

Math Stats Lecture Notes



Lecture 6. Sequential Decision Theory and Sufficiency

1

Lecture 14. Concentration Inequality

38

Lecture 7. Bias - Variance Trade-off.

7

Lecture 15. Random vectors in high dims.

47

Lecture 8. Basic Estimation Methods

14

Lecture 16. Norms of Random Matrices

53

Lecture 9. Law of Large Number and Estimation Consistency

19

Lecture 17. High-dim Statistical Phenomena

56

Lecture 12. Unbiased Estimation and UMVUE.

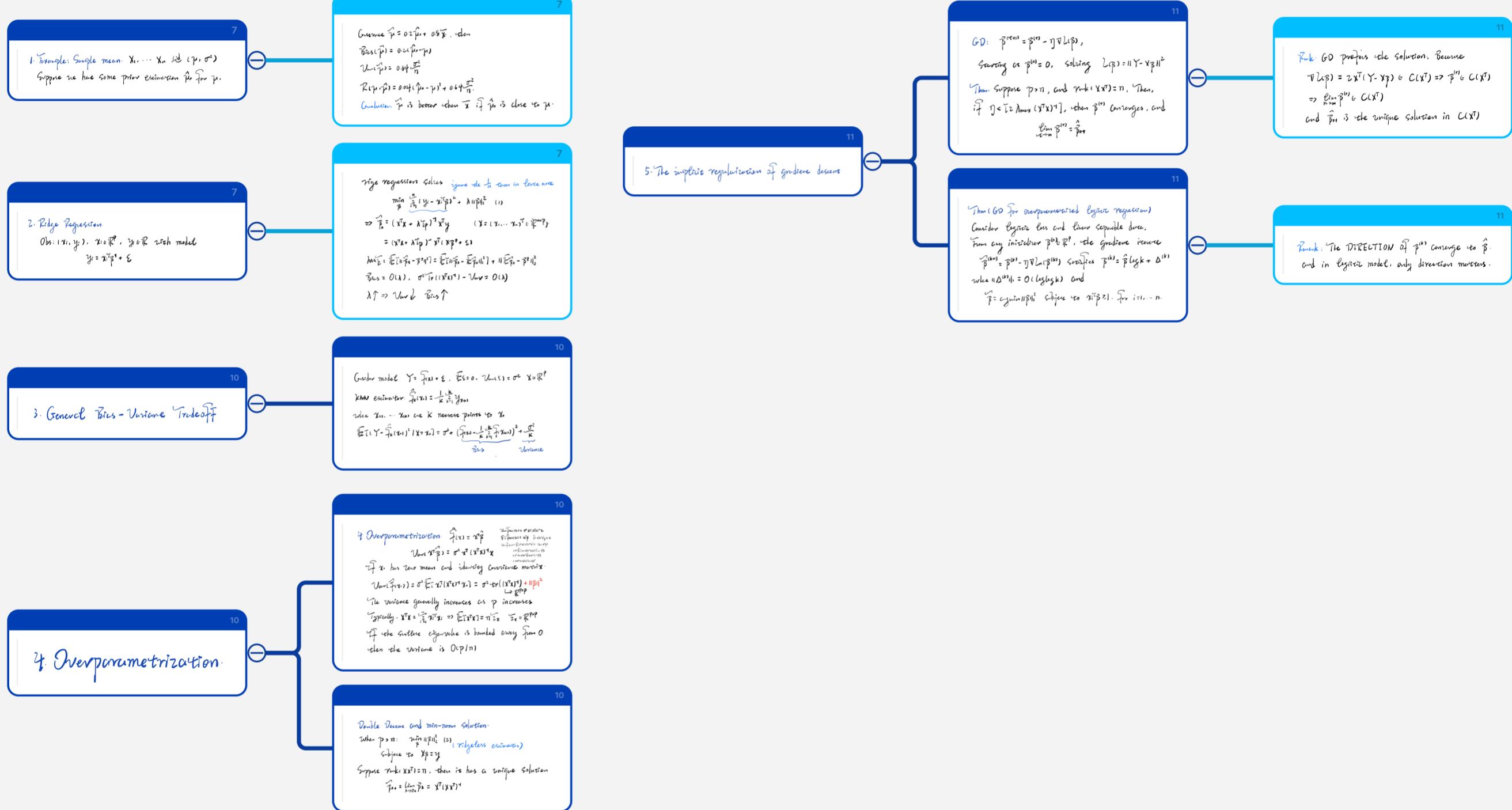
32

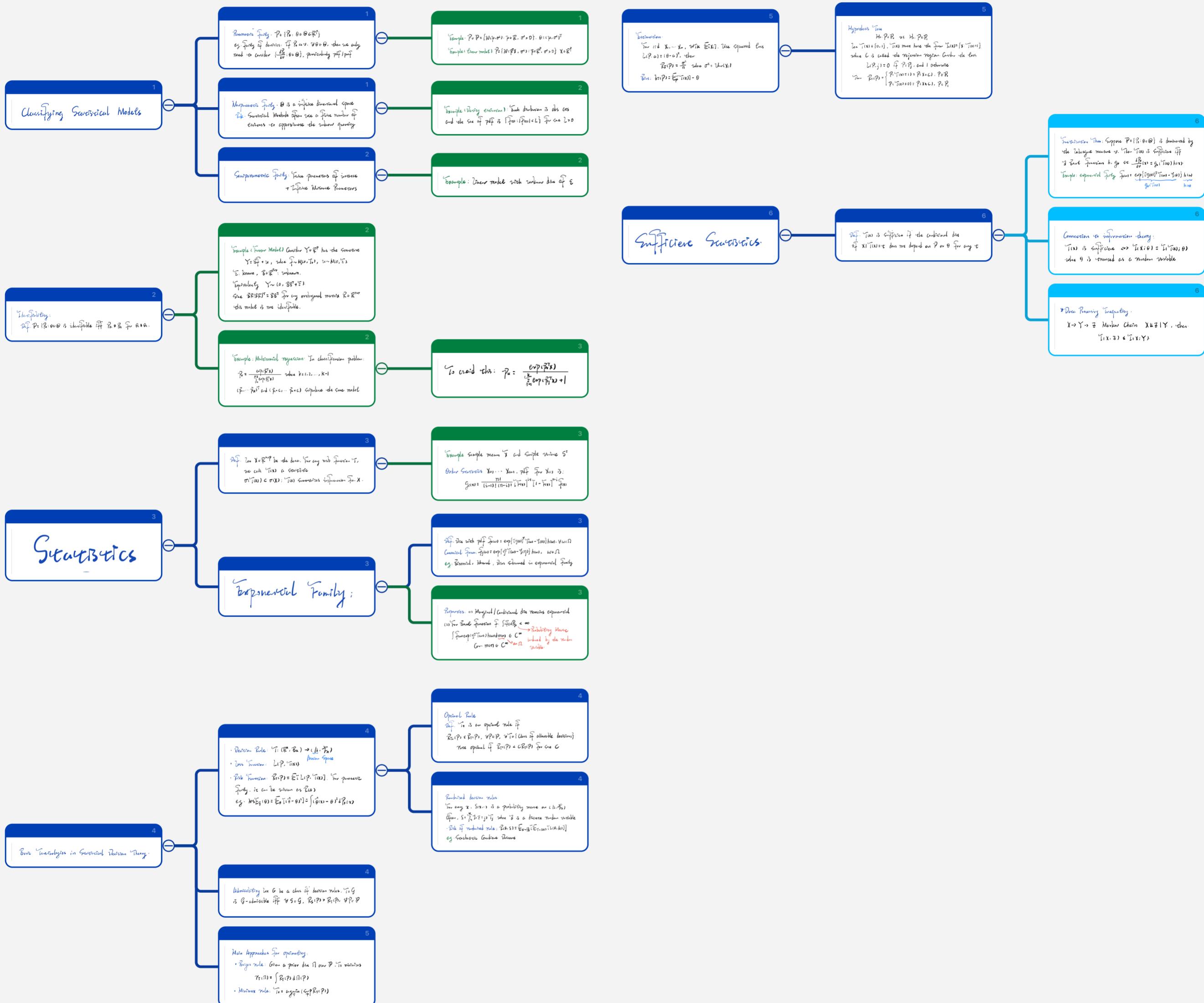
Lecture 18. Stein's phenomenon and shrinkage estimation.

63

Lecture 13: Fisher information and C-R Lower Bound.

34





1. Complete Statistic

Def. (Complete Statistics) $T(x)$ is complete for $P \in \mathcal{P}$ iff \forall Borel f , $E[f(T)] = 0 \Rightarrow f(T) = 0$ a.s. P

Rank: Requires the statistic has no redundant information.

Proposition. Let $\mathcal{P} = \{\tilde{P}_\eta, \eta \in \Xi\}$ be an exponential family of full rank (contains an open set) with pdf $\tilde{f}_\eta(x) = \exp\{\eta^\top T(x) - \tilde{\zeta}(\eta)\} h(x)$. Then $T(x)$ is sufficient and complete for $\eta \in \Xi$.

2. UMVUE

Def. An unbiased estimator $\hat{T}(x)$ of ν is called the UMVUE iff $\text{Var}(\hat{T}(x)) \leq \text{Var}(U(x))$ for any $P \in \mathcal{P}$ and any other unbiased estimator $U(x)$ of ν . Rank UMVUE does not always exist.

Then Suppose there exists a sufficient and complete statistic $T(x)$ for P . If ν is estimable then there exists unique UMVUE, which is of form $h(T)$ where h is a Borel function.

↳ Unbiased estimator

3. Construct UMVUE

Method 1. Find sufficient and complete statistic $T(x)$, then find $h(T)$ such that $E[h(T)] = \nu$, $\forall P \in \mathcal{P}$

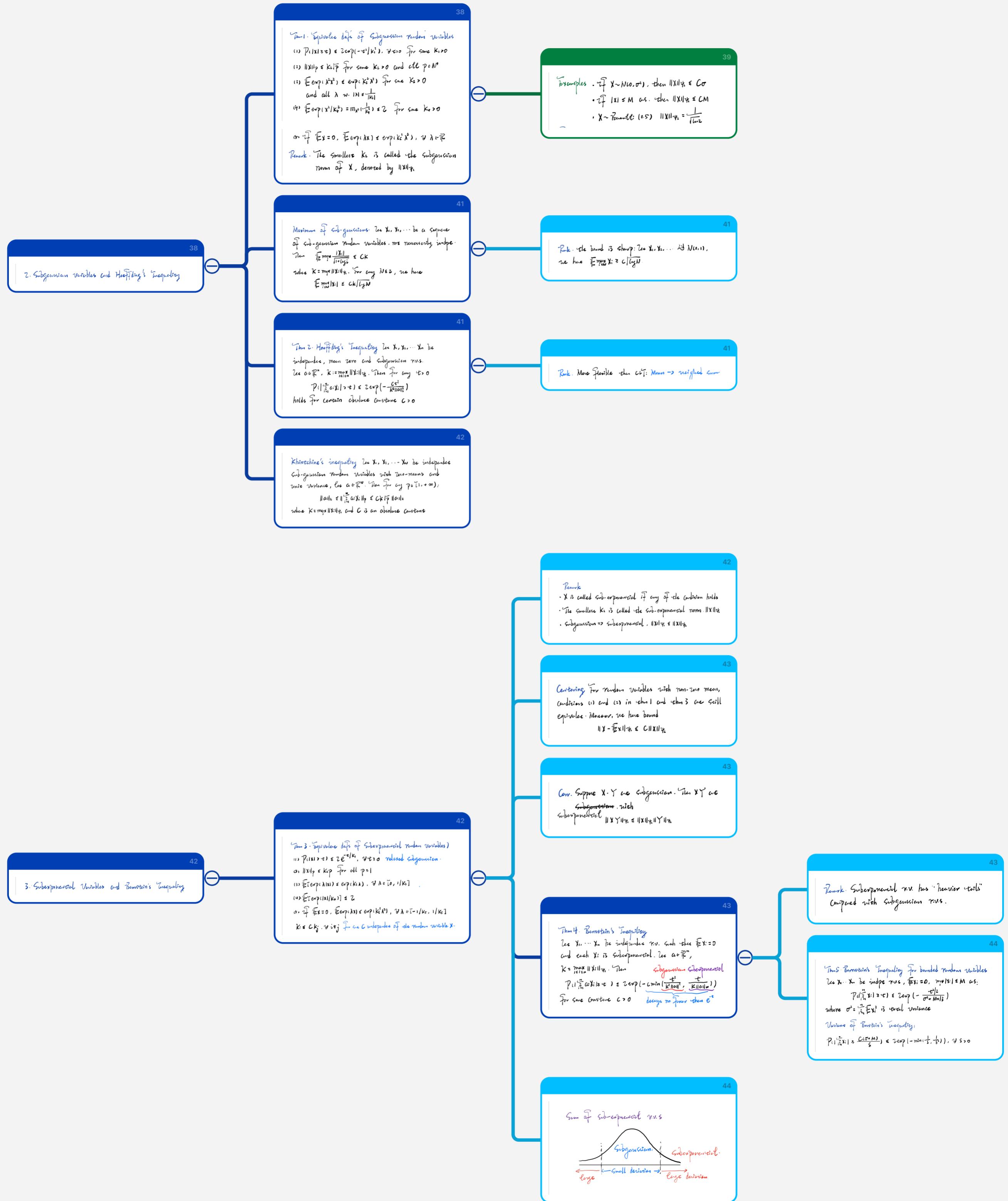
Method 2. Find sufficient, complete $T(x)$ and unbiased $U(x)$, then $E[U|T]$ is an UMVUE.

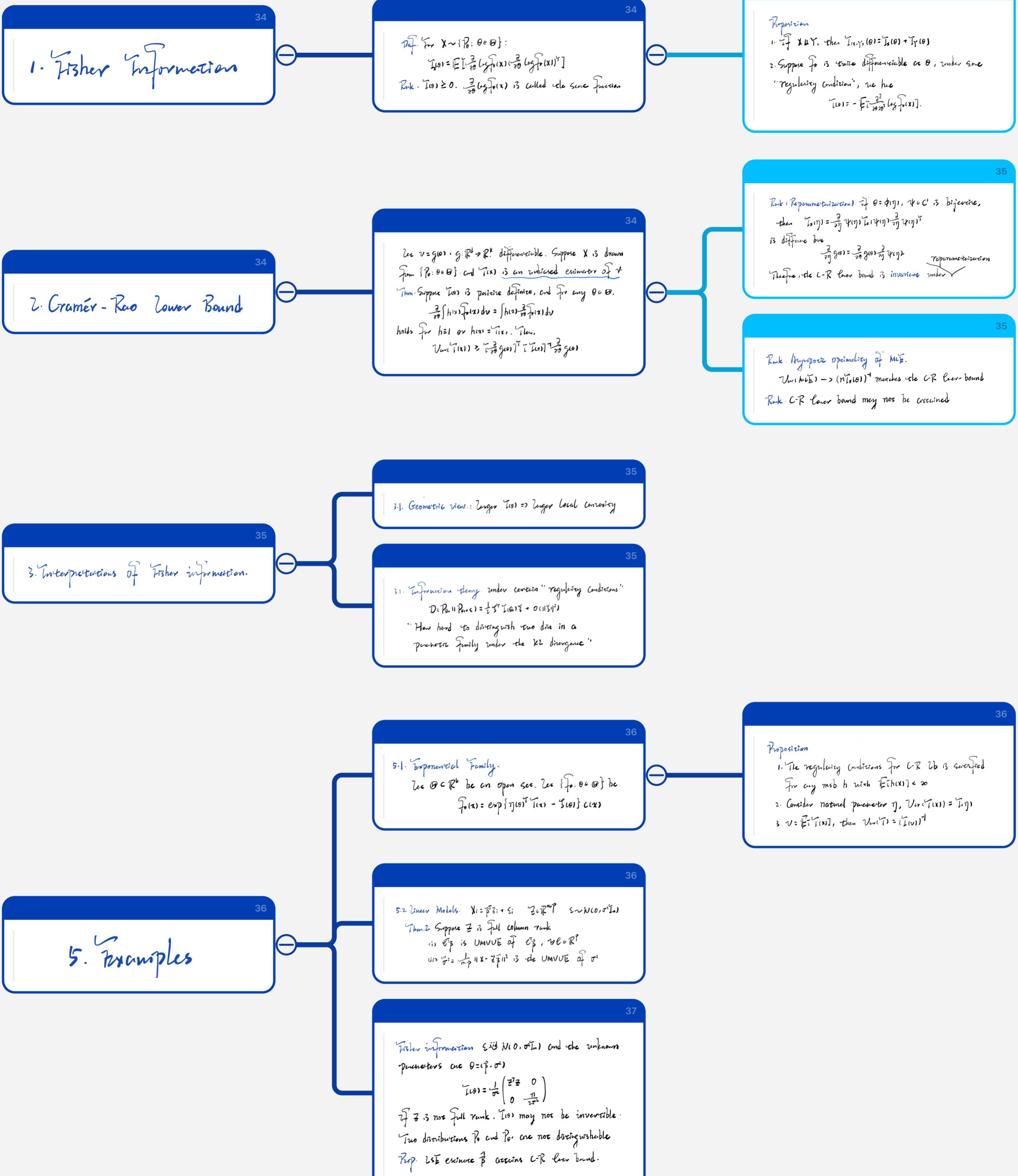
Method 3. Find UMVUE without knowing complete statistics

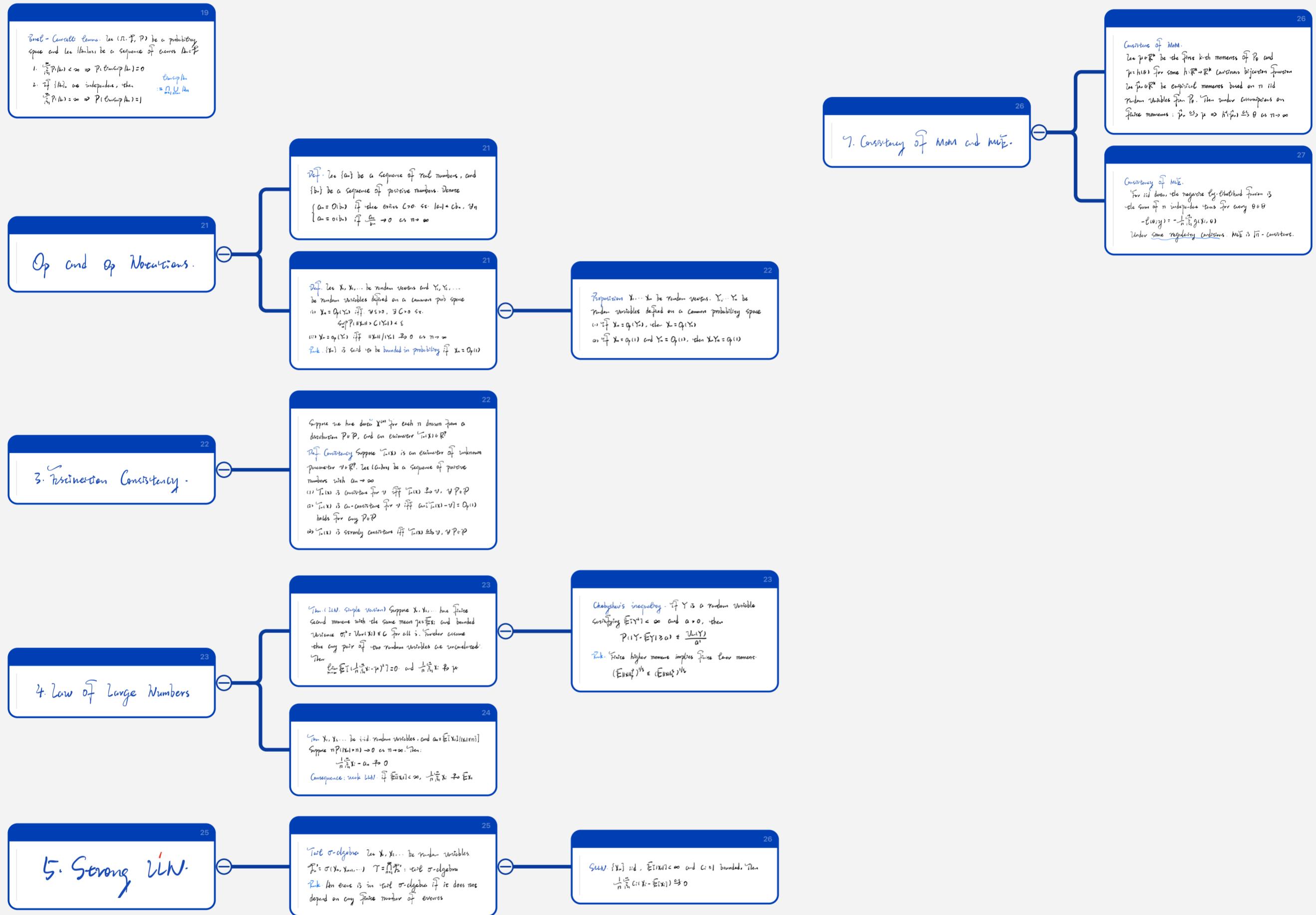
Theorem 2. Let $\mathcal{U} = \{U : E[U(X)] = 0, \text{Var}(U(X)) < \infty, \forall P \in \mathcal{P}\}$. T is unbiased for ν with $E[T(X)] < \infty$

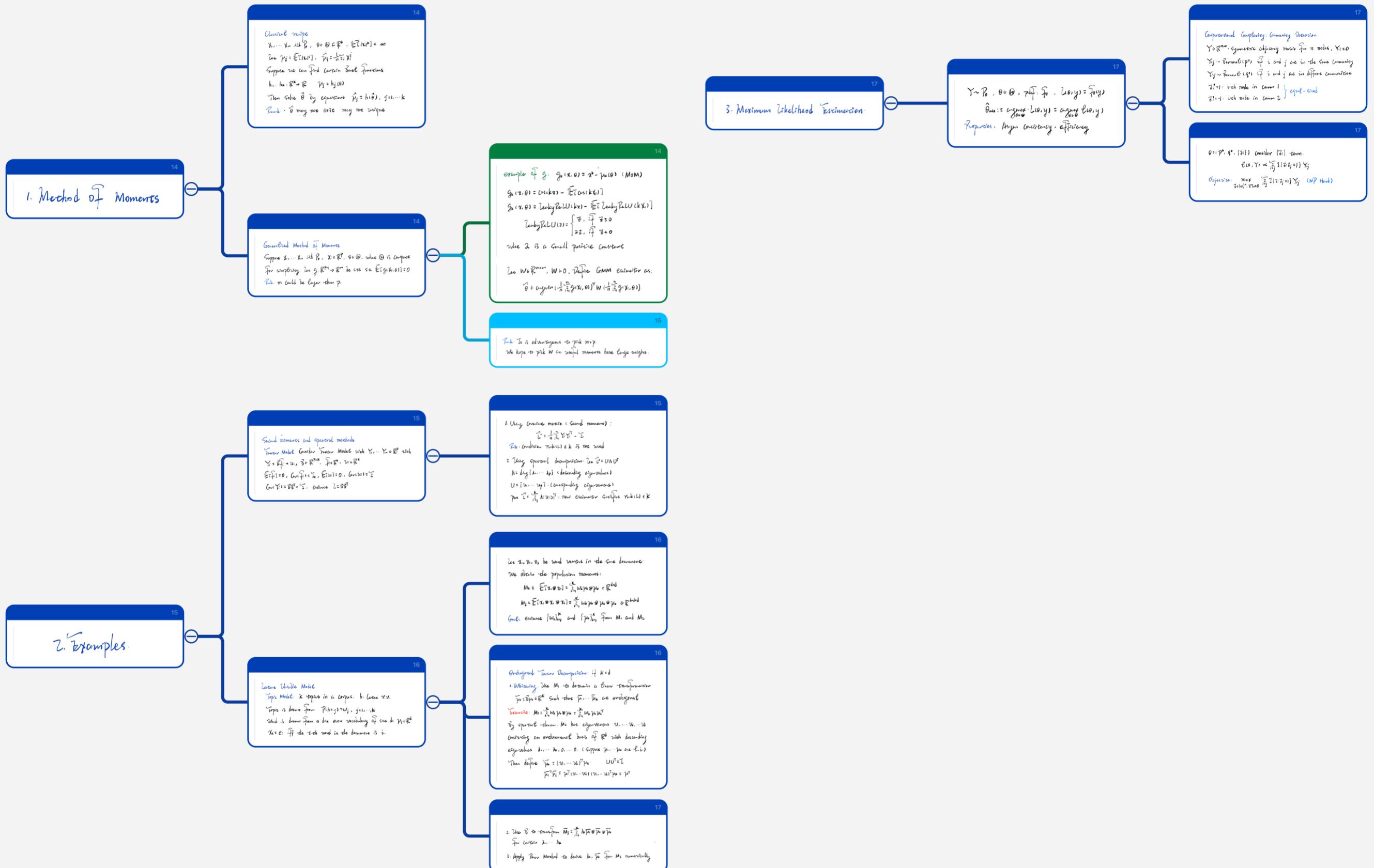
- $T(x)$ is UMVUE iff $E[T(X)U(X)] = 0, \forall U \in \mathcal{U}, \forall P \in \mathcal{P}$
- T is sufficient for P , i.e. $\tilde{U} = \mathcal{U} \cap \{g(T) : g \text{ Borel}\}$

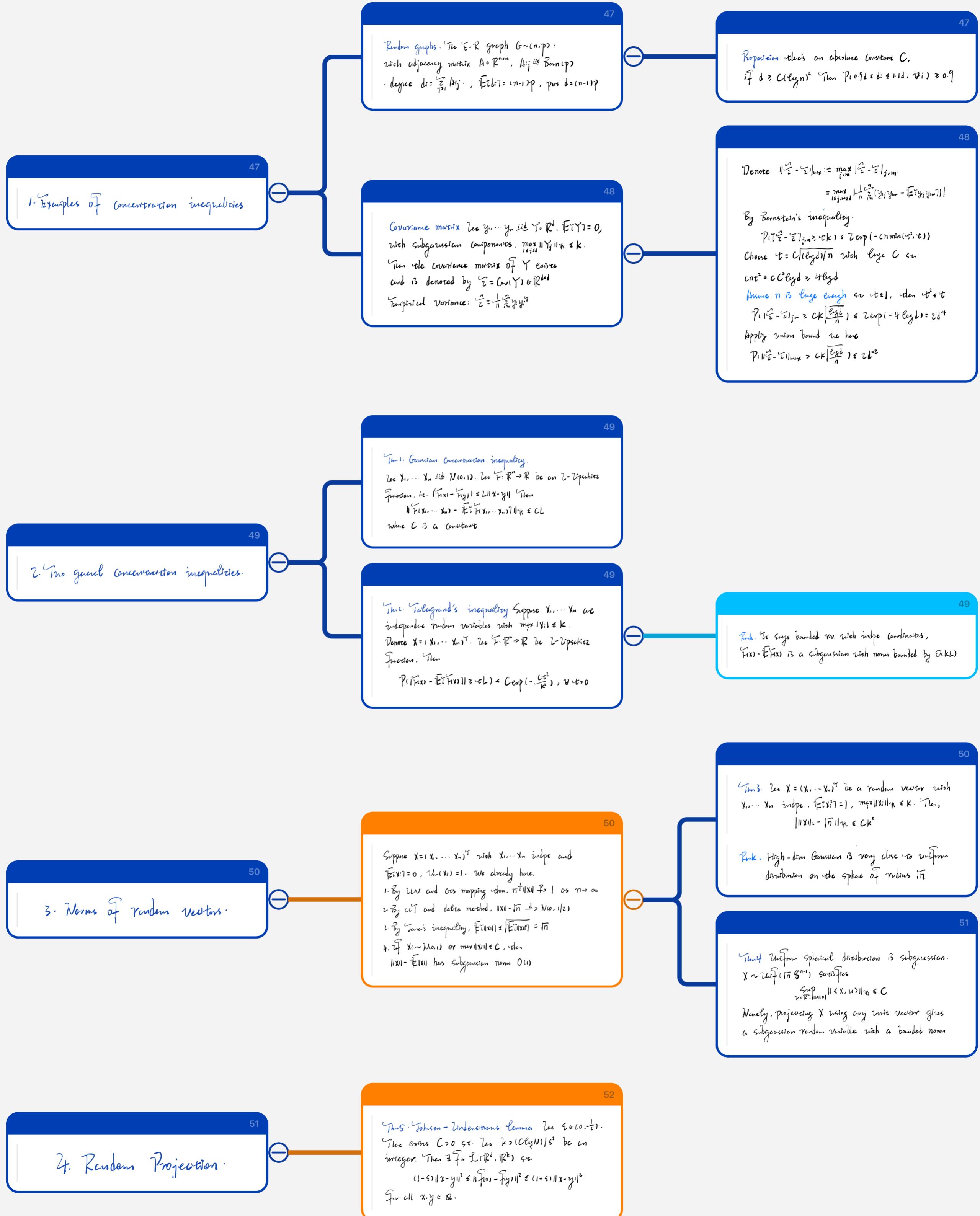
Then $T = h(\tilde{T})$ is UMVUE iff $E[T(X)U(X)] = 0, \forall U \in \tilde{\mathcal{U}}, \forall P \in \mathcal{P}$











54
2 Operator Norm of Subgaussian Random Matrix.

54
Thm 1. Let $A \in \mathbb{R}^{m,n}$ have independent subgaussian entries ($|A_{ij}|$ jointly indep. subgaussian with $\mathbb{E} A_{ij} = 0$). Then $\|A\| \leq Ck(\sqrt{m} + \sqrt{n} + k)$ with probability at least $1 - 2e^{-t^2}$ for certain constant C . and $k = \max_{ij} \|A_{ij}\|_{\infty}$.

54
Lemma 2. Let X_1, \dots, X_n be independent subgaussian random variables with $\mathbb{E} X_i = 0$, where $\sum_i \mathbb{E} |X_i|^p$ is also subgaussian with $\|\sum_i X_i\|_{\ell_p}^p \leq C \sum_i \mathbb{E} |X_i|^p$.

54
Rank $\|A\|_F$ is roughly of order $\sqrt{m} + \sqrt{n}$ with high probability. Here $n \geq m$, $\sigma^2 = \mathbb{E} A_{ij}^2$. Since $\|A\|_F^2 = \sum_{ij} (\mathbb{E} A_{ij}^2)^2$, we have $\frac{1}{m} \sum_{ij} \mathbb{E} A_{ij}^2 \geq \sigma^2$. The largest singular value is of order $\sqrt{n}\sigma$. with high prob \Rightarrow so does the largest singular value.

55
3. Covariance Truncation

55
Suppose $X_1, \dots, X_n \in \mathbb{R}^p$ iid random vectors with $\mathbb{E} X_i = 0$ and $\text{Cov}(X_i) = \Sigma$ exists. By LN, $\frac{1}{n} \sum_{i=1}^n X_i X_i^\top \xrightarrow{\text{P}} \Sigma$
Thm 2. Any $X \in \mathbb{R}^p$ is a subgaussian random vector, ($\sup_{x \in \mathbb{R}^p} \|\langle x, x \rangle\|_{\infty} \leq K$). Suppose X_1, \dots, X_n are iid copies of X , $\bar{X} = (X_1, \dots, X_n)^\top \in \mathbb{R}^{np}$. Then for every $t > 0$, with prob at least $1 - 2e^{-nt^2}$ $\|\frac{1}{n} \bar{X}^\top \bar{X} - \Sigma\| \leq \max\{\delta, \delta'\} \|\Sigma\|$ where $\delta = C \left(\sqrt{\frac{1}{n}} + \frac{t}{n} \right)$

55
Thm 2'. Covariance Truncation (HDP).
Let $X \in \mathbb{R}^p$ subgaussian ($\|\langle x, x \rangle\|_{\infty} \leq K \|\langle x, x \rangle\|_2$ for some $K \geq 1$, any $x \in \mathbb{R}^p$). Then $\mathbb{E} [\|\frac{1}{n} \bar{X}^\top \bar{X} - \Sigma\|] \leq CK \left(\sqrt{\frac{1}{n}} + \frac{2}{n} \right) \|\Sigma\|$.

1. Principle Component Analysis.

Suppose x_1, \dots, x_n iid $\mathcal{N}(0, I_p)$, $E[x_i] = 0$. $S = \text{Cov}(x_i)$ covers p_1, \dots, p_p be orthonormal eigenvectors of S with eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$. Then consider the empirical covariance $\tilde{S} = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$ with orthonormal eigenvectors $\tilde{p}_1, \dots, \tilde{p}_p$ and eigenvalues $\tilde{\lambda}_1 \geq \tilde{\lambda}_2 \geq \dots \geq \tilde{\lambda}_p \geq 0$.

1.1. Fixed Dimension

56

Theorem Consistency of PCA in fixed dimension. Suppose $\lambda_1 > \lambda_2$. Then as $n \rightarrow \infty$, we have $\max_{k \leq p} |\tilde{\lambda}_k - \lambda_k| \rightarrow 0$ and $\min_{S \subseteq \{1, \dots, p\}} \|S\tilde{u}_i - u_i\|_2 \rightarrow 0$. Rank $\lambda_1 > \lambda_2$ ensures λ_1 has multiplicity of 1.

Davis-Kahan Theorem (HDG) Let S and T be symmetric matrices with the same dimensions. Assume $\min_{j \neq i} |\lambda_j(S) - \lambda_j(T)| = \delta > 0$. Then the angle $\sin(\tilde{v}_i(S), \tilde{v}_i(T)) \leq \frac{2\|S-T\|}{\delta}$ where $v_i(T)$ gives the i -th largest eigenvector.

1.2. Growing Dimensions.

58

\tilde{p}_n, \tilde{u}_n depend on n : $x_i \sim \mathcal{N}(0, I_p)$. Spiked model: $\tilde{u}_n = (\lambda_1 - \lambda_2) u_1 u_1^T + \lambda_2 \tilde{u}_n$ eigenvalues are λ_1 and λ_2 with multiplicity $p-1$. Note that $\frac{1}{2} \min_{S \subseteq \{1, \dots, p\}} \|S\tilde{u}_i - u_i\|_2^2 = |-\langle \tilde{u}_i, u_i \rangle|$

Inconsistency of PCA in high dimensions: Assume $\|\tilde{u}_n\|_2 \leq C$ and $\lambda_1(\tilde{u}_n) - \lambda_2(\tilde{u}_n) \geq k > 0$ holds for certain constant C and k .

Rank population and sample principle direction are almost orthogonal in high dims.

Under Parameterized Regime

59

Proposition 1. Suppose x_1, x_2, \dots, x_n iid $\mathcal{N}(0, I_p)$. Assume $p = p_n$, $\frac{p_n}{n} = o(1)$ as $n \rightarrow \infty$, we have $E[\|\tilde{\beta} - \beta\|^2 | X] = (1 + o(1)) \frac{\sigma^2 p_n}{n}$, $E[(\tilde{\beta}^T \beta - \beta^T \beta)^2 | X] = (1 + o_p(1)) \frac{\sigma^2 p_n}{n}$ as $n \rightarrow \infty$. Rank $\tilde{\beta}$ increases approximately linear in p_n .

2. Linear Regressions in High Dims.

Consider (x_i, y_i) , $x_i \in \mathbb{R}^p$ and $i = 1, 2, \dots, n$. $y_i = x_i^T \beta + \epsilon_i$, $i = 1, 2, \dots, n$. $X = (x_1, \dots, x_n)^T \in \mathbb{R}^{np}$, $y = (y_1, \dots, y_n)^T \in \mathbb{R}^n$. $\epsilon = (\epsilon_1, \dots, \epsilon_n)^T \in \mathbb{R}^n$, $E[\epsilon \epsilon^T] = 0$, $\text{Var}(\epsilon_i) = \sigma^2$. Least squares: $\hat{\beta} = (X^T X)^{-1} X^T y$, Prediction: $\hat{f}(x_i) = x_i^T \hat{\beta}$.

Overparameterized Regime:

60

Minimum-Norm Interpolator: $\min_{\beta \in \mathbb{R}^p} \|\beta\|^2$ subject to $y = X\beta$ unbiased. Since $\text{rank}(X) = n$. The solution is $\hat{\beta} = X^T (X^T X)^{-1} y$. $\hat{\beta}^T \beta = \hat{\beta}^T \beta^*$, $\forall \beta \in N(X)$, $\text{rank}(X) = p$. Note that $\hat{\beta} \in C(X^T)$ and $C(X^T) \perp N(X)$. we have $\|\hat{\beta}^T \beta\|^2 = \|\hat{\beta}\|^2 + \|\beta\|^2 \geq \|\hat{\beta}\|^2$.

Proposition 2. Suppose $x_1, \dots, x_n \sim \mathcal{N}(0, I_p)$ and $p = p_n$ depends on n and $n/p_n = o(1)$ as $n \rightarrow \infty$. We have $\text{rank}(X) = n$ with probability $1 - o(1)$ and $E[\|\hat{\beta} - \beta\|^2 | X] = (1 + o_p(1)) \cdot [(1 - \frac{n}{p_n}) \|\beta\|^2 + \frac{\sigma^2 n}{p_n}]$, $E[(\hat{\beta}^T \beta - \beta^T \beta)^2 | X] = (1 + o_p(1)) \cdot [(1 - \frac{n}{p_n}) \|\beta\|^2 + \frac{\sigma^2 n}{p_n}]$. Bias Variance: Smaller when $p_n > n$.

1. Motivation

63
 1.1. High-D Perspective
 $\text{Let } \mathbf{x} = \mu + \mathbf{z}, \mathbf{z} \sim N_p(0, I_p)$
 $P(|\langle \mu, \mathbf{z} \rangle| > t) \leq \exp(-t^2/\|\mu\|^2)$
 When p is large:
 $\frac{\|\mu\|}{\sqrt{p}} = O(1) \Rightarrow |\langle \mu, \mathbf{z} \rangle| = O(\sqrt{p})$
 $\|\mathbf{x}\|^2 = \|\mu\|^2 + \|\mathbf{z}\|^2 + 2\langle \mu, \mathbf{z} \rangle \approx \|\mu\|^2 + p$
 $\|\mathbf{x}\| > \|\mu\|$ with high probability

63
 1.2. Bias-Variance trade-off perspective
 $R(p) = \mathbb{E}[\|\mathbf{x} - \mu\|^2 | \mathbf{X} \sim N(\mu, I_p)]$
 Consider $\hat{\mu} = (1-\varepsilon)\mathbf{x}$
 $\mathbb{E}[\|\hat{\mu} - \mu\|^2] = \varepsilon^2\|\mu\|^2 + (1-\varepsilon)^2 p$
 $= (1-2\varepsilon)p + \varepsilon^2(\|\mu\|^2 + p)$
 $= (1-2\varepsilon)p + O(\varepsilon)$

63
 Goal: $\min_{\varepsilon \in [0,1]} \{ \mathbb{E}[\|\hat{\mu}\|^2] + (1-\varepsilon)^2 p \}$.
 $\|\mu\|$ is unknown, but we can use $\|\mathbf{x}\|^2 \approx \|\mu\|^2 + p$ when p is large to estimate $\|\mu\|^2$

2. Shrinkage Estimators

64
 James-Stein Estimator: noise $\hat{\mu}_J = \mathbf{x}$
 $\hat{\mu}_{JS} = (1 - \frac{p-2}{\|\mathbf{x}\|^2})\mathbf{x}$ $\hat{\mu}_{SS} = (1 - \frac{p-2}{\|\mathbf{x}\|^2}) + \mathbf{x}$

64
 Rule: Shrinkage to 0 is not the only direction where it works. For any $c \in \mathbb{R}^p$, we can make
 $\hat{\mu}_S = (1 - \frac{p-2}{\|\mathbf{x}-c\|^2})(\mathbf{x}-c) + c$
 $\hat{\mu}_{SS} = (1 - \frac{p-2}{\|\mathbf{x}-c\|^2}) + (\mathbf{x}-c) + c$

3.2. Shrinkage under Superharmonic Function.

66
 Def. For $f \in C_c^2(\mathbb{R}^p)$, denote the Laplace operator by $\Delta f(x) = \sum_{i=1}^p \frac{\partial^2 f}{\partial x_i^2}(x)$. We say f is superharmonic iff $\Delta f(x) \leq 0, \forall x \in \mathbb{R}^p$.

66
 Thm 3: Let $\tilde{f} \in C_c^2(\mathbb{R}^p)$, $\tilde{f} > 0$, \tilde{f} is superharmonic iff $\mathbb{E}[f(\frac{1}{\tilde{f}(x)} \partial_x^2 \tilde{f}(x))] \leq 0$, $\mathbb{E}[\nabla \log \tilde{f}(x)]^2 \leq 0$
 we have $\mathbb{E}[\|\mathbf{x} + \nabla \log \tilde{f}(x) - \mu\|^2] = p + 4\mathbb{E}[\frac{\Delta \tilde{f}(x)}{\tilde{f}(x)}] \leq p$

66
 Rule: \mathbf{x} is a minimax estimator
 $\Rightarrow \mathbf{x} + \nabla \log \tilde{f}(x)$ is also a minimax estimator

4. Diffusion Model:

67
 Forward Process: $\mathbf{x}_t = \sqrt{1-\sigma_t^2} \mathbf{x}_{t-1} + \sigma_t \mathbf{z}_t$
 $\mathbf{z}_t \sim N(0, I_p)$.
 Score Function $\nabla \log p_t$ of \mathbf{x}_t
 $\mathbb{E}[\mathbf{x}_{t+1} | \mathbf{x}_t = \mathbf{x}_t] = \frac{1}{\sqrt{1-\sigma_t^2}} \mathbf{x}_t + \frac{\sigma_t^2}{\sqrt{1-\sigma_t^2}} \nabla \log p_t(\mathbf{x}_t)$
 $\mathbf{x}_{t+1} | \mathbf{x}_t = \mathbf{x}_t \sim N(\frac{1}{\sqrt{1-\sigma_t^2}} \mathbf{x}_t + \frac{\sigma_t^2}{\sqrt{1-\sigma_t^2}} \nabla \log p_t(\mathbf{x}_t), \sigma_t^2 I_p)$

67
 Reconstruction: $\mathbf{x}_1 = \frac{1}{\sqrt{1-\sigma_0^2}} \mathbf{x}_0 + \frac{\sigma_0^2}{\sqrt{1-\sigma_0^2}} S(\theta, \mathbf{x}_0) + \sigma_0 \mathbf{z}_0$
 $S(\theta, \mathbf{x}_0)$ is trained to approximate $\nabla \log p_0(\mathbf{x}_0)$:
 $\min_{\theta} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}_0, \mathbf{z}_0), \mathbf{x}_0 \sim p(\mathbf{x}_0, \mathbf{z}_0)} [\sigma_0 \|\nabla \log p_0(\mathbf{x}_0) - S(\theta, \mathbf{x}_0, \mathbf{z}_0)\|^2]$

Lecture 6. Statistical Decision Theory and Sufficiency

Basic Terminology

$$\mathcal{D}_{\text{obs}} \sim (\Omega, \mathcal{P}, P)$$

X_1, \dots, X_n : random variables $X = (X_1, \dots, X_n)^T \in \mathbb{R}^{n \times p}$

Example. Graph \mathcal{D}_{obs} I_2 -R model $G(n, p)$:

Graph G , n vertices, A : adjacency matrix.

$$A_{ij} \sim \text{Bernoulli}(p), \text{ then } \mathbb{E}[A] = p I_n I_n^T$$

$A \approx \mathbb{E}[A]$ under some conditions. $\text{I}(A) \approx \text{I}(\mathbb{E}[A])$

Classifying Statistical Models

Parametric Family. $\mathcal{P} \subseteq \{P_\theta : \theta \in \Theta \subset \mathbb{R}^p\}$

e.g. family of densities: if $P_\theta \ll \nu$. If $\theta \in \Theta$, then we only need to consider $\{\frac{dP_\theta}{d\nu} : \theta \in \Theta\}$, particularly pdf / pmf

Example. $\mathcal{P}_0 \mathcal{P} = \{N(\mu, \sigma^2) : \mu \in \mathbb{R}, \sigma^2 > 0\}$. $\theta = (\mu, \sigma^2)^T$

Example (linear model) $\mathcal{P} = \{N(\beta^T X, \sigma^2) : \beta \in \mathbb{R}^p, \sigma^2 > 0\}$ $X \in \mathbb{R}^p$

Nonparametric Family. Θ is a infinite dimensional space

rk. Statistical Methods often use a finite number of estimates to approximate the unknown quantity

Example (Blowing estimator): Each distribution is abs cont and the set of pdf is $\{f(x) : |f'(x)| < L\}$ for some $L > 0$

Semi-parametric family: True parameters of interest + Inverse Nuisance Parameters

Example: Linear model with unknown dist of ε

Identifiability:

Def. $P = \{\theta_0 + B\}$ is identifiable iff $P_{\theta_1} \neq P_{\theta_2}$ for $\theta_1 \neq \theta_2$.

Example (Factor Model) Consider $Y \in \mathbb{R}^p$ has the structure

$$Y = Bf + u, \text{ where } f \sim N(0, I_r), u \sim N(0, \Sigma).$$

Σ : known, $B \in \mathbb{R}^{p \times r}$: unknown.

equivalently $Y \sim (0, BB^T + \Sigma)$.

Since $BR(BR)^T = BB^T$ for any orthogonal matrix $R \in \mathbb{R}^{r \times r}$

this model is not identifiable.

Example: Multinomial regression. In classification problem:

$$\hat{P}_k = \frac{\exp(\hat{\beta}_k^T X)}{\sum_{j=1}^{k-1} \exp(\hat{\beta}_j^T X)} \quad \text{where } k=1, 2, \dots, K-1$$

$(\hat{\beta}_1, \dots, \hat{\beta}_K)^T$ and $(\hat{\beta}_1 + c, \dots, \hat{\beta}_K + c)$ specifies the same model

$$\text{To avoid this: } p_k = \frac{\exp(\beta_k^T x)}{\sum_{j=1}^k \exp(\beta_j^T x) + 1}$$

Statistics

Def. Let $X \in \mathbb{R}^{n \times p}$ be the data. For any nsb function T , we call $T(X)$ a statistic.

$\sigma(T(X)) \subset \sigma(X)$: $T(X)$ summarizes information from X .

Example Sample mean \bar{X} and Sample Variance S^2

Order Statistics $X_{(1)}, \dots, X_{(n)}$. pdf for $X_{(i)}$ is:

$$g_{(i)}(x) = \frac{n!}{(i-1)! (n-i)!} [F(x)]^{i-1} [1 - F(x)]^{n-i} f(x)$$

Exponential Family:

Def. Dst with pdf $f_\theta(w) = \exp\{\eta(\theta)^T T(w) - \zeta(\theta)\} h(w)$, $\forall w \in \Omega$

Canonical Form: $f_\eta(w) = \exp\{\eta^T T(w) - \zeta(\eta)\} h(w)$, $w \in \Omega$

e.g.: Binomial, Normal, Dsts subsumed in exponential family

Properties: (1) Marginal / Conditional dst remains exponential

(2) For Borel Function F : $\int |F| dP_\theta < \infty$

$\int f(w) \exp(\eta^T T(w)) h(w) d\nu(w) \in C^\infty$ Probability Measure induced by the random variable.

Con. $m(\tau) \in C^\infty_{\text{on } \Omega}$

Basic Terminologies in Statistical Decision Theory.

- Decision Rule: $T: (\mathbb{R}^k, \mathcal{B}_k) \rightarrow (\mathcal{A}, \mathcal{P}_{\mathcal{A}})$
- Loss Function: $L(P, T(x))$.
- Risk Function: $R_T(P) = \bar{\mathbb{E}}_P [L(P, T(x))]$. For parameter θ , it can be written as $R(\theta)$
e.g. $MS\bar{E}_\theta(\theta) = \bar{\mathbb{E}}_\theta [(\hat{\theta}(x) - \theta)^2] = \int (\hat{\theta}(x) - \theta)^2 dP_\theta(x)$

Optimal Rule

Def. T^* is an optimal rule if

$R_{T^*}(P) \leq R_T(P), \forall P \in \mathcal{P}, \forall T \in \mathcal{G}$ {Class of allowable decisions}

True optimal if $R_{T^*}(P) < CR_T(P)$ for some C

Randomized decision rules.

For every x , $S(x, \cdot)$ is a probability measure on $(\mathcal{A}, \mathcal{P}_{\mathcal{A}})$

Often, $S = \sum_j^n \mathbf{1}_{\{Z=j\}} T_j$ where Z is a discrete random variable

Risk of randomized rule: $R(S, P) = \bar{\mathbb{E}}_{X \sim P} [\bar{\mathbb{E}}_{T \sim S(x)} [L(T, A)]]$

e.g. Stochastic Gradient Descent

Admissibility Let G be a class of decision rules. $T \in G$

is G -admissible iff $\forall S \in G, R_S(P) \geq R_T(P), \forall P \in \mathcal{P}$

Main Approaches for optimality:

- Bayes rule: Given a prior dist π over \mathcal{P} , T_π minimizes

$$R_{T_\pi}(\pi) = \int R_T(p) d\pi(p)$$

- Minimax rule: $T_\pi = \arg\min_{T_\pi} (\sup_p R_T(p))$

Estimation.

For iid x_1, \dots, x_n , with \bar{x} . Use squared loss

$$L(p, a) = (\theta - a)^2, \text{ then}$$

$$\bar{R}_{\bar{x}}(p) = \frac{\sigma^2}{n} \text{ where } \sigma^2 = \text{Var}(x_i)$$

$$\text{Bias: } b_{T_\pi}(p) = \bar{E}_p[T_\pi(x)] - \theta$$

Hypothesis Test

$$H_0: P \in \mathcal{P}_0 \text{ vs. } H_1: P \in \mathcal{P}_1$$

Let $T_\pi(x) \in \{0, 1\}$, $T_\pi(x)$ must have the form $T_\pi(x) = \{x : T_\pi(x) = 1\}$.

where C is called the rejection region. Consider the loss

$$L(p, j) = 0 \text{ if } p \in P_j, \text{ and } 1 \text{ otherwise}$$

$$\text{Then } R_{T_\pi}(p) = \begin{cases} p(T_\pi(x) = 1) = p(x \in C), & p \in \mathcal{P}_0 \\ p(T_\pi(x) = 0) = p(x \notin C), & p \in \mathcal{P}_1 \end{cases}$$

Sufficient Statistics

Def. $T(x)$ is sufficient if the conditional dist
of $X | T(x) = t$ does not depend on P or θ for any t

Factorization Thm: Suppose $P = \{P_\theta : \theta \in \Theta\}$ is dominated by
the Lebesgue measure ν . Then $T(x)$ is sufficient iff

exists functions h, g_θ s.t. $\frac{dP_\theta}{d\nu}(x) = g_\theta(T(x)) h(x)$.

Example: exponential family $f_\theta(w) = \underbrace{\exp\{\langle \eta(\theta) \rangle^T T(w) - \psi(\theta)\}}_{g_\theta(T(x))} h(w)$

Connection to information theory:

$T(x)$ is sufficient $\Leftrightarrow I(X; \theta) = I(T(x); \theta)$

when θ is treated as a random variable

* Data Processing Inequality.

$X \rightarrow Y \rightarrow Z$ Markov Chain $X \perp\!\!\!\perp Z | Y$, then

$$I(X, Z) \leq I(X; Y).$$

Lecture 7. Bias-Variance Trade-off.

1. Example: Sample mean. $X_1, \dots, X_n \text{ iid } (\mu, \sigma^2)$

Suppose we have some prior estimation $\hat{\mu}_0$ for μ .

Convince $\hat{\mu} = 0.2\hat{\mu}_0 + 0.8\bar{x}$, then

$$\text{Bias}(\hat{\mu}) = 0.2(\hat{\mu}_0 - \mu)$$

$$\text{Var}(\hat{\mu}) = 0.64 \frac{\sigma^2}{n}$$

$$R(\mu, \hat{\mu}) = 0.04(\hat{\mu}_0 - \mu)^2 + 0.64 \frac{\sigma^2}{n}.$$

Conclusion: $\hat{\mu}$ is better than \bar{x} if $\hat{\mu}_0$ is close to μ .

2. Ridge Regression.

Obs: (X_i, y_i) , $X_i \in \mathbb{R}^p$, $y_i \in \mathbb{R}$ with model

$$y_i = X_i^\top \beta^* + \varepsilon$$

Ridge regression solves ignore the $\frac{1}{n}$ term in lecture note

$$\min_{\beta} \underbrace{\sum_{i=1}^n (y_i - X_i^\top \beta)^2}_{\text{MSE}} + \lambda \|\beta\|_2^2 \quad (1)$$

$$\Rightarrow \hat{\beta}_1 = (X^\top X + \lambda I_p)^{-1} X^\top y \quad (X = (X_1, \dots, X_n)^\top \in \mathbb{R}^{n \times p})$$

$$= (X^\top X + \lambda I_p)^{-1} X^\top (X \beta^* + \varepsilon)$$

$$\text{MSE} = \mathbb{E}[\|\hat{\beta}_1 - \beta^*\|^2] = \mathbb{E}[\|\hat{\beta}_1 - \mathbb{E}\hat{\beta}_1\|_2^2] + \|\mathbb{E}\hat{\beta}_1 - \beta^*\|_2^2$$

$$\text{Bias} = O(\lambda), \quad \sigma^2 \text{Tr}((X^\top X)^{-1}) - \text{Var} = O(\lambda)$$

$$\lambda \uparrow \Rightarrow \text{Var} \downarrow \text{Bias} \uparrow$$

Exercise:

Bias:

$$\begin{aligned}\tilde{\mathbb{E}}[\tilde{\beta}_\lambda] &= (X^T X + \lambda I_p)^{-1} X^T X \beta^* \\ &= (X^T X + \lambda I_p)^{-1} (X^T X + \lambda I_p - \lambda I_p) \beta^* \\ &= \beta^* + \lambda (X^T X + \lambda I_p)^{-1} \beta^*\end{aligned}$$

$$\text{Let } M = (X^T X + \lambda I_p) \in \mathbb{R}^{P \times P}$$

$$\begin{aligned}\|M^{-1}\| &= \max_{x \neq 0} \frac{\|M^{-1}x\|}{\|x\|} \\ &= \max_{x \neq 0} \frac{\|x\|}{\|Mx\|} \\ &= \left(\min_{x \neq 0} \frac{\|Mx\|}{\|x\|} \right)^{-1} \\ &= (\mu_{\min}^*)^{-1} \leq (\mu_{\min})^{-1}\end{aligned}$$

where μ_{\min}^* is the minimum eigenvalue of M .

μ_{\min} is the smallest eigenvalue of $X^T X$

$$\begin{aligned}\text{Bias} &:= \|\lambda (X^T X + \lambda I_p)^{-1} \beta^*\| \\ &\leq \lambda \|(\lambda (X^T X + \lambda I_p)^{-1})\| \|\beta^*\| \\ &\leq \lambda (\mu_{\min})^{-1} \|\beta^*\| \quad \Rightarrow \text{Bias} = O(\lambda)\end{aligned}$$

Variance

$$\hat{\beta}_\lambda = (X^T X + \lambda I_p)^{-1} X^T X \beta^* + (X^T X + \lambda I_p)^{-1} X^T \varepsilon$$

$$\hat{\beta}_\lambda - \tilde{\mathbb{E}}[\hat{\beta}_\lambda] = (X^T X + \lambda I_p)^{-1} X^T \varepsilon$$

$$\begin{aligned}\tilde{\mathbb{E}}[\|\hat{\beta}_\lambda - \tilde{\mathbb{E}}[\hat{\beta}_\lambda]\|^2] &= \tilde{\mathbb{E}}[\varepsilon^T X (X^T X + \lambda I_p)^{-2} X^T \varepsilon] \\ &= \sigma^2 \text{Tr}(X (X^T X + \lambda I_p)^{-2} X^T)\end{aligned}$$

Spectral Decomposition: $X^T X = U \Lambda U^T$. $X^T X + \lambda I_p = U (\Lambda + \lambda I_p) U^T$

$$\begin{aligned}
& X(X^T X + \lambda I_p)^{-2} X^T \\
&= U \Lambda U^T U (\Lambda + \Lambda^T \Lambda)^{-2} U^T U \Lambda U^T \\
&= U \Lambda (\Lambda + \Lambda^T \Lambda) \Lambda U^T \\
&= U \text{diag} \left\{ \frac{\mu_i}{(\mu_i + \lambda)^2} \right\} U^T
\end{aligned}$$

where μ_i are eigenvalues of $X^T X$

$$\text{Var}(\hat{\beta}_\lambda) = \sigma^2 \sum_{i=1}^n \frac{\mu_i}{(\mu_i + \lambda)^2} \quad \frac{\mu_i}{(\mu_i + \lambda)^2} =$$

$$\text{WIS: } \sigma^2 \text{Tr}((X^T X)^{-1}) - \text{Var} = O(\lambda)$$

3. General Bias-Variance Tradeoff

Consider model $Y = f(x) + \varepsilon$, $E\varepsilon = 0$, $\text{Var}(\varepsilon) = \sigma^2$ $x \in \mathbb{R}^p$

KNN estimator: $\hat{f}_k(x_0) = \frac{1}{k} \sum_{i=1}^k y_{x(i)}$

where $x_{(1)}, \dots, x_{(k)}$ are k nearest points to x_0

$$\mathbb{E}[(Y - \hat{f}_k(x_0))^2 | X = x_0] = \sigma^2 + \underbrace{\left(\hat{f}(x_0) - \frac{1}{k} \sum_{i=1}^k \hat{f}(x_{(i)}) \right)^2}_{\text{Bias}} + \underbrace{\frac{\sigma^2}{k}}_{\text{Variance}}$$

4. Overparametrization

$$\hat{f}(x) = x^\top \hat{\beta}$$

$$\text{Var}(x^\top \hat{\beta}) = \sigma^2 x^\top (X^\top X)^{-1} x$$

If x_0 has zero mean and identity covariance matrix.

$$\text{Var}(\hat{f}(x_0)) = \sigma^2 \mathbb{E}[x_0^\top (X^\top X)^{-1} x_0] = \sigma^2 \text{tr}((X^\top X)^{-1}) + \|\beta\|^2$$

$$\begin{aligned} \text{Var}(\hat{f}(x_0)) &= \sigma^2 \mathbb{E}[\sigma^2 x_0^\top (X^\top X)^{-1} x_0] + \text{Var}(x_0^\top \hat{\beta}) \\ \mathbb{E}[\hat{f}(x_0) | x_0] &= x_0^\top \hat{\beta} \quad \hat{\beta} = (X^\top X)^{-1} X^\top y \perp \!\!\! \perp x_0 \\ \text{Var}(\hat{f}(x_0)) &= \mathbb{E}[\sigma^2 x_0^\top (X^\top X)^{-1} x_0] + \text{Var}(x_0^\top \hat{\beta}) \\ &= \sigma^2 \mathbb{E}[\text{tr}((X^\top X)^{-1} x_0^\top x_0)] + \hat{\beta}^\top \hat{\beta} \\ &= \sigma^2 \text{tr}((X^\top X)^{-1}) + \|\hat{\beta}\|^2 \\ &= n\sigma^2 \text{tr}((X^\top X)^{-1}) + \|\beta\|^2 \end{aligned}$$

The variance generally increases as p increases

$$\text{Typically: } X^\top X = \sum_{i=1}^n x_i^\top x_i \Rightarrow \mathbb{E}[X^\top X] = n \sum_{i=1}^n z_i z_i^\top \in \mathbb{R}^{p \times p}$$

If the smallest eigenvalue is bounded away from 0.

then the variance is $O(p/n)$

Double Descent and min-norm solution.

When $p > n$: $\min_{\beta} \|\beta\|_2^2$ (2) (ridgeless estimator)
subject to $X\beta = y$

Suppose $\text{rank}(X^\top X) = n$, then it has a unique solution

$$\hat{\beta}_{0+} = \lim_{\lambda \rightarrow 0+} \hat{\beta}_\lambda = X^\top (X^\top X)^{-1}$$

5. The implicit regularization of gradient descent

$$GD: \hat{\beta}^{(t+1)} = \hat{\beta}^{(t)} - \eta \nabla L(\hat{\beta}),$$

Starting at $\hat{\beta}^{(0)} = 0$, solving $L(\hat{\beta}) = \|Y - X\hat{\beta}\|^2$

Thm. Suppose $p > n$, and $\text{rank}(X^T X) = n$. Then,

if $\eta < [2 \lambda_{\max}(X^T X)^{-1}]$, then $\hat{\beta}^{(t)}$ converges, and

$$\lim_{t \rightarrow \infty} \hat{\beta}^{(t)} = \hat{\beta}_{\text{opt}}$$

Rank. GD prefers the solution. Because

$$\nabla L(\hat{\beta}) = 2X^T(Y - X\hat{\beta}) \in C(X^T) \Rightarrow \hat{\beta}^{(t)} \in C(X^T)$$

$$\Rightarrow \lim_{n \rightarrow \infty} \hat{\beta}^{(t)} \in C(X^T)$$

and $\hat{\beta}_{\text{opt}}$ is the unique solution in $C(X^T)$

Thm (GD for overparameterized logistic regression)

Consider logistic loss and linear separable data.

From any initializer $\hat{\beta}^{(0)} \in \mathbb{R}^p$, the gradient iterate

$$\hat{\beta}^{(k+1)} = \hat{\beta}^{(k)} - \eta \nabla L(\hat{\beta}^{(k)}) \text{ satisfies } \hat{\beta}^{(k)} = \hat{\beta} \log k + \Delta^{(k)}$$

where $\|\Delta^{(k)}\|_2 = O(\log \log k)$ and

$$\hat{\beta} = \arg \min \|\hat{\beta}\|_2^2 \text{ subject to } x_i^T \hat{\beta} \geq 1 \text{ for } i=1, \dots, n.$$

Remark: The DIRECTION of $\hat{\beta}^{(k)}$ converge to $\hat{\beta}$.

and in logistic model, only direction matters.

Lesson 7 Exercises

Exercise 1: Solve ridge regression

$$\min_{\beta} \mathcal{L}(\beta) = \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \beta)^2 + \lambda \|\beta\|_2^2 \quad (1)$$

$$\begin{aligned} \frac{\partial \mathcal{L}(\beta)}{\partial \beta} &= \frac{1}{n} \sum_{i=1}^n 2(-\mathbf{x}_i)(y_i - \mathbf{x}_i^\top \beta) + 2\lambda \beta \\ &= \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i(y_i - \mathbf{x}_i^\top \beta) + 2\lambda \beta \end{aligned}$$

$$\begin{aligned} &= \frac{1}{n} (\sum_i \mathbf{x}_i y_i - \sum_i \mathbf{x}_i^\top \beta \mathbf{x}_i) + 2\lambda \beta \\ &= \frac{1}{n} (\mathbf{x}^\top \mathbf{y} - \mathbf{x}^\top \mathbf{x} \beta) + 2\lambda \beta \\ &= \frac{1}{n} \mathbf{x}^\top \mathbf{y} + (\frac{1}{n} \mathbf{x}^\top \mathbf{x} + 2\lambda) \beta \end{aligned}$$

$$\frac{\partial^2 \mathcal{L}(\beta)}{\partial \beta^2} = \frac{1}{n} \mathbf{x}^\top \mathbf{x} + 2\lambda I \succeq 0$$

$$\frac{\partial \mathcal{L}(\beta)}{\partial \beta} = 0 \Rightarrow \hat{\beta} = (\frac{1}{n} \mathbf{x}^\top \mathbf{x} + 2\lambda I)^{-1} \frac{1}{n} \mathbf{x}^\top \mathbf{y}$$

Exercise 2. Prove the ridge less solution

$$\hat{\beta}^* = \lim_{\lambda \rightarrow 0^+} \hat{\beta}_\lambda = X^T(XX^T)^{-1}y$$

From we prove two optimization problems are the same when $\lambda \rightarrow 0$

$$\min_{\beta} L(\beta) = \|Y - X\beta\|^2 + \lambda \|\beta\|^2 \quad \text{Set } \mathcal{B} = \{\beta : Y - X\beta = 0\}$$

For any $\beta \in \mathcal{B}^\perp$ and $\tilde{\beta} \in \mathcal{B}$

$$L(\beta) - L(\tilde{\beta}) = \|Y - X\beta\|^2 + \lambda(\|\beta\| - \|\tilde{\beta}\|^2) > 0 \quad \text{as } \lambda \rightarrow 0$$

Then the minimizer must satisfy $Y - X\beta = 0$

Then we solve: $\min_{\beta} \|\beta\|_2^2 \quad (2)$

subject to $X\beta = y$

pre $h(\beta) = X\beta - y$. Let $L(\beta, \lambda)$ be the lagrangian

$$L(\beta, \lambda) = \|\beta\|_2^2 - \lambda^T(X\beta - y) \quad \text{for some } \lambda \in \mathbb{R}^n$$

$$\nabla L(\beta, \lambda) = \begin{pmatrix} X\beta - \lambda^T \\ X\beta - y \end{pmatrix} = 0$$

$$\Rightarrow \beta = \frac{1}{2}X^T\lambda$$

$$X\beta = \frac{1}{2}XX^T\lambda - y = 0$$

$$\Rightarrow \lambda = (\frac{1}{2}XX^T)^{-1}y \quad \text{since } XX^T \text{ invertible.}$$

$$\text{Then we have } \beta = X^T(XX^T)^{-1}y$$

Lecture 8. Basic Estimation Methods

1. Method of Moments

Classical recipe

X_1, \dots, X_n iid P_θ , $\theta \in \Theta \subset \mathbb{R}^k$, $\mathbb{E}[|X_i|^k] < \infty$

Let $\mu_j = \mathbb{E}[X_i^j]$, $\hat{\mu}_j = \frac{1}{n} \sum_i X_i^j$

Suppose we can find certain Borel functions

$h_1, \dots, h_k: \mathbb{R}^k \rightarrow \mathbb{R}$ $\mu_j = h_j(\theta)$

Then solve $\hat{\theta}$ by equations $\hat{\mu}_j = h_j(\hat{\theta})$, $j=1, \dots, k$

Remark: $\hat{\theta}$ may not exist may not unique

Generalized Method of Moments

Suppose X_1, \dots, X_n iid P_θ , $X_i \in \mathbb{R}^p$, $\theta \in \Theta$, where Θ is compact

For simplicity. Let $g: \mathbb{R}^{p+1} \rightarrow \mathbb{R}^m$ be cts s.t. $\mathbb{E}[g(X_i, \theta)] = 0$

Rank m could be larger than p .

example of g : $g_k(x, \theta) = x^k - \mu_k(\theta)$ (MoM)

$$g_k(x, \theta) = \cos(kx) - \mathbb{E}[\cos(kx_i)]$$

$$g_k(x, \theta) = \text{LeakyReLU}(kx) - \mathbb{E}[\text{LeakyReLU}(kx_i)]$$

$$\text{LeakyReLU}(z) = \begin{cases} z, & \text{if } z \geq 0 \\ \alpha z, & \text{if } z < 0 \end{cases}$$

where α is a small positive constant

Let $W \in \mathbb{R}^{m \times m}$, $W > 0$. Define GMM estimator as:

$$\hat{\theta} \leftarrow \operatorname{argmin} \left(\frac{1}{n} \sum_{i=1}^n g(x_i, \theta) \right)^T W \left(\frac{1}{n} \sum_{i=1}^n g(x_i, \theta) \right).$$

Rank: It is advantageous to pick $m > p$.

We hope to pick W so useful moments have large weights.

2. Examples

Second moments and spectral methods

Factor Model: Consider Factor Model with $Y_1, \dots, Y_n \in \mathbb{R}^p$ with

$$Y_i = B f_i + u_i, \quad B \in \mathbb{R}^{p \times k}, \quad f_i \in \mathbb{R}^k, \quad u_i \in \mathbb{R}^p$$

$$\mathbb{E}[f_i] = 0, \quad \text{Cov}(f_i) = I_k, \quad \mathbb{E}[u_i] = 0, \quad \text{Cov}(u_i) = \Sigma$$

$$\text{Cov}(Y_i) = BB^T + \Sigma, \quad \text{estimate } L = BB^T$$

1. Using Covariance matrix (Second moments):

$$\hat{L} = \frac{1}{n} \sum_{i=1}^n Y_i Y_i^T - \Sigma$$

Rank: Condition $\text{rank}(L) \leq k$ is not used

2. Using Spectral decomposition: Let $\hat{L} = U \Lambda U^T$

$\Lambda = \text{diag}\{\lambda_1, \dots, \lambda_p\}$ (descending eigenvalues)

$U = [u_1, \dots, u_p]$ - (corresponding eigenvectors)

put $\tilde{L} = \sum_{i=1}^k \lambda_i u_i u_i^T$: new estimator satisfies $\text{rank}(L) \leq k$

Latent Variable Model

Topic Model: k topics in a corpus. h : Latent V.V.

w_{pj} is drawn from $P(h=j) = w_j$, $j=1, \dots, k$

Word is drawn from a dist over vocabulary of size d : $\mu_j \in \mathbb{R}^d$

$X_t = e_i$ iff the t -th word in the document is i .

Let x_1, x_2, x_3 be word vectors in the same document.

We obtain the population moments:

$$M_2 = \bar{E}[x_1 \otimes x_2] = \sum_{k=1}^K w_k \mu_k \otimes \mu_k \in \mathbb{R}^{d \times d}$$

$$M_3 = \bar{E}[x_1 \otimes x_2 \otimes x_3] = \sum_{k=1}^K w_k \mu_k \otimes \mu_k \otimes \mu_k \in \mathbb{R}^{d \times d \times d}$$

Grob: estimate $\{w_k\}_{k=1}^K$ and $\{\mu_k\}_{k=1}^K$ from M_1 and M_2

Orthogonal Tensor Decomposition: if $k < d$

1. Whitening: Use M_2 to determine a linear transformation

$\tilde{\mu}_k = B \mu_k \in \mathbb{R}^K$ such that $\tilde{\mu}_1, \dots, \tilde{\mu}_k$ are orthogonal

Exercise: $M_2 = \sum_{k=1}^K w_k \mu_k \otimes \mu_k = \sum_{k=1}^K w_k \mu_k \mu_k^T$

By spectral theorem. M_2 has eigenvalues $\lambda_1, \dots, \lambda_k, \dots, \lambda_d$

constructing an orthonormal basis of \mathbb{R}^d with descending eigenvalues $\lambda_1, \dots, \lambda_k, 0, \dots, 0$. (Suppose μ_1, \dots, μ_k are l.i.)

Then define $\tilde{\mu}_k = (\nu_1, \dots, \nu_d)^T \mu_k \quad UU^T = I$

$$\tilde{\mu}_i^T \tilde{\mu}_j = \mu_i^T (\nu_1, \dots, \nu_d) (\nu_1, \dots, \nu_d)^T \mu_k = \mu_i^T \mu_k$$

2. Use \mathcal{B} to transform $\bar{M}_3 = \sum_{i=1}^k \lambda_i \bar{\mu}_i \otimes \bar{\mu}_i \otimes \bar{\mu}_i$

for certain $\lambda_1, \dots, \lambda_k$

3. Apply Power Method to derive $\lambda_k, \bar{\mu}_k$ from M_3 numerically

3. Maximum Likelihood Estimation

$Y \sim P_\theta$, $\theta \in \Theta$, pdf: f_θ , $L(\theta, y) = f_\theta(y)$

$$\hat{\theta}_{MLE} := \underset{\theta \in \Theta}{\operatorname{argmax}} \cdot L(\theta, y) = \underset{\theta \in \Theta}{\operatorname{argmax}} \ell(\theta, y)$$

Properties: Asym Consistency, Efficiency

Computational Complexity: Community Detection

$Y \in \mathbb{R}^{n \times n}$: Symmetric adjacency matrix for n nodes, $Y_{ii} = 0$

$Y_{ij} \sim \text{Bernoulli}(p^*)$ if i and j are in the same community

$Y_{ij} \sim \text{Bernoulli}(q^*)$ if i and j are in different communities.

$Z_i^* = 1$: i th node in Comm 1
 $Z_i^* = -1$: i th node in Comm 2

} equal-sized

$\theta = (P^*, q^*, \{Z_i\})$ Consider $\{Z_i\}$ term:

$$\ell(\theta, Y) \propto \sum_{i < j} I\{Z_i Z_j = 1\} Y_{ij}$$

Objective: $\max_{Z \in \{-1, 1\}^n: Z^T Z = 0} \sum_{i < j} I\{Z_i Z_j = 1\} Y_{ij}$ (NP Hard)

Note that $\sum \{z_i z_j = 1\} = \frac{1 + \sum z_i z_j}{2}$

$$\max_{\mathbf{z} \in \{-1, 1\}^n : \mathbf{z}^\top \mathbf{1}_n = 0} \sum_{i,j} \sum \{z_i z_j = 1\} Y_{ij}$$

$$\Leftrightarrow \max_{\mathbf{z} \in \{-1, 1\}^n : \mathbf{z}^\top \mathbf{1}_n = 0} \sum_{i,j} z_i Y_{ij} z_j$$

$$\Leftrightarrow \max_{\mathbf{z} \in \{-1, 1\}^n : \mathbf{z}^\top \mathbf{1}_n = 0} \mathbf{z}^\top \mathbf{Y} \mathbf{z}$$

$$\xrightarrow{\text{relax}} \max_{\|\mathbf{z}\| = \sqrt{n}} \mathbf{z}^\top \mathbf{Y} \mathbf{z}$$

$\tilde{\mathbf{z}}$: longest eigenvector of \mathbf{Y} with $\|\tilde{\mathbf{z}}\| = \sqrt{n}$.

Discretization: $\hat{z}_i = 1$ if $\tilde{z}_i > 0$, otherwise, $\hat{z}_i = -1$

Lecture 9. Law of Large Number and Estimation Consistency

$\text{Ex: } X_n \xrightarrow{\text{a.s.}} X \Rightarrow X_n \xrightarrow{P} X.$

$$P, \omega: \lim_{n \rightarrow \infty} X_n(\omega) = X(\omega) = 1$$

$$P\left(\bigcup_{n=1}^{\infty} \bigcap_{n \geq N} \{||X_n(\omega) - X(\omega)|| < \varepsilon\}\right) = 1 \quad \text{for any } \varepsilon > 0.$$

$$\text{put } A_{n,\varepsilon} = \{\omega: ||X_n(\omega) - X(\omega)|| < \varepsilon\}$$

$$B_{N,\varepsilon} = \bigcap_{n \geq N} A_{n,\varepsilon} \quad B_{1,\varepsilon} \subseteq B_{2,\varepsilon} \subseteq \dots$$

$$\lim_{n \rightarrow \infty} P(B_{n,\varepsilon}) = P\left(\bigcup_{n=1}^{\infty} B_{n,\varepsilon}\right) = 1$$

$$B_{N,\varepsilon} \subseteq A_{n,\varepsilon} \Rightarrow \lim_{n \rightarrow \infty} P(A_{n,\varepsilon}) = 1$$

$$\Leftrightarrow \lim_{n \rightarrow \infty} P(||X_n(\omega) - X(\omega)||_2 > \varepsilon) = 0$$

Borel-Cantelli Lemma. Let (Ω, \mathcal{F}, P) be a probability space and let $\{A_n\}_{n \geq 1}$ be a sequence of events $A_n \in \mathcal{F}$

$$1. \sum_{n=1}^{\infty} P(A_n) < \infty \Rightarrow P(\limsup A_n) = 0$$

2. If $\{A_n\}_n$ are independent, then

$$\sum_{n=1}^{\infty} P(A_n) = \infty \Rightarrow P(\limsup A_n) = 1$$

$$\begin{aligned} \limsup A_n \\ := \bigcap_{m=1}^{\infty} \bigcup_{n \geq m} A_n \end{aligned}$$

$$\text{Pf. 1. } \text{Let } N(\omega) = \sum_{n=1}^{\infty} I_{A_n}(\omega)$$

$$\mathbb{E}[N(\omega)] \stackrel{\text{def}}{=} \sum_{n=1}^{\infty} \mathbb{E}[I_{A_n}] = \sum_{n=1}^{\infty} P(A_n) < \infty$$

$$\Rightarrow P(N = \infty) = 0$$

Alternatively:

$$P(\bigcap_{m \geq 1} \bigcup_{n \geq m} A_n) = \lim_{m \rightarrow \infty} P(\bigcup_{n \geq m} A_n) \leq \lim_{m \rightarrow \infty} \sum_{n=m}^{\infty} P(A_n) = 0$$

2. For $M < N$, $P(\bigcap_{n=M}^N A_n^c) = \prod_{n=M}^N (1 - P(A_n))$

$$\leq \prod_{n=M}^{\infty} e^{-P(A_n)}$$

$$= \exp\left(-\sum_{n=M}^{\infty} P(A_n)\right)$$

$\rightarrow 0$ as $N \rightarrow \infty$

$$\begin{aligned} P(\limsup_n A_n)^c &= P\left(1 - \bigcup_{m=1}^{\infty} \bigcap_{n=m}^{\infty} A_n^c\right) \\ &= \lim_{m \rightarrow \infty} P\left(\bigcap_{n=m}^{\infty} A_n^c\right) = 0 \end{aligned}$$

$$P(\bigcap_{n=M}^{\infty} A_n^c) = 0 \Rightarrow P(\bigcup_{n=M}^{\infty} A_n) = 1 \text{ for any } M$$

$$\Rightarrow P\left(\bigcap_{M=1}^{\infty} \bigcup_{n=M}^{\infty} A_n\right) = 1$$

Exercise: $\lim_{N \rightarrow \infty} \sum_{n=N}^{\infty} P(\|X_n - X\|_2 > \varepsilon) = 0, \quad \forall \varepsilon > 0$

$$\Rightarrow \lim_{N \rightarrow \infty} P\left(\bigcup_{n=N}^{\infty} \{ \|X_n - X\|_2 > \varepsilon \}\right) = 0$$

$$\Leftrightarrow P\left(\bigcap_{N=1}^{\infty} \bigcup_{n=N}^{\infty} \{ \|X_n - X\|_2 < \varepsilon \}\right) = 1 \quad \text{i.e. } X_n \xrightarrow{a.s.} X$$

Pf. $\lim_{N \rightarrow \infty} \sum_{n=N}^{\infty} P(\|X_n - X\|_2 > \varepsilon) = 0$

$$\Rightarrow \sum_{n=1}^{\infty} P(\|X_n - X\|_2 > \varepsilon) < \infty \quad \text{i.e. } X_n \xrightarrow{\text{if}} X$$

$$\Rightarrow P\left(\bigcap_{N=1}^{\infty} \bigcup_{n=N}^{\infty} \{ \|X_n - X\|_2 > \varepsilon \}\right) = 0 \Rightarrow P\left(\bigcup_{N=1}^{\infty} \bigcap_{n=N}^{\infty} \{ \|X_n - X\|_2 \leq \varepsilon \}\right) = 1$$

Op and O_p Notations.

Def. Let $\{a_n\}$ be a sequence of real numbers, and $\{b_n\}$ be a sequence of positive numbers. Denote

$$\begin{cases} a_n = O(b_n) & \text{if there exists } C > 0 \text{ s.t. } |a_n| \leq C b_n, \forall n \\ a_n = o(b_n) & \text{if } \frac{a_n}{b_n} \rightarrow 0 \text{ as } n \rightarrow \infty \end{cases}$$

Def. Let X_1, X_2, \dots be random vectors and Y_1, Y_2, \dots be random variables defined on a common prob space

(i) $X_n = O_p(Y_n)$ iff. $\forall \varepsilon > 0, \exists C > 0$ s.t.

$$\limsup_n P(|X_n| > C|Y_n|) < \varepsilon$$

(ii) $X_n = o_p(Y_n)$ iff. $|X_n|/|Y_n| \xrightarrow{P} 0$ as $n \rightarrow \infty$

Rank. $\{X_n\}$ is said to be bounded in probability if $X_n = O_p(1)$

Exercise. $X_n \sim N(0, \sigma_n^2), n=1, 2, \dots$, then

$$\lim_{n \rightarrow \infty} \sigma_n^2 = \sigma^2 < \infty \Leftrightarrow X_n = O_p(1)$$

Pf. $P(|X_n| \geq C) \leq \frac{\text{Var}(X_n)}{C^2}$ by Markov's inequality

$$\limsup_n P(|X_n| \geq C) \leq \frac{\sigma^2}{C^2}$$

We could choose C^2 large enough s.t. $\frac{\sigma^2}{C^2} < \varepsilon$

" \Leftarrow ". Suppose $\sigma_n^2 \rightarrow \infty$ as $n \rightarrow \infty$

Note that $\frac{X_n}{\sigma_n} \sim N(0, 1)$

$$P(|X_n| \geq c) \geq P(X_n \geq c) = (1 - \Phi(\frac{c}{\sigma_n}))$$

For any $c > 0$ we find σ_n large enough such that $\Phi(\frac{c}{\sigma_n}) \leq 0.8$.

$P(|X_n| \geq c) \geq 0.2$ for large enough n .

$\{|X_n|\}$ is not bounded

Proposition X_1, \dots, X_n be random vectors. Y_1, \dots, Y_n be random variables defined on a common probability space

(1) If $X_n = O_p(Y_n)$, then $X_n = O_p(Y_n)$

as if $X_n = O_p(1)$ and $Y_n = O_p(1)$, then $X_n Y_n = O_p(1)$

3. Estimation Consistency.

Suppose we have data $X^{(n)}$ for each n drawn from a distribution $P \in \mathcal{P}$, and an estimator $\hat{T}_n(x) \in \mathbb{R}^p$

Def. Consistency Suppose $\hat{T}_n(x)$ is an estimator of unknown parameter $v \in \mathbb{R}^p$. Let $(a_n)_{n \geq 1}$ be a sequence of positive numbers with $a_n \rightarrow \infty$

(1) $\hat{T}_n(x)$ is consistent for v iff $\hat{T}_n(x) \xrightarrow{P} v$, $\forall P \in \mathcal{P}$

(2) $\hat{T}_n(x)$ is a_n -consistent for v iff $a_n [\hat{T}_n(x) - v] = O_p(1)$ holds for any $P \in \mathcal{P}$

(3) $\hat{T}_n(x)$ is strongly consistent iff $\hat{T}_n(x) \xrightarrow{a.s.} v$, $\forall P \in \mathcal{P}$

4. Law of Large Numbers

Thm. (LLN. Simple version) Suppose X_1, X_2, \dots have finite second moment with the same mean $\mu = \mathbb{E}X_i$ and bounded variance $\sigma_i^2 = \text{Var}(X_i) \leq C$ for all i . Further assume that any pair of two random variables are uncorrelated.

Then:

$$\lim_{n \rightarrow \infty} \mathbb{E}[(\frac{1}{n} \sum_{i=1}^n X_i - \mu)^2] = 0 \text{ and } \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{P} \mu.$$

Chebychev's inequality - If Y is a random variable satisfying $\mathbb{E}[Y^2] < \infty$ and $a > 0$, then

$$P(|Y - \mathbb{E}Y| \geq a) \leq \frac{\text{Var}(Y)}{a^2}$$

Rem. Finite higher moment implies finite lower moment.

$$(\mathbb{E}\|X\|_p^p)^{1/p} \leq (\mathbb{E}\|X\|_q^q)^{1/q}$$

Exercise: Converse is false.

See X has density function $f(x) = \frac{2C^2}{x^{2+1}}$ for $x \geq C$, $1 < 2 < 2$

$$\int f(x) dx = 2C^2 \int_C^\infty x^{-(2+1)} dx = 2C^2 \left[-\frac{1}{2} x^{-2} \right]_C^\infty = 1$$

$$\mathbb{E}[X] = 2C^2 \int_C^\infty x^{-2} dx = \frac{2C}{2-1}$$

$$\mathbb{E}[X^2] = 2C^2 \int_C^\infty x^{-2+1} dx \text{ diverges}$$

\checkmark Then X_1, X_2, \dots be i.i.d. random variables, and $a_n = \bar{E}[X_1 \mathbb{1}_{\{|X_1| \leq n\}}]$

Suppose $nP(|X_n| > n) \rightarrow 0$ as $n \rightarrow \infty$. Then:

$$\frac{1}{n} \sum_{i=1}^n X_i - a_n \xrightarrow{P} 0$$

Consequence: weak LLN. If $\bar{E}|X_1| < \infty$, $\frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{P} \bar{E}X_1$.

\checkmark Pf (Exercise).

First we prove the consequence $a_n \rightarrow \bar{E}|X_1|$ by DCT.

$$nP(|X_n| > n) \leq n \frac{\bar{E}|X_n|}{n} = \bar{E}|X_1|$$

Suppose for contradiction that $nP(|X_n| > n)$ does not converge to 0. Then \exists subsequence $\{n_k\}$ s.t.

$n_k P(|X_{n_k}| > n_k) \geq c$ for all k and some $c > 0$

$$P(|X_{n_k}| > n_k) \geq \frac{c}{n_k}$$

$$\text{Then } \bar{E}[|X_1|] = \int_0^\infty P(|X_{n_k}| > t) dt$$

$$\geq \int_{n_k}^\infty P(|X_{n_k}| > t) dt$$

$$\geq \int_{n_k}^\infty \frac{c}{t} dt = c [\ln t]_{n_k}^\infty = \infty$$

Contradicts to $\bar{E}|X_1| < \infty$

Therefore, $nP(|X_n| > n) \rightarrow 0$

For the major part. Let $Y_{n_j} = X_j \mathbb{1}_{\{|X_j| \leq n\}}$. Then

$$P\left(\frac{1}{n} \sum_{j=1}^n X_j \neq \frac{1}{n} \sum_{j=1}^n Y_{n_j}\right) \leq \sum_{j=1}^n P(X_j \neq Y_{n_j}) = n P(|X_j| > n) \rightarrow 0$$

Thus, we only need to work on $\{Y_n\}$. Let $T_n = \frac{1}{n} \sum_{j \in n} Y_{nj}$

$$P(|T_n - \bar{E} T_n| \geq \varepsilon) \leq \frac{\text{Var}(Y_{nj})}{n \varepsilon^2} \leq \frac{\bar{E}[Y_{nj}^2]}{n \varepsilon^2}$$

$$\leq \frac{1}{n \varepsilon^2} \bar{E}[\min\{|X_j|, n\}^2]$$

$$= \frac{1}{n \varepsilon^2} \bar{E}\left[\int_0^n 2t \mathbb{1}_{(|X_j| \geq t)} dt\right]$$

$$(T_{\text{ubini}}) = \frac{1}{n \varepsilon^2} \int_0^n t \underbrace{P(|X_j| \geq t)}_{\text{if } |X_j| \geq t} dt \rightarrow 0$$

$\int_0^n 2t dt = |X_j|^2$

$$\int_0^n 2t dt = |X_j|^2$$

$\int_0^n 2t dt = |X_j|^2$

$$\int_0^n 2t dt = |X_j|^2$$

$\frac{1}{n} \int_0^n t P(|X_j| \geq t) dt$: average of $t P(|X_j| \geq t)$

5. Strong LLN: $\text{and } \tau P(|X_j| \geq t) \rightarrow 0$

Tail σ -algebra Let X_1, X_2, \dots be random variables.

$F'_n = \sigma(X_n, X_{n+1}, \dots)$ $T = \bigcap_{n=1}^{\infty} F'_n$: tail σ -algebra

Remk. An event B is in tail σ -algebra if it does not depend on any finite number of events

Kolmogorov's 0-1 law $A \in T$. $P(A) \in \{0, 1\}$

$$\text{? Exercise } A_1 = \left\{ \limsup_n \frac{X_1 + \dots + X_n}{n} \geq \bar{E} X_1 + \varepsilon \right\}$$

$$A_2 = \left\{ \limsup_n \frac{X_1 + \dots + X_n}{n} \leq \bar{E} X_1 - \varepsilon \right\}$$

$$P(A_1) = 1$$

Sum. $\{X_n\}$ iid, $\mathbb{E}|X_1| < \infty$ and $C_i \geq 1$ bounded. Then

$$\frac{1}{n} \sum_{i=1}^n C_i(X_i - \bar{E}[X_i]) \xrightarrow{a.s.} 0$$

7. Consistency of Mom and Var.

Ex: X_1, \dots, X_n iid $P \in \mathcal{P}$. Assume $\mathbb{E}|X_1| = \mu < \infty$

$$\text{Var}(X) = \sigma^2 < \infty, \text{ Then } S_n = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

and $\frac{n-1}{n} S_n$ are strongly consistent for σ^2

$$\begin{aligned} P \quad S_n &= \frac{1}{n-1} \left(\sum_{i=1}^n X_i^2 - 2 \sum_{i=1}^n X_i \bar{X} + \sum_{i=1}^n \bar{X}^2 \right) \\ &= \frac{1}{n-1} \left(\sum_{i=1}^n X_i^2 - 2n \bar{X}^2 + n \bar{X}^2 \right) \\ &= \frac{n}{n-1} \left(\frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}^2 \right) \xrightarrow{\text{a.s.}} \bar{E}[X^2] - \bar{E}[X]^2 = \text{Var}(X) \end{aligned}$$

$$\frac{n-1}{n} S_n \xrightarrow{\text{a.s.}} \text{Var}(X)$$

Consistency of Mom.

Let $\mu \in \mathbb{R}^k$ be the first k -th moments of P_θ and

$\mu = h(\theta)$ for some $h: \mathbb{R}^k \rightarrow \mathbb{R}^k$ continuous bijection function

Let $\hat{\mu}_n \in \mathbb{R}^k$ be empirical moments based on n iid

random variables from P_θ . Then under assumptions on

finite moments: $\hat{\mu}_n \xrightarrow{\text{a.s.}} \mu \Rightarrow h^{-1}(\hat{\mu}_n) \xrightarrow{\text{a.s.}} \theta$ as $n \rightarrow \infty$

Consistency of $\hat{\theta}$.

For iid data, the negative log-likelihood function is the sum of n independent terms for every $\theta \in \Theta$.

$$-\ell(\theta; y) = -\frac{1}{n} \sum_{i=1}^n g(x_i, \theta)$$

Under some regularity conditions, $\hat{\theta}$ is \sqrt{n} -consistent.

Lecture 10. Weak Convergence and CLT

1. Weak Convergence.

Equivalent definitions for $X_n \xrightarrow{d} X$:

- cdf $F_n(x) \rightarrow F(x)$ for each continuity point x of F
- $\lim_{n \rightarrow \infty} E[h(X_n)] = E[h(X)]$ for all bounded continuous functions h on \mathbb{R}^d
- Characteristic function $\lim_{n \rightarrow \infty} \phi_{X_n}(t) = \phi_X(t)$ for all $t \in \mathbb{R}^d$
- $X \xrightarrow{d} X$ iff $C^T X_n \xrightarrow{d} C^T X$, $\forall C \in \mathbb{R}^{d \times d}$

Thm. Let $\{X_i\}_{i \in \mathbb{Z}^+}$ be random vectors in \mathbb{R}^d

- (i) $X_n \xrightarrow{P} X$ then $X_n \xrightarrow{d} X$
- (ii) $X_n \xrightarrow{d} C$, $C \in \mathbb{R}^k$ then $X_n \xrightarrow{P} C$
- (iii). $X_n \xrightarrow{d} X$, then $X_n = O_p(1)$

2. Central Limit Theorem

Thm. Let X_1, X_2, \dots be iid random vectors in \mathbb{R}^d with finite second moment. Let $\Sigma = \text{Cov}(X_1)$, then

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n (X_i - E[X_i]) \xrightarrow{d} N(0, \Sigma)$$

Lecture 11. Asymptotic normality and Delta Method

1. Convergence of Transformations

Continuous Mapping Theorem Let X_1, X_2, \dots be random vectors

in \mathbb{R}^d defined on a common probability space and g be a measurable function from $(\mathbb{R}^d, \mathcal{B}^d)$ to $(\mathbb{R}^k, \mathcal{B}^k)$. Suppose

g is cts a.s. w.r.t. μ_X , then

$$(i) X_n \xrightarrow{\text{a.s.}} X \Rightarrow g(X_n) \xrightarrow{\text{a.s.}} g(X)$$

$$(ii) X_n \xrightarrow{P} X \Rightarrow g(X_n) \xrightarrow{P} g(X)$$

$$(iii) X_n \xrightarrow{d} X \Rightarrow g(X_n) \xrightarrow{d} g(X)$$

PF. (i). $A_0 = \{w: \lim_{n \rightarrow \infty} X_n = X\}$ $D = \{x \in \mathbb{R}^k, g \text{ cts at } x\}$.

$$\forall w \in A_0 \cap X^{-1}(D) \quad \lim g(X_n) = g(X)$$

$$A^c = A_0^c \cup X^{-1}(D^c) \quad P(A^c) \leq P(A_0^c) + P(X^{-1}(D^c)) = 0 \Rightarrow P(A) = 1$$

$$(ii). \forall \varepsilon > 0. \exists \delta > 0 \quad \|g(x) - g(y)\| \leq \varepsilon \quad \text{if } \|x - y\| < \delta$$

$$\text{Since } \lim P(|X_n - X| > \delta) = 0$$

$$P(\|g(X_n) - g(X)\| > \varepsilon) \leq P(|X_n - X| > \delta) \rightarrow 0 \quad \text{as } n \rightarrow \infty$$

$$(iii) X_n \xrightarrow{d} X \Rightarrow \exists Y_1, Y_2, \dots, Y_i \stackrel{d}{=} X_i \text{ and } Y_n \xrightarrow{\text{a.s.}} Y$$

Svesky's Theorem Let $X_1, X_2, \dots, Y_1, Y_2, \dots$ be random variables

on a probability space. Suppose $X_n \xrightarrow{d} X$, $Y_n \xrightarrow{d} c$. Then

(i) $X_n + Y_n \xrightarrow{d} X + c$

(ii) $X_n Y_n \xrightarrow{d} cX$

(iii) $X_n / Y_n \xrightarrow{d} X/c$ if $c \neq 0$

2. Delta Method

The Delta method Let X_1, X_2, \dots and Y be random vectors in \mathbb{R}^k satisfying $a_n(X_n - c) \xrightarrow{d} Y$. For $c \in \mathbb{R}^k$ and $\{a_n\}$ being a sequence of positive numbers with $\lim_{n \rightarrow \infty} a_n = \infty$. Let $g: \mathbb{R}^k \rightarrow \mathbb{R}$. If g is differentiable at c :

$$a_n [g(X_n) - g(c)] \xrightarrow{d} \nabla g(c)^T Y$$

3. Prevalence of asymptotic normality

Asymptotic Variance $\{a_n\}$. Positive sequence and either $a_n \rightarrow \infty$ or $a_n \rightarrow a > 0$. Assume

$$a_n [T_n(\bar{X}) - \theta] \xrightarrow{d} Y \text{ with } 0 < \mathbb{E} Y^2 < \infty$$

(i) The asymptotic variance of T_n is defined as $\frac{\text{Var}(Y)}{a_n^2}$

(ii) Let $T'_n(\bar{X})$ be another estimator. Asymptotic relative efficiency is defined to be the ratio between the

The Consistent Variable

$$\text{as } \mathbb{E}[Y^2] \in \liminf_{n \rightarrow \infty} \mathbb{E}[a_n^2(\bar{T}_n - \theta)^2]$$

Lecture 12: Unbiased Estimation and UMVUE.

1. Complete Statistic

Def. (Complete Statistics) $T(x)$ is complete for $P \in \mathcal{P}$

iff \forall Borel f , $E[f(T)] = 0 \Rightarrow f(T) = 0$ a.s. P

Rmk: Requires the statistic has no redundant information.

Proposition. Let $\mathcal{P} = \{P_\eta, \eta \in \Xi\}$ be an exponential family of full rank (contains an open set) with pdf

$$f_\eta(x) = \exp\{\eta^\top T(x) - \zeta(\eta)\} h(x)$$

Then $T(x)$ is sufficient and complete for $\eta \in \Xi$

2. UMVUE

Def. An unbiased estimator $T(x)$ of ν is called

the UMVUE iff $\text{Var}(T(x)) \leq \text{Var}(U(x))$

for any $P \in \mathcal{P}$ and any other unbiased estimator $U(x)$ of ν

Rmk UMVUE does not always exist

Thm: Suppose there exists a Sufficient and Complete statistic $T(x)$ for P . If v is estimable then there exists unique UMVUE, which is of form $h(T)$ where h is a Borel function.

Exists unbiased estimator

3. Construct UMVUE

Method 1: Find Sufficient and Complete statistic $T(x)$, then find $h(T)$ such that $\bar{E}[h(T)] = v$, $\forall P \in \mathcal{P}$

Method 2: Find Sufficient, complete $T(x)$ and unbiased $U(x)$, then $\bar{E}[U(T)]$ is an UMVUE.

Method 3: Find UMVUE without knowing complete statistics

Thm 2: Let $\mathcal{U} = \{U : \bar{E}[U(x)] = 0, \text{Var}(U(x)) < \infty, \forall P \in \mathcal{P}\}$

T is unbiased for v with $\bar{E}[T(x)] < \infty$

(1). $T(x)$ is UMVUE iff

$$\bar{E}[T(x)U(x)] = 0, \forall U \in \mathcal{U}, \forall P \in \mathcal{P}$$

(2) If T is sufficient for P , let

$$\tilde{\mathcal{U}} = \mathcal{U} \cap \{g(T) : g \text{ Borel}\}$$

Then $T = h(\tilde{T})$ is UMVUE iff

$$\bar{E}[T(x)U(x)] = 0, \forall U \in \tilde{\mathcal{U}}, \forall P \in \mathcal{P}$$

Method 4*: Variational Calculus

Lecture 13: Fisher information and C-R Lower Bound.

1. Fisher Information

Def. For $X \sim \{P_\theta : \theta \in \Theta\}$:

$$I(\theta) = E\left[\frac{\partial}{\partial \theta} \log f_\theta(x) \left(\frac{\partial}{\partial \theta} \log f_\theta(x)\right)^T\right]$$

Rank. $I(\theta) \geq 0$. $\frac{\partial}{\partial \theta} \log f_\theta(x)$ is called the score function

Proposition

1. If $X \perp\!\!\!\perp Y$, then $I_{(X,Y)}(\theta) = I_X(\theta) + I_Y(\theta)$

2. Suppose f_θ is twice differentiable at θ , under some "regularity condition", we have

$$I(\theta) = -E\left[\frac{\partial^2}{\partial \theta^2} \log f_\theta(x)\right].$$

2. Cramér-Rao Lower Bound

Let $v = g(\theta)$, $g: \mathbb{R}^d \rightarrow \mathbb{R}^k$ differentiable. Suppose X is drawn from $\{P_\theta : \theta \in \Theta\}$. And $\hat{T}(x)$ is an unbiased estimator of v .

Thm. Suppose $I(\theta)$ is positive definite, and for any $\theta \in \Theta$,

$$\frac{\partial}{\partial \theta} \int h(x) f_\theta(x) dv = \int h(x) \frac{\partial}{\partial \theta} f_\theta(x) dv$$

holds for $h \in \mathbb{I}$ or $h(x) = \hat{T}(x)$. Then,

$$\text{Var}(\hat{T}(x)) \geq \left[\frac{\partial}{\partial \theta} g(\theta) \right]^T I(\theta)^{-1} \left[\frac{\partial}{\partial \theta} g(\theta) \right]$$

Rmk (Reparameterization). If $\theta = \phi(\eta)$, $\psi \in C^1$ is bijective,

$$\text{then } \tilde{I}_x(\eta) = \frac{\partial}{\partial \eta} \psi(\eta) \tilde{I}_x(\psi(\eta)) \frac{\partial}{\partial \eta} \psi(\eta)^T$$

is different but

$$\frac{\partial}{\partial \eta} g(\theta) = \frac{\partial}{\partial \theta} g(\theta) \frac{\partial}{\partial \eta} \psi(\eta).$$

reparameterization

Therefore, the C-R lower bound is invariant under

Rmk. Asymptotic optimality of $\hat{\mu}_{\bar{E}}$.

$\text{Var}(\hat{\mu}_{\bar{E}}) \rightarrow (n \tilde{I}_x(\theta))^{-1}$ matches the C-R lower bound

Rmk. C-R lower bound may not be attained

3. Interpretations of Fisher information.

3.1. Geometric view.: Larger $\tilde{I}(\theta) \Rightarrow$ Larger local convexity

3.2. Information theory: under certain "regularity conditions"

$$D(P_\theta || P_{\theta+\xi}) = \frac{1}{2} \xi^T \tilde{I}(\theta) \xi + O(\|\xi\|^2)$$

"How hard to distinguish two dist in a parametric family under the K2 divergence"

5. Examples

5.1. Exponential Family.

$\mathcal{X} \subset \mathbb{R}^k$ be an open set. Let $\{\tilde{f}_\theta : \theta \in \Theta\}$ be

$$\tilde{f}_\theta(x) = \exp\{\eta(\theta)^T T(x) - \tilde{g}(\theta)\} c(x)$$

Proposition

1. The regularity conditions for C-R to be satisfied
for any msb h with $\mathbb{E}[h(x)] < \infty$
2. Consider natural parameter η , $\text{Var}[\tilde{T}(x)] = \tilde{I}(\eta)$
3. $V = \mathbb{E}[\tilde{T}(x)]$, then $\text{Var}[\tilde{T}] = (\tilde{I}(V))^{-1}$

5.2 Linear Models. $X_i = \beta^T z_i + \epsilon_i \quad z \in \mathbb{R}^{np} \quad \epsilon \sim N(0, \sigma^2 I_n)$

Thm 2 Suppose Z is full column rank

(i) $\hat{\beta}^T$ is UMVUE of β^T , $\forall \beta \in \mathbb{R}^p$

(ii) $\hat{\sigma}^2 = \frac{1}{n-p} \|X - Z\hat{\beta}\|^2$ is the UMVUE of σ^2

Pf (Sketch) $(Z\hat{\beta}, \|X - Z\hat{\beta}\|^2)$ is complete and sufficient
for $\theta = (\beta, \sigma^2)$. Further verify unbiasedness

Thm 3 Under assumptions in thm 2, $\hat{\beta}^T$ is independent
of $\hat{\sigma}^2$, moreover,

$$\hat{\beta}^T \sim N(\beta^T, \sigma^2 \hat{V}^T (Z^T Z)^{-1} \hat{V}), \quad \hat{\sigma}^2 \frac{n-p}{\sigma^2} \sim \chi_{n-p}^2$$

Fisher information $\{\text{iid } N(0, \sigma^2 I_n)\}$ and the unknown parameters are $\theta = (\beta, \sigma^2)$

$$I(\theta) = \frac{1}{\sigma^2} \begin{pmatrix} z^T z & 0 \\ 0 & \frac{n}{z^T z} \end{pmatrix}$$

If z is not full rank, $I(\theta)$ may not be invertible.

Two distributions P_θ and $P_{\theta'}$ are not distinguishable

Prop. LS estimate $\hat{\beta}$ attains C-R lower bound.

Lecture 14. Concentration Inequality

1. Motivation Non-asymptotic bounds (finite sample)

Gaussian Tail Inequality $\text{Let } G \sim N(0,1). \text{ Then}$

$$\left(\frac{1}{t} - \frac{1}{t^2}\right) \frac{1}{\sqrt{2\pi}} e^{-t^2/2} \leq P(G \geq t) \leq \frac{1}{t} \frac{1}{\sqrt{2\pi}} e^{-t^2/2} \quad \forall t > 0$$

$\text{Let } \tilde{F} \text{ and } \tilde{P} \text{ be the Gaussian CDF and PDF. Then}$

$$1 - \tilde{F}(t) = (1 + O(1)) \frac{1}{t} \tilde{P}(t)$$

Berry-Essen Inequality - Suppose X_1, \dots, X_n iid

with $\mathbb{E}[X_i] = 0$ and $\mathbb{E}[X_i^3] < \infty$, then

$$|P\left(\frac{1}{\sqrt{n}} \sum_{i=1}^n X_i \geq t\right) - P(G \geq t)| \leq \frac{C \mathbb{E}[X^3]}{\sqrt{n}}$$

Rmk: The general bound is too large.

2. Subgaussian variables and Hoeffding's Inequality

Theorem 1. Equivalent defns of Subgaussian random variables

(1) $P(|X| \geq t) \leq 2 \exp(-t^2/k_1^2)$, $\forall t > 0$ for some $K_1 > 0$

(2) $\|X\|_p \leq K_2 \sqrt{p}$ for some $K_2 > 0$ and all $p \in \mathbb{N}^*$

(3) $\mathbb{E} \exp(\lambda^2 X^2) \leq \exp(K_3^2 \lambda^2)$ for some $K_3 > 0$
and all λ w. $|M| \leq \frac{1}{|K_3|}$

(4) $\mathbb{E} \exp(X^2/k_4^2) = m_{X^2}(\frac{1}{k_4^2}) \leq 2$ for some $K_4 > 0$

(5) If $\mathbb{E}X=0$, $\mathbb{E}\exp(\lambda X) \leq \exp(k_2^2 \lambda^2)$, $\forall \lambda \in \mathbb{R}$

Remark. The smallest k_2 is called the subgaussian norm of X , denoted by $\|X\|_{\psi_2}$.

Examples. If $X \sim N(0, \sigma^2)$, then $\|X\|_{\psi_2} \leq C\sigma$

• If $|X| \leq M$ a.s. then $\|X\|_{\psi_2} \leq CM$

• $X \sim \text{Poisson}(0.5)$ $\|X\|_{\psi_2} = \sqrt{\ln 2}$

Pf. (1) \Rightarrow (2). Assume $K_1 = 1$:

$$\mathbb{E}|X|^p = \int_0^\infty P(|X|^p \geq u) du$$

$$= \int_0^\infty P(|X| \geq t) dt^p$$

$$\leq \int_0^\infty 2e^{-t^2} p t^{p-1} dt$$

$$= 2p \int_0^\infty e^{-v} v^{\frac{1}{2}(p-1)} dv^{\frac{1}{2}}$$

$$= p \int_0^\infty e^{-v} v^{\frac{1}{2}p-1} dv$$

$$= p P(\frac{1}{2}p)$$

$$\leq 3p(p/2)^{p/2} \Rightarrow \|X\|_p \leq \sqrt{\frac{3}{\pi}} p^{1/p} p^{1/2} \leq \sqrt{\frac{3}{\pi}} \sqrt{p}$$

we have (2) with $K_2 \leq 3$.

(2) \Rightarrow (1) Assume $K_3 = 1$.

$$\mathbb{E}\exp(\lambda^2 X^2) = \mathbb{E}[1 + \sum_{p=1}^{\infty} \frac{(\lambda^2 X^2)^p}{p!}] = 1 + \sum_{p=1}^{\infty} \frac{\lambda^{2p} \mathbb{E}[X^{2p}]}{p!}$$

Gamma Function:

$$P(z) = \int_0^\infty t^{z-1} e^{-t} dt$$

Stirling's Approximation

$$n! / \sqrt{2\pi n} \left(\frac{n}{e}\right)^n \rightarrow 1$$

$$n! \approx \sqrt{2\pi n} \left(\frac{n}{e}\right)^n e^{\frac{1}{12n}}$$

$$P(x) \sim x^{x-\frac{1}{2}} e^{-x} \sqrt{2\pi}$$

$$P(x) \approx 3x^x$$

Stirling's approximation

also yields $p! \geq (\pi/e)^p$

$$\leq 1 + \sum_{p=1}^{\infty} \frac{\lambda^p (2p)^p}{(p/e)^p} = \sum_{p=1}^{\infty} (2e\lambda^2)^p = \frac{1}{1-2e\lambda^2} \leq \exp(1.4e\lambda^2)$$

↓
Suppose $2e\lambda^2 < 1$

↓
 $2e\lambda^2 \in [0, \frac{1}{2}]$

For $|\lambda| \leq \frac{1}{2\sqrt{e}}$. and $K_3 = 2\sqrt{e}$

(3) \Rightarrow (4): \checkmark Trivial.

(4) \Rightarrow (1). Assume $K_4 = 1$.

$$P(X \geq t) = P(e^{X^2} \geq e^{t^2}) \leq \bar{E}[e^{X^2}] \leq 2e^{-t^2}$$

When we have $\bar{E}X = 0$:

(3) \Rightarrow (5): Assume $K_3 = 1$.

$$\begin{aligned} \bar{E}[e^{\lambda X}] &\leq \bar{E}[\lambda X + e^{\lambda X^2}] \\ &= \bar{E}[e^{\lambda^2 X^2}] \leq e^{\lambda^2} \quad \text{for } |\lambda| \leq 1. \end{aligned}$$

When $|\lambda| \geq 1$: $2\lambda X \leq \lambda^2 + X^2$

$$\bar{E}e^{\lambda X} \leq e^{\lambda^2/2} \bar{E}e^{X^2/2} \leq e^{\lambda^2/2} e^{1/2} \leq e^{\lambda^2}$$

(5) \Rightarrow (1): \checkmark for any $\lambda > 0$:

$$P(X \geq t) = P(e^{\lambda X} \geq e^{\lambda t})$$

$$\leq e^{-\lambda t} \bar{E}e^{\lambda X}$$

$$\leq e^{-\lambda t} e^{\lambda^2}$$

$$= e^{-\lambda t + \lambda^2}$$

$$\leq e^{-t^2/2}$$

Maximum of sub-gaussians. Let X_1, X_2, \dots be a sequence of sub-gaussian random variables, not necessarily indep.
 Then $\mathbb{E} \max_i \frac{\|X_i\|}{\sqrt{1+\log_i}} \leq CK$

where $K = \max_i \|X_i\|_{\psi_2}$. For any $N \in \mathbb{Z}$, we have

$$\mathbb{E} \max_{i \leq N} |X_i| \leq CK \sqrt{\log N}$$

Rank. The bound is sharp: Let X_1, X_2, \dots i.i.d. $N(0, 1)$, we have $\mathbb{E} \max_{i \leq N} X_i \geq C \sqrt{\log N}$.

Thm 2. Hoeffding's Inequality Let X_1, X_2, \dots, X_n be independent, mean zero and subgaussian r.v.s.

Let $a \in \mathbb{R}^n$, $K := \max_{1 \leq i \leq n} \|X_i\|_{\psi_2}$. Then for any $t > 0$

$$P\left(\left|\sum_{i=1}^n a_i X_i\right| > t\right) \leq 2 \exp\left(-\frac{ct^2}{K^2 \|a\|_2^2}\right)$$

holds for certain absolute constant $C > 0$

Rank. More flexible than CLT: Mean \rightarrow weighted sum

$$\begin{aligned} \text{Pf. } \mathbb{E}[\exp(\lambda \sum_{i=1}^n a_i X_i)] &= \prod_{i=1}^n \mathbb{E}[\exp(\lambda a_i X_i)] \\ &\leq \prod_{i=1}^n \exp(C \lambda^2 a_i^2 \|X_i\|_{\psi_2}^2) \\ &\leq \exp(C \lambda^2 \sum_{i=1}^n a_i^2 K^2) \\ &= \exp(C \lambda^2 \|a\|^2 K^2) \quad \text{Then use Thm 1 (1)} \end{aligned}$$

Khintchine's inequality Let X_1, X_2, \dots, X_n be independent sub-gaussian random variables with zero-means and unit variance, let $a \in \mathbb{R}^n$. Then for any $p \in [2, +\infty)$,

$$\|a\|_2 \leq \left\| \sum_{i=1}^n a_i X_i \right\|_p \leq C K \sqrt{p} \|a\|_2$$

where $K = \max_i \|X_i\|_{\psi_2}$ and C is an absolute constant

3. Subexponential Variables and Bernstein's Inequality

Thm 3. (equivalence defns of subexponential random variables)

(1) $P(|X| > t) \leq 2e^{-t/K_1}$, $\forall t \geq 0$ relaxed subgaussian.

(2) $\|X\|_p \leq K_2 p$ for all $p \geq 1$

(3) $\mathbb{E}[\exp(\lambda |X|)] \leq \exp(K_3 \lambda)$, $\forall \lambda \in [0, 1/K_3]$.

(4) $\mathbb{E}[\exp(|X|/K_4)] \leq 2$

(5) $\mathbb{E}[X] = 0$, $\mathbb{E}[\exp(\lambda X)] \leq \exp(K_5 \lambda^2)$, $\forall \lambda \in [-1/K_5, 1/K_5]$

$K_i \leq C k_j$. $\forall i \neq j$ for some C independent of the random variable X .

Remark

- X is called sub-exponential if any of the condition holds
- The smallest K_2 is called the sub-exponential norm. $\|X\|_{\psi_2}$.
- subgaussian \Rightarrow subexponential, $\|X\|_{\psi_1} \leq \|X\|_{\psi_2}$

Centering: For random variables with non-zero mean, conditions (1) and (2) in Thm 1 and Thm 3 are still equivalent. Moreover, we have bound

$$\|X - \bar{E}X\|_{\psi_2} \leq C \|X\|_{\psi_2}$$

Cor. Suppose X, Y are subgaussian. Then X, Y are subexponential with

$$\|XY\|_{\psi_2} \leq \|X\|_{\psi_2} \|Y\|_{\psi_2}$$

$$\begin{aligned} \mathbb{P}[\bar{E}((|X|^p |Y|^p)]^{1/p}] &\leq [\bar{E}|X|^p \bar{E}|Y|^p]^{1/p} \\ &\leq \|X\|_{\psi_2} \sqrt[p]{p} \|Y\|_{\psi_2} \sqrt[p]{p} \\ \frac{1}{p} [\bar{E}((|X|^p |Y|^p)]^{1/p}] &\leq \|X\|_{\psi_2} \|Y\|_{\psi_2}. \end{aligned}$$

Thm 4. Bernstein's Inequality

Let X_1, \dots, X_n be independent r.v. such that $\bar{E}X_i = 0$ and each X_i is subexponential. Let $a \in \mathbb{R}^n$,

$$K = \max_{1 \leq i \leq n} \|X_i\|_{\psi_2}.$$

subgaussian subexponential

$$\mathbb{P}\left(\left|\sum_{i=1}^n a_i X_i\right| \geq t\right) \leq 2 \exp\left(-\min\left(\frac{t^2}{K^2 \|a\|^2}, \frac{t}{K \|a\|_\infty}\right)\right)$$

for some constant $C > 0$

decays no faster than e^{-t}

Remark: Subexponential r.v. has "heavier tails" compared with subgaussian r.v.s.

Thm 5: Bernstein's Inequality For bounded random variables

Let X_1, \dots, X_n be indept r.v.s, $\mathbb{E}X_i = 0$, $\max|X_i| \leq M$ a.s.

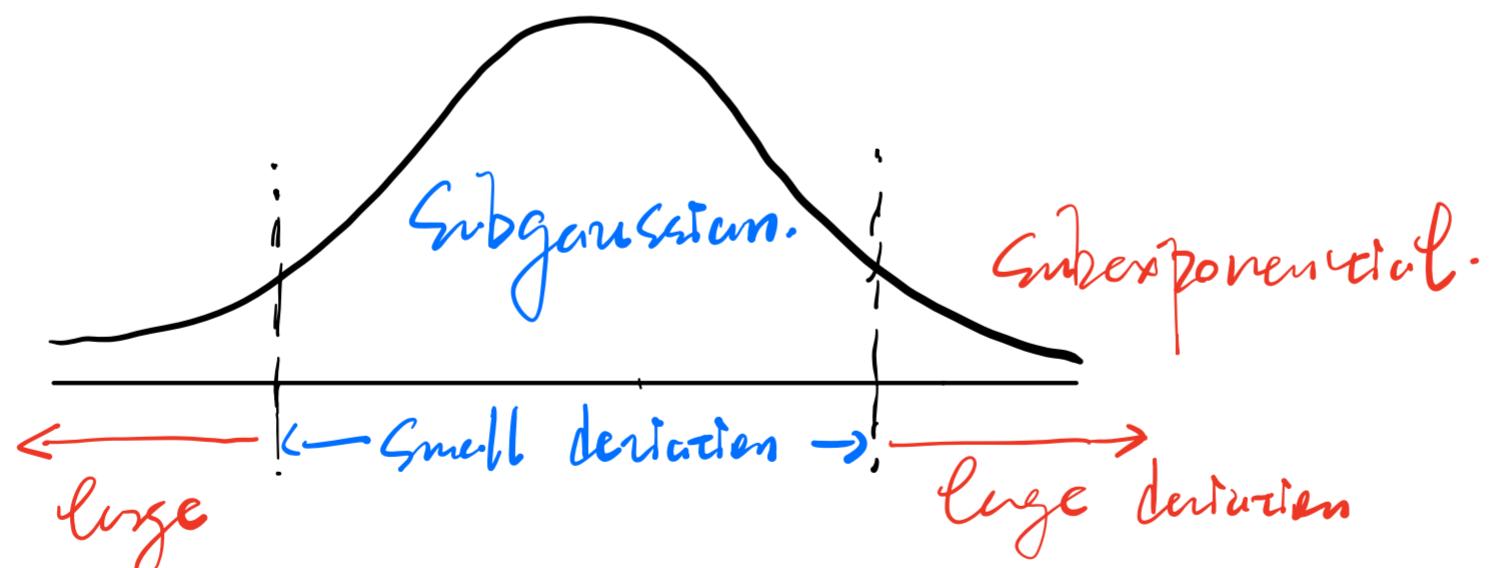
$$P(|\sum_{i=1}^n X_i| \geq t) \leq 2\exp\left(-\frac{t^2/2}{\sigma^2 + M^2/3}\right)$$

where $\sigma^2 = \sum_{i=1}^n \mathbb{E}X_i^2$ is total variance

Variants of Bernstein's Inequality:

$$P(|\sum_{i=1}^n X_i| \geq \frac{C(\sigma + M)}{\delta}) \leq 2\exp\left(-\min\left(\frac{1}{\delta}, \frac{1}{\delta^2}\right)\right), \forall \delta > 0$$

Sum of sub-exponential r.v.s



HDP 2.3. Chernoff's inequality

Thm (Chernoff's inequality) Let X_i be independent

Bernoulli random variables with parameters p_i . Consider

$$S_N = \sum_{i=1}^N X_i, \quad \mu := \mathbb{E} S_N. \quad \text{Then for any } t > \mu:$$

$$P(S_N > t) \leq e^{-\mu} \left(\frac{e^\mu}{t}\right)^t$$

For any $t < \mu$:

$$P(S_N \leq t) \leq e^{-\mu} \left(\frac{e^\mu}{t}\right)^t$$

(or. Poisson Tails: Let $X \sim \text{Poisson}(\lambda)$ For any $t > \lambda$:

$$P(X > t) \leq e^{-\lambda} \left(\frac{e^\lambda}{t}\right)^t$$

Thnk. using Stirling's Formula we have:

$$P(X=k) \approx \frac{1}{\sqrt{2\pi k}} e^{-\lambda} \left(\frac{e^\lambda}{k}\right)^k$$

So our bound on the entire tail of X has the same form as the probability of hitting one value k in the tail.

(or. (Small deviations). For some absolute constant $C > 0$

$$P(|S_N - \mu| > \delta \mu) \leq 2e^{-C \mu \delta^2}$$

(Con. (Poisson distribution near the mean) • Let $X \sim \text{Poisson}(\lambda)$.

For $c \in (0, \lambda]$, we have

$$P(|X - \lambda| \geq c) \leq 2c \exp\left(\frac{-ct^2}{\lambda}\right)$$

Rank.

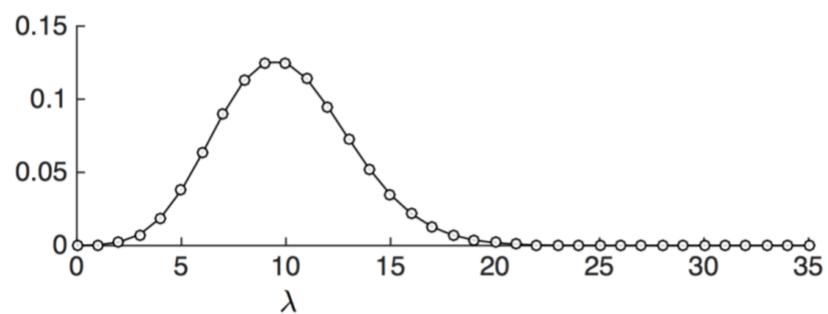


Figure 2.1 The probability mass function of the Poisson distribution $\text{Pois}(\lambda)$ with $\lambda = 10$. The distribution is approximately normal near the mean λ , but to the right from the mean the tails are heavier.

Normal Approximation to Poisson: $X \sim \text{Poisson}(\lambda)$

as $\lambda \rightarrow \infty$, we have

$$\frac{X - \lambda}{\sqrt{\lambda}} \xrightarrow{d} N(0, 1)$$

Lecture 15: Random vectors in high dims.

1. Examples of concentration inequalities

Random graphs. The G-R graph $G \sim (n, p)$:

with adjacency matrix $A \in \mathbb{R}^{n \times n}$, $A_{ij} \stackrel{iid}{\sim} \text{Bern}(p)$

• degree $d_i = \sum_{j \neq i} A_{ij}$, $\mathbb{E}[d_i] = (n-1)p$, put $d = (n-1)p$

Proposition there's an absolute constant C ,

if $d \geq C(\log n)^2$ Then $P(0.9d \leq d_i \leq 1.1d, \forall i) \geq 0.9$

Pf. Apply Bernstein's inequality. For any i , any $\delta > 0$

$$P(|d - d_i| \geq \frac{C_1(\sigma + M)}{\delta}) \leq 2 \exp(-\min(\frac{1}{\delta}, \frac{1}{\delta^2}))$$

where $\sigma = \sqrt{\text{Var}(d_i)}$, $M \leq 2$ being a simple bound.

$\text{Var}(d_i) = (n-1)p(1-p) \leq d$. Pick $\delta = (\alpha \log n)^{-1}$ for some $\alpha > 1$.

$$P(|d - d_i| \geq \frac{C_1(\sqrt{d} + 2)}{\delta}) \leq 2 \exp(-\min(\frac{1}{\delta}, \frac{1}{\delta^2})) = 2n^{-\alpha}$$

Note that $C_1(\sqrt{d} + 2)\alpha \log n \leq 0.1d$ for large C

Then, by union bound, we have

$$P(\exists i, |d_i - d| \geq 0.1d) \leq \sum_{i=1}^n P(|d_i - d| \geq \frac{C_1(\sigma + M)}{\delta}) \leq 2n^{1-\alpha}$$

Covariance matrix $\text{Cov } y_1, \dots, y_n \text{ iid } Y \in \mathbb{R}^d$. $E[YY^\top] = 0$,

with subgaussian components. $\max_{1 \leq j \leq d} \|Y_j\|_{\psi_2} \leq k$.

Then the covariance matrix of Y exists

and is denoted by $\Sigma = \text{Cov}(Y) \in \mathbb{R}^{d \times d}$

Empirical variance: $\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n y_i y_i^\top$

Denote $\|\hat{\Sigma} - \Sigma\|_{\max} := \max_{j,m} |\hat{\Sigma}_{j,m} - \Sigma_{j,m}|$.

$$= \max_{1 \leq j, m \leq d} \left| \frac{1}{n} \sum_{i=1}^n (y_{ij} y_{im} - \hat{\Sigma}_{j,m}) \right|$$

By Bernstein's inequality.

$$P(|\hat{\Sigma}_{j,m} - \Sigma_{j,m}| \geq t) \leq 2 \exp(-cn \min(t^2, t))$$

Choose $t = C \sqrt{\log d / n}$ with large C s.t.

$$cnt^2 = CC^2 \log d \geq 4 \log d$$

Assume n is large enough s.t. $t \leq 1$, then $t^2 \leq t$

$$P(|\hat{\Sigma}_{j,m} - \Sigma_{j,m}| \geq Ck \sqrt{\frac{\log d}{n}}) \leq 2 \exp(-4 \log d) = 2d^{-4}$$

Apply union bound we have

$$P(\|\hat{\Sigma} - \Sigma\|_{\max} \geq Ck \sqrt{\frac{\log d}{n}}) \leq 2d^{-2}$$

2. Two general concentration inequalities.

Thm 1. Gaussian Concentration inequality.

Let X_1, \dots, X_n iid $\mathcal{N}(0, 1)$. Let $F: \mathbb{R}^n \rightarrow \mathbb{R}$ be an L -Lipschitz function. i.e. $|F(x) - F(y)| \leq L\|x - y\|$. Then

$$\|F(X_1, \dots, X_n) - \mathbb{E}[F(X_1, \dots, X_n)]\|_{\psi_1} \leq CL$$

where C is a constant

Example 1. Suppose Y is a matrix of i.i.d standard normal entries. Let $\|\cdot\|$ be either Frobenius norm or matrix operator norm. Then

$$P(\|Y\| - \mathbb{E}\|Y\| > t) \leq e^{-ct^2}$$

Pf. Note that $|\|X\| - \|Y\|| \leq \|X - Y\|$, then apply Thm 1.

Thm 2. Talagrand's inequality. Suppose X_1, \dots, X_n are independent random variables with $\max_i |X_i| \leq K$.

Denote $X = (X_1, \dots, X_n)^T$. Let $F: \mathbb{R}^n \rightarrow \mathbb{R}$ be L -Lipschitz function. Then

$$P(|F(x) - \mathbb{E}[F(x)]| \geq ctL) \leq C \exp\left(-\frac{ct^2}{K^2}\right), \quad \forall t > 0$$

Rmk. It says bounded r.v. with indep coordinates,

$F(x) - \mathbb{E}[F(x)]$ is a subgaussian with norm bounded by $O(KL)$

3. Norms of Random Vectors.

Suppose $X = (X_1, \dots, X_n)^T$ with X_1, \dots, X_n indep and $\mathbb{E}[X_i] = 0, \text{Var}(X_i) = 1$. We already have:

1. By LN and CLT mapping thm, $n^{-\frac{1}{2}}\|X\| \xrightarrow{P} 1$ as $n \rightarrow \infty$
2. By CLT and delta method, $\|X\| - \sqrt{n} \xrightarrow{d} N(0, 1/2)$
3. By Jensen's inequality, $\mathbb{E}[\|X\|] \leq \sqrt{\mathbb{E}[\|X\|^2]} = \sqrt{n}$
4. If $X_i \sim N(0, 1)$ or $\max\|X_i\| \leq C$, then
 $\|X\| - \mathbb{E}\|X\|$ has subgaussian norm $O(1)$

Thm 3. Let $X = (X_1, \dots, X_n)^T$ be a random vector with X_1, \dots, X_n indep. $\mathbb{E}[X_i^2] = 1, \max\|X_i\|_{\psi_2} \leq K$. Then,

$$\|\|X\|_2 - \sqrt{n}\|_{\psi_2} \leq CK^2$$

Pf. Assume $K \geq 1$. Let $Y_i = X_i^2 - 1$, Y_i is subexponential with $\|Y_i\|_{\psi_2} \leq CK^2$. Applying Bernstein's inequality,

$$P\left(\left|\frac{1}{n}\|X\|^2 - 1\right| \geq u\right) \leq 2\exp\left(-\frac{Cn}{K^2} \min(u, u^2)\right), \text{ for all } u > 0$$

More note $|z - 1| \geq s \Rightarrow |z^2 - 1| \geq \max(s, s^2)$, $\forall s > 0$.

Then we obtain $P\left(\left|\frac{1}{n}\|X\|^2 - 1\right| \geq s\right)$

$$\leq P\left(\left|\frac{1}{n}\|X\|^2 - 1\right| \geq \max(s, s^2)\right)$$

$$\leq 2\exp\left(-\frac{Cn}{K^2} s^2\right)$$

Changing variables to $t = s\sqrt{n}$, we therefore obtain

$$P(\|X\|_2 - \sqrt{n} \geq t) \leq 2\exp\left(-\frac{Ct^2}{K^2}\right), \quad \forall t > 0$$

Rank. High-dim Gaussian is very close to uniform distribution on the sphere of radius \sqrt{n} .

Uniform spherical random variable

$r \mathbb{S}^{n-1}$: sphere in \mathbb{R}^d with radius r .

Thm 3: $N(0, I_n)$ is similar to $\text{Unif}(\sqrt{n} \mathbb{S}^{n-1})$ in high dimensions.

If $G \sim N(0, I_n)$, then

$$G = UR \text{ where } U = G/\|G\|, R = \|G\|.$$

Satisfies $U \sim \text{Uniform}(\mathbb{S}^{n-1})$, U and R are independent.

Thm 4. Uniform spherical distribution is subgaussian.

$X \sim \text{Unif}(\sqrt{n} \mathbb{S}^{n-1})$ satisfies

$$\sup_{u \in \mathbb{R}^n, \|u\|=1} \|\langle X, u \rangle\|_{\psi_2} \leq C$$

Namely, projecting X using any unit vector gives a subgaussian random variable with a bounded norm

4. Random Projection

A large collection high-dim points $Q := \{x_i\}_{i \in N} \subset \mathbb{R}^d$.

Goal: Dimensionality Reduction.

Thm 5 · Johnson-Lindenstrauss Lemma 2cc $\mathcal{E}_G(0, \frac{1}{2})$.

Take errors $C > 0$ s.t. $2cc k > (C \log N) / \varepsilon^2$ be an integer. Then $\exists \tilde{f} \in L(\mathbb{R}^d, \mathbb{R}^k)$ s.t.

$$(1 - \varepsilon) \|x - y\|^2 \leq \|\tilde{f}(x) - \tilde{f}(y)\|^2 \leq (1 + \varepsilon) \|x - y\|^2$$

for all $x, y \in Q$.

PF.

Lecture 16: Norms of Random Matrices

1. Linear Algebra Prep

Def. (ε -net). Let (\mathcal{T}, d) be a metric space,

then consider a subset $K \subset \mathcal{T}$ and some $\varepsilon > 0$,

then a subset $N \subset K$ is called a ε -net of K

if $\forall x \in K, \exists x_0 \in N$ s.t. $d(x_0, x) < \varepsilon$.

Def. Covering Number (the number of points in the smallest ε -net of K , denoted by $N(K, d, \varepsilon)$)

Prop 1. Consider unit ball $\bar{B}_2^n = \{x \in \mathbb{R}^n \mid \|x\|_2 \leq 1\} \subset \mathbb{R}^n$

$$(\frac{1}{\varepsilon})^n \leq N(\bar{B}_2^n, d_{L_2}, \varepsilon) \leq (1 + \frac{2}{\varepsilon})^n$$

Lemma 1. $\forall \varepsilon \in (0, 1)$, N is an ε -net of S^{n-1} . then

$$\sup_{x \in N} \|Ax\|_2 \leq \|A\|_{op} \leq \frac{1}{1-\varepsilon} \sup_{x \in N} \|Ax\|_2$$

If in addition there's a ε -net M of S^{m-1} , then

$$\sup_{x \in N, y \in M} \langle Ax, y \rangle \leq \|A\|_{op} \leq \frac{1}{1-2\varepsilon} \sup_{x \in N, y \in M} \langle Ax, y \rangle$$

Pf. $\forall x \in S^{n-1}, \exists x_0 \in N$ s.t. $\|x - x_0\| < \varepsilon$. Let x be the top eigenvector of $A^T A$, then

$$\|Ax - Ax_0\| \leq \|A\| \|x - x_0\| \leq \varepsilon \|A\|$$

$$\|Ax_0\| \geq \|Ax\| - \|Ax - Ax_0\| \geq \|A\| - \varepsilon \|A\|$$

$$\|A\| \leq \frac{1}{1-\varepsilon} \|Ax_0\|$$

Rank $\|A\|$ is approximately the same on a discretized space.

2 Operator Norm of Subgaussian Random Matrix.

Theorem 1. Let $A \in \mathbb{R}^{m \times n}$ have independent subgaussian entries

(A_{ij} jointly indep. subgaussian with $\mathbb{E} A_{ij} = 0$)

Then $\|A\| \leq Ck(\sqrt{m} + \sqrt{n} + t)$ with probability

at least $1 - 2e^{-t^2}$ for certain constant C .

$$\text{and } k = \max_{ij} \|A_{ij}\|_{\psi_2}$$

Lemma 2. Let X_1, \dots, X_n be independent subgaussian

random variables with $\mathbb{E} X_i = 0$, then $\sum_{i=1}^n X_i$ is

also subgaussian with

$$\left\| \sum_{i=1}^n X_i \right\|_{\psi_2}^2 \leq C \sum_{i=1}^n \|X_i\|_{\psi_2}^2$$

Rank $\|A\|$ is roughly of order $\sqrt{m} + \sqrt{n}$ with high

probability. Assume $n \geq m$, $\sigma^2 = \mathbb{E} \|A_{ij}\|^2$

$$\|A\|_F^2 = \sum_{i,j} \|A_{ij}\|^2 \xrightarrow{\text{P}} nm\sigma^2$$

Since $\|A\|_F^2 = \sum_{j=1}^m [\sigma_j(A)]^2$, we have $\frac{1}{m} \sum_{j=1}^m \sigma_j^2(A) \xrightarrow{\text{P}} n\sigma^2$

The averaged singular value is of order $\sqrt{n}\sigma$.

With high prob \Rightarrow so does the largest singular value.

3. Covariance Estimation

Suppose $X_1, \dots, X_n \in \mathbb{R}^p$ iid random vectors with $\mathbb{E}X_i = 0$ and $\text{Cov}(X_i) = \Sigma$ exists. By LN, $\frac{1}{n} \sum_{i=1}^n X_i X_i^\top \xrightarrow{P} \Sigma$

Theorem 2. Assume $X_0 \in \mathbb{R}^p$ is a subgaussian random vector, ($\sup_{u \in S^{n-1}} \| \langle X_0, u \rangle \|_{\psi_2} \leq K$) Suppose X_1, \dots, X_n

are iid copies of X_0 , $X = (X_1, \dots, X_n)^\top \in \mathbb{R}^{n \times p}$

Then for every $t \geq 0$, with prob at least $1 - 2e^{-4t^2}$

$$\left\| \frac{1}{n} X^\top X - \Sigma \right\| \leq \max\{\delta, \delta^2\} \|\Sigma\| \quad \text{where } \delta = C \sqrt{\frac{p}{n}} + \frac{4t}{\sqrt{n}}$$

Theorem 2': Covariance Estimation (HDP).

Let $X_0 \in \mathbb{R}^p$ subgaussian ($\|\langle X_0, x \rangle\|_{\psi_2} \leq K \|\langle X_0, x \rangle\|_{L^2}$ for some $K \geq 1$, any $x \in \mathbb{R}^p$). Then

$$\mathbb{E} \left[\left\| \frac{1}{n} X^\top X - \Sigma \right\| \right] \leq CK^2 \left(\sqrt{\frac{p}{n}} + \frac{2}{\sqrt{n}} \right) \|\Sigma\|.$$

Example: Subgaussian Random Vectors.

- $X \sim N(0, I_n)$
- $X \sim \text{Unif}(\{\pm 1\}^n)$
- $X \sim \text{Unif}(\mathbb{S}^{n-1})$

Lecture 17. High-dim Statistical Phenomena

1. Principle Component Analysis.

Suppose X_1, \dots, X_n iid \mathbb{X}_0 , $\mathbb{E}[Z|X_i] = 0$. $\Sigma = \text{Cov}(X_0)$ exists.

Let μ_1, \dots, μ_p be orthonormal eigenvectors of Σ with eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$.

Then consider the empirical covariance $\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n X_i X_i^\top$ with orthonormal eigenvectors $\hat{\mu}_1, \dots, \hat{\mu}_p$ and eigenvalues $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_p \geq 0$.

1.1. Fixed Dimension

Theorem 1 Consistency of PCA in fixed dimension.

Suppose $\lambda_1 > \lambda_2$. Then as $n \rightarrow \infty$, we have

$$\max_{k \leq p} |\hat{\lambda}_k - \lambda_k| \xrightarrow{P} 0 \quad \min_{S \subseteq \{1\}} \|S\hat{\mu}_1 - u_1\|_2 \xrightarrow{P} 0$$

Rank: $\lambda_1 > \lambda_2$ ensures λ_1 has multiplicity of 1.

Pf By LHW we have $\hat{\Sigma} \xrightarrow{P} \Sigma$

By Weyl's inequality

$$\max_{k \leq p} |\hat{\lambda}_k - \lambda_k| \leq \|\hat{\Sigma} - \Sigma\|_{\text{op}}$$

Since $\|\hat{\Sigma} - \Sigma\|_{\text{op}} \rightarrow 0$, $\max_{k \leq p} |\hat{\lambda}_k - \lambda_k| \rightarrow 0$

Thm 2 Davis-Kahan Sin- Θ Theorem (Simplified)

Let u_1, \tilde{u}_1 be the first normalized eigenvectors of 2 positive-definite matrices $\tilde{\Sigma}, \tilde{\Sigma} \in \mathbb{R}^{P \times P}$.

Assume $\lambda_1 > \lambda_2$. then

$$\min_{s \in \{-1\}} \|s\tilde{u}_1 - u_1\|_2 \leq \frac{2\|\tilde{\Sigma} - \Sigma\|_{op}}{\lambda_1 - \lambda_2} \text{ eigen gap}$$

Pf. Suppose $\lambda_1 - \lambda_2 < 2\|\tilde{\Sigma} - \Sigma\|_F$, then we have

$$LHS \leq \|2u_1\|_2 = \sqrt{2} \leq \sqrt{2} \cdot \frac{2\|\tilde{\Sigma} - \Sigma\|_{op}}{\lambda_1 - \lambda_2} = RHS$$

If $\lambda_1 - \lambda_2 > 2\|\tilde{\Sigma} - \Sigma\|_{op}$.

Pick $s_0 = 1$ if $\langle \tilde{u}_1, u_1 \rangle > 0$, $s_0 = -1$ if $\langle \tilde{u}_1, u_1 \rangle < 0$

$$\begin{aligned} \min_{s \in \{-1\}} \|s\tilde{u}_1 - u_1\|_2 &\leq \|s_0\tilde{u}_1 - u_1\|_2 = \sqrt{2 - 2\langle s_0\tilde{u}_1, u_1 \rangle} \\ &\leq \sqrt{2 - 2\langle s_0\tilde{u}_1, u_1 \rangle} \sqrt{1 + \langle s_0\tilde{u}_1, u_1 \rangle} \\ &= \sqrt{2} \sqrt{1 - \langle s_0\tilde{u}_1, u_1 \rangle^2} \\ &= \sqrt{2} \sqrt{1 - \langle \tilde{u}_1, u_1 \rangle^2} \end{aligned}$$

Since u_1, \dots, u_p is an orthonormal basis

$$\langle u_1, \tilde{u}_1 \rangle^2 + \sum_{k=2}^p \langle u_k, \tilde{u}_1 \rangle^2 = 1.$$

It suffice to show $\sum_{k=2}^p \langle u_k, \tilde{u}_1 \rangle^2 \leq \frac{4\|\tilde{\Sigma} - \Sigma\|_{op}^2}{(\lambda_1 - \lambda_2)^2}$

$$\langle u_k, \tilde{\Sigma} \tilde{u}_1 \rangle = \tilde{\lambda}_1 \langle u_k, \tilde{u}_1 \rangle$$

$$\langle \tilde{u}_1, \tilde{\Sigma} u_k \rangle = \lambda_k \langle u_k, \tilde{u}_1 \rangle$$

$$\Rightarrow \langle u_k (\tilde{\Sigma} - \Sigma) \tilde{u}_1 \rangle = (\tilde{\lambda}_1 - \lambda_k) \langle u_k, \tilde{u}_1 \rangle$$

$$\tilde{\lambda}_1 - \lambda_k \geq \tilde{\lambda}_1 - \lambda_2 \geq \lambda_1 - \lambda_2 - \frac{1}{2} \sum_{k=3}^n \| \tilde{z}_k \|_{\text{op}}^2 \geq \frac{1}{2} (\lambda_1 - \lambda_2)$$

where $\lambda_1 - \lambda_2 \geq 2 \| \tilde{z} \|_2^2 - \sum \| \tilde{z}_k \|_{\text{op}}^2$. Thus we have

$$\begin{aligned} \sum_{k=1}^n \langle \tilde{z}_k, \tilde{u}_k \rangle^2 &\leq \sum_{k=0}^n \frac{\langle \tilde{u}_k, (\tilde{z} - \tilde{z}) \tilde{u}_k \rangle^2}{(\tilde{\lambda}_1 - \lambda_k)^2} \leq \frac{4 \| (\tilde{z} - \tilde{z}) \tilde{u}_1 \|_2^2}{(\lambda_1 - \lambda_2)^2} \\ &\leq \frac{4 \| \tilde{z} - \tilde{z} \|_{\text{op}}^2}{(\lambda_1 - \lambda_2)^2} \end{aligned}$$

Davis-Kahan Theorem (HDPP) Let S and T be symmetric matrices with the same dimensions. Assume

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

Then the angle $\sin(\nu_i(S), \nu_i(T)) \leq \frac{2 \| S - T \|_F}{\delta}$
where $\nu_i(\cdot)$ gives col i-th largest eigenvector.

1.2. Growing Dimension.

P_n, \tilde{z}_n depend on n . $X_0 \sim N(0, \tilde{z}_n)$

Spiked model: $\tilde{z}_n = (\lambda_1 - \lambda_2) u_1 u_1^\top + \lambda_2 \tilde{L}_n$

eigenvalues are λ_1 and λ_2 with multiplicity $p-1$

Note that $\frac{1}{2} \min_{S \in \{I\}} \| S \tilde{u}_1 - u_1 \|_2^2 = 1 - | \langle \tilde{u}_1, u_1 \rangle |$

Inconsistency of PCA in high dimensions:

Assume $\| \tilde{z}_n \|_{\text{op}} \leq C$ and $\lambda_1(\tilde{z}_n) - \lambda_2(\tilde{z}_n) \geq k > 0$

holds for certain constants C and k

Suppose $\frac{P_n}{n} \rightarrow \gamma \in [0, +\infty)$

- If $\gamma = 0$, then there exists a constant $C_0 > 0$ such that $\min_{1 \leq i \leq n} \|S_{ii} - u_i\|_2 \xrightarrow{P} 0$ as $n \rightarrow \infty$
- If $\gamma > 0$, under sparsity model assumption we have $| \langle \hat{u}_i, u_i \rangle | \xrightarrow{P} \eta_\gamma$, where $\eta_\gamma < 1$. Moreover, $\exists \gamma^*$. if $\gamma > \gamma^*$, $| \langle \hat{u}_i, u_i \rangle | \not\rightarrow 0$

Rank population and sample principle direction are almost orthogonal in high dims.

2. Linear Regressions in High Dims.

Consider (x_i, y_i) , $x_i \in \mathbb{R}^p$ and $i = 1, 2, \dots, n$.

$$y_i = x_i^\top \beta + \varepsilon_i, \quad i = 1, 2, \dots, n.$$

$$X = (x_1, \dots, x_n)^\top \in \mathbb{R}^{n \times p} \quad y = (y_1, \dots, y_n)^\top \in \mathbb{R}^n$$

$$\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^\top \in \mathbb{R}^n, \quad \mathbb{E}[\varepsilon] = 0, \quad \text{Var}(\varepsilon) = \sigma^2 I$$

$$\cdot LS_{\mathbb{R}}: \hat{\beta} = (X^\top X)^{-1} X^\top y, \quad \text{Prediction: } \hat{f}(x_0) = x_0^\top \hat{\beta}$$

Under Parameterized Regime

Proposition 1. Suppose $x_0, x_1, \dots, x_n \sim N(0, \Sigma_p)$.

Assume $P = P_n$, $\frac{P_n}{n} = o(1)$ as $n \rightarrow \infty$, we have

$$\mathbb{E}_t [\|\hat{\beta} - \beta\|^2 | X] = (1 + o(1)) \frac{\sigma^2 p_n}{n}$$

$$\mathbb{E}_t [(x_0^\top \hat{\beta} - x_0^\top \beta)^2 | X] = (1 + o_p(1)) \frac{\sigma^2 p_n}{n} \quad \text{as } n \rightarrow \infty$$

Rank. The ms.e increases approximately linear in p_n

Overparameterized Regime:

Minimum-Norm Interpolator:

$$\min_{\beta \in \mathbb{R}^p} \|\beta\|^2$$

subject to $y = X\beta$

Assume $\text{rank}(X) = n$. The solution is $\hat{\beta} = \underbrace{X^\top (X X^\top)^{-1} y}_{\text{unbiased}}$

$\hat{\beta}\beta^* \approx X\beta^* = y$, $\exists v \in N(X)$ s.t. $\beta^* = \hat{\beta} + v$.

Note that $\hat{\beta} \in C(X^\top)$ and $C(X^\top) \perp N(X)$

$$\text{we have } \|\beta^*\|^2 = \|\hat{\beta}\|^2 + \|v\|^2 \approx \|\hat{\beta}\|^2$$

Proposition 2. Suppose $x_1, \dots, x_n \sim N(0, I_p)$ and $p = p_n$

depends on n and $n/p_n = o(1)$ as $n \rightarrow \infty$. We have

$\text{rank}(X) = n$ with probability $1 - o(1)$ and

$$\mathbb{E}_t [\|\hat{\beta} - \beta\|^2 | X] = (1 + o_p(1)) \cdot \left[(1 - \frac{n}{p_n}) \|\beta\|^2 + \frac{\sigma^2 n}{p_n} \right]$$

$$\mathbb{E}_t [(x_0^\top \hat{\beta} - x_0^\top \beta)^2 | X] = (1 + o_p(1)) \cdot \underbrace{\left[(1 - \frac{n}{p_n}) \|\beta\|^2 + \frac{\sigma^2 n}{p_n} \right]}_{\text{Bias}}$$

Variance:

Smaller when $p_n \gg n$

3. Community Detection $\tilde{z}^* \in \{\pm 1\}^n$: unknown membership

Stochastic Block Model. Let $A \in \mathbb{R}^{n \times n}$ be the adjacency matrix, $P_{ij} | A_{ij} = 1) = \begin{cases} p & \text{if } z_i^* = z_j^* \\ q & \text{if } z_i^* \neq z_j^* \end{cases}$ with $p > q$. A_{ij} jointly indep. Assume each community has equal size: $\mathbf{1}^\top \tilde{z}^* = 0$

Spectral Method.

Eigenvalues of A :

$$A^* = \mathbb{E}[A] = \frac{p-q}{2} Z^* Z^{*\top} + \frac{p+q}{2} \mathbf{1} \mathbf{1}^\top$$

It has only 2 positive eigenvalues:

$\lambda_1^* = \frac{(p+q)n}{2}$ and $\lambda_2^* = \frac{(p-q)n}{2}$ with corresponding eigenvectors $u_1^* = \mathbf{1}_n / \sqrt{n}$ and $u_2^* = Z^* / \sqrt{n}$

Spectral Clustering:

1. Calculate the second eigenvector $u = u_2$ of A
if p, q are known, use $A - \frac{p+q}{2} \mathbf{1} \mathbf{1}^\top$
2. Define $\tilde{z}_i = 1$ if $u_i \geq 0$, $\tilde{z}_i = -1$ if $u_i < 0$

3.1. Theoretical guarantee of spectral method.

Thm 5. For any $\varepsilon > 0$, we have

$$\Pr_{S \in \{ \pm 1 \}^n} \left(\frac{1}{\sqrt{n}} \| S \tilde{z} - \tilde{z}^* \|_2 > \varepsilon \right) = o(1)$$

Rank 2 says the mismatch ratio is arbitrarily small if $n \rightarrow \infty$

Lecture 18: Stein's phenomenon and Shrinkage estimation.

1. Motivation

1.1. High-D Perspective

$$\text{Let } \mathbf{x} = \boldsymbol{\mu} + \mathbf{z}, \quad \mathbf{z} \sim N_p(0, I_p)$$

Hoeffding:

$$P(|\langle \boldsymbol{\mu}, \mathbf{z} \rangle| > t) \leq 2 \exp\left(-C \frac{t^2}{\|\boldsymbol{\mu}\|^2}\right)$$

When p is large:

$$\text{if } \|\boldsymbol{\mu}\| = O(1) \Rightarrow \langle \boldsymbol{\mu}, \mathbf{z} \rangle = O_p(1)$$

$$\|\mathbf{x}\|^2 = \|\boldsymbol{\mu}\|^2 + \|\mathbf{z}\|^2 + 2\langle \boldsymbol{\mu}, \mathbf{z} \rangle \approx \|\boldsymbol{\mu}\|^2 + p$$

$\|\mathbf{x}\| \gg \|\boldsymbol{\mu}\|$ with high probability

1.2. Bias-Variance trade-off perspective

$$R_\varepsilon(\boldsymbol{\mu}) = \mathbb{E}[\|\mathbf{x} - \boldsymbol{\mu}\|^2] \quad \mathbf{x} \sim N(\boldsymbol{\mu}, I_p).$$

$$\text{Consider } \tilde{\boldsymbol{\mu}} = (1-\varepsilon)\mathbf{x}$$

$$\begin{aligned} \mathbb{E}[\|\tilde{\boldsymbol{\mu}} - \boldsymbol{\mu}\|^2] &= \varepsilon^2 \|\boldsymbol{\mu}\|^2 + (1-\varepsilon)^2 p \\ &= (1-2\varepsilon)p + \varepsilon^2(\|\boldsymbol{\mu}\|^2 + p) \\ &= (1-2\varepsilon)p + O(\varepsilon^2) \end{aligned}$$

$$\text{Goal: } \min_{\varepsilon \in [0, 1]} \{ \varepsilon \|\boldsymbol{\mu}\|^2 + (1-\varepsilon)^2 p \}.$$

$\|\boldsymbol{\mu}\|$ is unknown, but we can use $\|\mathbf{x}\|^2 \approx \|\boldsymbol{\mu}\|^2 + p$

when p is large to estimate $\|\boldsymbol{\mu}\|^2$

2. Shrinkage Estimators signal $\mu \uparrow$

James-Stein Estimator: noise $p \downarrow \Rightarrow \varepsilon \downarrow$

$$\hat{\mu}_{JS} = (1 - \frac{p-2}{\|x\|^2})x \quad \hat{\mu}_{JS+} = (1 - \frac{p-2}{\|x\|^2})_+ x$$

Rank · Shrinkage to 0 is not the only direction where it works. For any $c \in \mathbb{R}^p$, we can make

$$\hat{\mu}_{JS} = (1 - \frac{\kappa(p-2)}{\|x-c\|^2})(x-c) + c$$

$$\hat{\mu}_{JS+} = (1 - \frac{\kappa(p-2)}{\|x-c\|^2})_+ (x-c) + c$$

Theorem 1. If $\hat{\mu}$ is any of the two shrinkage estimators then for $\gamma \in (0, 2)$, when $p \geq 3$

$$R_{\hat{\mu}}(\mu) < R_x(\mu) \text{ for every } \mu \in \mathbb{R}^p$$

Moreover, for $\hat{\mu} = \hat{\mu}_{JS}$, we have

$$R_{\hat{\mu}}(\mu) = p - (2\gamma - \gamma^2)(p-2)^2 \bar{E}[\bar{E}[\|x-c\|^2]]$$

Theorem 2. Stein's lemma if $g \in C^1(\mathbb{R}^p)$ with $\bar{E}[\|\nabla g(z)\|] < \infty$

where $z \sim N(0, I_p)$, then

$$\bar{E}[z^T g(z)] = \bar{E}[z g(z)]$$

$$\text{Proof: } \bar{E}[z^T g(z)] = \int z^T g(z) \frac{1}{\sqrt{2\pi}} e^{-\|z\|^2/2} dz$$

$$(\text{integration by parts}) = - \int g(z) \frac{1}{\sqrt{2\pi}} z^T e^{-\|z\|^2/2} dz$$

$$= \int g(z) \frac{1}{\sqrt{2\pi}} z_i e^{-\|z\|^2/2}$$

$$= \bar{E}[z_i g(z)].$$

Rank General case: $g \in C^1(\mathbb{R}^p, \mathbb{R}^p)$, then

$$\bar{E}[D \cdot g(z)] = \bar{E}[z^T g(z)].$$

3. Tucke's Formula.

3.1. Bayesian Perspective.

Suppose μ has a prior distribution π_μ

$X|\mu \sim N(\mu, \Sigma_p)$. The marginal density of X is:

$$f_x(x) = \int_{\mathbb{R}^p} f_{x|\mu}(x|\mu) d\mu = \frac{1}{(2\pi)^{p/2}} \int_{\mathbb{R}^p} \exp(-\|x - \mu\|^2/2) \pi_\mu(\mu) d\mu.$$

Consider $\phi: \mathbb{R}^p \rightarrow \mathbb{R}^p$

$$\text{MSE} = \bar{E}\|\mu - \phi(x)\|^2 = \bar{E}_x[\bar{E}_{\mu|x=x}[\|\mu - \phi(x)\|^2]]$$

$$\bar{E}_{\mu|x=x}[\|\mu - \phi(x)\|^2] \geq \bar{E}_{\mu|x=x}[\|\mu - \bar{E}_{\mu|x=x}[\mu]\|^2], \forall x$$

We only need to choose $\phi(x) = \bar{E}[\mu | x=x]$

$$\bar{E}[\mu | x=x] = \frac{\int \mu \bar{f}_{x|\mu}(x|\mu) \pi_\mu(\mu) d\mu}{\int \bar{f}_{x|\mu}(x|\mu) \pi_\mu(\mu) d\mu}$$

$$\bar{E}[\mu | x=x] = \frac{\int (\mu - x) e^{-\|\mu - x\|^2/2} \pi_\mu(\mu) d\mu}{\int e^{-\|\mu - x\|^2/2} \pi_\mu(\mu) d\mu}$$

$$= \frac{\nabla \tilde{f}_x(x)}{\tilde{f}_x(x)} = \nabla \log \tilde{f}_x(x)$$

$$\psi_{\mu}(x) = \tilde{E}[\mu | x] = x + \nabla \log \tilde{f}_x(x)$$

Rank. From this perspective, shrinkage is a prior on the unknown μ .

3.2. Shrinkage under Superharmonic Function

Def. For $f \in C^2(\mathbb{R}^p)$, denote the Laplace operator by $\Delta f(x) = \sum_{j=1}^p \partial_j^2 f(x)$. We say f is superharmonic iff $\Delta f(x) \leq 0$, $\forall x \in \mathbb{R}^p$

Thm. Let $f \in C^2(\mathbb{R}^p)$, $f > 0$, if f is superharmonic

if $\tilde{E}\left[\frac{1}{f(x)} \sum_{j=1}^p |\partial_j^2 f(x)|\right] < \infty$, $\tilde{E}\left[\|\nabla \log f(x)\|^2\right] < \infty$

we have $\tilde{E}\left[\|x + \nabla \log f(x) - \mu\|^2\right] = p + 4\tilde{E}\left[\frac{\Delta \log f(x)}{f(x)}\right] \leq p$

Rank. x is a minimax estimator

$\Rightarrow x + \nabla \log f(x)$ is also a minimax estimator

4. Diffusion Model:

Forward Process: $X_t = \sqrt{1 - \sigma_t^2} X_{t-1} + \sigma_t Z_t$

$Z_t \text{ iid } N(0, I)$.

Score Function P_t : pdf of X_t

$$\tilde{E}[X_{t+1} | X_t = x_t] = \frac{1}{\sqrt{1 - \sigma_t^2}} x_t + \frac{\sigma_t^2}{\sqrt{1 - \sigma_t^2}} \underbrace{\nabla \log P_t(x_t)}_{\text{Score Function}}$$

$$X_{t+1} | X_t = x_t \sim N\left(\frac{1}{\sqrt{1 - \sigma_t^2}} x_t + \frac{\sigma_t^2}{\sqrt{1 - \sigma_t^2}} \nabla \log P_t(x_t), \sigma_t^2 I\right)$$

$$\text{Reconstruction: } \hat{X}_{t+1} = \frac{1}{\sqrt{1 - \sigma_t^2}} X_t + \frac{\sigma_t^2}{\sqrt{1 - \sigma_t^2}} S(\theta, X_t) + \sigma_t Z_t$$

$S(\theta, X_t)$ is trained to approximate $\nabla \log P_t(X_t)$:

$$\min_{\theta} \sum_{t=1}^T \mathbb{E}_{Z_t \sim N(0, I), X_0 \sim p_0, X_t \sim p(X_t | X_0)} [\sigma_t \| \nabla \log P_t(X_t) - S_\theta(X_t, t) \|^2]$$

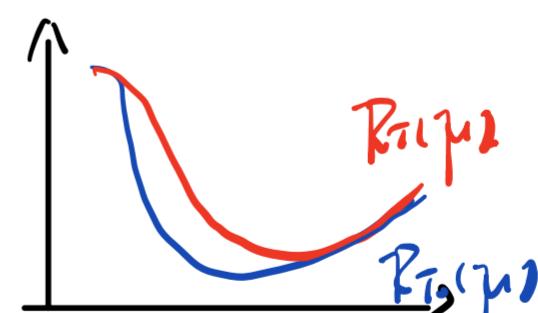
Result:

Admissibility:

$T(x)$ is admissible if there's no T_0 sc.

① $R_{T_0}(\mu) \leq R_T(\mu)$ for all μ

② $R_{T_0}(\mu) < R_T(\mu)$ for some μ



Minimality: $T(x)$ is minimax if for any $T_0(x)$

$$\sup_{\mu} R_{T_0}(\mu) \leq \sup_{\mu} R_T(\mu)$$

Setting $X \sim N(\mu, \Sigma_p)$, estimate μ

① Is $\hat{\mu}$ minimax? Yes.

The supremum of risk is controlled even if the risk is point-wise smaller than $\hat{\mu}$ -s estimator

② Is $\hat{\mu}$ admissible?

Yes when $p=1, 2$	}

Consider $r=1$. For any $p \geq 3$:

$$R_{\hat{\mu}_{SS}}(\mu) = p - (p-2)^2 \underbrace{E[\|X\|^{-2}]}_{\|X\|^2 \sim \Sigma_p} < R_{\hat{\mu}}(\mu) = p$$

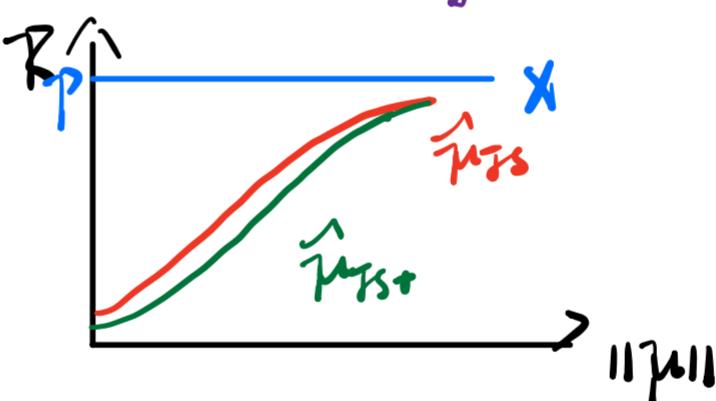
Two extreme case: $\|X\|^2 \sim \Sigma_p$ $\underbrace{E[\|X\|^{-2}]}_{\|X\|^2 \sim \Sigma_p} = \frac{1}{p-2}$.

$$R_{\hat{\mu}_{SS}}(0) = p - (p-2)^2 \frac{1}{p-2}$$

Turn why we need $p \geq 3$

$$R_{\hat{\mu}_{SS}}(\infty) = p$$

$$\cdot R_{\hat{\mu}_{SS}} > R_{\hat{\mu}}$$



$\Rightarrow \hat{\mu}_{SS}$ is inadmissible

Remark. $\hat{\mu}_{SS+}$ is also inadmissible

Empirical Bayes Perspective:

$$X \sim N(\mu, \Sigma_p) \quad \mu \sim N(0, \underbrace{\Sigma_p}_{hyperprior})$$

Bayes estimator. $\delta(x) = \arg \min_{\delta(x)} \mathbb{E}_{\mu \sim \Sigma_p} [R_{\delta}(x)] = (1 - \frac{1}{1 + \frac{1}{\mu^2}}) x$

$$P(x|\tau) = N(0, (1 + \tau^2)^{-1} I_p)$$

Empirical Bayes: estimate τ from data using $P(x|\tau)$

- UMVUE of $\frac{1}{1+\tau^2}$: $\frac{P-2}{\|x\|^2}$

Plug in: $\hat{g}(x) = (1 - \frac{P-2}{\|x\|^2})x = \hat{\mu}_{BS}$

- MLE of $1 + \tau^2$ is $\frac{\|x\|^2}{P}$

Plug in: $\hat{g}(x) = (1 - \frac{P}{\|x\|^2})_+ x \approx \hat{\mu}_{BS+}$

Lecture 19. Denoising via Shrinkage

1. Sparsity in data Analysis

1.2 Quantifying Sparsity.

Def. we say $x \in \mathbb{R}^n$ is s -sparse if

$$\|x\|_0 = \sum_{j=1}^n \mathbb{I}(x_j \neq 0) \leq s$$

We can also relax the definition using $\|x\|_q$, $q \in [0, \infty]$

Rank. If $q < 1$, $\|x\|_q$ is only a "Quasinorm". (Not a norm)

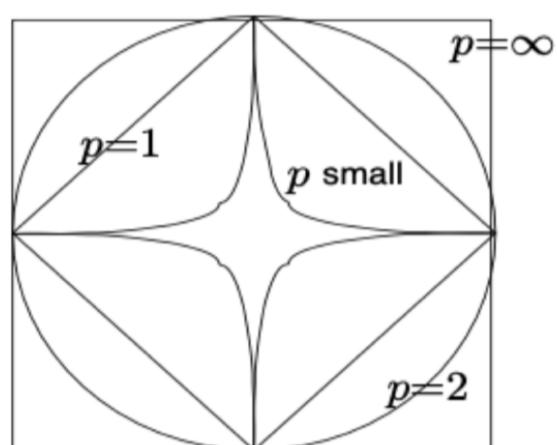


Figure 1.5 Contours of ℓ_p balls

3. Thresholding.

3.1. Gaussian sequence model.

Assume our observations follows $Y_k = \theta_k + \sigma Z_k$, $Z_k \sim N(0, 1)$

Raw measurements $X \in \mathbb{R}^n$

$$X_k = f(t_k) + \sigma W_k, \quad W_k \sim N(0, 1)$$

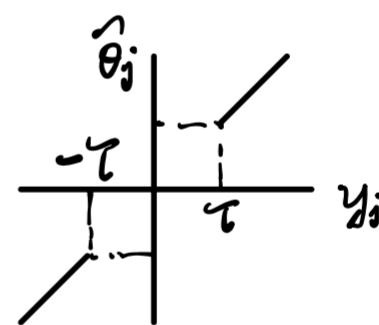
Apply orthogonal transform $Y = W X$ where W is an orthogonal matrix. $\theta = W f(t)$, $Y \sim N(\theta, \sigma^2 I_n)$

Sparcity Assumption: $\|\theta\|_0 = \sum_{j=1}^n \mathbb{I}(\theta_j \neq 0) \leq S$

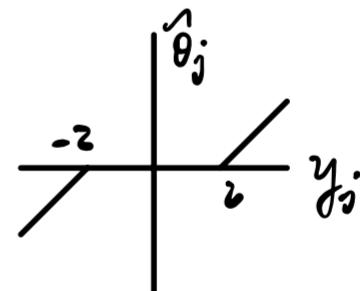
Unbiased estimator: \hat{Y} . ℓ_1 -s estimator: $(1 - \frac{(p-1)\sigma^2}{\|Y\|^2}) Y$.

3.2. Thresholding Functions.

Hard Thresholding: $\hat{\theta}_j = \begin{cases} Y_j & \text{if } |Y_j| > T \\ 0 & \text{if } |Y_j| \leq T \end{cases}$



Soft Thresholding: $T_\tau(u) = \text{Sign}(u)(|u| - \tau)_+$



Determining T : $\|\tilde{Z}\|_\infty \geq \sqrt{2 \log n}$

We may choose $T = \sqrt{2 \log n}$

Further, by Gaussian Concentration inequality: $\max_i \tilde{Z}_i$ is

close to $\sqrt{2 \log n}$, and

$$\sqrt{2 \log n} - 1 - \sqrt{2 \log 2} \leq \frac{1}{\sigma} \sqrt{\max_i \tilde{Z}_i} \leq \sqrt{2 \log n}$$

So choosing $\tau = (1+o(1))\sqrt{2\sigma^2 \log n}$ is guaranteed to remove all noise with high probability.

4. Denoising theory:

Oracle-aided ideal risk: $R(\theta, \sigma) = \min_{i=1}^n \{\theta_i^2, \sigma^2\}$

If θ is S -sparse: $R(\theta, \sigma) = S\sigma^2$

Theorem 1. Suppose $Y \sim N(\theta, \sigma^2 I_n)$. For soft-thresholding estimator $\hat{\theta}_{\tau_n} = T_{\tau_n}(Y)$ with rule level $\tau_n = \sqrt{2\sigma^2 \log n}$,

$$\mathbb{E}[\|\hat{\theta}_{\tau_n} - \theta\|^2] \leq (2\log n + 1)[\sigma^2 + R(\theta, \sigma)]$$

Rank. We don't know the indices of the large signals. The risk is increased by a log factor compared to the ideal risk

4.2 Oracle Inequalities for S -sparse Signals

Theorem 2. Suppose $S = S_n$ satisfies

$$\frac{S_n}{n} \rightarrow 0, \quad n \rightarrow \infty.$$

Then, as $n \rightarrow \infty$,

- (1) The soft-thresholding estimator $\hat{\theta}_{\tau_n}$ with $\tau_n = \sqrt{2\sigma^2 \log(n/S_n)}$,

$$\sup_{\theta \in \Theta_{n,s}} \mathbb{E}[\|\hat{\theta}_{\tau_n} - \theta\|^2] \leq (1 + o(1)) \cdot 2S_n \sigma^2 \log\left(\frac{n}{S_n}\right).$$

- (2) The soft-thresholding estimator matches the minimax lower bound:

$$\inf_{\hat{\theta}} \sup_{\theta \in \Theta_{n,s}} \mathbb{E}[\|\hat{\theta}_{\tau_n} - \theta\|^2] \geq (1 + o(1)) \cdot 2S_n \sigma^2 \log\left(\frac{n}{S_n}\right).$$

Final · 3 problems with subproblems, 16 subproblems in total

Most · No more than 3 lines of proofs

Lecture 20: Basis pursuit, compressed sensing

1. Basis Pursuit identifiable?

Find sparse solution θ^* to $Y = X\theta^*$, $X \in \mathbb{R}^{n \times p}$, $n < p$

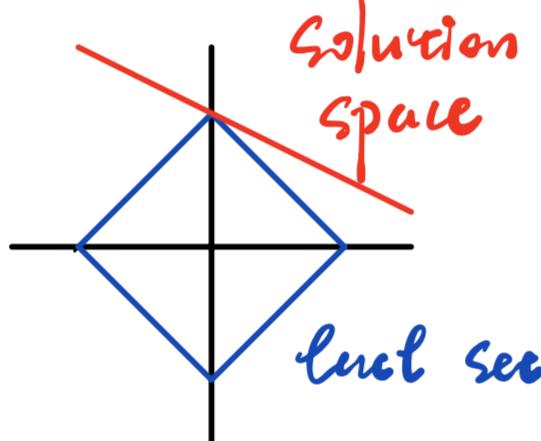
$$\min_{\theta \in \mathbb{R}^p} \|\theta\|_1$$

affine linear subspace

subject to $y = X\theta \Leftrightarrow \theta \in \theta^* + \text{Ker}(X)$

1.1. Restricted null space property.

Consider the level set $S_r := \{\theta : \|\theta\|_1 = r\}$.



Rank High-dim ℓ_1 balls are sharp.

In \mathbb{R}^p , the probability of failing to eliminate some parameters is $\frac{1}{2^{p-1}}$

RNS: Formulating geometric intuitions.

$\text{Zer } S = \{j \in p : \theta_j^* \neq 0\}$.

$$C(S) = \{\Delta \in \mathbb{R}^p : \|\Delta_{S^c}\|_1 \leq \|\Delta_S\|_1\}.$$

where $\Delta_S = (\Delta_j \mathbb{1}_{(j \in S)})_{j \in p}$

We say X satisfies restricted null space (RNS)

if $\text{Ker}(X) \cap C(S) = \{0\}$.

Proposition 1. X satisfies RNS, then BP recovers θ^*

1.2. Pairwise incoherence

Def Pairwise incoherence. Let $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_p) \in \mathbb{R}^{n \times p}$.

The pairwise incoherence is:

$$\delta_{pw}(\mathbf{X}) = \left\| \frac{1}{n} \mathbf{X}^T \mathbf{X} - \mathbf{I}_p \right\|_{max} = \max_{k,j} \left| \frac{1}{n} \langle \mathbf{x}_k, \mathbf{x}_j \rangle - \mathbb{I}(k=j) \right|$$

Proposition. Suppose $\mathbf{X} = (\mathbf{x}_{kj})_{k \in [n], j \in [p]}$, where

$\mathbf{x}_{kj} \sim \mathcal{N}(0, 1)$. Then, for any $\delta \in (0, 1)$, with

probability at least $1 - \delta$,

$$\delta_{pw}(\mathbf{X}) \leq C \frac{\sqrt{\log p + \log(1/\delta)}}{\sqrt{n}}$$

Rmk. Given exponentially many random normalized vectors, all pairs of vectors are nearly orthogonal.

Thm. If \mathbf{X} has pairwise incoherence, $\delta_{pw}(\mathbf{X}) < \frac{1}{2k}$.

and $|S| \leq k$, then \mathbf{X} satisfies RNS property.

Gr. Assume $\mathbf{x}_{kj} \sim \mathcal{N}(0, 1)$ and θ^* is k -sparse. Then BP recovers θ^* with high probability if

$$n \geq Ck^2 \log p$$

where C is an absolute constant