

# Affine Invariant Covariance Estimation for Heavy-Tailed Distributions

**Dmitrii M. Ostrovskii**   Alessandro Rudi  
INRIA Paris, Ecole Normale Supérieure

COLT 2019  
Phoenix, AZ  
June 28, 2019



# Covariance Estimation Problem

**Problem:** estimate the covariance matrix  $\mathbf{S} = \mathbf{E}[XX^\top]$  of a zero-mean random vector  $X \in \mathbb{R}^d$  from its  $n$  i.i.d. copies  $X_1, \dots, X_n$ .

**Empirical covariance estimator:**

$$\tilde{\mathbf{S}} = \frac{1}{n} \sum_{i=1}^n X_i X_i^\top.$$

- **easy to compute:**  $O(nd^2)$  time,  $O(d^2)$  memory.
- **statistically favorable** when data is **light-tailed**.
- **equivariant:**  $\tilde{\mathbf{S}}$  behaves as  $\mathbf{S}$  under linear transforms:  $\tilde{\mathbf{S}}' = \mathbf{A}\tilde{\mathbf{S}}\mathbf{A}^\top$ .

Equivariance is useful in applications – gives **affine invariant bounds**.

# Empirical Covariance Estimator: Background

- $\tilde{\mathbf{S}}$  is **statistically favorable** when data is **light-tailed**:

**Assumption: subgaussian marginals.** For any  $u \in \mathbb{R}^d$  and  $p \geq 2$ ,

$$\mathbf{E}^{1/p}[|\langle X, u \rangle|^p] \leq \kappa \sqrt{p} \mathbf{E}^{1/2}[\langle X, u \rangle^2],$$

where  $\kappa$  is a constant for any  $u$  and  $p$ .

- E.g., this holds with  $\kappa = 3$  when  $X \sim \mathcal{N}(0, \mathbf{S})$  with *arbitrary*  $\mathbf{S}$ .

**Theorem** (Koltchinskii and Lounici [2014]): with probability  $\geq 1 - \delta$ ,

$$\frac{\|\tilde{\mathbf{S}} - \mathbf{S}\|}{\|\mathbf{S}\|} \lesssim \kappa^2 \sqrt{\frac{r(\mathbf{S}) + \log(1/\delta)}{n}},$$

where  $\|\cdot\|$  is the operator norm,  $r(\mathbf{S}) = \frac{\text{Tr}(\mathbf{S})}{\|\mathbf{S}\|}$  the effective rank.

# Affine Invariant Bound via Equivariance

$$\frac{\|\tilde{\mathbf{S}} - \mathbf{S}\|}{\|\mathbf{S}\|} \lesssim \kappa^2 \sqrt{\frac{r(\mathbf{S}) + \log(1/\delta)}{n}}, \quad r(\mathbf{S}) = \frac{\text{Tr}(\mathbf{S})}{\|\mathbf{S}\|}.$$

The bound is **not** affine invariant, while the assumption is. **Let's fix it:**

- $\tilde{\mathbf{S}}$  is **equivariant**: behaves the same as  $\mathbf{S}$  under linear transforms.
- Consider (virtual) **decorrelated observations**  $Z_i = \mathbf{S}^{-1/2} X_i$  with **identity covariance** and empirical covariance estimator

$$\frac{1}{n} \sum_{i=1}^n Z_i Z_i^\top = \mathbf{S}^{-1/2} \tilde{\mathbf{S}} \mathbf{S}^{-1/2}.$$

Hence, using the previous result,

$$\|\mathbf{S}^{-1/2}(\tilde{\mathbf{S}} - \mathbf{S})\mathbf{S}^{-1/2}\| \lesssim \kappa^2 \sqrt{\frac{d + \log(1/\delta)}{n}} =: \varepsilon.$$

Equivalently,  $(1 - \varepsilon)\mathbf{S} \preceq \tilde{\mathbf{S}} \preceq (1 + \varepsilon)\mathbf{S}$ , so **relative-scale** eigenvalue bounds.

- Applications in random-design least-squares, noisy subspace iteration.

# Heavy-Tailed Distributions: Truncation

**Assumption:** marginals for any  $u \in \mathbb{R}^d$  have **kurtosis** bounded by  $\kappa$ :

$$\mathbf{E}^{1/4}[|\langle X, u \rangle|^4] \leq \kappa \mathbf{E}^{1/2}[\langle X, u \rangle^2].$$

Under this Asm., Minsker and Wei [2017] consider the truncation estimator

$$\widehat{\mathbf{S}}^{\text{MW}} = \frac{1}{n} \sum_{i=1}^n \tau_{\theta}(\|X_i\|^2 / \|\mathbf{S}\|) X_i X_i^{\top},$$

where  $\tau_{\theta}(\cdot)$  is the truncation map given by  $\tau_{\theta}(x) = \min(x, \theta)/x$ , and prove

$$\frac{\|\widehat{\mathbf{S}}^{\text{MW}} - \mathbf{S}\|}{\|\mathbf{S}\|} \lesssim \kappa^2 \sqrt{\frac{\mathbf{r}(\mathbf{S}) \cdot \log(2d/\delta)}{n}}.$$

- We would like **affine-invariant** bound, something like

$$\|\mathbf{S}^{-1/2}(\widehat{\mathbf{S}} - \mathbf{S})\mathbf{S}^{-1/2}\| \lesssim \kappa^2 \sqrt{\frac{d \cdot \log(2d/\delta)}{n}}.$$

- By equivariance, this bound would hold for the oracle “estimator”:

$$\widehat{\mathbf{S}}^* = \frac{1}{n} \sum_{i=1}^n \tau_{\theta}(\|\mathbf{S}^{-1/2} X_i\|^2) X_i X_i^{\top}.$$

**Theorem.** Under the kurtosis assumption, there exists an estimator  $\widehat{\mathbf{S}}$ , with time complexity  $O(nd^2 + d^3)$  and memory complexity  $O(d^2)$ , that satisfies,

$$\|\mathbf{S}^{-1/2}(\widehat{\mathbf{S}} - \mathbf{S})\mathbf{S}^{-1/2}\| \leq 48\kappa^2 \sqrt{\frac{d \cdot \log(2d/\delta)}{n}}$$

with probability at least  $1 - \delta$ , provided that

$$n \geq 96^2 \kappa^4 d \log(2d/\delta) \cdot \log(\text{cond}(\mathbf{S})).$$

- Extra factor  $\log(\text{cond}(\mathbf{S}))$  in the required  $n$ , but not in the rate.
- Similar cost as for the empirical covariance estimator when  $n \gg d$ .
- Estimator requires (loose) bounds on  $\|\mathbf{S}\|$  and  $\lambda_{\min}(\mathbf{S})$ .

# Key Idea

$$\frac{\|\widehat{\mathbf{S}}^{\text{MW}} - \mathbf{S}\|}{\|\mathbf{S}\|} \lesssim \kappa^2 \sqrt{\frac{\mathbf{r}(\mathbf{S}) \cdot \log(2d/\delta)}{n}}$$

$$\widehat{\mathbf{S}}^{\text{MW}} = \frac{1}{n} \sum_{i=1}^n \tau(\|X_i\|^2 / \|\mathbf{S}\|) X_i X_i^\top$$

$$\|\mathbf{S}^{-1/2}(\widehat{\mathbf{S}} - \mathbf{S})\mathbf{S}^{-1/2}\| \lesssim \kappa^2 \sqrt{\frac{d \cdot \log(2d/\delta)}{n}}$$

$$\widehat{\mathbf{S}}^* = \frac{1}{n} \sum_{i=1}^n \tau(\|\mathbf{S}^{-1/2} X_i\|^2) X_i X_i^\top$$

# Key Idea

$$\frac{\|\widehat{\mathbf{S}}^{\text{MW}} - \mathbf{S}\|}{\|\mathbf{S}\|} \lesssim \kappa^2 \sqrt{\frac{\mathbf{r}(\mathbf{S}) \cdot \log(2d/\delta)}{n}}$$

$$\widehat{\mathbf{S}}^{\text{MW}} = \frac{1}{n} \sum_{i=1}^n \tau(\|X_i\|^2 / \|\mathbf{S}\|) X_i X_i^\top$$

$$\|\mathbf{S}^{-1/2}(\widehat{\mathbf{S}} - \mathbf{S})\mathbf{S}^{-1/2}\| \lesssim \kappa^2 \sqrt{\frac{d \cdot \log(2d/\delta)}{n}}$$

$$\widehat{\mathbf{S}}^* = \frac{1}{n} \sum_{i=1}^n \tau(\|\mathbf{S}^{-1/2} X_i\|^2) X_i X_i^\top$$

- General problem for  $\lambda \geq 0$ :

$$\|(\mathbf{S} + \lambda \mathbf{I})^{-1/2}(\widehat{\mathbf{S}} - \mathbf{S})(\mathbf{S} + \lambda \mathbf{I})^{-1/2}\| \lesssim \kappa^2 \sqrt{\frac{d_\lambda(\mathbf{S}) \cdot \log(2d/\delta)}{n}},$$

where  $d_\lambda(\mathbf{S}) = \text{Tr}[\mathbf{S}(\mathbf{S} + \lambda \mathbf{I})^{-1}]$  is the effective dimension, with oracle

$$\widehat{\mathbf{S}}_\lambda^* = \frac{1}{n} \sum_{i=1}^n \tau(\|(\mathbf{S} + \lambda \mathbf{I})^{-1/2} X_i\|^2) X_i X_i^\top.$$



# Key Idea

$$\frac{\|\widehat{\mathbf{S}}^{\text{MW}} - \mathbf{S}\|}{\|\mathbf{S}\|} \lesssim \kappa^2 \sqrt{\frac{\mathbf{r}(\mathbf{S}) \cdot \log(2d/\delta)}{n}}$$

$$\widehat{\mathbf{S}}^{\text{MW}} = \frac{1}{n} \sum_{i=1}^n \tau(\|\mathbf{X}_i\|^2 / \|\mathbf{S}\|) \mathbf{X}_i \mathbf{X}_i^\top$$

$$\|\mathbf{S}^{-1/2}(\widehat{\mathbf{S}} - \mathbf{S})\mathbf{S}^{-1/2}\| \lesssim \kappa^2 \sqrt{\frac{d \cdot \log(2d/\delta)}{n}}$$

$$\widehat{\mathbf{S}}^* = \frac{1}{n} \sum_{i=1}^n \tau(\|\mathbf{S}^{-1/2} \mathbf{X}_i\|^2) \mathbf{X}_i \mathbf{X}_i^\top$$

- General problem for  $\lambda \geq 0$ :

$$\|(\mathbf{S} + \lambda \mathbf{I})^{-1/2}(\widehat{\mathbf{S}} - \mathbf{S})(\mathbf{S} + \lambda \mathbf{I})^{-1/2}\| \lesssim \kappa^2 \sqrt{\frac{d_\lambda(\mathbf{S}) \cdot \log(2d/\delta)}{n}},$$

where  $d_\lambda(\mathbf{S}) = \text{Tr}[\mathbf{S}(\mathbf{S} + \lambda \mathbf{I})^{-1}]$  is the effective dimension, with oracle

$$\widehat{\mathbf{S}}_\lambda^* = \frac{1}{n} \sum_{i=1}^n \tau(\|(\mathbf{S} + \lambda \mathbf{I})^{-1/2} \mathbf{X}_i\|^2) \mathbf{X}_i \mathbf{X}_i^\top.$$

- **Left:**  $\lambda = \|\mathbf{S}\|$ , oracle is  $\widehat{\mathbf{S}}^{\text{MW}}$  – available! **Right:**  $\lambda = 0$  – what we need.

# Key Idea

$$\frac{\|\widehat{\mathbf{S}}^{\text{MW}} - \mathbf{S}\|}{\|\mathbf{S}\|} \lesssim \kappa^2 \sqrt{\frac{\mathbf{r}(\mathbf{S}) \cdot \log(2d/\delta)}{n}}$$

$$\widehat{\mathbf{S}}^{\text{MW}} = \frac{1}{n} \sum_{i=1}^n \tau(\|\mathbf{X}_i\|^2 / \|\mathbf{S}\|) \mathbf{X}_i \mathbf{X}_i^\top$$

$$\|\mathbf{S}^{-1/2}(\widehat{\mathbf{S}} - \mathbf{S})\mathbf{S}^{-1/2}\| \lesssim \kappa^2 \sqrt{\frac{d \cdot \log(2d/\delta)}{n}}$$

$$\widehat{\mathbf{S}}^* = \frac{1}{n} \sum_{i=1}^n \tau(\|\mathbf{S}^{-1/2} \mathbf{X}_i\|^2) \mathbf{X}_i \mathbf{X}_i^\top$$

- General problem for  $\lambda \geq 0$ :

$$\|(\mathbf{S} + \lambda \mathbf{I})^{-1/2}(\widehat{\mathbf{S}} - \mathbf{S})(\mathbf{S} + \lambda \mathbf{I})^{-1/2}\| \lesssim \kappa^2 \sqrt{\frac{d_\lambda(\mathbf{S}) \cdot \log(2d/\delta)}{n}},$$

where  $d_\lambda(\mathbf{S}) = \text{Tr}[\mathbf{S}(\mathbf{S} + \lambda \mathbf{I})^{-1}]$  is the effective dimension, with oracle

$$\widehat{\mathbf{S}}_\lambda^* = \frac{1}{n} \sum_{i=1}^n \tau(\|(\mathbf{S} + \lambda \mathbf{I})^{-1/2} \mathbf{X}_i\|^2) \mathbf{X}_i \mathbf{X}_i^\top.$$

- **Left:**  $\lambda = \|\mathbf{S}\|$ , oracle is  $\widehat{\mathbf{S}}^{\text{MW}}$  – available! **Right:**  $\lambda = 0$  – what we need.
- **Construction:** start with  $\widehat{\mathbf{S}}^{(0)} = \widehat{\mathbf{S}}^{\text{MW}}$ , and approximate  $\widehat{\mathbf{S}}^*$  iteratively:

$$\widehat{\mathbf{S}}^{(t+1)} = \frac{1}{n} \sum_{i=1}^n \tau(\|(\widehat{\mathbf{S}}^{(t)} + \lambda_t)^{-1/2} \mathbf{X}_i\|^2) \mathbf{X}_i \mathbf{X}_i^\top$$

with  $\lambda_t = \|\mathbf{S}\| \cdot 2^{-t}$ . Proceed for  $\log(\text{cond}(\mathbf{S}))$  iterations, until  $\lambda \leq \lambda_{\min}(\mathbf{S})$ .

- Affine-Invariant bounds are important in applications.
- For equivariant estimators, they follow “automatically” from operator-norm bounds. However, without equivariance this is not so.

We construct an iterative procedure that results in estimators satisfying such bounds in the case of robust covariance estimation.

**Thanks!**

Koltchinskii, V. and Lounici, K. (2014). Concentration inequalities and moment bounds for sample covariance operators. *arXiv:1405.2468*.

Minsker, S. and Wei, X. (2017). Estimation of the covariance structure of heavy-tailed distributions. *arXiv:1708.00502*.

$$\widehat{\mathbf{S}}^{(t+1)} = \frac{1}{n} \sum_{i=1}^n \tau(\|(\widehat{\mathbf{S}}^{(t)} + \lambda_t)^{-1/2} \mathbf{X}_i\|^2) \mathbf{X}_i \mathbf{X}_i^\top, \quad 0 \leq t \leq \log(\text{cond}(\mathbf{S})).$$

Instead, we consider sample splitting:

$$\widehat{\mathbf{S}}^{(t+1)} = \frac{1}{b_t} \sum_{i=1}^{b_t} \tau(\|(\widehat{\mathbf{S}}^{(t)} + \lambda_t)^{-1/2} \mathbf{X}_i^{(t+1)}\|^2) \mathbf{X}_i^{(t+1)} [\mathbf{X}_i^{(t+1)}]^\top,$$

where  $\mathbf{X}_1^{(t+1)}, \dots, \mathbf{X}_{b_t}^{(t+1)}$  is a fresh batch of observations.

**Key lemma:** w.h.p. we have correct accuracy at step  $t + 1$ , i.e.,

$$\|(\mathbf{S} + \lambda_{t+1} \mathbf{I})^{-1/2} (\widehat{\mathbf{S}}^{(t+1)} - \mathbf{S}) (\mathbf{S} + \lambda_{t+1} \mathbf{I})^{-1/2}\| \lesssim \underbrace{\kappa^2 \sqrt{\frac{d \lambda_{t+1}(\mathbf{S}) \cdot \log(2d/\delta)}{n}}}_{\varepsilon_{t+1}},$$

provided **fixed accuracy**  $\varepsilon_t = 1/2$  at step  $t$ .

- Take  $b_t = \frac{n}{2 \log(\text{cond}(\mathbf{S}))}$  for  $t < \log(\text{cond}(\mathbf{S}))$ ;  $b = n/2$  in the end.

# Application: Ridge Regression with Heavy-Tailed Design

Fit  $Y = X^T w^*$  from i.i.d. sample  $(X_i, Y_i)_{i=1}^n$  with  $\mathbf{E}[X] = 0$ ,  $\mathbf{E}[XX^T] = \mathbf{S}$ .

- **Ridge regression** estimator of  $w^*$ :

$$\tilde{w}_\lambda = \frac{1}{n} \sum_{i=1}^n (\tilde{\mathbf{S}} + \lambda \mathbf{I})^{-1/2} X_i Y_i.$$

- Instead, consider

$$\hat{w}_\lambda = \frac{1}{n} \sum_{i=1}^n (\hat{\mathbf{S}} + \lambda \mathbf{I})^{-1/2} \hat{Z}_i,$$

where  $\hat{\mathbf{S}}$  is computed from a hold-out sample by our method;  $\hat{Z}_i$ 's are obtained by appropriately truncating  $Z_i = X_i Y_i$ 's in  $\|\cdot\|_{(\hat{\mathbf{S}} + \lambda \mathbf{I})^{-1}}$ -norm.

**Theorem.** With prob.  $1 - \delta$ ,

$$\|\hat{w}_\lambda - w^*\|_{\hat{\mathbf{S}}}^2 \lesssim \left[ (\kappa^4 + \kappa^2 \varkappa^2) \frac{v^2 d_\lambda(\mathbf{S}) \log(2d/\delta)}{n} + \lambda^2 \left\| (\mathbf{S} + \lambda \mathbf{I})^{-1/2} w^* \right\|^2 \right],$$

whenever  $X$  has  $\kappa$ -bounded marginal kurtoses,  $\mathbf{E}[Y^2] \leq v^2$ ,  $\mathbf{E}[Y^4] \leq \varkappa^4 v^4$ .