

Course: High Dimensional Probability:  
applications to statistical learning  
and data science.

Text(s): "High-Dimensional Probability"  
Roman Vershynin  
"High-Dimensional Statistics"  
Martin J. Wainwright

Geography: Let  $n$  be sample size  
 $d$  is the dimension of data

- Classical Asymptotics

$n \rightarrow \infty$ ,  $d$  fixed

Ex: Laws of large numbers  
Central limit theorem

- High dimensional Asymptotics

$(n, d) \rightarrow \infty$  while  $\Psi(n, d) \rightarrow d < \infty$ ,  
where  $\Psi$  is a "scaling function"

E.g.:  $\Psi(n, d) = \frac{d}{n}$

- Non-asymptotic bounds  $\leftarrow$  This class  
 $(n, d)$  is fixed. Probability of successful  
estimation/inference depends on  $(n, d)$

See ( $W:1$ ) for examples of the three viewpoints.

## Lecture 1 : Review ( $V:1$ )

Let  $X$  be a random variable on probability space

Expectation and Variance

$$EX, \text{Var}(X) = E(X - EX)^2 = E[X^2] - (EX)^2$$

Moment Generating Function :

$$M_X(t) = Ee^{tX}, t \in \mathbb{R}$$

$L^p$ -norm  $\|X\|_p = (\mathbb{E}|X|^p)^{\frac{1}{p}}, p \in (0, \infty)$

Banach Space :

$$L^p = \{X : \|X\|_p < \infty\}$$

Remark:  $L^2$  is a Hilbert space

$$\langle X, Y \rangle_2 = E[XY], \|X\|_2 = \sqrt{\mathbb{E}[X^2]}$$

Then

$$\|X - EX\|_2 = \sqrt{\text{Var}(X)} \quad \text{and}$$

$$\begin{aligned} \text{cov}(X, Y) &:= E(X - EX)(Y - EY) \\ &= \langle X - EX, Y - EY \rangle_2 \end{aligned}$$

# Limit Theorems

Thm: (Strong Law of Large Numbers)

Let  $X_1, X_2, \dots$  be i.i.d. random variables with  $\mathbb{E}X_i = \mu$ . Then  $S_n = \sum_{i=1}^n X_i$  satisfies

$\frac{1}{n} \cdot S_n \rightarrow \mu$  almost surely

$$\left[ \lim_{n \rightarrow \infty} \frac{S_n}{n} = \mu \text{ w.p. 1} \right]$$

Intuition:  $\text{Var}(S_n) = \frac{\text{Var}(X_i)}{n}$

Thm: (Central Limit Theorem)

Let  $X_1, X_2, \dots$  be i.i.d. with

$$\mathbb{E}X_i = \mu, \quad \text{Var } X_i = \sigma^2$$

Define  $Z_n = \frac{S_n - \mathbb{E}S_n}{\sqrt{\text{Var}(S_n)}} = \frac{1}{\sqrt{n}} \sum_{i=1}^n (X_i - \mu)$

Then  $Z_n \rightarrow N(0, 1)$  in distribution

$$\left[ \text{P}[Z_n \geq t] \rightarrow \frac{1}{\sqrt{2\pi}} \int_t^{\infty} e^{-x^2/2} dx \right]$$

# Chapter 1

## Concentration Inequalities

- Chernoff Bound
- Sub-Gaussian RV (Hoeffding)
- Sub-Exponential RV (Bernstein)
- Johnson-Lindenstrauss
- McDiarmid inequality
- Robust Mean Estimation

# Concentration Inequalities

Goal: Bounds of the form

$$\Pr\left[\frac{S_n}{n} \geq \mu + t\right] \leq \text{small}(n, t)$$

Prop: Let  $Z \sim N(\mu, \sigma^2)$ . Then

$$\Pr[Z \geq \mu + t] \leq e^{-t^2/2\sigma^2} \quad t > 0$$

So would hope  $\Pr\left[\frac{S_n}{n} \geq \mu + t\right] \leq C e^{-\frac{nt^2}{2\sigma^2}}$

Bad news from CLT:

$$\sup_t \left\{ |\Pr[S_n \geq t] - \Pr[Z \geq t]| \right\} \begin{cases} \text{can be } \sqrt{\frac{1}{n}} \\ \text{even for Bernoulli} \end{cases}$$

where  $Z \sim N(0, 1)$

[See (V: 2.1) for example]

Chernoff Method: {Sub-Gaussian, Sub-Exponential} RV's

Lemma (Markov) For any non-negative  $X$  and  $t \geq 0$ , we have

$$\Pr[X \geq t] \leq \frac{\mathbb{E}X}{t}$$

$$\text{pf: } \mathbb{E}X = \mathbb{E} X \mathbf{1}_{\{X \geq t\}} + \mathbb{E} X \mathbf{1}_{\{X < t\}}$$

$$\geq t \mathbb{E} \mathbf{1}_{\{X \geq t\}} + 0 = t \Pr[X \geq t]$$

Cor: (Chebychev) For any random variable  $X$  with  $\text{var}(X) \leq \sigma^2 < \infty$ , we have

$$\Pr[|X - \mathbb{E}[X]| \geq t] \leq \frac{\sigma^2}{t^2} \quad t \geq 0$$

More generally, suppose  $\mu = \mathbb{E}X < \infty$ . Then for all  $\lambda \geq 0$ , we have

$$\Pr[X - \mu \geq t] = \Pr[e^{\lambda(X-\mu)} \geq e^{\lambda t}] \stackrel{\text{Markov}}{\leq} \frac{\mathbb{E} e^{\lambda(X-\mu)}}{e^{\lambda t}}$$

$$\Rightarrow \log \Pr[X - \mu \geq t] \leq \inf_{\lambda \geq 0} \left\{ \log \mathbb{E} e^{\lambda(X-\mu)} - \lambda t \right\} \\ = - \sup_{\lambda \geq 0} \left\{ \lambda t - \log \mathbb{E} e^{\lambda(X-\mu)} \right\}$$

Define for any function  $\varphi: \mathbb{R} \rightarrow \mathbb{R} \cup \{\pm\infty\}$ , the Fenchel conjugate

$$\varphi^*(t) = \sup_{\lambda} \{ t\lambda - \varphi(\lambda) \}$$

Let's look at the main example

$$\varphi(\lambda) = \log \mathbb{E} e^{\lambda(X-\mu)}$$

For all  $\lambda \in \mathbb{R}$ , observe

$$\varphi(\lambda) = \log \mathbb{E} e^{\lambda(X-\mu)} \geq \log e^{\lambda(X-\mu)} = 0$$

So when  $\lambda < 0$  and  $t > 0$ , we have

$$t\lambda - \varphi(\lambda) \leq 0 = 0 \cdot t - \varphi(0).$$

Therefore for  $t \geq 0$ , equality holds:

$$\varphi_x^*(t) = \sup_{\lambda \geq 0} \{ t\lambda - \varphi(\lambda) \}$$

We therefore arrive at the bound:

Chernoff Bound:

$$\Pr[X \geq t] \leq \exp(-\psi_x^*(t))$$

where  $\psi_x(\lambda) := \log(E e^{\lambda(X-\mu)})$ .

Ex: Let  $X \sim N(\mu, \sigma^2)$ . Then

$$E e^{\lambda(X-\mu)} = e^{\frac{\sigma^2}{2}\lambda^2} \quad \forall \lambda \in \mathbb{R}.$$

So  $\Pr[X \geq \mu + t] \leq e^{-\frac{t^2}{2\sigma^2}} \quad t > 0$

Defn:  $X$  with mean  $\mu = E X$ , is

sub-Gaussian with parameter  $\sigma > 0$

$$E e^{\lambda(X-\mu)} \leq e^{\frac{\sigma^2}{2}\lambda^2} \quad \forall \lambda \in \mathbb{R}$$

If  $X$  is sub-Gaussian, so is  $-X$ .

$$\Rightarrow \Pr[|X-\mu| \geq t] \leq 2e^{-\frac{t^2}{2\sigma^2}}$$

Ex: (Rademacher)

Let  $\epsilon$  be a rademacher RV:

$$\mathbb{P}[\epsilon = 1] = \mathbb{P}[\epsilon = -1] = \frac{1}{2}.$$

$$\mathbb{E}[e^{\lambda\epsilon}] = \frac{1}{2}(e^{-\lambda} + e^{\lambda}) = \cosh(\lambda)$$

Exercise:  $\cosh(x) \leq \exp\left(\frac{x^2}{2}\right) \forall x \in \mathbb{R}$

So  $\epsilon$  is sub-Gaussian with  $\sigma=1$ .

Ex: (Bounded RV)

Suppose  $X$  is supported on  $[a, b]$ .

Then Jensen's inequality implies

$$\mathbb{E} e^{\lambda(X - \mathbb{E}_x X')} \leq \mathbb{E}_{x'} e^{\lambda(X - x')}$$

where  $X'$  is independent copy of  $X$ .

$$\text{Observe } X - X' \sim \mathcal{E}(X - \bar{X})$$

$$\text{So } \mathbb{E}_{x, x'} e^{\lambda(X - x')} = \mathbb{E}_{x, x'} \mathbb{E}_{\mathcal{E}} e^{\lambda \cdot \mathcal{E}(X - x')}$$

$$\leq \mathbb{E}_{x, x'} \exp\left(\lambda^2 (X - \bar{X})^2 / 2\right)$$

$$\leq \exp\left(\lambda^2 (b - a)^2 / 2\right)$$

A slightly more involved argument shows that  $X$  is sub-Gaussian with parameter  $\frac{b-a}{2}$

$$\Rightarrow \mathbb{P}[|X-\mu| \geq t] \leq 2 \exp\left(-\frac{2t^2}{(b-a)^2}\right)$$

Lemma: [Sum Rule]

$X_i$  are independent  $\sigma_i$ -sub-Gaussian  $\forall i=1, \dots, n \Rightarrow \sum_{i=1}^n X_i$  is  $\|\sigma\|_2$ -sub Gaussian

Cor: (Hoeffding)

Suppose  $X_1, \dots, X_n$  are independent with  $\mathbb{E} X_i = \mu_i$  and are  $\sigma_i$ -sub-Gaussian

Then

$$\mathbb{P}\left[\sum_{i=1}^n (X_i - \mu_i) \geq t\right] \leq \exp\left\{-\frac{t^2}{2\|\sigma\|_2^2}\right\}$$

$\Rightarrow$  If  $\mu_i = \mu$ ,  $\sigma_i = \sigma$ , then

$$\mathbb{P}\left[\frac{1}{n} \sum_{i=1}^n X_i \geq \mu + t\right] \leq \exp\left\{-\frac{nt^2}{2\sigma^2}\right\}$$

Subexponential RV:

Ex: Let  $Z \sim N(0, 1)$ . Let's compute

$$\mathbb{E}[e^{\lambda(Z^2 - 1)}] = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{\lambda(x^2 - 1)} e^{-x^2/2} dx$$
$$= \frac{e^{-\lambda}}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-(1-2\lambda)x^2/2} dx$$
$$= \begin{cases} \frac{e^{-\lambda}}{\sqrt{1-2\lambda}}, & \text{if } \lambda \leq \frac{1}{2} \\ +\infty, & \text{if } \lambda > \frac{1}{2} \end{cases}$$

Defn:  $X$  with mean  $\mu = \mathbb{E}X$  is

subexponential with parameters  $(\sigma^2, \mu)$

if  $\mathbb{E}[e^{\lambda(X-\mu)}] \leq e^{\frac{\sigma^2\lambda^2}{2}} \quad \forall |\lambda| \leq \frac{1}{\sigma}$

Back to example  $Z \sim N(0, 1)$

$$\mathbb{E}[e^{\lambda(Z^2 - 1)}] \leq \frac{e^{-\lambda}}{\sqrt{1-2\lambda}} \leq e^{\frac{4\lambda^2}{2}} \quad \forall |\lambda| \leq \frac{1}{4}$$

So  $Z^2$  is  $(2, 4)$ -subexponential.

Thm (Subexponential tail bound)  
 Let  $X$  be subexponential with  $(\sigma, \alpha)$ .

Then

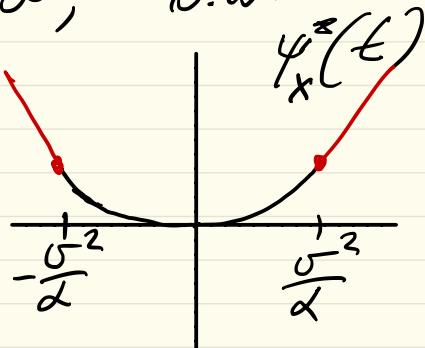
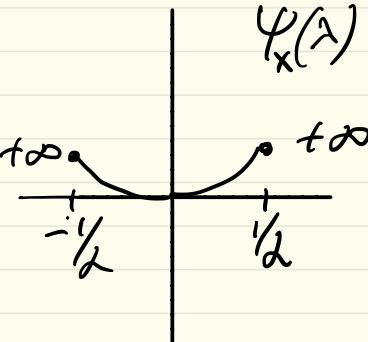
$$\mathbb{P}[X - \mu \geq t] \leq \begin{cases} e^{-\frac{t^2}{2\sigma^2}}, & \text{if } |t| \leq \frac{\sigma^2}{2} \\ e^{-\frac{t}{2\alpha}}, & \text{o.w.} \end{cases}$$

pf: Back to Chernoff

$$\log \mathbb{P}[X - \mu \geq t] \leq -\psi_x^*(t)$$

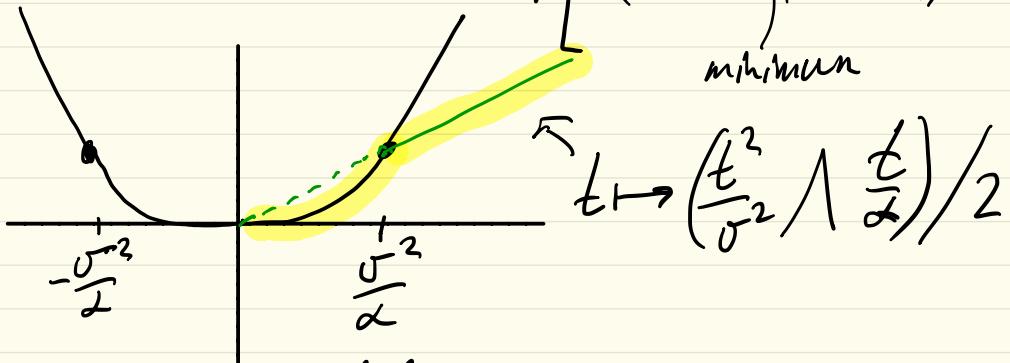
$$\text{where } \psi_x(\lambda) = \log \mathbb{E} e^{\lambda(X-\mu)}$$

$$= \begin{cases} \frac{\lambda^2}{2}, & \text{if } |\lambda| \leq \frac{1}{\alpha} \\ +\infty, & \text{o.w.} \end{cases}$$



Thm: (Bernstein) Let  $X$  be subexponential with parameter  $(\sigma, \alpha)$  and mean  $\mu = \mathbb{E}X$ . Then

$$\mathbb{P}[|X - \mu| \geq t] \leq 2 \exp \left[ - \left( \frac{t^2}{\sigma^2} \wedge \frac{t}{\alpha} \right) / 2 \right]$$



Lemma: [Sum Rule]

$X_i$  are  $(\sigma_i, \alpha_i)$ -subexp  $\forall i=1, \dots, n \Rightarrow \sum_{i=1}^n X_i$  is  $(\|\sigma\|_2, \|\alpha\|_\infty)$ -subexp

Thm (Bernstein for sums) Let  $X_1, \dots, X_n$  be independent sub-exponential with parameters  $(\sigma_i, \alpha_i)$ , and with mean  $\mu_i = \mathbb{E}X_i$ .

$$\mathbb{P}\left[\left|\sum_{i=1}^n (X_i - \mu_i)\right| \geq t\right] \leq 2 \exp \left[ - \frac{1}{2} \left( \frac{t^2}{\|\sigma\|_2^2} \wedge \frac{t}{\|\alpha\|_\infty} \right) \right]$$

Thm: (Improved Bernstein for bounded RV)

Suppose  $|X-\mu| \leq b$ ,  $\mathbb{E}(X-\mu)^2 = \sigma^2$ . Then

$$\mathbb{E} e^{\lambda(X-\mu)} \leq \exp\left(\frac{\lambda^2 \sigma^2}{2(1-b/\lambda)}\right) \quad \text{if } |\lambda| < \frac{1}{b}$$

Therefore

$$\textcircled{A} \quad \mathbb{P}[|X-\mu| \geq t] \leq 2e^{-\frac{t^2}{2(\sigma^2 + bt)}}$$

pf: Taylor expansion:

$$\begin{aligned} \mathbb{E} e^{\lambda(X-\mu)} &= \sum_{k=0}^{\infty} \lambda^k \frac{\mathbb{E}(X-\mu)^k}{k!} = 1 + \lambda \frac{\sigma^2}{2} + \sum_{k=3}^{\infty} \lambda^k \frac{\mathbb{E}(X-\mu)^k}{k!} \\ &\leq 1 + \sum_{k=2}^{\infty} \frac{\lambda^2 \sigma^2}{2 \cdot 3 \cdot 4 \cdots k} \frac{b^k \lambda^k}{k!} \\ &\leq 1 + \frac{\lambda^2 \sigma^2}{2} \cdot \frac{1}{1 - b/\lambda} \leq \exp\left(\frac{\lambda^2 \sigma^2}{2} \cdot \frac{1}{1 - b/\lambda}\right) \end{aligned}$$

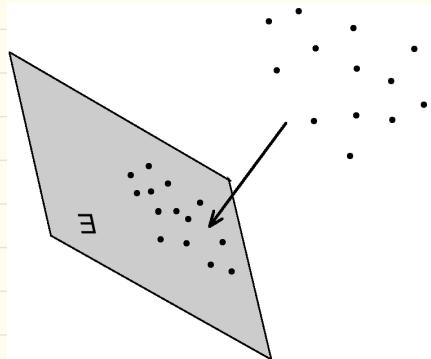
A Follows from Chernoff by

$$\text{setting } \lambda = \frac{t}{bt + \sigma^2}, \mathbb{E}\left[\frac{1}{b}\right]$$

This is superior to Hoeffding when  $\sigma^2 \ll b$ .

## Application: Dimensionality Reduction

Given  $u_1, \dots, u_m \in \mathbb{R}^d$  with  $m \ll d$ ,  
 can one map  $u_1, \dots, u_m$  into a  
 lower dimensional space with low distortion?



Thm: (Johnson-Lindenstrauss)

Fix  $\epsilon, \delta \in (0, 1)$ , a set  $\mathcal{U} \subseteq \mathbb{R}^d$  of  $m$  points and a number  $n > \frac{16 \ln(\frac{m^2}{\delta})}{\epsilon^2}$ .

Let  $X \in \mathbb{R}^{n \times d}$  consist of i.i.d  $N(0, 1)$  entries.

Then with probability  $1 - \delta$ , the map  $f(u) = \frac{1}{\sqrt{n}} X u$  satisfies

$$1 - \epsilon \leq \frac{\|f(u) - f(v)\|_2^2}{\|u - v\|_2^2} \leq 1 + \epsilon \quad \text{Hence.}$$

pf: Observe  $i^{\text{th}} \text{ row of } X$

$$\frac{\|X_{\text{all}}\|_2^2}{\|u\|_2^2} = \sum_{i=1}^n \underbrace{\left\langle \frac{x_i}{\|u\|_2}, \frac{u}{\|u\|_2} \right\rangle}_{\substack{i.i.d \\ N(0, 1)}}^2$$

$\Rightarrow \frac{\|X_{\text{all}}\|_2^2}{\|u\|_2^2}$  is  $(2\sqrt{n}, 4)$  - subexponential

$\Rightarrow$  Bernstein:

$$P\left[\left|\frac{\|X_{\text{all}}\|_2^2}{\|u\|_2^2} - 1\right| > \varepsilon\right] \leq 2 \exp\left[-\left(\frac{n\varepsilon^2}{8}\right) \wedge \left(\frac{n\varepsilon}{8}\right)\right]$$

So for any  $i, j$  get  $= 2 \exp\left(-\frac{n\varepsilon^2}{8}\right) \forall 0 \leq \varepsilon \leq 1$

$$P\left[\frac{\|f(u_i - u_j)\|_2^2}{\|u_i - u_j\|_2^2} \notin [1 - \varepsilon, 1 + \varepsilon]\right] \leq 2e^{-n\varepsilon^2/8}$$

Take union bound over  $\binom{m}{2}$  pairs of points  
 $2\binom{m}{2}e^{-n\varepsilon^2/8} \leq m^2 e^{-n\varepsilon^2/8} = S$   $\blacksquare$

Question: What if  $m = \infty$  but  $U$  has few "degrees of freedom"?

So far, we have focused on concentration of the average  $\frac{1}{n} \sum_{i=1}^n x_i$ .

Often one is interested in bounds

$$P[f(x_1, \dots, x_n) - E f(x_1, \dots, x_n) | > t] \leq \text{small}(n, t)$$

where  $x_i$  are independent and  
 $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is some function.

Useful insight:

As long as  $f(x_1, \dots, x_n)$  depends weakly on individual  $x_i$ , concentration holds

Thm: (McDiarmid) Suppose that  $f: \mathcal{X}^n \rightarrow \mathbb{R}$  has the bounded difference property:

$\exists L_1, \dots, L_n$  such that

$$|f(x_1, \dots, x_k, \dots, x_n) - f(x_1, \dots, x'_k, \dots, x_n)| \leq L_k \quad \forall x, x' \in \mathcal{X}^n$$

Then for independent R's  $X = (x_1, \dots, x_n)$  have

$$\mathbb{P}[|f(X) - \mathbb{E}f(X)| > t] \leq 2e^{-\frac{2t}{LK\|L\|_2^2}}$$

pf: We will use the Martingale method.

Define

$$y_0 = \mathbb{E}f(X) \text{ and } y_i = \mathbb{E}[f(X) | x_1, \dots, x_i] \quad \forall i$$

Observe

$$y_i = y_0 + \sum_{j=0}^{i-1} (y_{j+1} - y_j) =: D_{j+1} = y_0 + \sum_{j=1}^i D_j$$

and

$$\mathbb{E}[y_i | x_1, \dots, x_{i-1}] = \mathbb{E}[f(X) | x_1, \dots, x_{i-1}] = y_{i-1}$$

$$\Rightarrow \mathbb{E}[D_{j+1} | x_1, \dots, x_i] = 0$$

$$\Rightarrow \mathbb{E}[e^{\lambda(f(X) - \mathbb{E}f(X))}] = \mathbb{E}[e^{\lambda(y_n - y_0)}] =$$

$$\begin{aligned}
 \mathbb{E}[e^{\lambda(f(x) - \mathbb{E}f(x))}] &= \mathbb{E}[e^{\lambda(Y_n - \bar{Y})}] \\
 &= \mathbb{E}\left[e^{\lambda \sum_{j=1}^n D_j}\right] \\
 &= \mathbb{E}\left[e^{\lambda(Y_{n-1} - \bar{Y})} e^{\lambda D_n}\right] \\
 &= \mathbb{E}\left[e^{\lambda(Y_{n-1} - \bar{Y}_0)} \underbrace{\mathbb{E}\left[e^{\lambda D_n} \mid X_1, \dots, X_{n-1}\right]}_{\mathbb{E}\left[e^{\lambda D_i} \mid X_1, \dots, X_{i-1}\right]}\right]
 \end{aligned}$$

Let  $x'$  differ from  $x$  in  $x_i$ . Then

$$\mathbb{E}\left[e^{\lambda D_i} \mid X_1, \dots, X_{i-1}\right] = \mathbb{E}\left[e^{\lambda(Y_i - \bar{Y}_{i-1})} \mid X_1, \dots, X_{i-1}\right]$$

$$= \mathbb{E}\left[e^{\lambda \mathbb{E}[f(x) - f(x')] \mid X_1, \dots, X_i} \mid X_1, \dots, X_{i-1}\right]$$

$$\stackrel{\text{Jensen}}{\leq} \mathbb{E}\left[e^{\lambda(f(x) - f(x'))} \mid X_1, \dots, X_{i-1}\right]$$

↑ bounded by  $L_i$

$$\Rightarrow \mathbb{E}\left[e^{\lambda(f(x) - f(x'))} \mid X_1, \dots, X_{i-1}\right] \leq e^{\frac{\lambda^2 L_i^2}{8}}$$

$$\text{So } \mathbb{E}[e^{\lambda(f(x) - \mathbb{E}f(x))}] \leq e^{\frac{\lambda^2 \|L\|^2}{8}}. \text{ Apply Chernoff}$$

Recall if  $X_1, \dots, X_n$  are independent  $\zeta$ -subGaussian with  $\mathbb{E}X_i = \mu$ , the Hoeffding implies that  $\hat{x} = \frac{1}{n} \sum_{i=1}^n X_i$  satisfies

$$\Pr[|\hat{x} - \mu| \leq t] \geq 1 - 2\exp\left(-\frac{nt^2}{2\zeta^2}\right)$$

or equivalently

$$\Pr\left[|\hat{x} - \mu| \leq \sqrt{\frac{2\zeta^2 \ln\left(\frac{2}{p}\right)}{n}}\right] \geq 1 - p.$$

Can one achieve a similar guarantee without subGaussian assumption with a different estimator  $\hat{x}$ ?

Answer: yes, almost!

Thm: (Median of means)

Consider  $X \in \mathbb{R}$  with  $\mathbb{E}X = \mu$  and  $\text{Var}(X) = \sigma^2$ .  
Let  $X_1, \dots, X_n$  be i.i.d realizations of  $X$ .  
Subdivide into  $k = 18 \log(\frac{1}{\rho})$  bins and form the empirical means  $\hat{x}_j$  for  $j = 1, \dots, k$ .

Then  $\hat{x} = \text{median}(\hat{x}_1, \dots, \hat{x}_k)$  satisfies

$$P\left[|\hat{x} - \mu| \leq \frac{\sqrt{27\sigma^2 \log(\frac{1}{\rho})}}{n}\right] \geq 1 - \rho$$

pf: By Chebychev

$$P\left[|\hat{x}_i - \mu| \geq \frac{\sqrt{3\sigma^2 k}}{2n}\right] \leq \frac{\frac{\sigma^2}{n/k}}{\frac{3\sigma^2 k}{2n}} = \frac{2}{3} \quad \forall i.$$

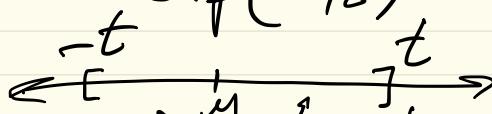
Let  $1_{i,j}$  be indicator of this event

Then by Hoeffding,

$$P\left[\sum_{i=1}^k 1_{i,j} > \frac{t}{2}\right] \geq 1 - \exp\left(-\frac{k}{18}\right)$$

In this event,

$$\Rightarrow |\hat{x} - \mu| \leq \sqrt{\frac{3\sigma^2 k}{2n}}$$



Notice that in contrast to sub-Gaussian case,  $\hat{x}$  depends on confidence level  $\rho$ .

# Chapter 2

Random Vectors in High Dimensions

- Concentration of the norm
- Isotropy
- Similarity of Normal and Spherical
- Sub-Gaussian and Sub exponential random vectors.

---

Two main results we'll prove:

- 1) Sub-Gaussian vectors concentrate around a sphere.
- 2) Two independent isotropic sub-Gaussian random vectors are nearly orthogonal in high dimensions.

We next investigate the behavior of random vectors in high dimensions.

### Concentration of the norm

Let  $X = (X_1, \dots, X_d) \in \mathbb{R}^d$  have independent  $\sigma$ -subGaussian coordinates with  $\mathbb{E} X_i = 0$  and  $\mathbb{E} X_i^2 = 1$ .

What can we expect from  $\|X\|_2^2$  and  $\|X\|_2$ ?

Lemma: Suppose  $y$  is  $\sigma$ -subGaussian. Then  $y^2$  is  $(\sigma^2, 4\sigma^2)$  subexponential.

pf sketch: Step 1: Estimate  $\mathbb{E}[|y|^r] \leq r^{1/2} \sigma \Gamma(r/2)$

$$\text{using } \mathbb{E}[|y|^r] = \int_{-\infty}^{\infty} \mathbb{P}[|y| > t]^{1/r} dt$$

Step 2: Use Taylor Expansion

$$\mathbb{E}[e^{\lambda(y^2 - \mathbb{E}y^2)}] \leq 1 + \sum_{r=2}^{\infty} \lambda^r 2^{r+1} \frac{\sigma^{2r}}{6^r} \leq 1 + \frac{8\lambda^2 \sigma^4}{1 - 2\lambda\sigma^2} \leq \exp(\dots) \quad \square$$

Cor: Let  $X = (X_1, \dots, X_d) \in \mathbb{R}^d$  have independent 6-subGaussian coordinates with

$$\mathbb{E} X_i = 0 \text{ and } \mathbb{E} X_i^2 = 1$$

$$\text{Then } \mathbb{P}\left[\|X\|_2^2 - d \geq t d\right] \leq 2 \exp\left(-\frac{d}{46^2}(t/t^2)\right)$$

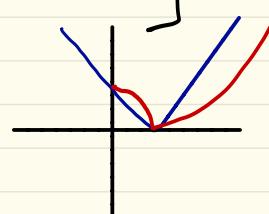
$$\mathbb{P}\left[\|X\|_2 - \sqrt{d} \geq t\sqrt{d}\right] \leq 2 \exp\left(-\frac{dt^2}{46^2}\right)$$

$$\text{pf: } \|X\|_2^2 = \sum_{i=1}^d X_i^2 \Rightarrow \begin{cases} \mathbb{E} \|X\|_2^2 = d \\ \|X\|_2^2 \text{ is } (6\sqrt{d}, 46^2) \text{ subexponential.} \end{cases}$$

$$\text{Bernstein} \Rightarrow \mathbb{P}\left[\left|\frac{1}{d}\|X\|_2^2 - 1\right| \geq t\right] \leq 2 \exp\left[-\frac{d}{46^2}(t/t^2)\right]$$

Observation: for any  $z \geq 0$ , have

$$|z-1| \geq t \Rightarrow |z^2-1| \geq t \sqrt{t^2}$$



$$\text{So } \mathbb{P}\left[\left|\frac{1}{\sqrt{d}}\|X\|_2 - 1\right| \geq t\right] \leq \mathbb{P}\left[\left|\frac{1}{d}\|X\|_2^2 - 1\right| \geq t \sqrt{t^2}\right] \\ \leq 2 \exp\left(-\frac{d}{46^2}t^2\right) \quad \square$$

## Isotropic Vectors

Recall for  $X \in \mathbb{R}^d$ , covariance

$$\text{cov}(X) = E((X - \mu)(X - \mu)^T)$$

where  $\mu = EX$ .

Defn: A random vector  $X \in \mathbb{R}^d$  with  $EX=0$  is isotropic if

$$\Sigma(X) := EXX^T = I_d$$

Remark: If  $\Sigma = \Sigma(X)$  is invertible, then

$$Z := \Sigma^{-\frac{1}{2}}(X - \mu) \text{ is isotropic.}$$

Lemma:  $X$  is isotropic iff

$$E\langle X, y \rangle^2 = \|y\|_2^2 \quad \forall y \in \mathbb{R}^d$$

pf:  $X$  is isotropic iff

$$EXX^T = I$$

$$\text{iff } y^T EXX^T y = y^T y$$

$$\text{iff } E y^T X X^T y = \|y\|_2^2$$

$$\text{iff } E\langle X, y \rangle^2 = \|y\|_2^2 \quad \square$$

Thus if  $\mathbb{E}X=0$ , then  $X$  is isotropic iff marginal  $\left\langle X, \frac{y}{\|y\|} \right\rangle$  has unit variance  $\forall y \in \mathbb{R}^d$

Lemma: Let  $X \in \mathbb{R}^d$  be isotropic. Then

$$\mathbb{E} \|X\|_2^2 = d.$$

Moreover, if  $X$  and  $y$  are two independent isotropic vectors, then

$$\mathbb{E} \langle X, y \rangle^2 = d$$

Pf: First

$$\|X\|_2^2 = X^T X = \text{trace}(XX^T)$$

$$\Rightarrow \mathbb{E} \|X\|_2^2 = \text{trace}(I_d) = d.$$

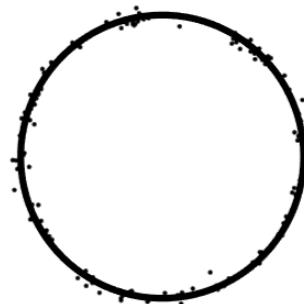
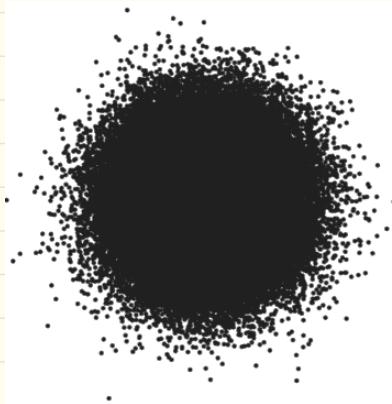
$$\begin{aligned} \text{Next } \mathbb{E} \langle X, y \rangle^2 &= \mathbb{E}_y \left[ \mathbb{E}_x \langle X, y \rangle^2 | y \right] \\ &= \mathbb{E}_y \|y\|_2^2 = d \end{aligned}$$

□

Let  $X$  and  $Y$  be independent and isotropic

Then  $\|X\| \sim \sqrt{d}$  and  $\left\langle \frac{X}{\|X\|}, \frac{Y}{\|Y\|} \right\rangle \sim \frac{1}{d}$ .  
 $\|Y\| \sim \sqrt{d}$   
 $\therefore$  Almost orthogonal.

Can be made rigorous by assuming  
light tails.



# Examples of isotropic RV:

- 1) Spherical  $X \sim \text{Unif}(\mathbb{S}^{d-1})$  HW
- 2) Symmetric Bernoulli:  $X \sim \text{Unif}(\{-1, 1\}^d)$
- 3) Any vector  $X = (X_1, \dots, X_d)$ , where  $X_i$  are independent, zero mean, unit variance.
- 4) Coordinate  $\text{Unif}(\{\sqrt{d}e_i\}_{i=1}^d)$
- 5) Gaussian  $g = (g_1, \dots, g_d) \sim N(0, I_d)$   
 Recall this means  $g_i$  are i.i.d.  $N(0, 1)$   
 $\Rightarrow$  Density  $\frac{d}{\prod_{i=1}^d p_i(x)} = \prod_{i=1}^d \frac{1}{\sqrt{2\pi}} e^{-x_i^2/2} = \frac{1}{(2\pi)^{d/2}} e^{-\|x\|^2/2}$ .  
 $\Rightarrow N(0, I_d)$  is rotation invariant.

## Exercise:

Let  $g \sim N(0, I_d)$ . Then  
 $r := \|g\|_2$  and  $\theta = \frac{g}{\|g\|_2}$   
 are independent random variables and  
 $\theta \sim \text{Unif}(\mathbb{S}^{d-1})$

Defns:  $X \in \mathbb{R}^d$  is  $\sigma$ -subGaussian if  
 $\langle X, u \rangle$  is  $\sigma$ -subGaussian  $\forall u \in \mathbb{S}^{d-1}$

Ex: Let  $X = (X_1, \dots, X_d)$  be RV  
with independent  $\sigma$ -subGaussian  $X_i$ .  
Then  $X$  is  $\sigma$ -subGaussian.

Ex: 1)  $N(0, I_d)$  is 1-subGaussian.  
2)  $\text{Unif}([-1, 1]^d)$  is 1-subGaussian.  
3)  $\text{Unif}(\{\sqrt{d}e_i\}_{i=1}^d)$  is  $\sigma$ -subGaussian  
with  $\sigma \asymp \sqrt{d \log(d)}$

Way too big to be useful

4)  $\text{Unif}(\sqrt{d}[-1, 1]^{d-1})$  is  $c$ -subGaussian

Q: How to get high probability bound on  
 $\|X - \mathbb{E}X\|$ ?

Idea:

$$\sup_{y \in \mathbb{R}^d} \langle X - \mathbb{E}X, y \rangle$$

# Chapter 3

## Uniform Laws:

- Glivenko-Cantelli!
- sample complexity of learning
- uniform law w\ Rademacher Complexity
- Upper Bounds on Rademacher Complexity:
  - Massart Lemma
  - Finite sample Glivenko-Cantelli!
  - VC dimension
  - Dimension Independent generalization bounds
    - linear models
    - convex losses

## Central Set-up:

Let  $\mathcal{F}$  be a family of integrable functions on some probability space  $(X, \Sigma, P)$ .

Goal: Estimate

$$\|P_n - P\|_{\mathcal{F}} := \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(x_i) - \mathbb{E}_{X \sim P}[f(X)] \right|$$

where  $x_1, \dots, x_n \stackrel{iid}{\sim} P$ .

Ex: (Empirical CDF)

Suppose we want to estimate the CDF

$$g(t) := P[X \leq t]$$

We can use the empirical CDF

$$g_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{(-\infty, t]}(x_i)$$

Want to control

$$\|g_n - g\|_{\infty} := \sup_t |g_n(t) - g(t)|$$

This is exactly

$$\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(x_i) - \mathbb{E}[f(X)] \right|$$

where  $\mathcal{F} = \left\{ \frac{1}{x-\alpha, t} : t \in \mathbb{R} \right\}$ .

Thm: (Glivenko-Cantelli): For any distribution

$$\|g_n - g\|_\infty \xrightarrow{\text{a.s.}} 0$$

[We'll prove a finite sample version soon]

Cor (Plug in Estimator) [HW]

Let  $\gamma(\cdot)$  be a  $\|\cdot\|_\infty$ -continuous functional of CDF's. Then

$$\gamma(g_n) \rightarrow \gamma(g) \text{ in probability}$$

Ex:  $\gamma$  quantiles median, mean  
uniform distance to hypothesis' CDF.

Ex: (Statistical Learning)

Consider a family of probability measures

$$\{P_\theta : \theta \in \Omega\}$$

where  $\Omega$  is some set.

Suppose we are given

$$X_1, \dots, X_n \stackrel{\text{iid}}{\sim} P_{\theta^*}$$

for some  $\theta^*$  [may not be in  $\Omega$ ]

Fix a "goodness of fit" function

$$\theta \mapsto L_\theta(X)$$

Goal is to minimize population risk

$$\min_{\theta \in \Omega} R(\theta) := \mathbb{E}_{X^n \sim P_{\theta^*}} [L_\theta(X)]$$

One approach is to instead minimize  
the Empirical Risk:

$$\min_{\theta \in \Omega} R_n(\theta) := \frac{1}{n} \sum_{i=1}^n L_\theta(X_i)$$



Goal: Bound the excess risk

$$E(\hat{\theta}_n) := R(\hat{\theta}_n) - \inf_{\theta \in \Omega} R(\theta)$$

where  $\hat{\theta}_n$  is the minimizer of  $R(\theta)$

Ex: (Maximum Likelihood)

Let  $p_\theta$  be the density of  $P_\theta$

Define

$$L_\theta(x) = \log \left( \frac{p_{\theta^*(x)}}{p_\theta(x)} \right)$$

Then

$$\hat{\theta}_n \in \arg\max_{\theta \in \Omega} \left\{ \frac{1}{n} \sum_{i=1}^n \log(p_\theta(x_i)) \right\}$$

Ex: (Binary Classification)

Data  $(X_i, Y_i) \in \mathbb{R}^d \times \{-1, +1\}$  according to  $P$ .

We want to estimate  $f: \mathbb{R}^d \rightarrow \{-1, +1\}$  solving

$$\min_{f \in T} \underbrace{P[f(X) \neq Y]}_{\text{Some function class}} = \min_{f \in T} \underbrace{E_{(X,Y) \sim P} \frac{1}{2} \{ f(X) \neq Y \}}_{R(f)}$$

Some function class

Try instead to solve

$$\min_{f \in T} \frac{1}{n} \sum_{i=1}^n \mathbb{I}[f(x_i) \neq y_i]$$

## Decomposition of Access Risk

$$E(\hat{\theta}_n, \theta^*) = R(\hat{\theta}_n) - \inf_{\theta \in \Sigma} R(\theta)$$

$$= R(\hat{\theta}_n) - R_n(\hat{\theta}_n)$$

$$+ R_n(\hat{\theta}_n) - R_n(\theta_0)$$

$$+ R_n(\theta_0) - R(\theta_0)$$

where  $\theta_0 = \arg \min_{\theta \in \Sigma} R(\theta)$

(yellow box) is approximation error. How to control?

(cyan box) is optimization error. If  $\hat{\theta}_n$  is true minimizer of  $R_n(\theta)$ , then

$$\leq 0$$

(green box) easy to control approximation error:

$$\mathbb{E}_{\theta} R_n(\theta_0) = R(\theta_0)$$

$$= R(\hat{\theta}_n) - R_n(\hat{\theta}_n)$$

Mark Difficult:  $\hat{\theta}_n$  is not fixed!  
But we can estimate

$$\leq \sup_{\theta \in \mathcal{S}} |R_n(\theta) - R(\theta)|$$

$$= \sup_{\theta \in \mathcal{S}} \left| \frac{1}{n} \sum_{i=1}^n L_\theta(x_i) - \mathbb{E}_{x \sim P_\theta} L_\theta(x) \right|$$

So

$$E(\hat{\theta}_n, \theta^*) \leq 2 \|P_n - P\|_F$$

where  $F := \{x \mapsto L_\theta(x) : \theta \in \mathcal{S}\}$

Remark: Function class  $F$  s.t.

$\|P_n - P\|_F \rightarrow 0$  in probability is  
called Glivenko-Cantelli

- If  $F$  is too big, it may not be Glivenko-Cantelli (see Ex 4.7 in W)

Our first approach will be based on Rademacher Complexity.

Defn:

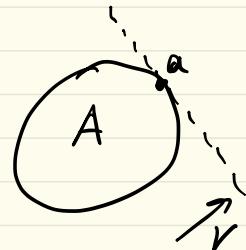
The Rademacher Complexity of a set  $A \subseteq \mathbb{R}^n$ , is the quantity

$$R(A) := \mathbb{E}_{\epsilon} \sup_{a \in A} \langle a, \epsilon \rangle$$

where  $\epsilon = (\epsilon_1, \dots, \epsilon_d)$  are i.i.d. Rademacher random variables

Recall that support function of a set  $A$  is

$$\sigma_A(v) := \sup_{a \in A} \langle v, a \rangle$$



So

$$R(A) = \mathbb{E}_{\epsilon} \sigma_A(\epsilon)$$

Defn: Consider a sequence of random <sup>independent</sup> variables  $(x_1, \dots, x_n) \subseteq \mathcal{X}^n$  and a class of functions  $\mathcal{F}$  on  $\mathcal{X}$ . The Rademacher complexity of  $\mathcal{F}$  is

$$\begin{aligned} R_n(\mathcal{F}) &:= \mathbb{E}_{\substack{x, \epsilon \\ f \in \mathcal{F}}} \sup \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(x_i) \right| \\ &= \mathbb{E}_x R(\mathcal{F}^{\pi}(x), \cdot) \\ &\leq 2 \mathbb{E}_x R(\mathcal{F}(x)/n) \end{aligned}$$

where  $\mathcal{F}(x) = \{(f(x_1), \dots, f(x_n)) : f \in \mathcal{F}\}$

Thm:  $\mathbb{E} \|P_n - P\|_{\mathcal{F}} \leq 2 R_n(\mathcal{F})$

Pf: Let  $y_1, \dots, y_n \stackrel{\text{iid}}{\sim} P$  and independent of  $x_1, \dots, x_n$ . Then

$$\begin{aligned} \mathbb{E} \|P_n - P\|_{\mathcal{F}} &= \mathbb{E}_x \left[ \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(x_i) - \mathbb{E}_y \left[ f(y_i) \right] \right| \right] \\ &= \mathbb{E}_x \left[ \sup_{f \in \mathcal{F}} \left| \mathbb{E}_y \left[ \frac{1}{n} \sum (f(x_i) - f(y_i)) \right] \right| \right] \end{aligned}$$

$$\begin{aligned}
&= \mathbb{E}_X \left[ \sup_{f \in \mathcal{F}} \left| \mathbb{E}_y \left[ \frac{1}{n} \sum (f(x_i) - f(y_i)) \right] \right| \right] \\
&\leq \mathbb{E}_X \mathbb{E}_y \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum (f(x_i) - f(y_i)) \right| \\
&= \mathbb{E}_{X, Y, \epsilon} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i (f(x_i) - f(y_i)) \right| \\
&\leq 2 \mathbb{E}_{X, \epsilon} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(x_i) \right| = 2 R_n(\mathcal{F})
\end{aligned}$$

Thm: Suppose  $\mathcal{F}$  is  $b$ -uniformly bounded, meaning

$$\|f\|_\infty = \sup_{x \in X} |f(x)| \leq b \quad \forall f \in \mathcal{F}.$$

Then for any  $n \in \mathbb{N}$ ,  $t \geq 0$ , it holds

$$\Pr \left( \|P_n - P\|_{\mathcal{F}} \leq 2 R_n(\mathcal{F}) + t \right) \geq 1 - e^{-\frac{n t^2}{2b^2}}$$

pf: All we have to do is show  
 the bounded difference property  
 for the function  $\|P_n - P\|_F$   
 with  $L_i = \frac{2b}{n}$ .

Define  $\bar{f}(x) = f(x) - E f(x)$ .

$$\text{Then } \|P_n - P\|_F = \sup_{f \in F} \left| \frac{1}{n} \sum_{i=1}^n \bar{f}(x_i) \right|$$

As before, let  $x'$  differ from  $x$  only in  $i$ 'th entry. Then

$$\left| \frac{1}{n} \sum_{i=1}^n \bar{f}(x_i) \right| - \sup_{g \in F} \left| \frac{1}{n} \sum_{i=1}^n \bar{g}(x'_i) \right| \\ \leq \left| \frac{1}{n} \sum_{i=1}^n \bar{f}(x_i) \right| - \left| \frac{1}{n} \sum_{i=1}^n \bar{f}(x'_i) \right| \leq \frac{2b}{n}$$

Take  $\sup_{f \in F}$  and exchange  $x, x'$  12

Rademacher complexity characterizes the asymptotics of  $\|P_n - P\|_{\mathcal{F}}$ .

$$\text{Prop } \frac{1}{2} R_n(\bar{\mathcal{F}}) \leq \mathbb{E}_X [\|P_n - P\|_{\mathcal{F}}] \leq 2 R_n(\mathcal{F})$$

where  $\bar{\mathcal{F}} = \{f - \mathbb{E}f : f \in \mathcal{F}\}$ .

[See Prop 4.1.1 in Wainwright]

Prop (Hw) Suppose  $\mathcal{F}$  is  $b$ -bounded  
 Then  $\forall n \in \mathbb{N}, t > 0$ , if holds  
 $\Pr_{\mathcal{F}} (\|P_n - P\|_{\mathcal{F}} \geq \frac{1}{2} R_n(\bar{\mathcal{F}}) - \frac{\sup_{f \in \mathcal{F}} |\mathbb{E}f(X)|}{2\sqrt{n}} - t) \geq 1 - e^{-\frac{nt^2}{2b^2}}$

Goal: Bound the Rademacher Complexity of interesting sets.

Lemma (Basic Calculus)

For any  $A, B \subseteq \mathbb{R}^n$ ,  $c \in \mathbb{R}$ ,  $a_0 \in \mathbb{R}^n$  it holds

- $R(cA) = |c|R(A)$
- $R(A+B) \leq R(A) + R(B)$   
with equality if  $A$  and  $B$  are convex
- $R(A+a_0) = R(A)$
- $R(A) = R(\text{conv}(A))$

Lemma: (Massart) Let  $A = \{a_1, \dots, a_N\} \subseteq \mathbb{R}^n$ .

Define  $\bar{a} = \frac{1}{N} \sum_{i=1}^N a_i$ . Then

$$R(A) \leq \max_{a \in A} \|a - \bar{a}\| \cdot \sqrt{2 \log(N)}$$

pf: WLOG assume  $\bar{a} = 0$ . Let  $\lambda > 0$  and  
 $A' = \{\lambda a_1, \dots, \lambda a_N\}$ . Then

$$R(A') = \mathbb{E}_\varepsilon \max_{a \in A'} \langle \varepsilon, a \rangle$$

$$= \mathbb{E}_\varepsilon \log \left( \max_{a \in A'} e^{\langle \varepsilon, a \rangle} \right)$$

$$\leq \mathbb{E}_\varepsilon \log \left( \sum_{a \in A'} e^{\langle \varepsilon, a \rangle} \right)$$

Jensen

$$\leq \log \left( \sum_{a \in A'} \mathbb{E}_\varepsilon e^{\langle \varepsilon, a \rangle} \right)$$

Observe

$$\mathbb{E}_\varepsilon e^{\langle \varepsilon, a \rangle} = \prod_{i=1}^d \mathbb{E}_\varepsilon e^{a_i \varepsilon_i} \leq \prod_{i=1}^d e^{\frac{a_i^2}{2}} = e^{\|a\|_2^2}$$

$$\leq \log \left( \sum_{a \in A'} e^{\|a\|_2^2} \right)$$

$$\leq \log (|A'| \max_{a \in A'} (e^{\|a\|_2^2}))$$

$$= \log (|A'|) + \max_{a \in A'} \frac{\|a\|_2^2}$$

$$\text{Thus } R(A) = \frac{1}{\lambda} R(A') \leq \frac{\log(|A|) + \max_{a \in A} \frac{\lambda^2 \|a\|^2}{2}}{\lambda}$$

Optimize over  $\lambda > 0$ . □

Defn:  $F$  has polynomial discrimination of order  $V \geq 1$ , if  $\forall n \in \mathbb{N}$  and

$x_1, \dots, x_n \in \mathcal{X}$ , we have

$$\text{card}(F(x_1, \dots, x_n)) \leq (n+1)^V.$$

Cor: Suppose  $F$  has polynomial discrimination of order  $V$ . Then

$$R_n(F) \leq 4 E \left[ \sup_{f \in F} \sqrt{\frac{1}{n} \sum_{i=1}^n f^2(x_i)} \right] \cdot \sqrt{V \log(n+1)}$$

pf: Fix  $x = (x_1, \dots, x_n)$ . Then Massart implies

$$R(F(x)/n) \leq 4 \cdot \frac{\max_{f \in F} \sqrt{\sum_{i=1}^n f_i^2(x)}}{n} \cdot \sqrt{V \log(n+1)}$$

□

In particular, if  $F$  is b-bounded, then

$$R_n(F) \leq 4b \sqrt{\frac{c \log(n+1)}{n}}$$

Cor: (Glivenko-Cantelli)

Let  $g(t) = P[X \leq t]$  be the CDF of  $X$ , and let  $g_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}[X_i \leq t]$  where  $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} P$ . Then

$$P\left[\|g_n - g\|_\infty \geq 8 \sqrt{\frac{\log(n+1)}{n}} + \delta\right] \leq \exp\left(-\frac{n\delta^2}{2}\right)$$

Pf: Let  $\mathcal{F} = \{\mathbb{1}_{(-\infty, t]} : t \in \mathbb{R}\}$ .

$$\begin{aligned} \text{Then } \|g_n - g\|_\infty &= \sup_{t \in \mathbb{R}} \left| \frac{1}{n} \sum_{i=1}^n \mathbb{1}[X_i \leq t] - P[X \leq t] \right| \\ &= \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - E[f(X)] \right| \end{aligned}$$

Since  $F$  is 1-bounded

$$P[\|g_n - g\|_\infty \leq 2R_n(F) + \delta] \geq 1 - \exp\left(-\frac{n\delta^2}{2}\right)$$

Observe  $F$  has poly discrimination with  $V=1$   $\square$

## Vapnik-Chervonenkis (VC) Theory

Method for bounding polynomial /  
discrimination of  $\{0,1\}$ -valued  $F$ .

Consider a class  $F$  of binary valued  
functions on  $X$ .

Def: We say that  $x = (x_1, \dots, x_n)$  is  
shattered by  $F$  if  $\text{card}(F(x)) = 2^n$ .

The VC dimension is

$$\text{VC}(F) = \sup \{ n \in \mathbb{N} : \exists x \in X^n \text{ shattered by } F \}$$

Notation: If  $F = \{ \bigcup_S : \text{some sets } S \}$

set  $S(x) := F(x)$  and  $\text{VC}(S) := \text{VC}(F)$

Ex:  $S_{\text{left}} = \{ (-\infty, a] : a \in \mathbb{R} \} \Rightarrow \text{VC}(S_{\text{left}}) = 1$

$S_{\text{two}} = \{ (a, b] : a, b \in \mathbb{R}, a < b \} \Rightarrow \text{VC}(S_{\text{two}}) = 2$



Thm (Sauer and Shelah)

For any  $x = (x_1, \dots, x_n)$  with  $n \geq VC(S)$ , we have

$$\text{card}(S(x)) \leq \sum_{i=0}^{VC(S)} \binom{n}{i} \leq (n+1)^{VC(S)}$$

Therefore  $F$  has polynomial discriminability of order  $VC(S)$  and

$$R_n(F) \leq 2 \sqrt{\frac{VC(S) \cdot \log(n+1)}{n}} \leftarrow \begin{matrix} \log(n+1) \\ \text{can be removed.} \end{matrix}$$

[See Prop 4.18 in Wainwright for a proof.]

Examples: Let  $G$  be a class of functions. For any  $g: X \rightarrow \mathbb{R}$  define

$$S_g = \{x \in X : g(x) \leq 0\}$$

$$S(G) = \{S_g : g \in G\}$$

Prop: Let  $G$  be a vector space of functions  $g: \mathbb{R}^d \rightarrow \mathbb{R}$  with  $\dim(G) < \infty$ .

Then

$$VC(S(G)) \leq \dim(G)$$

Pf: Set  $n = \dim(G) + 1$  and fix  $x = \{x_1, \dots, x_n\}$  with  $x_i \in \mathcal{X}$ . Define  $L: G \rightarrow \mathbb{R}^n$  by

$$L(g) = (g(x_1), \dots, g(x_n))$$

Since  $n > \dim(G)$ , there exists  $\alpha \neq \gamma \in \mathbb{R}^n$  s.t.  $\langle \alpha, L(g) \rangle = 0 \forall g \in G$ .

$$\Rightarrow \sum_{\{i : \gamma_i \leq 0\}} (-\gamma_i) g(x_i) = \sum_{\{i : \gamma_i > 0\}} \gamma_i g(x_i) \quad \forall g \in G$$

WLOG suppose  $\gamma_i > 0$  for some  $i$ .

Suppose there were  $g \in G$  such that  $S_g$  includes only  $\{x_i : \gamma_i \leq 0\}$ . Then

$$0 \geq \{LHS\} = \{RHS\} > 0$$

Contradiction  $\square$

Ex: (Halfspaces)  
Define  $S_{ab} = \{x \in \mathbb{R}^d : \langle a, x \rangle + b \leq 0\}$

$$S = \{S_{a,b} : a, b \in \mathbb{R}\}$$

$$G = \{x \mapsto \langle a, x \rangle + b : a, b \in \mathbb{R}\}$$

Then

$$VC(S) \leq \dim(G) = d+1.$$

actually equality

Ex: (Balls)

$$\text{Define } S_{a,b} = \{x \in \mathbb{R}^d : \|x - a\|_2 \leq b\}$$

$$S = \{S_{a,b} : a \in \mathbb{R}^d, b \geq 0\}$$

$$\text{Define } g_{a,b}(x) = \|x\|_2^2 - 2\langle a, x \rangle + \|a\|_2^2 - b^2$$

Trick: Define  $\Theta : \mathbb{R}^d \rightarrow \mathbb{R}^{d+2}$  by

$$\Theta(x) = (1, x_1, \dots, x_d, \|x\|_2^2)$$

$$g_c(x) = \langle c, \Theta(x) \rangle \text{ where } c \in \mathbb{R}^{d+2}$$

Then

$$\left\{ g_{a,b} : \begin{array}{l} a \in \mathbb{R} \\ b \geq 0 \end{array} \right\} \subseteq \left\{ g_c : c \in \mathbb{R}^d \right\}$$

vector space of dimension  $d+2$ .

$$\rightarrow VC(\mathcal{S}) \leq d+2$$

[Exact bound is  $d+1$ : harder to prove]

Rademacher Complexity and VC-dim often scale with the dimension of the ambient space  $x \in \mathcal{X}$ .

Thm: Consider  $\min_{w \in W \subseteq \mathbb{R}^d} f(w) = \mathbb{E}_z f(w, z)$  where

$\max_{w \in W} \|w\| \leq B$  and  $f(\cdot, z)$  is  $L$ -Lipschitz

Then 
$$\mathbb{E} \left[ \sup_{w \in W} \left( \frac{1}{n} \sum_{i=1}^n f(w, z_i) - \mathbb{E} f(w, z) \right) \right] \stackrel{\approx f(w)}{\leftarrow} \text{tight} \leq O \left( \sqrt{\frac{L^2 B^2 d \log(n)}{n}} \right)$$
  
[We'll prove this later!].

# Dimension Independent Bound for Generalization

- linear models
- convexity

## Linear Models:

Consider the problem

$$\min_{w \in W} \mathbb{E}_{(a,b) \sim P} l(\langle w, a \rangle, b)$$

To get generalization bounds we need to compute  $R_n(F)$  where

$$F = \{(a, b) \mapsto l(\langle w, a \rangle, b) : w \in W\}$$

Lemma: (Contraction)

Consider a set  $A \subseteq \mathbb{R}^n$  and let

$\phi_i : \mathbb{R} \rightarrow \mathbb{R}$  be a  $\beta$ -Lipschitz function.

Let  $A' = \{\phi_1(a_1), \dots, \phi_n(a_n) : a \in A\}$

Then  $R(A') \leq R(A)$

pf: WLOG, assume  $\beta = 1$ . It suffices  
to assume  $A' = \{(\varnothing, (a_1, a_2, \dots, a_m)) : a \in A\}$

Then

$$\mathcal{R}(A') = \mathbb{E}_{\varepsilon} \left[ \sup_{a \in A'} \sum_{i=1}^n \varepsilon_i a_i \right]$$

$$= \frac{1}{2} \mathbb{E}_{\varepsilon_2, \dots, \varepsilon_n} \left[ \sup_{a \in A} \left\{ \varnothing(a) + \sum_{i=2}^n \varepsilon_i a_i \right\} \right. \\ \left. + \sup_{\hat{a} \in A} \left\{ -\varnothing(\hat{a}) + \sum_{i=2}^n \varepsilon_i \hat{a}_i \right\} \right]$$

$$= \frac{1}{2} \mathbb{E}_{\varepsilon_2, \dots, \varepsilon_n} \left[ \sup_{\substack{a, \hat{a} \in A}} (\varnothing_i(a_i) - \varnothing_i(\hat{a}_i) + \sum_{i=2}^n \varepsilon_i (a_i + \hat{a}_i)) \right]$$

$$\leq \frac{1}{2} \mathbb{E}_{\varepsilon_2, \dots, \varepsilon_n} \left[ \sup_{\substack{a, \hat{a} \in A}} a_i - \hat{a}_i + \sum_{i=2}^n \varepsilon_i (a_i + \hat{a}_i) \right]$$

$$= \mathcal{R}(A)$$

□

# Rademacher Complexity of linear classes

Lemma: Consider a set of vectors  $x_1, \dots, x_n \in \mathbb{R}^d$  and define

$$A = \left\{ (\langle w, x_1 \rangle, \dots, \langle w, x_n \rangle) : \|w\|_2 \leq 1 \right\}$$

Then  $R(A) \leq \sqrt{\sum_{i=1}^n \|x_i\|_2^2}$

and therefore

$$R_n(\{x \mapsto (\langle w, x \rangle : \|w\|_2 \leq 1)\}) \leq 2 \frac{\max_{i=1, \dots, n} \|x_i\|_2}{\sqrt{n}}$$

Pf:

$$\begin{aligned} R(A) &= \mathbb{E}_\epsilon \sup_{w \in A} \sum_{i=1}^m \epsilon_i \cdot a_i \\ &= \mathbb{E}_\epsilon \sup_{\|w\|_2 \leq 1} \sum_{i=1}^m \epsilon_i \cdot \langle w, x_i \rangle \\ &= \mathbb{E}_\epsilon \sup_{\|w\|_2 \leq 1} \left\langle w, \sum_{i=1}^m \epsilon_i \cdot x_i \right\rangle \\ &= \mathbb{E}_\epsilon \left\| \sum_{i=1}^m \epsilon_i \cdot x_i \right\|_2 \\ &\leq \sqrt{\mathbb{E}_\epsilon \left\| \sum_{i=1}^m \epsilon_i \cdot x_i \right\|_2^2} \leq \sqrt{\sum_{i=1}^m \|x_i\|_2^2} \end{aligned}$$

□

Lemma: Consider a set of vectors  $x_1, \dots, x_m \in \mathbb{R}^d$  and define

$$A = \left\{ (\langle w, x_1 \rangle, \dots, \langle w, x_m \rangle) : \|w\|_1 \leq 1 \right\}$$

Then

$$R(A) \leq \sqrt{2n \log(2d)} \cdot \max_{i=1, \dots, n} \|x_i\|_\infty$$

and therefore

$$R_n(\{x \mapsto (\langle w, x \rangle : \|w\|_1 \leq 1)\}) \leq \frac{\sqrt{8R(A)}}{\sqrt{n}} \cdot \max_{i=1, \dots, n} \|x_i\|_\infty$$

$$\text{Pf: } R(A) = \mathbb{E}_\epsilon \sup_{a \in A} \sum_{i=1}^m \epsilon_i \cdot a_i$$

$$= \mathbb{E}_\epsilon \sup_{\|w\|_1 \leq 1} \sum_{i=1}^m \epsilon_i \cdot \langle w, x_i \rangle$$

$$= \mathbb{E}_\epsilon \sup_{\|w\|_1 \leq 1} \left\langle w, \sum_{i=1}^m \epsilon_i \cdot x_i \right\rangle$$

$$= \mathbb{E}_\epsilon \left\| \sum_{i=1}^m \epsilon_i \cdot x_i \right\|_\infty$$

$$= \mathbb{E}_\epsilon \sup_{V \in \{ \pm e_j \}_{j=1}^d} \sum_{i=1}^m \epsilon_i \langle x_i, V \rangle$$

$$\mathbb{E}_\varepsilon \sup_{V \in \{\pm e_i\}_{i=1}^d} \sum_{i=1}^n \varepsilon_i \langle x_i, v \rangle$$

$$= \mathbb{R} \left( (\langle x_1, v \rangle, \dots, \langle x_n, v \rangle) : v \in \{\pm e_i\} \right)$$

$$\leq \sqrt{2 \log(2d)} \cdot \max_{V \in \{\pm e_i\}} \|(\langle x_1, v \rangle, \dots, \langle x_n, v \rangle)\|_2$$

$$\leq \sqrt{2n \log(2d)} \cdot \max_{i=1, \dots, n} \|x_i\|_\infty. \quad \square$$

Back to  $\min_{w \in W} \mathbb{E}_{(a,b) \sim P} l(\langle w, a \rangle, b)$  where  $W = B_2$  or  $B_1$

So if  $l(\cdot, b)$  is  $\beta$ -Lipschitz  $\forall b$ , then the Rademacher bounds are

$$\underline{l_2\text{-case}}: P \cdot \frac{\mathbb{E} \max_{i=1, \dots, n} \|a_i\|_2}{\sqrt{n}}$$

$$\underline{l_1\text{-case}}: P \cdot \sqrt{\log(d)} \cdot \frac{\mathbb{E} \max_{i=1, \dots, n} \|a_i\|_\infty}{\sqrt{n}}$$

Convexity: Regret without uniform laws

Suppose we want to solve

$$\textcircled{X} \quad \min_{w \in W} f(w) = \mathbb{E}_{x \sim P} f(w, x)$$

Let  $S = \{x_1, \dots, x_n\}$  be iid from  $P$ , and let  $t(S)$  be an output of an algorithm which aims to approximately solve  $\textcircled{X}$ . Let  $S^i = (x_1, \dots, x_{i-1}, x', x_{i+1}, \dots, x_n)$  where  $x' \sim P$  independent of  $x_1, \dots, x_n$ . Which  $t(S)$  generalizes?

Intuition:  $f(A(S^i), x_i) - f(A(S), x_i)$  should not be big, otherwise overfitting

Thm:  $E_S [f(A(S)) - \frac{1}{n} \sum_{i=1}^n f(A(S), x_i)]$

$$= E_{\substack{(S, x') \sim P \\ i \sim U(n)}} [f(A(S^i), x_i) - f(A(S), x_i)]$$

Pf: for every  $i$ ,

$$\mathbb{E}_S f(A(S)) = \mathbb{E}_{S^{(2)}} f(A(S), x') = \mathbb{E}_{S^{(2)}} f(A(S^i), x_i)$$

Observe

$$\mathbb{E}_S \frac{1}{n} \sum_{i=1}^n f(A(S), x_i) = \mathbb{E}_{S, i} [f(A(S), x_i)]$$

Defn:  $A(\cdot)$  is leave-one-out stable with rate  $\epsilon(n)$  if

$$\mathbb{E}_{\substack{(S' x') \sim P^{n+1} \\ i \sim l(n)}} [f(A(S^i), x_i) - f(A(S), x_i)] \leq \epsilon(n).$$

Henceforth, fix  $\lambda > 0$  and we'll analyze

$$A(S) := \arg \min_{w \in W} \frac{1}{n} \sum_{i=1}^n f(w, x_i) + \frac{\lambda}{2} \|w\|^2.$$

Also suppose  $W$  and  $f(\cdot; x_i)$  are convex.

Defn: A function  $g: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  is  $\lambda$ -strongly convex if  $g - \frac{\lambda}{2} \|\cdot\|^2$  is convex.

Lemma: If  $g$  is  $\lambda$ -strongly convex, then it has a unique minimizer  $\bar{w}$  and

$$g(w) - g(\bar{w}) \geq \frac{\lambda}{2} \|w - \bar{w}\|^2 \quad \forall w$$

Thm: Suppose  $f$  is convex, and  $f(x)$  is convex and  $\rho$ -Lipschitz. Then the rule

$$A(S) = \arg \min_{w \in W} \frac{1}{n} \sum_{i=1}^n f(w, x_i) + \frac{\lambda}{2} \|w\|^2$$

is leave-one-out-stable with rate  $\frac{2\rho}{\lambda n}$

$$\text{Therefore } E_S [f(A(S)) - \frac{1}{n} \sum_{i=1}^n f(A(S), x_i)] \leq \frac{2\rho}{\lambda n}^2$$

Pf: Define  $f_S(w) := \frac{1}{n} \sum_{i=1}^n f(w, x_i) + \frac{\lambda}{2} \|w\|^2$

$$\Rightarrow f_S(w) - f_S(A(S)) \geq \frac{\lambda}{2} \|w - A(S)\|^2 \quad \forall w$$

For all  $w, v \in \mathbb{R}^n$  it holds

$$\begin{aligned} f_S(w) - f_S(v) &= \frac{1}{n} \sum_{x \in S} f(w, x) + \frac{\lambda}{2} \|w\|^2 \\ &\quad - \frac{1}{n} \sum_{x \in S} f(v, x) - \frac{\lambda}{2} \|v\|^2 \\ &= \frac{1}{n} \sum_{x \in S^c} f(w, x) + \frac{\lambda}{2} \|w\|^2 \\ &\quad - \frac{1}{n} \sum_{x \in S^c} f(v, x) - \frac{\lambda}{2} \|v\|^2 \\ &\quad + \frac{f(w, x_i) - f(v, x_i)}{n} + \frac{f(v, x^*) - f(w, x^*)}{n} \end{aligned}$$

Setting  $w = A(S^c)$ ,  $v = A(S)$  get

$$f_S(A(S^c)) - f_S(A(S)) \leq f_{S^c}(A(S^c)) - f_{S^c}(A(S))$$

$$\begin{aligned} &\leq -\frac{\lambda}{2} \|A(S^c) - A(S)\|^2 \\ &\quad + \frac{2\varphi \|A(S^c) - A(S)\|}{n} \end{aligned}$$

$$\Rightarrow \|A(S^c) - A(S)\| \leq \frac{2\varphi}{\lambda n} \quad \square$$

$$\text{Cor: } \mathbb{E}_S f(A(S)) \leq \min_{\bar{w}} f + \frac{\lambda}{2} \|\bar{w}\|^2 + \frac{2P^2}{n}$$

where  $\bar{w}$  is any minimizer of  $f$  on  $W$ .  
 Therefore under optimal choice  $\lambda = \sqrt{\frac{4P^2}{n\|\bar{w}\|^2}}$   
 get

$$\mathbb{E}_S (f(A(S))) \leq \min f + 2 \sqrt{\frac{P^2 \|\bar{w}\|^2}{n}}$$

$$\begin{aligned} \text{pf: } \mathbb{E}_S f(A(S)) &= \mathbb{E}_S \left[ \frac{1}{n} \sum_{x \in S} f(A(S), x) \right] \\ &\quad + \mathbb{E}_S [f(A(S)) - \frac{1}{n} \sum_{x \in S} f(A(S), x)] \\ &\leq \mathbb{E}_S \left[ \frac{1}{n} \sum_{x \in S} f(A(S), x) \right] + \frac{2P^2}{n} \end{aligned}$$

For any  $w$ , we have

$$\begin{aligned} \mathbb{E}_S \left[ \underbrace{\frac{1}{n} \sum_{x \in S} f(A(S), x)}_{f_S(A(S))} + \frac{\lambda}{2} \|A(S)\|^2 \right] &\leq \mathbb{E}_S f_S(w) \\ &\leq f(w) + \frac{\lambda}{2} \|w\|^2 \\ \Rightarrow \mathbb{E}_S \left[ \frac{1}{n} \sum_{z \in S} f(A(S), z) \right] &\leq \min f + \frac{\lambda}{2} \|\bar{w}\|^2 \end{aligned}$$

# 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  $(T, p)$  be a metric space

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

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

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

•  $L^2([0, 1])$  with

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

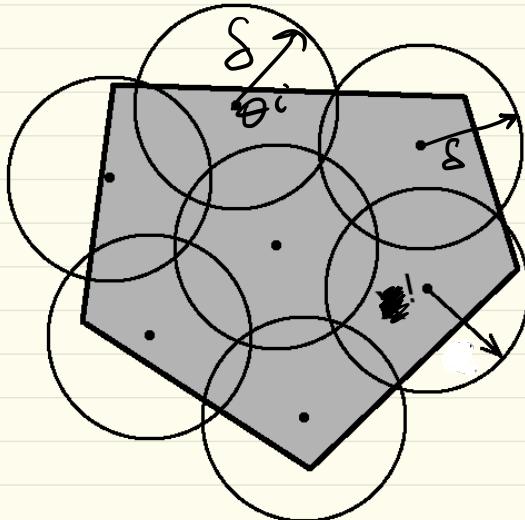
•  $C[0, 1]$  with

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

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

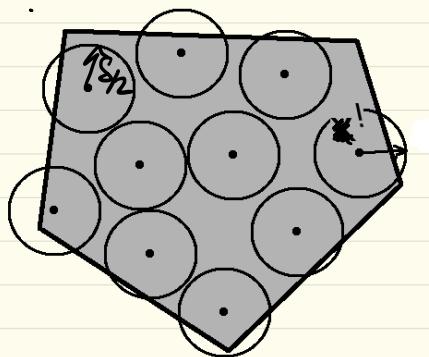
$\forall t \in T \exists i \in \{1, \dots, N\}$  s.t.  $p(t, \theta^i) \leq \delta$

The  $\delta$ -covering number  $N(\delta; T, p)$  is the cardinality of the smallest  $\delta$ -cover. The quantity  $\log(N(\delta; T, p))$  is called the metric entropy.



Defn: A  $\delta$ -packing of  $T$  with respect to  $p$  is a set  $\{\theta^1, \dots, \theta^m\} \subseteq T$  such that  $p(\theta^i, \theta^j) > \delta \quad \forall i, j \in \{1, 2, \dots, m\}$

The  $\delta$ -packing number  $M(\delta; T, p)$  is the cardinality of the largest  $\delta$ -packing.



Lemma:  $\overline{M}(2\delta; T, p) \stackrel{(a)}{\leq} N(\delta; T, p) \stackrel{(b)}{\leq} M(\delta; T, p)$

Pf: ⑥  $M(\delta; T, p)$       ⑦  $M(2\delta; T, p)$

$\xleftarrow{x} \xrightarrow{\theta} \xleftarrow{\delta}$

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

No  $\delta$  Ball can cover  $\theta$ .  $\theta^i$  simultaneously.  $\square$

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

$$B' = \{x: \|x\|' \leq 1\}$$

Then

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

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

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

pf of lemma: If  $\{\Theta_1, \dots, \Theta_n\}$  is a  $s$ -covering of  $B$ , then  $B \subseteq \bigcup_{i=1}^n \{\Theta_i + sB'\}$

$$\Rightarrow \text{Vol}(B) \leq Ns^d \sum_{i=1}^n \text{Vol}(B') \Rightarrow (a)$$

Let  $\{\Theta_1, \dots, \Theta_m\}$  be a maximal  $s/2$ -packing of  $B$ .  $\Rightarrow$  must also be  $s$ -covering.

Taking into account:

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

are disjoint

get:

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

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

■

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

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

Metric entropy can be very large...

Ex: Define

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

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

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

$$A = \sum_{i=1}^r S_i U_i V_i^T = \begin{bmatrix} S_1 & \dots & S_r \\ \vdots & \ddots & \vdots \end{bmatrix} \begin{bmatrix} U_1 & & \\ & \ddots & \\ & & U_r \end{bmatrix} \begin{bmatrix} V_1 \\ \vdots \\ V_r \end{bmatrix}$$

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

$$S_1 \geq S_2 \geq S_3 \geq \dots \geq S_n \geq 0$$

$$\bullet S_i(A) = \sqrt{\lambda_i(AA^T)} = \sqrt{\lambda_i(A^TA)} \quad i=1, \dots, r.$$

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

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

$$\|A\|_2 = \sup_{\|x\|_2 \leq 1} \|Ax\|_2 = S_1(A) = \|S(A)\|_\infty$$

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:

$$\|A\|_2 \lesssim \sqrt{m} + \sqrt{n}$$

with high probability.

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

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

p.f.:  $\textcircled{a}$  is clear. To see  $\textcircled{b}$ , suppose

$$\|A\|_2 = \langle A_{x,y} \rangle \text{ for some } x \in S^{n-1}, y \in S^{m-1}.$$

$$\exists \hat{x} \in N, \hat{y} \in M \text{ with } \|x - \hat{x}\| \leq \varepsilon, \|y - \hat{y}\| \leq \varepsilon$$

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

$$\leq \langle A_{\hat{x}, \hat{y}} \rangle + 2\varepsilon \|A\|_2 + \|A\|_2 \varepsilon^2$$

$$\Rightarrow \|A\|_2 \leq \frac{1}{1-2\varepsilon-\varepsilon^2} \langle A_{\hat{x}, \hat{y}} \rangle \quad \square$$

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

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

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

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

p.f: Step 1: Approximation

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

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

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

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

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

## Step 2: Concentration

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

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

$$\begin{aligned} \text{So Hoeffding } \Rightarrow & \frac{-u^2}{2 \max_{i,j} |A_{ij}|^2} \quad u > 0 \\ P[\langle Ax, y \rangle \geq u] &\leq e^{-\frac{u^2}{2 \max_{i,j} |A_{ij}|^2}} \\ \Rightarrow P\left[\max_{\substack{x \in X \\ y \in M}} \langle Ax, y \rangle \geq u\right] &\leq e^{-\frac{u^2}{2 \max_{i,j} |A_{ij}|^2}} \end{aligned}$$

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

Then

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

by choosing  $C$  large  $\blacksquare$

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

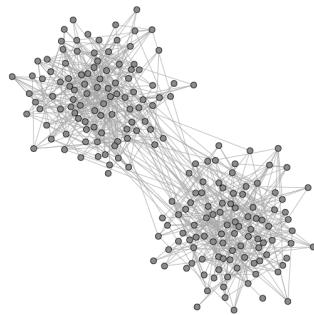
Obvious analogues hold for symmetric matrices.

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

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



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: (Hh)

Show that  $D$  has rank 2 and

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

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

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  $\|D\| = \lambda_1 \sim n$

$$\Pr \left[ \|R\| \leq C\sqrt{n} \right] \geq 1 - 4e^{-n}$$

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  $Q_1 = \{ i : [V_2(A)]_i > 0 \}$   
 $Q_2 = \{ i : [V_2(A)]_i \leq 0 \}$ .

## 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

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

Then

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

Remark: This implies

$$\min_{\theta \in \mathbb{R}} \|V_i(X) - \Theta V_i(\theta)\|_2 \leq \frac{2\|X - Y\|_2}{\delta^{\frac{3}{2}}}$$

Let's compute the eigengap

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

$\underbrace{\quad}_{\parallel}$

$\overset{M}{\vdots}$

$$\Rightarrow \min_{D \in \{I, S\}} \|V_2(D) - \Theta V_2(A)\|_2 \leq \frac{\|D - A\|_2}{\sqrt{mn}} = \frac{\|R\|_2}{\sqrt{mn}}$$

$$\Rightarrow P \left[ \max_{D \in \{I, S\}} \|V_2(D) - \Theta V_2(A)\|_2 \leq \frac{C}{\sqrt{mn}} \right] \geq 1 - 4e^{-n}$$

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

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

We next obtain two-sided bounds on the full spectrum

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

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  $\sigma_i$ -sub-Gaussian isotropic. Then

$$\sqrt{m} - CK^2(\sqrt{n} + t) \leq \sigma_n(A) \leq \sigma_1(A) \leq \sqrt{m} + CK^2(\sqrt{n} + t)$$

w.p.  $1 - 2\exp(-t^2)$ , where  $K = \max_i \sigma_i$ .  
 p.f. 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  $N$  of  $\mathbb{S}^{n-1}$ , with  $\text{card}(N) \leq 9^n$ .

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

and

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

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

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

Step 2: Concentration

Fix  $w \in V$  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,  $\sigma_i$ -subGaussian  
 with  $\mathbb{E} X_i^2 = 1$ . Therefore, Bernstein

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

$$2\exp\left[-\frac{1}{2}\left(\frac{mt^2}{K^4} - \frac{1}{4K^2}\right)\right] = 2\exp\left[-c_1 m \sqrt{\frac{t^2}{K^4} + \frac{t}{K^2}}\right]$$

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

$$\Rightarrow = 2\exp\{-c_1 m \delta^2\}$$

$\leq$

Step 3: Union bound

$$\begin{aligned} & \mathbb{P}\left[\max_{W \in \mathcal{N}} \left| \frac{1}{m} \|Aw\|_2^2 - 1 \right| \geq K \max\{\delta, \delta^2\} \right] \\ & \leq 9 \cdot 2 \exp\{-c_1 C^2 (n + t^2)\} \leq 2e^{-t^2} \end{aligned}$$

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  $1 - \delta \leq S_n(A) \leq S_1(A) \leq 1 + \delta$

□

$$\text{Cor: (HW)} \quad \mathbb{E} \left\| \frac{1}{m} A^T A - I \right\|_2 \leq C K^2 \left( \sqrt{\frac{n}{m}} + \frac{n}{m} \right)$$

## Covariance Estimation

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

Define  $\Sigma = \mathbb{E} X X^T$ . Let's estimate  $\Sigma$  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  $m > n$ . We'll see  $m = n$  suffice

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

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

for some  $K > 0$ .

Then

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

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

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

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

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

So  $Z$  and  $Z_i$  are  $K$ -sub-Gaussian

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

Define  $A$  whose rows are  $Z_i$ . Then

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



So can get

$$\mathbb{E} \left\| \Sigma_m - \Sigma \right\|_2 \leq \varepsilon \left\| \Sigma \right\|_2$$

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

Ex: (Hw) Actually get the high probability estimate

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

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

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

get

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

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  $x_i \sim N(\mu, I)$

- if tails choose  $x_i \sim N(-\mu, I)$

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  $\Sigma_m = \frac{1}{m} \sum_{i=1}^m x_i x_i^\top$

2) Compute eigenvector of  $\Sigma_m$  corresponding to  $\lambda_1(\Sigma_m)$ .

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.  $\|M\|_2 \geq C \sqrt{\log(\frac{1}{\delta})}$ .

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

pf sketch: Note  $X = E M + g$  where  $E$  is a Rademacher RV and  $g \sim N(0, I)$   
 $\Rightarrow E X X^T = E(M M^T + g g^T) = I + M M^T$

$\Rightarrow E X X^T$  has eigenvalues  $1 + \|M\|^2, 1, \dots, 1$   
Eigenvector corresponding to  $1 + \|M\|^2$  is  $M$   
So from the e-vector corresponding to the largest eigenvalue, you can "guess" at the communities by checking the sign of  $\langle x_i, M \rangle$ .  
You do the rest in your HW.

Our next goal is to investigate fails  
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
- $\mathcal{S}_+^d = \{ Q \in \mathcal{S}^d : \langle Qx, x \rangle \geq 0 \forall x \in \mathbb{R}^d \}$
- $\gamma_1(Q) \geq \gamma_2(Q) \geq \dots \geq \gamma_d(Q)$  are the eigenvalues of  $Q \in \mathcal{S}^d$
- $\mathcal{O}^d = \{ U \in \mathbb{R}^{d \times d} : U^T U = I \}$ .

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

$$f(Q) = U \operatorname{diag}(f(\gamma_1(Q)), \dots, f(\gamma_d(Q))) U^T$$

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

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

Note

$$\delta(f(Q)) = \{f(\delta_j(Q)) : j=1 \dots d\}$$

Matrix exponential:  $e^Q$

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

The Löwner ordering on

$$X \preceq Y \iff Y - X \in S_{++}^d$$

Lemma:  $(Hw) Q_2 \succcurlyeq Q_1, \forall 0 \Rightarrow \log(Q_2) \succcurlyeq \log(Q_1)$

[ $e^Q$  is not monotone]

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

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

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

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  $Q \in \mathbb{S}^d_+$ , is

the function  $\Psi_Q: \mathbb{R} \rightarrow \mathbb{S}^d$  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  $Q \in \mathbb{S}^d_+$  is

$V$ -subGaussian (with  $V \in \mathbb{S}^d_+$ ) if

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

Ex: Suppose  $Q = EA$  where  $E_i$  is a Rademacher RV and  $A \in \mathbb{S}^d$  is fixed.

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

So  $Q$  is  $A$ -subGaussian.

## Exercise (HW)

If  $Q = gA$ , where  $g \in \mathbb{R}$  is a symmetric 6-subGaussian variable, then  $Q$  is  $\sigma^2 A^2$  subGaussian.

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

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

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

Defn: A zero-mean random matrix  $Q$  is subexponential with parameters  $(V, \alpha)$  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  $\varepsilon$  is Rademacher  
and  $g \sim N(0, 1)$   
independent.

[You'll compute  $V$  and  $\lambda$  in HW]

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

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

$$\begin{aligned} \text{pf: } E e^{\lambda Q} &= \sum_{k=0}^{\infty} \frac{\lambda^k E Q^k}{k!} = \\ &= I + \frac{\lambda^2 \text{Var}(Q)}{2} + \sum_{j=3}^{\infty} \frac{\lambda^j E Q^j}{j!} \\ &\leq I + \frac{\lambda^2 \text{Var}(Q)}{2} + \sum_{j=3}^{\infty} \frac{\lambda^2 \lambda^{j-2} E Q^j b^{j-2}}{j!} \\ &\leq I + \frac{\lambda^2 \text{Var}(Q)}{2} + \sum_{j=3}^{\infty} \frac{\lambda^2 \text{Var}(Q) \sum_{i=1}^{\infty} \lambda^{j-2-i} b^{j-2-i}}{j!} \\ &= I + \frac{\lambda^2 \text{Var}(Q)}{2(1 - b|\lambda|)} \leq \exp\left(\frac{\lambda^2 \text{Var}(Q)}{2(1 - b|\lambda|)}\right) \end{aligned}$$

## Matrix Chernoff:

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

$$P[\gamma_1(Q) \geq t] \leq \text{tr}(\Psi_Q(\lambda)) e^{-\lambda t} \quad \forall \lambda \in [0, \alpha]$$

Consequently

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

pf: We begin as in scalar case:

$$\begin{aligned} P[\gamma_1(Q) \geq t] &= P[e^{\gamma_1(Q)} \geq e^{\lambda t}] \\ &= P[\gamma_1(e^{\lambda Q}) \geq e^{\lambda t}] \\ &\leq E[\gamma_1(e^{\lambda Q})] e^{-\lambda t} \\ &\leq E[\text{tr } e^{\lambda Q}] e^{-\lambda t} \\ &\leq \text{tr}[\Psi_Q(\lambda)] e^{-\lambda t} \quad \square \end{aligned}$$

Lemma: Let  $Q_1, \dots, Q_n$  be independent with  $\Psi_{Q_i}(\lambda)$  existing for all  $\lambda$  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

$$\text{Pr}\left[\left\|\frac{1}{n} \sum_{i=1}^n Q_i\right\|_2 \geq t\right] \leq 2 \text{tr}\left(e^{\sum_{i=1}^n \log \Psi_{Q_i}(\lambda)}\right) e^{-xtn}$$

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

Then

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

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

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

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

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

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

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  $Q_1, \dots, Q_n \in \mathbb{R}^d$  be zero-mean, independent, with  $\|Q_i\|_2 \leq b$ . Then

$$\begin{aligned} \mathbb{P}\left[\left\|\frac{1}{n} \sum_{i=1}^n Q_i\right\|_2 \geq t\right] &\leq 2 \text{rank}(V) \exp\left(-\frac{nt^2}{2(6^2 + bt^2)}\right) \\ &\leq 2 \text{rank}(V) \exp\left(-\frac{nt^2}{46^2} \wedge \frac{nt}{2b}\right) \end{aligned}$$

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

Cor: Under same assumptions,

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

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  $Q_i := \begin{bmatrix} O_{d_1 \times d_2} & A_i \\ A_i^T & O_{d_1 \times d_2} \end{bmatrix}$

$$\frac{\text{Lemma (Hh)}}{6^2} \|Q_i\|_2 = \|A_i\| \quad \text{and} \quad \left\| \frac{1}{n} \sum_{i=1}^n E[A_i A_i^T] \right\|_2 = \left\| \frac{1}{n} \sum_{i=1}^n E[A_i A_i^T] \right\|_2$$

and therefore

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

## Covariance Estimation for general distributions

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

$E[X] = 0$ . Define  $\Sigma = E[XX^T]$  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 subGaussian, we proved

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

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

Then

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

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

pf: Matrix Bernstein

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

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

Let's estimate  $\sigma^2$  and  $M$ .

$$\mathbb{E}(xx^T - \Sigma)^2 = \mathbb{E}(xx^T)^2 - \Sigma^2 \lesssim \mathbb{E}(xx^T)^2$$

and

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

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

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

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  $\Sigma$  is sparse, we should be able to estimate it using fewer observations. Goal: replace  $\|\Sigma_n - \Sigma\|_2 \leq \sqrt{\frac{C}{n}}$  by  $\|\Sigma_n - \Sigma\|_2 \leq \sqrt{\frac{\ell \log(d)}{n}}$  where " $\ell$ " 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

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

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

p.f.:  $A + (s-1)I$  is diagonally dominant and therefore positive semidefinite.  
 $\Rightarrow |\lambda_d(A)| \leq s-1$ .

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

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

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

it holds

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

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

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

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

Pf: 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 semi-definite,  $\|A^2\|_2 = \gamma(A^2)$   
 $\|B^2\|_2 = \gamma(B^2)$

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

$$\gamma(A^2) = \sup_{\|u\|_2 \leq 1} \sum (A^2)_{ij} u_i u_j \leq \sup_{\|u\|_2 \leq 1} \sum (B^2)_{ij} |u_i| |u_j| = \gamma(B^2)$$

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

$$\Rightarrow \|A\|_2^2 = \|A^2\|_2 \leq \|B^2\|_2 \leq \|B\|_2^2 \quad \boxed{\text{Pf}}$$

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

$$\Rightarrow \|\Sigma_n - \Sigma\|_2 \leq 2\lambda_n \|A\|_2$$

Lemma:

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

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

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

If  $i \neq j$ , then

$$2\Delta_n^{ij} = \frac{1}{n} \sum_{k=1}^n (\mathbf{x}_k)_i (\mathbf{x}_k)_j - 2\sum_{ij}$$

$$= \frac{1}{n} \sum_{k=1}^n \left[ (\mathbf{x}_k)_i + (\mathbf{x}_k)_j \right]^2 - \left[ \sum_{ii} + \sum_{jj} + 2\sum_{ij} \right]$$

$$- \left( (\mathbf{x}_k)_i^2 + (\mathbf{x}_k)_j^2 - \sum_{ii} - \sum_{jj} \right)$$

Apply Bernstein again + union bound.



Notice the bound

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

Can be  $\sim d\sqrt{\frac{\log d}{n}}$ , if  $\Sigma$  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.

pf: Write

$$A = \sum_{i \leq j} Z_{ij} \quad \text{where}$$

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

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

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

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

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

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

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  $\mathbb{E}_Z \sqrt{\left\| \sum_{i,j} Z_{ij} \right\|_2^2}$

Quick computation shows

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

diagonal.

$\Rightarrow$

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

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

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

$$\mathbb{E} \|A\|_2 \leq C \sqrt{\log(m+n)} (\mathbb{E}_{\max_i} \|A_{ii}\|_2 + \mathbb{E}_{\max_j} \|A_{jj}\|_2)$$

Application: Matrix Completion.

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

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

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

Thm: Let  $\hat{X}$  be a best rank- $r$  approximation to  $\bar{P}^{-1}Y$ . Then

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

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

So in order to ensure

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

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

pf: First bound error is operator norm

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

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

↑  
independent mean-zero

$$\Rightarrow \mathbb{E} \|Y - pX\|_2 \leq C \sqrt{\log(n)} \cdot \left( \mathbb{E}_{i \max} \| (Y - pX)_i \|_2 + \mathbb{E}_{j \max} \| (Y - pX)_j \|_2 \right)$$

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  $\mathbb{E}_{i \max} \sum_j (S_{ij} - p)^2 \leq Cpd$

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  $2r$

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

$\boxed{\square}$

## Chapter 5

### Quadratic forms and mean estimation

- Mean estimation for  $N(\mu, \Sigma)$
- Hanson-Wright inequality
- norm concentration w/out isotropy
- projections of random vectors.

As motivation, let's find the sample complexity of estimating the mean of a multivariate normal.

Thm: Suppose  $x_i \stackrel{\text{i.i.d.}}{\sim} N(\mu, \Sigma)$ . Then

$$P\left[\left\|\frac{1}{n} \sum_{i=1}^n x_i - \mu\right\|_2 \leq \sqrt{\frac{\text{tr}(\Sigma)}{n}} + \sqrt{\frac{2 \log\left(\frac{2}{\delta}\right) \|\Sigma\|_2}{n}}\right] \geq 1 - \delta$$

Read Lugosi-Mendelson survey! K530

pf: Fix a matrix  $A \in \mathbb{R}^{d \times d}$  and  $z \sim N(0, I)$

$$\text{Then } \|Az\|_2^2 = z^T A^T A z = \sum_{i,j} (A^T A)_{ij} z_i z_j$$

Write the eigenvector-decomposition

$$A = \sum_{i=1}^d \lambda_i u_i u_i^T$$

Then

$$z^T A^T A z = z^T \left( \sum_{i=1}^d \lambda_i^2 u_i u_i^T \right) z = \sum_{i=1}^d \lambda_i^2 \underbrace{\langle u_i, z \rangle}_{N(0, 1)}^2$$

Notice  $\{\langle u_i, z \rangle\}_{i=1}^m$  are independent.

Bernstein  $\Rightarrow$

$$\mathbb{P}\left[\left|\|A_2\|_2^2 - \|A\|_F^2\right| > t\right]$$

$$\leq 2 \exp\left[-\frac{1}{2} \left( \frac{t^2}{\sum_i \mathbf{x}_i^T \mathbf{x}_i} \wedge \frac{4t}{\|A\|_2^2} \right)\right]$$

$\sim N(0, I)$

Finally set  
 $A = \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{x}_i^{1/2}$

$$Z = \underbrace{\frac{1}{\sqrt{n}} \sum_{i=1}^n}_{\mathbf{1}} \sum_{i=1}^n -\frac{1}{2} (\mathbf{x}_i - \mu)^T$$

Then

$$\mathbb{P}\left[\left|\left\|\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i - \mu\right\|^2 - \frac{\text{tr}(\Sigma)}{n}\right| > t\right]$$

$$\leq 2 \exp\left[-\frac{1}{2} \left( \frac{t^2}{\|\Sigma_F\|^2/n} \wedge \frac{t}{\|\Sigma\|_2} \right) n\right]$$

◻

We will now generalize results of this type to sub-Gaussians. Namely we want to understand

$$|x^T A x - \mathbb{E} x^T A x|$$

where the coordinates of  $x = (x_1, \dots, x_n)$  are independent and subGaussian.

Then (Hanson-Wright),

Let  $x = (x_1, \dots, x_n) \in \mathbb{R}^n$  be a random vector with  $x_i$  independent, mean zero, and  $\sigma$ -subGaussian. Let  $A \in \mathbb{R}^{n \times n}$  be arbitrary.

Then

$$\mathbb{P}[|x^T A x - \mathbb{E} x^T A x| \geq t] \leq 2 \exp\left(-c \min\left(\frac{t^2}{6^4 \|A\|_F^2}, \frac{t}{6^2 \|A\|_2}\right)\right)$$

See Vershynin 6.1-6.2 for the proof

Cor: Let  $B \in \mathbb{R}^{m \times n}$  be a fixed matrix.  
 and let  $x = (x_1, \dots, x_n) \in \mathbb{R}^n$  have  
 independent  $\sigma$ -subGaussian coordinates  
 with  $\mathbb{E} x_i = 0$  and  $\mathbb{E} x_i^2 = 1$

Then

$$\mathbb{P}\left[\left|\|Bx\|_2 - \|B\|_F\right| > t\right] \leq 2e^{-\frac{ct^2}{6^4\|B\|_F^2}}$$

pf: Apply Hanson-Wright with  
 $A = B^T B$ . Compute

$$\begin{aligned} x^T A x &= \|Bx\|_2^2, \quad \mathbb{E} x^T A x = \mathbb{E} \text{tr}(A x x^T) = \text{tr}(B^T B) \\ &= \|B\|_F^2 \end{aligned}$$

$$\|A\|_2 = \|B\|_2^2, \quad \|A\|_F = \|B^T B\|_F \leq \|B\|_2 \cdot \|B\|_F$$

$$\Rightarrow \mathbb{P}\left[\left|\|Bx\|_2^2 - \|B\|_F^2\right| \geq u\right] \leq 2\exp\left(-\frac{c}{6^4} \left(\frac{u^2}{\|B\|_2^2 \|B\|_F^2} \right)\right)$$

$\Rightarrow$  continue like in the concentration  
 of the norm  $\|x\|$ .  $\square$

Cor: Let  $E \subset \mathbb{R}^n$  be a subspace with  $\dim(E) = d$ . Let  $x = (x_1, \dots, x_n)$  have independent  $\sigma$ -subGaussian coordinates with  $\mathbb{E}x_i = 0$ ,  $\mathbb{E}x_i^2 = 1$ .

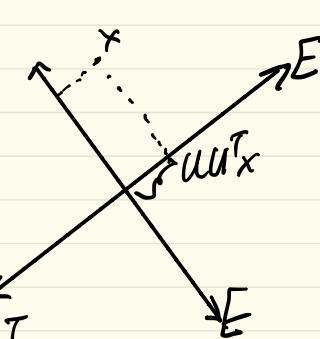
Then

$$\mathbb{P}[|d(x, E) - \sqrt{n-d}| > t] \leq 2 \exp\left(-\frac{ct^2}{\sigma^4}\right)$$

Pf: Let  $U \in \mathbb{R}^{n \times (n-d)}$  have as columns an orthonormal basis for  $E^\perp$

Then  $\text{Proj}_{E^\perp} = UU^\top$  and

$$\text{dist}(x, E) = \|UU^\top x\|_2$$



Then

$$\begin{aligned} \mathbb{E} \text{dist}^2(x, E^\perp) &= \mathbb{E} x^\top U U^\top U U^\top x \\ &= \mathbb{E} x^\top U U^\top x = \text{tr}(UU^\top) = n-d \end{aligned}$$

$$\Rightarrow \mathbb{E} \text{dist}(x, E^\perp) \leq \sqrt{n-d} \quad \text{Apply previous result } B = UU^\top \blacksquare$$

# Chapter 6

## Gaussian Processes

- Concentration of Lipschitz functions of the uniform spherical random vector
- Basics of Gaussian processes.
- Slepian's Inequality.
- Sudakov-Fernique Inequality

In this chapter, we explore a number of constructions built from gaussians. We begin with the following

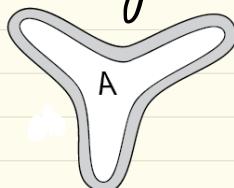
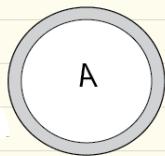
Thm: Let  $f: \mathbb{S}^{n-1} \rightarrow \mathbb{R}$  be Lipschitz with constant  $L_f$  and let  $x \sim \text{Unif}(\mathbb{S}^{n-1})$ . Then

$$P[|f(x) - E[f(x)]| \geq t] \leq 2 \exp\left(-\frac{ct^2}{L_f^2}\right)$$

The proof will be based on the isoparametric inequality

Thm: Let  $\epsilon > 0$ . Then among all sets  $A \in \mathbb{S}^{n-1}$  with prescribed area  $\mu^{n-1}(A)$ , the spherical caps minimize the area of the blow-up

$$A_\epsilon = \{x \in \mathbb{S}^{n-1} : \exists y \in A \text{ s.t. } \|x-y\|_2 \leq \epsilon\}$$



Recall that  $\mathbb{S}^{n-1}$  is a natural scaling of the sphere because

$\underline{X} \sim \text{Unif}(\mathbb{S}^{n-1})$  is isotropic

and  $c$ -sub Gaussian for a constant  $c$ .

[Vershynin Thm 3.4.6 :  $\underline{X} \sim \mathbb{S}^{n-1} \frac{\underline{g}}{\|\underline{g}\|_2}$  where  
 $\underline{g} \sim N(0, I)$  and  $\|\underline{g}\|_2 \sim \sqrt{n}$  w.h.p.]

Lemma: Fix  $A \subset \mathbb{S}^{n-1}$  with  $\mu^{n-1}(A) \geq \frac{1}{2}$   
 where  $\mu^{n-1}$  denotes the normalized area on  $\mathbb{S}^{n-1}$ . Then

$$\mu^{n-1}(A_t) \geq 1 - 2e^{-ct^2} \quad t \geq 0.$$

pf: Define

$$H := \{x \in \mathbb{S}^{n-1} : x_i \leq 0\}$$

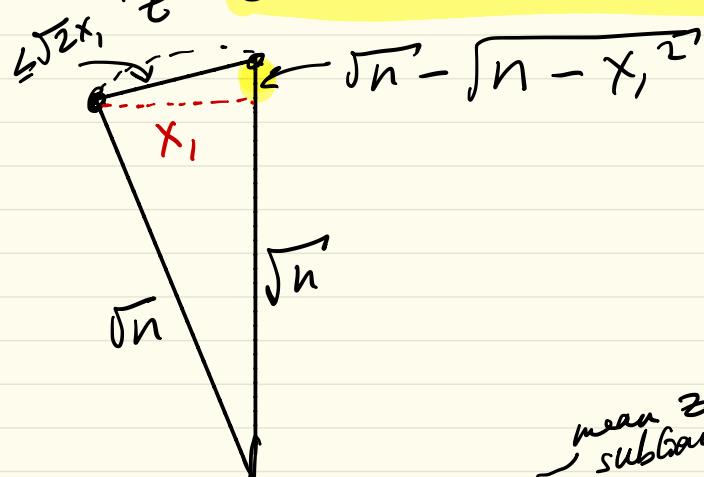
By assumption  $\mu^{n-1}(A) \geq \frac{1}{2} = \mu^{n-1}(H)$

$\Rightarrow \mu^{n-1}(A_t) \geq \mu^{n-1}(H_t)$  ← let's compute

$$M^{n-1}(H_t) = P[\sum_{i \in H_t} f^{n-1}_i]$$

Notice

$$H_t = \left\{ x \in \sum_{i=1}^n f^{n-1}_i : x_i \leq \frac{t}{\sqrt{n}} \right\}$$



mean zero  
subgaussian

$$\Rightarrow M^{n-1}(H_t) \geq P\left[X_i \leq \frac{t}{\sqrt{n}}\right] \geq 1 - 2 \exp(-ct^2)$$

pf of Thm ⑦:

Let  $M$  denote a median of  $f(X)$ , namely

$$P[f(X) \leq M] \geq \frac{1}{2} \text{ and } P[f(X) \geq M] \leq \frac{1}{2}$$

Define

$$A := \{x \in \mathbb{R}^n : f(x) \leq M\}$$

Since  $P[X \in A] \geq \frac{1}{2}$ , we just showed

$$P[X \in A_t] \geq 1 - 2e^{-ct^2}$$

Notice

$$P[X \in A_t] \leq P[f(X) \leq M + L_f t]$$

Since for  $\forall x \in A_t \exists y \in A$  s.t.

$$f(x) \leq f(y) + L_f \|x-y\| \leq M + L_f t$$

$$\text{So } P[f(X) \leq M + L_f t] \leq 1 - 2e^{-ct^2}$$

Replace  $f$  by  $-f$  proves the lower deviation inequality. Finally, the centering inequality (Vershynin 2.6.8)

$$\begin{aligned} \|\mathbb{E}[f(X)] - f(X)\|_{\ell_2} &= \|(\mathbb{E}[f(X)] - M) - (\mathbb{E}[f(X)] - M)\|_{\ell_2} \\ &\leq \|f(X) - M\|_{\ell_2} \leq cL_f \end{aligned}$$

QED

Cor: If  $f: S^{n-1} \rightarrow \mathbb{R}$  is Lipschitz and  $\underline{X} \sim \text{Unif}(S^{n-1})$ , then

$$\mathbb{P}[|f(\underline{X}) - \mathbb{E}f(\underline{X})| \geq t] \leq 2 \exp\left(-\frac{Cn t^2}{L_f^2}\right)$$

This result extends to many other settings

Thm: Let  $\epsilon > 0$ . Then among all sets  $A \in \mathbb{R}^n$  with prescribed Gaussian measure  $\gamma_n(A)$ , the half-spaces minimize the Gaussian measure  $\gamma_n(A_\epsilon)$

(HW)  
Thm: Consider a  $\underline{X} \sim N(0, I)$  and a Lipschitz  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ . Then

$$\mathbb{P}[|f(\underline{X}) - \mathbb{E}f(\underline{X})| > t] \leq \exp\left(-\frac{ct^2}{L_f^2}\right)$$

See Vershynin 5.2 for other examples of probability spaces on which functions concentrate.

Typical results require  $\mathbb{X}$  to be very symmetric. One important exception is the following.

Thm (Talagrand)

Consider a random  $\bar{\mathbf{X}} = (\bar{x}_1, \dots, \bar{x}_n) \in \mathbb{R}^n$  whose coordinates are independent and satisfy  $|\bar{x}_i| \leq 1$  almost surely. Then for any convex Lipschitz  $f: [0, 1]^n \rightarrow \mathbb{R}$ ; it holds

$$P\{|f(\bar{\mathbf{X}}) - \mathbb{E}f(\bar{\mathbf{X}})| \geq t\} \leq 2 \exp\left(-\frac{Cn t^2}{L_f^2}\right)$$

[See Boucheron, Lugosi, Massart for a proof]

Defn: A random process is a collection of random variables  $\{X_t : t \in T\}$  indexed by some set  $t \in T$ .

Let  $\{X_t\}_{t \in T}$  be a random process.  
Assume for simplicity  $E X_t = 0 \quad \forall t \in T$ .  
Define the covariance function

$$\Sigma(t, s) := \text{cov}(X_t, X_s) = E X_t X_s$$

and the increments

$$d(t, s) := \sqrt{E(X_t - X_s)^2}$$

Rem:  $d(\cdot, \cdot)$  is a pseudometric on  $T$ .

$$\text{Rem: } d^2(t, s) = \Sigma(t, t) + \Sigma(s, s) - 2 \Sigma(t, s)$$

Defn:  $\{X_t\}_{t \in T}$  is called a Gaussian process if for any finite  $T_0 \subset T$ , the random vector  $(X_t)_{t \in T_0}$  has a Gaussian distribution.

Main example is the canonical Gaussian Process indexed by  $T \subseteq \mathbb{R}^n$ :

$$X_t = \langle g, t \rangle \quad \forall t \in T$$

where  $g \sim N(0, I)$ .

Notice

$$d(t, s) = \sqrt{\mathbb{E} \langle g, t-s \rangle^2} = \sqrt{\mathbb{E} \text{tr}(t-s)(t-s)^T g g^T} = \|t-s\|_2$$

Our goal is to estimate

$$\mathbb{E} \sup_{t \in T} X_t$$

or more formally

$$\sup_{T \subseteq T \text{ finite}} \mathbb{E} \sup_{t \in T} X_t.$$

Thm (Slepian's inequality)  
 Let  $\{X_t\}_{t \in T}$ ,  $\{Y_t\}_{t \in T}$  be mean-zero Gaussian processes. Assume  $X_t, Y_t$  :

$$E X_t^2 = E Y_t^2 \text{ and } E(X_t - X_s)^2 \leq E(Y_t - Y_s)^2$$

Then

$$P\left[\sup_{t \in T} X_t \geq s\right] \leq P\left[\sup_{t \in T} Y_t \geq s\right]$$

Therefore

$$E \sup_{t \in T} X_t \leq E \sup_{t \in T} Y_t.$$

Proof Strategy:

Assume  $T$  is finite and set  $X = (X_t)$ ,  $Y = (Y_t)$ ,  
 $Z(u) = \sum X + \sqrt{1-u} Y \quad \text{for } u \in [0, 1].$

We will show that

$$u \mapsto P\left[\max_i Z_i(u) \geq s\right] \text{ is non-increasing}$$

Lemma (Hw)

Let  $X \sim N(0, \Sigma)$ . Then for any differentiable  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ , it holds:

$$\mathbb{E} X f(X) = \sum_i \mathbb{E} Df(x)$$

or equivalently

$$\mathbb{E} X_i f(X) = \sum_{j=1}^n \sum_{ij} \mathbb{E} \frac{\partial^2 f}{\partial X_j}(x), \quad i=1, \dots, n$$

Lemma: (Interpolation) Consider independent  $X \sim N(0, \Sigma^X)$ ,  $Y \sim N(0, \Sigma^Y)$ , and set

$$Z(u) = \sqrt{u} X + \sqrt{1-u} Y \quad \text{for } u \in [0, 1].$$

Then for any twice-differentiable  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  it holds:

$$\frac{d}{du} \mathbb{E} f(Z(u)) = \frac{1}{2} \left\langle \sum_i \frac{X_i}{\sqrt{u}} - \sum_i \frac{Y_i}{\sqrt{1-u}}, \mathbb{E} D^2 f(Z(u)) \right\rangle \quad \forall u \in (0, 1)$$

pf: Chain rule:

$$\frac{d}{du} \mathbb{E} f(Z(u)) = \mathbb{E} \left\langle \frac{d}{du} Z(u), Df(Z(u)) \right\rangle$$

$$= \frac{1}{2\sqrt{u}} \underbrace{\mathbb{E} \langle X, Df(Z(u)) \rangle}_{(1)} - \frac{1}{2\sqrt{1-u}} \underbrace{\mathbb{E} \langle Y, Df(Z(u)) \rangle}_{(2)}$$

$$\begin{aligned}
 ① &= \sum_{i=1}^n E X_i \cdot \frac{\partial f}{\partial x_i} (\sqrt{u} X + \sqrt{1-u} Y) \\
 &= E_y \sum_{i=1}^n E_x X_i \cdot \frac{\partial f}{\partial x_i} (\sqrt{u} X + \sqrt{1-u} Y) \\
 &= E_y \sum_{i=1}^n \left( \sum_{j=1}^n \sum_{i,j}^x E_x \sqrt{u} \cdot \frac{\partial^2 f}{\partial x_i \partial x_j} (Z(u)) \right) \\
 &= \sqrt{u} E \left\langle \sum_i^x, D^2 f(Z(u)) \right\rangle
 \end{aligned}$$

Similar computation gives

$$② = \sqrt{1-u} E \left\langle \sum_j^y, D^2 f(Z(u)) \right\rangle$$

Lemma: Let  $X \sim N(0, \Sigma^X)$ ,  $Y \sim N(0, \Sigma^Y)$   
be independent. Assume

$$E X_i^2 = E Y_i^2, \quad E (X_i - X_j)^2 \leq E (Y_i - Y_j)^2 \quad \forall i, j$$

Consider  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  with  $[D^2 f]_{ij} \geq 0 \quad \forall i, j$

Then  $E f(X) \geq E f(Y)$

pt: Notice

$$\sum_{ii}^X = \sum_{ii}^Y \text{ and } \sum_{ij}^X \geq \sum_{ij}^Y$$

We can assume  $X$  and  $Y$  are independent. Apply the previous lemma.  $\square$

Finally, we are ready to prove Slepian's inequality.

Thm (Slepian's inequality)

Let  $\{X_t\}_{t \in T}$ ,  $\{Y_t\}_{t \in T}$  be mean-zero Gaussian processes. Assume  $\forall t_1, t_2 \in T$ :

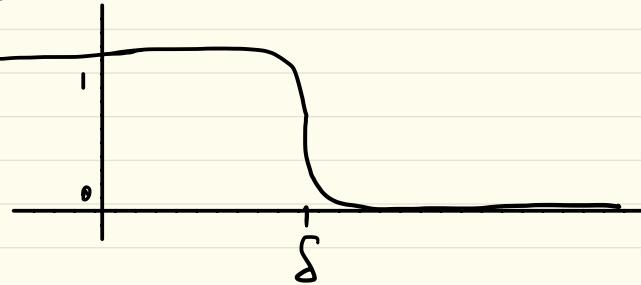
$$E X_t^2 = E Y_t^2 \text{ and } E(X_t - X_s)^2 \leq E(Y_t - Y_s)^2$$

Then

$$P\left[\sup_{t \in T} X_t \geq s\right] \leq P\left[\sup_{t \in T} Y_t \geq s\right] \wedge s$$

Therefore  $E \sup_{t \in T} X_t \leq E \sup_{t \in T} Y_t$ .

p.f.: Let  $h: \mathbb{R} \rightarrow [0,1]$  be the function



We can ensure  $h$  is a  $C^2$ -smooth approximation to  $\mathbb{1}_{(-\infty, s)}$ .

Define  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  by  $f(x) = h(x_1) \dots h(x_n)$   
[Note  $f$  approximates  $\mathbb{1}_{\{\max_i x_i \leq s\}}$ ]

Note  $\frac{\partial^2 f}{\partial x_i \partial x_j} = \mathbb{1}_0'(\mathbb{x}_i) \cdot \mathbb{1}_0'(\mathbb{x}_j) \prod_{k \neq i, j} h''(x_k) \geq 0$

So

$$Ef(x) \geq Ef(y)$$

Pass to the limit [e.g. dominated convergence]

□

Then (Sudakov-Fernique)

Let  $\{X_t\}_{t \in T}$ ,  $\{Y_t\}_{t \in T}$  be two mean-zero Gaussian processes. Assume

$$\mathbb{E}(X_t - X_s)^2 \leq \mathbb{E}(Y_t - Y_s)^2 \quad \forall s, t \in T.$$

Then  $\mathbb{E} \sup_{t \in T} X_t \leq \mathbb{E} \sup_{t \in T} Y_t$

p.f.: As before assume  $X \sim N(0, \Sigma^X)$   
 $Y \sim N(0, \Sigma^Y)$

Fix  $\beta > 0$  and define

$$f(x) = \frac{1}{\beta} \log \sum_{i=1}^n e^{\beta x_i}$$

You can check

$$\max_{i=1, \dots, n} x_i \leq f(x) \leq \max_{i=1, \dots, n} x_i + \frac{\log n}{\beta}$$

Apply Gaussian interpolation to deduce

$$\frac{d}{du} \mathbb{E} f(Z(u)) \leq 0 \quad (\text{HW})$$

◻

Cor: (Gaussian Contraction Inequality)  
 Fix a set  $T \subset \mathbb{R}^n$  and let  $g_1, \dots, g_n \stackrel{iid}{\sim} N(0, 1)$   
 Let  $\phi_i: \mathbb{R} \rightarrow \mathbb{R}$  be 1-Lipschitz.

Then

$$\mathbb{E} \sup_{t \in T} \sum_{i=1}^n g_i \phi_i(t_i) \leq \mathbb{E} \sup_{t \in T} \sum_{i=1}^n g_i t_i$$

Pf: Apply Sudakov-Fernique.  $\square$

We will now use Sudakov-Fernique to prove a useful lower bound on

$$\sup_{t \in T} X_t$$

$$\text{Recall: } d(t, s) = \sqrt{\mathbb{E}(X_t - X_s)^2}$$

defines a pseudo-metric on  $T$ .

Thm: (Sudakov's minoration inequality)  
 Let  $\{X_t\}_{t \in T}$  be a mean-zero Gaussian process. Then

$$\mathbb{E} \sup_{t \in T} X_t \geq C \sqrt{\log_2 N(T, d, \epsilon)} \quad \forall \epsilon > 0,$$

where  $N(T, d, \epsilon)$  denotes the covering number of  $T$  in  $d(\cdot, \cdot)$  metric.

pf: Suppose  $N := N(T, d, \epsilon)$  is finite.

You'll consider the infinite case for HW.

Let  $M$  be a maximal  $\epsilon$ -packing.  
 Then we know  $|M| \geq N$ . If

suffices to show

$$\mathbb{E} \sup_{t \in M} X_t \geq C \sqrt{\log N}$$

Define  $\{y_t\}_{t \in M}$  by  $y_t = \frac{\epsilon}{\sqrt{2}} g_t$   
 where  $g_t \sim N(0, 1)$ , independent.

Fix  $t, s \in M$ . Then

$$\mathbb{E}(X_t - X_s)^2 = d(t, s)^2 \geq \varepsilon^2$$

while

$$\mathbb{E}(Y_t - Y_s)^2 = \frac{\varepsilon^2}{2} \mathbb{E}(g_t - g_s)^2 = \varepsilon^2$$

Therefore Sudakov-Fernique

$$\Rightarrow \mathbb{E} \sup_{t \in M} X_t \geq \mathbb{E} \sup_{t \in M} Y_t \quad \text{HUI}$$

$$= \sqrt{\sum_{t \in M} \mathbb{E} \max_{t \in M} g_t} \geq C \sqrt{\log M}$$

We next study a particular type of a Gaussian process, which is an analogue of the Rademacher complexity.

Defn: The Gaussian width of a set

$$T \subseteq \mathbb{R}^n \text{ is } \omega(T) := \mathbb{E} \sup_{x \in T} \langle g, x \rangle$$

where  $g \sim N(0, I)$ .

## Properties:

(1)  $\omega(T) < \infty$ , iff  $T$  is bounded.

(2) Affine Invariance:

$$\omega(UT + y) = \omega(U)$$

for any orthogonal  $U \in \mathbb{R}^{n \times n}$  and  $y \in \mathbb{R}^n$ .

(3)  $\omega(\text{conv}(T)) = \omega(T)$

(4)  $\omega(T+S) = \omega(T) + \omega(S)$

$$\omega(aT) = |a|\omega(T) \quad \forall a \in \mathbb{R}.$$

(5)  $\omega(T) = \frac{1}{2} \omega(T-T) = \frac{1}{2} \sum_{x,y \in T} \langle g, x-y \rangle$

(6)  $\frac{1}{\sqrt{2\pi}} \text{diam}(T) \leq \omega(T) \leq \frac{\sqrt{n}}{2} \text{diam}(T)$

(7) For any 1-Lipschitz  $g: \mathbb{R} \rightarrow \mathbb{R}$  it holds  
 $\omega(g_0 T) \leq \omega(T)$

pf.: (1)-(4) follow from defn of support functions. To see (5) observe

$$\begin{aligned} \omega(T) &= \frac{1}{2} (\omega(T) + \omega(T)) = \frac{1}{2} (\omega(T) + \omega(-T)) \\ &\stackrel{(4)}{=} \frac{1}{2} \omega(T-T). \end{aligned}$$

To see (6), fix  $x, y \in T$ . Then

$$\begin{aligned}\omega(T) &\geq \frac{1}{2} \mathbb{E} \max(\langle x-y, g \rangle, \langle y-x, g \rangle) \\ &= \frac{1}{2} \mathbb{E} |\langle x-y, g \rangle| = \frac{1}{2} \sqrt{\frac{2}{\pi}} \|x-y\|\end{aligned}$$

Take  $\sup_{x, y \in T}$ .

Next

$$\omega(T) = \frac{1}{2} \mathbb{E}_g \sup_{x, y \in T} \langle g, x-y \rangle$$

$$\leq \frac{1}{2} \sup_{x, y \in T} \|x-y\|_2 \cdot \mathbb{E}_g \|g\| \leq \frac{\sqrt{n}}{2} \text{diam}(T)$$

(7) is the contraction inequality we already proved.

Lemma: For any  $A \in \mathbb{R}^{m \times n}$ , it holds

$$\omega(AT) \leq \|A\|_2 \omega(T)$$

[Hint: singular value decomposition + contraction]

Thm: (Gaussian vs Rademacher Complexity)  
 For any set  $T \subseteq \mathbb{R}^d$ , it holds:

$$\frac{\omega(T)}{2\sqrt{\log d}} \leq \mathcal{R}(T) \leq \sqrt{\frac{\pi}{2}} \omega(T)$$

$$\begin{aligned}
 \mathcal{R}(T) &= \mathbb{E}_\varepsilon \sup_{x \in T} \sum_{i=1}^d \varepsilon_i x_i \\
 &= \sqrt{\frac{\pi}{2}} \cdot \mathbb{E}_\varepsilon \sup_{x \in T} \sum_{i=1}^d \varepsilon_i \mathbb{E}[g_i] x_i \\
 &\leq \sqrt{\frac{\pi}{2}} \cdot \mathbb{E}_{\varepsilon, g} \sup_{x \in T} \sum_{i=1}^d \varepsilon_i |g_i| x_i \\
 &\stackrel{E_i |g_i| \leq g_i}{=} \sqrt{\frac{\pi}{2}} \cdot \mathbb{E}_g \sup_{x \in T} \sum_{i=1}^d g_i x_i \\
 &= \sqrt{\frac{\pi}{2}} \omega(T)
 \end{aligned}$$

Conversely

$$\begin{aligned}
 \omega(T) &= \mathbb{E}_g \sup_{x \in T} \sum_{i=1}^d g_i x_i \\
 &= \mathbb{E}_{g, \varepsilon} \sup_{x \in T} \sum_{i=1}^d \varepsilon_i g_i x_i
 \end{aligned}$$

$$= \mathbb{E}_g \mathcal{R}(g \circ T) \stackrel{\text{contraction}}{\leq} \mathbb{E}_{\substack{g \\ g_i \in T}} \max_{i=1 \dots d} \|g_i\|_1 \mathcal{R}(T)$$

$$\leq 2 \sqrt{\log^+ \mathcal{R}(T)} \quad \boxed{B}$$

There is a close cousin of Gaussian width called spherical width.

Defn: The spherical width of  $T \subseteq \mathbb{R}^d$

is  $s(T) = \mathbb{E}_{\Theta} \sup_{x \in T} \langle \Theta, x \rangle$

where  $\Theta \sim \text{Unif}(\mathbb{S}^{n-1})$

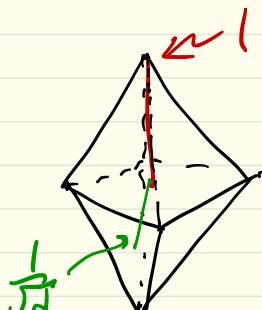
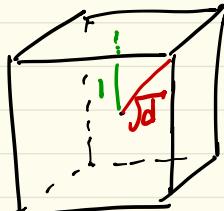
Lemma:  $(\gamma n - c) s(T) \leq w(T) \leq (\gamma n + c) s(T)$

$$\begin{aligned} w(T) &= \mathbb{E}_g \sup_{x \in T} \langle g, x \rangle = \mathbb{E}_g \sup_{x \in T} \langle \|g\|_2 \cdot \frac{g}{\|g\|_2}, x \rangle \\ &= \mathbb{E}_g \|g\|_2 \cdot s(T) \end{aligned}$$

Let's compare the sizes of  $B_2^d, B_1^d, B_\infty^d$ .

	$\log N(\cdot, \delta)$	$\omega(\cdot)$	$R_I$	$R_C$	# vertices
$B_2^d$	$d \log\left(\frac{1}{\delta}\right)$	$\sqrt{d}$	-	-	-
$B_1^d$	$\leq \frac{\log d}{\delta^2}$ ②	$\sqrt{\log d}$ ①	$\frac{1}{\sqrt{d}}$	1	$d$
$B_\infty^d$	$\leq \frac{d^2}{\delta^2}$ ④	$d$ ③	1	$\sqrt{d}$	$2^d$

- Pf:
- ①  $E \sup_{g \in B} \langle g, x \rangle = E_g \|g\|_\infty \approx \sqrt{\log d}$
  - ② Sudakov  $\Rightarrow \omega(B_1^d) \geq \max_s C s \sqrt{\log N(B_1^d, s)}$
  - ③  $\omega(B_\infty^d) = E_g \|g\|_1 = \sqrt{\frac{2}{\pi}} d$
  - ④ Sudakov



$$\omega(B_2^d) \approx \omega(\text{circumscribed ball}) \quad \omega(B_1^d) \approx \omega(\text{inscribed ball})$$

Algebraic Dimension is highly unstable

Defn: Define the stable dimension of  $T$  by

$$d(T) = \frac{h(T-T)}{\text{diam}(T)^2} \sim \frac{\omega(T)^2}{\text{diam}(T)^2}$$

where  $h(S) = \sqrt{\mathbb{E}_g \sup_{x \in S} \langle g, x \rangle^2}$

Lemma: (HW)  $2\omega(T) \leq h(T-T) \leq 2C\omega(T)$

Lemma:  $d(T) \leq \dim(\text{Span } T)$

pf: We can assume  $T \subseteq \text{Span}\{e_1, \dots, e_k\}$  where  $k = \dim(T)$ . Then

$$\begin{aligned} h(T-T)^2 &= \mathbb{E}_g \sup_{x \in T-T} \langle g, x \rangle^2 \leq \mathbb{E}_{g \sim N(0, I_n)} \|g\|^2 \\ &= k \cdot \text{diam}(T)^k \quad \blacksquare \end{aligned}$$

Exercise (HW): Let  $A \in \mathbb{R}^{m \times n}$ . Then

$$\|AB_2^d\| = \frac{\|A\|_F^2}{\|A\|_2^2}$$

- Stable rank  
of matrix
- changes gradually  
with perturbations

## Random Projections of Sets

Thm: Let  $T \subset \mathbb{R}^n$  and let  $P$  be a projection onto a uniformly random  $m$ -dimensional subspace.

Then w.p.  $1 - 2e^{-m}$ , we have

$$\text{diam}(PT) \leq C(S(T) + \sqrt{\frac{m}{n}} \text{diam}(T))$$

Lemma (HW): Let  $P$  be a projection in  $\mathbb{R}^n$  onto a uniformly random  $m$ -dimensional subspace. Let  $Q$  be an  $m \times n$ , orthogonal matrix drawn uniformly. Then

(a) For any  $x \in \mathbb{R}^n$ , it holds:

$\|Px\|_2$  and  $\|Qx\|_2$  have the same distribution.

(b) For any  $z$  with  $\|z\|=1$ , it holds  
 $Q^T z \sim \text{Unif}(S^{n-1})$

pf of them: Covering argument. Can replace  $P$  by  $Q$  by Lemma. Can assume  $\text{Diam}(T) \leq \frac{1}{m-1}$ .  
Approximation: Let  $N$  be  $\frac{1}{2}$ -net of  $S^{m-1}$ .  
 Then  $|N| \leq 5^m$ .

Then

$$\text{diam}(QT) \leq \sup_{X \in T-T} \|Q_X\|_2$$

$$= \sup_{X \in T-T} \sup_{Z \in \mathbb{S}^{n-1}} \langle Q_X, Z \rangle$$

$$\leq 2 \sup_{X \in T-T} \left[ \sup_{Z \in X} \langle X, Q^T Z \rangle \right]$$

$$\leq 2 \sup_{Z \in X} \left[ \sup_{X \in T-T} \langle X, Q^T Z \rangle \right]$$

Fix  $Z \in X$ . Then  $Q^T Z \sim \text{Unif}(\mathbb{S}^{n-1})$ .

Note:  $\int_Q \sup_{X \in T-T} \langle X, Q^T Z \rangle = S(T-T) = 2S(T)$

Concentration: The function  $\Theta \mapsto \sup_{X \in T-T} \langle \Theta, X \rangle$

is 1-Lipschitz on  $S^{n-1} \Rightarrow$  Concentration

$$P \left[ \sup_{X \in T-T} \langle X, Q^T Z \rangle \geq 2S(T) + t \right] \leq 2 \exp(-ct^2)$$

Union bound:

$$\mathbb{P}\left[\max_{Z \in \mathcal{X}} \sup_{x \in T-T} \langle Q_z^T, x \rangle \geq s(T) + t\right]$$

$$\leq 5^m \cdot 2 \exp(-c n t^2)$$

$$= 2 \exp(m \log(5) - c n t^2)$$

$$\text{set } t = e^{\sqrt{\frac{m}{n}}}$$

$$\leq e^{-m}.$$



Exercise (HW) For all bounded sets:

$$\mathbb{E} \text{diam}(PT) \geq C(s(T) + \sqrt{\frac{m}{n}} \text{diam}(T))$$

The phase transition:

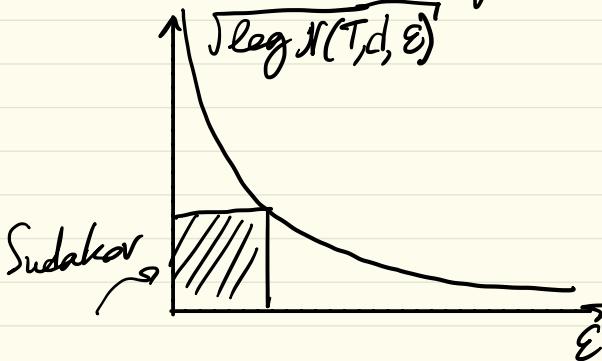
$$\text{diam}(PT) \leq C \begin{cases} \sqrt{\frac{m}{n}} \text{diam}(T), & \text{if } m \geq d(T) \\ s(T) & \text{, if } m \leq d(T) \end{cases}$$

Chaining

We now prove the following upper bound

$$\mathbb{E} \sup_{t \in T} X_t \leq CK \int_0^\infty \sqrt{\log N(T, d, \epsilon)} d\epsilon$$

on sub Gaussian processes.



Defn: Consider a random process  $\{X_t\}_{t \in T}$  on a metric space  $(T, d)$  if  $\exists K \geq 0$

s.t.

$$\|X_t - X_s\|_{\ell_2} \leq K d(t, s) \quad \forall t, s$$

**Sub Gaussian**

**norm**

Ex: If  $X_t$  is a <sup>canonical</sup> Gaussian process and

$$d(t, s) = \|t - s\|_2, \text{ then } K=1$$

Thm: (Discrete version)

Let  $\{X_t\}$  be a mean-zero subGaussian process on  $(T, d)$ . Then

$$\mathbb{E} \sup_{t \in T} X_t \leq CK \sum_{k \in \mathbb{Z}} 2^{-k} \sqrt{\log N(T, d, 2^{-k})}$$

[Idea: Let  $\Pi(\cdot)$  be projection on  $\epsilon$ -net repeat]

$$\Rightarrow \mathbb{E} \sup_{t \in T} X_t \leq \left[ \mathbb{E} \sup_{t \in T} X_{\Pi(t)} + \mathbb{E} \sup_{t \in T} [X_t - X_{\Pi_t}] \right]$$

p.f.: WLOG assume  $K=1$ ,  $T$  finite.

Set  $\epsilon_k = 2^{-k}$ ,  $k \in \mathbb{Z}$

and let  $T_k$  be an  $\epsilon_k$ -net of  $T$ .

Since  $T$  is finite  $\exists_{K_0, K_\infty}$  s.t.

$$T_{K_0} = \{t_0\} \text{ for some } t_0 \in T$$

$$T_{K_\infty} = T$$

For  $t \in T$ , let  $\Pi_k(t) \in T_k$  satisfy

$$d(t, \Pi_k(t)) \leq \epsilon_k.$$

Since  $\mathbb{E} \sum_{t_0} X_t = 0$ , we can write

$$\mathbb{E} \sup_{t \in T} X_t = \mathbb{E} \sup_{t \in T} [X_t - X_{t_0}]$$

Write

$$X_t - X_{t_0} = \sum_{k=k_0+1}^{K_0} (X_{T_k(t)} - X_{T_{k-1}(t)})$$

So

$$\mathbb{E} \sup_{t \in T} (X_t - X_{t_0}) \leq \sum_{k=k_0+1}^{K_0} \mathbb{E} \sup_{t \in T} (X_{T_k(t)} - X_{T_{k-1}(t)})$$

*Supremum is over*  
 $|T_k| \cdot |T_{k-1}| \leq |T_K|^2$

For fixed  $t$ , have

$$\begin{aligned} \|X_{T_k(t)} - X_{T_{k-1}(t)}\|_{\psi_2} &\leq d(T_k(t), T_{k-1}(t)) \\ &\leq d(T_k(t), t) + d(t, T_{k-1}(t)) \\ &\leq \varepsilon_k + \varepsilon_{k-1} \\ &\leq 2\varepsilon_{k-1} \end{aligned}$$

$$\text{So } \mathbb{E} \sup_{t \in T} (X_{\pi_n(t)} - X_{\pi_{n-1}(t)}) = C \varepsilon_{n-1} \sqrt{\log T \ln 1/\delta}$$

Exercise

2.5.10.  
(HW)

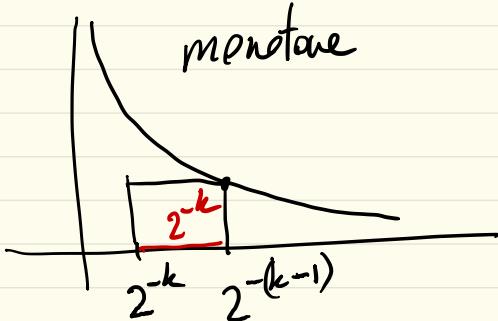
$$\text{So } \mathbb{E} \sup_{t \in T} |X_t - X_0| \leq C_1 \sum_{k=k_0}^{k_0-1} 2^{-k} \sqrt{\log N(T, d, 2^{-k})} \quad \square$$

Thm: (Dudley's Integral)

Let  $\{X_t\}_{t \in T}$  be a mean zero random process on  $(T, d)$  with subGaussian increments. Then

$$\mathbb{E} \sup_{t \in T} X_t \leq CK \int \sqrt{\log N(T, d, \varepsilon)} d\varepsilon$$

p.f.:



Thm (HW)

Let  $\{X_t\}_{t \in T}$  be a subGaussian process on  $(T, d)$ , possibly not centered. Then w.p.  $1 - 2e^{-\delta^2}$  have:

$$\sup_{t, s \in T} |X_t - X_s| \leq CK \left( \int_0^{\sqrt{\log N(T, d, \epsilon)}} de + \delta \cdot \text{diam}(T) \right)$$

Thm: (Two-sided Sudakov)

Let  $\{X_t\}_{t \geq 0}$  be a canonical Gaussian process on a set  $T$ . Define

$$\varphi(T) = \sup_{\epsilon \geq 0} \epsilon \sqrt{\log N(T, \epsilon)}$$

Then

$$c_2 \varphi(T) \leq W(T) \leq c_1 (\log n) \varphi(T)$$

[Read in Vershynin]

# Improved Uniform Laws

Let  $F$  be a class of indicator functions with VC-dimension  $V$ . We showed

$$R_n(F) \leq 2 \sqrt{\frac{V \log(n+1)}{n}}$$

Let's get a better bound on

$$R_n(F) = E_x E_{\varepsilon} \sup_{f \in F} \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(x_i)$$

Define  $Z_f := \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(x_i)$  and consider the process  $\{Z_f\}_{f \in F}$

Let's compute

$$\|Z_f - Z_g\|_2 = \frac{1}{\sqrt{n}} \left\| \sum_{i=1}^n \varepsilon_i (f(x_i) - g(x_i)) \right\|_2$$

$$\leq \frac{1}{\sqrt{n}} \sqrt{\sum_{i=1}^n (f(x_i) - g(x_i))^2} \leq \frac{1}{\sqrt{n}} \|f - g\|_2$$

So Audley

$$\Rightarrow \mathbb{E} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_i \varepsilon_i f(x_i) \right| \leq \frac{C}{\sqrt{n}} \int_0^2 \sqrt{\log N(\mathcal{F}, \delta, \| \cdot \|_2)} d\delta$$

$$\text{Thm: } N(\mathcal{F}, \delta, \| \cdot \|_2) \leq \left( \frac{2}{\delta} \right)^{C\sigma}$$

[Vershynin Thm 8.3.18]

So continuing

$$\leq \frac{C}{\sqrt{n}} \int_0^2 \sqrt{V \ln\left(\frac{2}{\delta}\right)} d\delta \leq C \sqrt{\frac{V}{n}}$$

$$\boxed{R_n(F) \leq C \sqrt{\frac{V}{n}}}$$

Thm: (Improved Glivenko-Cantelli)

Let  $X_1, \dots, X_n$  be iid random variables with CDF  $F$ . Define  $F_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{X_i \leq t\}}$

$$\text{Then } \mathbb{E} \|F_n - F\|_\infty \leq \frac{C}{\sqrt{n}}$$

Pf: Let  $\mathcal{F} = \{ \mathbb{1}_{(-\infty, t]} : t \in \mathbb{R} \}$

$$\Rightarrow \text{VC}(\mathcal{F}) = 1. \quad \blacksquare$$

Another Example:

Define

$$\mathcal{F} := \left\{ f: [0,1] \rightarrow [0,1] \mid \|f\|_{Lip} \leq L \right\}$$

Thm: Let  $X_1, \dots, X_n$  be iid random variables in  $[0,1]$ . Then

$$\mathbb{E} \sup_{f \in \mathcal{F}} \left( \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E}f(X) \right) \leq \frac{CL}{\sqrt{n}}$$

pf: WLOG assume  $L=1$ .

Again, let's upper bound

$$R_n(\mathcal{F}) = \mathbb{E}_X \underbrace{\mathbb{E}_{\varepsilon} \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \varepsilon_i \cdot f(X_i)}_{}$$

$$\text{Define } Z_f = \frac{1}{n} \sum_{i=1}^n \varepsilon_i \cdot f(X_i)$$

Then

$$\|Z_f - Z_g\|_{\ell_2} = \frac{1}{n} \left\| \sum_{i=1}^n \varepsilon_i (f(X_i) - g(X_i)) \right\|_{\ell_2}$$

$$\leq \frac{1}{n} \sqrt{\sum_{i=1}^n (f(X_i) - g(X_i))^2} \leq \frac{1}{\sqrt{n}} \|f - g\|_{\ell_\infty}$$

Dudley  $\Rightarrow$  diameter

$$\mathbb{E} \sup_{f \in \mathcal{F}} |Z_f| \leq \frac{C}{\sqrt{n}} \int_0^{\text{diameter}} \sqrt{\log N(\mathcal{F}, \|\cdot\|_\infty, \varepsilon)} d\varepsilon.$$

Lemma(HH):  $N(\mathcal{F}, \|\cdot\|_\infty, \varepsilon) \leq \left(\frac{2}{\varepsilon}\right)^{\frac{2}{\varepsilon}}$

$$\Rightarrow \mathbb{E} \sup_{f \in \mathcal{F}} |Z_f| \leq \frac{C}{\sqrt{n}} \int_0^1 \frac{1}{\varepsilon} \log\left(\frac{1}{\varepsilon}\right) d\varepsilon \\ \leq \frac{\hat{C}}{\sqrt{n}}$$

□

So far we know that

$$\sup_{T \in \mathcal{T}} \mathbb{E} [\log N(\mathcal{T}, \|\cdot\|_\infty)] \leq \mathbb{E} \sup_{t \in T} X_t \leq K \int_0^1 \sqrt{\log N(T, \varepsilon, d)} d\varepsilon$$

Gaussian

SubGaussian

Question: Is there a tighter bound that separates ① the subGaussian parameter?  
 ② size of  $T$ .  
 Yes!

Recall the key bound from chaining

$$\textcircled{X} \mathbb{E} \sup_{t \in T} X_t \leq C \sum_{k=k_0}^{\infty} \epsilon_{k-1} \sqrt{\log |T_k|}$$

where  $T_k$  are  $\epsilon_k$ -nets of  $T$ .

→ We chose  $\epsilon_n$  followed by  $T_n$ .

Let's switch the order.

Fix subset  $T_1, T_2, \dots, T_k, \dots \subset T$   
such that

$$|T_0|=1 \text{ and } |T_k| \leq 2^{2^k} \text{ for } k \geq 1$$

The sequence  $(T_k)_{k=0}^{\infty}$  is called admissible

Define  $\epsilon_k = \sup_{t \in T} d(t, T_k)$

Then  $T_k$  is an  $\epsilon_k$ -net of  $T$ .

$\textcircled{X}$  becomes

$$\mathbb{E} \sup_{t \in T} X_t \leq C \sum_{k=0}^{\infty} 2^{k\epsilon_k} \sup_{t \in T} d(t, T_k)$$

Defn (Talagrand's  $\gamma_2$  functional)

Let  $(T, d)$  be a metric space. Define

$$\gamma_2(T, d) = \inf_{(T_n)} \sup_{t \in T} \sum_{k=0}^{\infty} 2^{k/2} d(t, T_n)$$

where  $\inf_{(T_n)}$  is taken over all admissible sequences.

Thm (Generic chaining bound)

Let  $\{X_t\}_{t \in T}$  be a mean zero subGaussian process satisfying

$$\|X_t - X_s\|_{\psi_2} \leq k d(t, s).$$

Then

$$\mathbb{E} \sup_{t \in T} X_t \leq C K \gamma_2(T, d)$$

pf: WLOG assume  $k=1$ .

We do the chaining as before

$$t_0 \rightarrow \pi_{t_0}(t) \rightarrow \pi_{\pi_{t_0}(t)}(t) \rightarrow \dots \rightarrow \pi_{t_N}(t) = t$$

$$\Rightarrow X_t - X_0 = \sum_{k=1}^N (X_{\pi_k(t)} - X_{\pi_{k-1}(t)})$$

Fix  $k, t$ . We know

$$\|X_{\pi_k(t)} - X_{\pi_{k-1}(t)}\|_{\ell_2} \leq d(\pi_k(t), \pi_{k-1}(t))$$

$$\Rightarrow \Pr[|X_{\pi_k(t)} - X_{\pi_{k-1}(t)}| \leq cu^2 2^{k/2} d(\pi_k(t), \pi_{k-1}(t))] \\ \geq 1 - 2 \exp(-8u^2 2^k)$$

Take union bound over

$$|\mathcal{T}_k| |\mathcal{T}_{k-1}| \leq |\mathcal{T}_k|^2 = 2^{2^{k+1}} \quad k \in \mathbb{N}$$

So get uniform bound w.p.

$$1 - \sum_{k=1}^{\infty} 2^{2^k} \cdot 2 \exp(-8u^2 2^k) \stackrel{\text{check}}{\geq} 1 - 2 \exp(-cu^2)$$

for  $u > \hat{c}$ .

In this event, get

$$\begin{aligned} |X_t - X_0| &\leq C_u \sum_{k=1}^{\infty} 2^{k/2} d(\pi_k(t), \pi_{k-1}(t)) \\ &\leq C_u \underbrace{\sum_{k=1}^{\infty} 2^{k/2} d(t, T_k)}_{\delta_2(T, d)} \end{aligned}$$

$$\Rightarrow \sup_{t \in T} |X_t - X_0| \leq C u \delta_2(T, d)$$

Now integrate the tail  $\square$

Remarkably  $\delta_2(T, d)$  is also a lower bound

Thm (Talagrand) Let  $\{X_t\}_{t \in T}$  be a Gaussian Process on  $T$ . Define  $d(t, s) = \|X_t - X_s\|_2$ . Then

$$C \delta_2(T, d) \leq \mathbb{E} \sup_{t \in T} X_t \leq C \delta_2(T, d)$$

For us the most important consequence is the following generalization of

Sudakov - Fernique comparison inequality.

Thm: (Talagrand Comparison Inequality)

Let  $\{X_t\}_{t \in T}$  be a mean-zero random process on  $T$  and let  $\{Y_t\}_{t \in T}$  be a Gaussian process. Assume

$$\|X_t - X_s\|_{\ell_2} \leq K \|Y_t - Y_s\|_2$$

Then

$$\mathbb{E} \sup_{t \in T} X_t \leq CK \mathbb{E} \sup_{t \in T} Y_t.$$

p.f.: Define  $d(t, s) = \|Y_t - Y_s\|_2$ . Then

$$\mathbb{E} \sup_{t \in T} X_t \leq CK \gamma_2(T, d) \leq \hat{C} K \mathbb{E} \sup_{t \in T} Y_t.$$

upper  
bound  
in Talagrand

lower  
bound  
in Talagrand

Cor: Let  $\{X_t\}_{t \in T}$  be a mean-zero random process on  $T \subseteq \mathbb{R}^n$ . Assume

$$\|X_t - X_s\|_{\ell_2} \leq K \|t-s\|_2$$

Then  $\mathbb{E} \sup_{t \in T} X_t \leq CK w(T)$

Pf: Define  $y_t = \langle g, t \rangle$  where  $g \in N(0, I)$ . Apply previous thm.  $\square$

Thm: (SubGaussian Charvet Inequality)

Let  $A \in \mathbb{R}^{m \times n}$  be a random matrix with independent mean zero  $K$ -subGaussian parameters. Let  $T \subseteq \mathbb{R}^n$ ,  $S \subseteq \mathbb{R}^m$  be arbitrary bounded sets. Then

$$\mathbb{E} \sup_{x \in T, y \in S} \langle Ax, y \rangle \leq CK(w(T) \text{rad}(S) + \text{rad}(T) w(S))$$

pf: WLOG assume  $K=1$ . Define  
 $X_{u,v} := \langle Au, v \rangle$  for  $u \in T, v \in S$ .

Compute

$$\begin{aligned}
\|X_{uv} - X_{wz}\|_{\ell_2} &= \left\| \sum_{i,j} A_{ij} (u_i v_j - w_i z_j) \right\|_{\ell_2} \\
&\leq \left( \sum_{i,j} |u_i v_j - w_i z_j|^2 \right)^{1/2} \\
&= \|uv^T - wz^T\|_F \\
&= \|u(v-z)^T + (u-w)z^T\|_F \\
&\leq \|u(v-z)\|_F + \|(u-w)z^T\|_F \\
&= \|v-z\|_2 \|u\|_2 + \|u-w\|_2 \|z\|_2 \\
&\leq \|v-z\|_2 \cdot \text{rad}(T) + \|u-w\|_2 \cdot \text{rad}(S)
\end{aligned}$$

Define  $y_{uv} := \langle g, u \rangle \text{rad}(T) + \langle h, v \rangle \text{rad } S$

where  $g \sim N(0, I)$ ,  $h \sim N(0, I)$

$$\begin{aligned}
\|y_{uv} - y_{wz}\|^2 &= \|u-w\|_2^2 \text{rad}(T)^2 \\
&\quad + \|v-z\|_2^2 \text{rad}(S)^2
\end{aligned}$$

So

$$\|X_{uv} - X_{wz}\|_{\ell_2} \leq c \|y_{uv} - y_{wz}\|_2$$

Talagrand

$$\Rightarrow \mathbb{E} \sup_{u \in T, v \in S} X_{uv} \leq c \mathbb{E} \sup_{u \in T, v \in S} y_{uv}$$

$$= \mathbb{E} \sup_{u \in T} \langle g_u, u \rangle \text{rad}(S)$$

$$+ \mathbb{E} \sup_{v \in S} \langle h_v, v \rangle \text{rad}(T)$$

$$= \omega(T) \text{rad}(S) + \omega(S) \text{rad}(T).$$

□