# High Dimensional Statistics

Yanbo Tang

Imperial College London

March. 2025

# Content of this week

[a]

Covariance estimation X1/2 = ź 1/ź - ź1/
PCA

- PCA
- Matrix Regression  $Y = X @ f F_{-1} Waterx$ milliple

Some useful matrix inequalities 
$$A = \mathbb{Z} \times \mathcal{U}_{c}^{T} \mathcal{U}_{c}$$
  
 $\|A\|_{F}^{2} : \stackrel{?}{\underset{i=1}{\sum}} \chi_{c}(A)^{2} \quad (?)$   
 $\|W_{Y}\|_{F}^{2} : \stackrel{?}{\underset{i=1}{\sum}} \chi_{c}(B) \stackrel{?}{\underset{i=1}{\sum}} \|A \cdot B\|_{F}^{2} \quad (?)$   
 $\|H_{0}\|_{F}^{2} : \stackrel{?}{\underset{i=1}{\sum}} (\Lambda_{i}(B) - \Lambda_{i}(B)) \stackrel{?}{\underset{i=1}{\sum}} \|A \cdot B\|_{F}^{2} \quad (.2)$   
 $\|A \cdot B \rangle = \frac{1}{2} (\Lambda_{i}^{T}B) \quad (.2)$   
 $\|A \cdot B \rangle = \frac{1}{2} (\Lambda_{i}^{T}B) \quad (.2)$   
 $\|A \cdot B \rangle = \frac{1}{2} (\Lambda_{i}^{T}B) \quad (.2)$ 

0

# Matrix Bernstein

#### Theorem

(Matrix Bernstein's inequality). Let  $X_1, \ldots, X_N$  be independent, mean zero,  $n \times n$  symmetric random matrices, such that  $||X_i|| \le K$  almost surely for all *i*. Then, for every  $t \ge 0$ , we have

$$\mathbb{P}\left\{\left\|\sum_{i=1}^{N} X_{i}\right\| \geq t\right\} \leq 2n \exp\left(-\frac{t^{2}/2}{\sigma^{2} + Kt/3}\right).$$

Here  $\sigma^2 = \left\|\sum_{i=1}^{N} \mathbb{E}X_i^2\right\|$  is the norm of the matrix variance of the sum.

### Expectation

#### Lemma

(Matrix Bernstein's inequality: expectation). Let  $X_1, \ldots, X_N$  be independent, mean zero,  $n \times n$  symmetric random matrices, such that  $||X_i|| \leq K$  almost surely for all *i*.  $\int_{U}^{N} \mathcal{E} \left\| \sum_{i=1}^{N} X_i \right\| \lesssim \left\| \sum_{i=1}^{N} \mathbb{E} X_i^2 \right\|^{1/2} \sqrt{1 + \log n} + K(1 + \log n).$ 

### Covariance estimation

We can estimate the second moment matrix  $\Sigma = \mathbb{E}XX^T$  by its sample version

$$\sum_{\mu} = \frac{1}{m} \sum_{i=1}^{m} X_i X_i^{\mathsf{T}}.$$

Recall that if X has zero mean, then  $\Sigma$  is the covariance matrix of X and  $\Sigma_m$  is the sample covariance matrix of X.

# Covariance Estimation sub-Gaussian

#### Theorem

Let  $Y \in \mathbb{R}^d$  be a random vector such that  $\mathbb{E}[Y] = 0$ ,  $\mathbb{E}[YY^\top] = I_d$  and  $Y \sim subG_d(1)$ . Let  $X_1, \ldots, X_n$  be n independent copies of sub-Gaussian random vector  $X = \Sigma^{1/2}Y$ . Then  $\mathbb{E}[X] = 0$ ,  $\mathbb{E}[XX^\top] = \Sigma$  and  $X \sim subG_d(\|\overline{\Sigma}\|_{op})$ . Moreover,

$$\|\hat{\Sigma} - \Sigma\|_{\text{exp}} \lesssim \|\Sigma\|_{\text{exp}} \left(\sqrt{\frac{d + \log(1/\delta)}{n}} \bigvee \frac{d + \log(1/\delta)}{n}\right),$$

with probability  $1 - \delta$ .

### General covariance estimation

#### Theorem

(General covariance estimation). Let X be a random vector in  $\mathbb{R}^n$ ,  $n \ge 2$ . Assume that for some  $K \ge 1$ ,

$$\|X\|_2 \leq K(\mathbb{E}\|X\|_2)^{1/2}$$
 almost surely.

Then, for every positive integer *m*, we have  $\mathbb{E}\|\Sigma_m - \Sigma\| \lesssim \left(\sqrt{\frac{K^2 \# \log \#}{m}} + \frac{K^2 \# \log \#}{m} + \frac{K^2 \# \log \#}{m}\right) \|\Sigma\|.$ 

(5.16)

Proof  $E \leq ||\chi||^2 = f_{\mu}(\mathbb{Z})$  $\|\chi\|_{1}^{2} \leq k^{2} f(\mathbb{Z})$ ET 112-5/1] = 1 EU = CX:X7-5) 17 & 1 (of layed) + Mlayed)  $\mathcal{F}^{2} = h || E(r k^{7} - z)^{2} || = || = || = E(r k^{7} - z) ||$  $M : \|XX^{T} - \Xi\| \le M \quad 0.5.$ ~ <sup>2</sup>  $E_{Z}(Xx^{T}-\Xi)^{2}$  =  $E_{Z}(Xx^{T})^{2}$  -  $\Xi^{2} = E_{Z}(Xx^{T})^{2}$  $(\chi\chi^{r})^{2} = \|\chi\|^{2} \chi\chi^{r} \leq k^{2} f(z) \chi\chi^{T}$ => EE (XX T)2] - k2 f(E) ||E|| => 52 = k2 n. t(=) 1151

# Proof cont. $\|\chi\chi^{7} - \varepsilon\| \leq \|\chi\|_{2}^{2} \neq \|\Xi\|$

# A probability bound

#### Lemma

Under the same assumption as the previous theorem, we have

$$\|\hat{\Sigma} - \Sigma\| \lesssim \left(\sqrt{\frac{\kappa^2 \mathcal{O}(\log \mathcal{O} + \log(2/\delta))}{\mathcal{O} \mathcal{O}}} + \frac{\kappa^2 \mathcal{O}(\log \mathcal{O} + \log(2/\delta))}{\mathcal{O} \mathcal{O}}\right) \|\Sigma\|$$

with probability at least  $1 - \delta$ .



### Eckart-Young-Mirsky

#### Lemma

Let A be a rank-r matrix with singular value decomposition

$$A = \sum_{i=1}^r \lambda_i u_i v_i^{\top},$$

where  $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_r > 0$  are the ordered singular values of A. For any k < r, let  $A_k = \sum_{i=1}^k \lambda_i u_i v_i^{\top}$ . Then for any matrix B such that rank $(B) \leq k$ , it holds

$$\|A-A_k\|_F\leq \|A-B\|_F.$$

Moreover,

$$\mathcal{O} \ \|A - A_k\|_F^2 = \sum_{j=k+1}^r \lambda_j^2.$$

13/33

Proof Q is by definition wind  $A \cdot A_k = \sum_{\substack{j \in k \in I \\ j \in k \in I}}^{n} \lambda_j u_j V_i^T$   $3 \quad 1/A - A_k II_F^2 = \sum_{\substack{j \in k \in I \\ j \in k \in I}}^{n} \lambda_j^2$  $\forall B \text{ of } rank \neq k$  $\|A - B\|_{F}^{2} = \sum_{j=1}^{k} (\lambda_{j} - \delta_{j})^{2}$   $= \sum_{j=1}^{k} (\lambda_{j} - \delta_{j})^{2} + \sum_{j=k}^{k} \lambda_{j}^{2}$   $= \sum_{j=1}^{k} (\lambda_{j} - \delta_{j})^{2} + \sum_{j=k}^{k} \lambda_{j}^{2}$ 6,7.8,7% are singilar value of B

## Spiked covariance

For a fixed direction  $v \in S^{d-1}$ , and consider the sequence  $Y_1, \ldots, Y_n \sim N(0, I_d)$ , then the vectors  $v^{\top} Y_i v$  all live in the one dimensional space spanned by v. Typically we would have some noise in our observations so that they will be closer to:

$$X_i = \mathbf{v}^\top Y_i \mathbf{v} + Z_i,$$

for a noise vector  $Z_i \sim N(0, \sigma^2 I_d)$ . Note that the covariance matrix of  $X_i$  will be:

$$\Sigma = E\left[XX^{\top}\right] = vv^{\top} + \sigma^2 I_d.$$

# Spiked covariance

### Definition

A covariance matrix  $\Sigma \in \mathbb{R}^{d \times d}$  is a spiked covariance matrix if:

$$\Sigma = \theta v v^{\top} + I_d,$$

for  $\theta > 0$  and  $v \in S^{d-1}$ . The vector v is called the spike.

$$V_1(z) = v$$

# Davis Kahan reminder

#### Theorem

(Davis-Kahan) Let S and T be symmetric matrices with the same dimensions. Fix i and assume that the i-th largest eigenvalue of S is well separated from the rest of the spectrum:

$$\min_{j:j\neq i} |\lambda_i(S) - \lambda_j(S)| = \delta > 0.$$

Then the angle between the eigenvectors of S and T corresponding to the *i*-th largest eigenvalues (as a number between 0 and  $\pi/2$ ) satisfies

$$\sin \angle (v_i(S), v_i(T)) \leq rac{2\|S-T\|}{\delta}.$$

which implies

$$\exists heta \in \{-1,1\} : \|v_i(S) - heta v_i(T)\|_2 \le rac{2^{3/2} \|S - T\|}{\delta}$$

### Guarantees

#### Corollary

Let  $Y \in \mathbb{R}^d$  be a random vector such that  $\mathbb{E}[Y] = 0$ ,  $\mathbb{E}[YY^T] = I_d$  and  $Y \sim subG_d(1)$ . Let  $X_1, \ldots, X_n$  be n independent copies of sub-Gaussian random vector  $X = \Sigma^{1/2}Y$  so that  $\mathbb{E}[X] = 0$ ,  $\mathbb{E}[XX^T] = \Sigma$  and  $X \sim subG_d(\|\Sigma\|_{op})$ . Assume further that  $\Sigma = \theta vv^T + I_d$  satisfies the spiked covariance model. Then, the largest eigenvector  $\hat{v}$  of the empirical covariance matrix  $\hat{\Sigma}$  satisfies,  $\|\mathcal{G}_{\mathcal{F}}\| \leq \|$  $\min_{e \in \{\pm 1\}} \|e\hat{v} - v\|_2 \lesssim \frac{1+\theta}{\min(\theta, 1)} \left(\sqrt{\frac{d + \log(1/\delta)}{n}} \vee \frac{d + \log(1/\delta)}{n}\right)$ 

with probability  $1 - \delta$ .

18/33

Proof By Dans kalon 
$$\|V_{i}(\vec{z}) - V_{i}(\vec{z})\|_{2}^{2} = 2^{\frac{3^{2}}{2}} \frac{||\vec{z} - \vec{z}||}{\delta}$$
  
 $|fO_{i}, O_{i}, O_{i} \in \min(1, \Theta)$   
 $||\vec{z} - \vec{z}||$ 

### Sparse PCA

If we assume that v in the spiked covariance model is k-sparse:  $|v|_0 = k$ . Therefore, a natural candidate to estimate v is given by  $\hat{v}$  defined by

$$\hat{v}^{\top} \hat{\Sigma} \hat{v} = \max_{u \in S^{d-1}, |u|_0 = k} u^{\top} \hat{\Sigma} u.$$

It is the case that  $\lambda_{\max}^k(\hat{\Sigma}) = \hat{v}^\top \hat{\Sigma} \hat{v}$  is the largest of all leading eigenvalues among all  $k \times k$  sub-matrices of  $\hat{\Sigma}$  so that the maximum is indeed attained.



# Sparse PCA

#### Theorem

wi

Let  $Y \in \mathbb{R}^d$  be a random vector such that  $\mathbb{E}[Y] = 0$ ,  $\mathbb{E}[YY^{\top}] = I_d$  and  $Y \sim subG_d(1)$ . Let  $X_1, \ldots, X_n$  be n independent copies of sub-Gaussian random vector  $X = \Sigma^{1/2}Y$  so that  $\mathbb{E}[X] = 0$ ,  $\mathbb{E}[XX^{\top}] = \Sigma$  and  $X \sim subG_d(\|\Sigma\|_{op})$ . Assume further that  $\Sigma = \theta vv^{\top} + I_d$  satisfies the spiked covariance model for v such that  $|v|_0 = k \leq d/2$ . Then, the k-sparse largest eigenvector  $\hat{v}$  of the empirical covariance matrix satisfies,

$$\begin{split} \min_{\varepsilon \in \{\pm 1\}} \|\varepsilon \hat{v} - v\|_{2} \\ \lesssim \frac{1 + \theta}{\min(\theta, 1)} \left( \sqrt{\frac{k \log(ed/k) + \log(1/\delta)}{n}} \vee \frac{k \log(ed/k) + \log(1/\delta)}{n} \right) \\ th \text{ probability } 1 - \delta. \end{split}$$

Multiple regression  $\begin{pmatrix} \gamma_{11} & \gamma_{12} & \gamma_{23} \\ \gamma_{el} & \gamma_{21} & \gamma_{23} \\ \vdots & \vdots & \vdots \\ \gamma_{el} & \gamma_{el} & \gamma_{el} \end{pmatrix} \subset \chi \stackrel{(4)}{\leftarrow} + \mathcal{E}_{-7} M^{(4)} + \mathcal{E}_{-7} M^$  $Y_1 = X_1 + \mathcal{E}_1$  $Y_2 = X_2 + \mathcal{E}_2$ 

### Introduction

A simple question to ask is can we extend the classical regression set up to matrices? Specifically, is it useful to consider a model such that:

 $\mathbb{Y} = \mathbb{X}\Theta^{\star} + E,$ 

where  $Y \in \mathbb{R}^{n \times T}$ ,  $\mathbb{X} \in \mathbb{R}^{n \times d}$  and  $\Theta \in \mathbb{R}^{d \times T}$  is the unkown matrix of coefficients for some noise matrix  $E \sim subG_{n \times T}(\sigma)$ .

#### Definition

We call a  $n \times m$  matrix  $A subG_{n \times m}(\sigma)$  if for every  $x \in S^{m-1}$  and  $y \in S^{n-1}$ if:

 $y^{\top}Ax \sim subG(\sigma^2).$ 

### Direct observation model

Linew regression 
$$\int_{2}^{2} \gamma M_{1}^{2} \left( \frac{x}{2} \partial \sigma^{2} \right)$$

We are going to make our life simple (at first) and assume that we have an orthogonal design (ORT condition) for our design matrix, i.e.,  $X^{\top}X = nI_d$ . Under the ORT assumption,

$$\frac{1}{n}X^{\top}Y = \Theta^* + \frac{1}{n}X^{\top}E.$$

Which can be written as an equation in  $\mathbb{R}^{d \times T}$  called the *sub-Gaussian* matrix model (*sGMM*):

$$y = \Theta^* + F$$
,  
where  $y = \frac{1}{n} X^\top Y$  and  $F = \frac{1}{n} X^\top E \sim \text{subG}_{d \times T}(\sigma^2/n)$ .

If  $\Theta^*$  is sparse then it is possible to estimate it from a single observation. Consider the SVD of  $\Theta^*$ :

$$\Theta^* = \sum_j \lambda_j u_j v_j^\top.$$

and define  $\|\Theta^*\|_0 := |\lambda|_0$ . Therefore, if we knew  $u_j$  and  $v_j$ , we could simply estimate the  $\lambda_j$ s thresholding. It turns out that estimating the eigenvectors by themselves is sufficient. Consider the SVD of the observed matrix y:

$$y = \sum_{j} \hat{\lambda}_{j} \hat{u}_{j} \hat{v}_{j}^{\top}.$$

# Singular value tresholding

**Definition** The singular value thresholding estimator with threshold  $\not p_{\tau_n} \ge 0$  is defined by

$$\hat{\Theta}^{\mathsf{SVT}} = \sum_{j} \hat{\lambda}_{j} \mathbb{I}(|\hat{\lambda}_{j}| > \tau_{\mathsf{n}}) \hat{u}_{j} \hat{v}_{j}^{\top}.$$

### Singular value tresholding

**Definition** The singular value thresholding estimator with threshold  $2\tau \ge 0$  is defined by

$$\hat{\Theta}^{\mathsf{SVT}} = \sum_{j} \hat{\lambda}_{j} \mathbb{I}(|\hat{\lambda}_{j}| > \tau_{n}) \hat{u}_{j} \hat{v}_{j}^{\top}.$$

Lemma

Let A be a  $d \times T$  random matrix such that  $A \sim subG_{d \times T}(\sigma^2)$ . Then

$$\|A\|_{BMP} \leq 4\sigma \sqrt{\log(12)(d \vee T)} + 2\sigma \sqrt{2\log(1/\delta)}$$

with probability  $1 - \delta$ .

#### Theorem

Consider the multivariate linear regression model under the assumption **ORT** or, equivalently, the sub-Gaussian matrix model. Then, the singular value thresholding estimator  $\hat{\Theta}^{SVT}$  with threshold

$$\tau_n = 8\sigma \sqrt{\frac{\log(12)(d \vee T)}{n}} + 4\sigma \sqrt{\frac{2\log(1/\delta)}{n}},$$

satisfies

$$\begin{split} \frac{1}{n} \|\bar{X}\hat{\Theta}^{SVT} - \bar{X}\Theta^*\|_F^2 &= \|\hat{\Theta}^{SVT} - \Theta^*\|_F^2 \leq 36 \operatorname{rank}(\Theta^*)\tau_n^2 \\ &\lesssim \frac{\sigma^2 \operatorname{rank}(\Theta^*)}{n} \left( d \lor T + \log(1/\delta) \right). \end{split}$$
with probability  $1 - \delta$ .

Proof 
$$@^* \quad \lambda_1 \geq \lambda_2 \geq \dots$$
  
 $y \quad \lambda_r \geq \lambda_2 \geq \lambda_3 \geq \dots$   
 $Afme \quad S = Sj : \quad |\lambda_j| \geq 2\pi S \qquad ||F|| \leq 3K \quad u.p \quad 1 \rightarrow S$   
 $\mid \lambda_j^2 - \lambda_j \mid \leq ||F|| \leq 2\pi S \qquad ||F|| \leq 3K_{2} \quad y = 1 \rightarrow S$   
 $S \subset Sj : \quad |\lambda_j| \geq 2K S \qquad S \subset [j : |\lambda_j|] \leq 3K_{2} \quad y = 2K_{2} \quad (j = 1) \quad |\lambda_j| \leq 3K_{2} \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad (j = 1) \quad |\lambda_j| \leq 1 \rightarrow S \quad (j = 1) \quad (j = 1)$ 

# Without ORT

We can then consider penalizing the rank of the matrix instead of direct truncation in cases when we don't have the ORT assumption. Let  $\hat{\Theta}^{RK}$  be any solution to the following minimization problem:

$$\min_{\Theta \in \mathbb{R}^{d \times T}} \left\{ \frac{1}{n} \| Y - X\Theta \|_F^2 + \tau_n^2 \operatorname{rank}(\Theta) \right\}.$$

### Guarantees

#### Theorem

Consider the multivariate linear regression model (5.1). Then, the estimator by rank penalization  $\hat{\Theta}^{RK}$  with regularization parameter  $\tau_n^2$ , where  $\tau_n$  is defined in as in the previous theorem, satisfies

$$\frac{1}{n} \| X \hat{\Theta}^{RK} - X \Theta^* \|_F^2 \leq 2 \operatorname{rank}(\Theta^*) \tau^2 \lesssim \frac{\sigma^2 \operatorname{rank}(\Theta^*)}{n} \left( d\sqrt{T} + \log(1/\delta) \right),$$

with probability  $1 - \delta$ .

# Proof

Why can we solve this problem efficiently

# Concluding remarks

Tologrand - Concertration inequalities -> Issperinetic Takyoul e.y. Chernell bunk inequalities Masurt 2-700 Concontration un spheres. on granstion processes. - Comparit. Sup / IKE 11, Vershyain -up 29, f> feF -Rondom matrix theory - minimer rates levrence lao. Kiyollet Chyden 4/ mining the worst lase in Sop