

Lecture 6. Statistical Decision Theory and Sufficiency

Basic Terminology

$$\text{Data} \sim (\Omega, \mathcal{F}, P)$$

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

Example. Graph Data f_2 -R model $G(n, p)$:

Graph G , n vertices, A : adjacency matrix.

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

$A \approx E[A]$ under some conditions. $T(A) \approx T(E[A])$

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 $\Theta \subset \mathbb{R}$, then we only need to consider $\{\frac{dP_\theta}{d\nu} : \theta \in \Theta\}$, particularly pdf / pmf

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

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

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

Example (Density estimation): Each distribution is abs cont and the set of pdf is $\{f(x) : \int f(x) < L\}$ for some $L > 0$

Semi-parametric Family: True parameters of interest
+ Intrinsic Nonparametric Parameters

Example: Linear model with unknown dist of ε

Identifiability:

Def. $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$ 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 $(\hat{\beta}_1 + c, \dots, \hat{\beta}_K + c)$ stipulate 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:

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

Exponential Family:

Def. Dse 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, Dses subsumed in exponential family

Properties: (1) Marginal / Conditional dist 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.

Cor. $m(\eta) \in C^\infty$ on Ω

Basic Terminologies in Statistical Decision Theory.

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

Optimal Rule

Def. \bar{T}_* is an optimal rule if

$R_{T_*}(P) \leq R_T(P), \forall P \in \mathcal{P}, \forall T \in \{\text{class of allowable decisions}\}$

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

Randomized decision rules.

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

Often, $\delta = \sum_{j=1}^n I(Z=j) \bar{T}_j$ where Z is a discrete random variable

- Risk of randomized rule: $R(\theta, \delta) = \bar{E}_{x \sim P_\theta} [\bar{E}_{z \sim \delta(x)} [L(\theta, A)]]$

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 P , T_π minimizes

$$R_T(\pi) = \int R_T(p) d\pi(p)$$

- Minimax rule: $T_0 = \arg \min_{T_\pi} (\sup_p R_T(p))$

Estimation.

For i.i.d X_1, \dots, X_n , w/k $\bar{E}[X_i]$. Use squared loss

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

$$\bar{R}_T(P) = \frac{\sigma^2}{n} \text{ where } \sigma^2 = \text{Var}(X_i)$$

$$\text{Bias: } b_T(P) = \bar{E}_P[T(X)] - \theta$$

Hypothesis Test

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

Let $T(X) \in \{0, 1\}$, $T(X)$ must have the form $T_C(X) = \{X : T(X) = 1\}$.

where C is called the rejection region. Consider the loss

$$L(P, j) = 0 \text{ if } P \in \mathcal{P}_j, \text{ and } 1 \text{ otherwise}$$

$$\text{Then } R_T(P) = \begin{cases} P(T(X) = 1) = P(X \in C), & P \in \mathcal{P}_0 \\ P(T(X) = 0) = P(X \notin C), & P \in \mathcal{P}_1 \end{cases}$$

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

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

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

Connection to information theory:

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

where 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: Simple mean. $X_1, \dots, X_n \sim i.i.d. (\mu, \sigma^2)$

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 least square

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

$$\Rightarrow \hat{\beta}_\lambda = (X^\top X + \lambda I_p)^{-1} X^\top y \quad (X = (x_1, \dots, x_n)^\top \in \mathbb{R}^{n \times p})$$

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

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

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

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

Exercise:

Bias:

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

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

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

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

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

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

Variance

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

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

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

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

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

where μ_i are eigenvalues of $X^T X$

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

$\overline{\operatorname{Var}}(S) : \sigma^2 \operatorname{Tr}((X^T X)^{-1}) - \operatorname{Var} = O(\lambda)$

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: $\hat{f}_k(x_0) = \frac{1}{K} \sum_{k=1}^K y_{x_{(k)}}$

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

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

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 X :

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

The variance generally increases as p increases

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

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 $\text{rank}(X^T) = n$. Then,

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

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

Rank: GD prefers the solution. Because

$$\begin{aligned} \nabla L(\hat{\beta}) &= 2X^T(Y - X\hat{\beta}) \in C(X^T) \Rightarrow \hat{\beta}^{(t)} \in C(X^T) \\ &\Rightarrow \lim_{t \rightarrow \infty} \hat{\beta}^{(t)} \in C(X^T) \end{aligned}$$

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

Thm (GD for overparametrized 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

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

$$\begin{aligned} \frac{\partial L(\beta)}{\partial \beta} &= \frac{1}{n} \sum_{i=1}^n 2(-x_i)(y_i - x_i^\top \beta) + 2\lambda \beta \\ &= \frac{1}{n} \sum_{i=1}^n x_i(y_i - x_i^\top \beta) + 2\lambda \beta \\ &= \frac{1}{n} (\sum_i x_i y_i - \sum_i x_i^\top \beta x_i) + 2\lambda \beta \\ &= \frac{1}{n} (X^\top y - X^\top X \beta) + 2\lambda \beta \\ &= \frac{1}{n} X^\top y + (\frac{1}{n} X^\top X + 2\lambda) \beta \\ \frac{\partial^2 L(\beta)}{\partial \beta^2} &= \frac{1}{n} X^\top X + 2\lambda I \succeq 0 \end{aligned}$$

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

Exercise 2. Prove the ridge less solution

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

First we prove two optimization problem are the same when $\lambda \rightarrow 0$

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

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

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

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

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

$$\nabla L(\beta, \lambda) = \begin{pmatrix} 2\beta - X^T\lambda \\ X\beta - y \end{pmatrix} = 0$$

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

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

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

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

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

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

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

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

We hope to pick W s.t. 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
 $Y_i = Bf_i + u_i$, $B \in \mathbb{R}^{p \times k}$, $f_i \in \mathbb{R}^k$, $u_i \in \mathbb{R}^p$

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

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

1. Using Covariance matrix (second moment):

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

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

2. Using spectral decomposition. $\exists \tilde{L} = U \Lambda U^T$

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

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

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

Latent Variable Model

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

Topic j is drawn from $P(h=j) = w_j$, $j=1, \dots, K$

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

$X_t = i$: 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:

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

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

Goal: 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 eigenvectors $v_1, \dots, v_K, \dots, v_d$

considering 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 $\tilde{\mu}_k = (v_1 \dots v_d)^T \mu_k \quad UU^T = I$

$$\tilde{\mu}_i^T \tilde{\mu}_j = \mu_i^T (v_1 \dots v_d) (v_1 \dots v_d)^T \mu_k = \mu_i^T$$

2. Use B to transform $\bar{M}_3 = \sum_{i=k}^k \lambda_k \bar{\mu}_k \otimes \bar{\mu}_k \otimes \bar{\mu}_k$

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 , $L(\theta, y) = f_\theta(y)$

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

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)

$$\text{Note that } \mathbb{I}\{z_i z_j = 1\} = \frac{1 + z_i z_j}{2}$$

$$\max_{\mathbf{z} \in \{\pm 1\}^n : \mathbf{z}^\top \mathbf{z} = 0} \sum_{i,j} \mathbb{I}\{z_i z_j = 1\} Y_{ij}$$

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

Alternatively:

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

$$\begin{aligned} \text{2. For } M < N, \quad P(\bigcap_{n=M}^N A_n^c) &= \prod_{n=M}^N (1 - P(A_n)) \\ &\leq \prod_{n=M}^N e^{-P(A_n)} \\ &= e^{\sum_{n=M}^N -P(A_n)} \\ &\rightarrow 0 \text{ as } N \rightarrow \infty \end{aligned}$$

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

$$\begin{aligned} P\left(\bigcap_{n=M}^{\infty} A_n^c\right) = 0 \Rightarrow P\left(\bigcup_{n=M}^{\infty} A_n\right) &= 1 \quad \text{for any } M \\ \Rightarrow P\left(\bigcap_{M=1}^{\infty} \bigcup_{n=M}^{\infty} A_n\right) &= 1 \end{aligned}$$

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

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

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

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

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

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

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

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

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

Rmk. $\{X_n\}$ 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 $\Phi(\frac{C}{\sigma_n}) \leq 0.8$.

$$P(|X_n| \geq C) \geq 0.2 \text{ for large enough } n.$$

$\{|X_n|\}$ 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 $X_n = O_p(1)$

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:

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

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

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

Rmk. Finite higher moments implies finite lower moments.

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

Exercise: Converse is false.

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

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

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

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

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

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 $\bar{E}|X_1| < \infty$, $\frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{P} \bar{E}X_1$

Pf (Exercise)

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

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

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 $c > 0$

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

$$\begin{aligned} \text{Then } \bar{E}|X_1| &= \int_0^\infty P(|X_{n_k}| > t) dt \\ &\geq \int_{n_k}^\infty P(|X_{n_k}| > t) dt \\ &\geq \int_{n_k}^\infty \frac{c}{t} dt = c [\ln t]_{n_k}^\infty = \infty \end{aligned}$$

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

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

For the major part. Let $Y_{n_j} = X_j \mathbb{I}_{\{X_j \leq n\}}$. 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 $T_n = \frac{1}{n} \sum_{j \in n} Y_{nj}$

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

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

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

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

$\int_{X_j} |X_j| \leq t$

$$\int_{X_j} t dt = |X_j|^2$$

$\int_{X_j} |X_j| \geq t$

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

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

5. Strong UN. and $t P(|X_j| \geq t) \rightarrow 0$

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

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

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

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

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

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

$$P(A_1) = 1$$

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

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

7. Consistency of MoM and MLE.

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

$$\text{Var}(X_1) = \sigma^2 < \infty, \text{ Then } S_n = \frac{1}{n-1} \sum_{i=1}^{n-1} (X_i - \bar{X})^2$$

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

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

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

Consistent 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: $\hat{\mu}_n \xrightarrow{\text{a.s.}} \mu \Rightarrow h(\hat{\mu}_n) \xrightarrow{\text{a.s.}} \theta$ as $n \rightarrow \infty$

Consistency of $\hat{\mu}_k$.

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, $\hat{\mu}_k$ is \sqrt{n} -consistent.

Lesson 10. Weak Convergence and CLT

1. Weak Convergence

Equivalent definition 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 bdd ces functions h on \mathbb{R}^d
- Characteristic function $\lim_{n \rightarrow \infty} \phi_{X_n}(t) = \phi_X(t)$ for all $t \in \mathbb{R}^d$
- $X_n \xrightarrow{d} X$ iff $C^r X_n \xrightarrow{d} C^r X$, $\forall C \in \mathbb{R}^d$

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

- $X_n \xrightarrow{P} X$ then $X_n \xrightarrow{d} X$
- $X_n \xrightarrow{d} C$, $C \in \mathbb{R}^k$ then $X_n \xrightarrow{P} C$
- $X_n \xrightarrow{d} X$, 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 continuous w.r.t. μ_X , then

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

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

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

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

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

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

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

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

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

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

Slivskiy'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

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

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

(ii) Let $T_n(x)$ be an order estimator. Asymptotic relative efficiency is defined to be the ratio between the

The asymptotic Variance

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

Lecture 12. Unbiased Estimation and UMVUE.

1. Complete Statistic

Def. (Complete Statistics) $T(x)$ is complete for $P \in \mathcal{P}$
iff if and only if $E[f(T)] = 0 \Rightarrow f(T) = 0$ 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

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

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.

\checkmark gives unbiased estimator

3. Construct UMVUE

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

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

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

T is unbiased for v with $E[T(X)] < \infty$

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

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

(2) \tilde{T} is sufficient for P , let

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

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

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

Method 4*: Variational Calculus

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

1. Fisher Information

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

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

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

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

2. Gramer-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 $\tilde{v}(x)$ is an unbiased estimator of v .

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

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

holds for $h \in \mathcal{H}$ or $h(x) = \tilde{v}(x)$. Then,

$$\text{Var}(\tilde{v}(x)) \geq \left(\frac{\partial}{\partial \theta} g(\theta) \right)^T \tilde{I}(\theta) \left(\frac{\partial}{\partial \theta} g(\theta) \right)$$

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

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

is different but

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

reparameterization
under ✓

Therefore, the C-R lower bound is invariant under ✓

Rank. Asymptotic optimality of \hat{ML} .

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

Rank. C-R lower bound may not be attained

3. Interpretations of Fisher information.

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

3.2. Information theory. under certain "Regularity Conditions"

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

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

5. Examples

5.1. Exponential Family.

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

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

Proposition

1. The regularity conditions for C-R lb is satisfied

for any msb h with $E[h(x)] < \infty$

2. Consider natural parameter η , $\text{Var}(\tilde{T}(x)) = \tilde{I}(\eta)$

3. $V = E[\tilde{T}(x)]$, then $\text{Var}(V) = (\tilde{I}(V))^{-1}$

5.2 Linear Models. $X_i = Z^T \beta + \epsilon_i$ $Z \in \mathbb{R}^{n \times p}$ $\epsilon \sim N(0, \sigma^2 I_n)$

Thm 2 Suppose Z is full column rank

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

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

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

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

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

Fisher information $\{z \sim N(\theta, \sigma^2 I_n)\}$ 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 estimator $\hat{\beta}$ attains C-R lower bound.

Lecture 14. Concentration Inequality

1. Motivation Non-asymptotic bounds (finite sample)

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

$$\left(\frac{1}{t} - \frac{1}{t^3}\right) \frac{1}{\sqrt{2\pi}} e^{-t^2/2} \leq P(G \geq t) \leq \frac{1}{t} \frac{1}{\sqrt{2\pi}} e^{-t^2/2} \quad \forall t > 0$$

$\text{Let } \tilde{F}$ and \tilde{P} 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 $\bar{E}[X_i] = 0$ and $\bar{E}[X_i^3] < \infty$, then

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

WTS: Tail prob of \bar{X} is $O(e^{-nt^2/2})$

2. Subgaussian variables and Hoeffding's Inequality

Theorem 1. Subgaussian random variables

Let X be an r.v. with $\bar{E}X = 0$. Then the conditions

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

$$(2) \sup_{p \geq 1} \frac{1}{\sqrt{p}} (\bar{E}|X|^p)^{1/p} \leq k_2 \quad \text{for some } k_2 > 0$$

$$(3) \bar{E}[e^{\lambda X}] \leq e^{k_3 \lambda^2}, \quad \forall \lambda \in \mathbb{R} \quad \begin{matrix} \text{constants independent of} \\ \uparrow \text{random variable} \end{matrix}$$

are equivalent in the sense that $k_i < \underline{c}_{ij} k_j \quad \forall i \neq j$.

X is said to be a subgaussian r.v if they holds.

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 $\|X\|_{\psi_2} \leq CM$

Proof of Thm 1:

$$(1) \Rightarrow (2) \cdot \mathbb{E}|X|^p = \int_0^\infty p(|X| > t) p t^{p-1} dt$$

$$\begin{aligned} \text{w.l.o.g. } k_2 &= 1 & \leq \int_0^\infty 2pt^{p-1} e^{-t^2} \\ (\text{we can always rescale } X) & & = p \int_0^\infty u^{(p-1)/2} e^{-u} du \\ & & = p \Gamma(p/2) \leq 3p(\frac{p}{2})^{p/2} \end{aligned}$$

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

$$\text{Let } g(p) = \log(3p)^{1/p} = \frac{1}{p} \log(3p)$$

$$g'(p) = -\frac{1}{p^2} \log(3p) + \frac{1}{3p^2} = \frac{1 - 3 \log(3p)}{3p^2}$$

$$g(p) \leq g(\frac{1}{3}e^{\frac{1}{3}}) \Rightarrow \|X_p\| \leq \frac{\sqrt{2}}{2} \exp(g(\frac{1}{3}e^{\frac{1}{3}}))$$

(3) \Rightarrow (1). For $\lambda > 0$, $t > 0$

$$\begin{aligned} P(X \geq t) &\leq P(e^{\lambda X} \geq e^{\lambda t}) \leq e^{-\lambda t} \bar{E} e^{\lambda X} \\ &\leq e^{-\lambda t + \lambda^2} \leq e^{-t^2/4} \end{aligned}$$

Choose $\lambda = t/2$ to minimize $-\lambda t + \lambda^2$

Similarly, $P(-X \geq t) \leq e^{-t^2/4}$

Then, $P(|X| \geq t) \leq 2e^{-t^2/4}$

(2) \Rightarrow (3): W.L.O.G. $k_i = 1$. Let $\tilde{X} \triangleq X$ be an independent

copy of X . Then $X - \tilde{X} \triangleq \tilde{X} - X$ because

$$\phi_{X-\tilde{X}}(t) = \phi_X(t) \phi_{\tilde{X}}(-t) = \phi_{\tilde{X}}(t) \phi_X(-t) = \phi_{\tilde{X}-X}(t)$$

By Jensen's inequality we have $\bar{E} e^{\lambda X} \geq e^{\lambda \bar{E} X} = 1$

Therefore, $\bar{E} e^{\lambda X} \leq \bar{E} e^{\lambda X} \bar{E} e^{-\lambda \tilde{X}} = \bar{E} e^{\lambda(X-\tilde{X})}$

$$= \bar{E} \left[\sum_{k=0}^{\infty} \frac{(\lambda(X-\tilde{X}))^k}{k!} \right]$$

$$= 1 + \sum_{k: \text{even}}^{\infty} \frac{\lambda^k \bar{E}[(X-\tilde{X})^k]}{k!}$$

$$\leq 1 + \sum_{k: \text{even}}^{\infty} C_1 \frac{(2\lambda)^k k^{k/2}}{\sqrt{2\pi k/2} (k/(2e))^k} \quad (\text{by Stirling's}) ?$$

$$\leq 1 + \sum_{k: \text{even}}^{\infty} C_2 \frac{(2\sqrt{e}\lambda)^k}{\sqrt{2\pi k/2} (k/(2e))^{k/2}}$$

$$\leq 1 + \sum_{m=0}^{\infty} \frac{(4C_2 e \lambda^2)^m}{m!} \quad (m := k/2)$$

$$= e^{k^2 \lambda^2} \quad (\text{choose } k_2 = 4C_2 e)$$

Thm 2. Hoeffding's Inequality Let x_1, x_2, \dots, x_n be independent r.v.s s.t. $E x_i = 0$ and each x_i is subgaussian. Let $a \in \mathbb{R}^n$, $K := \max_{1 \leq i \leq n} \|x_i\|_{\psi_2}$. Then

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

holds for certain absolute constant $C > 0$

3. Subexponential Variables and Bernstein's Inequality

Thm 3. (Subexponential random variables) Let X be a random variable with $E X = 0$. Then

$$(i) P(|X| > t) \leq 2e^{-t/k_1}, \forall t \geq 0$$

$$(ii) \sup_{p \geq 1} \frac{1}{p} E[|X|^p]^{1/p} \leq k_2$$

$$(iii) E[e^{\lambda X}] \leq e^{k_3 \lambda^2}, \forall |\lambda| < 1/k_3$$

Constants independent
of the random variable X .

\uparrow

are equivalent in the sense that $k_3 \leq \underline{C}_{ij} k_j$

Remark

- X is called sub-exponential if any of the condition holds
- The smallest k_3 is called the sub-exponential norm. $\|X\|_{\psi_2}$
- Subgaussian \Rightarrow subexponential, $\|x_i\|_{\psi_1} \leq \|x_i\|_{\psi_2}$

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

$$\|XY\|_{\psi_1} \leq \|X\|_{\psi_1} \|Y\|_{\psi_1}$$

D.F. $[\bar{E}(|X|^p|Y|^p)]^{1/p} \leq [\bar{E}|X|^{2p} \bar{E}|Y|^{2p}]^{1/2p}$

$$\langle X, Y \rangle_p \leq \|X\|_{2p} \|Y\|_{2p} ?$$

Thm 4. Bernstein's Inequality

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

$K = \max_{1 \leq i \leq n} \|X_i\|_{\psi_1}$. Then

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

for some constant $c > 0$

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