

# Condition Number Shrinkage by Joint Distributionally Robust Covariance-Precision Estimation

Renjie Chen, Viet Anh Nguyen, Huifu Xu  
Speaker: Renjie Chen

Chinese University of Hong Kong, Department of Systems Engineering and Engineering Management

August, 2025 @ ICSP

# Sample Covariance Matrix and Sample precision Matrix

- Let  $\xi \in \mathbb{R}^p$  be a random vector with mean zero and **covariance matrix**  $\Sigma_0$ .
- The **precision matrix** (also called inverse covariance matrix) of  $\xi$  is  $\Sigma_0^{-1}$ .
- **Sample Covariance Matrix:**

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n \xi_i \xi_i^T$$

where  $\xi_1, \dots, \xi_n$  are i.i.d. samples of  $\xi$ .

- If  $\hat{\Sigma}$  is non-singular, a standard estimator of the precision matrix is  $\hat{\Sigma}^{-1}$ .
- However, in high-dimensional and data-deficient case,  $\hat{\Sigma}$  could be ill-conditioned or even singular, which makes the computation of  $\hat{\Sigma}^{-1}$  unstable or not feasible.

## Ill-conditioning of the Sample Covariance Matrix: Overestimation and Underestimation of Eigenvalues

- The **condition number** of the sample covariance matrix  $\hat{\Sigma}$  is defined as:

$$\kappa(\hat{\Sigma}) = \frac{\lambda_{\max}(\hat{\Sigma})}{\lambda_{\min}(\hat{\Sigma})}$$

where  $\lambda_{\max}(\hat{\Sigma})$  and  $\lambda_{\min}(\hat{\Sigma})$  are the largest and smallest eigenvalues of  $\hat{\Sigma}$ .

## Ill-conditioning of the Sample Covariance Matrix: Overestimation and Underestimation of Eigenvalues

- It is well-known that  $\lambda_{\max}(\cdot)$  is convex, the by Jensen's inequality,

$$\mathbb{E}[\lambda_{\max}(\hat{\Sigma})] \geq \lambda_{\max}(\mathbb{E}[\hat{\Sigma}]) = \lambda_{\max}(\Sigma_0).$$

The largest eigenvalue tends to be **overestimated**.

- $\lambda_{\min}(\cdot)$  is concave over  $\mathbb{S}_{++}^p$

$$\mathbb{E}[\lambda_{\min}(\hat{\Sigma})] \leq \lambda_{\min}(\mathbb{E}[\hat{\Sigma}]) = \lambda_{\min}(\Sigma_0).$$

The smallest eigenvalue tends to be **underestimated**.

- **Overestimated Condition Number:**

$$\kappa(\hat{\Sigma}) \text{ is often much larger than } \kappa(\Sigma_0).$$

# Ledoit Wolf's Linear shrinkage estimator<sup>1</sup>

- The linear shrinkage estimator

$$S = \rho \nu I + (1 - \rho) \hat{\Sigma},$$

where  $\rho$  is the shrinkage intensity and  $\nu I$  is the shrinkage target.

- The optimal parameter is chosen by minimizing the mean squared error of  $S$ :

$$\min_{0 \leq \rho \leq 1, \nu \geq 0} \mathbb{E} \left[ \|S - \Sigma_0\|^2 \right], \text{ s.t } S = \rho \nu I + (1 - \rho) \hat{\Sigma}$$

- The optimal shrinkage target is  $\nu^* = \frac{\text{tr}(\Sigma_0)}{p}$ , and the optimal intensity is

$$\rho^* = \frac{\mathbb{E}[\|\hat{\Sigma} - \Sigma_0\|^2]}{\mathbb{E}[\|\hat{\Sigma} - \nu^* I\|^2]}.$$

- **Backwards:** Linear structure is not expressive, and shrink too much (too conservative).

---

<sup>1</sup>Ledoit and Wolf 2004, "A well-conditioned estimator for large-dimensional covariance matrices".

# Distributionally Robust Covariance Estimation<sup>2</sup>

- Let  $\hat{\mathbb{P}}_n$  be the empirical distribution of i.i.d. samples  $\xi_1, \dots, \xi_n$ .
- The sample covariance matrix is  $\hat{\Sigma} = \mathbb{E}_{\hat{\mathbb{P}}_n}[\xi\xi^T]$ .
- Yue et al. proposed to view  $\hat{\Sigma}$  as the unique solution to the minimization problem

$$\hat{\Sigma} = \arg \min_{\Sigma \in \mathbb{S}_+^p} \|\Sigma\|_F^2 - 2\langle \Sigma, \hat{\Sigma} \rangle = \arg \min_{\Sigma \in \mathbb{S}_+^p} \underbrace{\|\Sigma\|_F^2 - 2\mathbb{E}_{\hat{\mathbb{P}}_n}[\langle \Sigma, \xi\xi^T \rangle]}_{\triangleq \text{FrobeniusLoss}(\Sigma, \hat{\mathbb{P}}_n)}.$$

---

<sup>2</sup>Yue et al. 2024, “A Geometric Unification of Distributionally Robust Covariance Estimators: Shrinking the Spectrum by Inflating the Ambiguity Set”.

# Distributionally Robust Covariance Estimation

- A moment-based ambiguity set centered at the nominal distribution  $\hat{\mathbb{P}}_n$  with radius  $\varepsilon > 0$ :

$$\mathcal{P}_\varepsilon(\hat{\mathbb{P}}_n) \triangleq \left\{ \mathbb{Q} : \mathbb{E}_{\mathbb{Q}}[\xi] = 0, \mathbb{E}_{\mathbb{Q}}[\xi\xi^\top] = S, D\left(S, \mathbb{E}_{\hat{\mathbb{P}}_n}[\xi\xi^\top]\right) \leq \varepsilon \right\},$$

where  $D$  is a divergence in the space of positive semidefinite matrices to measure the discrepancy between two matrices.

- The distributionally robust covariance estimation model:

$$\min_{\Sigma \in \mathbb{S}_{++}^p} \max_{\mathbb{Q} \in \mathcal{P}_\varepsilon(\hat{\mathbb{P}}_n)} \|\Sigma\|_F^2 - 2\mathbb{E}_{\hat{\mathbb{P}}_n}[\langle \Sigma, \xi\xi^\top \rangle]$$

# Distributionally Robust Covariance Estimation

- Under certain assumption on choice of divergence  $D$ , the distributionally robust covariance estimation model reduces to

$$\min_{D(S, \hat{\Sigma}) \leq \varepsilon} \|S\|_F^2.$$

- It shrinks all the eigenvalues towards 0, i.e., the shrinkage target of covariance matrix is zero matrix.
- The shrinkage intensity is decided by  $\varepsilon$ . The larger  $\varepsilon$ , the more shrinkage.
- **Backwards:** (1) shrinkage target 0 is not well-conditioned; (2) it underestimates  $\lambda_{\min}$  even worse.



# Distributionally Robust Precision Estimation<sup>3</sup>

- If  $\hat{\Sigma}$  is non-singular, the sample precision matrix estimator can be viewed as the optimal solution of:

$$\hat{X} = \arg \min_{X \in \mathbb{S}_{++}^p} -\log \det X + \langle X, \hat{\Sigma} \rangle = \arg \min_{X \in \mathbb{S}_{++}^p} \underbrace{-\log \det X + \mathbb{E}_{\hat{\mathbb{P}}_n} [\langle X, \xi \xi^T \rangle]}_{\triangleq \text{SteinLoss}(X, \mathbb{P}_n)}.$$

- The optimal solution is  $\hat{X} = \hat{\Sigma}^{-1}$ .

---

<sup>3</sup>Nguyen, Kuhn, and Mohajerin Esfahani 2022, “*Distributionally robust inverse covariance estimation: The Wasserstein shrinkage estimator*”.

# Distributionally Robust Precision Estimation

- Nguyen et al. consider the distributionally robust precision estimation model:

$$\min_{X \in \mathbb{S}_{++}^p} -\log \det X + \max_{Q \in \mathcal{P}_\varepsilon(\hat{\mathbb{P}}_n)} \mathbb{E}_Q[\langle X, \xi \xi^\top \rangle].$$

- The resulting distributionally robust precision matrix estimator is also a shrinkage estimator.
- It shrinks all the eigenvalues of precision matrix towards 0, i.e., the shrinkage target of precision matrix is zero matrix.
- It is equivalent to shrink the eigenvalues of covariance matrix towards  $+\infty$ .
- The shrinkage intensity is also control by radius  $\varepsilon$ .
- **Backwards:** It overestimates  $\lambda_{\max}$  even worse.

# Distributionally Robust Covariance-precision Estimation

- Consider the joint covariance-precision matrix estimation model:

$$\begin{aligned} \min_{\Sigma, X} \quad & \text{SteinLoss}(X, \hat{\mathbb{P}}_n) + \frac{\tau}{2} \text{FrobeniusLoss}(\Sigma, \hat{\mathbb{P}}_n) \\ \text{s.t.} \quad & \Sigma \in \mathbb{S}_+^p, \quad X \in \mathbb{S}_{++}^p, \quad X\Sigma = I. \end{aligned}$$

When  $\hat{\Sigma}$  is non-singular, the optimal solution is  $(\Sigma^* = \hat{\Sigma}, X^* = \hat{\Sigma}^{-1})$ .

# Distributionally Robust Covariance-precision Estimation

- Consider the moment-based ambiguity set centered at the nominal distribution  $\hat{\mathbb{P}}_n$  with radius  $\varepsilon > 0$  again:

$$\mathcal{P}_\varepsilon(\hat{\mathbb{P}}_n) \triangleq \left\{ \mathbb{Q} : \mathbb{E}_{\mathbb{Q}}[\xi] = 0, \mathbb{E}_{\mathbb{Q}}[\xi\xi^\top] = S, D\left(S, \mathbb{E}_{\hat{\mathbb{P}}_n}[\xi\xi^\top]\right) \leq \varepsilon \right\}.$$

- Define the pair of distributionally robust covariance-precision matrix estimator  $(\Sigma^*, X^*)$  by

$$(\Sigma^*, X^*) \triangleq \arg \min_{\substack{\Sigma, X \in \mathbb{S}_{++}^p \\ X\Sigma = I}} \max_{\mathbb{Q} \in \mathcal{P}_\varepsilon(\hat{\mathbb{P}}_n)} -\log \det X + \mathbb{E}_{\mathbb{Q}}[\langle X, \xi\xi^\top \rangle] + \frac{1}{2}\tau \left( \|\Sigma\|_F^2 - 2\mathbb{E}_{\mathbb{Q}}[\langle \Sigma, \xi\xi^\top \rangle] \right) \quad (\text{DRO})$$

# Distributionally Robust Covariance-precision Estimation

- Note that the objective function depends on  $\mathbb{Q}$  only through term  $\mathbb{E}_{\mathbb{Q}}[\xi\xi^T]$ .
- The DRO model reduces to the following robust optimization (RO) problem:

$$(\Sigma^*, X^*) = \arg \min_{\substack{\Sigma, X \in \mathbb{S}_{++}^p \\ X\Sigma = I}} \max_{S \in \mathbb{S}_+^p : D(S, \hat{\Sigma}) \leq \varepsilon} \left\{ -\log \det X + \langle X, S \rangle + \frac{1}{2} \tau (\|\Sigma\|_F^2 - 2\langle \Sigma, S \rangle) \right\}. \quad (\text{RO})$$

- Question:
- The constraint of (RO) is not convex. How can we solve it?
- Do we get a condition-number-shrunk optimal solution pair?
- If yes, what is the shrinkage target and intensity?
- How to choose parameter  $\tau$  and  $\varepsilon$ ?

# Solvable max-min problem

- Interchange the order of min-max of (RO) gives

$$\max_{D(S, \hat{\Sigma}) \leq \varepsilon} \left( \min_{\substack{X, \Sigma \\ \text{s.t.} \\ X \in \mathbb{S}_{++}^p, \Sigma \in \mathbb{S}_+^p, X\Sigma = I}} \langle \hat{\Sigma}, X \rangle - \log \det X + \frac{1}{2} \tau (\|\Sigma\|_F^2 - 2\langle \Sigma, S \rangle) \right).$$

- The inner min problem is solved by  $\Sigma^* = (X^*)^{-1} = S$ , when  $S$  is non-singular. Then it becomes

$$\max_{S \in \mathbb{S}_+^p : D(S, \hat{\Sigma}) \leq \varepsilon} \log \det S - \frac{1}{2} \tau \|S\|_F^2. \quad (\text{P-Mat})$$

- The following saddle-point theorem shows that the interchanging of min-max keeps the problem nature, i.e., optimal solution of (P-Mat) solves (RO).

# Saddle Point Theorem

## Theorem 1

Suppose that

1.  $D(\cdot, \hat{\Sigma})$  is convex, differentiable and  $\{S \in \mathbb{S}_{++}^p : D(S, \hat{\Sigma}) \leq \varepsilon\}$  is compact,
2.  $(\hat{\Sigma}, \hat{\Sigma}) \in \text{dom}(D)$ , and
3. Slater's condition holds, i.e., there exists  $S \in \mathbb{S}_+^p$  such that  $D(S, \hat{\Sigma}) < \varepsilon$ .

Let  $\tau > 0$  and  $S^*$  be an optimal solution defined as

$$S^* = \arg \max_{S \in \mathbb{S}_+^p : D(S, \hat{\Sigma}) \leq \varepsilon} \log \det S - \frac{1}{2} \tau \|S\|_F^2. \quad (\text{P-Mat})$$

Then the tuple  $(\Sigma = S^*, X = (S^*)^{-1}, S = S^*)$  is an optimal solution of (RO).

# Unbinding Constraint Case

Now we deal with the solvability of (P-Mat).

$$S^* = \arg \max_{S \in \mathbb{S}_+^p : D(S, \hat{\Sigma}) \leq \varepsilon} \log \det S - \frac{1}{2} \tau \|S\|_F^2. \quad (\text{P-Mat})$$

We first consider a relatively straightforward case where the constraint is unbinding.

## Proposition 1

*If  $\varepsilon \geq \varepsilon_{\max} \triangleq D\left(\sqrt{\frac{1}{\tau}}I, \hat{\Sigma}\right)$ , then  $S^* = \sqrt{\frac{1}{\tau}}I$  is the optimal solution to (P-Mat).*

In the following discussion, we assume that  $0 < \varepsilon < \varepsilon_{\max}$ , and show that under some assumption on choice of divergence  $D$ , (P-Mat) can be solved in quasi-closed form.



## Assumption 1: Convex Spectral Divergence $D$ I

The divergence function  $D : \mathbb{S}_+^p \times \mathbb{S}_+^p \rightarrow \mathbb{R}_+$  is non-negative, continuous, and satisfies the identity of indiscernibles, that is, for any  $X, Y \in \text{dom}(D)$  we have  $D(X, Y) = 0$  if and only if  $X = Y$ . For any  $Y \in \text{dom}(D)$ ,  $D(\cdot, Y)$  is convex, differentiable and for any  $Y \in \text{dom}(D)$ ,  $\varepsilon > 0$  the set  $\{S \in \mathbb{S}_{++}^p : D(S, Y) \leq \varepsilon\}$  is compact. In addition,  $D$  satisfies the following structural conditions.

- (i) (Orthogonal equivariance) For any  $X, Y \in \mathbb{S}_+^p$  and orthogonal matrix  $V$ , we have  $D(X, Y) = D(VXV^T, VYV^T)$ .
- (ii) (Spectrality) There exists a function  $d : \mathbb{R}_+ \times \mathbb{R}_+ \rightarrow \mathbb{R}_+$  such that

$$D(\text{diag}(x), \text{diag}(y)) = \sum_{i=1}^p d(x_i, y_i), \forall x, y \in \mathbb{R}_+^p$$

and  $d(a, b)$  is twice continuously differentiable in  $a$  for every  $b \geq 0$ . In the following, we refer to  $d$  as the generator of  $D$ .

## Assumption 1: Convex Spectral Divergence $D$ II

(iii) (Rearrangement property) For any  $x, y \in \mathbb{R}_+^p$  and  $V \in \mathcal{O}(p)$  we have

$$D(V \operatorname{diag}(x^\uparrow) V^\top, \operatorname{diag}(y^\uparrow)) \geq D(\operatorname{diag}(x^\uparrow), \operatorname{diag}(y^\uparrow)).$$

If its left side is finite, this inequality becomes an equality if and only if  $V \operatorname{diag}(x^\uparrow) V^\top = \operatorname{diag}(x^\uparrow)$ .

**Remark:** Consider  $(aI, bI) \in \operatorname{dom}(D)$ . Then by (ii), the domain of  $d$  is

$$\operatorname{dom}(d) = \{(a, b) \in \mathbb{R}_+^2 : (aI, bI) \in \operatorname{dom}(D)\},$$

and

$$D(aI, bI) = p \times d(a, b).$$

It implies that  $d$  inherits non-negativity, continuity, identity of indiscernibles, and convexity from  $D$ . Then we conclude that  $d(b, b) = 0$ ,  $\frac{\partial d}{\partial a}(b, b) = 0$  and thus  $b$  is the unique minimizer of the function  $d(\cdot, b)$  for any  $b > 0$ .

# Reduction of (P-Mat) to a Vector Space Problem

- Let  $\hat{\Sigma} = \hat{V} \text{diag}(\hat{\lambda}_1, \dots, \hat{\lambda}_p) \hat{V}^T$  be the spectral decomposition of  $\hat{\Sigma}$ .
- Under Assumption 1, we are able to show that (P-Mat) can be reduced to a problem that optimizes over all vectors in the non-negative orthant  $\mathbb{R}_+^p$ :

$$\begin{aligned} \max_{s_i \geq 0} \quad & \sum_{i=1}^p \log s_i - \frac{1}{2} \tau \sum_{i=1}^p s_i^2 \\ \text{s.t.} \quad & \sum_{i=1}^p d(s_i, \hat{\lambda}_i) \leq \varepsilon. \end{aligned} \tag{P-Vec}$$

## Proposition 2 (Equivalence of (P-Mat) and (P-Vec))

*If Assumption 1 holds, then problem (P-Mat) is equivalent to (P-Vec) in the following sense.*

- (i) *If  $s^*$  solves (P-Vec), then  $\hat{V} \text{diag}(s^*) \hat{V}^T$  solves (P-Mat).*
- (ii) *If  $S^*$  solves (P-Mat), then the vector of its eigenvalues  $\lambda(S^*)$  solves (P-Vec).*

## Quasi-closed Form of Solution of (P-Vec)

- By the standard Lagrange dual approach, the dual of (P-Vec) is equivalent to

$$\min_{\gamma \geq 0} \max_{s_i \geq 0} \sum_{i=1}^p \log s_i - \frac{1}{2} \tau \sum_{i=1}^p s_i^2 - \gamma \left( \sum_{i=1}^p d(s_i, \hat{\lambda}_i) - \varepsilon \right).$$

- Define a solution mapping  $\varphi : \mathbb{R}_+^3 \rightarrow \mathbb{R}$  via

$$\varphi(\tau, \gamma, b) = \text{the unique solution } a^* > 0 \text{ of the equation } 0 = \frac{1}{a^*} - \tau a^* - \gamma \frac{\partial d}{\partial a}(a^*, b).$$

- For fixed  $\gamma$ , the inner maximization problem is solved by

$$s_i^* = \varphi(\tau, \gamma, \hat{\lambda}_i).$$

- By complementary slackness condition, the optimal Lagrange multiplier  $\gamma^*$  must satisfy

$$\sum_{i=1}^p d(\varphi(\tau, \gamma^*, \hat{\lambda}_i), \hat{\lambda}_i) - \varepsilon = 0.$$

# Quasi-closed Form of Solution of (P-Vec)

## Proposition 3 (Solution of (P-Vec))

If  $\sum_{i=1}^p d\left(\sqrt{\frac{1}{\tau}}, \hat{\lambda}_i\right) > \varepsilon$ , then (P-Vec) admits a unique optimal solution  $s^*$  and

$$s_i^* = \varphi(\tau, \gamma^*, \hat{\lambda}_i),$$

where  $\gamma^*$  is the unique solution of the nonlinear equation

$$\sum_{i=1}^p d(\varphi(\tau, \gamma^*, \hat{\lambda}_i), \hat{\lambda}_i) - \varepsilon = 0.$$

# Uniqueness of $\gamma^*$

- Define  $F : \mathbb{R}_+ \rightarrow \mathbb{R}$  through  $F(\gamma) = \sum_{i=1}^p d(\varphi(\tau, \gamma, \hat{\lambda}_i), \hat{\lambda}_i)$ .

## Proposition 4 (Differentiable and strictly decreasing $F$ )

*The function  $F$  is differentiable and strictly decreasing over  $\mathbb{R}_+$ . If  $\varepsilon < \varepsilon_{\max} \triangleq \sum_{i=1}^p d\left(\sqrt{\frac{1}{\tau}}, \hat{\lambda}_i\right)$ , then  $\lim_{\gamma \downarrow 0} F(\gamma) > \varepsilon$  and  $\lim_{\gamma \rightarrow \infty} F(\gamma) < \varepsilon$ .*

The proposition reveals that

- $F(\gamma) = \varepsilon$  admits a unique positive root,
- the equation can be solved efficiently by bisection or Newton's method.

# Construction of Distributionally Robust Covariance-precision Estimators

Taking all steps so far:

$(\text{DRO}) \Leftrightarrow (\text{RO}) \Leftrightarrow (\text{P-Mat}) \Leftrightarrow (\text{P-Vec}) \Rightarrow$  quasi-closed form solution.

- Recall that  $\hat{\Sigma} = \hat{V} \text{diag}(\hat{\lambda}_1, \dots, \hat{\lambda}_p) \hat{V}^T$  is the spectral decomposition of  $\hat{\Sigma}$ .
- The theorem reveals how to construct the distributionally robust covariance-precision estimator.

## Theorem 2 (Construction of covariance estimator)

*Under certain assumptions, the optimal solution of (RO) is given by  $(\Sigma^* = \hat{V} \Phi(\tau, \gamma^*, \hat{\lambda}) \hat{V}^T, X = \Sigma^{*-1})$ , where*

$$\Phi(\tau, \gamma^*, \hat{\lambda}) \triangleq \text{diag}(\varphi(\tau, \gamma^*, \hat{\lambda}_1), \dots, \varphi(\tau, \gamma^*, \hat{\lambda}_p)),$$

*where  $\gamma^*$  is the unique positive root of the equation  $\sum_{i=1}^p d(\varphi(\tau, \gamma^*, \hat{\lambda}_i), \hat{\lambda}_i) - \varepsilon = 0$ .*

# Proof Structure of Theorem 2

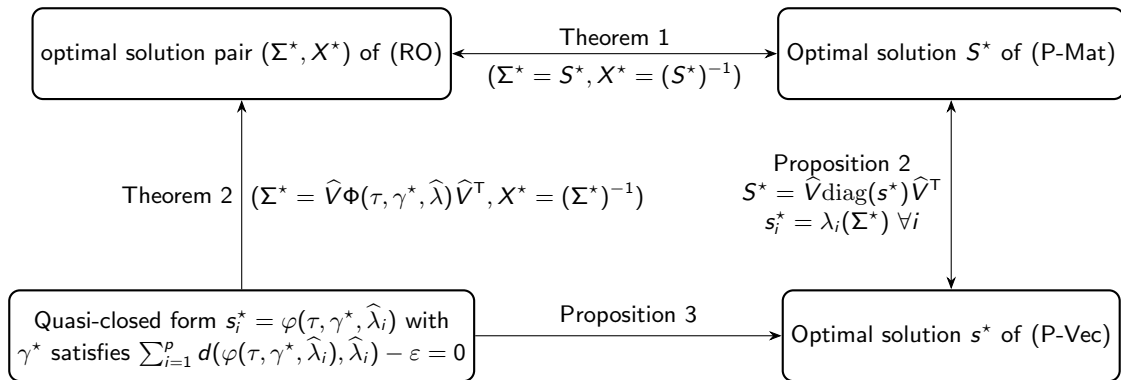


Figure 1: Structure of the proof of Theorem 2.



- Define the covariance matrix estimator by

$$\Sigma^*(\tau, \varepsilon) \triangleq \widehat{V} \text{diag}(\varphi(\tau, \gamma_{\widehat{\lambda}}(\varepsilon), \widehat{\lambda}_1), \dots, \varphi(\tau, \gamma_{\widehat{\lambda}}(\varepsilon), \widehat{\lambda}_p)) \widehat{V}^T, \varepsilon \in (0, \varepsilon_{\max}], \quad (1)$$

where  $\gamma_{\widehat{\lambda}}(\varepsilon)$  denote the unique root of  $\sum_{i=1}^p d(\varphi(\tau, \gamma^*, \widehat{\lambda}_i), \widehat{\lambda}_i) - \varepsilon = 0$ .

# Nonlinear Shrinkage and Improved Condition Number

## Theorem 3

If Assumptions hold, then for any  $\tau > 0$ ,  $\Sigma^*(\tau, \varepsilon)$  is continuous over  $\varepsilon \in (0, \varepsilon_{\max}]$  and  $\lim_{\varepsilon \rightarrow 0} \Sigma^*(\tau, \varepsilon) = \hat{\Sigma}$ ,  $\Sigma^*(\tau, \varepsilon_{\max}) = \sqrt{\frac{1}{\tau}} I$ . Moreover, if  $\lambda_{\min}(\hat{\Sigma}) < \sqrt{\frac{1}{\tau}} < \lambda_{\max}(\hat{\Sigma})$ , then we have

- (i)  $\lambda_{\min}(\hat{\Sigma}) < \lambda_{\min}(\Sigma^*(\varepsilon)) < \sqrt{\frac{1}{\tau}} < \lambda_{\max}(\Sigma^*(\varepsilon)) < \lambda_{\max}(\hat{\Sigma})$  for any  $\varepsilon \in (0, \varepsilon_{\max}]$ ,
- (ii)  $\kappa(\Sigma^*(\varepsilon))$  is strictly decreasing on  $\varepsilon$ , and thus  $\kappa(\Sigma^*(\varepsilon)) < \kappa(\hat{\Sigma})$  for any  $\varepsilon \in (0, \varepsilon_{\max}]$ .

- The proposed covariance-precision matrix estimator can be interpreted as a **nonlinear shrinkage estimator** with
  - $\sqrt{\frac{1}{\tau}} I$  is the shrinkage target,
  - $\varepsilon$  controls the shrinkage intensity.
- The nonlinear shrinkage improves the condition number strictly.

# Convex Spectral Divergences

## Theorem 4

*All divergences in Table 1 satisfy the Convex Spectral Divergence Assumption.*

Divergence function	$D(\Sigma, \hat{\Sigma})$	$d(a, b)$	Domain of $D$
Kullback-Leibler	$\frac{1}{2} \left( \text{tr}(\hat{\Sigma}^{-1}\Sigma) - p + \log \det(\hat{\Sigma}\Sigma^{-1}) \right)$	$\frac{1}{2} \left( \frac{a}{b} - 1 - \log \frac{a}{b} \right)$	$\mathbb{S}_{++}^p \times \mathbb{S}_{++}^p$
Wasserstein	$\text{tr} \left( \Sigma + \hat{\Sigma} - 2(\Sigma^{\frac{1}{2}}\hat{\Sigma}\Sigma^{\frac{1}{2}})^{\frac{1}{2}} \right)$	$a + b - 2\sqrt{ab}$	$\mathbb{S}_+^p \times \mathbb{S}_+^p$
Symmetrized Stein	$\frac{1}{2} \left( \text{tr}(\hat{\Sigma}^{-1}\Sigma) + \text{tr}(\Sigma^{-1}\hat{\Sigma}) - 2p \right)$	$\frac{1}{2} \left( \frac{b}{a} + \frac{a}{b} - 2 \right)$	$\mathbb{S}_{++}^p \times \mathbb{S}_{++}^p$
Squared Frobenius	$\text{tr} \left( (\Sigma - \hat{\Sigma})^2 \right)$	$(a - b)^2$	$\mathbb{S}_+^p \times \mathbb{S}_+^p$
Weighted Frobenius	$\text{tr} \left( (\Sigma - \hat{\Sigma})\hat{\Sigma}^{-1} \right)$	$\frac{(a-b)^2}{b}$	$\mathbb{S}_+^p \times \mathbb{S}_{++}^p$

Table 1: Divergence functions and their generators.

# Solution Mapping and Upper Bound of Dual Variable

Taking  $D$  as any divergence in Table 1, we can construct the optimal solution of (RO) using Theorem 2. The key steps are

- to find out the solution mapping  $\varphi$  as function of  $\tau$ ,  $\gamma$  and  $b$ ,
- to solve the equation  $\sum_{i=1}^P d(\varphi(\tau, \gamma^*, \hat{\lambda}_i), \hat{\lambda}_i) - \varepsilon = 0$  to find out dual variable  $\gamma^*$ .

# Solution Mapping and Upper Bound of Dual Variable

Divergence function	Solution mapping $\varphi(\tau, \gamma, b)$	Upper bound of $\gamma^*$
Kullback-Leibler	$\frac{-\gamma + \sqrt{\gamma^2 + 8\tau(2+\gamma)b^2}}{4\tau b}$	$\max \left\{ \frac{2p}{\varepsilon}, \frac{2\tau\hat{\lambda}_p^2 + 1}{e^{\varepsilon/p} - 1} \right\}$
Wasserstein	unique positive root $a$ of $\tau a^2 + \gamma a - \gamma\sqrt{b}\sqrt{a} - 1 = 0$	$\sqrt{\max \left\{ \frac{p(1-\tau\hat{\lambda}_1^2)^2}{\varepsilon\hat{\lambda}_1}, \frac{p(1-\tau\hat{\lambda}_p^2)^2}{\varepsilon\hat{\lambda}_1} \right\}}$
Symmetrized Stein	unique positive root $a$ of $2\tau ba^3 + \gamma a^2 - 2ba - \gamma b^2 = 0$	$\sqrt{\max \left\{ \frac{p(1-\tau\hat{\lambda}_p^2)^2}{4\varepsilon\hat{\lambda}_1^4}, \frac{p(1-\tau\hat{\lambda}_1^2)^2}{4\varepsilon\hat{\lambda}_1^4} \right\}}$
Squared Frobenius	$\frac{\gamma b}{\tau + 2\gamma} + \frac{\sqrt{\gamma^2 b^2 + \tau + 2\gamma}}{\tau + 2\gamma}$	$\sqrt{\max \left\{ \frac{p(1-\tau\hat{\lambda}_1^2)^2}{4\varepsilon\hat{\lambda}_1^2}, \frac{p(1-\tau\hat{\lambda}_p^2)^2}{4\varepsilon\hat{\lambda}_1^2} \right\}}$
Weighted Frobenius	$\frac{\gamma b}{\tau b + 2\gamma} + \frac{\sqrt{\gamma^2 b^2 + b(\tau b + 2\gamma)}}{\tau b + 2\gamma}$	$\sqrt{\max \left\{ \frac{p\hat{\lambda}_p(1-\tau\hat{\lambda}_1^2)^2}{4\varepsilon\hat{\lambda}_1^2}, \frac{p\hat{\lambda}_p(1-\tau\hat{\lambda}_p^2)^2}{4\varepsilon\hat{\lambda}_1^2} \right\}}$

**Table 2:** Solution mapping and upper bound of dual variable.  $\hat{\lambda}_1$  and  $\hat{\lambda}_p$  are the smallest and largest eigenvalues of  $\hat{\Sigma}$ .

# Parameter Tuning: Minimizing MSE not Working

- Our estimator  $\Sigma^*(\tau, \varepsilon)$  is parameterized by  $\tau$  and  $\varepsilon$ , where  $\tau$  determine the shrinkage target and  $\varepsilon$  is the shrinkage intensity.
- In Ledoit-Wolf's linear shrinkage estimator, they minimize MSE of the proposed estimator to find optimal parameters.
- If we borrow the same idea , we could solve

$$\min_{\tau \in \mathbb{R}_{++}, \varepsilon \in \mathbb{R}_+} \mathbb{E}[\|\Sigma^*(\tau, \varepsilon) - \Sigma_0\|_F^2]. \quad (2)$$

- However, the dependence of  $\Sigma^*(\tau, \varepsilon)$  on  $\tau$  and  $\varepsilon$  is nonlinear. It makes (2) hard to analyze.
- From another perspective, we hope the optimal solution of (P-Mat) is close to  $\Sigma_0$ .

# Optimal Shrinkage Target

- In the first step, we hope the shrinkage target is near to  $\Sigma_0$ .
- Recall that the objective of (P-Mat) is

$$\ell_\tau(S) \triangleq \log \det S - \frac{1}{2}\tau \|S\|_F^2,$$

and the shrinkage target  $\sqrt{\frac{1}{\tau}}I$  is the optimal solution of

$$\min_{S \in \mathbb{S}_{++}^p} \ell_\tau(S).$$

- We choose  $\tau$  such that  $\Sigma_0$  can nearly solve  $\min_{S \in \mathbb{S}_{++}^p} \ell_\tau(S)$  (and thus close to  $\sqrt{\frac{1}{\tau}}I$ ).
- We minimize the gradient norm of  $\ell_\tau(S)$  evaluated at  $S = \Sigma_0$ :

$$\min_{\tau \in \mathbb{R}_{++}} \|\nabla_S \ell_\tau(S)|_{S=\Sigma_0}\|_F = \min_{\tau \in \mathbb{R}_{++}} \|\Sigma_0^{-1} - \tau \Sigma_0\|_F.$$

# Optimal Shrinkage Target

- The above problem is solved by  $\tau = \tau^* \triangleq \frac{p}{\|\Sigma_0\|_F^2}$ , and the induced shrinkage target is

$$S = \frac{\|\Sigma_0\|_F}{\sqrt{p}} I.$$

- Note that  $\left\| \frac{\|\Sigma_0\|_F}{\sqrt{p}} I \right\|_F = \|\Sigma_0\|_F$ , and thus the target can be viewed as a uniform prior with the same scale as  $\Sigma_0$ , i.e., the total scale is distributed evenly to each dimension.



# Optimal Shrinkage Intensity I

- With the optimal target chosen as above, we tune the radius  $\varepsilon$  so that the true covariance  $\Sigma_0$  is near the optimal solution of

$$S_{\tau^*}^* \triangleq \arg \max_{S \in \mathbb{S}_+^p : D(S, \hat{\Sigma}_n) \leq \varepsilon} \log \det S - \frac{1}{2} \tau^* \|S\|_F^2. \quad (\text{P-Mat-}\tau^*)$$

- When the constraint  $D(S, \hat{\Sigma}_n) \leq \varepsilon$  is active, the first-order optimality condition of (P-Mat- $\tau^*$ ) dictates that  $S_{\tau^*}^*$  and the dual variable  $\gamma_{\tau^*}^*$  satisfies

$$(S_{\tau^*}^*)^{-1} - \tau^* S_{\tau^*}^* - \gamma_{\tau^*}^* D'(S_{\tau^*}^*, \hat{\Sigma}_n) = 0, \quad (3)$$

$$D(S_{\tau^*}^*, \hat{\Sigma}_n) - \varepsilon = 0, \quad (4)$$

$$\gamma_{\tau^*}^* > 0, S_{\tau^*}^* \in \mathbb{S}_+^p. \quad (5)$$

## Optimal Shrinkage Intensity II

- Recall that  $S = \Sigma^*(\tau^*, \varepsilon)$  uniquely solves (3)-(5). By the construction of  $\Sigma_n^*(\tau^*, \varepsilon)$ , equation (4) becomes

$$\sum_{i=1}^p d(\varphi(\tau^*, \gamma, \hat{\lambda}_i), \hat{\lambda}_i) - \varepsilon = 0.$$

- Then the optimal solution of (P-Mat- $\tau^*$ ) solves

$$S^{-1} - \tau^* S - \gamma_{\hat{\lambda}}(\varepsilon) D'(S, \hat{\Sigma}_n) = 0. \quad (6)$$

We call the left hand side of (6) as the *extended* gradient of problem (P-Mat- $\tau^*$ ).

- To measure the optimality of  $S$  on average, we define the average extended gradient by

$$g(S) \triangleq S^{-1} - \tau^* S - \gamma_{\lambda}(\varepsilon) \mathbb{E}_n[D'(S, \hat{\Sigma}_n)].$$

## Optimal Shrinkage Intensity III

- Then, again, to choose  $\varepsilon$  so that  $\Sigma_0$  can nearly solve (P-Mat- $\tau^*$ ) on average and thus be close to  $\Sigma_n^*(\tau^*, \varepsilon)$ , we minimize the norm of average extended gradient evaluated at  $S = \Sigma_0$ :

$$\varepsilon_n^* = \arg \min_{\varepsilon > 0} \|g(\Sigma_0)\|_F^2 = \arg \min_{\varepsilon > 0} \|\Sigma_0^{-1} - \tau^* \Sigma_0 - \gamma_\lambda(\varepsilon) \mathbb{E}_n D'(\Sigma_0, \hat{\Sigma}_n)\|_F^2. \quad (7)$$

- Though the optimal intensity defined by (7) is not achievable, since we have no access to precise evaluation of  $\Sigma_0$  and  $\mathbb{E}_n D'(\Sigma_0, \hat{\Sigma}_n)$ , it is possible for us to characterize the limit behavior of  $\varepsilon_n^*$ .

## Assumption 1 (Locally-quadratic divergence)

*For any  $b > 0$ , there exists positive constant  $C_{b,d}$  that depends on the divergence  $d$  and the limit point  $b$  such that*

$$\lim_{a \rightarrow b} \frac{d(a, b)}{(a - b)^2} = C_{b,d}.$$

- All the divergences in Table 1 satisfy Assumption 1.

## Assumption 2 (Non-degenerating divergence gradient)

*For all integers  $n$  sufficiently large, the derivative of divergence  $D$  satisfies  $\mathbb{E}_n[D'(\Sigma_0, \hat{\Sigma}_n)] \neq 0$  and there exists strictly positive  $C_1$  and  $C_2$  such that*

$$\lim_{n \rightarrow \infty} n \|\mathbb{E}_n[D'(\Sigma_0, \hat{\Sigma}_n)]\|_F = C_1$$

*and*

$$\lim_{n \rightarrow \infty} n \left\langle \Sigma_0^{-1} - \tau^* \Sigma_0, \mathbb{E}_n[D'(\Sigma_0, \hat{\Sigma}_n)] \right\rangle = C_2.$$

- Frobenius divergence does not satisfy the assumption, since  $D(\Sigma_0, \Sigma_n) = \|\Sigma_0 - \Sigma_n\|_F^2$ , and

$$\mathbb{E}[D'(\Sigma_0, \Sigma_n)] = \mathbb{E}[2(\Sigma_0 - \Sigma_n)] = 0.$$

# $1/n^2$ -order optimal radius

## Theorem 5 ( $1/n^2$ -order optimal radius)

*Under the above assumptions, the optimal intensity  $\varepsilon_n^*$  defined as (7) satisfies*

$$\lim_{n \rightarrow \infty} n^2 \varepsilon_n^* = \varepsilon^*,$$

*where the limit  $\varepsilon^*$  depends on divergence  $D$  and the underlying data-generating distribution.*

## Corollary: Optimal Shrinkage Intensity of Convex Divergences I

Let the optimal radius  $\varepsilon_n^*$  be defined as (7). Then under Assumptions, we have

(i) if  $D$  is taken as Kullback-Leibler divergence, then

$$\lim_{n \rightarrow \infty} n^2 \varepsilon_n^* = \frac{(p+1)^2 \|\Sigma_0^{-1}\|_F^4}{16 \left( \|\Sigma_0^{-1}\|_F^2 - \frac{p^2}{\|\Sigma_0\|_F^2} \right)^2} \sum_{i=1}^p (1 - \tau^* \lambda_i^2)^2, \quad (8)$$

(ii) if  $D$  is taken as symmetrized Stein divergence, then

$$\lim_{n \rightarrow \infty} n^2 \varepsilon_n^* = \frac{(p+1)^2 \|\Sigma_0^{-1}\|_F^4}{32 \left( \|\Sigma_0^{-1}\|_F^2 - \frac{p^2}{\|\Sigma_0\|_F^2} \right)^2} \sum_{i=1}^p (1 - \tau^* \lambda_i^2)^2, \quad (9)$$

## Corollary: Optimal Shrinkage Intensity of Convex Divergences II

(iii) if  $D$  is taken as Wasserstein divergence, then

$$\lim_{n \rightarrow \infty} n^2 \varepsilon_n^* = \frac{(p+1)^2 p^2}{256 \left( \text{tr}(\Sigma_0^{-1}) - \frac{p}{\|\Sigma_0\|_F^2} \text{tr}(\Sigma_0) \right)^2} \sum_{i=1}^p \frac{(1 - \tau^* \lambda_i)^2}{\lambda_i}. \quad (10)$$



- **Data Generation**

- True covariance

$$\Sigma_0 = V^\top \Lambda V \in \mathbb{S}_+^p, \quad \Lambda = \text{diag}(1, 2, \dots, p),$$

with  $V$  drawn uniformly from the orthogonal group. Thus  $\kappa(\Sigma_0) = p$ .

- Draw  $n$  samples from  $\mathcal{N}(0, \Sigma_0)$  and compute the sample covariance  $\hat{\Sigma}$ .

- **Estimators under Comparison**

$$W(\hat{\Sigma}), \quad L(\hat{\Sigma}), \quad NL(\hat{\Sigma}),$$

denoting the Wasserstein-based nonlinear shrinkage estimator, Ledoit–Wolf linear estimator, and Ledoit–Wolf nonlinear estimator, respectively.

- **Condition Number Error**

$$\varepsilon(M) = \mathbb{E}[|\kappa(M(\hat{\Sigma})) - \kappa(\Sigma_0)|], \quad M \in \{W, L, NL\}.$$

- **Hypothesis Tests** (paired  $t$ -test,  $\alpha = 0.05$ )

$$H_L^0 : \varepsilon(W) \geq \varepsilon(L) \quad \text{vs} \quad H_L^1 : \varepsilon(W) < \varepsilon(L), \quad (\text{I})$$

$$H_{NL}^0 : \varepsilon(W) \geq \varepsilon(NL) \quad \text{vs} \quad H_{NL}^1 : \varepsilon(W) < \varepsilon(NL). \quad (\text{II})$$

- Repeat sampling  $M$  times to obtain  $\hat{\varepsilon}(M)$ .
- Test statistic:  $\hat{\varepsilon}(W) - \hat{\varepsilon}(L)$  or  $\hat{\varepsilon}(W) - \hat{\varepsilon}(NL)$ .
- Reject  $H^0$  implies  $W$  significantly outperforms the competitor.

# Simulation Results

Sample size	Average of $\kappa(W(\hat{\Sigma}))$	Average of $\kappa(L(\hat{\Sigma}))$	$t$ -value	$p$ -value	Reject $H_L^0$ or not
300	320.60	3.00	-76.58	0.0	<b>Reject</b>
400	230.49	3.40	-306.72	0.0	<b>Reject</b>
600	189.82	4.16	-628.65	0.0	<b>Reject</b>
800	180.02	4.88	-578.36	0.0	<b>Reject</b>
1000	177.75	5.61	-642.50	0.0	<b>Reject</b>
1200	177.52	6.30	-652.57	0.0	<b>Reject</b>
1400	177.93	6.98	-690.35	0.0	<b>Reject</b>
1600	178.81	7.65	-710.80	0.0	<b>Reject</b>
1800	179.42	8.30	-736.08	0.0	<b>Reject</b>
2000	180.55	8.95	-780.41	0.0	<b>Reject</b>

Table 3: Hypothesis test for  $p = 200$ .

# Simulation Results

Sample size	Average of $\kappa(W(\hat{\Sigma}))$	Average of $\kappa(NL(\hat{\Sigma}))$	$t$ -value	$p$ -value	Reject $H_{NL}^0$ or not
300	320.60	63.96	-10.30	0.0	<b>Reject</b>
400	230.49	97.50	-69.99	0.0	<b>Reject</b>
600	189.82	128.77	-231.95	0.0	<b>Reject</b>
800	180.02	144.63	-252.25	0.0	<b>Reject</b>
1000	177.75	154.88	-187.61	0.0	<b>Reject</b>
1200	177.52	161.76	-138.53	0.0	<b>Reject</b>
1400	177.93	166.40	-106.77	0.0	<b>Reject</b>
1600	178.81	169.92	-86.35	0.0	<b>Reject</b>
1800	179.42	172.26	-77.84	0.0	<b>Reject</b>
2000	180.55	174.70	-69.77	0.0	<b>Reject</b>

Table 4: Hypothesis test for  $p = 200$ .