis, E. Truslow, M. Pieper, T. Cooley, and M. Brueggeman, "Detection algorithms in hyperspectral imaging systems: An overview of practical algorithms," *IEEE Signal Processing Magazine*, vol. 31, no. 1, pp. 24–33, Jan 2014.
- [11] D. Manolakis, R. Lockwood, T. Cooley, and J. Jacobson, "Is there a best hyperspectral detection algorithm?" *Proc. SPIE 7334*, p. 733402, 2009.
- [12] J. Frontera-Pons, F. Pascal, and J. P. Ovarlez, "False-alarm regulation for target detection in hyperspectral imaging," in *Computational Advances in Multi-Sensor Adaptive Processing (CAMSAP), 2013 IEEE 5th International Workshop on*, Dec 2013, pp. 161–164.
- [13] J. Frontera-Pons, M. A. Veganzones, S. Velasco-Forero, F. Pascal, J. P. Ovarlez, and J. Chanussot, "Robust anomaly detection in hyperspectral imaging," in *2014 IEEE Geoscience and Remote Sensing Symposium*, July 2014, pp. 4604–4607.
- [14] S. Matteoli, M. Diani, and G. Corsini, "A tutorial overview of anomaly detection in hyperspectral images," *Aerospace and Electronic Systems Magazine, IEEE*, vol. 25, no. 7, pp. 5–28, 2010.
- [15] I. S. Reed and X. Yu, "Adaptive multiple-band CFAR detection of an optical pattern with unknown spectral distribution," *IEEE Transactions on Acoustics Speech and Signal Processing*, vol. 38, pp. 1760–1770, Oct. 1990.
- [16] E. J. Kelly, "An adaptive detection algorithm," *IEEE Transactions on Aerospace and Electronic Systems*, vol. AES-22, no. 2, pp. 115–127, March 1986.
- [17] C.-I. Chang and S.-S. Chiang, "Anomaly detection and classification for hyperspectral imagery," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 40, no. 6, pp. 1314–1325, Jun 2002.
- [18] J. C. Harsanyi, *Detection and classification of subpixel spectral signatures in hyperspectral image sequences*. Ph.D. dissertation, Dept. Electr. Eng., Univ. Maryland, Baltimore, MD, USA, 1993.
- [19] J. Frontera-Pons, F. Pascal, and J. P. Ovarlez, "Adaptive nonzero-mean gaussian detection," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 55, no. 2, pp. 1117–1124, Feb 2017.
- [20] M. Pourahmadi, "Joint mean-covariance models with applications to longitudinal data: unconstrained parameterisation," *Biometrika*, vol. 86, no. 3, pp. 677–690, 1999.
- [21] P. J. Bickel and E. Levina, "Regularized estimation of large covariance matrices," *The Annals of Statistics*, vol. 36, no. 1, pp. 199–227, 2008.
- [22] A. J. Rothman, E. Levina, and J. Zhu, "Generalized thresholding of large covariance matrices," *Journal of the American Statistical Association*, vol. 104, no. 485, pp. 177–186, 2009.
- [23] G. Cao and C. Bouman, "Covariance estimation for high dimensional data vectors using the sparse matrix transform," in *Advances in Neural Information Processing Systems 21*. Curran Associates, Inc., 2009, pp. 225–232.
- [24] W. B. Wu and M. Pourahmadi, "Nonparametric estimation of large covariance matrices of longitudinal data," *Biometrika*, vol. 90, no. 4, p. 831, 2003. [Online]. Available: <http://dx.doi.org/10.1093/biomet/90.4.831>
- [25] J. Z. Huang, N. Liu, M. Pourahmadi, and L. Liu, "Covariance matrix selection and estimation via penalised normal likelihood," *Biometrika*, vol. 93, no. 1, pp. 85–98, 2006.
- [26] E. J. Kelly, "An adaptive detection algorithm," *Aerospace and Electronic Systems, IEEE Transactions on*, vol. 23, no. 1, pp. 115–127, November 1986.
- [27] J. Fan and R. Li, "Variable selection via nonconcave penalized likelihood and its oracle properties," *Journal of the American Statistical Association*, vol. 96, no. 456, pp. 1348–1360, 2001.
- [28] P. Gong, C. Zhang, Z. Lu, J. Huang, and J. Ye, "Gist: General iterative shrinkage and thresholding for non-convex sparse learning," 2013.
- [29] J. Barzilai and J. M. Borwein, "Two-point step size gradient methods," *IMA Journal of Numerical Analysis*, vol. 8, no. 1, pp. 141–148, 1988.
- [30] T. Zhang, "Analysis of multi-stage convex relaxation for sparse regularization," *J. Mach. Learn. Res.*, vol. 11, pp. 1081–1107, Mar. 2010. [Online]. Available: <http://dl.acm.org/citation.cfm?id=1756006.1756041>
- [31] ——, "Multi-stage convex relaxation for feature selection," *Bernoulli*, vol. 19, no. 5B, pp. 2277–2293, 11 2013. [Online]. Available: <http://dx.doi.org/10.3150/12-BEJ452>
- [32] C. Z. P. Gong, J. Ye, "Multi-stage multi-task feature learning," in *Advances in Neural Information Processing Systems 25*. Curran Associates, Inc., 2012, pp. 1988–1996.
- [33] E. Candès, M. B. Wakin, and S. P. Boyd, "Enhancing sparsity by reweighted ℓ_1 minimization," *Journal of Fourier Analysis and Applications*, vol. 14, no. 5, pp. 877–905, 2008.
- [34] C.-H. Zhang, "Nearly unbiased variable selection under minimax concave penalty," *Ann. Statist.*, vol. 38, no. 2, pp. 894–942, 04 2010. [Online]. Available: <http://dx.doi.org/10.1214/09-AOS729>
- [35] J. Frontera-Pons, M. A. Veganzones, F. Pascal, and J. P. Ovarlez, "Hyperspectral anomaly detectors using robust estimators," *Selected Topics in Applied Earth Observations and Remote Sensing, IEEE Journal of*, vol. 9, no. 2, pp. 720–731, Feb 2016.
- [36] official website of the MUSE Project, "<http://muse.univ-lyon1.fr/>"
- [37] Y. Zhang, B. Du, and L. Zhang, "A sparse representation-based binary hypothesis model for target detection in hyperspectral images," *IEEE TGRS*, vol. 53, no. 3, pp. 1346–1354, March 2015.