# High Dimensional Statistics

Yanbo Tang

Imperial College London

March. 2025

## Content of this week

- LASSO sparse set recovery (PDW)
- Matrix concentration ||A||<sub>op</sub>
  Spectral clustering on Rodom glaphs
  Matrix Bernstein || Z A c ||<sub>op</sub>
  Covariance estimation = X:X

### LASSO

#### Theorem

Consider an S-sparse linear regression model for which the design matrix satisfies Assumptions 2 and 3. Then for any choice of regularization parameter such that

$$\lambda_n \ge \frac{2}{1-\alpha} \left\| X_{\mathcal{S}}^{\mathsf{T}} \mathsf{\Pi}_{\mathcal{S}^{\perp}}(\mathsf{X}) \frac{\epsilon}{n} \right\|_{\infty},\tag{1}$$

the Lasso program has the following properties:

- (a) Uniqueness: There is a unique optimal solution  $\hat{\theta}$ .
- (b) No false inclusion: This solution has its support set  $\hat{S}$  contained within the true support set S.
- (c)  $\ell_{\infty}$ -bounds: The error  $\hat{\theta} \theta^*$  satisfies

$$\|\hat{\theta}_{S} - \theta_{S}^{*}\|_{\infty} \leq \left\| \left(\frac{X_{S}^{T} X_{S}}{n}\right)^{-1} X_{S}^{T} \frac{\epsilon}{n} \right\|_{\infty} + \left\| \left(\frac{X_{S}^{T} X_{S}}{n}\right)^{-1} \right\|_{\infty} \lambda_{n},$$

$$\tag{2}$$

where  $||A||_{\infty} = \max_{i=1,...,s} \sum_{j} |A_{i,j}|$  is the matrix  $\ell_{\infty}$ -norm.

(d) No false exclusion: The Lasso includes all indices  $i \in S$  such that  $|\theta_i^*| > B(\lambda_n; X)$ , and hence is variable selection consistent if  $\min_{i \in S} |\theta_i^*| > B(\lambda_n; X)$ .

#### Corollary

For a S-sparse linear model based on a noise vector  $\epsilon$  with zero-mean i.i.d.  $\sigma$ -sub-Gaussian entries, and a deterministic design matrix X that satisfies Assumptions 2 and 3, as well as the C-column normalization condition  $\max_{j=1,...,d} \|X_j\|_2 / \sqrt{n} \leq C$ . Suppose that we solve the Lasso program with regularization parameter

$$\lambda_n = \frac{2C\sigma}{1-\alpha} \left\{ \sqrt{\frac{2\log(d-k)}{n}} + \delta \right\}$$

for some  $\delta > 0$ . Then the optimal solution  $\hat{\theta}$  is unique with its support contained within S, and satisfies the  $\ell_{\infty}$ -error bound

$$\|\widehat{\theta}_{S} - \theta_{S}^{*}\|_{\infty} \leq \frac{\sigma}{\sqrt{c_{\min}}} \left( \sqrt{\frac{2\log s}{n}} + \delta \right) + \left\| \left( \frac{\mathsf{X}_{S}^{T}\mathsf{X}_{S}}{n} \right)^{-1} \right\|_{\infty} \lambda_{n}, \qquad (3)$$

all with probability at least  $1 - 4e^{-n\delta^2/2}$ .

### Proof

Conver function f: 12 dr-7/R Sub-gradients We say ZE/Rd is a sub-gradient "1 f at (8) f(OtA) Zf(0) t <Z, ~> Ha C/Rd then ZEZ (10) (be song a poir (0,2) is primal-Dual optimal if 1) & 15 a minimizer 2) ZEZ11011 Z; EI-1, 17  $\frac{1}{2}\chi^{T}(\chi\hat{G}-\chi)+\lambda_{T}\hat{Z}=0$ 

# Primal Dual Witness

### Definition

### Primal-dual witness (PDW) construction:

• Set 
$$\hat{\theta}_{S^c} = 0$$
.

2 Determine  $(\widehat{\theta}_S, \widehat{z}_S) \in \mathbb{R}^s \times \mathbb{R}^s$  by solving the *oracle subproblem* 

$$\widehat{\theta}_{S} \in \arg\min_{\theta_{S} \in \mathbb{R}^{s}} \left\{ \underbrace{\frac{1}{2n} \|\mathbf{y} - \mathbf{X}_{S} \theta_{S}\|_{2}^{2} + \lambda_{n} \|\theta_{S}\|_{1}}_{=:f(\theta_{S})} \right\},$$
(4)

and then choosing  $\widehat{z}_{S} \in \partial \|\widehat{\theta}_{S}\|_{1}$  such that  $\nabla f(\theta_{S})|_{\theta_{S} = \widehat{\theta}_{S}} + \lambda_{n}\widehat{z}_{S} = 0$ .

**③** Solve for  $\hat{z}_{S^c} \in \mathbb{R}^{d-s}$  via the zero-subgradient equation, and check whether or not the *strict dual feasibility* condition  $\|\hat{z}_{S^c}\|_{\infty} < 1$  holds.

Witness for LASSO 
$$\lambda_{min}\left(\frac{\chi_{5}^{7}\chi_{5}}{m}\right) > C_{min} > 0$$

#### Lemma

If the lower eigenvalue condition holds, then success of the PDW construction implies that the vector  $(\hat{\theta}_S, 0) \in \mathbb{R}^d$  is the unique optimal solution of the Lasso.

State if 
$$\vec{O}$$
  $||\vec{O}||_{1} \leq \chi \neq 1, \vec{O} \geq \leq ||\vec{Z}||_{\infty} ||\vec{O}||_{1} = |\vec{I}\vec{O}||_{1}$   
 $\vec{Z}_{i} = -1, 1 \qquad \langle \neq 1, \vec{O} \rangle = ||\vec{O}||_{1}$   
 $\vec{z}_{i} = -1, 1 \qquad \langle \neq 2, \vec{O} \rangle = ||\vec{O}||_{1}$   
 $\vec{z}_{i} = -1, 1 \qquad \langle \neq 2, \vec{O} \rangle = ||\vec{O}||_{1}$   
 $\vec{z}_{i} = -1, 1 \qquad \langle \neq 2, \vec{O} \rangle = ||\vec{O}||_{1}$   
 $\vec{z}_{i} = -1, 1 \qquad \langle \neq 2, \vec{O} \rangle = ||\vec{O}||_{1}$   
 $\vec{z}_{i} = -1, \vec{D} \qquad \langle \neq 2, \vec{O} \rangle = ||\vec{O}||_{1}$ 

### Proof

### Proof of main theorem

$$\frac{2}{Sc} = \frac{-1}{2} X_{S}^{T} X_{S} (\theta_{S} - \theta_{S}^{*}) \neq X_{S}^{T} (\frac{s}{2\pi n}) \quad (i)$$

$$\frac{2}{Sc} = \frac{-1}{2} X_{S} (X_{S})^{4} X_{S}^{T} \varepsilon - \lambda_{n} n (X_{S}^{T} X_{S})^{-1} \frac{s}{2s} \quad (2)$$

$$\frac{2}{Sc} = \frac{1}{2} \frac{1}{S} \frac{1}{S} \frac{1}{S} (X_{S})^{4} \frac{1}{2s} \varepsilon + \frac{1}{2s} \frac{1}{S} \frac{1}{S} \frac{1}{2s} \frac{$$

Yanbo Tang (Imperial College London)

# Matrix Concentration

### Some review

### The Courant Fisher min-max theorem which states that:

$$\lambda_i(A) = \max_{dim(E)=i} \min_{x \in S(E)} x^\top A x$$

where the maximum is taken over all *i*-dimensional subspaces E of  $\mathbb{R}^n$ .

### Some review

The Courant Fisher min-max theorem which states that:

$$\lambda_i(A) = \max_{dim(E)=i} \min_{x \in S(E)} x^\top A x$$

where the maximum is taken over all *i*-dimensional subspaces E of  $\mathbb{R}^n$ . Using this we can characterize the operator norm or the maximum singular value as follows:

$$\|A\|_{\mathcal{V}} := \max_{x \in \mathbb{R}^n \setminus \{0\}} \frac{\|Ax\|_2}{\|x\|_2} = \max_{x \in S^{n-1}} \|Ax\|_2.$$

Equivalently, the operator norm of A can be computed by maximizing the quadratic form  $\langle Ax, y \rangle$  over all unit vectors x, y:

$$\|A\| = \max_{x \in S^{n-1}, y \in S^{m-1}} \langle Ax, y \rangle.$$

# Controlling the operator norm

#### Lemma

Let A be an  $m \times n$  matrix and  $\varepsilon \in [0,1)$ . Then, for any  $\varepsilon$ -net  $\mathcal{N}$  of the sphere  $S^{n-1}$ , we have

$$\begin{split} \sup_{x \in \mathcal{N}} \|Ax\|_{2} &\leq \|A\| \leq \frac{1}{1 - \varepsilon} \sup_{x \in \mathcal{N}} \|Ax\|_{2}. \\ \hline \text{Proof:} \quad fix \quad x \quad for \quad which \quad \|A\| = \|Ax\|_{2} \\ \hline \text{Choose} \quad \chi_{0} \in \mathcal{N} \quad where \quad \|X - X_{0}\|_{2} \leq \varepsilon \\ \quad \|Ax - Ax_{0}\|_{2} \leq \|A\| \quad \|X - X_{0}\|_{2} \leq \varepsilon \\ \quad \|Ax - Ax_{0}\|_{2} \leq \|A\| \quad \|X - X_{0}\|_{2} \leq \varepsilon \\ \quad \|Ax - Ax_{0}\|_{2} \leq \|A\| \quad \|X - X_{0}\|_{2} \leq \varepsilon \\ \quad \|A\| \leq \|A\| \\ \quad \|A\| \leq \sup_{x \in \mathcal{N}} \frac{\|Ax\|}{1 - \varepsilon} \\ \quad \|Ax\| \\ \quad \|A\| \leq \sup_{x \in \mathcal{N}} \frac{\|Ax\|}{1 - \varepsilon} \\ \end{split}$$

### Proof

## The actual lemma

#### Lemma

Let A be an  $m \times n$  matrix and  $\varepsilon \in [0,1)$ . Then, for any  $\varepsilon$ -net  $\mathcal{N}$  of the sphere  $S^{n-1}$  and any  $\varepsilon$ -net  $\mathcal{M}$  of the sphere  $S^{m-1}$ , we have

$$\sup_{x \in \mathcal{N}} \sup_{y \in \mathcal{M}} y^{\top} A x \leq \|A\|_{\mathscr{B}} \leq \frac{1}{1 - 2\varepsilon} \sup_{x \in \mathcal{N}} \sup_{y \in \mathcal{M}} y^{\top} A x.$$

### Random sub-Gaussian matrices

#### Theorem

Let A be an  $m \times n$  random matrix whose entries  $A_{ij}$  are independent, mean zero,  $\sigma$ -sub-gaussian random variables. Then, for any t > 0 we have

$$\|A\| \le C\sigma \left(\sqrt{m} + \sqrt{n} + t\right)$$

with probability at least  $1 - 2\exp(-t^2)$  for some constant C > 0.

Step 1: Approximation. Choose  $\mathcal{E} = V_{4}$  then have  $\mathcal{E} - \text{Net } \mathcal{N}(x)$   $\mathcal{E} - \text{Net } \mathcal{M}(y)$   $|\mathcal{N}| \leq q^{n}$   $|\mathcal{M}| \leq q^{m}$  $||\mathcal{A}|| \leq 2 \max_{\substack{K \in N_{r} \\ Y \in M}} \mathcal{L}A_{K,Y}^{2}$ 

 $\left( \begin{array}{c} \overbrace{z} \overbrace{z} \overbrace{z} A_{ij} A_{ij} A_{ij} \end{array} \right)$ Step 2: Concentration.  $fix \quad x \in \mathcal{N} \quad cil \quad y \in \mathcal{N}$   $fix \quad x \in \mathcal{N} \quad cil \quad y \in \mathcal{N}$   $LAx_{i}y \quad 7^{2} \quad z = z = A_{ij} \quad \lambda_{i} \quad \lambda_{j} \quad is \quad a \quad s = c_{j} \quad independent \\ cil \quad j^{2} \quad i \quad c^{2} \quad j^{2} \quad c^{2} \quad s = b - Gaassians.$   $fix \quad x \in \mathcal{N} \quad x \in \mathcal{N} \quad x \in \mathcal{N}$ p(2Ar, y) = 12) ¿ exp(- (F))

Step 3: Union bound.

 $if \max_{X \in \mathcal{N}, y \in \mathcal{M}} \mathcal{L}Ax, y \geq \mathcal{I}$   $p(\max_{X \in \mathcal{N}, y \in \mathcal{M}} \mathcal{L}Ax, y \geq \mathcal{I}) \leq \mathcal{I} \mathcal{P}(\mathcal{L}Ax, y \geq \mathcal{I})$   $p(\max_{X \in \mathcal{N}, y \in \mathcal{M}} \mathcal{L}Ax, y \geq \mathcal{I}) \leq \mathcal{I} \mathcal{I} \mathcal{I}$  $\leq q^{hfm} \cdot 2 \exp\left(\frac{-\mu^2}{c\sigma^2}\right)$ u= CB (Vatum ff) then solve to get  $\leq 2 \exp(-\epsilon^2)$ 

Application stochastic block model A12 Vi A27 Vi A27 Vi S G(n,p) G2(n,p;q) E[A] = [pp qq pp qq n/2 per group [qq pp] qq pp] 2 groups IELA] || ≈ N, IR || ≤ CIA A = EZAIFR  $\lambda_{1} = (\underline{p+q}) \cdot \mathcal{N} \qquad \qquad \mathcal{A}_{1} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}, \quad \mathcal{A}_{2} = \begin{bmatrix} 1 \\ 1 \\ -1 \\ -1 \end{bmatrix}$ EGAl X3 20

### Illustration



Figure 1: Taken from High Dimensional Probability by Roman Vershynin, with n = 200, p = 1/20 and q = 1/200.

# Cont. Pavis kahan.

#### Theorem

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

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

The conclusion of the Davis-Kahan theorem implies that the unit eigenvectors  $v_i(S)$  and  $v_i(T)$  are close to each other up to a sign, namely

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

Cont.  $\int = mor(\frac{p}{2}, q) \cdot n =: \mu \cdot n$  $70 \in \{-1, 1\}$  :  $||V_2(E(A)) - OV_2(A)||_2 \leq \frac{Con}{n - n} = \frac{C}{n \sqrt{n}}$ /- expc-n)  $\| V_{\Delta}(E(A)) - \Theta V_{\Delta}(A) \| \leq \frac{C}{\mathcal{H}}$  $f_{f_{a}}$  number of mirden  $\leq \frac{C}{\mu^2}$ 

## Bernstein inequality for matrices

X,2MO,2)

Recall that we were able to use the fact that

$$E[\exp(t(X+Y))] = E[\exp(tX)\exp(tY)],$$

and use Chernoff's method if X and Y are real valued random variables. To generalize this approach we need to define what a *matrix exponential* means:

### Bernstein inequality for matrices

Recall that we were able to use the fact that

$$E[\exp(t(X+Y))] = E[\exp(tX)\exp(tY)],$$

and use Chernoff's method if X and Y are real valued random variables. To generalize this approach we need to define what a *matrix exponential* means:

### Definition

For a function  $f : \mathbb{R} \to \mathbb{R}$  and an  $n \times n$  symmetric matrix

$$X=\sum_{i=1}^n\lambda_iu_iu_i^{\top},$$

define

$$f(X) := \sum_{i=1}^n f(\lambda_i) u_i u_i^{\top}.$$

Yanbo Tang (Imperial College London)

**Probability for Statistics** 

### Matrix power series

For a convergent power series expansion of f about  $x_0$ :

$$f(x) = \sum_{k=1}^{\infty} a_k (x - x_0)^k.$$

It is the case that series of matrix terms converges, and

$$f(X) = \sum_{k=1}^{\infty} a_k (X - x_0 I)^k.$$

As an example, for each  $n \times n$  symmetric matrix X we have

$$e^{X} = I + X + \frac{X^{2}}{2!} + \frac{X^{3}}{3!} + \cdots$$

$$A \cdot B \leq 13 \cdot A$$

$$A \cdot B \leq 13 \cdot A$$

# Generalization of exponential inequality

### Theorem

(Golden-Thompson inequality). For any  $n\times n$  symmetric matrices A and B, we have

$$tr(e^{A+B}) \leq tr(e^A e^B).$$

Unfortunately, Golden-Thompson inequality does not hold for three or more matrices: in general, the inequality  $tr(e^{A+B+C}) \leq tr(e^A e^B e^C)$  may fail.

#### Theorem

(Lieb's inequality). Let H be an  $n \times n$  symmetric matrix. Define the function on matrices

$$f(X) := tr \exp(H + \log X).$$

Then f is concave on the space on positive definite  $n \times n$  symmetric matrices.

### Cont.

#### Lemma

(Lieb's inequality for random matrices). Let H be a fixed  $n \times n$  symmetric matrix and Z be a random  $n \times n$  symmetric matrix. Then

 $\mathbb{E}tr\exp(H+Z) \leq tr\exp(H+\mathbb{E}Z).$ 

This follows by using Jensen's inequality.

## 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.

# Moment generating function of random matrices $\begin{bmatrix} 0 \\ 0 \\ 2 \end{bmatrix} \neq \begin{bmatrix} 2 \\ 0 \\ 2 \end{bmatrix}$

#### Lemma

(Moment generating function). Let X be an  $n \times n$  symmetric mean zero random matrix such that  $||X|| \leq K$  almost surely. Then

$$\mathbb{E} \exp(\lambda X) \preceq \exp(g(\lambda)\mathbb{E}X^2)$$
 where  $g(\lambda) = \frac{\lambda^2/2}{1-|\lambda|K/3}$ ,

provided that  $|\lambda| < 3/K$ .



Matrix Bernsetein - Step 1: Reduction to MGF  $\int := \sum_{i=1}^{N} \chi_{i}$ ||S|| = max | 2.(5) = max ( 2max (S), - 5/2/10)  $\frac{\|S\| - m_{r}}{(z)} = p(e^{\lambda \lambda nax(s)} = e^{\lambda t})$   $= e^{-\lambda t} E E e^{\lambda \lambda naa(s)} A$ do max ? E = EEMmax(e<sup>AS</sup>)] = EE tr esp (AS)7

EZA7 Step 2: Application of Lieb's inequality EZ EZA 1877  $E \leq E I f_{T} \exp\left(\frac{\sum_{i=1}^{N-1} \lambda X_{i} + \lambda X_{i}}{\sum_{i=1}^{N-1} \lambda X_{i} + \lambda X_{i}}\right)$ Condition on  $(X_i)_{i=1}^{N-1}$  apply lower with  $H := \sum_{i=1}^{n-1} \chi_{X_i}$ ont puls z: => XN ≤ E tr eng( E AA; + lag EzeAM)]) = thep I = by ( EE exis)] = tr exp(g(N) Z) Z = E EZX: ] 2 n. Tack ( exp(g(x) 2)) = n · exp (g() /my(z))  $= n - exp(g(\lambda) \sigma^2)$ then play back into \$ out avining get the

Yanbo Tang (Imperial College London)

### Step 3: Using the MGF bound

### 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*.

$$\mathbb{E}\left\|\sum_{i=1}^{N} X_{i}\right\| \lesssim \left\|\sum_{i=1}^{N} \mathbb{E} X_{i}^{2}\right\|^{1/2} \sqrt{1 + \log n} + \mathcal{K}(1 + \log n).$$

### General covariance estimation

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

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

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.

### General covariance estimation

#### Theorem

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

$$\|X\|_2 \le K(\mathbb{E}\|X\|_2)^{1/2}$$
 almost surely. (5.16)

Then, for every positive integer m, we have

$$\mathbb{E}\|\Sigma_m - \Sigma\| \leq C\left(\sqrt{\frac{K^2 n \log n}{m}} + \frac{K^2 n \log n}{m}\right) \|\Sigma\|.$$

### Proof

### Proof cont.