# **Optimal Compressive Covariance Sketching via Rank-One Sampling**

Wenbin Wang and Ziping Zhao School of Information Science and Technology, ShanghaiTech University, Shanghai, China

### Background and Motivation

- **High-dimensional Streaming Data:** Each time a snapshot of a data vector  $x \in \mathbb{R}^d$  is generated, with d large.
- Challenges in Modern Data Acquisition:
  - Data generation at unprecedented rate: data samples are 1) not observable due to privacy or security constraints; 2) distributed at multiple locations; 3) online generated on the fly and can only be accessed once.
  - Limited processing power at sensor platforms: 1) time-sensitive: impossible to obtain a complete snapshot of the system; 2) storage-limited: cannot store the whole data set; 3) power-hungry: minimize the number of observations.
- **Covariance Sketching:** Key Observation: The covariance structure can be recovered without measuring the whole data stream.
- Contribution: We propose a nonconvex problem with an efficient algorithm, which can be proved to attain the oracle statistical rate.

# Sampling Model

Consider *n* independent observations  $\{\boldsymbol{x}_t\}_{t=1}^n$ , each drawn from a zero-mean random vector  $\boldsymbol{x}$  with covariance matrix  $\boldsymbol{\Sigma}^*$ . Given *m* sensing vectors  $\{\boldsymbol{a}_i\}_{i=1}^m$ , the quadratic measurement measurement  $y_i$ ,  $i = 1, \ldots, m$ , is given by  $\boldsymbol{y} = \boldsymbol{A}_{\otimes} \operatorname{vec}(\boldsymbol{S}) + \boldsymbol{\eta} = \mathcal{A}(\boldsymbol{S}) + \boldsymbol{\eta}$ ,

#### Main Results

#### **Essential Assumptions**:

The true covariance matrix  $\Sigma^*$  satisfies  $0 < \frac{1}{\kappa} \le \lambda_{\min}(\Sigma^*) \le \lambda_{\max}(\Sigma^*) \le \kappa < \infty,$ 

where  $\boldsymbol{y} [\boldsymbol{y}_1, \dots, \boldsymbol{y}_m]^\top$ ,  $\boldsymbol{\eta} [\eta_1, \dots, \eta_m]^\top$  are additive measurement noises,  $\boldsymbol{A}_{\otimes} = [(\boldsymbol{a}_1 \otimes \boldsymbol{a}_1) \cdots (\boldsymbol{a}_m \otimes \boldsymbol{a}_m)]^\top$  and  $\operatorname{vec}(\boldsymbol{S})$  denotes the vectorization of  $\boldsymbol{S}$  obtained by stacking its columns, and  $\mathcal{A} : \mathbb{R}^{d \times d} \mapsto \mathbb{R}^m$  is a linear operator.

- Assumptions:
  - The sensing vectors { a<sub>i</sub>}<sup>m</sup><sub>i=1</sub> are i.i.d. sub-Gaussian random variables with zero mean and identity covariance.
  - The measurement noises  $\{\eta_i\}_{i=1}^m$  are i.i.d. sub-exponential random variables with mean 0 and variance proxy  $\sigma^2$ .

## Problem Formulation

We propose to estimate the sparse covariance matrices from quadratic measurements using the non-convex penalty

$$\min_{\boldsymbol{\varSigma}\succ \mathbf{0}} \left\{ \frac{1}{2m} \| \boldsymbol{y} - \mathcal{A}(\boldsymbol{\varSigma}) \|_{\mathrm{F}}^2 - \tau \log \det \boldsymbol{\varSigma} + \sum_{i,j} p_{\lambda}(|\boldsymbol{\varSigma}_{ij}|) \right\}.$$

• Assumptions on the non-convex penalty function  $p_{\lambda}(\cdot)$ :

(i) p<sub>λ</sub>(t) is non-decreasing on [0, +∞) with p<sub>λ</sub>(0) = 0 and is differentiable on (0, +∞);
(ii) 0 ≤ p'<sub>λ</sub>(t<sub>1</sub>) ≤ p'<sub>λ</sub>(t<sub>2</sub>) ≤ λ for all t<sub>1</sub> ≥ t<sub>2</sub> ≥ 0 and lim<sub>t→0</sub> p'<sub>λ</sub>(t) = λ;
(iii) There exists an α > 0 such that p'<sub>λ</sub>(t) = 0 for t ≥ αλ.

Prototypical examples of the non-convex penalty function  $p_{\lambda}(\cdot)$ :

 $\int \lambda$ , for  $0 < t \leq \lambda$ ,

- for some constant  $\kappa \geq 1$ .
- $\blacksquare$  There exist universal constants  $\alpha$  and  $\mu$  such that

$$\|\boldsymbol{\Sigma}_{\mathcal{S}}^{\star}\|_{\min} = \min_{(i,j)\in\mathcal{S}} |\boldsymbol{\Sigma}_{ij}^{\star}| \ge (\alpha + \mu) \lambda,$$
  
where  $\mu \in (0, \alpha)$  satisfies  $p_{\lambda}'(\mu\lambda) \ge \frac{\lambda}{2}$ .

Shanghai

**Theorem 1**: Let 
$$S = \{(i,j) \mid \Sigma_{ij}^* \neq 0\}$$
 and  $|S| = s$ .  
Define  $f(\Sigma) = \frac{1}{2m} ||\mathbf{y} - \mathcal{A}(\Sigma)||_{\mathrm{F}}^2 - \tau \log \det \Sigma$ .

Under some standard assumptions, the  $\varepsilon$ -optimal solution  $\widetilde{\Sigma}^{(k)}$   $(1 \leq k \leq K)$  satisfies the following contraction property:

$$\left\|\widetilde{\boldsymbol{\Sigma}}^{(k)} - \boldsymbol{\Sigma}^{\star}\right\|_{\mathrm{F}} \leq \frac{1}{\rho} \left( \underbrace{\left\| (\nabla f(\boldsymbol{\Sigma}^{\star}))_{\mathcal{S}} \right\|_{\mathrm{F}}}_{\text{oracle rate}} + \underbrace{\varepsilon \sqrt{s}}_{\text{optimization error}} \right) + \underbrace{\delta \left\| \widetilde{\boldsymbol{\Sigma}}^{(k-1)} - \boldsymbol{\Sigma}^{\star} \right\|_{\mathrm{F}}}_{\text{contraction}},$$

where  $\delta \in (0, 1)$  is the contraction parameter.

**Theorem 2**: Under some standard assumptions, let  $\boldsymbol{x}$  be a sub-Gaussian random vector with mean zero and covariance  $\boldsymbol{\Sigma}^*$  and  $\{\boldsymbol{x}_i\}_{i=1}^n$  be a collection of i.i.d. samples drawn from  $\boldsymbol{x}$ , if  $\lambda \asymp \sqrt{\frac{\log d}{mn}}$ ,  $\tau \lesssim \sqrt{\frac{1}{mn}} \left\| (\boldsymbol{\Sigma}^*)^{-1} \right\|_{\max}^{-1}$ ,  $\varepsilon \lesssim \sqrt{\frac{1}{mn}}$ , and  $K \gtrsim \log(\lambda \sqrt{mn}) \gtrsim \log \log d$ , then the  $\varepsilon$ -optimal solution  $\boldsymbol{\Sigma}^{(K)}$  satisfies

SACD: 
$$p'_{\lambda}(t) = \begin{cases} \frac{b\lambda-t}{b-1}, & \text{for } \lambda \leq t \leq b\lambda, \\ 0, & \text{for } t \geq b\lambda, \end{cases}$$

where b > 2 is an additional tuning parameter.

MCP: 
$$p_{\lambda}(t) = \operatorname{sign}(t) \lambda \cdot \int_{0}^{|t|} \left(1 - \frac{z}{\lambda b}\right)_{+} dz$$

for some b > 0.

#### Algorithm

Algorithm 1: Majorization-Minimization Based Multistage Convex Relaxation

for 
$$k = 1, 2, ..., K$$
 do  
update  $\Lambda_{ij}^{(k-1)} = p'_{\lambda}(|\widetilde{\Sigma}_{ij}^{(k-1)}|);$   
obtain  $\widetilde{\Sigma}^{(k)}$  by solving  

$$\min_{\Sigma \succ 0} \left\{ \frac{1}{2m} \| \mathbf{y} - \mathcal{A}(\Sigma) \|_{\mathrm{F}}^{2} - \tau \log \det \Sigma + \sum_{i,j} p_{\lambda}(|\Sigma_{ij}|) \right\},$$

$$k = k + 1;$$
  
and for.

Algorithm 2: Proximal Newton Algorithm With Back-tracking Line Search

Input  $\Sigma^{k-1}$ ,  $\Lambda^k$ ,  $\varepsilon$ Initialize t = 0,  $\Sigma_t = \Sigma^{k-1}$ ,  $\mu = 0.8$ ,  $\alpha = 0.3$   $\left\|\widetilde{\boldsymbol{\Sigma}}^{(K)}-\boldsymbol{\Sigma}^{\star}\right\|_{F}\lesssim\sqrt{\frac{s}{mn}}$ 

with high probability.

#### Numerical Simulations

We use the MCP penalty, defined as

 $p_{\lambda}(t) = \operatorname{sign}(t)\lambda \cdot \int_{0}^{|t|} \left(1 - \frac{z}{\lambda b}\right)_{+} dz,$ 

with b = 2 across all trials. The regularization parameters  $\tau$  and  $\lambda$  are selected via five-fold cross-validation. The ground-truth covariance matrix  $\Sigma^*$  is generated using the built-in sprandsym function in MATLAB with *s* nonzero entries. We draw n = 50 independent samples from the multivariate normal distribution  $\mathcal{N}(0, \Sigma^*)$ , and the noise variables  $\eta_i$  are sampled from a sub-exponential distribution with scale parameter  $\gamma$ , i.e.,  $\eta_i \sim \gamma \cdot \mathcal{N}(0, 1)$ . To evaluate recovery performance, we measure the success probability as visualized in the color-coded matrix in Fig. (a). To reduce the impact of limited sample size, we directly apply sketching to the true covariance matrix  $\Sigma^*$ . A recovery is considered successful if the relative Frobenius error satisfies

$$rac{\left\| \boldsymbol{\varSigma} - \boldsymbol{\varSigma}^{\star} 
ight\|_{\mathrm{F}}}{\left\| \boldsymbol{\varSigma}^{\star} 
ight\|_{\mathrm{F}}} \leq 10^{-3}.$$

Fig. (b) compares the proposed estimator with the  $\ell_1$ -norm-based method under a

Repeat

$$\begin{split} \boldsymbol{\Sigma}_{t+\frac{1}{2}} \in \arg\min_{\boldsymbol{\Sigma}\succ\boldsymbol{0}} \widetilde{f}_t(\boldsymbol{\Sigma}) + \|\boldsymbol{\Lambda}\odot\boldsymbol{\Sigma}\|_1; \\ \boldsymbol{\Delta}_t &= \boldsymbol{\Sigma}_{t+\frac{1}{2}} - \boldsymbol{\Sigma}_t; \\ \boldsymbol{\delta}_t &= \langle \nabla f(\boldsymbol{\Sigma}_t), \boldsymbol{\Delta}_t \rangle - \|\boldsymbol{\Lambda}^k \odot \boldsymbol{\Sigma}_t\|_1 + \|\boldsymbol{\Lambda}^k \odot (\boldsymbol{\Sigma}_t + \boldsymbol{\Delta}_t)\|_1; \\ \boldsymbol{\beta} &= 1, \ \boldsymbol{q} = 0; \\ \textbf{Repeat} \\ \boldsymbol{\beta} &= \mu^q, \ \boldsymbol{q} &= \boldsymbol{q} + 1; \\ \textbf{If } \boldsymbol{\Sigma}_t + \boldsymbol{\beta}\boldsymbol{\Delta}_t \leq \boldsymbol{0} \text{ then} \\ \text{ continue;} \\ \textbf{end} \\ \textbf{until } \overline{F}(\boldsymbol{\Sigma}_t + \boldsymbol{\beta}\boldsymbol{\Delta}_t) \leq \overline{F}(\boldsymbol{\Sigma}_t) + \alpha \boldsymbol{\beta} \boldsymbol{\delta}_t \\ \boldsymbol{\Sigma}_{t+1} &= \boldsymbol{\Sigma}_t + \boldsymbol{\beta}\boldsymbol{\Delta}_t; \\ \boldsymbol{t} &= t + 1; \\ \textbf{until } \max_{i,j} \left| (\nabla f(\boldsymbol{\Sigma}_{t+1}) + \boldsymbol{\Lambda}^k \odot \boldsymbol{\Xi}^k)_{ij} \right| \leq \varepsilon \\ \textbf{Output: } \boldsymbol{\Sigma}^K &= \boldsymbol{\Sigma}_{t+1} \end{split}$$

consistent noise level  $\gamma = 10^{-1}$ . As the number of measurements increases, the recovery error decreases, and our method consistently outperforms the  $\ell_1$ -based estimator.





Wenbin Wang and Ziping Zhao