

Yihong Wu, High dimensional statistics. Lec. 1.

Examples.

1) Sparse linear regression (compressed sensing) $\beta \in \mathbb{R}^p$, some underlying vector
 Observe: $y = A\beta + \text{noise}$
 If β is k -sparse, then we need $\frac{m}{k} \log p$ samples $\left. \begin{matrix} \text{ } \\ \text{ } \end{matrix} \right\} p \times m$

2) Cov. matrix $X_1, \dots, X_n \stackrel{iid}{\sim} N(0, \Sigma)$; $\Sigma: p \times p$, psd.
 $n \times p$ data matrix

Test: $H_0: \Sigma = I_p$.
 $H_1: \Sigma \neq I_p$ or its variant $\|\Sigma - I_p\| \geq \epsilon$.

3) Species problem: Ball-urn model.
 k balls in urn
 coloured, red, black, blue, etc. #. colours unknown.
 Can only draw from urn. Estimate # colours.
 Ref: Fisher's work in the 1940s, butterflies in the Malayan islands.

I: Fundamentals

II: Unstructured problems - primarily for lower bounds.

III: Structured problems - sparsity/smoothness/shape

IV: Unstructured functional estimation:

- E.g., instead of estimating Σ completely, focus on the principal component.
- instead of a distribution, estimate entropy.

V: Computational issues and barriers.

- E.g., community detection.
- Illustrative relaxation of combinatorial procedures.

I: Fundamentals: Decision theoretic framework
 Focus on minimax risk.

• Experiment a) $\{P_\theta, \theta \in \Theta\}$ prob distributions on alphabet \mathcal{X} .
 Θ : parameter space. (Model)

b) Data: $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} P_\theta$.

c) Estimator: $\hat{\theta}: \mathcal{X}^n \rightarrow \Theta$
 $(x_1, \dots, x_n) \mapsto \hat{\theta}(x^n)$.

Loss function: $l: \Theta \times \Theta \rightarrow \mathbb{R}$
 $(\theta, \hat{\theta}) \mapsto l(\theta, \hat{\theta})$, e.g. $\|\theta - \hat{\theta}\|$

$\inf_{\hat{\theta}} \sup_{\theta \in \Theta} \mathbb{E}_{\theta} [l(\hat{\theta}, \theta)]$: worst-case expected loss. $:= R_n(\Theta)$.

Questions: How does it depend on n, p , sparsity parameter, Θ .

Worst case Bayes: $\theta \sim \pi$.

$$\max_{\pi} \inf_{\hat{\theta}} \mathbb{E}_{\theta \sim \pi} [l(\hat{\theta}, \theta)] \leq \inf_{\hat{\theta}(\cdot)} \sup_{\theta \in \Theta} \mathbb{E}_{\theta} [l(\hat{\theta}, \theta)] = R_n(\Theta)$$

A convex program.

↳ Relax to randomized strategies

$P_{\hat{\theta}} | x^n$.

Relaxed version is also a convex program.

Remarks: 1) All lower bounds go this way; pick a ~~smart~~ π smartly.

2) Inspired good procedures (Aggregation, exponential weighting)

3) But often, impossible to compute, esp. in high dimensions.

Objective) Classical asymptotics: Fixed p , large n , $\Theta \subset \mathbb{R}^p$.

How does $R_n(\Theta) \xrightarrow{n \rightarrow \infty} 0$; $\frac{1}{n} \times \text{const}$

Asymp. normality $\text{Var} = \frac{1}{n} \cdot \text{Fisher information}$

In high dimensions: $p \gg n$ or $p, n \rightarrow \infty$.

So we want nonasympt. characterisations of $R_n(\Theta)$ with universal constant factors.

Minimax rate vs. minimax risk

$$\Psi_n(\Theta) \approx R_n(\Theta)$$

Notation: $X_n \gtrsim Y_n \iff \frac{X_n}{Y_n} \geq C_1$, where C_1 is some abs. constant

$Y_n \lesssim X_n \iff \frac{X_n}{Y_n} \geq C_2$

When both, we write $X_n \asymp Y_n$.

Objective: Characterise $R_n(\Theta)$ as $\Psi_n(\Theta)$ as above.

Lecture 2

Remarks: Is it wise to focus on $\mathbb{E}_\theta[l(\cdot)]$? Why not high prob. statements?

* Qn: $\Psi_n(\Theta)$ as a function of l ?
Structure theorems on the class of l 's?

Focus on:

• Θ - normed linear space
 $l(\theta, \hat{\theta})$, a norm of $\|\theta - \hat{\theta}\|$.

← set of sparse cov. matrices
set of sparse vectors
set of distr.

Sample complexity: Min sample size for which $\exists \hat{\theta}$ s.t. $\sup_{\theta} \mathbb{E}_\theta[l(\hat{\theta}, \theta)] \leq \epsilon$
 $:= n^*(\epsilon)$.

(inverse of $\mathbb{E} n \rightarrow R_n(\Theta)$).

Qn: Is independence of sampling important?

Example: Scalar Gaussian Location Model.

$$\theta \in \Theta = \mathbb{R}$$

$$X_i = \theta + Z_i, \quad Z_i \stackrel{iid}{\sim} N(0, 1)$$

2) Many hypotheses. Fano's inequality

Given $\{P_1, P_2, \dots, P_M\}$ on \mathcal{X} .

Test $\hat{J}: \mathcal{X} \rightarrow [M] = \{1, 2, \dots, M\}$

$$\max_{j \in [M]} P_j[\hat{J} \neq j] \geq \frac{1}{M} \sum_j P_j[\hat{J} \neq j] \geq 1 - \frac{I(U; X) + \log 2}{\log M}$$

References: 1) A. Tsybakov Introduction to nonparametric statistics.

2) Le Cam Asymptotically efficient ...

3) Le Cam & Yang. 2nd book.

$$\inf_{\hat{\theta}} \sup_{\theta \in \Theta} \mathbb{E}_{\theta} [\|\hat{\theta} - \theta\|^2] \geq \sup_{\theta_0, \theta_1 \in \Theta} \left[\frac{\|\theta_0 - \theta_1\|^2}{8} (1 - \text{TV}(P_{\theta_0}, P_{\theta_1})) \right]$$

$$\begin{aligned} \text{TV}(P_{\theta_0}, P_{\theta_1}) &= \text{TV} \left(N(0, I_k), N \left(\begin{pmatrix} 0 \\ \vdots \\ \|\theta_1 - \theta_0\| \end{pmatrix}, I_k \right) \right) \\ &\stackrel{\text{here}}{=} \text{TV} (N(0, 1), N(\|\theta_1 - \theta_0\|, 1)) = 2Q\left(\frac{1}{2}\|\theta_1 - \theta_0\|\right) \end{aligned}$$

$$\begin{aligned} \text{Thus } \inf_{\hat{\theta}} \sup_{\theta \in \Theta} [\dots] &\geq \sup_{\theta_0, \theta_1 \in \mathbb{R}^k} \frac{\|\theta_1 - \theta_0\|^2}{8} \left(1 - 2Q\left(\frac{\|\theta_1 - \theta_0\|}{2}\right)\right) \\ &= \textcircled{H} (1) \quad (\text{Order } n \text{ disappears}) \end{aligned}$$

How to fix this issue?

$$1) \hat{\theta} : \mathbb{R}^k \rightarrow \mathbb{R}^k$$

$$\hat{\theta}_1 : \mathbb{R}^k \rightarrow \mathbb{R}$$

$$\vdots$$

$$\hat{\theta}_k : \mathbb{R}^k \rightarrow \mathbb{R}$$

$$\begin{aligned} \sum_{i=1}^k \inf_{\hat{\theta}_i} \sup_{\theta \in \Theta} \mathbb{E}_{\theta} [(\hat{\theta}_i - \theta_i)^2] &\geq \sum_{i=1}^k \inf_{\hat{\theta}_i} \sup_{\theta_i} \mathbb{E}_{\theta_i} [(\hat{\theta}_i - \theta_i)^2] \\ &= k \cdot \textcircled{H} (1). \end{aligned}$$

2) Many hypotheses. Fano's inequality

Given $\{P_1, P_2, \dots, P_M\}$ on \mathcal{X} .

Test $\hat{J} : \mathcal{X} \rightarrow [M] = \{1, 2, \dots, M\}$

$$\max_{j \in [M]} P_j[\hat{J} \neq j] \geq \frac{1}{M} \sum_j P_j(\hat{J} \neq j) \geq 1 - \frac{I(U; \mathcal{X}) + \log 2}{\log M} \quad \leftarrow \text{distr. that is unif. }([M]).$$

References: 1) A. Tsybakov Introduction to nonparametric statistics.

2) Le Cam Asymp. theory of ... statistics.

3) Le Cam & Yang. 2nd book.

Yihong Wu Lecture 3.

Recapitulation.

- $\{P_\theta : \theta \in \Theta\}$ model

Typically $l(\theta, \hat{\theta}) = \|\theta - \hat{\theta}\|^2$

$$\inf_{\hat{\theta}} \sup_{\theta \in \Theta} \mathbb{E}_\theta [\|\theta - \hat{\theta}\|^2]$$

Running example.

$$P_\theta = N(\theta, \frac{I_k}{n}) ; x = \theta + \frac{z}{\sqrt{n}}$$

min-max risk is $\frac{k}{n}$.

$$R_n(\Theta) = \inf_{\hat{\theta}} \sup_{\theta \in \Theta} \mathbb{E}_\theta [\|\theta - \hat{\theta}\|^2] = \max_{\pi} \inf_{\hat{\theta}} \mathbb{E}_{\theta \sim \pi} [\|\theta - \hat{\theta}\|^2]$$

Pick θ_0, θ_1 , and concentrate π on it.

Provides a l.b.

This is Le Cam's method.

$$R_n(\Theta) \geq \sup_{\theta_0, \theta_1} \frac{\|\theta_1 - \theta_0\|^2}{8} (1 - TV(P_{\theta_0}, P_{\theta_1}))$$

$$= \Theta(1).$$

So, we need a better π .

- Two fixes

a) Component-wise distances

b) Many hypotheses and Fano inequality.

- Recall Fano inequality.

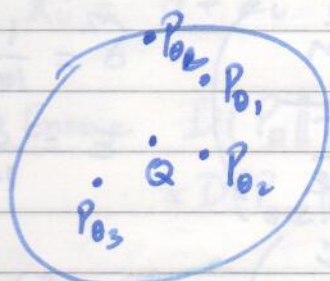
$$P_{\text{err, avg}} = \frac{1}{M} \sum_{j=1}^M P_j(\hat{J} \neq j) \geq 1 - \frac{I(U; X) + \log 2}{\log M}$$

$$I(U; X) = D(P_{X|U} \| P_X | P_U) = \inf_Q D(P_{X|U} \| Q | P_U) = \inf_Q \frac{1}{M} \sum_{j=1}^M D(P_j \| Q)$$

Note: $C = \max_{P_U} I(U; X) = \max_{P_U} \inf_Q D(P_{X|U} \| Q | P_U) \leq \inf_Q \max_{P_U} D(P_{X|U} \| Q | P_U) = \inf_Q \max_j D(P_j \| Q)$.

$$\leq \max_{i,j} D(P_j \| P_i).$$

Thus $P_{\text{err,avg}} \geq 1 - \frac{\max_{i,j} D(P_j \| P_i) + \log 2}{\log M}.$



• Suppose $\theta_1, \theta_2, \dots, \theta_M \in T \subset \mathbb{H}$ are such that $\|\theta_i - \theta_j\| \geq \varepsilon.$

• Estimator $\hat{J} = \arg \min_j \|\theta_j - \hat{\theta}\|^2.$

$$\approx \varepsilon^2 \left(1 - \log 2 \frac{C(T) + \log 2}{\log M(T, \varepsilon)} \right), \text{ where } C(T) = \max_{U \in T} I(U; X).$$

$M(T, \varepsilon)$: packing number. (T will be defined soon as a convex set)
Metric entropy

• $T \subset (\mathbb{H}, d)$

ε -packing number of T : $M(T, \varepsilon) = \max \# \text{ disjoint balls of radius } \varepsilon \text{ (centres in } T \text{) that can be packed in } T.$

ε -covering number of T : $N(T, \varepsilon) = \text{small number of balls of radius } \varepsilon \text{ that cover } T.$

• Fact: $M(T, \varepsilon) \geq N(T, 2\varepsilon) \geq M(T, 2\varepsilon).$

Fact: $M(T, \varepsilon) \leq \text{Vol}$

$$N(T, 2\varepsilon) \geq \frac{\text{Vol}(T)}{\text{Vol}(B(0, 2\varepsilon))}$$



• Let $T = B(0, \delta)$, L^2 -ball

$C(T)$ for the example is $\max_{\|\theta\| \leq \delta} I(\theta; \theta + \frac{1}{\sqrt{n}} \Gamma_k) \leq \max_{\|\theta\| \leq \delta} I(\theta; \theta + \frac{1}{\sqrt{n}} \Gamma_k)$

$$= \frac{k}{2} \log \left(1 + \frac{\delta^2 n}{k} \right) \leq \text{const} \cdot \frac{1}{2} \delta^2 n.$$

Yihong Wu Lecture 3

Recap $\log M(T, \epsilon) \geq \log \frac{\text{Vol}(T)}{\text{Vol}(B(0, 2\epsilon))} = \log \frac{\delta^k}{(2\epsilon)^k} = k \log \frac{\delta}{2\epsilon}$

Now lower bound: $\epsilon^2 \left(1 - \frac{\text{const.} \cdot \frac{1}{2} \delta^2 n}{k \log \frac{\delta}{2\epsilon}} \right) \quad \delta^2 = \frac{1}{100} \frac{k}{n}$
 $\epsilon = \frac{\delta}{100}$

$\epsilon^2 \left(1 - \frac{\text{const.} \cdot \frac{k^2 \cdot n}{nk}}{k} \right) = O(1) \cdot \epsilon^2 = O\left(\frac{k}{n}\right)$
 controlled

$(X; U) I_{T=U} = (1) I_{T=U}$

So we need a better ϵ

Two fixed $(k, \epsilon) \subset T$
 T is called (k, ϵ) -separated if $\forall x, y \in T, x \neq y, \|x - y\| \geq \epsilon$

Recall Fano inequality: $I(U; X) \leq H(p) - \sum_{j=1}^M p_j H(p_j | P_j)$

$I(U; X) = D(P_{X|U} \| P_X | P_U) = \sum_{j=1}^M p_j D(P_j \| P_X)$

$I(U; X) \geq \max_{j \in \{1, \dots, M\}} I(U; X | U=j) \geq \max_{j \in \{1, \dots, M\}} \log \frac{1}{p_j} = \log M$

Lec 4.

- 1) What if noise is non-Gaussian?
- 2) What if loss is not mean-square.

1) $X_i = \theta + Z_i \rightarrow$ non-Gaussian

$$\text{Ensemble } D(P_\theta \| P_{\theta'}) \approx \frac{1}{2} \|\theta - \theta'\|^2$$

$$\propto D(P_{\theta+Z} \| P_Z) \approx \frac{1}{2} \|\theta\|^2$$

$$\frac{C(T)}{M(T, \epsilon)} \leq \frac{\max_{\theta, \theta' \in T} D(P_\theta^{\otimes n} \| P_{\theta'}^{\otimes n})}{M(T, \epsilon)} \approx \frac{n\sigma^2}{M(T, \epsilon)}$$

- 2) Diff. loss for $\|\cdot\|$ norm on \mathbb{R}^n .

$$\|x\|_* = \sup_{\|y\| \leq 1} \langle x, y \rangle$$

We will see $\frac{k}{n} \leq R_n \leq \frac{E[\|Z_n\|_*^2]}{n}$

$$\frac{1}{n} \left(\frac{k}{E[\|Z\|_*]} \right)^2$$

Take $\|\cdot\|_1$. Then $E[\|Z\|_\infty] \approx \sqrt{\log n}$

Thus $\frac{k^2}{n \cdot \log n}$: off by a factor $\log n$.

Use the volume method:

$$\log M(B_2(\delta), \epsilon) \geq \log \frac{\text{Vol}(B_2(\delta))}{\text{Vol}(B_{\|\cdot\|}(\epsilon))}$$

Need upper bound on $\text{Vol}(B_{\|\cdot\|}(\epsilon))$.

Fact: Let $K \subset \mathbb{R}^d$ convex body with a nonempty interior.

Then $\text{Vol}(K) \geq \frac{1}{\mathbb{E}[\|Z\|_K]}$, where $Z \sim N(0, I_d)$

$$\|x\|_K = \inf \{ \alpha > 0 \mid x \in \alpha K \}$$

$$= \inf \{ \alpha > 0 \mid \frac{x}{\alpha} \in K \}$$

Now lower bound:

If $K = B_{\|\cdot\|}$, then $\|\cdot\|_K = \|\cdot\|$.

Fact: $\text{Vol}(K) \leq \frac{w(K)}{\sqrt{d}}$, where $w(K) = \mathbb{E} \sup_{y \in K} \langle Z, y \rangle$, Gaussian width of K .

Thus $\text{Vol}(K)^{1/d} \leq \delta \text{Vol}(B_{\|\cdot\|}) \leq \frac{\delta}{\sqrt{d}} \mathbb{E} \left[\sup_{\|y\| \leq 1} \langle Z, y \rangle \right]$.

• Covariance matrix estimation.

X_1, X_2, \dots, X_n iid $N(0, \Sigma_{k \times k})$.

Suff. statistic $S = \frac{1}{n} \sum_{i=1}^n X_i X_i^T = \frac{1}{n} X X^T$.

Is it minimax? Stein: minimax must satisfy a group invariance property. This S does not.

$$\inf_{\Sigma} \sup_{\|\Sigma\|_{\text{op}} \leq \lambda} \mathbb{E} \left[\|\hat{\Sigma} - \Sigma\|^2 \right] \asymp \left(\frac{k \wedge 1}{n} \right) \lambda^2 \cdot \|I_k\|^2$$

\downarrow
 $\sigma_{\max}(\Sigma)$

as long as the norm is unitarily invariant.

Fact: S attains the upper bound to within constant factors.

$$D(N(0, \Sigma_0) \| N(0, \Sigma_1)) = \frac{1}{2} \text{Tr} \left(\Sigma_0^{-1} \Sigma_1 - I \right) - \frac{1}{2} \log \frac{\det \Sigma_1}{\det \Sigma_0}$$

$$= \frac{1}{2} \sum_k \left(\sigma_i(\Sigma_0^{-1} \Sigma_1) - 1 - \log \sigma_i(\Sigma_0^{-1} \Sigma_1) \right)$$

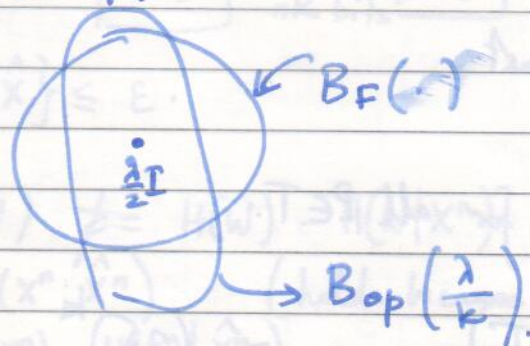
Take $T = \left\{ \Sigma \mid D(N(0, \Sigma) \parallel N(0, \Sigma_0)) \leq \dots \right\}$.

When $\sigma_i(\Sigma_0^{-1} \Sigma_1) \approx 1$, then $D(\cdot \parallel \cdot)$ is roughly quadratic and

we get
$$\frac{1}{2} \sum_i (\sigma_i(\Sigma_0^{-1} \Sigma_1 - I))^2 = \frac{1}{2} \|\Sigma_0^{-1} \Sigma_1 - I\|_F^2$$

$$\leq \frac{1}{2} \|\Sigma_0^{-1}\|_F^2 \cdot \|\Sigma_0 - \Sigma_1\|_F^2$$

Use this to pick T :



Basic idea: $\frac{\text{Vol}(\text{KL inbd})}{\text{Vol}(\text{ball coming from loss } f_{\text{KL}})}$

- DMS $P_X = P$.
 X_1, X_2, \dots, X_n iid P
 P : a distr. on $[k] = \{1, 2, \dots, k\}$.

$$\inf_{\hat{P}} \sup_{P \in M_k} \mathbb{E} \text{TV}(P, \hat{P}) = R(n, k) \asymp \sqrt{\frac{k-1}{n}} \wedge 1.$$

Sufficient statistic empirical distr. $\frac{N(j | X^n)}{n}, j \in [k].$

Focus on $k \geq 2$

$$\mathbb{E}[\text{TV}(\hat{P}, P)] \leq \mathbb{E}[\|P - \hat{P}\|_1]$$

$$= \sum_{j=1}^k |p_j - \hat{p}_j| \leq \sum_{j=1}^k \sqrt{\frac{p_j(1-p_j)}{n}} \leq \frac{1}{\sqrt{n}} \sum_{j=1}^k \sqrt{p_j} \leq \sqrt{\frac{k}{n}}$$

Lower bound: $D(P||Q) = \sum_{j=1}^k p_j \log \frac{p_j}{q_j} = \sum_{j=1}^k p_j \left(-\frac{q_j}{p_j} + 1 - \left(\frac{q_j}{p_j} - 1\right)^2 + \dots \right)$

$$\approx \sum_{j=1}^k p_j \left(\frac{q_j}{p_j} - 1\right)^2 + \dots$$

$$Q = \left(\frac{1}{2(k-1)}, \dots, \frac{1}{2(k-1)}, \frac{1}{2} \right)$$

$$T = \left\{ P \mid |p_i - q_i| \leq \delta \quad \forall i=1, \dots, k-1, \quad p_k = 1 - \sum_{i=1}^{k-1} p_i \right\} \subset M_k$$

prob-complex.

$$\delta = \left(\frac{1}{k-1} \wedge \frac{1}{\sqrt{n(k-1)}} \right) \text{ const}$$

$D(P||Q) \leq \delta^2 (k-1)^2$ follows for all $P \in T$.

$$\log M(T, \varepsilon) \geq \frac{\text{Vol}(T)}{\text{Vol}(B_{k-1}(\varepsilon))} \geq \frac{(2\delta)^k}{\frac{(2\varepsilon)^{k-1}}{(k-1)!}} \geq \frac{\delta(k-1)^{k-1}}{2\varepsilon e}$$

$\rightarrow l_1$ -distance

Summary:

$$\inf_P \sup_{P \in M_k} \mathbb{E}[D(P||\hat{P})] \asymp \frac{k-1}{n} \quad (\text{for ML})$$

Need to change the estimator \hat{P} to get better constants

$$\left. \begin{array}{l} \text{add 1 estimator (Laplace)} \\ \text{add } \frac{1}{2} \text{ estimator ex.} \end{array} \right\} \frac{k-1}{2n}$$

Discussion.

① Open questions

Does metric entropy related bounds govern statistic noise?

(*) convex body \mathcal{C}^n , $X = \theta + \frac{1}{\sqrt{n}} Z$

$\|\cdot\|_2^2$ loss $f_{\mathcal{C}}$.

How does $R_n(\odot)$ behave as a $f_{\mathcal{C}}$ of (*).

-11- Can say something if \odot is characterised by sparsity.

Day 3

Lecture 5

Yihong Wu

$$\bullet R_n(\mathbb{H}) \geq \varepsilon^2 \left(1 - \frac{C(T) + \log 2}{\log M(T, \varepsilon)} \right)$$

A connection: Set ε to be pairwise ^{min} distance

$$M(T, \varepsilon) = \max \{ M \mid \exists \theta_1, \theta_2, \dots, \theta_M \text{ s.t. } \|\theta_i - \theta_j\| \geq \varepsilon \forall i, j \}$$

We already saw

$$M(T, \varepsilon) \geq N(T, \frac{\varepsilon}{2}) \geq \frac{\text{Vol}(T)}{\text{Vol}(B_{\mathbb{H}}(\frac{\varepsilon}{2}))}$$

$$T = B_2(1) \subset \mathbb{R}^k$$

$$\left(\frac{1+\frac{\varepsilon}{2}}{\varepsilon} \right)^k = \left(\frac{1+\varepsilon/2}{\varepsilon/2} \right)^k \geq M(T, \varepsilon) \geq N(T, \frac{\varepsilon}{2}) \geq \left(\frac{1}{\varepsilon/2} \right)^k$$

$$X = \theta + Z \text{ where } Z \sim N(0, I_k), \theta \in \mathbb{R}^k$$

Criterion min. $\|\theta - \hat{\theta}\|^2$

$$\text{Know: Min MSE} = \frac{k}{n} \times \text{const.}$$

$$T = B_2(c_1 \sqrt{k})$$

$$\varepsilon = c_2 \sqrt{k}$$

$$C(T) = \max_{\theta, \theta' \in T} D(P_\theta \| P_{\theta'}) = \max_{\theta, \theta' \in T} \|\theta - \theta'\|^2 = \text{const. } k$$

$$\log M(T, \varepsilon) = \text{const. } k$$

Now, structure.

Dennoising with sparsity.

$$X = \theta + Z, \quad Z \sim N(0, I_p), \quad \theta \in \mathbb{R}^p$$

$$\text{But now } \theta \in \mathbb{H} = \{k\text{-sparse vectors}\} \subset \mathbb{R}^p$$

$$= \{ \theta \in \mathbb{R}^p \mid \sum_{i=1}^p \mathbb{1}_{\{\theta_i \neq 0\}} \leq k \} = \{ \theta \in \mathbb{R}^p \mid \ell_0\text{-norm of } \theta \leq k \}$$

$$= B_0(k)$$

Aside: $Z_1, Z_2, \dots, Z_p \stackrel{iid}{\sim} N(0, 1)$

$$\max_{1 \leq i \leq p} Z_i = \sqrt{2 \log p} + o_p(1) : \Pr\{Z_i > a\} = Q(a) \sim \frac{1}{a} \varphi(a) = \frac{\text{const.}}{a} e^{-a^2/2}$$

$$\Pr\left\{ \max_{1 \leq i \leq p} Z_i > \sqrt{2 \log p} \right\} \leq \frac{p}{\sqrt{2 \log p}} e^{-\log p} = o(1).$$

Lower bound: $\inf_{\hat{\theta}} \sup_{\|\theta_0\|_0 \leq k} \mathbb{E}_0 \|\hat{\theta} - \theta\|_2^2 \geq k \left(\frac{1}{k} \right)$ (because even if told where the nonzero locations are, need to estimate θ).

extra for sparsity locations.

$$\log \frac{ep}{k} \quad \left(\begin{array}{l} \forall k \in [p] \\ \forall p \in \mathbb{N} \end{array} \right)$$

- $T = B_2(\delta) \cap B_0(k)$
- $C(T) \leq \delta^2$.
- Choose (anticipating above result)
- $\varepsilon = \text{const.} \sqrt{k \log \frac{ep}{k}}$, $\delta = \text{const.} \cdot \varepsilon$.

$$\Rightarrow C(T) \leq k \log \frac{ep}{k}$$

Goal: Find $\theta_1, \theta_2, \dots, \theta_M \in T$ s.t. $\|\theta_i - \theta_j\| \geq \varepsilon$

First, find $b_1, b_2, \dots, b_M \in \{0, 1\}^p$ s.t.

a) $w_H(b_i) = k$

b) $d_H(b_i, b_j) \geq \text{const.} \cdot k =: m$

Define: $\theta_i = b_i \sqrt{\frac{\delta}{k}}$, $i = 1, \dots, M$.

Now $\|\theta_i - \theta_j\|_2^2 = d_H(b_i, b_j) \cdot \frac{\delta}{k} \geq \text{const.} \cdot \delta = \varepsilon$.

Number of such binary vectors: $M \geq \frac{\binom{p}{k}}{\left| \text{Binary}(m) \right|} = \frac{\binom{p}{k}}{1 + \binom{p}{1} + \binom{p}{2} + \dots + \binom{p}{m}} = \frac{\binom{p}{k}}{\binom{p}{\leq m}}$

Know: $\frac{\binom{p}{k}}{\binom{p}{\leq m}} \leq \left(\frac{ep}{k} \right)^k$

Thus $M \geq \frac{\binom{p}{k}}{\left(\frac{ep}{k} \right)^k} \Rightarrow \log M \geq \frac{k \log p/k - m \log \left(\frac{ep}{k} \right)}{\text{const.} \cdot k \cdot \log \left(\frac{ep}{k} \right)}$

MLE: $x \sim P_\theta = N(\theta, I_p)$
 $P_{x|\theta} \propto \exp(-\frac{1}{2} \|x - \theta\|_2^2)$

$\hat{\theta}_{MLE} = \operatorname{argmin}_{\theta: \|\theta\|_0 \leq k} \|x - \theta\|_2^2 \rightarrow$ Keep the largest k

Let $h = \hat{\theta}_{MLE} - \theta$

Observe that the ground truth is a feasible point.

Thus $\|x - \hat{\theta}\|_2^2 \leq \|x - \theta\|_2^2$

$\|z - h\|_2^2 \leq \|x - \theta\|_2^2 = \|z\|_2^2$ (Ground truth)
 $= \|z\|_2^2 + \|h\|_2^2 - 2\langle z, h \rangle$

$\Rightarrow \|h\|_2^2 \leq 2\langle z, h \rangle = 2\|h\| \cdot \langle z, \frac{h}{\|h\|} \rangle$

$\Rightarrow \|h\| \leq 2\langle z, \frac{h}{\|h\|} \rangle$

- ① $u = \frac{h}{\|h\|}$ has $\|u\| = 1$
- ② h is difference of 2 k -sparse vectors so u is $2k$ sparse.

Thus $\|h\| \leq 2 \cdot \sup_{u: \|u\|_0 \leq 2k} \langle z, u \rangle =: B$

- Can analyse this directly (chi-squared entries)
- Instead use metric entropy approach.

$u \in S^{p-1} \cap B_0(2k) \subseteq B_2(1) \cap B_0(2k)$

Let $\theta_1, \theta_2, \dots, \theta_N$ be an ϵ -net of $B_2(1) \cap B_0(2k)$

$\forall u, \exists i$ s.t. $u = \theta_i + \epsilon_i$, where $\|\epsilon_i\| \leq \epsilon$.

$\langle z, u \rangle = \langle z, \theta_i \rangle + \langle z, \epsilon_i \rangle$

\downarrow
 $N(0,1)$

\downarrow
 ϵ_i is $4k$ sparse.

$\sup_{\substack{\|u\|_2=1 \\ \|u\|_0 \leq 4k}} \langle z, u \rangle \leq \sqrt{2} B$

Thus $B \leq \max_{i \in [N]} \langle z, \theta_i \rangle + \epsilon \sqrt{2} B$

Thus $B \leq \max_{i \in [N]} \langle z, \theta_i \rangle \leq \sqrt{\log N}$

But we know $\sqrt{\log \left[\left(1 + \frac{2}{\epsilon}\right)^k \binom{p}{k} \right]} \asymp \sqrt{k \log \frac{ep}{k}}$

Remark: What if we don't know k ?

We will have two more procedures that give $\sqrt{k \log p}$ instead of $\sqrt{k \log \frac{ep}{k}}$.
But k need not be known.

Lecture 6.

Two procedures with simpler analyses.

$$1) \arg \min_{\theta} \|\theta\|_0 \quad \text{s.t.} \quad \|x - \theta\|_\infty \leq \tau = \sqrt{2 \log p} = \hat{\theta}$$

Solution: hard thresholding $\hat{\theta}_i = \mathbb{1}\{|x_i| > \tau\}$.

We will ^{analyse} it differently.

w.h.p., $E \|\hat{\theta} - \theta\|_2^2 \leq k \log p$

Also, ground truth θ is feasible w.h.p.

$$\left(\frac{\|z + \theta - \theta\|_\infty}{x} = \|z\|_\infty \stackrel{\text{w.h.p.}}{\leq} \sqrt{2 \log p} \Rightarrow \theta \text{ is feasible w.h.p.} \right)$$

$$\|\hat{\theta}\|_0 \leq \|\theta\|_0 \leq k$$

$\Rightarrow h = \hat{\theta} - \theta$ is at most $2k$ -sparse.

So $\|x - \hat{\theta}\|_\infty \leq \tau$

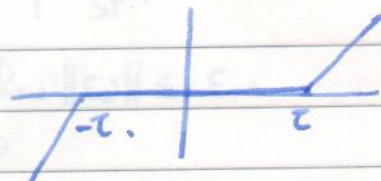
$\hookrightarrow \|z - h\|_\infty \leq \tau$. Since $\|z\| \leq \tau$, we have $\|h\|_\infty \leq 2\tau$.

Now $\|h\|_2^2 \leq \|h\|_\infty^2 \|h\|_0 \leq 4\tau^2 \cdot 2k$ which is $k \log p$.

However, we can get to $k \log \frac{p}{k}$ by thresholding exactly the top level differently, the second level differently, and so on.

$$2) \arg \min_{\theta} \|\theta\|_1 \quad \text{s.t.} \quad \|x - \theta\|_\infty \leq \tau = \sqrt{2 \log p} \quad \text{soft thresholding} \quad \hat{\theta}$$

$$\hat{\theta}_i = (|x_i| - \tau)_+$$



We will analyse it differently, as before.

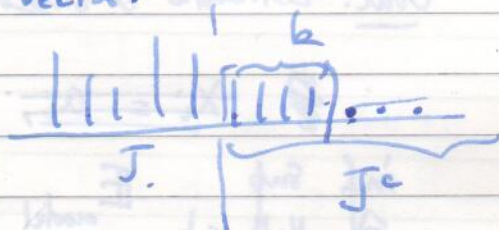
a) Some θ is feasible.

b) $\hat{\theta}$ has lower $\|\theta\|_1$ norm. $\|\hat{\theta}\|_1 \leq \|\theta\|_1$; $h = \hat{\theta} - \theta$

$$h = h_J + h_{J^c}, \text{ where } J = \text{supp}(\theta) \\ |J| \leq k.$$

$$\|h_{J^c}\|_1 = \|\theta\|_1 - \|\theta + h_J\|_1, \checkmark \\ \leq \|\theta\|_1 - \|\theta\|_1 + \|h_J\|_1 \\ \leq \|h_J\|_1.$$

Step 3. h is well-approx. by a $2k$ -sparse vector.
 $K = k$ largest indices of $|h_i|$, $i \notin J$.
 $|J \cup K| \leq 2k$.



Step 4: $\|h\|_2^2 \leq \|h_{J \cup K}\|_2^2 \times 2 \leq 2k \|h\|_\infty^2$ to be proved.
 $\leq \|h\|_\infty \cdot \|h_{J \cup K}\|_0$
 $= 2k \|h\|_\infty^2$

More than half of energy is in $2k$ -sparse components

$$\|h_{J^c}\|_1 \leq \|h_J\|_1 \implies \|h_{J \cup K}\|_2 \geq \|h_{(J \cup K)^c}\|_2$$

$$\forall i \in (J \cup K)^c, |h_{(i)}| \leq \frac{1}{i} \|h_{J^c}\|_1$$

$$\sum_{\substack{i \in (J \cup K)^c \\ i \geq k+1}} |h_{(i)}|^2 \leq \|h_{J^c}\|_1^2 \cdot \sum_{i \geq k+1} \frac{1}{i^2} \leq \|h_{J^c}\|_1^2 \cdot \frac{1}{k}$$

$$\leq \frac{1}{k} \|h\|_1^2 \stackrel{\text{CS ineq.}}{\leq} \|h_J\|_2^2 \leq \|h_{J \cup K}\|_2^2$$

Sparse linear regression. (Linear inverse problem, compressed sensing)

$$X = A\theta + Z \quad Z \sim N(0, I_n)$$

$\begin{matrix} & \xrightarrow{nxp} & \xrightarrow{px1} & \xrightarrow{nx1} \\ \begin{matrix} \text{ } \\ \text{ } \\ \text{ } \end{matrix} & \text{ } & \text{ } & \text{ } \\ \text{ } & \text{ } & \text{ } & \text{ } \end{matrix}$

Gaussian design matrix with ^{iid} entries $N(0, \frac{1}{n})$

Goal: Estimate θ based on X and A .

~~$x_i = a_i$~~ \rightarrow row of A .

$$\inf_{\hat{\theta}} \sup_{\|\theta\|_0 \leq k} \mathbb{E}_{\text{model}} \left[\|\hat{\theta} - \theta\|_2^2 \right] \gtrsim k \log \frac{ep}{k}$$

is a lower bound even if we see θ in noise directly.

over $P_{A,Z|\theta}$

Point, we can achieve $k \log \frac{ep}{k}$ inefficiently
 can achieve $k \log p$ efficiently
 when $n \gtrsim k \log \frac{ep}{k}$

Open: Why $n \gtrsim k \log \frac{ep}{k}$?

Lower bound:

$$D(P_{\theta_0} \| P_{\theta_1}) = D(P_{A,X|\theta_0} \| P_{A,X|\theta_1}) = \mathbb{E}_A \left[D(N(A\theta_0, I_n) \| N(A\theta_1, I_n)) \right]$$

$$= \frac{1}{2} \mathbb{E} \left[(\theta_0 - \theta_1)' A' A (\theta_0 - \theta_1) \right]$$

$$= \frac{1}{2} (\theta_0 - \theta_1)' \mathbb{E}[A' A] (\theta_0 - \theta_1)$$

$\xrightarrow{pxp} I_p$

$$= \frac{1}{2} \|\theta_0 - \theta_1\|_2^2$$

\rightarrow Same analysis as before
 carries over now

(It is as good as A is absent.)

Achievability MLE

$$\hat{\theta}_{MLE} = \arg \min_{\theta: \|\theta\|_0 \leq k} \|X - A\theta\|_2^2$$

Discussion: 1) Connection between randomised versus the worst-case settings we are studying in these lectures.

2) Sequential approaches?

3) θ iid $P_x = (1-\alpha)\delta_0 + \alpha N(0,1)$.

$$\frac{n}{p} = \alpha.$$

$n = \alpha p$ suffices.

Message passing style algorithms incorporating priors

Spatial coupling codes (entries of the matrix not necessarily iid Enable successive cancellations).

Achieves information theoretic lower bound.

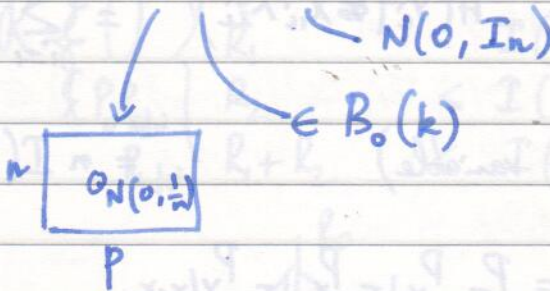
Day 4 23/07/2015

Yihong Wu Lec 7.

Recall: Sparse linear regression.

$$X = A\theta + Z$$

minimum MSE $\gtrsim k \log \frac{ep}{n}$



We will see that $R(k, p, n) \asymp k \log \frac{ep}{k}$, provided $n \gtrsim k \log \frac{ep}{k}$

$$\hat{\theta}_{MLE} = \arg \min_{\|\theta\|_0 \leq k} \|X - A\theta\|_2^2$$

$$\text{First } \|X - A\hat{\theta}\|_2^2 \leq \|X - A\theta\|_2^2 \leq \|Z\|_2^2$$

$$\text{Write } h = \hat{\theta} - \theta, \text{ then } \|Ah\|_2 \leq -\langle Ah, Z \rangle \leq \|h\|_2 \langle \frac{Ah}{\|h\|_2}, Z \rangle$$

Step 1. $\sup_{\substack{\|u\|_2=1 \\ \|u\|_0 \leq 2k}} \langle Au, Z \rangle$

Note that $\langle Au, Z \rangle = \langle A'z, u \rangle$

We will show $\sup \langle Au, Z \rangle \leq \sqrt{k \log \frac{ep}{k}}$

Step 2: $\|h\|_2^2 \approx \|Ah\|_2^2$

Discussion: $A'z = \left(\frac{A'z}{\|z\|} \right) \cdot \|z\| = \frac{\|z\|_2}{\sqrt{n}} \cdot N(0, I_p)$ vector.

$Ah =$ $2k$

conditional law, but unconditioned also yields $N(0, I_p)$.

$$\Pr \left\{ \min_{\theta} \left(\sum_{j \in J} A_j \theta_j \right) \leq 1 - \sqrt{\frac{m}{n}} - \frac{t}{\sqrt{n}} \right\} \leq e^{-t^2/2}$$

$N(0, \frac{1}{n})$

Since $\binom{p}{2k} \approx e^{2k \log \frac{ep}{k}}$

Choose $t = \sqrt{k \log \binom{p}{2k}}$

Recall:

We already saw $A = I_p$ case.

$$\left. \begin{array}{l} \min \|\theta\|_0 \\ \|\theta - y\|_\infty \leq \tau = \sqrt{2 \log p} \end{array} \right\} \text{ both were considered.}$$

Now A is a fat matrix.

$$\hat{\theta}_{\text{LS}} = \arg \min \|\theta\|_1 \quad (\text{Dantzig-Selctor})$$

s.t. $\|A'(x - A\theta)\|_\infty \leq \tau = \sqrt{2 \log p}$

Again, as before $h = \hat{\theta} - \theta$, $J = \text{supp}(\theta)$.

Then

$$\|h_J\|_1 \leq \|h\|_1$$

We will soon see that θ is also feasible. (Indeed, $A'z = \frac{\|z\|_2}{\sqrt{n}} N(0, I_n)$. Choose $z = \theta$.)

So $\|h_{J^c}\|_2^2 \geq \frac{1}{2} \|h\|_2^2$, as before, $K = \text{set of } k \text{ largest magnitudes outside } J$.

• By linear algebra + ϵ -net argument

$$\|h\|_2^2 \lesssim \|Ah\|_2^2 \leq 4\sqrt{k} \|h\|_2^2$$

near isometry argument

We will argue only the second inequality:

$$\|A'(x - A\theta)\|_\infty \leq \tau, \quad \|A'(x - A\theta)\|_\infty \leq \tau$$

$$\Rightarrow \|A'A h\| \leq 2\tau \text{ by triangle inequality}$$

Day 4 23/07/2015

Yihong Wu Lec 7. $(\frac{1}{2}, 0)h$

$$\|Ah\|^2 = \langle Ah, Ah \rangle \leq \langle A'A h, h \rangle$$

$$\leq \|A'A h\|_{\infty} \cdot \|h\|_1$$

$$\leq 2\epsilon \cdot 2 \|h\|_2 \leq 4\epsilon \sqrt{k} \|h\|_2$$

$$\leq 4\epsilon \sqrt{k} \|h\|_2$$

Functional Estimation. (E.g., Entropy estimation)

- $X = M + \frac{1}{\sqrt{n}} Z$
 $\left. \begin{array}{l} \text{pxp, with } N(0,1) \text{ components} \\ \text{estimate } M. \end{array} \right\}$

$$\inf_{\hat{M}} \sup_M \|\hat{M} - M\|_{op}^2 = \frac{1}{n} E[\|Z\|_{op}^2] \text{ by Anderson's lemma}$$

- What if we want only $\|M\|_{op}$?

$$\left| \|\hat{M}\|_{op} - \|M\|_{op} \right| \leq \|\hat{M} - M\|_{op}$$

$\Rightarrow \frac{1}{n}$ error is ach. via plug-in. \neq

It will turn out, we will ~~have~~ ^{have} $\frac{1}{n}$ here. So plug-in is best. No so for other functionals.

$$\underline{\text{Thm:}} \inf_{\hat{M}} \sup_M E \left[\left| \frac{1}{n} - \|M\|_{op} \right| \right] \asymp \frac{1}{n}$$

For lower bound, Le Cam's method works.

$$H_0: \|M\|_{op} = 0 \quad (M=0)$$

$$H_1: \|M\|_{op} = \lambda$$



Try $\lambda = \sqrt{p}$.

For simplicity, let $n=1$ (scaling)

$H_0: M=0 \quad (X=Z)$

$H_1: M = \lambda uu', \quad (u \text{ uniform } S^{p-1}), \quad X = \lambda uu' + Z$

$TV(Z, \lambda uu' + Z) \leq \dots$ is of interest since this arises in the LRT for the above test.

Let's bound $\sqrt{\chi^2(\lambda uu' + Z \| Z)}$ $(\chi^2(P \| Q) = \text{Var}_Q(\frac{P}{Q}) = \int \frac{(P-Q)^2}{Q})$

Observations:

$\chi^2(P_x \| P_z) = \int \frac{p_x^2}{p_z} - 1$

$p_x(z) = \mathbb{E}_M(\phi(x-M)) : \chi^2 = \int \frac{(\mathbb{E} \phi(x-M))^2}{\phi(x)} dx - 1$

Now have M, \tilde{M} iid

$= \mathbb{E} \int \frac{\phi(x-M)\phi(x-\tilde{M})}{\phi(x)} - 1$

$= \mathbb{E} [\exp \langle M, \tilde{M} \rangle] - 1$

$\langle M, \tilde{M} \rangle = \langle (\lambda uu'), (\lambda \tilde{u} \tilde{u}') \rangle = c.p. \langle u, \tilde{u} \rangle^2$

$\langle u, \tilde{u} \rangle = \langle \frac{Z}{\|Z\|}, \tilde{u} \rangle = O_p(\sqrt{p})$

Sparse M: Estimate $\|M\|_{op}$. $k \sim p^{1/3}$ is a threshold.

Lec. 8

Estimating entropy.

$X_1, X_2, \dots, X_n \stackrel{iid}{\sim} p$ on $[k]$

Estimate $H(p)$.

samples $\frac{k}{\log k}$ f(accuracy).

• If k is fixed, provably the best.

• If $k \sim n$, bias = $\mathbb{E}[\hat{H}_{\text{plug-in}} - H(p)] = -\mathbb{E}[D(\hat{P}||P)] = -\mathbb{E}[D(P||P)] \leq 0$.

Severely underbiased

Taylor expansion of bias $\mathbb{E}[\text{bias}] = -\frac{S(p)-1}{2n} + \frac{1}{12n^2} (\dots)$

support. But we may not know it.

Paminsky $\exists n_k = o(k)$ s.t. $R^*(k, n_k) \rightarrow 0$ as $k \rightarrow \infty$.

Valiant & Valiant: sharp scaling $n \sim \frac{k}{\log k}$

Result. $n \gtrsim \frac{k}{\log k}$

$\lesssim \frac{k}{\log k}$

$$R^*(k, n) \sim \left(\frac{k}{n \log k}\right)^2 + \frac{\log^2 k}{n}$$

no consistent est exists, i.e., $R^*(k, n) \gtrsim 1$

Compare with plugin: $\left(\frac{k}{n}\right)^2 + \dots$

Do unbiased estimators exist?

• Bernoulli r.v.'s. Estimate $f(p)$.

Unbiased est exists $\Leftrightarrow f(p) = \text{polynomial of degree } \leq n$.

$$\mathbb{E}[f(Y)] = \sum_{k=0}^n \hat{f}(k) \binom{n}{k} p^k (1-p)^{n-k}$$

- Use a large degree polynomial.

Bias $\rightarrow 0$. But variance \rightarrow increases because # parameters is large.

Find a sweet spot for the tradeoff.

- $h_j = \sum_{i=1}^k \mathbb{1}\{N_i=j\}$: histogram of histogram.

$$\hat{H} = \sum_{j=1}^k \alpha_j h_j, \text{ significantly faster than LP-based approaches.}$$

- Impossibility result.

Best polynomial approx.

$$\epsilon_L(f, I) = \inf_{\deg(p) \leq L} \sup_{x \in I} |f(x) - p(x)|$$

Dual: moment matching, upto L^{th} moments are same.

$$\epsilon_L(f, I) = \sup \mathbb{E}[f(U) - f(U')]$$

s.t. $U, U' \in I$

$$\mathbb{E}[U^j] = \mathbb{E}[U'^j], j=1, \dots, L$$

Fact:

$$\epsilon_L(f, I) = 2 \epsilon_L(f, I)$$

Idea. (LeCam's ~~idea~~)

Construct two ensembles of entropy distributions, one from H_0 : Entropy ≤ 3 bits, other from H_1 : Entropy ≥ 10 bits,

yet very difficult to distinguish between the two.

- One more example: How many butterflies, $S(P)$, support size estimation.

Community detection in networks

Summary

Community scales linearly, sharp recovery via SDP
 Sublinear, computationally hard via l.b.'s.

Two-equally sized communities, + -
 indep. w.p. $p = a \frac{\log n}{n}$ if ++ or --
 $q = b \frac{\log n}{n}$ if +- or -+

Assume $p > q$

MLE: $\max_{\sigma} \langle A, \sigma \sigma^T \rangle$
 s.t. $\sigma_i \in \{+1, -1\}, i \in [n]$
 $\sigma^T \mathbf{1} = 0$

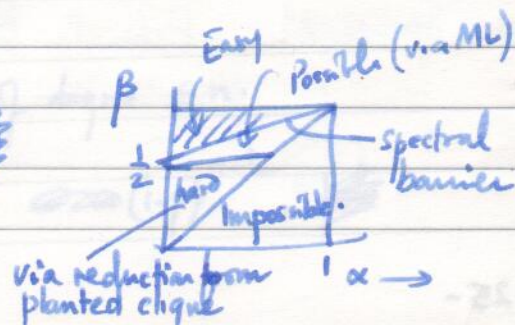
\Leftrightarrow $\max_Y \langle A, Y \rangle$
 s.t. $\text{rank } Y = 1$
 $Y_{ii} = 1$
 $\langle J, Y \rangle = 0$, $J = \text{all-1 matrix}$

Relax $\text{rank } Y = 1$ to $Y \succeq 0$.
 (Similar to Frieze-Jerrum 1995 for MAX BISECTION).

Info theoretic: 1) If $(\sqrt{a} - \sqrt{b})^2 > 2 \rightarrow$ recovery in poly
 $< 2 \rightarrow$ recovery is impossible.

2) SDP suffices in the recovery possible case.

Finding a single community
 (Planted clique problem) \rightarrow Planted dense subgraphs
 $K = pn$: SDP achieves sharp threshold.



-26- $K = n^p$

-24-

-LST

References

- Le Cam's book
- A. Tsybakov - Intro to nonparametric statistics.
lower bound (2nd chapter)
- Ibragimov & Khasminski - Statistical decision theory
- Iain Johnstone - Gaussian seq. models.
(On his web page).
- Inster & Sushina - Minimax \approx quadratic.
- ~~Quadratic~~ Quadratic \hat{f} estimator - papers

$$\frac{(U, X|Y)_n}{(X, X|Y)_n} - (U, X|Y)_n = (U, X|Y; X|Y)_n$$

Random access: binary source. See definitions of $\mathcal{P}^T - \mathcal{I}^T$

$$\frac{1}{2} \log \frac{(U, X|Y)_n}{(X, X|Y)_n} \geq \frac{1}{2} \log \frac{(U, X|Y)_n}{(X, X|Y)_n} \geq \frac{1}{2} \log \frac{(U, X|Y)_n}{(X, X|Y)_n}$$

$$\frac{1}{2} \log \frac{(U, X|Y)_n}{(X, X|Y)_n} \geq \frac{1}{2} \log \frac{(U, X|Y)_n}{(X, X|Y)_n}$$

a) CDMA: Treat interference as noise.
b) P_1, P_2 (y) $\frac{1}{2} \log \frac{(U, X|Y)_n}{(X, X|Y)_n}$