

01/18. 36709 Lec 1

not cover everything

pick a paper. connection/report the paper
how to prove. how cool idea
scribing. everyone

Non-asymptotic high-dimensional statistics

1. High-dimensional (linear) classification

2. Estimating high-dim matrices

3. Non-parametric regression

1. High-dim classifier.

Fisher's Linear Discriminant Analysis

$$y=1, X_1 \dots X_{n_1} \sim N(\mu_1, \Sigma)$$

$$y=2, X_{n_1+1} \dots X_{n_1+n_2} \sim N(\mu_2, \Sigma)$$

If know μ_1, μ_2, Σ Assume balanced classes

Bayes-optimal classifier

$$X, (x - \mu_1)^T \Sigma^{-1} (x - \mu_1) \leq (x - \mu_2)^T \Sigma^{-1} (x - \mu_2) \text{ then Label 1, otherwise Label 2}$$

$$\text{Replace } \mu_1 \text{ by } \hat{\mu}_1 = \frac{1}{n_1} \sum_{i=1}^{n_1} X_i, \mu_2 \text{ by } \hat{\mu}_2, \Sigma \text{ by } \hat{\Sigma} = \frac{1}{2} \left[\sum_{i=1}^{n_1} (x_i - \hat{\mu}_1)(x_i - \hat{\mu}_1)^T + \sum_{i=n_1+1}^{n_1+n_2} (x_i - \hat{\mu}_2)(x_i - \hat{\mu}_2)^T \right]$$

$$\rightarrow x \rightarrow \text{pred}_y, \max p(y|x=x)$$

$$\text{Classifier: } \hat{\psi}(x) = \underbrace{\langle \mu_1 - \mu_2, \hat{\Sigma}^{-1} (x - \frac{\mu_1 + \mu_2}{2}) \rangle}_{\hat{\psi}_N} \text{, Label 1 if } \hat{\psi}(x) > 0 \quad \text{Label 2 if } \hat{\psi}(x) \leq 0$$

$$\text{Err}(\hat{\psi}) = \frac{1}{2} P_{X \sim N(\mu, \Sigma)} [\hat{\psi}(x) < 0] + \frac{1}{2} P_{X \sim N(\mu_2, \Sigma)} [\hat{\psi}(x) > 0]$$

$$\text{err} = P(\hat{\psi}(x) \neq y) = P(\hat{\psi}(x) \neq 1 | y=1) P(y=1) + P(\hat{\psi}(x) \neq 2 | y=2) P(y=2)$$

Story: *classical asymptotics $n \rightarrow \infty$ d fixed

\rightarrow consistent estimate $\hat{\mu}_1, \hat{\Sigma}, \hat{\mu}_2$

$\rightarrow \text{Err}(\hat{\psi}) \xrightarrow{P} \text{Err}(\psi)$ $\xrightarrow{\text{Bayes-opt}}$

$$\text{Err}(\psi) = \Phi\left(-\frac{\|\mu_1 - \mu_2\|}{2}\right), \gamma = \|\mu_1 - \mu_2\|$$



* high-dim. asympt. Kolmogorov.

$$\hookrightarrow n \rightarrow +\infty, \left[\frac{n_1}{d}, \frac{n_2}{d} \right] \rightarrow \Delta, \frac{\log d}{n} \rightarrow 0, \frac{d}{n} \rightarrow c, \frac{d^2}{n} \rightarrow c$$

proportional. $\Delta \approx \exp\left(\frac{n}{\log d}\right)$

$\sum_{i=1}^d$ Kolmogorov. $\text{Err}(\hat{\psi}) \xrightarrow{P} \Phi\left(-\frac{\gamma^2}{2\sqrt{\gamma^2 + 2\Delta}}\right), \gamma = |\mu_1 - \mu_2|/\sqrt{\frac{n_1}{d} + \frac{n_2}{d}}$

Exercise. $\hat{\mu}_1 = \mu_1 + \bar{z}_1, \bar{z}_1 \sim N(0, \frac{1}{n_1}), \hat{\mu}_2 = \mu_2 + \bar{z}_2, \bar{z}_2 \sim N(0, \frac{1}{n_2})$

$$Q(N) = \langle \mu_1 - \mu_2, \bar{z}_1 - \bar{z}_2, \gamma - \frac{\mu_1 + \mu_2}{2} - \frac{\bar{z}_1 + \bar{z}_2}{2} \rangle, \text{ prove Kolmogorov. } \Phi\left(\frac{\gamma^2}{2\sqrt{\gamma^2 + 2\Delta}}\right)$$

2. High-dim. matrices

$$X_1, \dots, X_n \sim N(0, \Sigma), \text{ want to est. } \Sigma. \quad \hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n X_i X_i^T$$

$$\|\hat{\Sigma} - \Sigma\|_{\text{op}} := \lambda_{\max}(\hat{\Sigma} - \Sigma) = \sup_{\|X\|=1} \sqrt{(\hat{\Sigma} - \Sigma)V} \quad (\text{largest eigenvalue})$$

operator norm

$$\frac{\gamma_1}{d} \rightarrow \Delta \text{ or } \frac{\gamma_1}{d} \rightarrow 0 \quad d \text{ fixed } \Rightarrow \begin{cases} \text{classically} & \|\hat{\Sigma} - \Sigma\|_{\text{op}} \rightarrow 0 \\ & d \text{-fixed. } n \rightarrow \infty \end{cases}$$

$d/n \rightarrow \Delta, \hat{\Sigma} - \Sigma \rightarrow \text{eigenvalues. don't all converge to 0.}$
Top eigenvector \vec{v} of $\hat{\Sigma}$ won't be "close" \vec{v} of Σ

Eigenvalues distribution: Maroko-Pastur law. random mat

↳ covariance: $f(x) = \frac{1}{\pi} \int_{-\sqrt{\delta}}^{\sqrt{\delta}} \frac{1}{x - t} dt$

$$x = \text{left} + \sqrt{\delta} \cdot \delta, \delta = (-\delta)^2$$

3. Non-parametric regression.

d can be fixed.

est. Lipschitz funct.

$$n^{-\frac{2}{2+d}}, d = 1, \sqrt{n}^{-\frac{1}{6}}, \text{ much slower than } \sqrt{n}$$

error in parameter

Bias-Variante regression $y_i = f(x_i) + \varepsilon_i$



average:

y-values

of bins

$(\frac{1}{d})^d$

exponential
curve of dim

"Bad things" in high-dim

structure/structural assumptions

↓ sparsity of vector. ↓ pattern in non-zero - low rank - sparse

↓ functions. smoothness. additive decompositions

$$f(x) = \sum_j f_{i,j}(x_j)$$

↳ only depends on 5 variables.

Lec 2 Chpt 5

Recap:

1. Contrasting (1) classical asymptotic: $\frac{d}{n} \rightarrow 0$

(2). high-dim asym., $\frac{d}{n} \rightarrow d$

(3) high-dim Mar asym $\frac{d}{n} \rightarrow +\infty$

2 Example.

(1) vector estimation (classification)

(2) matrix estimate (PCA/cov. matrices)

(3) function estimation. (L -Lips)

3. Structure.

(1) Sparsity

(2) low-rankness

(3) Smoothness, sparsity

Motivating example:

$$y = \theta^* + \varepsilon, \quad \Sigma \sim N(0, I_d), \text{ estimate } \hat{\theta}^* \text{ (goal)} \\ \hat{\theta}^* \in K, \text{ (may be sparse)} \quad \text{encode structure} \quad \xrightarrow{\text{convex set}} \\ \text{Maximize likelihood: } \hat{\theta} = \arg \max_{\theta \in K} \left\{ \frac{1}{2} \exp \left\{ -\frac{(y - \theta)^T (y - \theta)}{\Sigma} \right\} \right\} \\ \begin{aligned} &= \arg \min_{\theta \in K} \frac{1}{2} \|y - \theta\|^2 \\ &\quad \theta \in X \end{aligned}$$



Basic inequality

$$\frac{1}{2} \|y - \hat{\theta}\|^2 \leq \frac{1}{2} \|y - \theta^*\|^2, \quad \hat{\theta} \text{: best estimate}$$

$$\Delta = \hat{\theta} - \theta^*, \quad y = \theta^* + \varepsilon, \quad \frac{1}{2} \|\hat{\theta} - \theta^* - \varepsilon\|^2 \leq \frac{1}{2} \|\varepsilon\|^2 \Rightarrow \frac{1}{2} \|\Delta\|^2 \leq \langle \Delta, \varepsilon \rangle$$

→ If Δ was fixed, $\langle \Delta, \varepsilon \rangle \sim N(0, \|\Delta\|^2)$

→ Δ depends on ε ,

$$\langle \Delta, \varepsilon \rangle \leq \sup_{\Delta \in K, \theta^*} \langle \Delta, \varepsilon \rangle = \sup_{\Delta \in K} \langle \Delta, \varepsilon \rangle \quad \xrightarrow{\text{bound on this}} \text{study this} \quad (\text{upper bound})$$

→ One valid upper bd: $\langle \Delta, \varepsilon \rangle \leq \sup_{\Delta \in K, \theta^*} \langle \Delta, \varepsilon \rangle \quad \xrightarrow{\Delta \in E} \text{just an example}$

Gaussian complexity of $K(\Sigma)$

$$\ell_g(K) = \mathbb{E} \sup_{\theta \in K} \langle \theta, \varepsilon \rangle$$

Measure of size of K , useful for statistical applications

Examples. $K = B_2(R) = \{ \theta : \|\theta\|_2 \leq R \}$

$\ell_g(K) \leq \mathbb{E} \|\theta\|_2 \|\varepsilon\|_2$ by Cauchy-Schwarz

$$\leq R \mathbb{E} \|\theta\|_2 \sqrt{\mathbb{E} \|\varepsilon\|_2^2} \leq R \sqrt{d}$$

Metric entropy

$K, \ell \stackrel{\text{(metric)}}{=} \{ \theta \in \mathbb{R} : p(\theta, \theta') \geq 0, \text{ only if } \theta = \theta' \}$
 $\forall \theta, \theta' \in K, p(\theta, \theta') \leq p(\theta, \tilde{\theta}) + p(\tilde{\theta}, \theta')$
 $\exists \tilde{\theta} \in K, p(\theta, \theta') = p(\theta, \tilde{\theta}) + p(\tilde{\theta}, \theta')$

$K \subseteq \mathbb{R}^d$, metric $p(\theta, \theta') = \|\theta - \theta'\|$
Hamming metric on $\{0,1\}^d$, $p(\theta, \theta') = \frac{1}{d} \sum_{j=1}^d \mathbf{1}_{\{\theta_j \neq \theta'_j\}}$

Covering number:

$N(\delta; K, \ell)$: minimum # of radius δ balls needed to "cover" K
 $\exists ! \dots \theta^j, \dots \theta^N$, minimal N , covering numbers
 $\forall \theta \in K, \exists j, p(\theta, \theta^j) \leq \delta$

Metric entropy: $(\log N(\delta; K, \ell), (\delta^{-d}, N))$

Example: $K = [-1, 1] \subseteq \mathbb{R}$, $p(\theta, \theta') = |\theta - \theta'|$

Upper bounds on cover: one valid cover

$$\theta^j = -1 + 2(j-1)\delta, N = \left\lceil \frac{1}{\delta} \right\rceil + 1$$

$$N(\delta; [-1, 1], \delta) \leq \left\lfloor \frac{1}{\delta} \right\rfloor + 1 \leq \frac{1}{\delta} + 1 \Rightarrow \log(N)$$

$$N(\delta; [-1, 1]^d, \delta) \leq (\frac{1}{\delta} + 1)^d \Rightarrow d \log(N)$$

Packing numbers

δ -packing of K , $\{\theta^1, \dots, \theta^M\} \subseteq K$

$$p(\theta^i, \theta^j) > \delta$$

$\rightarrow M(\delta; K, \ell)$ maximal δ -packing of K

Lemma: $M(2\delta; K, \ell) \leq N(\delta; K, \ell) \leq M(\delta; K, \ell)$

$\theta \in K : \theta^j \in \text{2st. of } \theta^i$
 $\theta^i \in \text{2st. of } \theta^j$
add to packing
(maximal)
if $\dots \theta^M$

$\theta^i \in \text{2st. of } \theta^j$
 $\theta^j \in \text{2st. of } \theta^i$
claim also
a δ -cover

Take any $v \in \text{covering}$

$$\theta^* = \arg \inf \{p(v, \theta^*)\}$$

$\forall \theta \in K, p(v, \theta) \geq p(v, \theta^*) - p(\theta^*, \theta)$ If $p(v, \theta^*) \geq 2\delta$, then $p(v, \theta) \geq \delta$
 $\geq 2\delta - \delta = \delta$

$N(\delta; K, \delta)$ upper bds. constant cover

{ Lower bd. construct 2δ packing, many & }

Unit cube

$K = [-1, 1] \subseteq \mathbb{R}, \ell = \|\theta - \theta'\|$
 $\left\lfloor \frac{2}{\delta} \right\rfloor + 1 \leq N(\delta; K, \ell) \leq \left\lceil \frac{3}{\delta} \right\rceil + 1$
Valid 2δ packing

$\cup B(\theta^i, \delta) \subseteq K$
 $\# \theta^i \leq \frac{1}{\delta}$

Recap

1. Estimating a structured vector:

$$y = \theta^* + \varepsilon, \theta^* \in K$$

$$\hat{\theta} = \arg \min_{\theta \in K} \|y - \theta\|_2$$

Basic inequality: $\|\hat{\theta} - \theta^*\|_2 \leq 2\|\hat{\theta} - \theta^*\|_2 \leq 2 \sup_{\theta \in K} \langle \theta, \varepsilon \rangle$

$\|\hat{\theta} - \theta^*\|_2 \leq 2 E(\sup_{\theta \in K} \langle \theta, \varepsilon \rangle)$ Gaussian (subth) (K - θ^*) complexity

2. Packing & Covering:

$N(\delta; K, l)$ → min # of δ -balls in metric l that are needed to cover K
 $M(\delta, K, l)$ → max # of points that can be placed on K such that $l(\theta^i, \theta^j) \geq \delta$

Metric entropy = $\log N(\delta, K, l)$

$$M(2\delta; K, l) \leq N(\delta, K, l) \leq M(\delta, K, l)$$

Metric entropy of B , $\|\cdot\|$ using balls B' , $\|\cdot\|'$

→ Examples $\|\cdot\|$ & $\|\cdot\|'$ are $\|\cdot\|_2$

$$B \subseteq \bigcup_{i=1}^n B'(\theta^i, \delta) \{ \theta^1, \dots, \theta^n \} \rightarrow \text{covers } B$$

Lower bd covering number

$$N \geq \frac{\text{vol}(B)}{\text{vol}(\delta B')} = \frac{1}{\delta^d} \frac{\text{vol}(B)}{\text{vol}(B')}$$

Upper bd covering number:

$$N \leq M(\delta, K, \|\cdot\|')$$

$$\text{In } \text{vol}\left(\frac{1}{2}B'\right) \leq \text{vol}\left(B + \frac{\delta}{2}B'\right) \text{ Minkowski sum. C+C': } \{x+y : x \in C, y \in C'\}$$

$$\frac{1}{\delta^d} \frac{\text{vol}(B)}{\text{vol}(B')} \sim N \leq \frac{\text{vol}\left(B + \frac{\delta}{2}B'\right)}{\text{vol}\left(\frac{1}{2}B'\right)} = \left(\frac{2}{\delta}\right)^d \frac{\text{vol}(B + \frac{\delta}{2}B')}{\text{vol}(B')}$$

Examples. $\|\cdot\|$, same norm

$$B + \frac{\delta}{2}B' \leq 2B, N \leq \left(\frac{2}{\delta}\right)^d \frac{\text{vol}(2B)}{\text{vol}(B)} \leq \left(\frac{4}{\delta}\right)^d$$

Covers unit sphere with δ balls $(\frac{\delta}{\sqrt{d}})^d$, cubes of side-length δ to cover unit cubes $\{\|\cdot\|_\infty \leq 1\}$

office

F-mom.

T-approx.

Takeaway: $\log N(\delta; \|\cdot\|, \|\cdot\|) \leq d \log(\frac{L}{\delta})$

Lipschitz functions on $[0,1]$, $\mathcal{F}_L = \{f(x) = 0, |f(x) - f(x')| \leq L|x-x'|, \forall (x,x') \in [0,1]\}$
 $\|f-g\|_\infty = \sup_{x \in [0,1]} |f(x) - g(x)|$, if $\exists \dots \exists N$ such that $f \in \mathcal{F}_L, \exists j, \|f-f_j\|_\infty \leq \delta$

Takeaway $\log N(\delta, \mathcal{F}_L, \|\cdot\|_\infty) \asymp (\frac{L}{\delta})^d$ and dims $\asymp (\frac{L}{\delta})^d$

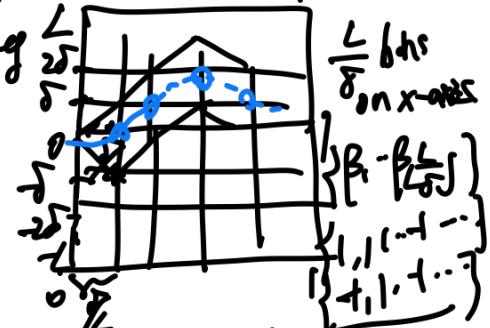
↳ Lipschitz \gg "unit balls", much much (use metric entropy for \mathcal{F}_L)

Cover bound: suffices to construct 2δ packing

each binary string $\rightarrow f_j$

$$\|f_j - f_k\| \leq 2\delta$$

$\hookrightarrow 2^{Ld}$ binary strings, $\log N \asymp \frac{Ld}{\delta}$



Covering argument:

\rightarrow at grid points $\{\frac{\delta}{2}, \frac{2\delta}{2}, \dots\}$, $\|f_j - f_k\| \leq \delta$

\rightarrow at non-grid points $|f_j(x_1) - f_k(x_1)| \leq |f_j(x_1) - f(x_1)| + |f(x_1) - f_k(x_1)|$
 (No: nearest gridpt) $\leq 3\delta$

$f_j \cdot \text{binary string}$

$$f_j = \sum_{i \in \{0,1\}^d} f_i \cdot i \cdot (\text{L-factor})$$

indexed by style
param $\theta \in \mathbb{R}$

Example: parametric class

$$f_\theta(x) = (-\exp(-\theta x), \theta \in [0,1])$$



$$\log N \leq d \log(\frac{L}{\delta}) \quad (\theta \in \mathbb{R}^d)$$

High-level: to cover \mathcal{F}_θ , just need to "cover" θ . intuition: change θ a little bit
 f_θ changes a little bit

$$P(f_\theta, f_{\theta'}) \leq P(\theta, \theta')$$

$\sup_{x \in [0,1]} |f_\theta(x) - f_{\theta'}(x)| \leq C(\theta - \theta')$
 \mathcal{F}_θ smooth function, θ older smooth. $\int_0^1 \theta^\alpha d\theta = \frac{1}{\alpha+1}$, $\alpha \in (0,1)$, $\delta = 1, \alpha = \theta \rightarrow$ Lipschitz

Smooth function, θ older smooth. $\int_0^1 \theta^\alpha d\theta = \frac{1}{\alpha+1}$, $\alpha \in (0,1)$, $\delta = 1, \alpha = \theta \rightarrow$ Lipschitz

$$|f^{(\alpha)}(x)| \leq L, |f_\theta^{(\alpha)}(x) - f_{\theta'}^{(\alpha)}(x)| \leq L|x - x'|^\alpha$$

$$\log N \leq (\frac{L}{\delta})^{\frac{d}{\alpha+1}}$$

$H \ll L$
 $|f_\theta(x) - f_{\theta'}(x)| \leq L$

↳ Cam's equation.

Rate of convergence.

Rate of convergence of $\|\cdot\|$

$$n \geq \log N(\varepsilon, \mathcal{F}_\theta, \|\cdot\|), \varepsilon: \text{rate of convergence}$$

$\varepsilon \downarrow \rightarrow n \geq \log N(\varepsilon, \mathcal{F}_\theta, \|\cdot\|) \asymp (\frac{L}{\varepsilon})^d$, typical non-param rate

$$n \geq \frac{d \log(\frac{L}{\varepsilon})}{d \log(\frac{L}{\delta})}, \varepsilon \asymp \sqrt{\frac{d \log(\frac{L}{\delta})}{d \log(\frac{L}{\varepsilon})}}, \text{Large metric complexity} \rightarrow \text{small converge rate}$$

Recap.

1. Metric entropy: $\log N(\delta; K, C)$ (P(1)).

2. For $B = \{\theta : \| \theta \| \geq 1\}$ & $B' = \{\theta : \| \theta \|' \leq 1\}$
 $\left(\frac{1}{\delta}\right)^d \frac{\text{Vol}(B)}{\text{Vol}(B')} \leq N(\delta, B, \| \cdot \|) \leq \left(\frac{2}{\delta}\right)^d \frac{\text{Vol}(B + \frac{\delta}{2}B')}{\text{Vol}(B')}$

Roughly, $N(\delta; B, \| \cdot \|) = \left(\frac{1}{\delta}\right)^d \frac{\text{Vol}(B)}{\text{Vol}(B')}$

$\rightarrow N(\delta; B_2(1), \| \cdot \|_2) \asymp \left(\frac{1}{\delta}\right)^d \rightarrow$ 

$N(\delta; B_\infty(1), \| \cdot \|_\infty) \asymp \left(\frac{1}{\delta}\right)^d \rightarrow$ 

3. $F_C = \{f(x) = 0, |f(x) - f(x')| \leq C|x-x'|, x \in [0, 1]\}$

$\log N(\delta; F_C, \| \cdot \|_\infty) \asymp \frac{1}{\delta}$

4. "Parametric functions" $\log N(\delta, F, \| \cdot \|_\infty) \asymp \log(\frac{1}{\delta})$

5. Le Cam: $\eta \Sigma^2 = \log N(\Sigma, F, \| \cdot \|)$

Elliptics: $f(x) = \sum_{i=1}^{\infty} \theta_i p_i(x) \quad x \in [0, 1], \int_0^1 p_i(x) dx = 1, \int_0^1 p_i(x)p_j(x) dx = \delta_{ij}$

Example. Fourier basis
orthogonal function

orthogonal series estimators, θ_j .
smooth functions: "most of the high-order coefficients are nearly zero"

coefficient decay: Fourier basis, $\sum_{j=1}^{\infty} j^{2d} \theta_j^2 \leq 1$ (Holderd functions)

Lipschitz functions: Fourier basis, $\sum_{j=1}^{\infty} j^{2d} \theta_j^2 \leq 1$ (2-times diff)

$d=1, 2, 3, \dots, \theta_j \leq \frac{1}{j^{2d}}$ θ_j decays faster α denotes

L_2 functions $\sum_{j=1}^{\infty} \theta_j^2 \leq \infty$, $\alpha > 1$, θ_j decays faster

Lipschitz $\sum_{j=1}^{\infty} j^2 \theta_j^2 \leq L$ (related to α)

$y_i = \theta_i + \varepsilon_i, \varepsilon_i \sim N(0, 1)$, orthogonal series estimators

$\sum_{j=1}^{\infty} \theta_j^2 j^{2d} \leq 1$, orthogonal series estimators

axis of ellipses, $j^{-d}, -d, 2^{-d}, 3^{-d}, \dots$ high-dim collapse...

Naive est. of $E g = 0$, variance \propto bias \rightarrow var(bias) \propto Var(y)

Truncated series estimator: pick T

Bias $\sum_{j=T+1}^{\infty} \theta_j^2 \leq T^{-2d} \sum_{j=1}^T \theta_j^2 \leq T^{-2d}$, Variance: $\frac{1}{n} \sum_{j=1}^T E(y_i - \theta_j)^2 = \sum_{j=1}^T V(\varepsilon_j) = \frac{1}{n} \sum_{j=1}^T$

$$\leq T^{-2d}$$

$$MSE \leq T^{-\frac{2}{2+1}} + \frac{1}{n} \asymp n^{-\frac{2}{2+1}}, \quad T \asymp n^{\frac{1}{2+1}}$$

$$y_i = f(x_i) + \varepsilon_i \quad f(x_i) = \sum_{j=1}^m \theta_j p_j(x_i)$$

$$z_1 = \frac{1}{n} \sum_{i=1}^n y_i p_1(x_i), \quad \mathbb{E} z_1 = \theta_1 \quad \int_0^1 f(x) p_1(x) dx = \theta_1$$

$$z_j = \frac{1}{n} \sum_{i=1}^n y_i p_j(x_i), \quad \mathbb{E} z_j = \theta_j$$

Metric entropy of $\{ \sum_{j=1}^m \theta_j p_j(x) \mid \theta \in K \}$, $\sum_{j=1}^m \theta_j^2 \leq 1$ infinite metric entropy

$$\log N(\delta, K \| \cdot \|_2) \asymp \left(\frac{1}{\delta} \right)^{1/2}$$

- truncate ellipse at some T can "ignore" all the remaining axis' s ...
- cover finite dimensional ellipse
- cover sub-Gaussian processes.

Metric entropy & sub-Gaussian processes

$$1. X \text{ sub-Gaussian if } \mathbb{E} \exp(tX) \leq \exp\left(\frac{t^2 \sigma^2}{2}\right)$$

$$2. \theta \in K, \rightarrow X_\theta - \bar{X}_\theta \text{ mean 0}$$

$$\text{Sub-Gaussian process w.r.t. } \rho \quad \mathbb{E} \exp(t(X_\theta - \bar{X}_\theta)) \leq \exp\left(\frac{t^2 \sigma^2}{2}\right), \forall t \in \mathbb{R}, \theta, \theta' \in K$$

\rightarrow canonical EIP

$$\text{Bound. } \mathbb{E} \sup_{\theta \in K} \langle \theta, \varepsilon \rangle, \varepsilon \sim N(0, I_d)$$

$$\langle \theta, \varepsilon \rangle = X_\theta \quad X_\theta: \text{sub-GP w.r.t. } \rho(\theta, \theta') = (\theta - \theta')_2$$

$$X_\theta, \mathbb{E} \sup_{\theta \in K} X_\theta$$

1-step discretization.

If N Gaussians

$$\mathbb{E} \sup_{\theta \in K} X_\theta \asymp \sqrt{\log N}$$

$$\mathbb{E} \sup_{\theta \in K} X_\theta \leq 2 \mathbb{E} \left(\sup_{\theta, \theta' \in K} (X_\theta - X_{\theta'}) \right) + 4 \sqrt{D^2 \log N(\delta, k, \rho)}$$

$$\mathbb{E} \sup_{\theta \in K} X_\theta = \mathbb{E} \sup_{\theta} (X_\theta - X_{\theta_0}) \quad (\mathbb{E} X_{\theta_0} = 0) \quad D = \sup_{\theta, \theta'} \rho(\theta, \theta')$$

$$\leq \mathbb{E} \sup_{\theta, \theta'} (X_\theta - X_{\theta'})$$

f. net points θ_i such that $\rho(\theta, \theta_i) \leq \delta$

$$X_\theta - X_{\theta_0} = X_\theta - X_{\theta_i} + X_{\theta_i} - X_{\theta_0}$$

$$\leq \sup_{\theta, \theta' \in K} (X_\theta - X_{\theta'}) + \max_{i=1 \dots N} |X_{\theta_i} - X_{\theta_0}|$$

$$X_\theta - X_{\theta_0} \leq 2 \sup_{\substack{\theta, \theta' \in K \\ \rho(\theta, \theta') \leq \delta}} (X_\theta - X_{\theta'}) + 2 \max_{j=1 \dots N} |X_{\theta_j} - X_{\theta_0}|$$

$$\mathbb{E} \sup_{\theta, \theta'} (X_\theta - X_{\theta'}) \leq 2 \mathbb{E} \left[4 \sqrt{D^2 \log N(\delta, k, \rho)} \right] + 2 \mathbb{E} \max_{j=1 \dots N} |X_{\theta_j} - X_{\theta_0}|$$

$$\max_{i=1 \dots N} |X_{\theta^i} - X_{\theta^1}|$$
$$E \exp t(\underbrace{X_{\theta^i} - X_{\theta^1}}) \leq \exp \left(\frac{t^2 \rho^2(\theta^i, \theta^1)}{2} \right) \leq \exp \left(\frac{t^2 D^2}{2} \right)$$

$X_1 \dots X_N$, σ -sub-Gaussian

$$E \max |X_i| \leq 2\sigma \sqrt{\log N}$$

Reap.

1. Gaussian complexity of K

$$G(K) = \mathbb{P} \sup_{\theta \in K} \theta^T \Sigma, \Sigma \sim N(0, I_d) \quad \begin{matrix} \text{1/2 Lemma} \\ \text{Sub-Gaussian} \\ \text{process} \\ \text{with } L_2 \end{matrix}$$

2. Sub-Gaussian process $\{X_\theta, \theta \in K\}$ mean 0

$$\mathbb{E} \exp(t(X_\theta - X_0)) \leq \exp\left(\frac{t^2 C(\theta, \theta')}{2}\right), \text{ e.g. } X_\theta = \theta^T \Sigma, \mathbb{E} \exp(t(X_\theta - X_0)) \leq \exp\left(\frac{t^2 C(\theta, \theta')}{2}\right)$$

1.b Function spaces.

$$f = \sum_{i=1}^n \theta_i \psi_i, \sum_i \theta_i^2 < \infty$$

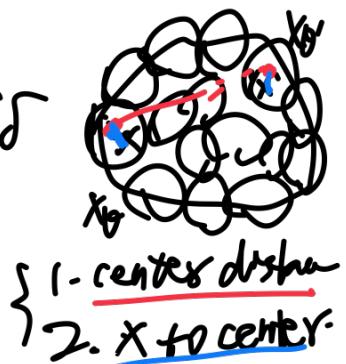
$$\text{Smoothness } \sum_{i=1}^n i^{2d} \theta_i^2 \leq 1 \rightarrow \log N(f, K, 1/\epsilon) \asymp \left(\frac{1}{\epsilon}\right)^{\frac{1}{2d}}$$

$$y_i = \theta_i \psi_i \text{ then } \mathbb{E}[\theta_i - \theta'_i]^2 \sim \frac{1}{I+T} \xrightarrow{\text{bits}} n^{-\frac{2d}{2d+1}}$$

\downarrow at origin fractioned sequences variance

3. One-step. $\mathbb{P} \sup_{\theta \in K} X_\theta \leq \mathbb{E} \sup_{\theta \in K} |X_\theta| \leq \mathbb{E} \sup_{\theta \in K} f(\theta) \leq \delta$

$$X_\theta = X_\theta' + \sqrt{D^2 \log N(\delta, K, \epsilon)}$$



Expected operator norm of random matrix

$$\text{h. W. } W_{ij} \sim N(0, 1) \quad \mathbb{E} \|W\|_{op} = \mathbb{E} \sup_{\|v\|=1} \langle W, v \rangle = \sum_{ij} W_{ij} a_i v_j, a \in \mathbb{R}^n \text{ fixed}$$

$$\mathbb{E} \|W\|_{op} = \mathbb{E} \sup_{\substack{\|M\| \leq 1 \\ \|M\|_F \leq 1}} \langle W, M \rangle = G(M) \quad \begin{matrix} \text{=} \mathbb{E} \sup_{\|M\| \leq 1} \langle W, M \rangle \\ \text{sub-Gaussian with metric } \sqrt{\text{rank}(M) + n - 1}/C \end{matrix}$$

$$D = \sup_{\text{diam } M} \|M - M'\|_F \leq 2$$

$$\mathbb{E} \|W\|_{op} \leq \mathbb{E} \sup_{\substack{M \in \mathbb{R}^{n \times n} \\ \|M - M'\|_F \leq D}} \langle W, M \rangle + \sqrt{\log N(\delta, M, \|F\|_F)} \leq \mathbb{E} \|W\|_{op} \frac{\|M\|_{\text{operator}}}{{\text{rank}(M) + n - 1}} \leq \sqrt{\sum_i \sigma_i^2(M)} \leq \sqrt{\sum_i \sigma_i^2(M)} \mathbb{E} \|W\|_{op}$$

$$\mathbb{E} \sup_{\theta \in K} \{X_1, \dots, X_n\} \leq \sqrt{\log n}$$

$$U = \{u_1, \dots, u_m\} \subset R^n \quad \forall u_i, \|u_i\|_2 = 1, \{u_i^T u_j\}_2 \leq 1$$

$$V = \{v_1, \dots, v_m\} \subset R^d$$

$\{(u_i v_j^T, \dots)\}_{i \in \{-m, \dots, m\}, j \in \{-N, \dots, N\}}, \text{ covering } M \text{ using } U, V.$

U^* covers U , V^* covers V

$$\|UV^T - U^*V^{*T}\|_F = \|UV^T - UV^{*T} + UV^{*T} - U^*V^{*T}\|_F \\ \leq \|U\|_2 \|V^*\|_2 + \|U - U^*\|_2 \|V^*\|_2 \leq \delta$$

$$M \times N : \left(\frac{2}{\delta}\right)^n \times \left(\frac{2}{\delta}\right)^d$$

$$\mathbb{E}\|W\|_{op} \leq \sqrt{2} \delta \mathbb{E}\|W\|_{op} + C \sqrt{n+d} \log(\frac{1}{\delta})$$

intuition: $U, V, U^T W V \sim \mathcal{N}(0, I)$
union bound $(\frac{1}{\delta})^{n+d}$

Non-parametric regressions
 $y_i = f^*(x_i) + \epsilon_i$ $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$

$$f \in \mathcal{F}_L \Leftrightarrow f(0) = 0, |f(x) - f(y)| \leq L|x-y|$$

Non-param. L.S. Represent them.
 $f = \arg\min \|y - f\|_2^2$ finite-dim opt. prob.

Basic ineq: $\frac{1}{n} \sum_i (f(x_i) - y_i)^2 \leq \frac{1}{n} \sum_i (f^*(x_i) - y_i)^2$

$$\mathbb{E} \left(\frac{1}{n} \sum_i (f(x_i) - f^*(x_i))^2 \right) \leq \frac{1}{n} \sum_i (f(x_i) - f^*(x_i))^2 \xrightarrow{\text{Gaussian complexity}} \\ (\text{in-sample error}/\text{fixed design error}) / \|f - f^*\|_{L^2(\rho)}^2 \int (f - f^*)^2 d\rho = \|f - f^*\|_{L^2(\rho)}^2$$

Now Gaussian complexity

$$\mathbb{E} \left(\sum_i (f(x_i) - f^*(x_i))^2 \right) \leq \mathbb{E} \sup_{f \in \mathcal{F}_L} \sum_i (f(x_i) - f^*(x_i))^2 \mathbb{E} \epsilon_i^2$$

Natural metric $\|\bar{f} - f\| = \sup_{f' \in \mathcal{F}_L} \sum_i \langle f - f', \epsilon_i \rangle$ $\mathcal{F}' = \{f(x_1), \dots, f(x_n), f^*\}$

$$\leq \mathbb{E} \sup_{f, f' \in \mathcal{F}_L} \sum_i \langle f - f', \epsilon_i \rangle + \sqrt{D \cdot \log N(\mathcal{F}_L)}$$

$$\|f' - f\|_{L^2(\rho)} \leq \|f' - f\|_\infty$$

$$\|f - f'\|_{L^2(\rho)} \leq \sqrt{D}$$

$$\leq \sqrt{n} \cdot \sqrt{D} \delta + \sqrt{D}$$

$$\leq \delta \sqrt{n} \cdot \sqrt{D} \quad \delta \asymp \sqrt{\frac{1}{n}}$$

$$\leq \sup_{f' \in \mathcal{F}_L} \|f - f'\|_{L^2(\rho)} \quad (\text{can drop } f)$$

Recap

1. Sub-Gaussian process $\{X_\theta : \theta \in K\}$

$$\mathbb{E} \exp(t(X_\theta - X_{\theta'})) \leq \exp(t^2 \frac{\rho(\theta, \theta')}{2})$$

2. One-step discretization. mean $X_\theta = 0$

$$\mathbb{E} \sup_{\theta \in K} X_\theta \leq \mathbb{E} \sup_{\theta, \theta'} X_\theta - X_{\theta'} + \sqrt{D \log N(f, K, \rho)}$$

Applications

1. $\mathbb{E} \|M\|_{op} \leq \sqrt{n} + \sqrt{d}$, $\mathcal{G}\{f M : \text{rank}(M) = 1, \|M\|_F = 1\}$
Nonparametric local smoothing $\rightarrow 35705$

2. Lipschitz fn estimation
N-pn LS. $\hat{f} = \arg \min_{f \in \mathcal{F}} \|y - f\|_2$. $\frac{1}{n} \sum_{i=1}^n (f(x_i) - f^*(x_i))^2 \leq \frac{2}{n} \sum_{i=1}^n (f(x_i) - f(x_i^*))^2$
 $\mathbb{E} \frac{1}{n} \sum_{i=1}^n (f(x_i) - f^*(x_i))^2 \leq \frac{1}{n} \mathbb{E} \sup_{f \in \mathcal{F}_L} \langle \hat{f} - f^*, \varepsilon \rangle$, $\rho(f, f^*) = \|f - f^*\|_{L^2(\mu_p)} = \sqrt{\frac{2}{n} \sum_{i=1}^n (f(x_i) - f^*(x_i))^2} \leq \|f - f^*\|_\infty$
 $\leq \frac{1}{n} \mathbb{E} \sup_{f \in \mathcal{F}_L} \frac{\langle \hat{f} - f^*, \varepsilon \rangle}{\sqrt{n}} + \underbrace{\text{(athy) check...?}}_{\|f - f^*\|_{L^2(\mu_p)}} \Rightarrow \text{do } \log N \frac{\sqrt{n}}{\delta}$

scale up with D $\leq \frac{1}{\sqrt{n}} (\|f\|_\infty \delta) + \sqrt{f} \sqrt{\frac{D \delta^2}{n}}$
 ambient bound, now increasing the resolution of covering, $\begin{aligned} & \delta & 2\delta & 4\delta \\ & |x_0 - x_1 + x_1 - x_2 + x_2 - x_3| & |x_0 - x_1 + x_1 - x_2 + x_2 - x_3| & |x_0 - x_1 + x_1 - x_2 + x_2 - x_3| \\ & |x_0 - x_1| + |x_1 - x_2| + |x_2 - x_3| & |x_0 - x_1 + x_1 - x_2 + x_2 - x_3| & |x_0 - x_1 + x_1 - x_2 + x_2 - x_3| \\ & \text{distance. } \sqrt{2(y_0)^2 + (w)^2} \log(2\delta) \dots & \end{aligned}$

Daddy's hand

$$\mathbb{E} \sup_{\theta} X_\theta \leq \mathbb{E} \sup_{\theta, \theta' \in \{0, 1\}^D} X_\theta - X_{\theta'} + \sqrt{\int_0^D \log N(u, K_d) du} \text{ entropy integral}$$

$$\mathbb{E} \frac{1}{n} \sum_{i=1}^n (f^* - f(x_i))^2 \leq \frac{1}{\sqrt{n}} \int_0^D \frac{1}{\sqrt{u}} du \leq \frac{1}{\sqrt{n}}, \text{ better than } n^{-\frac{1}{2}}, \text{ but if } \text{ad. change}$$

$$D = \sup_{f, f'} \|f - f'\|_{L^2(\mu_p)} \leq L \quad \text{lipschitz smthg} \quad \int_0^D \frac{1}{\sqrt{u}} du \text{ ad. change}$$

Density estimate $X_1 \dots X_n \sim p$ on $[0, 1]$ \rightarrow make assumptions on $p: \rightarrow$ non-param. method.
 \rightarrow make assumptions on metric weaker γ_0 (local smthg).

$$\begin{aligned} \mathbb{P}^G, F_N \in \mathcal{F}^G, \int_0^1 p(x) dx, L(\hat{f}^G, p) &= \sup_x |F_N^W - f(x)| \\ &= \frac{1}{n} \sum_{i=1}^n |X_i - \hat{x}| \end{aligned}$$

Wasserstein - 1: norm

$$W_1(p, q) = \sup_{\|f\|_{\infty} \leq 1} \left| \mathbb{E}_p f - \mathbb{E}_q f \right|$$

$$\hat{p}^{\text{un}} = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{X_i \leq x\}}$$

$$\mathbb{E} W_1(p, p^*) = \sup_{\|f\|_{\infty} \leq 1, f(0)=0} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E}_p f \right| = \sup_{f \in F_p} X_f, X_f = \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E}_p f$$

$$X_f - X_g = \frac{1}{n} \sum_{i=1}^n (f(X_i) - g(X_i)) - (\mathbb{E}_p f - g) \rightarrow \text{Average of bounded RVs . MGF}$$

$X_f - X_g$ changes by at most $\frac{2\|f-g\|_{\infty}}{n}$

Azuma's bound

$$\mathbb{E} \exp(t(X_f - X_g)) \leq \exp\left(\frac{ct^2\|f-g\|_{\infty}^2}{n}\right)$$

Hoeffding
 $(t-\alpha)^2 \leq (X_f - X_g)^2$

$\therefore X_f$ is a sub-Gaussian process
 with $\rho(f, g) = \|f-g\|_{\infty}$

Apply Dudley bound + Lipschitz

$$\mathbb{E} W_1(p, p^*) \leq \frac{1}{\sqrt{n}} : \|f\|_{\infty} \sum_{i=1}^n f(X_i) - \mathbb{E}_p f = O\left(\frac{1}{\sqrt{n}}\right)$$

for all $f \in \mathcal{H}$
 d-dim, rate $n^{-\frac{1}{2d}}$

VC classes

$$\Sigma(\mathcal{C}) = \sup_{C \in \mathcal{C}} \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{X_i \in C\}} \rightarrow \mathbb{P}(\mathcal{C})$$

\mathcal{C} has VC-dim = d

$$\mathbb{E} \Sigma \leq \sqrt{\frac{d \log n}{n}}$$

Covering number for VC collections of sets

$$(VC(d, \rho, n), \mathbb{P}(C)) \leq d^{2d} \left(\frac{1}{\rho}\right)^d$$

sub-Gaussian
 $\leq \sqrt{d}$

Recap.

1. Dudley's chaining bound

$$\mathbb{E} \sup_{\theta \in K} X_\theta \leq \mathbb{E} \sup_{\substack{\theta, \theta' \in K \\ P(\theta, \theta') < \delta}} (X_\theta - X_{\theta'}) + \int_0^D \sqrt{\log N(u, K, \delta)} du$$

Entropy Integral

2. Applications

1. Lipschitz regression $\mathbb{E} \frac{1}{n} \sum_i (f(x_i) - f(x_j))^2 \leq \frac{1}{n}, \geq \frac{1}{n}$

tight! 2. W_1 -distance between P_n & P : $\mathbb{P} W_1(P_n, P) \leq \frac{1}{n}$ in 1-D
 3. VC-classes $\mathbb{E} \sup_{C \in \mathcal{C}} |P_n(C) - P(C)| \leq \sqrt{\frac{d}{n}}, \text{ (0.5 } \sqrt{\frac{d \log n}{n}})$

1. Gaussian comparison ineq. $N(x_1, \dots, x_n)$

$$\mathbb{E} \sup_{\theta} X_\theta \quad \xrightarrow[\text{Gaussian process}]{\{Y_1, \dots, Y_n\}} \mathbb{E} \sup_i X_i \text{ v/s } \mathbb{E} \sup_i Y_i$$

Sudakov-Fernique. Intrahole metric $\rho_x(i, j) = \mathbb{E}(X_i - X_j)^2$

$$\rho_x(i, j) \leq \rho_g(i, j) \quad \forall i, j \Rightarrow \mathbb{E} \sup_i X_i \leq \mathbb{E} \sup_i Y_i$$

Gaussian contraction ineq.

$$\mathbb{E} \sup_{\theta \in K} \langle \phi(\theta), \varepsilon \rangle \quad \phi(\theta) = (\phi_1(\theta_1), \dots, \phi_d(\theta_d))$$

$\hookrightarrow \phi(K)$

$$\text{v.s. } \mathbb{E} \sup_{\substack{\theta \in \Phi(K) \\ \theta \in K}} \langle \phi(\theta), \varepsilon \rangle$$

$$G(F_L^n) = \mathbb{E} \sup_{f \in F_L} \langle f, \varepsilon \rangle, \quad H(F_L^n) = \mathbb{E} \sup_{f \in F_L} \langle f^2, \varepsilon \rangle$$

If ϕ_1, \dots, ϕ_d are L -Lipschitz & $\phi_i(0) = 0$ then $\mathbb{E} \sup_{\theta \in K} \langle \phi(\theta), \varepsilon \rangle \leq \mathbb{E} \sup_{\theta \in K} \|\theta\|_2$

$$\text{eg. } \|f\|_{\infty} \leq b, \quad \phi_i(x) = \begin{cases} \frac{x^2}{2b} & |x| \leq b \\ b & \text{else} \end{cases} \quad H(F_L^n) \leq 2b \cdot G(F_L^n) \quad ?$$

prop. $\mathbb{E}(Q_i - Q_j)^2 \geq \mathbb{E}(\phi(Q_i) - \phi(Q_j))^2$
 1-Lipschitz

Sudakov ineq.

$$\mathbb{E} \sup_{\theta \in K} X_\theta \geq \sup_{\theta \in K} \frac{1}{2} \sqrt{\log N(\frac{1}{2}, K, \delta)}, \quad \rho_x(\theta, \theta') = \mathbb{E}(X_\theta - X_{\theta'})^2$$

Proof: $Y_\theta \sim \mathcal{N}(0, \delta^2/2)$

$$\mathbb{E}(X_{\theta'} - X_\theta)^2 \geq \delta^2$$

$$\mathbb{E} (Y_\theta - Y_{\theta'})^2 = \delta^2$$



$$\mathbb{E} \sup_{\theta \in \Theta} \sum_{i=1}^n \sum_{j=1}^m X_{ij} Z_i \geq \mathbb{E} \sup_{i \in \{1, \dots, m\}} T_{i,i} \geq \delta \sqrt{\log n}$$

Metric on an L_1 -ball

$$P_1(\gamma) = \mathbb{E} \sup_{\|\theta\|_1 \leq 1} \langle \theta, \varepsilon \rangle = \mathbb{E} \| \theta \|^1 / \sum_{i=1}^n \sum_{j=1}^m \log$$

$$\log N(\delta, k \cdot \| \cdot \|_2) \sim \frac{\log \frac{1}{\delta}}{\| \cdot \|_2}$$

$$\log N(\delta, \| \cdot \|_2 \cdot \| \cdot \|_1) \sim d \log \frac{1}{\delta}$$

L_1 -ball, $\mathbb{E} P_1$

Matrix estimation.

- 1. Covariance matrix est. $X_1, \dots, X_n \sim P, \Sigma = E(X_i X_i^T)$
- 2. Matrix completion $\tilde{Y} = M + S, M_{ij} \xrightarrow{\text{observed}} 0, S_{ij} \xrightarrow{\text{observed}} 0$
- 3. Matrix $X = M + W$

$g_i = \text{tr}(M^T X_i) + \sum_j \beta_j^* X_{ij} + \varepsilon_i$

Matrix completion $X \in \mathbb{R}^{d_1 \times d_2}$

want to estimate M^* , observe \tilde{Y}

Stochastic Block Models

$$M_{ij}^* = \begin{cases} p & i=j \\ q & i \neq j \end{cases}, Y_{ij} = \begin{cases} 1 & \text{with prob. } M_{ij}^* \\ 0 & \text{with prob. } 1 - M_{ij}^* \end{cases}$$

$$E(Y_{ij}) = M_{ij}^*, Y = M^* + W, \text{ mean } 0$$

Reminder - Low rank matrix est./sparse-vec est. 705

hard thresholding $\xrightarrow{\text{w.p. } \sum_i \| \beta_i - \hat{\beta}_i \|_2^2 \leq \sum_i m_i \|\beta_i^*\|^2, \frac{\sigma^2 \log d}{n}}$

$\hat{\beta}^*$ is sparse, $\frac{\log d}{n}$

β^* is L_1 sparse

(Matrix analog, symmetric) $M^* = \sum_{i=1}^n \alpha_i v_i v_i^T, Y = \sum_{i=1}^n \beta_i w_i w_i^T, M = \sum_{i=1}^n \beta_i \frac{\sqrt{R_i \log n}}{\sqrt{R_i}} v_i v_i^T$

$$Y = M^* + W, \text{ If } (W)_{\text{top}} \leq t/2, \| M - M^* \|_F^2 \leq \sum_{i=1}^n m_i \|\beta_i^*\|^2, t^2 \}$$

Recap:

1. Sudakov-Fernique. $\{X_1 \dots X_n\} \{Y_1 \dots Y_n\}$ Gaussian process.

$$\text{If } \mathbb{E}(X_i - X_j)^2 \leq \mathbb{E}(Y_i - Y_j)^2 \forall i, j$$

Then. $\mathbb{E} \sup_i X_i \leq \mathbb{E} \sup_j Y_j$

(a) Grammatical contraction.

$\mathbb{E} \sup_{\theta} \langle \varphi(\theta), \varepsilon \rangle \leq \mathbb{E} \sup_{\theta} \langle \varphi, \varepsilon \rangle$, If $\varphi(\theta)$ is a contraction
e.g. $\varphi(\theta) = \varphi_1(\theta) \dots \varphi_n(\theta)$

(b) Sudakov minimization. $\sup_{\theta} \mathbb{E} \log M(\theta, \mu_\theta)$ μ_θ is $1 - \frac{1}{2} \mathbb{E} \|\theta\|^2$

$$\mathbb{E} \sup_{\theta} X_\theta \geq \sup_{\theta} \frac{d}{2} \int \log M(\theta, \mu_\theta)$$

$$\text{where } \rho(X_\theta, \theta) = \mathbb{E}(X_\theta - \mu_\theta)^2$$

2. Lévy's inequality.

If f is L -Lipschitz. $\mathbb{P}(|f(z) - \mathbb{E} f(z)| \geq t) \leq 2 \exp\left(-\frac{t^2}{2L^2}\right)$

Matrix Est. Z is std-normal. $\sup_{\theta} |X_\theta - \mathbb{E} \sup_{\theta} X_\theta|$

$\hat{Y} = M^* + W$ \hookrightarrow indep. sub-Gaussian mean deviations.

$$T = \sum_{i=1}^n \beta_i V_i V_i^T, M = \sum_{i=1}^n \beta_i \mathbf{1}_{\{\beta_i \geq 2t\}} V_i V_i^T, M^* = \sum_{i=1}^n \alpha_i U_i U_i^T$$

Suppose $\|M\|_{\text{op}} \leq t$, then $\|\hat{M} - M^*\|_F^2 \leq \sum_{i=1}^n \min\{\beta_i^2, 2t^2\} \Rightarrow$ prove it!

$\beta_1 \leq \beta_2 \leq \dots \leq \beta_n$, $\alpha_1 \leq \alpha_2 \leq \dots \leq \alpha_n$.

$S = \{i : |\beta_i| \geq 2t\}$, Weyl's Law. (matrix perturbation) $\Rightarrow \alpha_i - \beta_i \leq \|W_i\|_F$

$$\|\hat{M} - M^*\|_F \leq \|\hat{M} - M_S^*\|_F + \|M_S^* - M^*\|_F \quad M_S^* = \sum_{i \in S} \alpha_i U_i U_i^T$$

$$\|\hat{M} - M^*\|_F \leq \sqrt{2S} \|\hat{M} - M_S^*\|_{\text{op}} \quad \sqrt{\sum_{i \in S} \alpha_i^2} \leq \sqrt{\sum_{i \in S} \beta_i^2}$$

$$\sqrt{2S} \|\hat{M} - M_S^*\|_{\text{op}} \leq \sqrt{2S} \|\hat{M} - T\|_{\text{op}} + \|\hat{M} - T\|_{\text{op}} + \|\hat{M} - M_S^*\|_{\text{op}} \leq t = \max_{i \in S} \{|\alpha_i| \leq 3t\}$$

$\text{rank } S \leq 2t \sqrt{S} \leq t \sqrt{|S|}$ by Weyl's

$$\|\hat{M} - M^*\|_F \leq \sqrt{t^2 |S| + \sum_{i \in S} \alpha_i^2} \leq \sqrt{\sum_{i=1}^n m_i h_i^2 \alpha_i^2} \quad \begin{cases} \text{on } S, |\alpha_i| \leq 3t \\ \text{on } S^c, |\alpha_i| \leq 3t \end{cases}$$

Suppose: $\mathbb{E} \|W\|_{\text{op}} \leq \sqrt{n}$, $t \asymp \sqrt{n}$

$$\|\hat{M} - M^*\|_F^2 \leq \sum_{i=1}^n m_i h_i^2 \alpha_i^2 \leq n r \quad \text{Suppose } M^* \text{ has rank } r / \|\hat{M} - M^*\|_F^2 \leq nr$$

$\frac{n}{r} \leq \frac{1}{2} \text{rk } R \quad \frac{(n-r)^2}{r} \leq D$
eigenvalues

Oracle Ineq. form
 $\|M - M^*\|_F^2 \leq \inf_{\theta} \|M^* - \theta\|_F^2 + n \text{rank}(\theta) \Big\]$, M not necessarily low-rank
 θ is low rank

prove: $\|M - M^*\|_F^2 \leq \|M^* - M\theta\|_F^2 + \gamma \times n$
 $\leq \sum_{i=1}^n m_i(\theta_i, n) \leq \sum_{i \in S^c} \theta_i + \sum_{i \in S}$

Universal SVT

Strongly stochastically transitive matrices (SST) Nihar. March. Sila
 $\{1 \dots n\} P(i \succ j)$. parameter model: $\theta_1 \dots \theta_n$, Bradley-Terry law
 non-parametric Π^* decrease: best $\xrightarrow{\text{increas}} \xleftarrow{\text{decreas}}$ worst $\xrightarrow{\text{increas}} \xleftarrow{\text{decreas}}$
 If i is better than j . $y = \frac{1}{1 + \exp(\theta_i - \theta_j)}$
 $P(i \succ k) \geq P(j \succ k), \forall k.$ Bernoulli noise
 out measure/low rank

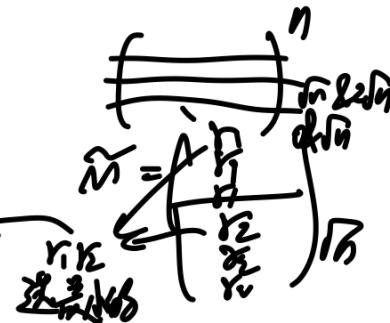
Sommer Chatterjee

$$\|M^* - M_{\sqrt{n}}\|_F^2 \leq n\sqrt{n} \Rightarrow \|M - M^*\|_F^2 \leq n\sqrt{n}$$

$$\|M^* - \tilde{M}\|_F^2 \leq \sum_{i=1}^n \sum_{j=1}^n (M_{ij}^* - \tilde{M}_{ij})^2$$

positive

$$\leq n \sqrt{n}$$



W. with prob $\geq 1 - \delta$, $\|W\|_F \leq \sqrt{n} + \sqrt{\log(1/\delta)}$

Symmetric

Gaussian. $\|W\|_F = \sup_{\|u\|=1} u^T W u \xrightarrow{\text{N} \frac{1}{2} \text{ rot of Sphere. } S^{n-2}} \sqrt{\frac{2}{n}}$

$\sup_{\|u\|=1} u^T \theta \leq \sup_{u \in \mathbb{R}^n} u^T \theta \leq \sqrt{n} \max_{u \in \mathbb{R}^n} u^T \theta$ \hat{u} closest not point

$$\hat{u}^T \theta \geq \frac{u^T \theta}{\sqrt{n}}$$

$$\geq -\frac{\|u^T \theta\|_2}{\sqrt{n}}$$

$$u^T \theta \leq 2\hat{u}^T \theta \leq 2\sup_{u \in \mathbb{R}^n} u^T \theta$$

Sub-Gauss
 $P(\sup_{u \in \mathbb{R}^n} u^T W u \geq t) \leq 2^n \exp(-t^2)$

Recap:

Chapter 6.

- Matrix estimation, signal noise model

$$Y = M^* + W \quad \text{independent mean } 0, \text{ sub-gaussian}$$

- Singular value thresholding.

$$\hat{M} = \sum_{i=1}^n \beta_i \mathbf{u}_i \mathbf{v}_i^T, \text{ where } Y = \sum_{i=1}^n \beta_i \mathbf{u}_i \mathbf{v}_i^T$$

- Guarantee If $\|W\|_{\text{op}} \leq t$, $\|\hat{M} - M^*\|_F \leq \sum_{i=1}^n m_i \{\alpha_i^2 + t^2\}$

$$\text{where } M^* = \sum_{i=1}^n \alpha_i \mathbf{u}_i \mathbf{v}_i^T \approx \min_{\theta} \|M^* - \theta I_F\|^2 + \lambda \text{rank}(\theta)$$

$$\Rightarrow \|W\|_{\text{op}} \geq 1 - \delta$$

$$\|W\|_{\text{op}} \lesssim \sqrt{n + (\log \frac{1}{\delta})}$$

Oracle
inequality
increasing

approx. error est. error

- Example. SST. matrix \downarrow decreases, \sqrt{n} close to \sqrt{n} -rank matrix

Gaussian covariance matrix

$$X_1, \dots, X_n \sim N(0, \Sigma) \rightarrow \frac{1}{n} \mathbf{X}^T \mathbf{X}, \Sigma = \frac{1}{n} \mathbf{X}^T \mathbf{X}, \text{ how close } \Sigma \& \Sigma^*$$

$$\text{Thm. } P\left(\max\left(\frac{X}{\sqrt{n}}\right) \geq (1+\delta)\sqrt{\Sigma} \text{ s.t. } \sqrt{\frac{\text{Tr}(\Sigma)}{n}} \leq \exp\left(\frac{-n\delta^2}{2}\right)\right) \xrightarrow{\text{proof}}$$

$$P\left(\max\left(\frac{X}{\sqrt{n}}\right) \leq (1-\delta)\lambda_{\min}(\Sigma) - \sqrt{\frac{\text{Tr}(\Sigma)}{n}}\right) \leq \exp\left(-\frac{n\delta^2}{2}\right)$$

$$X_1, \dots, X_n \sim N(\mu, \Sigma), E\|\hat{\mu} - \mu\|_2^2 = E\|\mathcal{N}(0, \frac{\Sigma}{n})\|_2^2 = \frac{\text{Tr}(\Sigma)}{n}, \text{ if } \Sigma \text{ is full rank of } \Sigma$$

Σ is identity. with prob $\geq 1 - 2\exp\left(-\frac{n\delta^2}{2}\right)$

$$\lambda_{\max}\left(\frac{\mathbf{X}^T \mathbf{X}}{n}\right) \leq \left(1 + \sqrt{1 + \frac{\text{Tr}(\Sigma)}{n}}\right)^2 \leq (1 + \Sigma^2 + 2\Sigma) \xrightarrow{\text{if } \Sigma \text{ is full rank of } \Sigma} \|\frac{\mathbf{X}^T \mathbf{X}}{n} - I\|_{\text{op}} \leq \Sigma^2 + 2\Sigma$$

$$\lambda_{\min}\left(\frac{\mathbf{X}^T \mathbf{X}}{n}\right) \geq 1 + \Sigma^2 - 2\Sigma$$

$$X = \sqrt{n} \mathcal{Z}, W \sim N(0, I_d), X_i \sim N(0, \Sigma)$$

$$\|\hat{\Sigma} - \Sigma\|_{\text{op}} = \|\Sigma^{\frac{1}{2}} \left(\frac{W^T W}{n} - I\right) \Sigma^{\frac{1}{2}}\|_{\text{op}} \leq \|\Sigma^{\frac{1}{2}}\|_{\text{op}} \left\| \frac{W^T W}{n} - I \right\|_{\text{op}} \xrightarrow{\text{Corollary}} \sqrt{\frac{\text{Tr}(\Sigma)}{n}} \leq \sqrt{\frac{\text{Tr}(\Sigma)}{n}} + \frac{d}{n} + \delta$$

$$\sqrt{\frac{\text{Tr}(\Sigma)}{n}} \leq \Sigma^2 + 2\Sigma$$

$$\text{w.p. } \geq 1 - 2\exp\left(-\frac{n\delta^2}{2}\right), \left\| \frac{W^T W}{n} - I \right\|_{\text{op}} \leq \sqrt{\frac{d}{n}} + \frac{d}{n} + \delta$$

$P\left(\frac{\sigma_{\max}(\tilde{X})}{\sqrt{n}} \geq \mathbb{E}\frac{\sigma_{\max}(\tilde{X})}{\sqrt{n}} + t\right) \leq \exp\left(-\frac{n t^2}{2}\right)$
 $\left(\left|\frac{\sigma_{\max}(\tilde{X})}{\sqrt{n}} - \mathbb{E}\frac{\sigma_{\max}(\tilde{X})}{\sqrt{n}}\right| \leq \frac{\|X - Y\|_{op}}{\sqrt{n}} \text{ s.t. } \frac{\|X - Y\|_{op}}{\sqrt{n}} = \frac{1}{\sqrt{n}} \text{ (psd)}\right)$
 $\hookrightarrow \text{Weglis Ineqn.}$
 $\mathbb{E}\frac{\sigma_{\max}(X)}{\sqrt{n}} = \mathbb{E}\frac{\sigma_{\max}(W\sqrt{\Sigma})}{\sqrt{n}} = \mathbb{E}\sup_{\|u\|=1} \frac{u^T W \sqrt{\Sigma} v}{\sqrt{n}} = \mathbb{E}\sup_{\|u\|=1} u^T W v$
 $Z_{u,v} = u^T W v, \mathbb{E}(Z_{u,v} - \tilde{Z}_{u,v})^2 = \mathbb{E}(u^T W v - \tilde{u}^T W v)^2 = \mathbb{E}(u_i v_j - \tilde{u}_i \tilde{v}_j)^2$
 cross term.
 $= \text{tr}(W(u - \tilde{u})^T W(u - \tilde{u})^T)$
 $= (\tilde{u}^T W - u^T W)(W^T W) \leq 0$
 $\text{since } W^T W \text{ is psd}$
 $\leq \|W(u - \tilde{u})\|^2$
 $\leq \|W\|_F \|u - \tilde{u}\|_2^2$
 $\leq \|W\|_F \|u - \tilde{u}\|_2^2 + \|W\|_F \|v - \tilde{v}\|_2^2$
 $\leq \|u - \tilde{u}\|_2^2 + \left(\sum_i \|\tilde{u}_i - u_i\|_2^2\right)$
 $\text{using } \|u\|_2 = \text{using } \|\tilde{u}\|_2$
 $\mathbb{E} X_{u,v} = \mathbb{E} \langle u, g \rangle + \mathbb{E} \langle v, h \rangle$
 $g \sim N(0, I), h \sim N(0, I)$
 $\mathbb{E} (X_{u,v} - \tilde{X}_{u,v}) = \|u - \tilde{u}\|_2^2 + \left(\sum_i \|\tilde{u}_i - u_i\|_2^2\right)$
 $\mathbb{E} \sup_{u,v} Z_{u,v} \leq \mathbb{E} \sup_{u,v} X_{u,v} \leq \|W\|_F \left(\mathbb{E} \sup_{u \in \mathcal{U}} \langle u, g \rangle + \mathbb{E} \sup_{v \in \mathcal{V}} \langle v, h \rangle \right) \leq \sup_{u \in \mathcal{U}} \sup_{v \in \mathcal{V}} \langle u, g \rangle + \sup_{u \in \mathcal{U}} \sup_{v \in \mathcal{V}} \langle v, h \rangle$
 $= \|W\|_F \sqrt{\sum_{i=1}^n \mathbb{E} \langle u_i, g \rangle^2} + \sqrt{\sum_{i=1}^n \mathbb{E} \langle v_i, h \rangle^2} \leq \|W\|_F \sqrt{\sum_{i=1}^n \mathbb{E} \langle u_i, g \rangle^2} + \sqrt{\sum_{i=1}^n \mathbb{E} \langle v_i, h \rangle^2}$
 $\mathbb{E} \frac{\sigma_{\max}(W\sqrt{\Sigma})}{\sqrt{n}} \leq \|W\|_F \sqrt{\sum_{i=1}^n \mathbb{E} \langle u_i, g \rangle^2} + \sqrt{\frac{\text{tr}(\Sigma)}{n}}$

Qqwww

Recap,

1. Gaussian cov est. $X_1 \cdot X_n \sim N(\Sigma)$

$$\Sigma = \frac{1}{n} \sum_{i=1}^n X_i X_i^T \text{ with probability } 1 - 2 \exp\left(-\frac{n\delta^2}{2}\right)$$

$$\frac{\|\Sigma - \Sigma\|_{op}}{\|\Sigma\|_{op}} \leq 2\sqrt{\frac{d}{n}} + 2\delta + \left(\delta + \sqrt{\frac{d}{n}}\right)^2$$

Alternatively, w.p. $1 - \delta$ if $n \geq d$, $\frac{\|\Sigma - \Sigma\|_{op}}{\|\Sigma\|_{op}} \leq \sqrt{\frac{d}{n}} + \sqrt{\frac{(d+1)\delta}{n}}$

2. (Not covered)

Sub-Gaussian case. suppose each X_i is sub-Gaussian ($\frac{d}{n} + \dots$)

$$\frac{\|\Sigma - \Sigma\|_2}{\sigma^2} \leq \sqrt{\frac{d}{n}} + \sqrt{\frac{\log d}{n}} + \frac{d}{n} \quad (\forall d) \rightarrow 2\text{-step discretization}$$

Q_1, \dots, Q_n symmetric, mean 0, matrices, study $\|\frac{1}{n} \sum_i Q_i\|_{op}$

IV) Chernoff method.

$$P(y \geq t) = P(\exp(sy) \geq \exp(st)) \leq \inf_{s > 0} \frac{\exp(st)}{\exp(sy)} \xrightarrow{s \rightarrow \infty} \frac{\exp(t)}{\exp(y)} \xrightarrow{\text{special.}} \alpha_1 = X X^T \leq \frac{\beta - \Sigma\|_{op}}{\beta}$$

$$[BS] \log(\mathbb{E} \exp(\lambda)) \cdot M = U^T S U, \exp(M) = U \mathbb{E} \left(\exp(S) \mathbb{E} \exp(S) \right) U$$

$$\bar{Q} = \frac{1}{n} \sum_i Q_i, \mathbb{E}(\exp(\lambda \bar{Q})) = \mathbb{E} \|\exp(\lambda \bar{Q})\|_{op} \leq \mathbb{E}[\text{tr}(\exp(\lambda \bar{Q}))] = \text{tr}[\mathbb{E}(\exp(\lambda \bar{Q}))]$$

Sub-Gaussian, $\sqrt{V}, \mathbb{E} \exp(\lambda \bar{Q}) \leq \exp\left(\frac{\lambda^2 V}{2}\right), \forall \lambda, (\mathbb{E} \exp(X) \leq \exp(\frac{\lambda^2 E[X]}{2}))$

Q : Def: \exists square symmetric similar to vector X, S

B symmetric, $Q = \Sigma B, \Sigma = \{1, -1\} - Q_i$'s B^2 sub Gaussian
 \rightarrow if b ~ b' equal probability b-sub-Gaussian.

$$\mathbb{E} \exp(\lambda B) = \frac{1}{2} [\exp(\lambda B) + \exp(-\lambda B)], \exp(\lambda B) = \sum_{j=0}^{\infty} \frac{\lambda^j}{j!} B^j$$

$$\mathbb{E} \exp(\lambda \Sigma B) \leq \sum_{j=0}^{\infty} \frac{(\lambda \Sigma B)^j}{(2j)!} \leq \sum_{j=0}^{\infty} \left(\frac{\lambda^2 B^2}{2}\right)^j j! \leq \exp\left(\frac{\lambda^2 B^2}{2}\right) \Rightarrow \Sigma B$$

$$\text{tr}(\mathbb{E} \exp(\lambda \bar{Q}))$$

$$\bar{Q} = \frac{1}{n} \sum_i Q_i, Q_i$$

For vector x : $\mathbb{E} \exp(\lambda x) = \mathbb{E} \exp\left(\frac{\lambda x}{n}\right)$
~~but not true for matrix~~
 $\text{tr} \mathbb{E} \exp(\lambda B + C) \neq \text{tr} \mathbb{E} \exp(\lambda B) \mathbb{E} \exp(C)$

Lieb's theorem. If symmetric H

$$f(A) = \text{tr}(\exp(H + \log A)), f \text{ is a concave fn of } A$$

$$\begin{aligned} & \text{tr}(E[\exp(S\bar{Q})]) \\ &= E[\text{tr}(\exp(S \sum_{i=1}^n Q_i) + \log E[\exp(SQ_i)])] \\ &\stackrel{\text{linearity}}{\leq} E[\text{tr}(\exp(S \sum_{i=1}^n Q_i) + \log E[\exp(SQ_i)])] \\ &\stackrel{(n \geq 1)}{\leq} \text{tr} \exp\left(\sum_{i=1}^n \log E[\exp(SQ_i)]\right) \end{aligned}$$

Matrix Hoeffding \$Q_1, \dots, Q_n\$ symmetric mean 0, sub-Gaussian \$Q_i\$

$$P\left(\frac{1}{n} \sum_i Q_i \geq t\right) \leq 2 \exp\left(\frac{-nt^2}{2\sigma^2}\right), \sigma^2 = \frac{1}{n} \sum_i V_i \text{ lop}$$

Chernoff.

$$\begin{aligned} P\left(\lambda \max\left(\frac{1}{n} \sum_i Q_i\right) \geq t\right) &\leq \inf_{s > 0} \exp(st) \text{tr}(E[\exp(S\bar{Q})]) \\ \sum_{i=1}^n \log E[\exp\left(\frac{S\bar{Q}_i}{n}\right)] &\leq \sum_{i=1}^n \log \exp\left(\frac{S^2 V_i}{2\sigma^2}\right) = \frac{S^2 \sum_i V_i}{2\sigma^2} \\ \text{tr} \exp\left(\frac{S^2 \sum_i V_i}{2\sigma^2}\right) &\leq d \exp\left(\frac{\sigma^2 \sigma^2}{2n}\right), \sigma^2 = \frac{1}{n} \sum_i V_i \text{ lop} \end{aligned}$$

$$P(-st) \leq \inf_{s > 0} \exp(st) \exp\left(\frac{S^2 \sigma^2}{2n}\right) \leq \exp\left(\frac{-nt^2}{2\sigma^2}\right)$$

$$d \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix} \Sigma_1 \text{ lop of the.}$$

$$d \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} \Sigma_2, \bar{Q} = \frac{1}{d} \sum_{i=1}^d \Sigma_i e_i e_i^T, \| \bar{Q} \|_{\text{lop}} = \frac{1}{d}$$

$$V_i = e_i e_i^T, \sigma^2 = \| \bar{Q} \|_{\text{lop}}^2$$

$$P(\|\bar{Q}\|_{\text{lop}} \geq t) \leq d \exp\left(\frac{-dt^2}{2}\right), \text{ w.p. } \frac{2}{3} \|\bar{Q}\|_{\text{lop}} \leq \sqrt{\frac{\log d}{2}}$$

If \$\Sigma_i\$ is given \$N(0, I)\$, \$\bar{Q} = \frac{1}{d} \sum_{i=1}^d g_i g_i^T \Rightarrow \|\bar{Q}\|_{\text{lop}} \leq \sqrt{\log d}\$, tight upper bound

Matrix Bern. Nrm. \$Q_1, \dots, Q_n\$. \$(Q_i)_i \leq b\$ \$Var(Q_i) \leq \frac{1}{n} \sum_i \frac{1}{n} \sum_j Q_i Q_j^T\$

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n Var(Q_i) \text{ lop}$$

$$\left| \frac{1}{n} \sum_{i=1}^n Q_i \text{ lop} \geq t \right| \leq 2d \exp\left(\frac{-nt^2}{2\sigma^2}\right), \sigma^2 = \frac{1}{n} \sum_{i=1}^n V_i \text{ lop}$$

sigma lop
Bernoulli
ext Hoeffding
tight

Recap

1. Matrix Hoeffding: $Q_1 \dots Q_n$ symmetric
 $P\left(\left\|\frac{1}{n} \sum_i Q_i\right\|_{\text{op}} \geq t\right) \leq 2d \exp\left(\frac{-nt^2}{2d^2}\right)$, $\sigma^2 = \frac{1}{n} \sum_i \text{Var}(Q_i)$
- $\sum_i \mathbb{E}[\text{Var}(Q_i)] \leq \frac{\sigma^2}{n}$
 X_i is β -sub-Gaussian
 $\sigma^2 \leq \frac{\sigma^2}{n}$
2. Matrix Bernstein. Suppose $\|Q_i\|_{\text{op}} \leq b$, Q_i is γ -sub-Gaussian
 $P\left(\left\|\frac{1}{n} \sum_i Q_i\right\|_{\text{op}} \geq t\right) \leq 2d \exp\left(\frac{-nt^2}{2(bt + \sigma^2)}\right)$
- $\mathbb{E}[\exp(tQ_i)] \leq \exp(tV_i)$
 $\text{Var}[Q_i] \leq \frac{t^2}{n} \sum_i \text{Var}(Q_i)$
- $\sigma^2 = \left\| \frac{1}{n} \sum_i \text{Var}(Q_i) \right\|_{\text{op}}$
- $X_1 \dots X_n \sim P$. $Q_i = X_i X_i^T - \mathbb{E}[X_i X_i^T]$, $\|X_i\|_2 \leq \sqrt{b}$
- mean Σ Then $P\left(\left\|\frac{1}{n} \sum_i \Sigma\right\|_{\text{op}} \geq t\right) \leq 2d \exp\left(\frac{-nt^2}{b^2(\sum_i \text{Var}(X_i))}\right)$
- cov. Σ $\left\|\Sigma\right\|_{\text{op}} \leq \sup_{\|X_i\|=1} \mathbb{E}[X_i X_i^T] \leq b$
- proof: To use Bernstein: $\left\|Q_i\right\|_{\text{op}} \leq \|X_i X_i^T\|_{\text{op}} \leq \|X_i\|_2^2 \leq b$, $\left\|\Sigma\right\|_{\text{op}} \leq b$
- $\text{Var}(Q_i) = \text{Var}(X_i X_i^T) = \mathbb{E}(X_i X_i^T)^2 - \Sigma^2 \leq \mathbb{E}[X_i X_i^T X_i X_i^T] - \Sigma^2 \leq b^2$
- $\left\|\Sigma\right\|_{\text{op}} \leq \sup_{\|X_i\|=1} \mathbb{E}[X_i X_i^T] \leq b$
- $\sigma^2 = b \left\|\Sigma\right\|_{\text{op}}$
- $X_i \rightarrow \mathbb{R}^d$ -sphere, $\|X_i\|_2 = \sqrt{d}$, $\mathbb{E}(X_i X_i^T) = I_d$, $\left\|\Sigma\right\|_{\text{op}} = 1$
- C-subGaussian: $\left\|\Sigma\right\|_{\text{op}} \leq \sqrt{d} + O\left(\frac{1}{n}\right)$, $\left\|\Sigma\right\|_{\text{op}} \leq \sqrt{\frac{d \log d}{n}} + O\left(\frac{\sqrt{d \log d}}{n}\right)$ lower order
- \hookrightarrow tighter uniform bound
- $X_i \rightarrow \mathbb{R}^d$ e.g., $j \sim \text{unif}(d)$, $\mathbb{E}(X_i X_i^T) = I_d$, d -sub-Gaussian, $\frac{d^{3/2}}{\sqrt{n}} \rightarrow$ worse than Bernstein, $\frac{d}{\sqrt{n}}$ water bound \leq
- Estimating Structured Cov Matrix (Sparse Matrix).
- Σ is sparse: $A_{ij} = \begin{cases} 1 & \text{if } \Sigma_{ij} \neq 0 \\ 0 & \text{otherwise} \end{cases}$, $\|A\|_{\text{op}} \approx \text{small}$, $\|A\|_F \approx \text{small}$.
- use $\frac{1}{n}$ threshold in entries
- with prob $\geq 1 - \delta$, $\left\|\frac{1}{n} \sum_i \Sigma_{ij}\right\|_\infty \leq 2$, $\sum_{j,k} \Sigma_{jk} = \sum_{i,j} X_{ij} X_{ik}$, X_{ij} is σ -sub Gaussian
- $\left(P\left(\left\|\sum_{j,k} \Sigma_{jk}\right\|_{\text{op}} \geq t\right) \leq \exp\left(\frac{-nt^2}{2\sigma^2}\right) \right)$
- $\Rightarrow \left\|\frac{1}{n} \sum_i \Sigma_{ij}\right\|_\infty \leq \frac{\sigma^2 \sqrt{\log(1/\delta)}}{n}$, pick t to be $t = \sqrt{\log(1/\delta)}$

$\tilde{\Sigma} = \overline{I}_{\leq t}(\Sigma)$ condition on $\|\tilde{\Sigma} - \Sigma\|_{\infty} \leq t$, $t \geq C \sigma^2 \sqrt{\frac{\log d}{n}}$

} zero coordinate $\tilde{\Sigma}_{jj} = 0$

} non-zero coordinate. $|\tilde{\Sigma}_{jj} - \Sigma_{jj}| \leq |\tilde{\Sigma}_{jj} - \frac{1}{t} \tilde{\Sigma}_{ii}| + |\tilde{\Sigma}_{jj} - \Sigma_{jj}| \leq 3t$

$|\tilde{\Sigma} - \Sigma| \leq 3tA$
↳ elementwise

$\|\tilde{\Sigma} - \Sigma\|_{\text{op}} \leq 3t \|A\|_{\text{op}}$

② 1) Σ have $t \leq$ threshold Σ ,
 $(\tilde{\Sigma})$

通过 $\tilde{\Sigma}$ 去 estimate Σ sparse

不相似 Σ and structured.
 $\sqrt{\log d}$ scale up
 $\hat{\Sigma} \uparrow$ sparse

Recap.

1. Applications of matrix Bernstein to Covariance est.

(a). $x_i \sim \text{unif}(\sqrt{d}\text{-sphere})$, union bound better

(b). $x_i \sim t\sqrt{d}e_j$, $j \sim \text{unif}(d)$, union bound worse

$$\|\sum - I\|_2 \leq \sqrt{\frac{d \log d}{n}}$$

2. Estimating sparse covariance matrices

$X_{ij} \rightarrow \sigma$ -sub Gaussian

$$\|\sum - \Sigma\|_\infty \leq \sigma \sqrt{\log d / n} \|T(\sum) - \Sigma\|_2 \leq \sigma \sqrt{\log d / n} \|A\|_{\text{op.}}$$

Sparse linear models compressed sensing

$$y = X\theta^* + w, \theta^* \text{ is sparse } \|B\|_q \leq q \cdot \|\theta^*\|_1$$

sensing vector $\leftarrow x_1 \cdots x_n$, $n \ll d$, $X: n \times d$

$$\theta^* = \psi^T B^*$$

min $\|\theta\|_0$ s.t. $y = X\theta$, θ^* is unique

min $\|\theta\|_1$ s.t. $y = X\theta$, ℓ_1 -norm linear interpolation

(1) Basis pursuit

$$\theta^* \rightarrow S \subseteq \{1, \dots, d\}, \theta \in \text{null}(X), (X\theta = 0) \quad \theta^* + \theta \text{ is also a solution}$$

$$C(S) = \{\|\theta_S^c\|_1 \leq \|\theta_S\|_1\}. \text{ Restricted null space. } C(S) \cap \text{null}(X) = \{0\}$$

Thm. 1.1. Basis pursuit works for any θ^* with support S , ($\hat{\theta} = X\theta^*$ unique solution. $\|\hat{\theta}\|_1$)

(2) Restricted null space RNS holds for S

$$\text{Proof: } 2 \Rightarrow 1. \hat{\theta} \text{ is a soln to BP. } \|X\hat{\theta}\|_1 \leq \|\theta^*\|_1, \Delta = \hat{\theta} - \theta^*$$

$$\|\theta^*\|_1 \geq \|\Delta + \theta^*\|_1 = \|(\Delta + \theta^*)S\|_1 + \|\Delta S^c\|_1 \geq \|\theta^*\|_1 + \|\Delta\|_1 + \|S^c\|_1$$

$$\Rightarrow \|\Delta_S\|_1 \leq \|\Delta\|_1, \Delta \in C(S) \therefore RNS \Rightarrow \Delta = 0$$

$1 \Rightarrow 2$
If RNS fails, θ^* with supp $S \not\subseteq \hat{S}$ which solves BP

$$\theta \in \text{null}(X) \cap C(S), y = X\theta_S \text{ solves BP (y, X)} \xrightarrow{\|\theta_S\|_1, \|\theta_S^c\|_1} \theta_S \neq 0, X\theta_S^c = -X\theta_S^c$$

X so that $RN(S)$ holds for $|S| \leq s$, BP will exactly reconstruct any θ^* , if $\theta^* \in S$

1. Pairwise incoherence δ , $\left\| \frac{X^T X}{n} - I \right\|_F \leq \delta_{PW}$

If $\delta_{PW} < \frac{1}{25}$ then $RN(S)$ holds for any $|S| \leq s$
 $(\text{if } X_{ij} \sim N(0, 1), \left\| \frac{1}{n} \sum_{j=1}^n X_{ij}^2 \right\|_F \leq \frac{1}{25}, n \geq 5^2 \log d \text{ without noise}).$

Proof: $\theta \in \text{null}(X)$, $\left\| \frac{X \theta_S}{\sqrt{n}} \right\|_2^2 = \theta_S^T \frac{X^T X}{n} \theta_S \geq -\theta_S^T (I - \frac{X^T X}{n}) \theta_S + \left\| \theta_S \right\|_2^2$

$$\begin{aligned} \left\| \frac{X \theta_S}{\sqrt{n}} \right\|_2^2 &= \left\langle \frac{X \theta_S}{\sqrt{n}}, -\frac{X \theta_S^c}{\sqrt{n}} \right\rangle = -\theta_S^T \frac{X^T X}{n} \theta_S^c \\ &= -\theta_S^T (I - \frac{X^T X}{n}) \theta_S^c + \theta_S^T \theta_S^c \\ &\leq \left\| \theta_S \right\|_2 \left\| \theta_S^c \right\|_2 \delta_{PW} \end{aligned}$$

$\geq \left\| \theta_S \right\|_2^2 - \left\| I - \frac{X^T X}{n} \right\|_F \left\| \theta_S \right\|_2^2$
 $\geq \left\| \theta_S \right\|_2^2 (1 - \delta_{PW} \delta)$
 $\stackrel{\text{use } \frac{\delta}{1-\delta} \leq \frac{1}{25}}{\leq} \frac{\left\| \theta_S \right\|_2^2}{\max_{j \in [d]} |\theta_S^j|}$

Combine upper bound & lower bound

$$\left\| \theta_S \right\|_2 \leq \frac{\sqrt{s} \left\| \theta_S \right\|_2 \delta \left\| \theta_S^c \right\|_1}{1 - \delta \delta}$$

$$\left\| \theta_S \right\|_1 \leq \sqrt{s} \left\| \theta_S \right\|_2 \leq \frac{\delta s}{1 - \delta} \left\| \theta_S^c \right\|_1, \quad \left\| \theta_S \right\|_1 \leq \left\| \theta_S^c \right\|_1$$

$\delta \leq \frac{1}{25}$

Restricted Isometry Prop (RIP) $\not\in C(S)$

$$\left\| \frac{X^T X}{n} - I \right\|_F \leq \delta_S, \text{ for all } |S| \leq s$$

$\rightarrow \delta_{2s} \leq \delta_s$. Then RN holds for all $|S| \leq s$

$$\text{e.g. } X_{ij} \sim N(0, 1), \delta_S \leq \sqrt{\frac{s + \log d}{n}} \leq \sqrt{\frac{\log d}{n}}, n \geq 5 \log d$$

Recap, $y = X\theta^*$, $X \in \mathbb{R}^{n \times d}$ and $d > n$

1. Basis pursuit with $\|\theta\|_1$, θ^* is sparse. (motivation)
finds uniquely
2. Restricted null space $C(S) \cap \{X\theta \mid \|\Delta_S \theta\|_1 \leq \|\Delta_S \theta^*\|_1\}$ supports
 RNS holds if $C(S) \cap \text{null}(X) = \{0\}$
3. $RN \& BP$. $\exists S$ RN holds for some S
(2). BP uniquely recovers any θ^* with support S
(1) \Leftrightarrow 2
4. Certify RN . (a). Parseval incoherence $\left\| \frac{X^T X - I}{n} \right\|_F \leq \delta_{PW}$
If $(n \geq S^2 \log d)$ $\delta_{PW} \leq \frac{1}{3S}$ then $RN(S)$ holds for any
(b). Restricted isometry $RIP(S)$ holds if $|S| \leq s$
 $\left\| \frac{X_S^T X_S - I}{n} \right\|_F \leq \delta_{RIP, s} \quad \forall |S| \leq s$
If $\delta_{RIP, s} \leq \frac{c}{3}$ then $RN(S)$ holds for any $|S| \leq s$
 $\rightarrow n \geq S \log d$

$\theta \in \text{null}(X)$. want to show $\theta \notin C(S)$, $\|\theta_S\|_1 < \|\theta_S^*\|_1$, it's

$$\theta = \theta_0 + \theta_1 + \dots + \theta_n$$

(largest entry) θ_0 (others are 0)

(second largest entry) θ_1

(third largest entry) θ_2

$\left\| \frac{X \theta_0}{\sqrt{n}} \right\|_2^2 = \theta_0^T \frac{X^T X}{n} \theta_0 = (\|\theta_0\|_2^2 + \theta_0^T \left(\frac{X^T X - I}{n} \right) \theta_0 \geq \|\theta_0\|_2^2 - \delta \|\theta_0\|_2^2 \geq (-\delta) \|\theta_0\|_2^2$

$$\left\| \frac{X \theta_0}{\sqrt{n}} \right\|_2^2 = -\theta_0^T \frac{X^T \theta_0^C}{n}, (X \theta_0 = -X \theta_0^C)$$

$$= -\theta_0^T \left(\frac{X^T - I}{n} \right) \sum_{j \neq 0} \theta_j \leq \delta \|\theta_0\|_2 \sum_{j \neq 0} \|\theta_j\|_2 \leq \frac{\delta \|\theta_0\|_2}{\sqrt{n}} \sum_{j \neq 0} \|\theta_j\|_1$$

$$\Rightarrow \|\theta_0\|_1 \leq \sqrt{s} \|\theta_0\|_2 \leq \frac{\sqrt{s} \|\theta_0\|_2}{\sqrt{n}} \sum_{j \neq 0} \|\theta_j\|_1 = \frac{\delta}{\sqrt{n}} \left(\|\theta_0\|_1 + \|\theta_0\|_1 \right)$$

Now: $y = X\theta^* + w$, w mean 0, sub-Gaussian (X : nxd, $d > n$)

1. Relaxed BP: min $\|\theta\|_1$
 $\frac{1}{2n} \|y - X\theta\|_2^2 \leq \beta$

2. Lagrangian LASSO $\min_{\theta} \frac{1}{2n} \|y - X\theta\|_2^2 + \lambda \|\theta\|_1$

3. $\min_{\theta} \frac{1}{2n} \|y - X\theta\|_2^2, \|\theta\|_1 \leq R$ constrained LASSO

1. Estimation error $\|\hat{\theta} - \theta^*\|_2$

2. Prediction error: in-sample $\frac{1}{n} \|X\hat{\theta} - X\theta^*\|_2^2$

3. Support recovery: $S \sim \text{supp}(\theta^*)$

4. Inference: pick coordinate j , C_d $P(\hat{\theta}_j^* \in C_d) \geq 1 - \alpha$, $C_d \subset \mathbb{R}$

Restricted eigenvalue cond. $\exists \lambda_1, C_d(S) \leq \|\Delta_S^c\|_1 \leq \lambda_1 \|\Delta_S^c\|_1$ Cone.

RE(LC) $\|\frac{X\Delta}{\sqrt{n}}\|_2^2 \geq k \lambda_{\min}, \forall \Delta \in C_d(S)$ if Fisher Infor. matrix is ~~carried~~ well
stronger than RN(S)

1. Estimation error $\|\hat{\theta} - \theta^*\|_2$. Constrained LASSO, $R = \|\theta^*\|_1$

$$\Delta = \hat{\theta} - \theta^* \\ \|\theta^*\|_1^2 / \|\Delta + \theta^*\|_1 = \|\Delta + \theta^*\|_1 / (\|\Delta\|_1 + \|\theta^*\|_1) \geq \|\Delta\|_1 / (\|\Delta\|_1 + \|\theta^*\|_1)$$

$$B.T.: \frac{1}{2n} \|Y - X\hat{\theta}\|_2^2 \leq \frac{1}{2n} \|Y - X\theta^*\|_2^2 \Rightarrow \|\Delta_S^c\|_1 \leq \|\Delta_S^c\|_1$$

$$\frac{K\|\Delta\|_2^2}{2} \leq \frac{1}{2n} \|X^T w\|_2^2 \leq \frac{1}{2n} \|X\Delta\|_2^2 \leq \frac{1}{n} \|X^T w\|_2 \Rightarrow \|\Delta\|_2 \leq 4 \frac{\|X^T w\|_2}{n} \frac{\sqrt{K}}{2}$$

Formal. $R = \|\theta^*\|_1, X \rightarrow P_B(1, K)$

$$\text{then } \|\Delta\|_2 \leq 4 \frac{\|X^T w\|_2}{n} \frac{\sqrt{K}}{2}, \text{ (use } \|\cdot\|_1 \leq \|\cdot\|_2 \text{)}$$

X : column normalized, $\|X_j\|_2 \leq c, \frac{\|X^T w\|_2}{\sqrt{n}} \sim \sqrt{N(0, \sigma_n^2)}, \frac{\|X^T w\|_2}{\sqrt{n}} \sim \frac{\|X^T w\|_2}{\sqrt{S_{\text{cond}}}}$

$$\text{Low-dim. } \frac{\|X\Delta\|_2^2}{2n} \leq \frac{1}{n} \|X^T w\|_2^2$$

$$\text{and } \frac{\|X\Delta\|_2^2}{n} \leq \frac{\|X^T w\|_2^2}{n} \leq \frac{\|\Delta\|_2^2}{n} \frac{\|X^T \Sigma\|_2}{\sqrt{n}}$$

$$\Rightarrow \|\Delta\|_2 \leq \sqrt{n}$$

Recomp. (i) d.

1. BP method, If $\|x\theta^*\|_2 \leq \sqrt{3}$ then RNN holds $\|\theta\|_1 \leq \epsilon$
 $y = X\theta$

2. LASSO. $y = X\theta^* + w$ \rightarrow Lagrangian $\hat{\theta} = \arg \min_{\theta} \frac{1}{2n} \|y - X\theta\|_2^2 + \lambda \|\theta\|_1$
 Constrained $\hat{\theta} = \arg \min_{\theta} \frac{1}{2n} \|y - X\theta\|_2^2$

Relaxed BP. $\hat{\theta} = \arg \min_{\theta} \|\theta\|_1, \frac{1}{2n} \|y - X\theta\|_2^2 \leq b$

3. Restricted eigenvalue (Q, k) $\frac{X^T X}{n} \geq I_k$, Δ eigen vector.
 $K \|\Delta\|_2^2 \leq \frac{\|X\Delta\|^2}{n}, \forall \Delta, \|\Delta_S\|_1 \leq 2\|\Delta\|_1, \forall S \subseteq S$

4. Result for constant LASSO. If $X \sim RF(k, k)$. $R = \|\theta^*\|_1$,

$$\text{Then } \frac{\|\Delta\|_2}{\|\theta^*\|_2} \leq \left\| \frac{X^T w}{n} \right\|_2 \leq \frac{\sqrt{k}}{R}$$

$\downarrow \sqrt{\log n}$

Lagrangian LASSO. $\lambda \geq \frac{\alpha}{n} \|X^T w\|_1$. $X \sim RF(k, k)$, then $\|\Delta\|_2 \leq \frac{\lambda \sqrt{k}}{R}$

$$\frac{1}{2n} \|y - X\hat{\theta}\|_2^2 \leq \frac{1}{2n} \|y - X\theta^*\|_2^2 + \lambda \|\hat{\theta} - \theta^*\|_1, \Delta = \hat{\theta} - \theta^*$$

$$\begin{aligned} \frac{\|\Delta\|_2^2}{2} &\leq \frac{1}{2n} \|X\Delta\|_2^2 \leq \frac{1}{n} \|X^T w\|_2 \|\Delta\|_1 + \lambda \|\hat{\theta}\|_1 - \|\theta^*\|_1 - \|\Delta_S\|_1 - \|\Delta_{S^c}\|_1 \\ &\leq \left\| \frac{X^T w}{n} \right\|_2 \|\Delta\|_1 + \lambda \|\Delta_S\|_1 + \|\Delta_{S^c}\|_1 \\ &\leq \frac{\lambda}{2} \left\{ 3\|\Delta_S\|_1 + \|\Delta_{S^c}\|_1 \right\} \end{aligned}$$

$$\text{Upper condition } \frac{\|\Delta_S\|_1}{\|\Delta_{S^c}\|_1} \leq 3 \|\Delta_S\|_1, \Rightarrow \frac{K \|\Delta\|_2^2}{\sum} \leq \frac{3\lambda}{2} \|\Delta\|_1 \leq 3\lambda \sqrt{k} \|\Delta\|_2$$

Get $\lambda \geq \frac{\alpha}{n} \|X^T w\|_1$, noise b large enough, eigen vector \rightarrow sparse

RPB. $b^2 \geq \frac{\|w\|^2}{2n}$ $RF(k, k)$, then $\|\Delta\|_2 \leq \left\| \frac{X^T w}{n} \right\|_2 \frac{\sqrt{k}}{R} + \frac{1}{R} \sqrt{b^2 + \frac{\|w\|^2}{2n}}$

$$\frac{1}{2n} \|y - X\theta\|_2^2 \leq \|\hat{\theta}\|_1 \leq \|\theta^*\|_1 \leq \|\Delta_S\|_1 + \|\Delta_{S^c}\|_1$$

$\hookrightarrow \theta^*$ is feasible

$$\frac{1}{2n} \|y - X\hat{\theta}\|_2^2 \leq \frac{1}{2n} \|y - X\theta^*\|_2^2 + b^2 - \frac{\|w\|^2}{2n}$$

$$\frac{1}{2n} \|X\Delta\|_2^2 \leq \left\| \frac{X^T w}{n} \right\|_2 \|\Delta\|_1 + b^2 - \frac{\|w\|^2}{2n}$$

$$(\|\Delta\|_1 \leq 2\|\Delta_S\|_1 \leq 2\beta \|\Delta\|_2) \downarrow$$

$$\frac{k\|\theta\|^2}{n} \leq 2\sqrt{n} \|X^T w\|_2 \|\Delta\|_2 + b^2 \frac{\|w\|^2}{n}$$

Prediction error: $\frac{\|X\hat{\theta} - X\theta^*\|_2^2}{n}$

$$\frac{1}{n} \|X\hat{\theta} - X\theta^*\|_2^2 \leq \frac{1}{n} \|X(X\theta^* - \Delta)\|_2^2 + \lambda \|\theta^*\|_1 - \|\theta\|_1$$

$$\frac{1}{n} \|X\Delta\|_2^2 \leq \frac{\|X^T w\|_2}{n} \|\Delta\|_1 \times 2 + 2\lambda \|\theta^*\|_1 - \|\theta\|_1$$

$$(p < \lambda \geq \frac{\|X^T w\|_2}{n})$$

$$\leq \lambda (\|\theta^*\|_1 + \|\theta\|_1 + 2\|\theta^*\|_1 - 2\|\theta\|_1) \quad (\|\theta\|_1 \leq 3\|\theta^*\|_1)$$

$$\leq 3\lambda \|\theta^*\|_1$$

$$\Rightarrow \frac{1}{n} \|X\Delta\|_2^2 \leq \lambda \|\theta^*\|_1 \leq \|\theta^*\|_1 \sqrt{\frac{\log d}{n}}$$

$$X \sim RE, \theta^* = \sigma, \text{sparse: } \|\Delta\|_2^2 \leq \frac{\|X\Delta\|_2^2}{nK}$$

$$\frac{1}{n} \|X\Delta\|_2^2 \leq \lambda \sqrt{\|X\Delta\|_2} \leq \sqrt{\frac{\lambda}{nK}} (\|X\Delta\|_2 / \sqrt{n})$$

$$\leq \frac{\lambda}{\sqrt{K}} \leq \frac{\text{sluggish}}{n\sqrt{K}} \rightarrow \text{fast rate}$$

$$Y = X\theta^* + w, \theta_j^*, \hat{\theta} = (X^T X)^+ X^T Y = \theta^* + (X^T X)^+ X^T w$$

$$w \sim N(0, \sigma^2 I)$$

$$\theta_j^* \sim N(0, (X^T X)^+ (X^T X)^T \sigma^2)$$

$$\theta_j^* \pm \sigma \sqrt{\frac{1}{n} \sum_{i=1}^n \hat{\theta}_i^2}$$

LASSO: not work, sparse, biased. related to λ .

Recap.

1. LASSO theory. $\hat{\theta}^*$ is s -sparse.

- Estimation error X satisfies RE(k, 3)

$$\|\hat{\theta} - \theta^*\|_2 \leq \sqrt{\frac{\log p}{n k}} \quad \text{column normalization.}$$

2. Prediction error. - slow rate - X satisfies col norm

$$\frac{\|X\hat{\theta} - X\theta^*\|_2^2}{n} \lesssim \sqrt{\frac{\log p}{n}} \|\theta^*\|_1$$

- Fast rate $X + RE + \theta^*$ is s -sparse

$$\frac{\|X\hat{\theta} - X\theta^*\|_2^2}{n} \lesssim s \frac{\log p}{n}$$

Debiasing the LASSO, $y = X\theta^* + w$. CI for θ_j^* , $w \sim N(0, \sigma^2_\epsilon)$

$$\hat{\theta} : \text{LASSO est.} \quad \hat{\theta} = \hat{\theta} + M X^T (y - X\hat{\theta}), \quad M \approx \left(\frac{X^T X}{n}\right)^{-1}$$

$$\text{debias: } \tilde{\theta} = \hat{\theta} + M X^T \frac{(y - X\hat{\theta})}{n} + M X^T X \theta^* + M X^T w$$

$$\therefore \tilde{\theta} = \theta^* - M X^T \frac{\hat{\theta}}{n} + M X^T \frac{X \theta^*}{n} + M X^T \frac{w}{n} \\ \approx \theta^* + \text{noise} \quad (\text{put } M = \left(\frac{X^T X}{n}\right)^{-1} \text{ into it})$$

$$\tilde{\theta} = \theta^* + \left(I - M \frac{X^T X}{n}\right) (\hat{\theta} - \theta^*) + M \frac{X^T w}{n}$$

$$\tilde{\theta} \stackrel{d}{=} N(\theta^*, \frac{M X^T X M^T \sigma^2_\epsilon}{n^2}) \text{ i.i.d. var}$$

$$\text{err} \leq \left\| I - \frac{M X^T X}{n} \right\| \approx \|\hat{\theta} - \theta^*\|_1 \lesssim \sqrt{\frac{\log p}{n}} \sqrt{s} \underbrace{\sqrt{\frac{\log p}{n}}}_{\approx \frac{\log p}{\sqrt{n}}} \underbrace{\|\hat{\theta} - \theta^*\|_1}_{\|\hat{\theta} - \theta^*\|_1 \leq \sqrt{\frac{\log p}{n}} \| \theta^* \|_1} \quad \begin{matrix} \text{under RE + col norm + } \theta^* \\ \text{s-sparse} \end{matrix}$$

$$N(0, 1) \quad \text{then } M = I$$

if $\frac{\log p}{n} \ll \frac{1}{m}$, then CI is correct

Support recovery. $y = X\theta^* + w$. $\text{supp}(\theta^*)$

$$\text{supp}(\hat{\theta}) = \text{supp}(\theta^*) \text{ under cond.}$$

$$\|\hat{\theta} - \theta^*\|_2 \leq \sqrt{\frac{\log p}{n}} \rightarrow \|\hat{\theta} - \theta^*\|_\infty \leq \sqrt{\frac{\log p}{n}} \text{ will recover any } |\theta_j^*| \geq \sqrt{\frac{\log p}{n}}$$

Prob cond. $\{|\theta_j^*| \text{ non-zero} > \sqrt{\frac{\log p}{n}}\} \Rightarrow \hat{\theta}_j \text{ will be non-zero}$

+ threshold $\hat{\theta} \rightarrow \text{at } \sqrt{\frac{\log p}{n}}$

(Imagine $X = I$, $N(T(x) \sqrt{\frac{\log p}{n}})$
will recover $\text{supp}(\theta^*)$ if $|\theta_j^*| > \sqrt{\frac{\log p}{n}}$)

Mutual incoherence

$$\max_{j \in S^c} \|(\mathbf{X}_S^T \mathbf{X}_S)^{-1} \mathbf{X}_S^T \mathbf{x}_j\|_1 < 1, \text{ product } \mathbf{x}_j \text{ from } \mathbf{X}_S \\ \|\beta_j\|_1 < 1$$

If $\beta_m \neq 0$: Mutual incoherence \Rightarrow support recovery

$$\text{KKT of LASSO: } \vec{\mathbf{x}} = (\mathbf{y} - \mathbf{X}\theta) + \lambda \text{sign}(\theta), \\ \frac{1}{2} \mathbf{x}^T (\mathbf{y} - \mathbf{X}\theta) + \lambda \text{sign}(\theta) = 0 \quad \text{sign}(\theta) = \begin{cases} 1 & \theta > 0 \\ 0 & \theta = 0 \\ -1 & \theta < 0 \end{cases}$$

$\hat{\theta}$ has support S

$(\hat{\theta}, \text{sign}(\theta))$ which solves KKT

primal
sol
dual
sol

Minimax lower bounds.

$$\theta(P), \hat{\theta}, \mathbb{E}_{\theta} P(\hat{\theta}, \theta(P))$$

semi-metric

minimax rate

$$SN(P, P) = \inf_{\hat{\theta}} \sup_{P \in \mathcal{P}} \mathbb{E}_P P(\hat{\theta}, \theta(P))$$

From estimation to testing, $\{\theta^1, \dots, \theta^M\}$, $P(\theta^i, \theta^j) \geq 2\delta$

testing error $\underbrace{\mathbb{P}(\psi(z) \neq j)}_{\psi(z) \rightarrow \hat{M}}$ $\xrightarrow{j \sim \text{unif}(M), z \sim p_{0j}} \text{find which } \theta^M?$ (testing)

$$\mathbb{P}(z | J=j) = p_{0j}, \quad Q(J=j) = \frac{1}{M}$$

$$\mathbb{P}(\psi(z) \neq j) \leq \inf_{\psi} \mathbb{P}(\psi(z) \neq j)$$

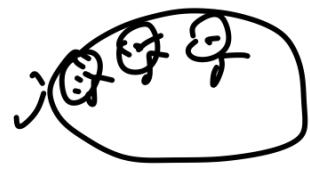


$$\sup_{P \in \mathcal{P}} \mathbb{E}_P \mathbb{P}(\psi(\theta, \theta(P)) \geq \psi(\theta)) \geq \mathbb{P}(\psi(\theta, \theta(P)) \geq \psi(\theta))$$

$$\geq \sup_{P \in \mathcal{P}} \mathbb{P}(\psi(\theta, \theta(P)) \geq \delta) \geq \mathbb{P}(\psi(\theta, \theta(P)) \geq \delta) \geq \sum_{i=1}^M \mathbb{P}_{0i} (P(\theta^i, \theta) \geq \delta)$$

$$\hat{\theta}, \hat{\psi}(\hat{\theta}) = \arg \min_{\theta} P(\theta, \theta)$$

If $c(\theta, \theta^j) \leq \delta$, $\Rightarrow \varphi(z) = j$
 $\dots \geq \underline{c}(\theta) \cdot R(\psi(z) \neq j)$
 $\inf_{\theta} \sup_{\rho} R \geq \underline{c}(\theta) \inf_{\rho} R(\psi(z) \neq j)$



Recap.

1. LASSO theory:

$$\text{est. error: } \|\hat{\theta} - \theta^*\|_2 \leq \sqrt{\frac{\log p}{n}} \text{ under } R\hat{\theta} + \theta^* + S\text{-sparse}$$

$$\text{pred. error: } \frac{\|X\hat{\theta} - X\theta^*\|_2^2}{n} \leq \frac{\log p}{n} \leq \sqrt{\frac{\log p}{n}} \|\theta^*\|_1$$

Variable selection / supp recovery. under col-norm.

$$\max_i \|(\mathbf{X}_S^T \mathbf{X}_S)^{-1} \mathbf{X}_S^T \mathbf{x}_j\|_q \leq 1$$

$$\text{Prmt } (\theta^* \text{ nh } \geq \sqrt{\frac{\log p}{n}} \text{ s.t. } \text{supp}(\hat{\theta}) = \text{supp}(\theta^*)) \rightarrow 1$$

Def. bias.

$$\text{if } \|A - M \frac{\mathbf{X}^T \mathbf{X}}{n}\|_2 \|\hat{\theta} - \theta^*\|_1 = O(\sqrt{n})$$

then. $\hat{\theta} = \theta^* + \frac{M \mathbf{X}^T (y - A\theta^*)}{n}$ is asymptotically mean θ^* Gaussian

Min Max.

$$E_{\theta} \Phi_P(\theta, \theta(P)) \quad M(P, \hat{\theta}, \theta) = \inf_{\theta} \sup_{P \in \mathcal{P}} R(\hat{\theta}, \theta(P))$$

$$= P(\hat{\theta} | \theta(P)) \quad \left\{ \theta^1, \dots, \theta^M \right\} \rightarrow P(\theta^i, \theta) \geq 2 \sum \underbrace{\min_{j \neq i} \theta_j}_{\text{testing error}} \cdot \theta^i$$

$$M \geq \overline{\Phi}(P) \inf_{\theta} Q(\theta | \theta \neq \bar{\theta}) \text{ testing error; lower bound}$$

Divergence measures.

$$1. TV(P, Q) = \sup_A |P(A) - Q(A)| = \frac{1}{2} \int |P(x) - Q(x)| dx \text{ of est. error}$$

$$2. KL(P|Q) = E_{x \sim P} \log \frac{P(x)}{Q(x)}$$

$$3. H^2(P, Q) = \frac{1}{2} \int (\sqrt{P(x)} - \sqrt{Q(x)})^2 dx$$

$$\text{Properties: } H^2 \leq TV \leq H \leq \sqrt{KL}$$

$$\text{Tensorization: } KL(P^n, Q^n) = n KL(P, Q) \quad x_1, \dots, x_n \text{ i.i.d.} \quad P, Q \text{ i.i.d.}$$

$$H^2(P^n, Q^n) \leq n H^2(P, Q)$$

$$H^2(P, Q) = 1 - \int \sqrt{PQ} \text{ Hellinger affinity}$$

$$H^2(P^n, Q^n) = 1 - \left(\int \sqrt{PQ} \right)^n = 1 - (1 - H^2)^n \leq n H^2$$

Mutual information:

$$I(X; Y) = KL(p(XY), p(X)p(Y)) \rightarrow \int_{x_1} \int_{x_2} \int_{x_3} \sqrt{p(x_1)q(x_1|x_2)} \sqrt{p(x_2)q(x_3|x_2)} dx_1 dx_2 dx_3$$

Le Cam's $M \geq \underline{I}(\theta) \text{ inf } Q(\theta_1, \theta_2)$

$$\underline{I} = \frac{1}{2} P_{\theta_1} + \frac{1}{2} P_{\theta_2} \geq \frac{\underline{I}(\theta)}{2} (-TV(P_{\theta_1}, P_{\theta_2}))$$

$X_1, \dots, X_n \sim N(\theta, \sigma^2)$, $R(\theta, \theta^*) \geq \frac{\sigma^2}{n}$

$$\theta_1, \dots, \theta_1 + 2\delta \quad R(\theta_1, \theta^*) = \underline{I}(\theta - \theta^*)$$

$$M \geq \delta \left(1 - TV(N(\theta_1, \sigma^2), N(\theta_1 + 2\delta, \sigma^2)) \right) \rightarrow \theta_1, \dots, \theta_1 + 2\delta \text{ i.i.d.}$$

$$TV \leq k \sqrt{n} K \sqrt{N(\theta_1 + 2\delta, \sigma^2)} \quad TV \leq \sqrt{\frac{n}{2}} \frac{\sigma \delta^2}{\sigma^2}$$

$$\Rightarrow \text{Le Cam: } M \geq \delta \left(1 - \sqrt{\frac{2\sigma^2}{\sigma^2}} \right), \delta = \frac{\sigma}{\sqrt{n}}, M \geq \frac{\sigma}{\sqrt{n}}$$

$X_1, \dots, X_n \sim U(\theta, \theta + 1)$

$$\hat{\theta}_{MLE} = \min_{i=1}^n X_i, |\hat{\theta}_{MLE} - \theta| \lesssim \frac{1}{n}$$

$$\theta_1: U(\theta, \theta + 1), \theta_{1+2\delta}: U(\theta + 2\delta, \theta + 2\delta + 1)$$

$$H^2(P_{\theta_1}, P_{\theta_1 + 2\delta}) = \frac{1}{2}(4\delta) = 2\delta$$

$$H^2(P^n, P^n) \leq 2^n \delta$$

$$M \geq \frac{\delta}{2} \left(1 - \sqrt{2^n \delta} \right), \delta \approx \sqrt{n}, M \geq \frac{1}{\sqrt{n}}$$

Lipschitz fkt at a point $x_1, \dots, x_n, y_1, \dots, y_n \sim U(0,1)$

$$g_i = f(x_i) + \sum_{j \neq i} \cdot |f(x_j) - f(y_j)| \leq 4|x-y|$$

$$\theta = f(0, \cdot)$$

$$P_1, f_1 = 0 \quad P_2: f_2 = \frac{\Delta^{1/4}}{4\pi}$$

$$KL(P_1(x, y), P_2(x, y)) = E_P \log \frac{P_1(y|x)}{P_2(y|x)} = \int_x P_1(x) f((P_1(y|x), P_2(y|x)))$$

$$= \int_x P_1(x) (f_1(x) - f_2(x))^2$$

Recap.

1. Est. of testing

$$M(\Phi \circ \rho, P) = \inf_{\theta \in \Theta} \mathbb{E}_{P \in \mathcal{P}} \Phi \circ \rho(\theta, \theta_P)$$

$\rightarrow M \geq \frac{1}{n} \sum_{j=1}^n I(\psi(z_j) \neq j), \theta_1, \dots, \theta_m$ are separated.
 $J \sim \text{unif}[M], z_j | J \sim P_{\theta_j}$

3. Law of Large Numbers

$$M \geq \frac{\Phi(1)}{2} \left(1 - TV(P_{\theta^1}, P_{\theta^2})\right)$$

$$4. \text{Req. } H^2 \leq TV \leq H \leq \sqrt{KL} \leq \sqrt{X^2}$$

$$KL(P^n, Q^n) = n KL(P, Q), H^2(P^n, Q^n) \leq n H^2(P, Q)$$

5. Examples. Unif Normal Location $\geq \frac{G}{T\sqrt{n}}$
 Uniform location $\geq \frac{1}{\sqrt{n}}$. Lipschitz fn. est. at const $\geq \frac{f}{\sqrt{n}}$

$$\text{Fano's inequality: } \inf_{\psi} I(\psi(z) \neq j) \geq 1 - \frac{I(z, J)}{\log M} + \frac{\log 2}{M} \quad J \sim [M]$$

$$I(z, J) = KL(Q(z, J) || Q(z|J))$$

$$\left(\frac{1}{n} \sum_{j=1}^n P_{\theta_j} \right) = \mathbb{E}_{J \sim [M]} \left(\log \frac{Q(z|J) Q(z|J)}{Q(z)} \right) = \frac{1}{n} \sum_{j=1}^n \mathbb{E}_{z|J} \log \frac{Q(z|J)}{Q(z)}$$

$$= \frac{1}{n} \sum_{j=1}^n KL(P_{\theta_j} || Q)$$

convexity of KL ($\lambda p_1 + (1-\lambda)p_2 || \lambda q_1 + (1-\lambda)q_2 \leq \lambda KL(p_1 || p_2) + (1-\lambda)KL(q_1 || q_2)$)

$$I(z, J) \leq \frac{1}{M^2} \sum_{j=1}^M \sum_{k=1}^n KL(P_{\theta_j} || P_{\theta_k}) \Rightarrow \left(\frac{1}{M} \sum_{k=1}^n P_{\theta_k} \right) \leq \sup_{\theta_j \in \Theta} KL(P_{\theta_j} || P_{\theta_k})$$

$$X_1, \dots, X_M \sim N(\theta, \sigma^2 I_d) \quad (\text{local packings})$$

$$KL(P_{\theta_j} || P_{\theta_k}) = \frac{n(\theta_j - \theta_k)^2}{2\sigma^2} \leq \frac{2n\mu^2}{\sigma^2}$$

$$\log M \geq \log \left(\frac{1}{\delta} \right)^d \geq d \log \left(\frac{1}{\delta} \right) \geq d.$$

$$\Rightarrow \text{Let } \frac{2n\mu^2}{\sigma^2} \leq d \Rightarrow \frac{2n\mu^2}{\sigma^2} \Rightarrow M \geq \frac{\sigma^2 d}{n}$$



$$KL(C) = \frac{(1-\theta)\ln(1-\theta)}{2\sigma^2}$$

Lipschitz f.n. est. with ℓ_2 loss

$X \sim \text{Uniform}(0,1)$ $y = f(x) + \epsilon$



$$\Phi(x) = \begin{cases} x, & x \in [0, \frac{1}{n}] \\ 1-x, & x \in (\frac{1}{n}, 1] \end{cases}$$

$$\varphi_j = h \Phi\left(\frac{x-j}{h}\right)$$

$$f_{\alpha}(x) = \sum_{j=0, h \geq 1}^{\lfloor \frac{1}{h} \rfloor} w_j \varphi_j(x), w_j \in \mathbb{R}$$

$$KL(P_{w_1}, P_{w_2}) = \frac{n \int (f_{w_1}(x) - f_{w_2}(x))^2 dx}{2\sigma^2} \lesssim \frac{n L^2 h^3}{\sigma^2} \times d_H(w_1, w_2)$$

\downarrow Hamming distance.
number of disagree
entries

$$L \cdot P(f_{w_1}, f_{w_2}) \gtrsim L^2 h^3 d_H(w_1, w_2)$$

Vardamov-Gilbert bound.

$$\exists \Sigma \subseteq \{1, 0\}^d, w_1, w_2 \in \Sigma$$

$$KL \lesssim \frac{n L^2 h^3}{2\sigma^2} \text{ and } \Sigma^2 \gtrsim L^2 h^3 \times \frac{1}{h} \times \frac{1}{h} \times h^2$$

$$\log M \gtrsim \frac{1}{h}$$

$$\Rightarrow n \frac{L^2 h^2}{\sigma^2} \lesssim \frac{1}{h}, h \lesssim n^{-\frac{1}{3}}, h^2 \gtrsim n^{\frac{2}{3}}$$

Recap

1. Fano's inequality
 $\inf_{\Psi} \mathbb{P}(\Psi(Z; J)) \geq 1 - \frac{I(Z; J) + \log 2}{\log M}$
 $\Rightarrow: P_{\theta^1}, \dots, P_{\theta^M}$ such that $P(\theta^i, \theta^j) \geq \delta$, $I(Z; J) \leq \log M$
 then $M \geq \underline{\Phi}(\delta)$

Boundary $I(Z; J)$ mutual information

$$\begin{aligned} I(Z; J) &= \frac{1}{M} \sum_{j=1}^M KL(P_{\theta^j} \| \frac{1}{M} \sum_{i=1}^M P_{\theta^i}) \\ &\leq \frac{1}{M} \sum_j \sum_k KL(P_{\theta^j} \| P_{\theta^k}) \\ &\leq \sup_{J, k} KL(P_{\theta^j} \| P_{\theta^k}) \quad (\text{local packings}) \end{aligned}$$

\circlearrowleft

Varshamov-Gilber:

$$\begin{aligned} S_2 &\subseteq \{-1, 1\}^d, |S_2| \geq 2^{\frac{d}{2}} \\ &\& d_{HTW}(V, W) \geq \frac{d}{2} \quad \text{if } V \neq W \\ &\& d(V, W') \in \mathbb{N} \end{aligned}$$

$\begin{aligned} &\text{if } P(\theta^j, \theta^k) \neq 0 \text{ risk} \\ &KL(P_{\theta^j}, P_{\theta^k}) \text{ risk} \\ &\text{if } \lambda \text{ lower bound} \\ &\text{of packmax.} \end{aligned}$

Eg: $X_1, \dots, X_n \sim N(\theta^*, \sigma^2 I)$
 Support recovery \rightarrow Est. $\|\theta\|_2 \leq S \rightarrow$ Est. $\|\theta\|_2 \leq R$

1) hard thresholding. $HT_\lambda(\lambda) = \sqrt{\frac{\log d}{n}}$

$$\hat{\theta} = HT_\lambda(\bar{x}), \text{supp}(\hat{\theta}) = \text{supp}(\theta^*) \text{ if } |\theta^*|_{\min} \geq \sqrt{\frac{\log d}{n}}$$

$$\begin{aligned} \mathbb{E}[\|\hat{\theta} - \theta^*\|_2^2] &\leq \frac{s \log d}{n} \\ &\leq \Omega R \sqrt{\frac{\log d}{n}} \end{aligned}$$

$\mathbb{P}\{\text{supp}(\hat{\theta}) \neq \text{supp}(\theta)\}, M = \inf_{\hat{\theta}} \sup_{\theta^*} \mathbb{P}(\text{supp}(\hat{\theta}) \neq \text{supp}(\theta^*)), \text{min max risk}$

$$\theta^* = \begin{pmatrix} \theta^*_{\min} \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \theta^j = \begin{pmatrix} 0 \\ \theta^*_{\min} \\ \vdots \\ 0 \end{pmatrix}, \dots, \delta = 1, \mathbb{P}(\theta^j, \theta^j) \geq 1$$

$$KL(P_{\theta^j} \| P_{\theta^2}) = \frac{2n \theta_{\min}^2}{2\sigma^2} \Rightarrow \frac{n \theta_{\min}^2}{(KL)^2} \leq \frac{\log d}{\log n}, M \geq C$$

lower bound min max

$$2) \alpha, w \in \{1, 0, -1\}^d, \theta^i = \alpha w^i \stackrel{(+) \rightarrow}{\text{if } (+)} \sum_{t=1}^n (\frac{\alpha}{S})^t \rightarrow (\frac{\alpha}{S})^n$$

$$\text{Spontaneous Varnish - Gilbert model. } S \subseteq H$$

$$d_H(w_i, w_j) \geq \frac{1}{8}, w_i, w_j \in S, |S| \geq C \sum_{t=1}^n (\frac{\alpha}{S})^t$$

$$\frac{1}{8} \sum_{i,j} d_H(w_i, w_j) \geq \frac{1}{8} \sum_{i,j} \frac{1}{8} = \frac{1}{64} \sum_{i,j} 1 = \frac{1}{64} n(n-1) \geq \frac{1}{8} \sum_{i,j} 1 = \frac{1}{8} n(n-1)$$

$$\|\theta^i - \theta^j\|_2^2 \geq \frac{1}{8} \sum_{i,j} 1 = \frac{1}{8} n(n-1), KL(P_{\theta^i} || P_{\theta^j}) \leq \frac{n \alpha^2}{2 \sigma^2}$$

Let $KL \leq \log \text{of packing.}$

$$\frac{n \alpha^2}{2 \sigma^2} \leq S \log \frac{d}{S} \Rightarrow \alpha \leq \sqrt{\frac{4 \log \frac{d}{S}}{n}}, \Rightarrow M \geq \sqrt{\frac{2 \sigma^2}{n} \log \frac{d}{S}}$$

$$3) \theta^i = \frac{R}{K} w^i, \|\theta^i - \theta^j\|_2^2 = \frac{R^2}{K^2} \sum_{i,j} 1, KL(P_{\theta^i} || P_{\theta^j}) \leq \frac{n K \left(\frac{R}{K}\right)^2}{2 \sigma^2}$$

$$w \in \{-1, 0, 1\}^d \quad = \delta^2 \quad \Rightarrow KL \leq \log \frac{d}{\delta} \Rightarrow K \geq \sqrt{\frac{1}{2} \frac{R^2}{\delta^2} \log \frac{d}{\delta}}$$

$$M \geq \sqrt{\frac{2 \sigma^2}{n} \log \frac{d}{\delta}} = \sigma R \sqrt{\frac{\log \frac{d}{\delta}}{n}}$$

Yang - Barron

$$I(Z, J) \leq \sup_{i,j} KL(P_{\theta^i} || P_{\theta^j})$$

$$I(Z, J) \leq \log M$$

$$\text{YB: } I(Z, J) \leq \inf_{\Sigma \gg} \left\{ \sum_{i,j} \log \left(\frac{P_{\theta^i}}{P_{\theta^j}} \right) \right\} \quad N_{KL} = \{P_{\theta^1}, \dots, P_{\theta^M}\}^P, \text{ for any } P \in \mathcal{P}$$

$$KL(P_{\theta^i} || P_{\theta^j}) \leq \epsilon^2 \text{ for some } \epsilon$$

$\hookrightarrow \Sigma \text{ is packing of } \mathbb{R}^d, \text{ i.e. } T_A \text{ upper bound}$

Roadmap:

$$1. \text{Find } \Sigma, \text{ s.t. } \Sigma^2 \geq \log N_{KL}(\Sigma)$$

$$\text{Smallest: } I \leq 2 \Sigma^2$$

$$2. \log M \gg 4 \cdot \Sigma^2 \xrightarrow{\text{find } \delta} \text{global packing}$$

$$3. \text{Apply Fano: } M \geq \frac{1}{2} \delta^2$$

local packing (before)

$\hookrightarrow \Sigma \leq KL \text{ number}$

$$\text{Non-parametric regression: } f \in \mathcal{F}, \log(\Sigma \mathcal{F}_1 \| \mathcal{F}_2) \asymp (\frac{1}{\Sigma})^{1/2}$$

$$g = f + \epsilon \sim \text{Uniform}(0, 1), KL(P_{f+} || P_{f-}) = \frac{n \|f^+ - f^-\|_2^2}{2 \sigma^2}, \sqrt{KL} \leq \Sigma \|f^+ - f^-\|_2 / \sqrt{n}$$

$$\Rightarrow KL \text{ covering number} \quad \log N_{KL}(\Sigma) = \left(\frac{\sqrt{n}}{\sigma \Sigma} \right)^{1/2}$$

$$\Sigma^2 \geq \left(\frac{\sqrt{n}}{\sigma \Sigma} \right)^{1/2}$$

$$1. I \leq \frac{1}{n^{1/(2d+1)}} \left(\frac{1}{\delta} \right)^{d/2} \geq \frac{d}{n^{1/(2d+1)}}, \delta \leq n^{1/(2d+1)} \Rightarrow M \geq n^{\frac{2d}{2d+1}}$$

$$\text{YB: } \{P_{\theta^1}, \dots, P_{\theta^M}\}, \{P_{\theta^1}, \dots, P_{\theta^M}\}, I(Z, J) \leq \sum_{i,j} KL(P_{\theta^i} || P_{\theta^j})$$

$$\text{pick any } \theta, \text{ then } I(Z, J) \leq \frac{1}{M} \sum_{i,j} KL(P_{\theta^i} || P_{\theta^j})$$

$$= \frac{1}{M} \sum_{i,j} E_{P_{\theta^i}} \log \frac{P_{\theta^i}}{P_{\theta^j}} = \frac{1}{M} \sum_{i,j} \left(E_{P_{\theta^i}} \log \frac{P_{\theta^i}}{P_{\theta^j}} + E_{P_{\theta^j}} \log \frac{P_{\theta^j}}{P_{\theta^i}} \right) = I(Z, J) + KL(\overline{\theta} || \theta)$$

$$\begin{aligned}
 \chi^2 &= \frac{1}{N} \sum_{j=1}^N P_{Yj}, \quad I \leq \frac{1}{M} \sum_{j=1}^M KL(P_{0jj} \| \frac{1}{N} \sum_{k=1}^N P_{Yk}) \\
 &\leq \mathbb{E}_{P_{0j}} \log \frac{P_{0j}}{\frac{1}{N} P_Y} = \overline{\mathbb{E}_{P_{0j}} \log \frac{P_{0j}}{P_Y}} + \log N \\
 &\leq \Sigma^2 + \log N \overline{KL(\cdot | \cdot)}
 \end{aligned}$$

Recap:

1. Tang-Raman. $I(\mathbb{Z}, \mathcal{J}) \leq \inf_{\varepsilon > 0} \left\{ \varepsilon^2 \log N_{\mathcal{F}_L}(\varepsilon) \right\}$
 $\mathcal{S} \in \mathbb{R}^d, \mathcal{J} \subset \mathbb{R}^m, \theta^1, \dots, \theta^m$
2. Strategy. Global padding. $J - \text{Unif}(n)$
 - a. Find smallest ε that $\varepsilon^2 \geq \log N_{\mathcal{F}_L}(\varepsilon)$
 - b. Find largest δ such that $\log M(\delta) \geq 4\varepsilon^2$
 - c. Conclude $M \geq \Phi(\delta)$

$$n\varepsilon^2 = \log N(\varepsilon)$$

Non-parametric max Likelihood.

$$\mathcal{P}, p^* \in \mathcal{P}, x_1, \dots, x_n \sim p^*$$

$$\hat{p} = \arg \max_{p \in \mathcal{P}} f_N(p) = \prod_{i=1}^n P(x_i)$$

Hellinger. $h^2(\hat{p}, p^*)$?

Informal. find $n \varepsilon^{*2} \asymp \log N(\hat{p}, \varepsilon^{*2})$

$$h^2(\hat{p}, p^*) \leq e^{*2}$$

Finite collection $\{p_1, \dots, p_N\}$

Lemma. Wong-Shen. If $h^2(p, p^*) \geq \delta^2$, $P\left(\frac{\ln p}{\ln p^*} \geq \exp\left(\frac{n\delta^2}{4}\right)\right) \leq \exp\left(-\frac{n\delta^2}{8}\right)$

$$\text{Proof: } P\left(\sqrt{\frac{\ln p}{\ln p^*}} \geq \exp\left(\frac{n\delta^2}{8}\right)\right) \leq \frac{E\left(\sqrt{\frac{\ln p}{\ln p^*}}\right)}{\exp\left(\frac{n\delta^2}{8}\right)}$$

$$\int \prod_{i=1}^n \frac{1}{p^*(x_i)} p^*(x_i) = \left(\int \sqrt{p(x)p^*(x)} dx \right)^n = (1-t^2)^n$$

$$H = h(\hat{p}, p^*) \geq \delta^2 \leq \exp(-n(t^2))$$

$$\left\{ p_1, \dots, p_n \right\} \quad P \left(\frac{\ln(p)}{\ln(p^*)} \geq 1, h^2(\beta, p^*) \geq \delta^2 \right) \\ \leq P \left(\sup_{p \in P, h^2(p, p^*) \geq \delta^2} \frac{\ln(p)}{\ln(p^*)} \geq 1 \right)$$

why
\$N \exp(-n\delta^2)\$

$$\therefore \text{If } \frac{n\delta^2}{16} \geq \log N, \text{ then } P(h^2(\beta, p^*) \geq \delta^2) \leq \exp\left(\frac{-n\delta^2}{16}\right)$$

Sieve MLE 

\rightarrow Sieve Hellinger cover, i.e. for any $p \in \mathcal{P}$, $\tilde{p} \in \mathcal{N}$, s.t.

$$\text{Assumption. 1. } \exists \tilde{p} \in \mathcal{N}, \sup_{x \in X} \frac{p(x)}{\tilde{p}(x)} \leq \frac{11}{8} \quad 2. \chi^2(p^*, \tilde{p}) \leq \frac{\delta^2}{16}$$

$$\text{Proof } P \left(\frac{\ln(p)}{\ln(p^*)} \geq 1, h^2(\beta, p^*) \geq \delta^2 \right) = P \left(\frac{\ln(p)}{\ln(p^*)} \cdot \frac{\ln(p^*)}{\ln(p^*)} \geq 1, \dots \right)$$

$u \rightarrow \text{with } p \in \mathcal{N} \text{ in } H^2 \text{ and } \chi^2$
 $u \in p^* + \delta^2 \text{ in } H^2$

$$\leq P \left(\frac{\ln(p)}{\ln(p^*)} \geq \exp\left(\frac{-n\delta^2}{8}\right), h^2(p, p^*) \geq \delta^2 \right) \\ + P \left(\frac{\ln(p^*)}{\ln(p^*)} > \exp\left(\frac{n\delta^2}{8}\right) \right)$$

$$\Rightarrow P(h^2(\beta, p^*) \geq \delta^2) \leq N \exp\left(-\frac{n\delta^2}{8}\right)$$

If $n\delta^2 \geq \log N$

$$\text{then } P(h^2(\beta, p^*) \geq \delta^2) \leq \exp\left(-\frac{n\delta^2}{16}\right)$$

$$(1). H^2(p^*, u) \leq \varepsilon^2, \text{ where } \frac{p^*}{u} \in \mathcal{E} \rightarrow \chi^2(p^*, u) \leq \varepsilon^2$$

$$\leq \left(N \exp\left(-\frac{n\delta^2}{8}\right) + \exp\left(-\frac{n\delta^2}{8}\right) \left[\frac{p^* - p^*}{u} \right]^n \right)^n \\ \leq \exp\left(\frac{n\delta^2}{8}\right) \exp\left(n\chi^2(p^*, u)\right) = \exp\left(-\frac{n\delta^2}{16}\right)$$

$$\int \underbrace{\frac{(p^* - \sqrt{u})^2 (\sqrt{p^*} + \sqrt{u})^2}{u}}_{\text{why?}} \frac{du}{n} \leq 2 \int \underbrace{\left(\frac{p^*}{u} + 1 \right) (\sqrt{p^*} - \sqrt{u})^2}_{\leq 2(1 + \varepsilon^2)} du \leq 2(1 + \varepsilon^2)$$

Total variation density est.

$$X_1, \dots, X_m \sim p^* \quad H_0: p = p_1, H_1: p = p_2, A = \{p_1 > p_2\}, \hat{p} = \arg \min_{p \in \mathcal{P}(p_1, p_2)} |P(A) - P_{\hat{p}}(A)|$$

$$X_1, \dots, X_m \sim p_1, \quad P_n(A) \hat{P}_1(A) = \left| p_1(A) - \frac{1}{n} \sum_{i=1}^n I\{X_i \in A\} \right|$$

$$P_1 \left(|A| - P_n(A) \geq \sqrt{\frac{\log(1/\delta)}{n}} \right) \leq \delta$$

$$P_1(A) - P_2(A) = TV(P_1, P_2), \quad D(P) = |P(A) - P_0(A)| = |P(A) - P_{\text{true}}(A)|$$

If $TV(P, P^*) \geq \max\left\{\sum_i \sqrt{\frac{\log k_i}{n}}, \sqrt{\frac{\log k^*}{n}}\right\}$

$$P(|P - P^*| \leq \delta) \leq P\left(|P - P^*| > \sqrt{\frac{\log k^*}{n}}\right) + P\left(|P - P^*| \leq \sqrt{\frac{\log k^*}{n}}\right)$$

TV-version · Wong & Shen. $\leq \delta$ $\exp(-n\delta^2)$ Hoeffding-

Recap.

1. Nonparametric MLE
 $X_1, \dots, X_n \sim p^* \in \mathcal{P}$, $\hat{p} = \arg\max_{p \in \mathcal{P}} \prod_{i=1}^n p(X_i)$

2. Wong-Shen. If $h^2(\hat{p} - p^*) \geq \varepsilon$

$$P\left(\prod_{i=1}^n \frac{p(X_i)}{p^*(X_i)} \geq \exp\left(-\frac{n\varepsilon^2}{h}\right)\right) \leq \exp\left(\frac{-n\varepsilon^2}{8}\right)$$

3. If $n\varepsilon^2 \geq \log(N/\varepsilon) + \chi^2$ approx enough

then $P(h^2(\hat{p}, p^*) \geq \varepsilon^2) \leq \exp(-cn\varepsilon^2)$

where $\hat{p} = \arg\max_{N(\varepsilon, \hat{p}, h)} \prod_{i=1}^n p(X_i)$

$X_1, \dots, X_n \sim p^*, p_1, p_2 \Rightarrow A = \{p_1 > p_2\}$
 $\Delta p = [p_n(A) - p(A)]$, $\hat{p} = \arg\min_{p \in \mathcal{P}(p_1, p_2)} \Delta(p)$

$\{p_1, \dots, p_n\} \rightarrow Yatkov's$ method. $A_{i,j} = \{p_i > p_j\} \rightarrow A$

$$\Delta(p) = \sup_{A \in \mathcal{A}} |p_n(A) - p(A)|$$

$$\Rightarrow \hat{p} = \arg\min_{p \in \mathcal{P}(p_1, p_n)} \Delta(p)$$

$$\begin{aligned} P(TV(\hat{p}, p^*) \geq \varepsilon) &= P(\Delta(\hat{p}) \leq \varepsilon), TV(\hat{p}, p^*) \geq \varepsilon \\ &\leq P(\Delta(\hat{p}) \leq \frac{\varepsilon}{2}, TV(\hat{p}, p^*) \geq \varepsilon) + P(TV(\hat{p}, p^*) \geq \frac{\varepsilon}{2}) \\ &\stackrel{\text{union bound}}{\leq} \sum_{i=1}^n P(A_{i,j}) \leq \frac{\varepsilon}{2}, TV(\hat{p}, p^*) \geq \varepsilon \quad \xrightarrow{\text{Hoeffding}} \exp\left(-\frac{2n\varepsilon^2}{4}\right) \cdot \binom{n}{2} \\ &\leq \sum_{i=1}^n P(|p_n(A_{i,j}) - p(A_{i,j})| \leq \frac{\varepsilon}{2}, TV(\hat{p}, p^*) \geq \varepsilon) \quad \xrightarrow{\text{Hoeffding}} \binom{n}{2} \\ &\leq \sum_{i=1}^n P(|p_n(A_{i,j}) - p^*(A_{i,j})| \geq \frac{\varepsilon}{2}) \quad \xrightarrow{\text{TV}(\hat{p}, p^*) \geq \varepsilon} \\ &\leq \sum_{i=1}^n P(|p_n(A_{i,j}) - p^*(A_{i,j})| \geq \frac{\varepsilon}{2}) + \binom{n}{2} \exp(-n\varepsilon^2) \\ &\leq n \exp(n\varepsilon^2) + \binom{n}{2} \exp(-n\varepsilon^2) \end{aligned}$$

Hoeffding
 $2t \cdot n\varepsilon^2 \geq \log N$
 $\text{then } P(TV(\hat{p}, p^*) \geq \varepsilon) \leq \exp(-n\varepsilon^2), \varepsilon \geq \sqrt{\frac{\log N}{n}}$

$$\rightarrow \gamma$$

$n\varepsilon^2 \geq \log N(P, \varepsilon)$
then $\mathbb{P}(TV(\beta, p^*) \geq \varepsilon) \leq \exp(-n\varepsilon^2)$

$TV(p^*, u) \leq \varepsilon$, $\mathbb{P}(\Delta\beta \leq \varepsilon, TV(\beta, p^*) \geq (1+\varepsilon)\varepsilon) \geq \underline{\mathbb{P}(TV(\beta, p^*) \geq (1+\varepsilon)\varepsilon)}$
 $\leq \mathbb{P}(\Delta\beta \leq \varepsilon, TV(\beta, p^*) \geq (1+\varepsilon)\varepsilon) + \mathbb{P}(\Delta u \geq \varepsilon)$

$$\mathbb{P}\left(\sup_{A \in \Sigma} |P_{\beta}(A) - P_u(A)| \geq 5\varepsilon\right) \leq \mathbb{P}\left(\sup_{A \in \Sigma} |P_{\beta}(A) - P^{**}(A)| \geq 4\varepsilon\right) \binom{N}{2} \exp(-cn\varepsilon^2)$$

robust estimation.

$$n\varepsilon^2 \geq \log N(\varepsilon)$$

$$\mathbb{P}(CTV(\beta, p^*) \leq 3 \text{ but } TV(u, p^*) \geq \varepsilon) \leq \exp(-n\varepsilon^2)$$

$$\text{Habermehl: } p^* \rightarrow (\beta) p + \ell \beta$$

$$TV(p^*, \beta) \leq \varepsilon$$

$$X_1, \dots, X_n \sim N(\mu, \sigma^2 I_d), \quad \mathbb{E} \| \hat{\mu} - \mu \|^2 \leq \frac{\sigma^2 d}{n}$$

$$\text{Nonparametric L-S} \quad y_i = f(x_i) + \varepsilon_i \sim N(0, 1)$$

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n (y_i - f(x_i))^2$$

$$\mathbb{E} \|\hat{f} - f^*\|_n^2 = \mathbb{E} \hat{f} \sum_{i=1}^n (f(x_i) - f^*(x_i))^2$$

$$\text{Local Gaussian width: } g(\mathcal{F}, F) = \mathbb{E} \left[\sup_{g \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \mathbb{E}_x g(x_i) \right]$$

$$\mathcal{F}^* = \{f - f^* : f \in \mathcal{F}\}$$

$$\text{Critical Eq.: } \mathbb{E}_x g(\mathcal{F}, \mathcal{F}^*) \leq \frac{\delta^2}{2}$$

Then if \mathcal{F}^* is star-shaped, for any $f \geq \mathcal{F}^*$

$$\mathbb{P}(\|\hat{f} - f^*\|_n^2 \geq (b + \delta)^2) \leq \exp\left(\frac{-n(b + \delta)^2}{2}\right)$$

$\mathbb{E}[\|\hat{f} - f^*\|_n^2 \leq \|\mathcal{F}^*\|^2 + b]$ → e.g. $n^{-1/3}$ one-step vs $n^{-1/2}$ -change.

star-shaped:
if $f \in \mathcal{F}$ then $\alpha f + (1-\alpha)f^* \in \mathcal{F}$, $\alpha \in (0, 1)$ → f^* is the origin.
If \mathcal{F}^* is not star-shaped, $\mathbb{E}_x(g(\mathcal{F}, \text{star}(\mathcal{F}^*)))$

Recap:

1. Estimation of TV. $X_1, \dots, X_n \sim p^* \in \mathcal{P}$
construct \hat{p} such that $TV(\hat{p}, p^*)$ is small.

2. Yatracos method.

→ find smallest ε $n\varepsilon \geq \text{say} TV(\hat{p}, \varepsilon)$

→ construct p_1, \dots, p_N which is an ε -cover of \mathcal{P}

→ $A = A_{ij} = \{p_i > p_j\} \quad i, j \in \{1, \dots, N\}$

$$\Delta(p) = \sup_{A \in \mathcal{A}} |p_n(A) - p(A)|$$

$$\hat{p} = \arg \inf_{p \in \{p_1, \dots, p_N\}} \Delta(p)$$

3. Givens that $TV(\hat{p}, p^*) < \sqrt{\dots}$ robust.

$y_i = f(x_i) + w_i$, $x_i \sim P_x$, $f^* \in \mathcal{F}$

$$f = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2 \rightarrow \left(\frac{1}{n} \sum_{i=1}^n (f(x_i) - f^*(x_i))^2 \right)^{1/2}$$

Local Gaussian width.

$$l_g(\delta, \mathcal{F}) = \mathbb{E}_{f \in \mathcal{F}} \sup_{i=1}^n \left| \sum_{j=1}^n w_{ij} g(x_j) \right| \quad \|g\|_F = \sqrt{\sum_{j=1}^n g(x_j)^2}$$

Critical radius, finds δ^* $l_g(\delta, \mathcal{F}) \leq \frac{\sigma}{2}$

\mathcal{F} is star-shaped around f^* $f^* \in \{f - f^* \mid f \in \mathcal{F}\}$. $\delta^* \cdot G(\delta, \mathcal{F}^*) \leq \frac{\sigma}{2}$
then $t \geq \delta^*$, $P(\|f - f^*\|_h \geq t) \geq \exp(-\frac{n \delta^2}{2})$

$$\|f - f^*\|_h \leq \delta^2 + \frac{1}{h}$$

Lemma. If \mathcal{F} is star-shaped, then $l_g(\delta, \mathcal{F})$ is \downarrow fn of δ

Proof: $l_g(\delta, \mathcal{F}) = \mathbb{E}_{\substack{f \in \mathcal{F} \\ \text{star-shaped}}} \sup_{i=1}^n \left| \sum_{j=1}^n w_{ij} h(x_j) \right| = \mathbb{E}_{\substack{h \in \mathcal{G} \\ \text{star}}} \sup_{i=1}^n \left| \sum_{j=1}^n w_{ij} h(x_j) \times \frac{1}{t} \right|$
 $\|h\|_h \leq 1$

$$\leq h(f, F) \frac{\delta}{t} \Rightarrow G(\delta, F) \leq \frac{h(F)}{t} \text{ decreasing..}$$



$$\frac{1}{n} \sum_{i=1}^n (f(x_i) - f^*(x_i))^2 \leq \frac{2}{n} \sum_{i=1}^n w_i (f(x_i) - f^*(x_i))$$

denote. $\Delta = f - f^*$.

$$\frac{1}{n} \sum_{i=1}^n (\Delta(x_i))^2 \leq \frac{1}{n} \sum_{i=1}^n w_i \Delta(x_i) \leq \sup_{\Delta \in \mathcal{D}} \frac{1}{n} \sum_{i=1}^n w_i g(x_i)$$

$$E|\Delta|^2 = \delta^2, \sum \Delta^2 \leq G(F, \delta)$$

$$\Rightarrow G(F, F^*) \leq \frac{\delta^2}{2}, \delta \geq \delta_0, A(u) = \{g \in F, |g|_h \geq u, \frac{1}{n} \sum w_i g(x_i) \geq 2|g|_h x u\}$$

$$P(A(u)) \leq \exp(-nu^2)$$

$$t \geq \delta, u = \sqrt{t\delta} \geq \delta.$$

$$|\Delta|_h \leq \sqrt{t\delta} \rightarrow \text{choose } |\Delta|_h \geq \sqrt{t\delta} \text{ cond. on } A(u) \subset |\Delta|_h \leq 4\sqrt{t\delta}$$

$$\frac{1}{n} \sum_{i=1}^n (\Delta(x_i))^2 \leq \frac{1}{n} \sum_{i=1}^n w_i (\Delta(x_i)). \frac{1}{n} \sum_{i=1}^n w_i \Delta(x_i)^2 \leq 2|\Delta|_h x u \leq 4u \Rightarrow \text{on } A(u) \quad (\frac{1}{n} \sum_{i=1}^n w_i \Delta(x_i)^2 \leq 16u^2)$$

$$P(|\Delta|^2 \geq 16u^2) \leq \exp(-\frac{n u^2}{2})$$

$$\tilde{g} = \frac{g}{|\Delta|_h} x u, \frac{1}{n} \sum_{i=1}^n w_i \tilde{g}(x_i) \geq 2u^2$$

$$Z(u) = \sup_{\tilde{g} \in A(u)} \frac{1}{n} \sum_{i=1}^n w_i \tilde{g}(x_i), P(A(u)) \leq P(Z(u) \geq 2u^2)$$

$$E(Z(u)) \leq u^2$$

$$\sup_{\tilde{g} \in A(u)} \frac{1}{n} \sum_{i=1}^n w_i \tilde{g}(x_i) - \sup_{\tilde{g} \in A(u)} \frac{1}{n} \sum_{i=1}^n w_i g(x_i)$$

$$(P(A(u)) \leq \exp(-\frac{n u^2}{2})) \leq \frac{1}{n} \sum_{i=1}^n (w_i - w_i) |g(x_i)| \leq \frac{1}{\sqrt{n}} |g(x)|$$

$Z(u) \xrightarrow{\text{a.s.}} 0$. L-Lipschitz.

$$L = L \cdot P = P(Z(u) - E(Z(u))^2 t) \leq \exp\left(-\frac{t^2}{2}\right)$$

$$P(A(u)) \leq P(Z(u) \geq 2u^2) \leq P(Z(u) - E(Z(u))^2 \geq u^2) \leq \exp\left(-\frac{n - 4u^2}{2u^2}\right) \leq \exp\left(-\frac{n}{2u^2}\right)$$

$$h(F, F) \leq \frac{1}{m} \int_{\frac{\delta}{4}}^{\delta} \sqrt{\log N(\delta, F)} d\delta + \frac{\delta^2}{4}$$

$$\frac{1}{m} \int_{\frac{\delta}{4}}^{\delta} \sqrt{\log N(\delta, F)} d\delta \leq \frac{N}{4}$$

e.g. L-Lipschitz fun.

$$\log N(F, \delta) \leq \log N_{\text{ad}}(F, \delta) \leq \frac{L}{\delta}, E|\Delta - f^*|^2 \leq \frac{1}{m} \int_0^L \frac{L}{\delta} d\delta \leq \frac{1}{m} L$$

$$E|\Delta - f^*|^2 \leq \frac{1}{m^2}$$

$$\sqrt{\frac{1}{m^2} \cdot \frac{L^2}{\delta^2}} \leq \frac{L}{4}, \frac{\sqrt{\frac{3}{2}} \cdot \frac{L}{\sqrt{m}}}{\delta} \geq \frac{1}{m^2}$$

Another example.. $\sqrt{d \log(k)}$

Recap

1. Nonparametric L.S. $\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_i (y_i - f(x_i))^2$

2. Local Gaussian width

$$G(\mathcal{L}, P) = \mathbb{E}_{\substack{f \sim \mathcal{F} \\ \|f\|_2 \leq \delta}} \frac{1}{n} \sum_{i=1}^n g(f(x_i)) w_i$$

3. Critical radius δ^{*}

$$\text{Any positive soln to } G(\mathcal{L}, P) \leq \frac{\delta^2}{2}$$

4. Main result: $\mathbb{E} \|\hat{f} - f^*\|_2^2 \leq \delta^{*2} + \frac{1}{n}$

5. Can upper bound $G(\mathcal{L}, P) \leq \frac{16}{\sqrt{n}} \int_{\mathcal{L}} \sqrt{\log(M_P(\mathcal{F}), \|f\|_2)} d\mu_f$
Sufficient to find δ^* , $\int_{\mathcal{L}} \sqrt{\frac{1}{4}} d\mu_f$

Estimating functional function

$X_1 \dots X_n \sim f$, Entropy $= \int f \log f$. Expected density $= f^2$

Given $y_1 \dots y_n \sim g$. Goodness-of-fit $\int (f - f_0)^2$

$$\hookrightarrow \text{KL: } \int f \log \frac{f}{g} + (\frac{1}{2} \int (f - f_0)^2)$$

Causal Inf: $(X, A, Y) \mid X \sim P_X$, $P(X = 1) \propto \text{prop score}(A, Y \mid A, X \sim N(\mu, 1))$
 $A \perp E = \int (y \mid A=1, X=x) P_X dx$

Plug-in estimator: $\hat{f} = T \hat{f} + \int f^2 - \hat{f}^2 = \int \hat{f}^2$

$$R = \mathbb{E}(\hat{f} - T)^2 = \int \hat{f}^2 - T^2 = \left| \int f^2 - T^2 \right| \leq \left(\int (T^2 - f^2) \right) \leq (T^2 - \int f^2)$$

$(\hat{f} - f) \leq \frac{\beta}{\sqrt{n}}$ assume $X_1 \dots X_n \sim f$, $(0, 1)^d$, Holder- β

Not obvious \rightarrow parameter rate not achievable

$$\int f^2 = \int \hat{f}^2 + \int 2f(\hat{f} - f) + \int (\hat{f} - f)^2 = - \int \hat{f}^2 + 2 \int f \hat{f} + \int (\hat{f} - f)^2$$

$$T = \int f^2 = - \int \hat{f}^2 + \frac{2}{n} \sum_{i=1}^n \hat{f}(x_i) \quad (x_1, x_2, \dots, x_n, \hat{f} \text{ is fixed})$$

$$(\mathbb{E} \hat{T} - T) = \int \hat{f}^2 + 2 \int f \hat{f} - \int f^2 = \int (\hat{f} - f)^2$$

Plug $\leq n^{-\frac{2\beta}{2\beta+1}}$, $R_{\text{first-order}} \leq \frac{1}{n} + n^{-\frac{4\beta}{2\beta+1}}$ if $\beta > \frac{1}{2}$ then $\frac{1}{n}$
if $\beta < \frac{1}{2}$ then nonparametric

$$R \leq \frac{1}{n} + n^{-\frac{8\beta}{4\beta+1}} \quad \beta > \frac{1}{4} \text{ (parametric)} \quad \beta < \frac{1}{4} \text{: nonparametric}$$

$$\begin{aligned}
 & \text{Parametric: } X_1, \dots, X_d \sim N(\theta^*, I_d) \\
 & T(\theta^*) = \|\theta^*\|_2^2 \quad \hat{T} = \frac{1}{n} \sum_{i=1}^n X_i^T X_i = \bar{X}^T \bar{X} \quad \bar{X} \cdot \bar{Y} \sim N(\theta^* \cdot \bar{Y}, \frac{\sigma^2}{n}) \\
 & E\hat{T} - T = 0 \quad \text{Var}(\hat{T}) = \sum_{j=1}^d \text{Var}(\bar{X}_j, \bar{Y}_j) = \sum_{j=1}^d \left[E(\bar{X}_j, \bar{Y}_j)^2 - \theta_j^{*2} \right] = \sum_{j=1}^d \left[E(\theta_j^* + \tilde{\epsilon}_j)(\tilde{x}_j + \tilde{\epsilon}_j) \right. \\
 & \quad \left. - \theta_j^{*2} \right] \\
 & \leq \sum_{j=1}^d \left(\eta^2 + \frac{\theta_j^{*2}}{\eta} \right) \leq \frac{\|\theta^*\|^2}{\eta} + \frac{d}{\eta^2}
 \end{aligned}$$

$$R_f \leq \frac{1}{\eta} + \frac{d}{\eta^2}. \quad \text{If d cl. parameter} \downarrow \text{D rate}$$

$$\mathbb{E}[\hat{\theta}_j - \theta_j^*]^2 \leq \frac{d}{n}$$

$$\text{with decay: } \theta^*: \sum_{j=1}^{\infty} j^{-2/\alpha} \theta_j^{*\alpha} \leq C$$

$$\text{Truncate J: Bias: } \frac{\sum_{j=J+1}^{\infty} \theta_j^{*\alpha}}{J}, \text{Var: } \frac{\sum_{j=J+1}^{\infty} \theta_j^{*\alpha}}{J} / n^2$$

$$\text{Now how good is the estimator?} \quad \hat{T} = \sum_{j=1}^J \bar{X}_j \bar{Y}_j \quad R \leq \left(\sum_{j=1}^J \theta_j^{*\alpha} \right)^2 + \frac{C}{n} + \frac{J}{n^2}$$

$$\hat{T} = \sum_{j=1}^{\infty} \theta_j \psi_j. \quad \int f^2 \sum_{j=1}^{\infty} \theta_j^2 \leq J \sum_{j=1}^J \left(\frac{\theta_j}{J} \right)^2 + \frac{C}{n} + \frac{J}{n^2} \quad J \leq n^{\frac{2d}{d+2\alpha}}$$

Truncated version \leq Romax