

# Lecture 17

## Last time

- ▷ Covariance Estimation
- ▷ Clustering a Gaussian Mixture.

## Today

- ▷ Matrix Calculus
- ▷ Matrix Chernoff
- ▷ Matrix Hoeffding

## Matrix Calculus

Our goal today is to show the following extension of Hoeffding's inequality for matrices

$$P\left(\left\| \sum_{i=1}^n a_i \mathbf{x}_i \right\|_p \geq t\right) \leq \dots$$

↑ Ind. mean-zero

We'll need some notation.

Def (Functions of matrices): Given a function  $f: \mathbb{R} \rightarrow \mathbb{R}$  and a  $X \in \mathbb{S}^d$  with spectral decomposition

$$X = \sum_{i=1}^d \lambda_i u_i u_i^\top \quad \text{define } f(X) := \sum_{i=1}^d f(\lambda_i) u_i u_i^\top.$$

### Examples:

- ▷ If  $f(t) = t^{-1} \Rightarrow f(X) = X^{-1}$
- ▷ If  $f(t) = \sum_{p=1}^{\infty} \alpha_p t^p \Rightarrow f(X) = \sum_{p=1}^{\infty} \alpha_p X^p$

We can Taylor expand!

Just as with scalars, symmetric matrices also have an order.

Def (Löwner ordering): Given  $X, Y \in S^d$  we say that  $X \leq Y$  if  $Y - X \in S_+^d$ .  $\dashv$

Lemma (\*) Suppose that  $X \leq Y$ , then:

1) For all  $A \in \mathbb{R}^{d \times k}$ ,

$$A^T X A \leq A^T Y A.$$

2) All eigenvalues satisfy

*i*th largest  $\lambda_i(X) \leq \lambda_i(Y)$

3) Let  $f: \mathbb{R} \rightarrow \mathbb{R}$  be nonincreasing, then  
 $\text{tr}(f(X)) \leq \text{tr}(f(Y))$ .  $\dashv$

Given the last item it is natural to wonder whether

(\*)  $X \leq Y \Rightarrow f(X) \leq f(Y)$

for nonincreasing  $f$ . In general this is not the case. Indeed,  
 $t \mapsto e^t$

is a counterexample (check!)

Def (Matrix monotone) A function is matrix monotone if (1) holds  $\forall X, Y \succcurlyeq$

Lemma (2) The functions  $t \mapsto t^{-1}$ ,  $t \mapsto t^{1/2}$ , and  $t \mapsto \log t$  are matrix monotone.

## Matrix Chernoff

Recall the strategy we use for scalar random variables:

- ① Sub-Gaussian MGF  $\Rightarrow$  Tail bounds  
Markov's
- ② Sum of iid sub-Gaussian  $\Rightarrow$  Sum is also sub-Gaussian.  
Peeling argument

In turn, the peeling argument is delicate. Define the MGF of a random matrix  $a$  as  $\Psi_a : \mathbb{R} \rightarrow \mathbb{S}^d$  given by

$$\Psi_a(\gamma) = \mathbb{E}[e^{\gamma a}] = \sum_{k=0}^{\infty} \frac{\gamma^k}{k!} \mathbb{E}[a^k].$$

Just as in the scalar case this controls concentration.

**Lemma (Matrix Chernoff Method):** Let  $A$  a random sym. matrix with MGF defined in some interval  $(-\alpha, \alpha)$ . Then for all  $t > 0$ , we have

$$\mathbb{P}(\lambda_1(A) \geq t) \leq \text{tr}(\Psi_A(\gamma)) e^{-\gamma t} \quad \forall \gamma \in [0, \alpha].$$

As a consequence

$$\mathbb{P}(\|A\|_{\text{op}} \geq t) \leq 2 \text{tr}(\Psi_A(\lambda)) e^{-\lambda t} \quad \forall \lambda \in [\max\{\lambda_1(A), -\lambda_n(A)\}],$$

**Proof:** Taking scalar exponentials

$$\begin{aligned} \mathbb{P}(\lambda_1(A) \geq t) &= \mathbb{P}(e^{\lambda_1(A)} \geq e^{\gamma t}) \\ &= \mathbb{P}(\lambda_1(e^{\gamma A}) \geq e^{\gamma t}) \\ \text{Markov's} \quad \text{tr}(A) = \sum \lambda_i(A) &\leq \mathbb{E} \lambda_1(e^{\gamma A}) e^{-\gamma t} \\ &\leq \mathbb{E} \text{tr}(e^{\gamma A}) e^{-\gamma t} \\ &= \text{tr} \Psi_A(\gamma) e^{-\gamma t}. \end{aligned}$$

□

This seems to suggest that we want  $\Psi_A(\gamma)$  to be nicely controlled by a "Gaussian tail" as before.

Def: A random symmetric matrix is V-sub-Gaussian if

$$\Psi_Q(\lambda) \leq e^{\lambda^2/2} \quad \forall \lambda \in \mathbb{R}. \quad \dagger$$

Example: Suppose  $Q = \epsilon B$  with  $\epsilon \sim \text{Unif}\{\pm 1\}$  and  $B \in S^d$  fixed.

Then  $\mathbb{E} Q^{2k+1} = 0$  and  $\mathbb{E} Q^{2k} = B^{2k}$ . So,

$$\mathbb{E} e^{\lambda Q} = \sum_{k=0}^{\infty} \frac{\lambda^{2k}}{(2k)!} B^{2k} \stackrel{?}{\leq} \sum_{k=0}^{\infty} \frac{1}{k!} \left( \frac{\lambda^2 B^2}{2} \right)^k = e^{\lambda^2 B^2 / 2}. \quad \dagger$$

(why?)

Unlike before we don't have a sum rule because

$$e^{A+B} \neq e^A e^B$$

in general (it only holds for commuting matrices). To "fix" the peeling argument we will use a deep result from analysis.

Theorem (Lieb inequality):

Let  $H \in S^d$ , and define  $f: S_+^d \rightarrow \mathbb{R}$  given by

$$f(X) = \text{tr} \exp(H + \log X).$$

Then,  $f$  is concave on  $S_+^d$ . +

We will not prove this result.

Lemma Let  $a_1, \dots, a_n \in S^d$  be independent with  $\psi_{a_i}(\cdot)$  defined over an interval  $J \subseteq \mathbb{R}$ . Let

$$S_n = \sum_{i=1}^n a_i. \text{ Then,}$$

$$\text{tr}(\psi_{S_n}(\gamma)) \leq \text{tr}(\exp(\sum_{i=1}^n \log \psi_{a_i}(\gamma)))$$

Consequently, <sup>Chernoff method</sup>  $\forall \gamma \in J$ .

$$\mathbb{P}\left(\|\frac{1}{n} \sum a_i\|_{op} \geq t\right) \leq 2 \text{tr}(e^{\sum \log \psi_{a_i}(\gamma)}) e^{-\gamma t}.$$

Proof: Expanding

$$\begin{aligned} \text{tr}(\psi_{S_n}(\gamma)) &= \text{tr} \mathbb{E} e^{\gamma S_n} \\ &= \text{tr} \mathbb{E} e^{\gamma S_{n-1} + \log \exp(\gamma a_n)} \\ &= \mathbb{E}_{S_{n-1}} \mathbb{E}_{a_n} \text{tr} e^{\gamma S_{n-1} + \log \exp(\gamma a_n)} \\ \xrightarrow{\text{Lieb + Jensen's}} &\leq \mathbb{E}_{S_{n-1}} \text{tr} e^{\gamma S_{n-1} + \log \psi_{a_n}(\gamma)} \end{aligned}$$

Recurising  
the same  
argument

$$\stackrel{\leq}{\dots} \leq \text{tr } e^{\sum_{i=1}^n \log \Psi_{\alpha_i}(\gamma)}.$$
□

**Theorem (Hoeffding):** Suppose  $\alpha_1, \dots, \alpha_n$  are zero-mean,  $V_i$ -sub-Gaussian random matrices in  $S^d$ . Then,

$$\begin{aligned} \mathbb{P}\left(\left\|\frac{1}{n} \sum_{i=1}^n \alpha_i\right\|_{\text{op}} \geq t\right) &\leq 2 \text{rank}(\sum V_i) e^{-\frac{n t^2}{2\sigma^2}} \\ &\leq 2d e^{-\frac{n t^2}{2\sigma^2}} \quad \forall t > 0, \end{aligned}$$

where  $\sigma^2 = \left\|\frac{1}{n} \sum V_i\right\|_{\text{op}}$ .

**Proof:** Let  $V = \sum_{i=1}^n V_i$ . From Lemma (10), it suffices to bound  $\text{tr}(\exp(\sum \log \Psi_{\alpha_i}(\lambda)))$ . By sub-Gaussianity and the monotonicity of the matrix  $\log$  (Lemma (5))

$$\sum_{i=1}^n \log \Psi_{\alpha_i}(\gamma) \leq \frac{\gamma^2}{2} \sum_{i=1}^n V_i.$$

Moreover since  $t \mapsto e^t$  is increasing

Lemma (\*) gives

$$\text{tr}(\exp(\sum_{i=1}^n \log \varphi_{a_i}(\gamma))) \leq 2 \text{tr}(e^{\frac{\gamma^2}{2} V}).$$

Thus, by Lemma (1),

$$P\left(\left\|\frac{1}{n} \sum a_i\right\|_{\text{op}} \geq t\right) \leq 2 \text{tr}(e^{\frac{\gamma^2 V}{2}}) e^{-\gamma n t}.$$

Note that

$$\text{tr}(e^A) \leq \text{rank}(A) e^{\|A\|_{\text{op}}},$$

moreover  $\left\|\frac{\gamma^2 V}{2}\right\|_{\text{op}} = \frac{\gamma^2}{2} n \sigma^2$ . So

$$P\left(\left\|\frac{1}{n} \sum a_i\right\|_{\text{op}} \geq t\right) \leq 2 \text{rank}(V) e^{\frac{\gamma^2 n \sigma^2 - \gamma n t}{2}}.$$

The best bound is given by taking  $\gamma = t/\sigma^2$ , which yields the claim.  $\square$

**Remark:** The additional  $d$  factor in the bound is in general unavoidable.