

Chapter 4

Metric Entropy and Matrix Concentration

- Nets, covering numbers, and metric entropy
- Eigenvalues and Singular values of random matrices
- Matrix Concentration
- Operator norm vs. row/column norms
- Applications: community detection,
 - (sparse) covariance matrix estimation
 - matrix completion

Let (T, ρ) be a metric space

Ex: $(\mathbb{R}^d, \|\cdot\|_2)$ with $\rho(x, y) = \|x - y\|_2$

• $(\{0, 1\}^d, \rho_H)$ with Hamming metric

$$\rho_H(\theta, \tilde{\theta}) = \frac{1}{d} \sum_{j=1}^d \mathbb{1}[\theta_j \neq \tilde{\theta}_j]$$

• $L^2(\mu, [0, 1])$ with

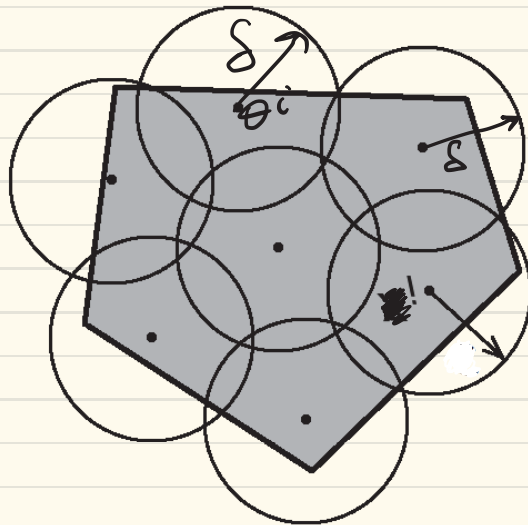
$$\|f - g\|_2 = \sqrt{\int_0^1 (f(x) - g(x))^2 d\mu(x)}$$

• $C[0, 1]$ with

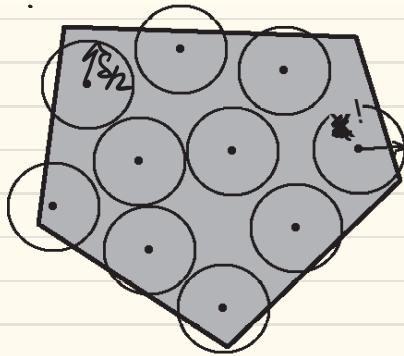
$$\|f - g\|_\infty = \sup_{x \in [0, 1]} |f(x) - g(x)|$$

Defn: A δ -cover of a set T with respect to a metric ρ is a set $\{\theta^1, \dots, \theta^N\} \subset T$ such that

$\forall \theta \in T \exists i \in \{1, \dots, N\}$ s.t. $\rho(\theta, \theta^i) \leq \delta$
The δ -covering number $N(\delta; T, \rho)$ is the cardinality of the smallest δ -cover. The quantity $\log(N(\delta; T, \rho))$ is called the metric entropy.



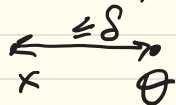
Defn: A δ -packing of T with respect to ρ is a set $\{\theta^1, \dots, \theta^m\} \subseteq T$ such that $\rho(\theta^i, \theta^j) > \delta \quad \forall i, j \in \{1, 2, \dots, m\}$.
 The δ -packing number $M(\delta; T, \rho)$ is the cardinality of the largest δ -packing.



Lemma:

$$M(2\delta; T, \rho) \stackrel{(a)}{\leq} N(\delta; T, \rho) \stackrel{(b)}{\leq} M(\delta; T, \rho)$$

pf: (b) $M(\delta; T, \rho)$



$$= M(\delta; T, \rho) \geq N(\delta; T, \rho)$$

(a) $M(2\delta; T, \rho)$

No δ Ball can cover θ, θ^i simultaneously. δ

Lemma: Let $\|\cdot\|, \|\cdot\|'$ be norms on \mathbb{R}^d
 and set $B = \{x: \|x\| \leq 1\}$
 $B' = \{x: \|x\|' \leq 1\}$

Then

$$\left(\frac{1}{\delta}\right)^d \frac{\text{vol}(B)}{\text{vol}(B')} \stackrel{\textcircled{a}}{\leq} N(\delta; B, \|\cdot\|) \stackrel{\textcircled{b}}{\leq} \frac{\text{vol}\left(\frac{2}{\delta}B + B'\right)}{\text{vol}(B')}$$

Remark: If $B' \subseteq B$, then simplifies

$$N(\delta; B, \|\cdot\|) \leq \left(1 + \frac{2}{\delta}\right)^d \frac{\text{vol}(B)}{\text{vol}(B')}$$

pf of lemma: If $\{\theta^1, \dots, \theta^n\}$ is a δ -covering
 of B , then $B \subseteq \bigcup_{i=1}^n \{\theta^i + \delta B'\}$

$$\Rightarrow \text{vol}(B) \leq N \delta^d \text{vol}(B') \Rightarrow \textcircled{a}$$

Let $\{\theta^1, \dots, \theta^m\}$ be a maximal $\delta/2$ -packing
 of B . \Rightarrow must also be δ -covering.

Taking into account:

$$\left\{ \emptyset + \frac{\delta}{2} B' \right\} \subseteq B + \frac{\delta}{2} B'$$

are disjoint

get:

$$M \text{vol} \left(\frac{\delta}{2} B' \right) \leq \text{vol} \left(B + \frac{\delta}{2} B' \right)$$

$$\left(\frac{\delta}{2} \right)^d M \text{vol}(B') \leq \left(\frac{\delta}{2} \right)^d \text{vol} \left(\frac{2}{\delta} B + B' \right)$$

~~□~~

Cor: If $\|\cdot\|' = \|\cdot\|$, then

$$d \log \left(\frac{1}{\delta} \right) \leq \log N \left(\frac{1}{\delta}; B, \|\cdot\| \right) \leq d \log \left(1 + \frac{2}{\delta} \right)$$

Metric entropy can be very large...

Ex: Define

$$\mathcal{F}_L([0,1]^d) = \left\{ f: [0,1]^d \rightarrow \mathbb{R} : f(0) = 0, |f(x) - f(x')| \leq L \|x - x'\|_\infty \right\}$$

Then $\log N \left(\frac{1}{\delta}; \mathcal{F}_L([0,1]^d), \|\cdot\|_\infty \right) \sim \left(\frac{L}{\delta} \right)^d$

See ex 5.10 in Wainwright.

One of the main uses of covering numbers is for random matrix theory.

Review: Any $A \in \mathbb{R}^{m \times n}$ can be written as

$$A = \sum_{i=1}^r s_i u_i v_i^T = \begin{bmatrix} u_1 & \dots & u_r \\ \vdots & & \vdots \end{bmatrix} \begin{bmatrix} s_1 & \dots & s_r \\ & & \\ & & \end{bmatrix} \begin{bmatrix} v_1^T \\ \vdots \\ v_r^T \\ \vdots \\ \vdots \end{bmatrix}$$

where $r = \text{rank}(A)$, $s_i \geq 0$, $\|u_i\|_2 = \|v_i\|_2 = 1$

• Decomposition is called a singular value decomposition.

• s_i are the singular values. We define $s_{r+1}, \dots, s_n = 0$ and arrange

$$s_1 \geq s_2 \geq s_3 \geq \dots \geq s_n \geq 0$$

• $s_i(A) = \sqrt{\lambda_i(AA^T)} = \sqrt{\lambda_i(A^T A)}$ $\forall i=1, \dots, r$.

The set $\mathbb{R}^{m \times n}$ admits the inner-product

$$\langle A, B \rangle := \sum_{i,j} A_{ij} B_{ij} = \text{tr}(A^T B)$$

and the norms

$$\|A\|_F := \sqrt{\sum_{i,j} A_{ij}^2} = \sqrt{\langle A, A \rangle} = \|s(A)\|_2$$

$$\|A\|_2 = \sup_{\|x\|_2=1} \|Ax\|_2 = s_1(A) = \|s(A)\|_\infty$$

Remark: $s_1(A)$ and $s_n(A)$ are respectively the smallest \bar{m} and largest \underline{m} with

$$\underline{m} \|x\|_2 \leq \|Ax\|_2 \leq \bar{m} \|x\|_2 \quad \forall x \in \mathbb{R}^n$$

Thm (Eckart-Young-Mirsky)

Let $A = \sum_{i=1}^n s_i u_i v_i^T$ be an SVD of A

Then the matrix $A' = \sum_{i=1}^r s_i u_i v_i^T$ solves the problem

$$\min_{A': \text{rank}(A') \leq r} \|A - A'\|$$

where $\|\cdot\|$ is either $\|\cdot\|_F$ or $\|\cdot\|_2$.

We next aim to show that for $A \in \mathbb{R}^{m \times n}$ with independent subGaussian entries, it holds:

with high probability, $\|A\|_2 \leq \sqrt{m} + \sqrt{n}$.

Lemma: Let $A \in \mathbb{R}^{m \times n}$ and $\epsilon \in [0, \sqrt{2}-1]$. Then for any ϵ -nets \mathcal{N} of S^{n-1} and \mathcal{M} of S^{m-1} , it holds:

$$\sup_{x \in \mathcal{N}, y \in \mathcal{M}} \langle Ax, y \rangle \stackrel{\textcircled{a}}{\leq} \|A\|_2 \stackrel{\textcircled{b}}{\leq} \frac{1}{1-2\epsilon-\epsilon^2} \sup_{x \in \mathcal{N}, y \in \mathcal{M}} \langle Ax, y \rangle$$

pt. \textcircled{a} is clear. To see \textcircled{b} , suppose

$\|A\|_2 = \langle Ax, y \rangle$ for some $x \in S^{n-1}, y \in S^{m-1}$.
 $\exists \hat{x} \in \mathcal{N}, \hat{y} \in \mathcal{M}$ with $\|x - \hat{x}\| \leq \epsilon, \|y - \hat{y}\| \leq \epsilon$

$$\begin{aligned} \Rightarrow \langle Ax, y \rangle &= \langle A\hat{x} - A(\hat{x} - x), \hat{y} - (\hat{y} - y) \rangle \\ &= \langle A\hat{x}, \hat{y} \rangle - \langle A\hat{x}, \hat{y} - y \rangle - \langle A(\hat{x} - x), \hat{y} \rangle \\ &\quad + \langle A(\hat{x} - x), \hat{y} - y \rangle \end{aligned}$$

$$\begin{aligned} &\leq \langle A\hat{x}, \hat{y} \rangle + 2\epsilon \|A\|_2 + \|A\|_2 \epsilon^2 \\ \Rightarrow \|A\|_2 &\leq \frac{1}{1-2\epsilon-\epsilon^2} \langle A\hat{x}, \hat{y} \rangle \quad \square \end{aligned}$$

Thm: Let $A \in \mathbb{R}^{m \times n}$, where A_{ij} are independent mean-zero σ_{ij} -subgaussian.

Then \exists constant $C > 0$ s.t.

$$\mathbb{P}[\|A\|_2 \leq C \cdot (\max_{i,j} \sigma_{ij}) \cdot (\sqrt{m} + \sqrt{n} + t)] \geq 1 - 2e^{-t^2}$$

[Remark: when $A_{ij} \sim N(0,1)$, optimal]
constant is $C=1$]

pt: Step 1: Approximation

We know $N(S^{n-1}, \epsilon, \|\cdot\|_2) \leq \left(1 + \frac{2}{\epsilon}\right)^n$

$$N(S^{m-1}, \epsilon, \|\cdot\|_2) \leq \left(1 + \frac{2}{\epsilon}\right)^m$$

\Rightarrow Set $\epsilon = \frac{1}{4}$. There are nets

$$\text{card}(\mathcal{N}) \leq 9^n, \quad \text{card}(\mathcal{M}) \leq 9^m$$

\Rightarrow Lemma $\|A\|_2 \leq 3 \cdot \max_{x \in \mathcal{N}, y \in \mathcal{M}} \langle Ax, y \rangle$

Step 2: Concentration

The RV $\langle A_{x,y} \rangle = \sum_{i,j} A_{ij} x_i y_j$ satisfies

$$\begin{aligned} e^{\lambda \langle A_{x,y} \rangle} &= \prod_{i,j} e^{(\lambda x_i y_j) A_{ij}} \leq \prod_{i,j} e^{\lambda^2 \sigma_{ij}^2 x_i^2 y_j^2} \\ &= e^{\lambda^2 \left(\max_{i,j} \sigma_{ij} \right)^2 \sum_{i,j} x_i^2 y_j^2} = e^{\lambda^2 \max_{i,j} \sigma_{ij}^2} \end{aligned}$$

So Hoeffding \Rightarrow
$$P[\langle A_{x,y} \rangle \geq u] \leq e^{-\frac{u^2}{2 \max_{i,j} \sigma_{ij}^2}} \quad \forall u > 0$$

$$\Rightarrow P \left[\max_{\substack{x \in \mathcal{X} \\ y \in \mathcal{Y}}} \langle A_{x,y} \rangle \geq u \right] \leq 9^{n+m} e^{-\frac{u^2}{2 \max_{i,j} \sigma_{ij}^2}}$$

Choose $u = C \left(\max_{i,j} \sigma_{ij} \right) (\sqrt{n} + \sqrt{m} + t)$

Then

$$\begin{aligned} \text{RHS} &\leq 9^{n+m} \cdot e^{-C(n+m+t^2)} \\ &= e^{\ln(9)(n+m) - C(n+m+t^2)} \leq e^{-t^2} \end{aligned}$$

by choosing C large \square

Cor: (HW) $\mathbb{E} \|A\|_2 \leq C (\max_{i,j} \sigma_{ij}) (\sqrt{m} + \sqrt{n})$

Obvious analogues hold for symmetric matrices.

Application: Community Detection

Defn: (Stochastic Block Model)

Divide n vertices into two sets Q_1, Q_2 of size $n/2$ each. Build the graph

$$\mathbb{P}[ij \text{ is an edge}] = \begin{cases} p & \text{if } i, j \in Q_1 \text{ or } i, j \in Q_2 \\ q & \text{o.w} \end{cases}$$



The goal is to find the communities from the realization of the graph.

Let A be the adjacency matrix.
Write it as

$$A = D + R \quad \text{where}$$

$$D = EA = \left[\begin{array}{cc|cc} p & p & q & q \\ p & p & q & q \\ \hline q & q & p & p \\ q & q & p & p \end{array} \right]$$

Exercise: (Hw)

Show that D has rank 2 and

$$\lambda_1 = \left(\frac{p+q}{2}\right)n, \quad u_1 = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} / \sqrt{n}$$

$$\lambda_2 = \left(\frac{p-q}{2}\right)n, \quad u_2 = \begin{bmatrix} 1 \\ \vdots \\ -1 \end{bmatrix} / \sqrt{n}$$

If we knew u_2 , we would be able to identify the communities.

Signal to noise ratio (SNR): $\frac{\|D\|_2}{\|R\|_2}$

We know $\|D\|_2 = \lambda_1 \sim n$

$$\mathbb{P}[\|R\| \leq C\sqrt{n}] \geq 1 - 4e^{-n}$$

So $SNR \sim \sqrt{n}$

\Rightarrow problem should be solvable.

A natural idea is to use spectral clustering.

Alg: (Spectral Clustering)

Input: graph G

1) Compute adjacency matrix A

2) Compute eigenvector $v_2(A)$ corresponding to second largest eigenvalue of A .

3) Set $Q_1 = \{i : [v_2(A)]_i \geq 0\}$
 $Q_2 = \{i : [v_2(A)]_i \leq 0\}$.

Perturbation Results:

Thm: (Weyl's inequality)

For any symmetric $X, Y \in \mathbb{R}^{n \times n}$, it holds

$$\max_{i=1, \dots, n} |\lambda_i(X) - \lambda_i(Y)| \leq \|X - Y\|_2$$

Thm (Davis-Kahan)

Let $X, Y \in \mathbb{R}^{n \times n}$ be symmetric. Fix an index i and assume

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

Then $\sin \angle(v_i(X), v_i(Y)) \leq \frac{2\|X - Y\|_2}{\delta}$

Remark: This implies

$$\min_{\Theta \in \mathbb{R}^{\pm 1}} \|v_i(X) - \Theta v_i(\Theta)\|_2 \leq \frac{2\|X - Y\|_2^{3/2}}{\delta}$$

Let's compute the eigengap

$$\delta = \min \{ \lambda_2, \lambda_1 - \lambda_2 \} = \min \left\{ \frac{p-q}{2}, q \right\} \cdot n$$

$$\Rightarrow \min_{\Theta \in \{\pm 1\}} \|V_2(D) - \Theta V_2(A)\|_2 \leq \frac{\|D - A\|_2}{\mu} = \frac{\|R\|_2}{\mu n}$$

$$\Rightarrow \mathbb{P} \left[\exists_{\Theta \in \{\pm 1\}} \|V_2(D) - \Theta V_2(A)\|_2 \leq \frac{c}{\mu \sqrt{n}} \right] \geq 1 - 4e^{-n}$$

Since all coordinates of $\Theta V_2(D)$ are $\frac{\pm 1}{\sqrt{n}}$, the signs can only disagree in at most $\frac{c^2}{\mu^2}$ coordinates.

Thm: Set $\mu = \min \{ q, \frac{p-q}{2} \}$. Then with probability $1 - 4e^{-n}$, the spectral clustering algorithm identifies the communities correctly up to $\left(\frac{c}{\mu}\right)^2$ misclassified vertices.

We next obtain two-sided bounds on the full spectrum

$$\sqrt{m} - C\sqrt{n} \leq \lambda_i(A) \leq \sqrt{m} + C\sqrt{n}$$

and we relax independence of entries to independence of the rows.

Thm: Let $A \in \mathbb{R}^{m \times n}$ be a matrix with rows A_i that are independent and σ_i -sub-Gaussian isotropic. Then

$$\sqrt{m} - CK^2(\sqrt{n} + t) \leq \lambda_n(A) \leq \lambda_1(A) \leq \sqrt{m} + CK^2(\sqrt{n} + t)$$

w.p. $1 - 2\exp(-t^2)$, where $K = \max_i \sigma_i$.
pt: We first show

$$\left\| \frac{1}{m} A^T A - I_n \right\|_2 \leq K^2 \max\{\delta, \delta^2\} \text{ where}$$

$$\delta = C \left(\sqrt{\frac{n}{m}} + \frac{t}{\sqrt{m}} \right)$$

Step 1: As before, there exist a $\frac{1}{4}$ -cover \mathcal{N} of S^{n-1} with $\text{card}(\mathcal{N}) \leq 9^n$.

For $B = \frac{1}{m} A^T A - I$ have $\|B\|_2 = \sup_{v \in \mathcal{S}^{h-1}} \langle Bv, v \rangle$
 If v achieves sup. Then $\exists w$ with $\|w-v\| \leq \frac{1}{4}$
 and

$$\begin{aligned} \|B\|_2 &= \langle Bw - B(w-v), w - (w-v) \rangle \\ &= \langle Bw, w \rangle - 2\langle Bw, w-v \rangle + \langle B(w-v), w-v \rangle \\ &\leq \sup_{w \in \mathcal{R}} \langle Bw, w \rangle + \frac{1}{4} \|B\|_2 + \frac{1}{16} \|B\| \end{aligned}$$

$$\Rightarrow \|B\|_2 \leq 2 \sup_{w \in \mathcal{R}} \langle Bw, w \rangle$$

$$= 2 \left(\sup_{w \in \mathcal{R}} \frac{1}{m} \|Aw\|_2^2 - 1 \right)$$

Step 2: Concentration

Fix $w \in \mathcal{R}$ and write

$$\|Aw\|_2^2 = \sum_{i=1}^m \langle A_i, w \rangle^2 =: \sum_{i=1}^m X_i^2$$

$\Rightarrow X_i$ are independent, ϕ_i -subGaussian with $\mathbb{E} X_i^2 = 1$. Therefore, Bernstein

$$P\left[\left| \frac{1}{m} \|Aw\|_2^2 - 1 \right| \geq t \right] \leq 2 \exp\left[-\frac{1}{2} \left(\frac{mt^2}{K^4} \wedge \frac{mt}{4K^2} \right) \right]$$

$$2 \exp\left[-\frac{1}{2} \left(\frac{mt^2}{K^4} \wedge \frac{mt}{4K^2} \right)\right] = 2 \exp\left[-c, m \left[\frac{t^2}{K^4} \wedge \frac{t}{K^2} \right]\right]$$

$$\text{Set } t = K^2 \max\{\delta, \delta^2\}$$

$$\rightarrow = 2 \exp\{-c, m \delta^2\}$$

\leq

Step 3: Union bound

$$\begin{aligned} & \mathbb{P}\left[\max_{W \in \mathcal{N}} \left| \frac{1}{m} \|AW\|_2^2 - 1 \right| \geq K^2 \max\{\delta, \delta^2\}\right] \\ & \leq 9 \cdot 2 \exp\{-c, C^2(n+t^2)\} \leq 2e^{-t^2} \end{aligned}$$

if C is big enough.

Finally result follows from

Lemma: Suppose $A \in \mathbb{R}^{m \times n}$ and $\delta > 0$ satisfy

$$\|A^T A - I\|_2 \leq \max\{\delta, \delta^2\}$$

Then

$$1 - \delta \leq \lambda_n(A) \leq \lambda_1(A) \leq 1 + \delta$$

□

Cor: (HW)

$$\mathbb{E} \left\| \frac{1}{m} A^T A - I \right\|_2 \leq CK^2 \left(\sqrt{\frac{n}{m}} + \frac{n}{m} \right)$$

Covariance Estimation

Assume $X \in \mathbb{R}^n$ and $\mathbb{E}X = 0$

Define $\Sigma = \mathbb{E}XX^T$. Let's estimate Σ

by the sample covariance

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

Clearly for $\Sigma - \Sigma_m$ to be small we need $m \gg n$. We'll see $m \sim n$ suffice

Thm: Suppose $X \in \mathbb{R}^n$ satisfies

$$\mathbb{E}_x e^{\lambda \langle X, w \rangle} \leq e^{K^2 \lambda^2 \langle \Sigma w, w \rangle} \quad \forall w \in \mathbb{R}^n$$

for some $K \geq 0$.

Then

$$\mathbb{E} \left\| \frac{\Sigma_m - \Sigma}{\|\Sigma\|_2} \right\|_2 \leq CK^2 \left(\sqrt{\frac{n}{m}} + \frac{n}{m} \right)$$

pf: Fix random X, X_1, \dots, X_m . Define

$$Z = \Sigma^{-1/2} X \quad Z_i = \Sigma^{-1/2} X_i$$

Recall Z, Z_i are independent and isotropic. Moreover for any $w \in \mathbb{S}^{n-1}$, have

$$\begin{aligned} \mathbb{E} e^{\lambda \langle Z, w \rangle} &= \mathbb{E} e^{\lambda \langle \Sigma^{-1/2} X, w \rangle} \\ &= \mathbb{E} e^{\lambda \langle X, \Sigma^{-1/2} w \rangle} \leq e^{\frac{1}{2} \lambda^2 K^2 \langle \Sigma^{-1/2} w, \Sigma^{-1/2} w \rangle} \\ &= e^{\frac{1}{2} \lambda^2 K^2} \end{aligned}$$

So Z and Z_i are K -sub-Gaussian

$$\begin{aligned} \Rightarrow \|\Sigma_m - \Sigma\|_2 &= \left\| \Sigma^{1/2} \left(\frac{1}{m} \sum_{i=1}^m Z_i Z_i^T - I_n \right) \Sigma^{1/2} \right\|_2 \\ &\leq \|\Sigma\|_2 \cdot \left\| \frac{1}{m} \sum_{i=1}^m Z_i Z_i^T - I \right\|_2 \end{aligned}$$

Define A whose rows are Z_i . Then

$$\frac{1}{m} A^T A - I = \frac{1}{m} \sum_{i=1}^m Z_i Z_i^T - I. \quad \square$$

So can get

$$\mathbb{P} \|\Sigma'_m - \Sigma\|_2 \leq \varepsilon \|\Sigma\|_2$$

as long as $m \sim \frac{n}{\varepsilon^2}$

Exc: (HW) Actually get the high probability estimate

$$\|\Sigma'_m - \Sigma\|_2 \leq CK^2 \left(\sqrt{\frac{n+t}{m}} + \frac{n+t}{m} \right) \|\Sigma\|_2$$

with probability $1 - 2e^{-t}$

Remark: In particular, setting $t = \delta m$
get

$$\frac{\|\Sigma'_m - \Sigma\|_2}{\|\Sigma\|_2} \leq CK^2 \left(\sqrt{\frac{n}{m}} + \frac{n}{m} + \delta \right)$$

with probability $1 - 2e^{-\delta m}$

Clustering (Gaussian Mixture Model)

Generate m random points in \mathbb{R}^n as follows.

Flip a fair coin

- if heads choose $x_i \sim \mathcal{N}(\mu, I)$
- if tails choose $x_i \sim \mathcal{N}(-\mu, I)$

Suppose we are given m points drawn from the Gaussian Mixture Model.

Goal: Identify which points belong to which cluster.

Spectral Clustering:

Input: $x_1, \dots, x_m \in \mathbb{R}^n$

Output: a partition of points into 2 clusters.

- 1) Compute $\Sigma_m = \frac{1}{m} \sum_{i=1}^m x_i x_i^T$
- 2) Compute eigenvector of Σ_m corresponding to $\lambda_1(\Sigma)$.
- 3) If $\langle v, x_i \rangle > 0$, put x_i in first community otherwise, put x_i in second.

Thm: Let $\delta > 0$ be s.t. $\|u\|_2 \geq C \sqrt{\log\left(\frac{n}{\delta}\right)}$.

Suppose $m \geq \left(\frac{n}{\|u\|_2}\right)^c$ where c is a constant. Then with probability $1 - 4e^{-n}$, the algorithm identifies communities correctly up to δm misclassified vertices.

pf sketch: Note $x = \epsilon u + g$ where

ϵ is a Rademacher RV and $g \sim \mathcal{N}(0, I)$

$$\Rightarrow \mathbb{E} x x^T = \mathbb{E} (u u^T + g g^T) = I + u u^T$$

$\Rightarrow \mathbb{E} x x^T$ has eigenvalues $1 + \|u\|^2, 1, \dots, 1$

Eigenvector corresponding to $1 + \|u\|^2$ is u

So from the e-vector corresponding to the largest eigenvalue, you can "guess" at the communities by checking the sign of $\langle x_i, u \rangle$.

You do the rest in your HW \square

Our next goal is to investigate tails of matrices generalizing bounds

$$\left\| \frac{1}{m} \sum_{i=1}^m x_i x_i^T - \Sigma \right\|_2 \text{ from before.}$$

Notation:

- \mathcal{S}^d are $d \times d$ symmetric matrices
- $\mathcal{S}_+^d = \{ Q \in \mathcal{S}^d : \langle Q, x \rangle \geq 0 \ \forall x \in \mathbb{R}^d \}$
- $\lambda_1(X) \geq \lambda_2(X) \geq \dots \geq \lambda_d(X)$ are the eigenvalues of $X \in \mathcal{S}^d$
- $\mathcal{O}^d = \{ U \in \mathbb{R}^{d \times d} : U^T U = I \}$

Any function $f: \mathbb{R} \rightarrow \mathbb{R}$ defines a function $f: \mathcal{S}^d \rightarrow \mathcal{S}^d$ by

$$f(Q) = U \operatorname{diag}(f(\lambda_1(Q)), \dots, f(\lambda_d(Q))) U^T$$

where $Q = U \operatorname{diag}(\lambda_1(Q)) U^T$ is any eigenvalue decomposition of Q .

[The choice of U does not matter (check!)]

Note

$$\chi(f(Q)) = \{f(\chi_j(Q)) : j=1, \dots, d\}$$

Matrix exponential: e^Q

Matrix Logarithm: $\log Q$ defined on S_{++}^d

The Löwner ordering on

$$X \preceq Y \Leftrightarrow Y - X \in S_+^d$$

Lemma: (HW) $Q_2 \succeq Q_1, \gamma > 0 \Rightarrow \log(Q_2) \preceq \log(Q_1)$

[e^Q is not monotone]

Lemma: If $f: \mathbb{R} \rightarrow \mathbb{R}$ is continuous and non-decreasing, then

$$Q \preceq R \Rightarrow \text{tr}(f(Q)) \leq \text{tr}(f(R))$$

Lemma: $I + A \succeq e^A$

Thm (Lieb): For any $H \in S^d$, the function $f: S_{++}^d \rightarrow \mathbb{R}$ given by $f(A) = \text{tr}(e^{H + \log(A)})$ is concave.

Defn: Moment generating function of a random matrix $Q \in \mathcal{S}^d$ is the function $\Psi_Q: \mathbb{R} \rightarrow \mathcal{S}^d$ given by

$$\Psi_Q(\lambda) = \mathbb{E} e^{\lambda Q} = \sum_{k=0}^{\infty} \frac{\lambda^k}{k!} \mathbb{E}[Q^k]$$

Defn: A zero-mean $Q \in \mathcal{S}_+^d$ is V -subGaussian (with $V \in \mathcal{S}_+^d$) if

$$\Psi_Q(\lambda) \preceq e^{\frac{\lambda^2 V}{2}} \quad \forall \lambda \in \mathbb{R}$$

Ex: Suppose $Q = \varepsilon A$ where ε is a Rademacher RV and $A \in \mathcal{S}^d$ is fixed.

$$\begin{aligned} \mathbb{E} e^{\lambda Q} &= \sum_{k=0}^{\infty} \frac{\lambda^k}{k!} \mathbb{E}(\varepsilon A)^k = \sum_{k=0}^{\infty} \frac{\lambda^{2k}}{(2k)!} A^{2k} \\ &\preceq \sum_{k=1}^{\infty} \frac{1}{k!} \left(\frac{\lambda^2 A^2}{2} \right)^k \\ &= e^{\frac{\lambda^2 A^2}{2}} \end{aligned}$$

So Q is A^2 -subGaussian.

Exercise (HW)

If $Q = gA$, where $g \in \mathbb{R}$ is a symmetric ^{distribution} σ -subGaussian variable, then Q is $\sigma^2 A^2$ subGaussian.

Ex: Suppose $Q = \varepsilon A$ where ε is Rademacher and A is random with $\|A\|_2 \leq b$. Then (why?)

$$\Psi_Q(\lambda) = \mathbb{E} e^{\lambda Q} = \mathbb{E}_A \mathbb{E}_\varepsilon e^{\varepsilon A} \leq \mathbb{E}_A e^{\frac{\lambda^2 A^2}{2}} \leq e^{\frac{\lambda^2 b^2}{2}}$$

So Q is $b^2 I$ -subGaussian.

Defn: A zero-mean random matrix Q is subexponential with parameters (V, α)

if

$$\Psi_Q(\lambda) \leq e^{\frac{\lambda^2 V}{2}} \quad \forall |\lambda| \leq \frac{1}{\alpha}.$$

Ex: $Q = \varepsilon g^2 A$ where ε is Rademacher
and $g \sim \mathcal{N}(0,1)$
independent.

[You'll compute V and α in HW]

Lemma: Suppose zero-mean $Q \in \mathbb{S}^d$ satisfies
 $\|Q\|_2 \leq b$. Then we have

$$\Psi_Q(\lambda) \leq \exp\left(\frac{\lambda^2 \text{Var}(Q)}{2(1-b|\lambda|)}\right) \quad \forall |\lambda| \leq \frac{1}{b}$$

pt: $E e^{\lambda Q} = \sum_{k=0}^{\infty} \frac{\lambda^k E Q^k}{k!} =$

$$= I + \frac{\lambda^2 \text{Var}(Q)}{2} + \sum_{j=3}^{\infty} \frac{\lambda^j E Q^j}{j!}$$

$$\leq I + \frac{\lambda^2 \text{Var}(Q)}{2} + \sum_{j=3}^{\infty} \frac{\lambda^2 \lambda^{j-2} E Q^{j-2} b^{j-2}}{j!}$$

$$\leq I + \frac{\lambda^2 \text{Var}(Q)}{2} + \frac{\lambda^2 \text{Var}(Q)}{2} \sum_{j=3}^{\infty} \frac{j!}{\lambda^{j-2} b^{j-2}}$$

$$= I + \frac{\lambda^2 \text{Var}(Q)}{2(1-b|\lambda|)} \leq \exp\left(\frac{\lambda^2 \text{Var}(Q)}{2(1-b|\lambda|)}\right)$$

Matrix Chernoff:

Lemma: Let $Q \in \mathbb{S}^d$ be zero mean with Ψ_Q existing on $(-a, a)$. Then for any $t > 0$, it holds:

$$\mathbb{P}[\chi_1(Q) \geq t] \leq \text{tr}(\Psi_Q(\lambda)) e^{-\lambda t} \quad \forall \lambda \in (0, a)$$

Consequently

$$\mathbb{P}[\|Q\|_2 \geq t] \leq 2 \text{tr}(\Psi_Q(\lambda)) e^{-\lambda t} \quad \forall \lambda \in (0, a)$$

pf: We begin as in scalar case:

$$\begin{aligned} \mathbb{P}[\chi_1(Q) \geq t] &= \mathbb{P}[e^{\chi_1(\lambda Q)} \geq e^{\lambda t}] \\ &= \mathbb{P}[\chi_1(e^{\lambda Q}) \geq e^{\lambda t}] \\ &\leq \mathbb{E}[\chi_1(e^{\lambda Q})] e^{-\lambda t} \\ &\leq \mathbb{E}[\text{tr} e^{\lambda Q}] e^{-\lambda t} \\ &\leq \text{tr}[\Psi_Q(\lambda)] e^{-\lambda t} \quad \square \end{aligned}$$

Lemma: Let $Q_1, \dots, Q_n \in \mathbb{S}^n$ be independent with $\Psi_{Q_i}(\lambda)$ existing for all λ in an interval J . Define $S_n = \sum_{i=1}^n Q_i$.

Then $\text{tr}(\Psi_{S_n}(\lambda)) \leq \text{tr}\left(e^{\sum_{i=1}^n \log \Psi_{Q_i}(\lambda)}\right) \quad \forall \lambda \in J$

[Note: together with the previous lemma, we deduce $\mathbb{P}\left[\left\|\frac{1}{n} \sum_{i=1}^n Q_i\right\|_2 \geq t\right] \leq 2 \text{tr}\left(e^{\sum_{i=1}^n \log \Psi_{Q_i}(\lambda)} \cdot e^{-\lambda t n}\right)$]

pf: Define $G(\lambda) := \text{tr}(\Psi_{S_n}(\lambda))$. Then

$$\begin{aligned} G(\lambda) &= \text{tr}\left(\mathbb{E} e^{\lambda S_{n-1} + \log \exp(\lambda Q_n)}\right) \\ &= \mathbb{E}_{S_{n-1}} \mathbb{E}_{Q_n} \text{tr}\left(e^{\lambda S_{n-1} + \log \exp(\lambda Q_n)}\right) \\ \text{Lieb} &\leq \mathbb{E}_{S_{n-1}} \text{tr}\left(e^{\lambda S_{n-1} + \log \Psi_{Q_n}(\lambda)}\right) \\ &\leq \dots \leq \text{tr}\left(e^{\sum_{i=1}^n \log \Psi_{Q_i}(\lambda)}\right) \quad \square \end{aligned}$$

Thm (Hoeffding)

Let $Q_1, \dots, Q_n \in \mathcal{S}^d$ be independent, zero-mean, and V_i -subGaussian

Then

$$\begin{aligned} \mathbb{P}\left[\left\|\frac{1}{n}\sum_{i=1}^n Q_i\right\|_2 \geq t\right] &\leq 2 \operatorname{rank}\left(\sum_{i=1}^n V_i\right) e^{-\frac{nt^2}{2\sigma^2}} \\ &\leq 2d e^{-\frac{nt^2}{2\sigma^2}} \quad \forall t > 0, \end{aligned}$$

where $\sigma^2 = \left\|\frac{1}{n}\sum_{i=1}^n V_i\right\|_2$.

pf: Define $V = \sum_{i=1}^n V_i$. Suppose first $\operatorname{rank} V = d$. Then

$$\sum_{i=1}^n \log \Psi_{Q_i}(\lambda) \leq \frac{\lambda^2}{2} \sum_{i=1}^n V_i$$

because $\log(\cdot)$ is matrix monotone.

The function $t \mapsto e^t$ is increasing

$$\Rightarrow \operatorname{tr}\left(e^{\sum_{i=1}^n \log \Psi_{Q_i}(\lambda)}\right) \leq \operatorname{tr}\left(e^{\frac{\lambda^2}{2} \sum_{i=1}^n V_i}\right)$$

Chernoff

$$\Rightarrow P\left[\left\|\frac{1}{n}\sum_{i=1}^n Q_i\right\|_2 \geq t\right] \leq 2 \operatorname{tr}\left(e^{\frac{\lambda^2}{2} V}\right) e^{-\lambda t}$$

Note $\operatorname{tr}(e^A) \leq d e^{\|A\|_2}$ for any $A \in \mathbb{R}^d$

$$\Rightarrow \leq 2d e^{\frac{\lambda^2}{2} n \sigma^2 - \lambda n t}$$

Optimize over $\lambda \Rightarrow \lambda = \frac{t}{\sigma^2}$. Done.

If $r := \operatorname{rank}(V) < d$, we can form an eigenvalue decomposition

$$V = U D U^T \text{ where } U \in \mathbb{R}^{d \times r}$$

Then replace Q by

$$\hat{Q} = U^T Q U \in \mathbb{R}^{r \times r}. \quad \square$$

Thm: Let $Q_1, \dots, Q_n \in \mathcal{S}^d$ be zero-mean, independent, with $\|Q_i\|_2 \leq b$. Then

$$\begin{aligned} \mathbb{P}\left[\left\|\frac{1}{n}\sum_{i=1}^n Q_i\right\|_2 \geq t\right] &\leq 2 \operatorname{rank}(V) \exp\left(-\frac{nt^2}{2\sigma^2 + bt}\right) \\ &\leq 2 \operatorname{rank}(V) \exp\left(-\frac{nt^2}{4\sigma^2} \wedge \frac{nt}{2b}\right) \end{aligned}$$

where $V = \sum_{i=1}^n \mathbb{E}Q_i^2$ and $\sigma^2 = \frac{1}{n} \|V\|_2$.
(Hw)

Cor: Under same assumptions,

$$\mathbb{E}\left\|\frac{1}{n}\sum_{i=1}^n Q_i\right\|_2 \leq \sqrt{\frac{\sigma^2 \log(d)}{n}} + \frac{b \log(d)}{n}$$

The Bernstein bound can be extended to rectangular matrices.

Idea: Given a sequence $A_1, \dots, A_n \in \mathbb{R}^{d_1 \times d_2}$
form $Q_i := \begin{bmatrix} 0_{d_1 \times d_1} & A_i \\ A_i^T & 0_{d_2 \times d_2} \end{bmatrix}$

Lemma: $\|Q_i\|_2 = \|A_i\|_2$ and

$$\sigma^2 \leq \max \left\{ \left\| \frac{1}{n} \sum_{i=1}^n \mathbb{E}[A_i A_i^T] \right\|_2, \left\| \frac{1}{n} \sum_{i=1}^n \mathbb{E}[A_i A_i^T] \right\|_2 \right\}$$

and therefore

$$\mathbb{P} \left[\left\| \frac{1}{n} \sum_{i=1}^n A_i \right\|_2 \geq t \right] \leq 2(d_1 + d_2) e^{-\frac{nt^2}{2\sigma^2 + bt}}$$

Covariance Estimation for general distributions

Let $x \in \mathbb{R}^d$ be a random vector with

$\mathbb{E}x = 0$. Define $\Sigma = \mathbb{E}xx^T$ and

$\Sigma_n = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$ where x_i are

i.i.d realizations of x .

[If x_i were Σ -K subGaussian, we proved]

$$\frac{\mathbb{E} \|\Sigma_n - \Sigma\|_2}{\|\Sigma\|_2} \leq K \left(\sqrt{\frac{d}{n}} + \frac{d}{n} \right)$$

Thm: Assume for some $K > 0$, have
 $\|X\|_2 \leq K \sqrt{\text{tr } \Sigma}$ [$= \sqrt{\mathbb{E} \|X\|_2^2}$]

Then

$$\frac{\mathbb{E} \|\Sigma_n - \Sigma\|_2}{\|\Sigma\|_2} \leq C \left(\sqrt{\frac{K^2 + \log(d)}{n}} + \frac{K^2 r \log(d)}{n} \right)$$

pf: Matrix Bernstein where $r = \frac{\text{tr}(\Sigma)}{\|\Sigma\|_2} \leq d$

$$\begin{aligned} \mathbb{E} \|\Sigma_n - \Sigma\|_2 &= \mathbb{E} \left\| \frac{1}{n} \sum_{i=1}^n (x_i x_i^T - \Sigma) \right\|_2 \\ &\leq C \left(\sqrt{\sigma^2 \log(d)/n} + M \log(d)/n \right) \end{aligned}$$

where $\sigma^2 = \|\mathbb{E}(x x^T - \Sigma)^2\|_2$ and $M \geq \|x x^T - \Sigma\|_2$

Let's estimate σ^2 and M .

$$\mathbb{E}(x x^T - \Sigma)^2 = \mathbb{E}(x x^T)^2 - \Sigma^2 \preceq \mathbb{E}(x x^T)^2$$

and

$$(x x^T)^2 = \|x\|_2^2 x x^T \preceq K^2 (\text{tr } \Sigma) x x^T$$

Take expectations $\rightarrow \mathbb{E}(x x^T)^2 \preceq K^2 \text{tr}(\Sigma) \Sigma$

$$\rightarrow \sigma^2 \leq K^2 \text{tr}(\Sigma) \|\Sigma\|_2$$

Finally

$$\begin{aligned}\|xx^T - \Sigma\|_2 &\leq \|x\|_2^2 + \|\Sigma\|_2 \\ &= K^2 \text{tr} \Sigma + \|\Sigma\|_2 \\ &\leq 2K^2 \text{tr} \Sigma =: M\end{aligned}$$

So

$$\mathbb{E} \|\Sigma_n - \Sigma\|_2 \leq C \left(\sqrt{\frac{K^2 \text{tr}(\Sigma) \|\Sigma\|_2 \log(d)}{n}} + \frac{2K^2 \text{tr}(\Sigma) \log(d)}{n} \right)$$

Simplifying completes the proof.

□

Sparse covariance estimation

If we know a priori that Σ is sparse, we should be able to estimate it using fewer observations. Goal: replace

$\|\Sigma_n - \Sigma\|_2 \leq \sqrt{\frac{d}{n}}$ by $\|\hat{\Sigma}_n - \Sigma\|_2 \leq \sqrt{\frac{l \log(d)}{n}}$
where "l" is the sparsity level.

For $\lambda > 0$, define the hard-thresholding operator

$$T_\lambda(u) := \begin{cases} u, & \text{if } |u| > \lambda \\ 0, & \text{o.w.} \end{cases}$$

As a surrogate for sparsity we will use $\|A\|_2$ where

$$A_{ij} = \begin{cases} 1, & \text{if } i=j \text{ or } \sum_{i,j} \neq 0 \\ 0, & \text{o.w.} \end{cases}$$

Lemma: If Σ has at most s nonzero entries per row, then $\|A\|_2 \leq s$.

pt: $A + (s-1)I$ is diagonally dominant and therefore positive semidefinite.

$$\Rightarrow |\lambda_j(A)| \leq s-1.$$

Let u be an eigenvector of A corresponding to $\lambda_j(A)$ and let j index its maximal entry. WLOG $u_j > 0$.
Then $\lambda_j u_j = (Au)_j = \langle A_j, u \rangle \leq s u_j$ \square

Thm: Let $\{x_i\}_{i=1}^n$ be i.i.d zero mean with covariance Σ , and suppose x_{ij} is σ -subGaussian. Then as long as $n > \log(d)$ setting

$$\lambda_n = \sigma^2 \left(8 \sqrt{\frac{\log(d)}{n}} + \delta \right)$$

it holds

$$\mathbb{P} \left[\frac{\|T_{\lambda_n}(\Sigma_n) - \Sigma\|_2}{\sigma^2} \geq 2 \|A\|_2 \left(8 \sqrt{\frac{\log(d)}{n}} + \delta \right) \right] \leq 8 e^{-\frac{n}{16} \min\{\delta, \delta^2\}}$$

We proceed with the following lemmas.

Lemma: Suppose $\|\Sigma_n - \Sigma\|_{\max} \leq \lambda_n$ (coordinate wise)

Then $|T_{\lambda_n}(\Sigma_n) - \Sigma| \leq 2\lambda_n A$ (coordinate wise)

pf: $A^{ij} = 0 \Rightarrow |\Sigma_n^{ij}| \leq \lambda_n \Rightarrow T_{\lambda_n}(\Sigma_n^{ij}) = 0$

$$\begin{aligned} A^{ij} \neq 0 \Rightarrow |T_{\lambda_n}(\Sigma_n^{ij}) - \Sigma^{ij}| &\leq |T_{\lambda_n}(\Sigma_n^{ij}) - \Sigma_n^{ij}| \\ &\quad + |\Sigma_n^{ij} - \Sigma^{ij}| \\ &\leq 2\lambda_n \quad \square \end{aligned}$$

Lemma: For any symmetric matrices $0 \leq A \leq B$,
it holds $\|A\|_2 \leq \|B\|_2$.

pt: Observe

$$\langle a_i, a_j \rangle \leq \langle b_i, b_j \rangle$$

$$\Rightarrow A^2 \leq B^2$$

Since A^2 and B^2 are positive
semidefinite, $\|A^2\|_2 = \sigma_1(A^2)$

$$\|B^2\|_2 = \sigma_1(B^2)$$

The Perron-Frobenius theorem guarantees
that there exists an eigenvector $u \geq 0$
corresponding to $\sigma_1(A^2)$. In this symmetric
setting, this follows directly from:

$$\sigma_1(A^2) = \sup_{\|u\| \leq 1} \sum (A^2)_{ij} u_i u_j \leq \sup_{\|u\| \leq 1} \sum (A^2)_{ij} |u_i| |u_j| \leq \sigma_1(A^2)$$

Therefore

$$\sigma_1(A) \leq \langle A^2 u, u \rangle \leq \sup_{\|v\| \leq 1} \langle B^2 v, v \rangle = \|B^2\|_2$$

$$\Rightarrow \|A\|_2^2 = \|A^2\|_2 \leq \|B^2\|_2 \leq \|B\|_2^2. \quad \square$$

$$\text{So } |\Sigma'_n - \Sigma|_{\max} \leq \lambda_n$$

Lemma: $\Rightarrow \left\| \frac{T(\Sigma_n)}{\lambda_n} - \Sigma \right\|_2 \leq 2\lambda_n \|A\|_2$

$$\mathbb{P} \left[\frac{|\Sigma_n - \Sigma|_{\max}}{\sigma^2} \geq t \right] \leq 8d e^{-\frac{n}{16} t^2}$$

pt: Set $\Delta_n := \Sigma_n - \Sigma$. Observe

$$\Delta_n^{ii} = \frac{1}{n} \sum_{k=1}^n (x_k)_i^2 - \mathbb{E} x_i^2. \quad \text{Apply Bernstein}$$

If $i \neq j$, then

$$2\Delta_n^{ij} = \frac{1}{n} \sum_{k=1}^n (x_k)_i (x_k)_j - 2\Sigma_{ij}$$

$$= \frac{1}{n} \sum_{k=1}^n \left(\begin{aligned} & [(x_k)_i + (x_k)_j]^2 - [\Sigma_{ii} + \Sigma_{jj} + 2\Sigma_{ij}] \\ & - (x_k)_i^2 - (x_k)_j^2 - \Sigma_{ii} - \Sigma_{jj} \end{aligned} \right)$$

Apply Bernstein again + union bound. \square

Notice the bound

$$\|T_{\chi_n}(\Sigma_n) - \Sigma\|_2 \leq C \|A\|_2 \sqrt{\frac{\log(d)}{n}}$$

can be $\sim d \sqrt{\frac{\log d}{n}}$, if Σ is dense, with small entries. This is much worse than our bounds from before.

Operator vs row norm:

For any $A \in \mathbb{S}^d$, it holds

$$\|A\|_2 \geq \max_i \|A_i\|_2$$

Thm: Let $A \in \mathbb{S}^d$ whose entries on and above the diagonal are independent mean zero. Then

$$\mathbb{E} \|A\|_2 \leq C \sqrt{\log(d)} \cdot \mathbb{E} \max_i \|A_i\|_2$$

The proof is based on Bernstein.

pt: Write

$$A = \sum_{i \leq j} Z_{ij} \quad \text{where}$$
$$Z_{ij} = \begin{cases} A_{ij}(e_i e_j^T + e_j e_i^T), & \text{if } i \neq j \\ A_{ii} & \text{if } i = j \end{cases}$$

$$\begin{aligned} \mathbb{E} \|A\|_2 &= \mathbb{E}_{\varepsilon} \left\| \sum_{i \leq j} Z_{ij} - \mathbb{E} \sum_{i \leq j} Z'_{ij} \right\|_2 \\ &\leq \mathbb{E}_{\varepsilon} \left\| \sum_{i \leq j} \varepsilon_{ij} (Z_{ij} - Z'_{ij}) \right\|_2 \\ &\leq 2 \mathbb{E}_{\varepsilon} \left\| \sum_{i \leq j} \varepsilon_{ij} Z_{ij} \right\|_2 \end{aligned}$$

Recall $\varepsilon_{ij} Z_{ij}$ is Z_{ij}^2 -subGaussian.
Therefore with $\sigma^2 = \left\| \sum_{i \leq j} Z_{ij}^2 \right\|_2$ get

$$\begin{aligned} \mathbb{P} \left[\left\| \sum_{i \leq j} \varepsilon_{ij} Z_{ij} \right\|_2 \geq t \right] &\leq 2d \cdot \exp \left(-\frac{t^2}{2\sigma^2} \right) \\ &= \exp \left(\log(2d) - \frac{t^2}{2\sigma^2} \right) \end{aligned}$$

set $t = \sqrt{2\sigma^2 \log(2d)} + u$

$$\Rightarrow \mathbb{P} \left[\left\| \sum_{i \leq j} \varepsilon_{ij} Z_{ij} \right\|_2 \geq \sqrt{2\sigma^2 \log(2d)} + u \right] \leq e^{-\frac{u^2}{2\sigma^2}}$$

Integrating, get

$$\mathbb{E}_\varepsilon \left\| \sum_{i,j} \varepsilon_{ij} z_{ij} \right\| \leq C \sqrt{\log d} \cdot \sqrt{\left\| \sum_{i \leq j} z_{ij}^2 \right\|_2}$$

Next evaluate $\mathbb{E}_z \sqrt{\left\| \sum_{i,j} z_{ij}^2 \right\|_2}$

Quick computation shows

$$z_{ij}^2 = \begin{cases} A_{ij}^2 (e_i e_i^T + e_j e_j^T), & i < j \\ A_{ii}^2 e_i e_i^T, & i = j. \end{cases}$$

↑
diagonal.

$$\Rightarrow \sum_{i \leq j} z_{ij}^2 = \begin{bmatrix} \|A_1\|_2^2 & & & 0 \\ & \|A_2\|_2^2 & & \\ & & \ddots & \\ 0 & & & \|A_d\|_2^2 \end{bmatrix}$$

$$\Rightarrow \sqrt{\left\| \sum_{i \leq j} z_{ij}^2 \right\|_2} = \max_{i=1, \dots, d} \|A_i\|_2 \quad \square$$

Cor (HW) Let $A \in \mathbb{R}^{m \times n}$ whose entries are independent mean-zero. Then

$$\mathbb{E} \|A\|_2 \leq C \sqrt{\log(m+n)} \left(\mathbb{E} \max_i \|A_i\|_2 + \mathbb{E} \max_j \|A^j\|_2 \right)$$

Application: Matrix Completion.

Consider $X \in \mathbb{R}^{d \times d}$ with $\text{rank}(X) = r \ll d$.

Suppose each entry X_{ij} is revealed independently with probability $p \in (0, 1)$

So we get to see $Y \in \mathbb{R}^{d \times d}$ with

$$Y_{ij} = S_{ij} X_{ij} \text{ where } S_{ij} \sim \text{Ber}(p)$$

Thus we get to see

$$m := p \cdot d^2 \text{ entries on average}$$

Thm: Let \hat{x} be a best rank- r approximation to $p^{-1}y$. Then

$$\mathbb{E} \frac{1}{d} \|\hat{X} - X\|_F \leq C \sqrt{\frac{rd \log(d)}{m}} \|X\|_{\max}$$

as long as $m \geq \log(d)$.

So in order to ensure

$$\mathbb{E} \sqrt{\frac{1}{d^2} \sum_{i,j} (\hat{X}_{ij} - x_{ij})^2} \leq \epsilon \|X\|_{\max},$$

it suffices to see $m = \frac{1}{\epsilon^2} rd \log(d)$ entries on average.

pf: First bound error is operator norm

$$\begin{aligned} \|\hat{X} - X\|_2 &\leq \|\hat{X} - p^{-1}y\|_2 + \|p^{-1}y - X\|_2 \\ &\leq 2\|p^{-1}y - X\|_2 = \frac{2}{p} \underbrace{\|y - pX\|_2}_{\text{easy to understand}} \end{aligned}$$

$$(Y - pX)_{ij} = (S_{ij} - p) X_{ij}$$

independent mean-zero

$$\Rightarrow \mathbb{E} \|Y - pX\|_2 \leq C \sqrt{\log(n)} \cdot (\mathbb{E} \max_i \| (Y - pX)_i \|_2 + \mathbb{E} \max_j \| (Y - pX)^j \|_2)$$

Observe

$$\begin{aligned} \| (Y - pX)_i \|_2^2 &= \sum_j (Y - pX)_{ij}^2 = \sum_j (S_{ij} - p)^2 X_{ij}^2 \\ &\leq \left(\sum_j (S_{ij} - p)^2 \right) \|X\|_{\max}^2 \end{aligned}$$

Can bound $\mathbb{E} \max_i \sum_j (S_{ij} - p)^2 \leq Cp d$
(HW)

using concentration and union bound.

$$\Rightarrow \mathbb{E} \|Y - pX\|_2 \leq C \sqrt{pd \log(d)} \|X\|_{\max}$$

$$\Rightarrow \mathbb{E} \|\hat{X} - X\|_2 \leq C \sqrt{\frac{d \log(d)}{p}} \|X\|_{\max}$$

Step 2: Pass to Frobenius.

$$\|\hat{X} - X\|_F \leq \sqrt{2r} \|\hat{X} - X\|_2$$

rank $2r$

$$\Rightarrow \mathbb{E} \frac{1}{d} \|\hat{X} - X\|_F \leq C \sqrt{\frac{rd \log(d)}{pd^2}} \|X\|_{\max}.$$

\square