

# Robust Online Covariance and Sparse Precision Estimation Under Arbitrary Data Corruption

Tong Yao and Shreyas Sundaram

**Abstract**—Gaussian graphical models are widely used to represent correlations among entities but remain vulnerable to data corruption. In this work, we introduce a modified trimmed-inner-product algorithm to robustly estimate the covariance in an online scenario even in the presence of arbitrary and adversarial data attacks. At each time step, data points, drawn nominally independently and identically from a multivariate Gaussian distribution, arrive. However, a certain fraction of these points may have been arbitrarily corrupted. We propose an online algorithm to estimate the sparse inverse covariance (i.e., precision) matrix despite this corruption. We provide the error-bound and convergence properties of the estimates to the true precision matrix under our algorithms.

## I. INTRODUCTION

Graph modeling is at the core of modern statistical learning, with applications spanning a wide range of disciplines, including finance and economics [1], neuroscience [2], and health and social science [3]. These problems address the challenges of estimating complex relationships between multiple variables to gain insight into the underlying interactions.

One way to represent and quantify these relationships is by learning the covariance and inverse covariance matrix through the collected multivariate data, which captures the degree to which the variables change together. Inverse covariance estimation is a crucial task, as the inverse covariance matrix (also known as the precision matrix) can reveal the underlying conditional independence structure between variables. By estimating the sparse inverse covariance matrix, one can identify the most relevant connections between variables, leading to more interpretable and efficient models.

Assuming that the underlying relationships follow Gaussian distributions, techniques such as graphical lasso [4]–[6] have been developed to tackle the problem by incorporating sparsity-promoting penalties. Different optimization methods have been proposed to solve the graphical lasso problem, including coordinate descent [5], proximal methods [7], [8], alternating minimization methods [9], [10], and Newton-conjugate gradient methods [11].

In many real-world scenarios, data is continuously generated and collected, making these traditional batch-processing methods infeasible or computationally expensive. More recently, graphical lasso has been extended to learn the static and dynamic relationships between variables in an online manner [12]–[16]. Online estimation methods offer several

advantages, including scalability, adaptability to changes, real-time decision-making, and reduced computational cost.

However, real-world data often contain outliers, corruptions, and even maliciously poisoned data, which severely impact the performance of statistical estimators. Robust estimation aims to develop methods that are less sensitive to such outliers, providing reliable and accurate estimates in the presence of contaminated data. Traditional robust mean estimators, such as the median-of-means [17] and the trimmed mean [18], [19] have been explored in the literature. Early works on robust statistics attempt to estimate the mean given an outlier model [20], [21], and recent advancements consider stronger contamination models [22]–[24]. For a comprehensive overview of robust mean estimation, readers are directed to a survey [25]. Classical robust covariance estimation techniques include Minimum Covariance Determinant [26] [27] to find a given proportion of uncorrupted observations and compute their empirical covariance matrix and the truncated inner product [28] which filters out data with large absolute values to mitigate the impact of arbitrary corruption. These traditional robust estimators were developed for batch datasets and are not suitable for data arriving in a sequence.

Motivated by the streaming nature and the potential corruption of modern datasets, we propose online and robust covariance and sparse inverse covariance estimation algorithms that enable efficient and robust updating of the estimates as new potentially corrupted and noisy data arrive. We show through theoretical performance guarantees and experimental simulations that our algorithms are effective and less sensitive to data corruption. By addressing these problems, we can build more reliable, accurate, and interpretable models, leading to a better understanding of complex systems and more informed decision-making across various domains.

## II. PROBLEM DEFINITION

Consider a set of  $p$  random variables  $X = [x_1 \ x_2 \ \cdots \ x_p]^T$ , that are jointly Gaussian with zero mean and covariance  $S^*$ . These variables can be represented by a graph  $G = (V, E)$ , where  $V = \{v_1, \dots, v_p\}$  is the set of nodes, with each node  $v_i$  representing the random variable  $x_i$ . An edge  $(v_i, v_j) \in E$  indicates that variable  $x_j$  is conditionally dependent on  $x_i$ , given all other random variables. Conversely, if  $(v_i, v_j) \notin E$ ,  $v_j$  is conditionally independent of  $v_i$ , given all the other variables. This lack of an edge corresponds to a zero-entry in the precision matrix  $\theta^* = (S^*)^{-1}$  [9], [29].

Tong Yao and Shreyas Sundaram are with Elmore Family School of Electrical and Computer Engineering at Purdue University. Email: {yao127, sundara2}@purdue.edu

This research was supported by NSF CAREER award 1653648.

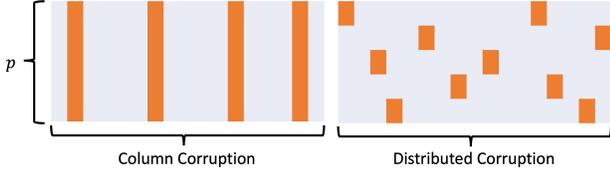


Fig. 1: Distribution of corruptions for data matrix, where orange indicates corrupted data points. In these examples, the number of corrupted products of observations is  $\eta t = 4$ .

These relationships (i.e., graph structure) between the variables are unknown *a priori*; the goal is to infer the edges of the graph based on samples from random variables. More specifically, at each time step  $t \in \{1, 2, \dots\}$ , the network generates a data vector  $X_t = [x_{1,t}, x_{2,t}, \dots, x_{p,t}]^T \in \mathbb{R}^p$ . We assume that each  $X_t$  is independently and identically sampled from the underlying Gaussian distribution  $X_t \sim \mathcal{N}(0, S^*)$ .

We consider the existence of an adversary who can inspect all data up to time  $t$  and replace them with arbitrary values. Let  $\mathcal{X} = [X_1, \dots, X_t] \in \mathbb{R}^{p \times t}$  be the data matrix containing all data points received up to time  $t$ . We say that the dataset  $\mathcal{X}$  is  $\eta$ -corrupted if, for each pair of variables  $x_i$  and  $x_j$ , where  $i, j \in \{1, \dots, p\}$ , the product of their observations  $x_{i,1}x_{j,1}, \dots, x_{i,t}x_{j,t}$  contains at most  $\eta t$  corrupted entries, where the corruption parameter  $\eta \in (0, 1/32)$ . We will explain the range of  $\eta$  in Section IV.

In the  $\eta$ -corruption scenario, the corrupted entries can be located in the same columns (observations) of  $\mathcal{X}$  across rows (variables). This *column corruption* allows at most  $\eta t$  corrupted data points for each variable. The  $\eta$ -corruption also implies that the corrupted data can be distributed anywhere in  $\mathcal{X}$ . This is defined as a *distributed corruption* scenario where each variable has at most  $\eta/2$  arbitrary corrupted observations. We show the corruption models in Fig. 1.

We apply the tilde notation to any potentially corrupted data or dataset, e.g.,  $\tilde{x}_t$  is a potentially corrupted data sample,  $\tilde{\mathcal{X}}$  is the  $\eta$ -corrupted data matrix containing observations up to time  $t$ , and  $\tilde{X}_t = [\tilde{x}_{1,t}, \tilde{x}_{2,t}, \dots, \tilde{x}_{p,t}]^T$  is the potentially corrupted observation vector at time  $t$ .

Given a sequence of potentially corrupted observations  $\{\tilde{X}_1, \tilde{X}_2, \dots\}$ , our objective is to perform real-time (i.e., online) robust estimation of the sample covariance matrix  $S^*$  and infer the conditional dependencies between the variables through estimating the sparse precision matrix  $\theta^*$ .

### III. BACKGROUND

We build our approach on an alternating minimization algorithm proposed in [10] (batch) and [16] (online) to solve problem (1). The online algorithm was shown in [16] to achieve a similar result with fewer iterations in real-time settings. We provide the background of the algorithm from [16] and subsequently discuss our modification to account for arbitrary data corruption in Section IV.

Given a set of uncorrupted data  $\{X_1, X_2, \dots, X_t\}$  up to time  $t$ , the maximum likelihood estimation problem for

recovering  $\theta$  is given by

$$\underset{\theta_t \in \mathcal{S}_{++}^p}{\text{minimize}} \quad -\log \det \theta_t + \text{tr}(S_t \theta_t) + \lambda |\theta_t|_{l_1}, \quad (1)$$

where the set of  $p \times p$  positive definite matrices is denoted by the set  $\mathcal{S}_{++}^p$ , and  $S_t = \frac{1}{t} \sum_{i=1}^t X_i X_i^T$  is the sample covariance matrix constructed from all the data up to the time step  $t$ . The terms  $-\log \det \theta_t + \text{tr}(S_t \theta_t)$  are derived from the Gaussian log-likelihood function [4], where  $\text{tr}$  denotes trace, and the term  $|\theta|_{l_1} = \sum_{i,j=1}^p |\theta_{ij}|$  is the element-wise  $l_1$  norm encouraging sparsity of the solution regulated by the penalty parameter  $\lambda \geq 0$ .

The approaches in [10] and [16] formulate the primal and dual objective functions for problem (1). The primal of (1) is:

$$\begin{aligned} & \underset{\Omega_t \in \mathcal{S}_{++}^p, \Phi_t \in \mathcal{S}_{++}^p}{\text{minimize}} \quad -\log \det \Omega_t + \text{tr}(S_t \Phi_t) + \lambda |\Phi_t|_{l_1} \\ & \text{subject to } \Phi_t = \Omega_t. \end{aligned} \quad (2)$$

The dual of (1) is given by

$$\begin{aligned} & \underset{\Gamma_t \in \mathcal{S}_{++}^p}{\text{minimize}} \quad -\log \det \Gamma_t - p \\ & \text{subject to } |\Gamma_{ij,t} - S_{ij,t}| \leq \lambda \quad \forall i, j, \end{aligned} \quad (3)$$

where the symmetric positive definite matrix  $\Gamma_t$  is the dual variable, and  $A_{ij,t}$  denotes the  $(i, j)$ -th element of matrix  $A$  at time  $t$ .

Given the sample covariance matrix  $S_t$ , the alternating minimization follows the iterative sequence of updates, where each iteration is indexed by the variable  $k \in \mathbb{N}$ :

$$\Omega_t^{k+1} = \underset{\Omega \in \mathcal{S}_{++}^p}{\text{argmin}} \quad -\log \det \Omega + \text{tr}(\Gamma_t^k \Omega), \quad (4)$$

$$\begin{aligned} \Phi_t^{k+1} = \underset{\Phi \in \mathcal{S}_{++}^p}{\text{argmin}} \quad & \text{tr}(S_t \Phi) + \lambda |\Phi|_{l_1} - \text{tr}(\Gamma_t^k \Phi) \\ & + \frac{\zeta_t^k}{2} \|\Omega_t^{k+1} - \Phi\|_F^2, \end{aligned} \quad (5)$$

$$\Gamma_t^{k+1} = \Gamma_t^k + \zeta_t^k (\Omega_t^{k+1} - \Phi_t^{k+1}). \quad (6)$$

In the above equations,  $\zeta_t^k$  is a step size, and  $\|A\|_F = \sqrt{\text{tr}(AA^T)}$  denotes the Frobenius norm of a given matrix  $A$ . The step size  $\zeta_t^k$  is chosen to guarantee convergence of the estimates  $\Gamma_t^k$  to their desired quantities and can be set as a constant  $\zeta = \zeta_t^k < a^2$  for some constant  $a$  [16]. The problem of recovering  $\theta_t$  reduces to solving (4)–(6). We provide robust estimation details in the next section.

### IV. ONLINE ROBUST COVARIANCE AND PRECISION ESTIMATION

In this section, we will introduce the robust covariance estimator and describe the modification of the sparse precision estimation.

#### A. Online Trimmed Inner Product

We describe an online trimmed inner product estimator based on the trimmed mean estimator in [23] (batch) and [24] (online). We modify the estimator for robust online estimation of covariance. The robust estimator can mitigate

the influence of arbitrary corruption and we will show the theoretical performance guarantees in Section V.

At the beginning of the algorithm, the system designer will select an initialization time step  $t_0 \in \mathbb{N}$ , a desired confidence interval  $\delta \in (0, 1)$ , and the corruption level  $\eta$ . Note that  $\eta$  can be an estimate of the upper bound of the corruption rate.

At each time step, a potentially corrupted data vector  $\tilde{X}_t \in \mathbb{R}^p$  arrives. We compute the product  $\tilde{s}_t = \tilde{X}_t \tilde{X}_t^T$ , where the  $(i, j)$ -th entry  $\tilde{s}_{ij,t} = \tilde{x}_{i,t} \tilde{x}_{j,t}$ .

For  $t < t_0$ , we temporarily store the product matrices in a data vector for each  $(i, j)$ -th entry  $\tilde{Y}_{ij,t} = [\tilde{s}_{ij,1}, \dots, \tilde{s}_{ij,t}]$ .

At  $t = t_0$ , we start with the set of potentially corrupted products  $\tilde{Y}_{ij,t_0} = [\tilde{s}_{ij,1}, \dots, \tilde{s}_{ij,t_0}]$ . Given the corruption parameter  $\eta$  and the desired confidence level  $\delta$ , define

$$\epsilon = 8\eta + 12 \frac{\log(4/\delta)}{t_0}. \quad (7)$$

Note that for sufficiently large  $t$ ,  $\epsilon < 0.25$  since  $\eta < 1/32$ . For each  $\tilde{Y}_{ij,t_0}$ , let  $\tilde{s}_{ij,1}^* \leq \tilde{s}_{ij,2}^* \leq \dots \leq \tilde{s}_{ij,t_0}^*$  be a non-decreasing rearrangement of  $\tilde{s}_{ij,1}, \dots, \tilde{s}_{ij,t_0}$ . For each  $i, j \in \{1, \dots, p\}$ , define  $\alpha_{ij} = \tilde{s}_{ij,\epsilon t}^*$  and  $\beta_t = \tilde{s}_{ij,(1-\epsilon)t}^*$  to be *trimming values*. For each  $\alpha_{ij} \leq \beta_{ij}$ , and  $s \in \mathbb{R}$ , define the *trim estimator*

$$\phi_{\alpha_{ij}, \beta_{ij}}(s) = \begin{cases} \beta_{ij} & \text{if } s > \beta_{ij}, \\ s & \text{if } s \in [\alpha_{ij}, \beta_{ij}], \\ \alpha_{ij} & \text{if } s < \alpha_{ij}. \end{cases} \quad (8)$$

We then initialize the  $(i, j)$ -th entry of the robust estimation of sample covariance

$$\hat{S}_{ij,t_0} = \frac{\sum_{k=1}^{t_0} \phi_{\alpha_{ij}, \beta_{ij}}(\tilde{s}_{ij,k})}{t_0}. \quad (9)$$

For all  $t > t_0$ , the process recursively updates the estimate using the previous estimate  $\hat{S}_{ij,t-1}$  and the new product  $\tilde{s}_{ij,t}$

$$\hat{S}_{ij,t} = \frac{(t-1)\hat{S}_{ij,t-1} + \phi_{\alpha_{ij}, \beta_{ij}}(\tilde{s}_{ij,t})}{t}. \quad (10)$$

We include the pseudo-code implementation for the online robust covariance estimation in Algorithm 1.

### B. Online Alternating Minimization Algorithm

At the initialization time step  $t_0$ , we obtain  $\hat{S}_{t_0}$  from (9) and initialize  $\Gamma_{t_0} = \hat{S}_{t_0} + \lambda I_p$ , where  $I_p \in \mathbb{R}^{p \times p}$  is an identity matrix, with data up to a user-defined  $t_0 \in \{1, 2, \dots\}$ . The sparsity penalty parameter  $\lambda$  ensures that the initial dual variable is positive definite, i.e.,  $\lambda > \max\{0, -\lambda_{\min}(\hat{S}_{t_0})\}$ .

We apply the robust covariance estimator and simplify the alternating minimization algorithm by allowing one optimization iteration per time step. Taking the derivatives of the expressions for  $\Omega_t$  and  $\Phi_t$  and equating them to 0, we obtain the closed-form updates of (4) and (5):

$$\Omega_t = (\Gamma_t)^{-1}, \quad (11)$$

$$\Phi_t = \frac{1}{\zeta_{t-1}} \mathcal{S}_\lambda(\zeta_{t-1} \Omega_t - \hat{S}_t + \Gamma_t). \quad (12)$$

Here,  $\mathcal{S}_\lambda(x) = \text{sign}(x)(\max(|x| - \lambda, 0))$  is the soft-thresholding operator (applied element-wise to a matrix argument). Following these update rules,  $\Omega_t$  is interpreted as an

---

### Algorithm 1: Online Trimmed Inner Product

---

**Parameter:**  $t_0, \delta, \eta$

**Input:**  $\tilde{X}_t = [\tilde{x}_{1,t}, \dots, \tilde{x}_{p,t}]$

**for**  $i, j \in \{1, 2, \dots, p\}$  **do**

    Compute  $\tilde{s}_{ij,t} = \tilde{x}_{i,t} \tilde{x}_{j,t}$

**if**  $t < t_0$  **then**

        Store data in  $\tilde{Y}_{ij,t_0}$

**if**  $t = t_0$  **then**

        Compute  $\epsilon = 8\eta + 12 \frac{\log(4/\delta)}{t_0}$

        Sort data in  $\tilde{Y}_{ij,t_0}$  and determine trimming thresholds  $\alpha_{ij}, \beta_{ij}$

$\hat{S}_{ij,t_0} = \frac{\sum_{k=1}^{t_0} \phi_{\alpha_{ij}, \beta_{ij}}(\tilde{s}_{ij,k})}{t_0}$

**if**  $t > t_0$  **then**

        update trimmed inner product

$\hat{S}_{ij,t} = \frac{(t-1)\hat{S}_{ij,t-1} + \phi_{\alpha_{ij}, \beta_{ij}}(\tilde{s}_{ij,t})}{t}$

**return**  $\hat{S}_t$

---



---

### Algorithm 2: Online Graphical Alternating Minimization Algorithm (O-GAMA)

---

**Parameter:**  $\lambda, t_0$

**Input:** Stream of potentially corrupted multivariate data  $\tilde{X}_1, \tilde{X}_2, \dots \in \mathbb{R}^p$

**for**  $t \in \{1, 2, \dots\}$  **do**

$\hat{S}_t = \text{online-trimmed-inner-product}(\tilde{X}_t)$

**if**  $t = t_0$  **then**

$\Gamma_{t_0} = \hat{S}_{t_0} + \lambda I_p$

        Choose  $\zeta_t \in (0, (\lambda_{\min}(\Gamma_{t_0}))^2)$

**else if**  $t > t_0$  **then**

$\Gamma_t = \mathcal{C}_\lambda(\Gamma_t - \hat{S}_t + \zeta_{t-1}(\Gamma_{t-1})^{-1}) + \hat{S}_t$

        Choose  $\zeta_t \in (0, (\lambda_{\min}(\Gamma_t))^2)$

        Update  $\Omega_t$  and  $\Phi_t$  as per (11) and (12)

**return** Sequence of sparse precision estimates

$\Phi_{t_0+1}, \Phi_{t_0+2}, \dots$

---

approximately sparse precision matrix and  $\Phi_t$  is interpreted as the estimate of the sparse precision matrix. Substituting (11) and (12) for the variables in (6), and using the clip function  $\mathcal{C}_\lambda(x) = \min(\max(x, -\lambda), \lambda)$  with property  $x = \mathcal{S}_\lambda(x) + \mathcal{C}_\lambda(x)$ , the dual update (6) can be written as:

$$\Gamma_t = \mathcal{C}_\lambda(\Gamma_{t-1} - \hat{S}_t + \zeta_{t-1}(\Gamma_{t-1})^{-1}) + \hat{S}_t. \quad (13)$$

The algorithm first iterates through (13) to obtain the new dual variable. The updated dual variable is then used to update (11) and (12).

We provide the pseudo-code implementation of the robust sparse precision estimation algorithm in Algorithm 2.

## V. THEORETICAL ANALYSIS

In this section, we provide theoretical guarantees of the quality of the robust sample covariance estimator and then analyze the quality of the sparse precision estimation.

### A. Robust Sample Covariance

Let  $tS = \sum_{k=1}^t X_k X_k^T$  be a  $p \times p$  random matrix, where each  $X_t$  is i.i.d. sampled from  $\mathcal{N}(0, S^*)$ . The matrix  $tS$  follows the Wishart distribution  $tS \sim \mathcal{W}(S^*, t)$  which arises as the distribution of the sample covariance matrix for a sample of a multivariate normal distribution [30]. Using results from the Wishart distribution, each  $(i, j)$ -th entry of  $S$ , denoted  $S_{ij}$ , is a real-valued random variable with mean  $S_{ij}^*$  and finite variance  $\sigma_{S_{ij}}^2 = S_{ij}^{*2} + S_{ii}^* S_{jj}^*$ .

We set up the analysis for each  $(i, j)$ -th entry of  $S$ . For simplicity, we omit the subscript  $ij$  of  $S$  for the following analysis. Define  $\bar{S} = S - S^*$ . For  $0 < q < 1$ , define the quantile

$$Q_q(\bar{S}) = \sup\{M \in \mathbb{R} : \mathbb{P}(\bar{S} \geq M) \geq 1 - q\}. \quad (14)$$

We have  $\mathbb{P}(\bar{S} \geq Q_q(\bar{S})) = 1 - q$  and from Chebyshev's inequality,

$$1 - q = \mathbb{P}(\bar{S} \geq Q_q(\bar{S})) \leq \mathbb{P}(|\bar{S}| \geq Q_q(\bar{S})) \leq \frac{\sigma_S^2}{Q_q^2(\bar{S})}.$$

As a result,

$$|Q_q(\bar{S})| \leq \frac{\sigma_S}{\sqrt{1-q}}. \quad (15)$$

The following result upper bounds the trimming values using quantiles.

*Lemma 1 ([23]):* Consider the corruption-free samples  $y_1, \dots, y_t$ . With probability at least  $1 - 4e^{-\epsilon t/12}$ , the inequalities

$$\begin{aligned} |\{i : y_i \geq S^* + Q_{1-2\epsilon}(\bar{S})\}| &\geq 3/2\epsilon t \\ |\{i : y_i \leq S^* + Q_{1-\epsilon/2}(\bar{S})\}| &\geq (1 - (3/4)\epsilon)t \\ |\{i : y_i \leq S^* + Q_{2\epsilon}(\bar{S})\}| &\geq 3/2\epsilon t \\ |\{i : y_i \geq S^* + Q_{\epsilon/2}(\bar{S})\}| &\geq (1 - (3/4)\epsilon)t \end{aligned}$$

hold simultaneously. We denote the event when the above four inequalities hold as event  $A$ . On event  $A$ , since corruption  $\eta \leq \epsilon/8$ , the following inequalities also hold

$$Q_{1-2\epsilon}(\bar{S}) \leq \beta - S^* \leq Q_{1-\epsilon/2}(\bar{S}), \quad (16)$$

$$Q_{\epsilon/2}(\bar{S}) \leq \alpha - S^* \leq Q_{2\epsilon}(\bar{S}). \quad (17)$$

With the background mentioned above, we provide the following result on the error of the robust estimation of covariance. Note that the proofs of all results are included in the Appendix.

*Theorem 1:* Let  $t_0 > \max(\frac{3 \log(8/\delta)}{2\eta}, \frac{12 \log(4/\delta)}{0.25-8\eta})$  and fix  $\delta \in [4e^{-t_0}, 1)$ . Following the procedures of Algorithm 1, with probability at least  $1 - \delta$ , the estimates satisfy

$$\begin{aligned} |\hat{S}_{ij,t} - S_{ij}^*| &\leq \left( \sqrt{2} + \frac{\sqrt{6}}{9} \right) \sigma_{S_{ij}} \sqrt{\frac{\log(4/\delta)}{t}} \\ &\quad + \frac{43\sqrt{2}}{12} \sigma_{S_{ij}} \sqrt{\epsilon}, \forall t \geq t_0, \end{aligned} \quad (18)$$

where  $\sigma_{S_{ij}} = \sqrt{S_{ij}^{*2} + S_{ii}^* S_{jj}^*}$ .

The theorem illustrates that for each  $(i, j)$ -th entry of the sample covariance, the estimation error consists of an initialization error derived from setting the trimming thresholds and a convergence error influenced by the variance of  $x_i, x_j$ , as well as the covariance between these variables. Intuitively, robust covariance estimation becomes more challenging in the presence of corruption when  $x_i$  and  $x_j$  exhibit large variances, or when their covariance is large, implying considerable joint variation.

We provide the following result on the convergence of the robustly estimated sample covariance matrix.

*Corollary 1:* Let  $t_0 > \max(\frac{3 \log(8/\delta)}{2\eta}, \frac{12 \log(4/\delta)}{0.25-8\eta})$  and fix  $\delta \in [4e^{-t_0}, 1)$ . Following the procedures of Algorithm 1, there exists a set of sample paths of measure 1, such that for each sample path in that set, there exists a finite time  $\bar{t}$ , such that for all  $t \geq \bar{t}$ , the robust sample covariance satisfies

$$\begin{aligned} \|\hat{S}_t - S^*\|_F &\leq \left( \sqrt{2} + \frac{\sqrt{6}}{9} \right) p \sigma_S \sqrt{\frac{\log(4/\delta)}{t}} \\ &\quad + \frac{43\sqrt{2}}{12} p \sigma_S \sqrt{\epsilon}, \forall t \geq t_0. \end{aligned} \quad (19)$$

where  $\sigma_S = \max_{ij} \sqrt{S_{ij}^{*2} + S_{ii}^* S_{jj}^*}$ .

Using the definition of  $\epsilon$  in (7) with Corollary 1, we obtain the result below on the convergence of estimates.

*Corollary 2:* Let  $t_0 > \max(\frac{3 \log(8/\delta)}{2\eta}, \frac{12 \log(4/\delta)}{0.25-8\eta})$  and fix  $\delta \in [4e^{-t_0}, 1)$ . Following the procedures of Algorithm 1, the robust estimate of sample covariance satisfies the following inequality almost surely,

$$\limsup_{t \rightarrow \infty} \|\hat{S}_t - S^*\|_F \leq \frac{43}{6} p \sigma_S \sqrt{4\eta + 6 \frac{\log(4/\delta)}{t_0}}, \quad (20)$$

where  $\sigma_S = \max_{ij} \sqrt{S_{ij}^{*2} + S_{ii}^* S_{jj}^*}$ .

From the above results, we observe that the estimation error converges to the error introduced by initialization. Thus, it requires a relatively large  $t_0$  to give reasonable estimates.

### B. Robust Sparse Precision Matrix

In this subsection, we provide error bounds for the estimates of sparse precision estimates. It was shown in [10], that  $\Omega^* = \theta^*$  is the optimal solution of (1) given  $S^*$ , if and only if  $\Gamma^* = (\Omega^*)^{-1}$  is a fixed point of (13) given  $S^*$  (i.e., with access to the true covariance matrix). We will analyze the dual variable  $\Gamma_t$  given  $t$  data points and as the number of data points increases.

We provide the following theorem to show the lower and upper bound of eigenvalues of the dual variables  $\Gamma_t$ . The robustly estimated sample covariance matrix is real and symmetric. However,  $S_t$  may no longer be positive semidefinite. For the inverse covariance estimation, we will use the penalty parameter  $\lambda$  to enforce positive definiteness in the estimates.

*Theorem 2:* Let  $t_0 > \max(\frac{3 \log(8/\delta)}{2\eta}, \frac{12 \log(4/\delta)}{0.25-8\eta})$  and fix  $\delta \in [4e^{-t_0}, 1)$ . Define  $b = \max_{t \geq t_0} \lambda_{\max}(\hat{S}_t) + p\lambda$  and  $a = \min_{t \geq t_0} e^{g(t)} b^{1-p}$ , where  $g(t) = \log \det \Gamma_{t_0} - \sum_{k=t_0}^{t-1} \Delta_k$ ,

and  $\Delta_t = \text{tr}((\Gamma_{t+1} - \Gamma_t)(-\Omega_t)) + \frac{1}{2\zeta_t} \|\Gamma_{t+1} - \Gamma_t\|_F^2$ . For any finite  $t \geq t_0$ , there exists a sufficiently large  $\lambda$  such that iterates of  $\Gamma_t$  of Algorithm 2 satisfy  $0 \prec aI_p \preceq \Gamma_t \preceq bI_p$ .

Note that here we provide the result of boundedness for finite time steps, and we can always select a sufficiently large penalty parameter  $\lambda$  such that estimates are bounded for finite time steps. We note that  $\hat{S}_t$  converges to a range of  $S^*$  almost surely and  $\Gamma_t$  converges to a range of  $\Gamma^*$  almost surely. The upper bound  $b$  is finite almost surely. For the lower bound, we observe from experiments that  $\sum_{k=t_0}^t \Delta_k$  is bounded for all  $t \geq t_0$  (see Fig. 4 in the Appendix). As a result, the appropriate value for the penalty parameter is also bounded, i.e.,  $\lambda \in (0, \infty)$ .

To analyze the finite-time performance of robust sparse precision estimation, we introduce the following theorem. The proof is similar to the online graphical alternating minimization algorithm [16], and we incorporate the robust sample covariance estimation error.

*Theorem 3:* Let  $t_0 > \max(\frac{3 \log(8/\delta)}{2\eta}, \frac{12 \log(4/\delta)}{0.25-8\eta})$  and fix  $\delta \in [4e^{-t_0}, 1)$ . Assume  $\forall t \geq t_0$ , the quantities  $\Gamma_t$  satisfy  $aI \preceq \Gamma_t \preceq bI$  and  $\zeta_t = \zeta$ , and assume  $\Gamma^*$  satisfies  $aI_p \preceq \Gamma^* \preceq bI_p$ . Then at time-step  $t+1$ , the dual variable  $\Gamma_{t+1}$  satisfies the following condition:

$$\|\Gamma_{t+1} - \Gamma^*\|_F \leq r^{t+1-t_0} \|\Gamma_{t_0} - \Gamma^*\|_F + 2 \sum_{k=t_0}^t r^{t-k} \|\hat{S}_{k+1} - S^*\|_F, \quad (21)$$

where  $r = \max\left\{|1 - \frac{\zeta}{a^2}|, |1 - \frac{\zeta}{b^2}|\right\}$ .

We defer discussing the implications of the result until the asymptotic analysis is presented later in this section.

For analysis, we require the following lemma.

*Lemma 2:* Let  $0 < r < 1$  and let  $\{\rho_t\}$  be a positive scalar sequence. Assume that  $\lim_{t \rightarrow \infty} \rho_t \leq \bar{\rho}$ . Then

$$\limsup_{t \rightarrow \infty} \sum_{l=0}^t r^{t-l} \rho_l \leq \frac{\bar{\rho}}{1-r}. \quad (22)$$

*Corollary 3:* Let  $t_0 > \max(\frac{3 \log(8/\delta)}{2\eta}, \frac{12 \log(4/\delta)}{0.25-8\eta})$  and fix  $\delta \in [4e^{-t_0}, 1)$ . Assume that there exist constants  $0 < a < b$  such that, for all  $t \geq t_0$ , the quantities  $\Gamma_t$  computed in Algorithm 2 satisfy  $aI \preceq \Gamma_t \preceq bI$  and  $\zeta_t = \zeta < a^2$ . Then, as the number of data points  $t \rightarrow \infty$ , the result  $\Gamma_t$  converges to a range of the optimal solution  $\Gamma^*$  almost surely:

$$\limsup_{t \rightarrow \infty} \|\Gamma_{t+1} - \Gamma^*\|_F \leq \frac{43}{6(1-r)} \sigma_S \sqrt{4\eta + 6 \frac{\log(4/\delta)}{t_0}}, \quad (23)$$

where  $r = \max\left\{|1 - \frac{\zeta}{a^2}|, |1 - \frac{\zeta}{b^2}|\right\}$  and  $\sigma_S = \max_{ij} \sqrt{S_{ij}^{*2} + S_{ii}^* S_{jj}^*}$ .

From the above results, the estimation error is introduced by the initialization of the dual variable and the initialization of the trimming threshold. Since the error introduced by dual initialization converges exponentially to 0 (for  $0 < r < 1$ ), the convergence of the dual variable is dominated by the convergence of the sample covariance matrices.

Asymptotically, the estimation error is introduced only by the trimming threshold determined during the initialization process. Notably, the error correlates with the corruption rate  $\eta$  but not to any severity of the corruption.

## VI. EXPERIMENTS

In this section, we demonstrate the effectiveness of robust online estimation algorithms through experimental simulations.

### Experimental Setup

We first generate a sparse Erdos-Renyi network with  $p$  nodes following the steps in [31]. We generate a  $p \times p$  symmetric matrix  $A$  by setting the probability of two nodes having no edge as 0.95, and the probability of two nodes having an edge as 0.05; in the latter case, the value of the corresponding entry is chosen to be uniformly distributed within certain intervals:

$$A_{ij} = \begin{cases} 0 & \text{Pr} = 0.95 \\ \text{Unif}([-0.6, -0.3] \cup [0.3, 0.6]) & \text{otherwise} \end{cases}$$

To ensure that the covariance matrix is positive definite, we set  $\theta^* = (S^*)^{-1} = A + (\xi + |\lambda_{\min}(A)|)I$ , and adjust  $\xi$  to make  $\lambda_{\min}(\theta^*) = 1$ .

We then generate a clean data matrix  $X \in \mathbb{R}^{p \times t}$  from a Gaussian distribution  $\mathcal{N}(0, \theta^*)$ . For each row in the data matrix, we randomly select  $\eta t$  of them to be corrupted, where  $\eta = 0.03$ . The corruption data are generated by two normal distributions  $\mathcal{N}(\mu, \sigma)$ , representing small and large corruptions, respectively, as follows:

$$(1) \mu = 1, \sigma = 2; \quad (2) \mu = 1, \sigma = 5.$$

### Performance of Covariance and Sparse Precision Estimates

For this set of experiments, we let  $p = 10, t_0 = 100, \delta = 0.9$ , and  $\lambda = 0.15$ . We include the simulation results for robust covariance estimation in Fig. 2 and robust sparse precision estimation in Fig. 3.

In Fig. 2, we let  $\hat{S}_t$  be zero matrices for all  $t \leq t_0$ . The performance of the two robust estimators is similar and the deviation curves overlap in the plot. We observe that the robust covariance estimator is effective against large corruption and shows improvement over small corruption. Asymptotically, the performance of the robust estimator is influenced by the initialization error and stays within a bounded range of the true covariance.

In Fig. 3, we observe that both small and large corruptions have a significant effect on the estimates of the sparse precision matrix. However, with the proposed robust estimator, the estimates are significantly improved and are bounded within a close range of the estimates without any corruption.

## VII. CONCLUSION AND FUTURE WORK

In this work, we proposed online robust covariance and sparse precision estimators. In [24], it was shown that by trading off computation and memory complexity, the initial estimation error can be eliminated. We will include such modifications to reduce initialization errors in future work.

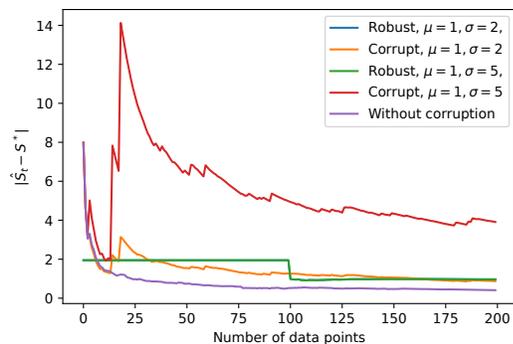


Fig. 2: Estimate of sample covariances.

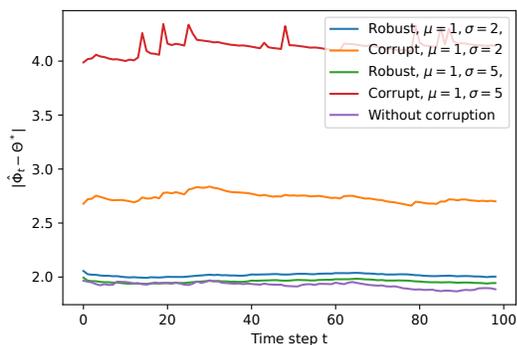


Fig. 3: Estimate of sparse inverse covariances.

## REFERENCES

[1] Jianqing Fan, Yuan Liao, and Martina Mincheva. High-dimensional covariance matrix estimation in approximate factor models. *The Annals of Statistics*, 39(6):3320–3356, 2011.

[2] Stephen M. Smith, Karla L. Miller, Gholamreza Salimi-Khorshidi, Matthew Webster, Christian F. Beckmann, Thomas E. Nichols, Joseph D. Ramsey, and Mark W. Woolrich. Network modelling methods for fmri. *NeuroImage*, 54(2):875–891, 2011.

[3] Bernd Fellinghauer, Peter Bühlmann, Martin Ryffel, Michael von Rhein, and Jan D. Reinhardt. Stable graphical model estimation with random forests for discrete, continuous, and mixed variables. *Computational Statistics and Data Analysis*, 64:132–152, 2013.

[4] Ming Yuan and Yi Lin. Model selection and estimation in the Gaussian graphical model. *Biometrika*, 94(1):19–35, 2007.

[5] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. Sparse inverse covariance estimation with the graphical lasso. *Biostatistics*, 9(3):432–441, 2008.

[6] Onureena Banerjee, Laurent El Ghaoui, and Alexandre d’Aspremont. Model selection through sparse maximum likelihood estimation for multivariate Gaussian or binary data. *Journal of Machine Learning Research*, 9(Mar):485–516, 2008.

[7] Benjamin Rolfs, Bala Rajaratnam, Dominique Guillot, Ian Wong, and Arian Maleki. Iterative thresholding algorithm for sparse inverse covariance estimation. In *Advances in Neural Information Processing Systems*, pages 1574–1582, 2012.

[8] Cho-Jui Hsieh, Mátyás A Sustik, Inderjit S Dhillon, Pradeep K Ravikumar, and Russell Poldrack. BIG & QUIC: Sparse inverse covariance estimation for a million variables. In *Advances in Neural Information Processing Systems*, pages 3165–3173, 2013.

[9] Katya Scheinberg, Shiqian Ma, and Donald Goldfarb. Sparse inverse covariance selection via alternating linearization methods. In *Advances in Neural Information Processing Systems*, pages 2101–2109, 2010.

[10] Onkar Dalal and Bala Rajaratnam. Sparse Gaussian graphical model

estimation via alternating minimization. *Biometrika*, 104(2):379–395, 2017.

[11] Salar Fattahi, Richard Y. Zhang, and Somayeh Sojoudi. Linear-time algorithm for learning large-scale sparse graphical models. *IEEE Access*, 7:12658–12672, 2019.

[12] Shuheng Zhou, John Lafferty, and Larry Wasserman. Time varying undirected graphs. *Machine Learning*, 80(2-3):295–319, 2010.

[13] Ricardo Pio Monti, Peter Hellyer, David Sharp, Robert Leech, Christoforos Anagnostopoulos, and Giovanni Montana. Estimating time-varying brain connectivity networks from functional MRI time series. *NeuroImage*, 103:427–443, 2014.

[14] Mladen Kolar, Le Song, Amr Ahmed, Eric P Xing, et al. Estimating time-varying networks. *The Annals of Applied Statistics*, 4(1):94–123, 2010.

[15] David Hallac, Youngsuk Park, Stephen Boyd, and Jure Leskovec. Network inference via the time-varying graphical lasso. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 205–213, 2017.

[16] T. Yao and S. Sundaram. Online estimation of sparse inverse covariances. In *Proceedings of the American Control Conference (ACC)*, pages 1935–1940, 2021.

[17] Noga Alon, Yossi Matias, and Mario Szegedy. The space complexity of approximating the frequency moments. *Journal of Computer and System Sciences*, 58(1):137–147, 1999.

[18] John W Tukey and Donald H McLaughlin. Less vulnerable confidence and significance procedures for location based on a single sample: Trimming/winsorization 1. *Sankhyā: The Indian Journal of Statistics, Series A*, pages 331–352, 1963.

[19] Elvezio M Ronchetti and Peter J Huber. *Robust statistics*. John Wiley & Sons, 2009.

[20] John W Tukey. A survey of sampling from contaminated distributions. *Contributions to probability and statistics*, pages 448–485, 1960.

[21] Peter J Huber. Robust estimation of a location parameter. In *Breakthroughs in statistics*, pages 492–518. Springer, 1992.

[22] Ilias Diakonikolas, Daniel M. Kane, and Ankit Pensia. Outlier robust mean estimation with subgaussian rates via stability. In *Advances in Neural Information Processing Systems*, volume 33, pages 1830–1840. Curran Associates, Inc., 2020.

[23] Gabor Lugosi and Shahar Mendelson. Robust multivariate mean estimation: the optimality of trimmed mean. *The Annals of Statistics*, 49(1):393–410, 2021.

[24] Tong Yao and Shreyas Sundaram. Robust online and distributed mean estimation under adversarial data corruption. In *Proceedings of the 2022 IEEE 61st Conference on Decision and Control (CDC)*, pages 4193–4198, 2022.

[25] Gábor Lugosi and Shahar Mendelson. Mean estimation and regression under heavy-tailed distributions: A survey. *Foundations of Computational Mathematics*, 19(5):1145–1190, 2019.

[26] Peter J. Rousseeuw. Least median of squares regression. *Journal of the American Statistical Association*, 79(388):871–880, 1984.

[27] Peter J. Rousseeuw and Katrien Van Driessen. A fast algorithm for the minimum covariance determinant estimator. *Technometrics*, 41(3):212–223, 1999.

[28] Lingxiao Wang and Quanquan Gu. Robust Gaussian graphical model estimation with arbitrary corruption. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 3617–3626. PMLR, 06–11 Aug 2017.

[29] Daphne Koller and Nir Friedman. *Probabilistic graphical models: principles and techniques*. MIT press, 2009.

[30] Alvin C. Rencher and William F. Christensen. *Methods of multivariate analysis*. Wiley, Hoboken, New Jersey, 2012.

[31] Karthik Mohan, Palma London, Maryam Fazel, Daniela Witten, and Su-In Lee. Node-based learning of multiple Gaussian graphical models. *The Journal of Machine Learning Research*, 15(1):445–488, 2014.

## A. Proof of Theorem 1

First, we state the following lemma to complete the proof.

*Lemma 3:* (Bernstein's inequality) Let  $S_1, \dots, S_t$  be independent random variables, such that  $|S_i - \mathbb{E}[S_i]| \leq a$  for all  $i$ . Then for any  $\delta \in (0, 1)$  and  $t \in \mathbb{N}$ , we have

$$\mathbb{P}\left(\left|\frac{1}{t} \sum_{k=1}^t (S_k - \mathbb{E}[S_k])\right| \leq \frac{2a}{3t} \log(2/\delta) + \sqrt{\frac{2\sigma^2 \log(2/\delta)}{t}}\right) \geq 1 - \delta, \quad (24)$$

where  $\sigma^2 = \frac{1}{t} \sum_{k=1}^t \sigma_{S_k}^2$ .

To simplify notation, we omit the subscript  $ij$  that denotes the  $(i, j)$ -th entry in the proof and let  $\phi_{t_0} = \phi_{\alpha_{t_0}, \beta_{t_0}}$ . From the triangle inequality, we can separate

$$\begin{aligned} \left| \frac{1}{t} \sum_{k=1}^t \phi_{t_0}(\tilde{s}_k) - S^* \right| &\leq |\mathbb{E}[\phi_{t_0}(S)] - S^*| \\ &+ \left| \frac{1}{t} \sum_{k=1}^t \phi_{t_0}(s_k) - \mathbb{E}[\phi_{t_0}(S)] \right| \\ &+ \left| \frac{1}{t} \sum_{k=1}^t \phi_{t_0}(\tilde{s}_k) - \frac{1}{t} \sum_{k=1}^t \phi_{t_0}(s_k) \right|. \end{aligned} \quad (25)$$

The error between the expected values of the trimming operator and the true mean (i.e.,  $S^* = \mathbb{E}[S]$ ) is introduced by data beyond the trimming thresholds

$$\begin{aligned} |\mathbb{E}[\phi_{t_0}(S)] - S^*| &\leq |\mathbb{E}[(\alpha_{t_0} - S)\mathbf{1}_{S \leq \alpha_{t_0}}] \\ &+ \mathbb{E}[(S - S)\mathbf{1}_{\alpha_{t_0} < S < \beta_{t_0}}] + \mathbb{E}[(\beta_{t_0} - S)\mathbf{1}_{S \geq \beta_{t_0}}]|. \end{aligned} \quad (26)$$

Notice that on the right-hand side of (26), the first term is positive, the second term is zero, and the third term is negative.

For  $t_0 \geq \frac{3 \log(8/\delta)}{2\eta}$ , the probability of event  $A$  in Lemma 1 satisfies  $4e^{-\epsilon t_0/12} \leq 4e^{-\frac{2}{3}\eta t_0} \leq \delta/2$ . With probability at least  $1 - \delta/2$ , on event  $A$ , the trimming thresholds satisfy (16) and (17), and recall  $\bar{S} = S - S^*$ ,

$$\begin{aligned} |\mathbb{E}[\phi_{t_0}(S)] - S^*| &\leq \max \left\{ \left| \mathbb{E}[(Q_{2\epsilon}(\bar{S}) - \bar{S})\mathbf{1}_{\bar{S} \leq Q_{2\epsilon}(\bar{S})}] \right|, \right. \\ &\left. \left| \mathbb{E}[(\bar{S} - Q_{1-2\epsilon}(\bar{S}))\mathbf{1}_{\bar{S} \geq Q_{1-2\epsilon}(\bar{S})}] \right| \right\}. \end{aligned} \quad (27)$$

The first term on the right-hand side of (27) can be upper

bounded as

$$\begin{aligned} &|\mathbb{E}[(Q_{2\epsilon}(\bar{S}) - \bar{S})\mathbf{1}_{\bar{S} \leq Q_{2\epsilon}(\bar{S})}]| \\ &\leq |\mathbb{E}[Q_{2\epsilon}(\bar{S})\mathbf{1}_{\bar{S} \leq Q_{2\epsilon}(\bar{S})}]| + |\mathbb{E}[\bar{S}\mathbf{1}_{\bar{S} \leq Q_{2\epsilon}(\bar{S})}]| \\ &\stackrel{(a)}{\leq} |Q_{2\epsilon}(\bar{S})|\mathbb{P}(\bar{S} \leq Q_{2\epsilon}(\bar{S})) \\ &\quad + \sqrt{\mathbb{E}[\bar{S}^2]\mathbb{E}[(\mathbf{1}_{\bar{S} \leq Q_{2\epsilon}(\bar{S})})^2]} \\ &\stackrel{(b)}{\leq} \frac{\sigma_S}{\sqrt{1-2\epsilon}} 2\epsilon + \sigma_S \sqrt{\mathbb{P}(\bar{S} \leq Q_{2\epsilon}(\bar{S}))} \\ &\stackrel{(c)}{\leq} \frac{\sigma_S}{\sqrt{2\epsilon}} 2\epsilon + \sigma_S \sqrt{\mathbb{P}(\bar{S} \leq Q_{2\epsilon}(\bar{S}))} \\ &= \sigma_S \sqrt{8\epsilon}, \end{aligned}$$

where (a) applies the Cauchy–Schwarz inequality, and (b) follows from the definition and upper bound of quantile in (14) and (15), and (c) is given by  $\epsilon < 0.25$  by the assumption  $t_0 > \frac{12 \log(4/\delta)}{0.25-8\eta}$ . The proof is similar for the second term on the right-hand side of (27), and we obtain

$$|\mathbb{E}[\phi_{t_0}(S)] - S^*| \leq \sigma_S \sqrt{8\epsilon}. \quad (28)$$

Using (16) and (17), and the upper bound on  $\mathbb{E}[\phi_{t_0}(S)]$  derived from (28), we have that any data point  $s_k$ ,

$$\begin{aligned} |\phi_{t_0}(s_k) - \mathbb{E}[\phi_{t_0}(S)]| &\leq \max\{|Q_{\epsilon/2}(\bar{S})|, |Q_{1-\epsilon/2}(\bar{S})|\} + \sigma_S \sqrt{8\epsilon} \\ &\leq \frac{\sigma_S}{\sqrt{\epsilon/2}} + \sigma_S \sqrt{8\epsilon}, \end{aligned} \quad (29)$$

where the last inequality follows from (15) and  $\epsilon < 0.25$ .

Applying Bernstein's inequality in Lemma 3, where  $a$  is the right-hand side of (29), with probability at least  $1 - \delta/2$ , we bound the second term of (25) as

$$\begin{aligned} \left| \frac{1}{t} \sum_{k=1}^t \phi_{t_0}(s_k) - \mathbb{E}[\phi_{t_0}(S)] \right| &\leq \sigma_S \sqrt{\frac{2 \log(4/\delta)}{t}} \\ &+ \frac{2\sigma_S}{\sqrt{\epsilon/2}} \frac{\log(4/\delta)}{3t} + \frac{2\sigma_S \sqrt{8\epsilon} \log(4/\delta)}{3t}. \end{aligned} \quad (30)$$

Expanding the second term of on the right-hand-side of (30), we obtain

$$\begin{aligned} \frac{2\sigma_S}{\sqrt{\epsilon/2}} \frac{\log(4/\delta)}{3t} &\stackrel{(a)}{\leq} \frac{2\sigma_S \log(4/\delta)}{3t \sqrt{4\eta + 6 \frac{\log(4/\delta)}{t_0}}} \\ &\stackrel{(b)}{\leq} \sigma_S \frac{2\sqrt{t_0}}{3t} \sqrt{\frac{\log(4/\delta)}{6}} \stackrel{(c)}{\leq} \sigma_S \frac{2}{3} \sqrt{\frac{\log(4/\delta)}{6t}}, \end{aligned} \quad (31)$$

where (a) and (b) are derived from the definition of  $\epsilon$  in (7), and (c) follows from  $t_0/t \leq 1, \forall t \geq t_0$ . Noting that  $\log(4/\delta)/t \leq 1$  by the assumption that  $\delta \geq 4e^{-t_0}$ , we have

$$\begin{aligned} \left| \frac{1}{t} \sum_{k=1}^t \phi_{t_0}(s_k) - \mathbb{E}[\phi_{t_0}(S)] \right| &\leq \left( \sqrt{2} + \frac{\sqrt{6}}{9} \right) \sigma_S \sqrt{\frac{\log(4/\delta)}{t}} + \frac{4\sqrt{2}}{3} \sigma_S \sqrt{\epsilon}. \end{aligned} \quad (32)$$

For corrupted data satisfying  $\phi_{t_0}(\tilde{s}_t) \neq \phi_{t_0}(s_t)$  at time  $t$ , the gap is bounded

$$|\phi_{t_0}(\tilde{s}_t) - \phi_{t_0}(s_t)| \leq |Q_{\epsilon/2}(\bar{S})| + |Q_{1-\epsilon/2}(\bar{S})|,$$

and it follows that,

$$\begin{aligned} & \frac{1}{t} \left| \sum_{k=1}^t \phi_{t_0}(s_k) - \sum_{k=1}^t \phi_{t_0}(\tilde{s}_k) \right| \\ & \stackrel{(a)}{\leq} \eta (|Q_{\epsilon/2}(\bar{S})| + |Q_{1-\epsilon/2}(\bar{S})|) \\ & \stackrel{(b)}{\leq} \eta \left( \frac{\sigma_S}{\sqrt{1-\epsilon/2}} + \frac{\sigma_S}{\sqrt{\epsilon/2}} \right) \\ & \stackrel{(c)}{\leq} \frac{\epsilon}{8} \left( \frac{2\sigma_S}{\sqrt{\epsilon/2}} \right) \leq \frac{\sigma_S \sqrt{2\epsilon}}{4}, \end{aligned} \quad (33)$$

where (a) reflects the potential presence of up to  $\eta t$  corrupted samples at time  $t$ , (b) follows from (15), and (c) follows from  $\eta \leq \epsilon/8$  by definition in (7) and  $\epsilon < 0.25$ .

Using (28), (32), and (33) in (25), with probability at least  $1 - \delta$ ,

$$|\hat{S}_t - S^*| \leq \left( \sqrt{2} + \frac{\sqrt{6}}{9} \right) \sigma_S \sqrt{\frac{\log(4/\delta)}{t}} + \frac{43\sqrt{2}}{12} \sigma_S \sqrt{\epsilon}. \quad (34)$$

### B. Proof of Corollary 1

Let  $\rho_{ij,t} = \left( \sqrt{2} + \frac{\sqrt{6}}{9} \right) \sigma_{S_{ij}} \sqrt{\frac{\log(4/\delta)}{t}} + \frac{43\sqrt{2}}{12} \sigma_{S_{ij}} \sqrt{\epsilon}$ . From Theorem 1, we have

$$\mathbb{P}[|\hat{S}_{ij,t} - S_{ij}^*| \geq \rho_{ij,t}] \leq \delta_t. \quad (35)$$

Applying the union bound, we have

$$\mathbb{P}[\cup_{ij} \{|\hat{S}_{ij,t} - S_{ij}^*| \geq \rho_{ij,t}\}] \leq p^2 \delta_t. \quad (36)$$

This result holds for any fixed  $t$ . Note that from Theorem 1, the range for the selection of  $\delta$  is time-dependent. We use  $\delta_t$  to represent this dependency and each  $\delta_t \in (4e^{-t_0}, 1/p^2)$ . We will now extend this to bound the deviation of  $\hat{S}_t$  from  $S^*$  for all sufficiently large  $t$ .

Define a bad event at time  $t$  as the event that there exists a pair of  $(i, j)$  such that  $|\hat{S}_{ij,t} - S_{ij}^*| > \rho_{ij,t}$ . Define the random variable  $B_t$ , with  $B_t = 1$  if the bad event occurs at the given  $t$ , and 0 otherwise. Define  $B = \sum_{k=t_0}^t B_k$  as the number of bad events up to time  $t$ . We can always select a constant  $c$ , such that  $\delta_t = ce^{-t} \in (4e^{-t_0}, 1/p^2)$ . Summing up the probability of bad events in (36), we have

$$\sum_{k=t_0}^t \mathbb{P}[B_k] \leq cp^2 \sum_{k=t_0}^t e^{-k} < \infty.$$

From the Borel-Cantelli lemma, the probability of infinitely many bad events occurring is 0. Thus, for a set of sample paths of measure 1, there exists a sample-path dependent finite time  $\bar{t}$  such that for  $t \geq \bar{t}$ , no bad event occurs. In other words, for all  $t \geq \bar{t}$ ,  $|\hat{S}_{ij,t} - S_{ij}^*| \leq \rho_{ij,t} \forall i, j$  almost surely, and  $\max_{ij} (|\hat{S}_{ij,t} - S_{ij}^*|) \leq \max_{ij} \rho_{ij,t}$ . Using the matrix norm inequality,  $\|\hat{S}_t - S^*\|_F \leq p\rho_t$ , where

$\rho_t = \max_{ij} \rho_{ij,t}$ . Thus, for each sample path in a set of measure 1, there exists a finite time  $\bar{t}$  such that  $\forall t \geq \bar{t}$ ,

$$\|\hat{S}_t - S^*\|_F \leq p\rho_t = p \max_{ij} \rho_{ij,t}. \quad (37)$$

### C. Proof of Theorem 2

Consider the update for  $\Gamma_t$  given by (13). According to Weyl's inequality and the upper bound of the eigenvalues of the clipping operator (see [16]), the eigenvalues of  $\Gamma_t$  for fixed  $t$  satisfy

$$\lambda_{\max}(\Gamma_t) \leq p\lambda + \lambda_{\max}(\hat{S}_t).$$

A symmetric matrix with positive diagonal elements has at least one positive eigenvalue, thus,  $\lambda_{\max}(\hat{S}_t) > 0, \forall t$ . Define  $b = p\lambda + \max_{t \geq t_0} \lambda_{\max}(\hat{S}_t)$ . We have  $0 < \lambda_{\max}(\Gamma_t) \leq b$  for all finite  $t$ .

Now we show the lower bound for eigenvalues. The step size  $\zeta_t$  is chosen to guarantee convergence as shown in Theorem 3. More specifically, it can be chosen for each iteration by backtracking line search such that for the next iteration,  $\Gamma_t$  satisfies the sufficient descent condition

$$f(\Gamma_t) \leq D_{\zeta_{t-1}}(\Gamma_t, \Gamma_{t-1}), \quad (38)$$

where  $f(\Gamma_t) = -\log \det \Gamma_t$  is the dual objective from (3) and  $D_{\zeta_{t-1}}$  is the quadratic approximation of the dual objective (3) around  $\Gamma_{t-1}$  given by

$$D_{\zeta_{t-1}}(\Gamma_t, \Gamma_{t-1}) = f(\Gamma_{t-1}) + \text{tr}[(\Gamma_t - \Gamma_{t-1})\nabla f(\Gamma_{t-1})] + \frac{1}{2\zeta_{t-1}} \|\Gamma_t - \Gamma_{t-1}\|_F^2. \quad (39)$$

Using the approximation and computing the derivative  $\nabla f$ , we have the sufficient descent condition

$$-\log \det \Gamma_t \leq -\log \det \Gamma_{t-1} + \Delta_{t-1}, \quad (40)$$

where  $\Delta_t = \text{tr}[(\Gamma_{t+1} - \Gamma_t)(-\Omega_t)] + \frac{1}{2\zeta_t} \|\Gamma_{t+1} - \Gamma_t\|_F^2$ .

The step size can be selected small enough to satisfy the above conditions, and it was shown in [10] that  $\zeta_t \leq \min(\lambda_{\min}(\Gamma_{t-1}, \Gamma_t))^2$  is sufficient to satisfy the condition.

Iterating through the above condition, we obtain

$$\log \det \Gamma_t \geq \log \det \Gamma_{t_0} - \sum_{k=t_0}^{t-1} \Delta_k. \quad (41)$$

Let  $a_t$  be the smallest eigenvalue of  $\Gamma_t$ , and let  $g(t) = \log \det \Gamma_{t_0} - \sum_{k=t_0}^{t-1} \Delta_k$ , we have

$$\log a_t + (p-1) \log b \geq \log \det \Gamma_t \geq g(t) \quad (42)$$

$$a_{t+1} \geq e^{g(t)} b^{1-p}. \quad (43)$$

For a finite time step  $t \geq t_0$ , we can always select a large enough  $\lambda$ , such that  $g(t) > -\infty$  for all  $t_0 \leq \dots \leq t$  and  $a_t \geq a$ , where  $a = \min_{t_0 \leq t} (a_t)$ .

Thus, we have  $aI_p \preceq \Gamma_t \preceq bI_p$  for all finite  $t \geq t_0$ .

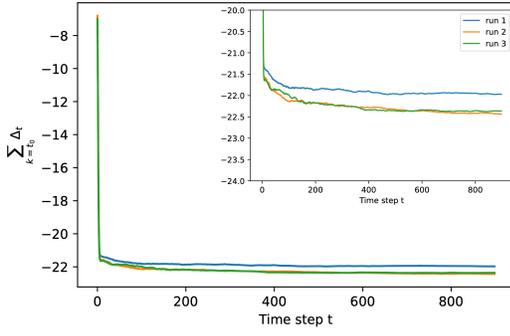


Fig. 4: Boundedness of the sum of error terms  $\Delta$ .

#### D. Proof of Lemma 2

Let  $\xi > 0$  be arbitrary. Since  $\lim_{t \rightarrow \infty} \rho_t \leq \bar{\rho}$ , we let  $\bar{t}$  to be the index such that  $\rho_t \leq \bar{\rho} + \xi$  for all  $t \geq \bar{t}$ .

For all  $t \geq \bar{t} + 1$ , we can partition

$$\sum_{l=0}^t r^{t-l} \rho_l = \sum_{l=0}^{\bar{t}} r^{t-l} \rho_l + \sum_{l=\bar{t}+1}^t r^{t-l} \rho_l. \quad (44)$$

Using properties of geometric series, we have

$$\begin{aligned} \sum_{l=0}^{\bar{t}} r^{t-l} \rho_l &\leq \max_{0 \leq k \leq \bar{t}} \rho_k \sum_{l=0}^{\bar{t}} r^{t-l} = \max_{0 \leq k \leq \bar{t}} \rho_k r^{t-\bar{t}} \sum_{t=0}^{\bar{t}} r^t \\ &\leq \max_{0 \leq k \leq \bar{t}} \rho_k \frac{r^{t-\bar{t}}}{1-r}, \end{aligned} \quad (45)$$

and

$$\sum_{l=\bar{t}+1}^t r^{t-l} \rho_l \leq (\bar{\rho} + \xi) \sum_{l=\bar{t}+1}^t r^{t-l} \leq \frac{\bar{\rho} + \xi}{1-r}. \quad (46)$$

Therefore, by combining the results, we have

$$\sum_{l=0}^t r^{t-l} \rho_l \leq \max_{0 \leq k \leq \bar{t}} \rho_k \frac{r^{t-\bar{t}}}{1-r} + \frac{\bar{\rho} + \xi}{1-r}. \quad (47)$$

Since  $\xi$  is arbitrary and  $r \in (0, 1)$ , we have the result

$$\limsup_{t \rightarrow \infty} \sum_{l=0}^t r^{t-l} \rho_l \leq \frac{\bar{\rho}}{1-r}. \quad (48)$$

#### E. Proof of Corollary 3

From Theorem 3, if  $\forall t \geq t_0$ ,  $\zeta_t^k = \zeta < a^2$ , then  $0 < r < 1$ . The first term  $r^{t+1-t_0} \|\Gamma_{t_0} - \Gamma^*\|_F$  in (21) converges exponentially to zero. The second term in (21) is an instance of Lemma 2, with  $r < 1$  and  $\rho_t \rightarrow \bar{\rho}$  almost surely as  $t \rightarrow \infty$  from Corollary 2.

#### F. Boundedness of Eigenvalues for Sparse Cases

In Fig. 4, we demonstrate the boundedness of error terms defined in (40), by running 3 independent experiments where  $p \gg t_0$ . The subplot with an adjusted  $y$ -axis shows the values with a higher resolution. The total number of data points is 1000,  $t_0 = 100$ ,  $p = 500$ ,  $\lambda = 0.5$ , and other parameters are identical to the previous experiment setup.