

① Decision theory

{
minimax
admissibility
James-Stein

② Bayesian Statistics

{ empirical Bayes
EM
Variational Bayes
BVM

③ Model Selection

{ BIC
AIC
LASSO

④ Hypothesis testing

{ Neyman-Pearson lemma

Chernoff information

Le Cam's method

multiple comparison

{ global null
FwR
FDR

① E. Lehmann Theory of point estimation - minimality, admissibility

② A van der Vaart Asymptotic Statistics

③ I. Johnstone Gaussian sequence (asymptotic)

④ E. Candès Stat 302 C Lecture notes

Decision Theory (A. Wald)

$X_1, \dots, X_n \stackrel{\text{iid}}{\sim} P_\theta$ $\theta \in \Theta \leftarrow$ parameter space
 \hookrightarrow parameter of interest

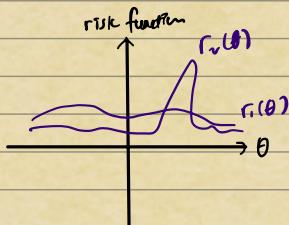
$\hat{\theta} = \hat{\theta}(X_1, \dots, X_n) \leftarrow$ estimator

Loss function $L(\hat{\theta}, \theta)$ e.g. $\|\hat{\theta} - \theta\|^2$

risk $R(\hat{\theta}, \theta) = E(L(\hat{\theta}, \theta))$

$$= \int L(\hat{\theta}(x), \theta) P_\theta(x) dx$$

two estimators: $\hat{\theta}_1, \hat{\theta}_2$
 $R_1(\theta) = R(\hat{\theta}_1, \theta)$ $R_2(\theta) = R(\hat{\theta}_2, \theta)$



idea: average risk $\int R(\theta, \theta) \pi(\theta) d\theta$
 maximum risk $\sup_{\theta \in \Theta} R(\hat{\theta}, \theta)$

def. ① $\hat{\theta}$ is a Bayes estimator if $\hat{\theta} = \arg \min_{\theta} \int R(\hat{\theta}, \theta) \pi(\theta) d\theta$
 $\Leftrightarrow V \hat{\theta} \int R(\hat{\theta}, \theta) \pi(\theta) d\theta \leq \int R(\hat{\theta}, \theta) \pi(\theta) d\theta$

② $\hat{\theta}$ is a minimax estimator if $\hat{\theta} = \arg \max_{\theta \in \Theta} R(\hat{\theta}, \theta)$

$$\Leftrightarrow V \hat{\theta}, \sup_{\theta \in \Theta} R(\hat{\theta}, \theta) \leq \sup_{\theta \in \Theta} R(\hat{\theta}, \theta)$$

Bayes estimator $\hat{\theta}_{\pi} = \arg \min_{\theta} \int R(\hat{\theta}, \theta) \pi(\theta) d\theta$

$$\int R(\hat{\theta}, \theta) \pi(\theta) d\theta = \int \int L(\hat{\theta}(\theta), \theta) P_{\theta}(x) dx \pi(\theta) d\theta$$

$$= \int \int L(\hat{\theta}(x, \theta) \underbrace{P(x|\theta)}_{\text{posterior}} \pi(\theta) dx d\theta$$

$$p(x|\theta) \pi(\theta) = \frac{p(x|\theta) \pi(\theta)}{\int p(x|\theta) \pi(\theta) d\theta} \int p(x|\theta) \pi(\theta) d\theta$$

$$= \pi(\theta|x) m(x)$$

\swarrow posterior \nwarrow marginal density

$$\int R(\hat{\theta}, \theta) \pi(\theta) d\theta = \underbrace{\int \int L(\hat{\theta}(x, \theta) \pi(\theta|x) dx)}_{\text{a function of } x} m(x) d\theta$$

claim: $\hat{\theta}_{\pi}(x) = \arg \min_{\theta} \int L(\theta, \theta) \pi(\theta|x) d\theta$

pf) wts $V \hat{\theta}$.

$$\int R(\hat{\theta}_{\pi}, \theta) \pi(\theta) d\theta \leq \int R(\hat{\theta}, \theta) \pi(\theta) d\theta$$

$$\begin{aligned} \int R(\hat{\theta}_{\pi}, \theta) \pi(\theta) d\theta &= \int \int L(\hat{\theta}_{\pi}(x, \theta), \theta) \pi(\theta|x) m(x) dx d\theta \\ &\leq \int \int L(\hat{\theta}(x, \theta), \theta) \pi(\theta|x) m(x) dx d\theta \\ &= \int R(\hat{\theta}, \theta) \pi(\theta) d\theta \end{aligned}$$

an important example $\Theta \subseteq \mathbb{R}$, $L(\theta, \theta) = (\hat{\theta} - \theta)^2$

$$\hat{\theta}_{\pi}(x) = \arg \min_{\theta} \int (\hat{\theta} - \theta)^2 \pi(\theta|x) d\theta$$

$$= \underset{\theta}{\operatorname{argmin}} E((\hat{\theta} - \theta)^2 | X)$$

$$= E(\hat{\theta}(X)) \quad (\text{since } E((\hat{\theta} - \theta)^2 | X) = (E(\hat{\theta}(X)) - \theta)^2 + \operatorname{Var}(\hat{\theta}(X)))$$

e.g. $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Ber}(p)$ for $\hat{p} = \frac{1}{n} \sum_{i=1}^n X_i$
 prior $\pi = \text{Beta}(\alpha, \beta)$ $\pi(p) \propto p^{\alpha-1} (1-p)^{\beta-1}$
 $p|X_1, \dots, X_n \sim \text{Beta}\left(\sum_{i=1}^n X_i + \alpha, \sum_{i=1}^n (1-X_i) + \beta\right)$
 Bayes estimator $\hat{p} = E(p|X_1, \dots, X_n) = \frac{\sum_{i=1}^n X_i + \alpha}{n + \alpha + \beta}$

$$R(\hat{p}, p) = E_p(\hat{p} - p)^2 = \operatorname{Var}(\hat{p}) + (E(\hat{p}) - p)^2$$

$$= \frac{(n)}{n+\alpha+\beta} \frac{n(1-p)}{n} + \left(\frac{\alpha+p}{\alpha+\beta+n}\right)^2 \left(\frac{\alpha}{\alpha+\beta} - p\right)$$

Risk of MLE Expectation of prior mean

minimax estimator $\hat{\theta}_{\text{minimax}} = \underset{\hat{\theta}}{\operatorname{argmin}} \sup_{\theta \in \Theta} R(\hat{\theta}, \theta)$

Theorem If for some π , $\hat{\theta}$ satisfies

$$\sup_{\theta \in \Theta} R(\hat{\theta}, \theta) = \inf_{\tilde{\theta}} \int R(\tilde{\theta}, \theta) \pi(\theta) d\theta$$

then $\hat{\theta}$ is minimax.

Proof A

$$\begin{aligned} \sup_{\theta \in \Theta} R(\hat{\theta}, \theta) &\geq \int R(\hat{\theta}, \theta) \pi(\theta) d\theta \\ &\geq \inf_{\tilde{\theta}} \int R(\tilde{\theta}, \theta) \pi(\theta) d\theta \\ &= \sup_{\theta \in \Theta} R(\hat{\theta}, \theta) \end{aligned}$$

Corollary If $\hat{\theta} = \hat{\theta}_\pi$ for some π and $R(\hat{\theta}_\pi, \theta)$ is constant over $\theta \in \Theta$

then $\hat{\theta}$ is minimax.

$$\begin{aligned} \text{Proof: } \sup_{\theta \in \Theta} R(\hat{\theta}, \theta) &= \int R(\hat{\theta}, \theta) \pi(\theta) d\theta \\ &= \inf_{\tilde{\theta}} \int R(\tilde{\theta}, \theta) \pi(\theta) d\theta \end{aligned}$$

$\hat{\theta}$ is minimax.

e.g. $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(p)$

for $(\hat{p} - p)^2$

$$R(\hat{P}, P) = \left[\left(\frac{\alpha+\beta}{\alpha+\alpha+\beta} \right)^2 - \frac{1}{n} \left(\frac{n}{\alpha+\alpha+\beta} \right)^2 \right] P^2 + \left[\frac{1}{n} \left(\frac{n}{\alpha+\alpha+\beta} \right)^2 - \left(\frac{\alpha+\beta}{\alpha+\alpha+\beta} \right)^2 \frac{2\alpha}{\alpha+\beta} \right] P + \left(\frac{\alpha+\beta}{\alpha+\alpha+\beta} \right)^2 \left(\frac{\alpha}{\alpha+\beta} \right)^2$$

$$\left\{ \begin{array}{l} (\alpha+\beta)^2 = n \\ 2\alpha(\alpha+\beta) = n \end{array} \right. \Rightarrow \alpha = \beta = \frac{\sqrt{n}}{2}$$

$$\hat{P}_{\text{unbiased}} = \frac{\sum x_i + \frac{\sqrt{n}}{2}}{n + \sqrt{n}}$$

≈ same

Lecture 2

lower bound $\sup_{\theta \in \Theta} R(\hat{\theta}, \theta) \geq \int R(\hat{\theta}, \theta) \pi(\theta) d\theta \quad \forall \pi$

Thm. If $\exists \pi$ s.t. $\hat{\theta}$ satisfies $\sup_{\theta \in \Theta} R(\hat{\theta}, \theta) = \inf_{\hat{\theta}} \int R(\hat{\theta}, \theta) \pi(\theta) d\theta$ then $\hat{\theta}$ is minimax

Cor. if $\hat{\theta}$ is Bayes estimator w.r.t. some π_θ and $R(\theta, \cdot)$ is constant
then $\hat{\theta}$ is minimax.

e.g. $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Ber}(p)$ loss $(\hat{p} - p)^2$ $\pi(p) = \text{Beta}(\alpha, \beta)$

$$\hat{p}_\pi = E(p | X_1, \dots, X_n) = \frac{\sum X_i + \alpha}{n + \alpha + \beta}$$

$$R(\hat{p}_\pi, p) = []p^2 + []p + []$$

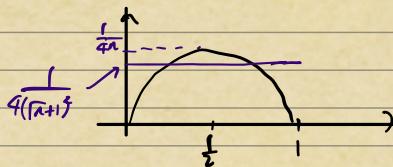
$\hookrightarrow \alpha \quad \hookrightarrow \beta \quad \Rightarrow \begin{cases} \alpha = \frac{n}{2} \\ \beta = \frac{n}{2} \end{cases}$

$$\hat{p}_{\text{minimum}} = \frac{\sum X_i + \frac{\sqrt{n}}{2}}{n + \alpha}$$

$$R(\hat{p}_{\text{minimum}}, p) = \frac{1}{4(n+1)^2}$$

$$\hat{p}_{MLE} = \bar{x} = \frac{1}{n} \sum_{i=1}^n X_i \quad R(\hat{p}_{MLE}, p) = E(\hat{p}_{MLE} - p)^2 = \text{Var}(\hat{p}_{MLE}) = \frac{p(1-p)}{n}$$

$$\sup_{p \in [0, 1]} R(\hat{p}_{MLE}, p) = \frac{1}{4n} > \frac{1}{4(n+1)^2} = R(\hat{p}_{\text{minimum}}, p)$$



e.g. $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bern}(p)$ loss $I(\hat{p}, p) = \frac{(\hat{p} - p)^2}{p(1-p)}$
just extra weight

Bayes estimator w.r.t. $\pi(p) = 1$

$$\hat{p}(x) = \arg \min_p \int \frac{(p - \hat{p})^2}{p(1-p)} \pi(p|x) dp$$

$$= \arg \min_p \int (p - \hat{p})^2 \frac{\pi(p|x)}{p(1-p)} dp$$

$$\frac{\pi(p|x)}{p(1-p)} \propto p^{\sum_{i=1}^n x_i - 1} (1-p)^{\sum_{i=1}^n (1-x_i) - 1} = \text{Beta}(\sum x_i, \sum (1-x_i))$$

$$\hat{p}(x) = \frac{\sum x_i}{\sum x_i + \sum (1-x_i)} = \frac{\sum x_i}{n} = \bar{x}$$

$$R(\hat{p}, p) = E\left(\frac{(\hat{p} - p)^2}{p(1-p)}\right) = \frac{\text{Var}(\hat{p})}{p(1-p)} = \frac{1}{n}$$

constant

\bar{X} is minimax

e.g. $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ $\mu \in \mathbb{R}$ loss $(\hat{\mu} - \mu)^2$

$$Q. \text{ Is } \bar{X} \text{ minimax? } R(\bar{X}, \mu) = E(\bar{X} - \mu)^2 = \frac{\sigma^2}{n}$$

Consider $\pi = N(0, \tau^2)$

$$\pi(\mu | x) \propto \pi(\mu) \prod_{i=1}^n p(x_i | \mu) \propto e^{-\frac{\mu^2}{2\tau^2}} e^{-\frac{1}{2} \sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma^2}}$$

$$f(\mu) = \frac{\mu}{\tau^2} + \sum \frac{(x_i - \mu)}{\sigma^2} \quad f'(\mu) = \frac{2\mu}{\tau^2} + \frac{1}{\sigma^2} \sum 2(\mu - x_i) = 0$$

$$\Leftrightarrow \frac{\mu}{\tau^2} + \frac{\sum x_i}{\sigma^2} = \frac{\sum x_i}{\tau^2}$$

$$\Rightarrow \hat{\mu} = \frac{\frac{1}{\tau^2} \sum x_i}{\frac{1}{\tau^2} + \frac{n}{\sigma^2}} = \frac{\frac{1}{\tau^2} \sum x_i}{\frac{1}{\tau^2} + \frac{n}{\sigma^2}} \bar{x}$$

$$R(\hat{\mu}, \mu) = E(\hat{\mu} - \mu)^2 = \text{Var}(\hat{\mu}) + E(\hat{\mu} - \mu)^2$$

$$= \left(\frac{\frac{1}{\tau^2}}{\frac{1}{\tau^2} + \frac{n}{\sigma^2}} \right)^2 \frac{\sigma^2}{n} + \left(\frac{\frac{1}{\tau^2}}{\frac{1}{\tau^2} + \frac{n}{\sigma^2}} \right) \mu^2$$

$$\int R(\hat{\mu}, \mu) \pi(\mu) d\mu = \left(\frac{\frac{1}{\tau^2}}{\frac{1}{\tau^2} + \frac{n}{\sigma^2}} \right)^2 \frac{\sigma^2}{n} + \left(\frac{\frac{1}{\tau^2}}{\frac{1}{\tau^2} + \frac{n}{\sigma^2}} \right)^2 \tau^2$$

$$= \frac{1}{\frac{1}{\tau^2} + \frac{n}{\sigma^2}}$$

$$\text{Recall } R(\bar{X}, \mu) = \frac{\sigma^2}{n}$$

Proof that \bar{X} is minimax:

$$\sup_{\mu} R(\bar{X}, \mu) = \frac{\sigma^2}{n} \quad \forall \hat{\mu}$$

$$\sup_{\mu} R(\hat{\mu}, \mu) \stackrel{\text{Want}}{\geq} \frac{\sigma^2}{n}$$

$$\sup_{\mu} R(\hat{\mu}, \mu) \geq \int R(\hat{\mu}, \mu) \pi(\mu) d\mu$$

$$\geq \inf_{\mu} \int R(\hat{\mu}, \mu) \pi(\mu) d\mu$$

$$= \frac{1}{\frac{1}{\tau^2} + \frac{n}{\sigma^2}} \xrightarrow{\tau \rightarrow \infty} \frac{\sigma^2}{n}$$

$$\underline{\text{Thm}} \quad \exists \{\pi_m\} \text{ s.t. } \sup_{\theta \in \Theta} R(\hat{\theta}, \theta) = \liminf_{m \rightarrow \infty} \int R(\hat{\theta}, \theta) \pi_m(\theta) d\theta$$

then $\hat{\theta}$ is minimax

Previously... \uparrow more general

Thm. If $\exists \pi$ s.t. $\hat{\theta}$ satisfies $\sup_{\theta \in \Theta} R(\hat{\theta}, \theta) = \inf_{\hat{\theta}} \int R(\hat{\theta}, \theta) \pi(\theta) d\theta$ then $\hat{\theta}$ is minimax

Admissibility

Def. $\hat{\theta}$ is inadmissible if $\exists \tilde{\theta}$ s.t.

$$R(\tilde{\theta}(n, \theta)) \leq R(\hat{\theta}(n, \theta)) \quad \forall \theta \in \Theta$$

$$R(\tilde{\theta}(n, \theta)) < R(\hat{\theta}(n, \theta)) \quad \text{for some } \theta_0 \in \Theta$$

Thm If $\hat{\theta}$ is Bayes, it is admissible

Suppose $\hat{\theta}$ is inadmissible $\exists \tilde{\theta}$

$$\textcircled{1} \quad R(\tilde{\theta}, \theta) \leq R(\hat{\theta}, \theta) \quad \forall \theta \in \Theta$$

$$\textcircled{2} \quad R(\tilde{\theta}, \theta_0) < R(\hat{\theta}, \theta_0) \quad \text{for some } \theta_0 \in \Theta \quad (\text{continuity assumption})$$

$\rightarrow \exists$ open set $\Theta_0 \ni \theta_0$ and $\varepsilon > 0$ such that

$$R(\tilde{\theta}, \theta) < R(\hat{\theta}, \theta) \quad \forall \theta \in \Theta_0$$

$$\begin{aligned} \Rightarrow \int R(\tilde{\theta}, \theta) \pi(\theta) d\theta &= \int_{\Theta_0} R(\tilde{\theta}, \theta) \pi(\theta) d\theta + \int_{\Theta_0^c} R(\tilde{\theta}, \theta) \pi(\theta) d\theta \\ &< \int_{\Theta_0} R(\hat{\theta}, \theta) \pi(\theta) d\theta + \int_{\Theta_0^c} R(\hat{\theta}, \theta) \pi(\theta) d\theta \\ &= \int R(\hat{\theta}, \theta) \pi(\theta) d\theta \quad \text{contradiction} \end{aligned}$$

Complete class theorem (Brown, 1986)

$\hat{\theta}$ is admissible $\Rightarrow \exists \pi_m$ s.t. $\hat{\theta}_{\pi_m} \rightarrow \hat{\theta}$

e.g. $X_1, \dots, X_n \stackrel{iid}{\sim} N(\theta, \sigma^2)$ $\theta \in \mathbb{R}$ $(\hat{\theta} - \theta)^2$

Q: Is \bar{X} admissible? Wald (1939) wrong prof. Colin Blyth

$$X_1, \dots, X_n \stackrel{iid}{\sim} N(\theta, I_p) \quad \theta \in \mathbb{R}^p \quad \|\hat{\theta} - \theta\|^2 = \sum_{j=1}^p (\hat{\theta}_j - \theta_j)^2$$

Q. Is \bar{X} admissible? (minimum, limit of Bayes, invariance property)

$p=1$ admissible Blyth

$p=2$ admissible Stein

$p \geq 3$ inadmissible Stein

$$\hat{\theta}_{JS} = \left(1 - \frac{p-2}{n\|\bar{X}\|}\right) \bar{X} \quad \text{James-Stein estimator}$$

Larry Brown 1971 Corresponding stochastic process recurrent or not
random walks,

Lecture 3 $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\theta, I_p)$ $L(\hat{\theta}, \theta) = \|\hat{\theta} - \theta\|$ Q. Is \bar{X} admissible or not?

Thm. $p=1$, admissible

pf) (Blyth's method)

Suppose $\hat{\theta} = \bar{X}$ is inadmissible

$$\exists \hat{\theta} \text{ s.t. } \textcircled{1} R(\hat{\theta}, \theta) \leq \frac{1}{n} \quad \forall \theta \in \mathbb{R}$$

$$\textcircled{2} R(\hat{\theta}, \theta_0) < \frac{1}{n} \quad \exists \theta_0 \in \mathbb{R}$$

$$\exists a < b \quad \varepsilon > 0 \text{ s.t. } (a, b) \ni \theta_0$$

$$R(\hat{\theta}, \theta) < \frac{1}{n} - \varepsilon \quad \forall \theta \in (a, b)$$

$$\text{Consider } \pi_m = N(0, m), \quad \int R(\hat{\theta}_{\pi_m}, \theta) \pi_m(d\theta) d\theta = \frac{1}{m+n}$$

$$\frac{1}{n} - \int R(\hat{\theta}_{\pi_m}, \theta) \pi_m(d\theta) d\theta = \frac{1}{n} - \frac{1}{m+n} = \frac{1}{n} \left(1 - \frac{n}{m+n} \right)$$

$$= \frac{1}{n} \frac{\frac{1}{m}}{\frac{m}{m+n}} \times \frac{1}{m} \quad (n \text{ fixed})$$

$$\begin{aligned} \frac{1}{n} - \int_{(a,b)} R(\hat{\theta}, \theta) \pi_m(d\theta) d\theta &= \frac{1}{n} - \int_{(a,b)} R(\hat{\theta}, \theta) \pi_m(d\theta) d\theta - \int_{(a,b)^c} R(\hat{\theta}, \theta) \pi_m(d\theta) d\theta \\ &= \int_{(a,b)} \frac{1}{n} - R(\hat{\theta}, \theta) \pi_m(d\theta) d\theta + \int_{(a,b)^c} \frac{1}{n} - R(\hat{\theta}, \theta) \pi_m(d\theta) d\theta \end{aligned}$$

$$\begin{aligned} > \sum_{(a,b)} \int \pi_m(d\theta) d\theta &= \sum_{(a,b)} P(a < N(0, m) < b) \\ &= \sum_{(a,b)} P\left(\frac{a}{\sqrt{m}} < N(0, 1) < \frac{b}{\sqrt{m}}\right) \asymp \frac{1}{\sqrt{m}} \end{aligned}$$

$$\exists m \text{ large s.t. } \int R(\hat{\theta}_{\pi_m}, \theta) \pi_m(d\theta) d\theta > \int R(\hat{\theta}_{\pi_m}, \theta) \pi_m(d\theta) d\theta$$

Contradiction $\Rightarrow \hat{\theta} = \bar{X}$ is admissible

$p=2$ \bar{X} is admissible (Blyth does not work)

proved by Stein

$p=3$ Stein \bar{X} is inadmissible

Jones - Stein estimator

$$\hat{\theta}_{JS} = \left(1 - \frac{p-2}{n\|\bar{X}\|^2}\right) \bar{X}$$

\rightarrow Stein Shrinkage

an empirical Bayes perspective (Efron)

$$x_1, \dots, x_n \stackrel{\text{IID}}{\sim} N(\theta, I_p) \quad \theta \sim N(0, \tau^2 I_p)$$

$$\hat{\theta}_{\text{Bayes}} = \mathbb{E}(\theta | x_1, \dots, x_n) = \frac{n}{n + \frac{1}{\tau^2}} \bar{x} = \left(1 - \frac{\frac{1}{\tau^2}}{n + \frac{1}{\tau^2}}\right) \bar{x}$$

Empirical Bayes framework $P(x|\theta) \pi(\theta|\tau^2)$

$$m(x|\tau^2) = \int P(x|\theta) \pi(\theta|\tau^2) d\theta$$

$$\text{estimate } \tau^2 \text{ from } m(x|\tau^2) \quad \xrightarrow{\text{"marginal likelihood"}} \\ x_1, \dots, x_n | \theta \stackrel{\text{IID}}{\sim} N(\theta, I_p)$$

$$\theta | \tau^2 \sim N(0, \tau^2 I_p)$$

$$x_1, \dots, x_n | \tau^2 \sim ?$$

$$x_i \sim N(\theta, I_p) \Leftrightarrow x_i = \theta + z_i, \quad z_i \stackrel{\text{IID}}{\sim} N(0, I_p) \quad \left. \begin{array}{l} \text{independent} \\ \theta \sim N(0, \tau^2 I_p) \Leftrightarrow \theta = \tau w \quad w \sim N(0, I_p) \end{array} \right\}$$

$$\Rightarrow x_i = \tau w + z_i, \quad i=1, \dots, n \quad (\text{dependence})$$

estimate τ^2 $\left\{ \begin{array}{l} \text{MLE} \\ \text{Method of moments} \Rightarrow \text{James-Stein} \end{array} \right.$

$$\bar{x} = \tau w + \bar{z} \sim N(0, (\tau^2 + \frac{1}{n}) I_p)$$

$$\Rightarrow \frac{\|\bar{x}\|^2}{\tau^2 + \frac{1}{n}} \sim \chi_p^2 \quad \frac{\tau^2 + \frac{1}{n}}{\|\bar{x}\|^2} \sim \text{Inv-}\chi_p^2$$

method of moment

$$\mathbb{E}\left(\frac{\tau^2 + \frac{1}{n}}{\|\bar{x}\|^2}\right) = \frac{1}{p-2}$$

$$\Rightarrow \mathbb{E}\left(\frac{p-2}{\|\bar{x}\|^2}\right) = \frac{1}{\tau^2 + \frac{1}{n}}$$

We can estimate $\frac{1}{\tau^2 + \frac{1}{n}}$ by $\frac{p-2}{\|\bar{x}\|^2}$

$$\text{Thm } p \geq 3 \quad \mathbb{E} \|\hat{\theta}_{JS} - \theta\|^2 \leq \frac{p}{n} = \mathbb{E} \|\bar{x} - \theta\|^2 \quad \forall \theta \in \mathbb{R}^p$$

Lemmas (Stein's Identity) $z \sim N(0, I_p)$

(can prove CFT)

sum
variables $\leftarrow \mathbb{E} Zg(z) = \mathbb{E} g'(z)$

$$\text{pf) } \phi(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} \quad \phi'(x) = -x \phi(x)$$

$$\mathbb{E} Zg(z) = \int x g(x) \phi(x) dx = - \int \phi'(x) g(x) dx$$

$$= - (g(x) \phi(x))|_{-\infty}^{\infty} - \int g(x) \phi(x) dx$$

$$= \mathbb{E} g'(z)$$

$$z \sim N(0, I_p) \quad g: \mathbb{R}^p \rightarrow \mathbb{R}^p \quad E(\langle z, g(z) \rangle) = E(\langle \nabla g(z), z \rangle) = \sum_{j=1}^p \frac{1}{\sqrt{n}} \frac{\partial g_j}{\partial z_j}(z)$$

Pf of Thm)

$$\begin{aligned} E_\theta \| \hat{\theta}_{JS} - \theta \|^2 &= E_\theta \| \left(1 - \frac{P-2}{n\|\bar{x}\|^2} \right) \bar{x} - \theta \|^2 \\ &= E_\theta \| \bar{x} - \theta - \frac{P-2}{n\|\bar{x}\|^2} \bar{x} \|^2 \quad \bar{x} \sim N(\theta, \frac{1}{n} I_p) \\ &\stackrel{\text{def}}{=} E_\theta \| \frac{1}{n} \bar{Z} - \frac{P-2}{n\|\bar{x}\|^2} \frac{1}{n} (\mu + z) \|^2 \quad \bar{x} = \theta + \frac{1}{n} z \quad z \sim N(0, I_p) \\ &= \frac{1}{n} E \| z - \frac{P-2}{n\|\bar{x}\|^2} (\mu + z) \|^2 \quad \mu = \sqrt{n}\theta \\ &= \frac{1}{n} \left(E \| z \|^2 + \frac{(P-2)^2}{n\|\bar{x}\|^2} - 2(P-2) \langle z, \frac{\mu + z}{\|\bar{x}\|^2} \rangle \right) \\ &= \frac{1}{n} \left(P + \frac{(P-2)^2}{n\|\bar{x}\|^2} - 2(P-2) \langle z, \frac{\mu + z}{\|\bar{x}\|^2} \rangle \right) \end{aligned}$$

Analysis of $E(\langle z, \frac{\mu + z}{\|\bar{x}\|^2} \rangle)$ $\hat{g}(z) = \begin{pmatrix} g_1(z) \\ \vdots \\ g_P(z) \end{pmatrix}$ $\hat{g}_j(z) = \frac{\mu_j + z_j}{\|\bar{x}\|^2}$

$$\frac{1}{\sqrt{n}} \hat{g}_j(z) = \frac{\|\mu + z\|^2 - 2(\mu_j + z_j)^2}{\|\bar{x}\|^4}$$

$$\begin{aligned} E \langle z, \frac{\mu + z}{\|\bar{x}\|^2} \rangle &= \sum_{j=1}^P E \frac{1}{\sqrt{n}} \hat{g}_j(z) \\ &= E \sum_{j=1}^P \frac{\|\mu + z\|^2 - 2(\mu_j + z_j)^2}{\|\bar{x}\|^4} = E \frac{P \|\mu + z\|^2 - 2 \|\mu + z\|^2}{\|\bar{x}\|^4} \\ &= E \frac{P-2}{\|\bar{x}\|^2} \end{aligned}$$

$$E_\theta \| \hat{\theta}_{JS} - \theta \|^2 = \frac{1}{n} \left(P + E \frac{(P-2)^2 - 2(P-2)^2}{\|\bar{x}\|^2} \right)$$

$$= \frac{P}{n} - \frac{1}{n} E \frac{(P-2)^2}{\|\bar{x}\|^2} \underset{n \rightarrow \infty}{\sim} 0$$

$$\begin{aligned} R(\hat{\theta}_{JS}, \theta) &= \frac{1}{n} \left(P - (P-2)^2 \frac{1}{\|\bar{x}\|^2} \right) \\ &= \frac{2}{n} \quad (\text{does not depend on } P) \end{aligned}$$



by symmetry, can show $R(\hat{\theta}_{JS}, \theta)$ is a function of $\|\theta\|$

Also show θ_w is minmax.

You can only break Cover-Rao lower bound upto measure 0 set.

Lecture 4

Shrinkage estimator $\hat{\theta}_c = c\bar{x}$

$$R(\hat{\theta}_c, \theta) = E_{\theta} \|c\bar{x} - \theta\|^2$$

$$= E_{\theta} \|c\bar{x} - c\theta\|^2 + \|c\theta - \theta\|^2 \\ = c^2 \frac{P}{n} + (c-1)^2 \|\theta\|^2 = f(c)$$

$$f'(c) = 2c \frac{P}{n} + 2(c-1) \|\theta\|^2 = 0$$

$$c^* = \frac{\|\theta\|^2}{\frac{P}{n} + \|\theta\|^2}$$

$$\hat{\theta}_{c^*} = \frac{\|\theta\|^2}{\frac{P}{n} + \|\theta\|^2} \bar{x} = \left(1 - \frac{\frac{P}{n}}{\frac{P}{n} + \|\theta\|^2}\right) \bar{x}$$

$$R(\hat{\theta}_{c^*}, \theta) = \frac{b}{a+b}^2 a + \frac{a}{a+b}^2 b \quad \begin{cases} a = \frac{P}{n} \\ b = \|\theta\|^2 \end{cases}$$

$$= \frac{b^2 a + a^2 b}{(a+b)^2} \\ = \frac{ab}{a+b} \leq \min(a, b)$$

$$= \min\left(\frac{P}{n}, \|\theta\|^2\right)$$

$$R(\hat{\theta}_{js}, \theta) = \frac{1}{n} (P - (P-\nu)^2 E \frac{1}{\|z+u\|^2}) \quad \begin{cases} z \sim N(0, I_p) \\ u \sim \mathcal{N}(0, \sigma^2 I_p) \end{cases}$$

$$\|z+u\|^2 \sim \chi^2_{p+1, \mu^2}$$

$$\text{an important fact} \stackrel{d}{=} \chi^2_{p+2N} \quad N \sim \text{pois}(\frac{\|u\|^2}{2})$$

$$E \frac{1}{\|z+u\|^2} = E \frac{1}{\chi^2_{p+2N}} = E(E(\frac{1}{\chi^2_{p+2N}} | N))$$

random

$$= E\left(\frac{1}{p+2N-2}\right)$$

Jensen

$$\geq \frac{1}{p+2N-2}$$

$$R(\hat{\theta}_{js}, \theta) \leq \frac{1}{n} (P - (P-\nu)^2 \frac{1}{p+2N-2})$$

$$= \frac{1}{n} \frac{P^2 + P\|u\|^2 - 2P - \nu^2 + 4\nu - 4}{p+2N-2}$$

$$= \frac{1}{n} \frac{P\|u\|^2 + 2P - \nu}{p+2N-2}$$

$$= \frac{2}{n} + \frac{(P-2) \|M\|^2/n}{\|M\|^2 + (P-2)} \quad M = \bar{x}\theta$$

$$= \frac{2}{n} + \frac{(P-2) \|\theta\|^2/n}{\|\theta\|^2 + (P-2)} / n$$

$$\leq \frac{2}{n} + \frac{\frac{P-2}{n} \|\theta\|^2}{1 + \frac{P-2}{n}} = \frac{2}{n} R(\hat{\theta}_0, \theta)$$

Theorem (oracle inequality)

$$R(\hat{\theta}_{JS}, \theta) \leq \inf_c R(\hat{\theta}_c, \theta^*) + \frac{2}{n}$$

dimension free

density estimation

$$x_1, \dots, x_n \stackrel{iid}{\sim} f, \quad f \in S_\alpha(R) \subseteq L^2[0, 1]$$

$$\text{loss} : L(\hat{f}, f) = \|\hat{f} - f\|^2 = \int_0^1 (\hat{f}(x) - f(x))^2 dx = \|\hat{\theta} - \theta\|^2$$

$$\text{Fourier expansion} \quad f(x) = a_0 + \sum_{j=1}^{\infty} (a_j \cos(2\pi j x) + b_j \sin(2\pi j x))$$

$$= \sum_{j=1}^{\infty} \theta_j \phi_j(x) \quad \int_0^1 \phi_j(x) \phi_k(x) dx = \begin{cases} 1 & j=k \\ 0 & j \neq k \end{cases}$$

$$f'(x) = \sum_{j=1}^{\infty} [-a_j(2\pi j) \sin(2\pi j x) + b_j(2\pi j) \cos(2\pi j x)]$$

$$\|f'\|^2 = \int_0^1 (f'(x))^2 dx = \sum (2\pi j)^2 (a_j^2 + b_j^2)$$

$$\|f''\|^2 = \sum_{j=1}^{\infty} (2\pi j)^4 (a_j^2 + b_j^2)$$

$$\text{for } \alpha > 0, \quad \|f^{(\alpha)}\|^2 = \sum_{j=1}^{\infty} (2\pi j)^{2\alpha} (a_j^2 + b_j^2)$$

$$S_\alpha(R) = \left\{ f = \sum_{j=1}^{\infty} \theta_j \phi_j : \sum_{j=1}^{\infty} j^{2\alpha} \theta_j^2 \leq R^2, \int f = 1, f \geq 0 \right\}$$

$$\inf_{\hat{f}} \sup_{f \in S_\alpha(R)} E_f \| \hat{f} - f \|^2 = (1 + o(1)) C_{\alpha, R} n^{-\frac{2\alpha}{2\alpha+1}} : \text{"non parametric rate"} \text{ cannot reach "parametric rate"}$$

\hookrightarrow finite constant

(minimax risk)

$$E_{\theta} \| \bar{x} - \theta \|^2 = \frac{P}{n} \quad n^{-1} : \text{"parametric rate"} \text{ goes to 0 faster}$$

If $P \rightarrow \infty$, $\frac{P}{n} \rightarrow 0$. In the above, due to regularity condition, it did not blow up.

Q. How to achieve minimax rate?

empirical Fourier coefficient

$$\hat{\theta}_j = \int_0^1 f \phi_j = E(\phi_j(x)) \approx \hat{\theta}_j = \frac{1}{n} \sum_i \hat{f}_i \phi_j(x_i)$$

$$\text{by CLT: } \frac{\sqrt{n}(\hat{\theta}_j - \theta_j)}{\sqrt{\text{Var}(\hat{\theta}_j)}} \sim N(0, 1)$$

$$\text{approximately } \hat{\theta}_j = \theta_j + \frac{1}{\sqrt{n}} \bar{z}_j \quad j=1, \dots, \infty$$

Infinite dimensional version of Gaussian location model
sufficient stat $\bar{X} \sim N(\theta, \frac{1}{n} I_p)$

\uparrow difference is minor

$$\bar{X} = \theta + \frac{1}{\sqrt{n}} \bar{Z} \quad \bar{Z} \sim N(0, I_p)$$

Gaussian sequence Model (Johnstone)

$$X_j = \theta_j + \frac{1}{\sqrt{n}} \bar{z}_j \quad j=1, \dots, \infty \Leftrightarrow X_j \sim N(\theta_j, \frac{1}{n})$$

$$\bar{z}_j \sim N(0, 1)$$

$$\begin{aligned} \text{equivalence} \quad & X_j \sim N(\theta_j, 1) \\ & \uparrow \\ & X \sim N(\bar{X}, 1) \\ & X \text{ (sufficient)} \end{aligned}$$

(Nussbaum, 1996)

nontrivial asymptotic equivalence between two nonparametric models (previously only for parametric models)

density estimation $\xrightarrow{\text{asymptotically equivalent}}$ Gaussian sequence model

sample $\xleftarrow{\text{cutter}} \text{sample}$

\Rightarrow go to connect problems!

non parametric regression $\xrightarrow{\text{density estimator}}$ white noise model

$$\Theta \in S_\alpha(R) = \left\{ \theta : \sum_{j=1}^{\infty} j^{-2d} \theta_j^2 \leq R^2 \right\}$$

$$L(\hat{\theta}, \theta) = \| \hat{\theta} - \theta \|^2 = \sum_{j=1}^{\infty} (\hat{\theta}_j - \theta_j)^2 \quad \text{We can pay the price of bias at the cost of variance } (\frac{1}{n} \text{ vs. } 0)$$

\Downarrow set $\sigma \rightarrow 0$ to use the regularity

$$\hat{\theta}_j = \begin{cases} X_j & j \leq k \\ 0 & j > k \end{cases} \quad \begin{matrix} \text{variance} & \hookrightarrow \\ \text{where} & \text{bias} \end{matrix}$$

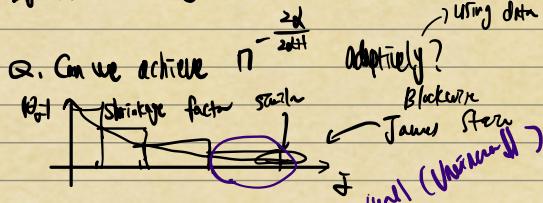
$$\begin{aligned} E_\theta \| \hat{\theta} - \theta \|^2 &= \sum_{j=1}^{\infty} E(\hat{\theta}_j - \theta_j)^2 = \sum_{j=1}^k E(X_j - \theta_j)^2 + \sum_{j=k+1}^{\infty} \theta_j^2 \\ &= \frac{k}{n} + \sum_{j=k+1}^{\infty} \theta_j^2 \leq \frac{k}{n} + R^2 k^{-2d} \end{aligned}$$

$$\text{bias}^2 = \sum_{j=k+1}^{\infty} \theta_j^2 = \sum_{j=k+1}^{\infty} j^{-2d} \hat{\theta}_j^2 \leq R^2 k^{-2d}$$

$$\text{choose } k \text{ by } \min_k \left(\frac{k}{n} + R^2 k^{-2d} \right)$$

$$\frac{k}{n} \times k^{-2d} \rightarrow k \propto n^{-\frac{1}{2d+1}} \quad (\text{effective dimension of nonparametric problem})$$

$$E_\theta \| \hat{\theta} - \theta \|^2 \leq C n^{-\frac{2d}{2d+1}} \quad \text{minimax rate optimal}$$



$$R(\hat{\theta}_{JS}, \theta) = \frac{1}{n} (P - (P-2)^2 E \frac{1}{(Z+M)^2}) \quad \begin{cases} Z \sim N(0, I_p) \\ M = \sqrt{n} \theta \end{cases}$$

large θ : risk increases

Lec 5

Leave JS when 0

Gaussian sequence model

$$X_j = \theta_j + \frac{1}{n} \sum_{i=1}^n \epsilon_i \sim N(0, 1) \quad \| \hat{\theta} - \theta \|^2$$

$$\Theta \in \mathbb{H}_d(\mathbb{R}) = \left\{ \theta : \sum_{j=1}^{2d} \theta_j^2 \leq R^2 \right\}$$

$$\inf_{\hat{\theta}} \sup_{\theta \in \mathbb{H}_d(\mathbb{R})} E \| \hat{\theta} - \theta \|^2 \asymp n^{-\frac{2d}{2d+1}} \quad \text{rate-optimal estimator} \quad \hat{\theta}_j = \begin{cases} X_j & j \leq L \\ 0 & j > L \end{cases}$$

$$\text{choose } k \propto n^{\frac{1}{2d+1}} \quad \sup_{\theta \in \mathbb{H}_d(\mathbb{R})} E_\theta \| \hat{\theta} - \theta \|^2 \lesssim n^{-\frac{2d}{2d+1}}$$

$$\textcircled{1} \quad R(\hat{\theta}_{JS}, \theta) \leq \frac{P}{n} \quad \forall \theta \in \mathbb{R}^P$$

$$\textcircled{2} \quad \text{oracle inequality} \quad R(\hat{\theta}_{JS}, \theta) \leq \inf_c R(c\bar{X}, \theta) + \frac{2}{n} \quad \forall \theta \in \mathbb{R}^P$$

where $\inf_c R(c\bar{X}, \theta) \leq \min \left(\frac{P}{n}, \frac{\| \theta \|^2}{2} \right)$

$$R(\bar{X}, \theta) \xrightarrow{\text{better than both}}$$

 $|B_j|$ 

$$\{1, \dots, n\} = \bigcup_{j=1}^m B_{\theta_j} \quad B_1 = \{1, 2, \dots\} \quad B_2 = \{4, 5, \dots, 7\} \quad B_3 = \{10, \dots, 27\}$$

$$B_{\theta_j} = \{3^{j-1} + i, 3^j\} \quad |B_{\theta_j}| = 3^j - 3^{j-1} = \frac{2}{3} 3^j$$

$$\hat{\theta} = \begin{pmatrix} \hat{\theta}_{B_1} \\ \hat{\theta}_{B_2} \\ \vdots \\ \hat{\theta}_{B_m} \end{pmatrix} \quad B_m = \{\dots, 3^m\} \quad m = \log_3 n$$

for $j=1 \dots m$

$$\hat{\theta}_{B_j} = \left(1 - \frac{|B_j| - 1}{n \| X_{B_j} \|^2} \right) X_{B_j}$$

$$\hat{\theta}_i = 0 \quad \text{if } i > n$$

$$E_\theta \| \hat{\theta} - \theta \|^2 = \sum_{j=1}^m E \| \hat{\theta}_{B_j} - \theta_{B_j} \|^2 + \sum_{i>n} \theta_i^2$$

$$\leq \sum_{j=1}^m \left[\min \left(\frac{|B_j|}{n}, \| \theta_{B_j} \|^2 \right) + \frac{2}{n} \right] + \sum_{i>n} \theta_i^2$$

$$= \sum_{j=1}^m \min \left(\frac{|B_j|}{n}, \| \theta_{B_j} \|^2 \right) + \frac{2 \cdot g_3 n}{n} + \sum_{i>n} \theta_i^2$$

$$\textcircled{3} \quad \sum_{i>n} \theta_i^2 = \sum_{i>n} i^{-2d} i^{-2d} \theta_i^2 \leq n^{-2d} \sum_{i>n} \theta_i^2 \leq R^2 n^{-2d} = O(n^{-\frac{2d}{2d+1}})$$

$$\textcircled{4} \quad \theta = O(n^{-\frac{2d}{2d+1}})$$

$$\textcircled{5} \quad \sum_{j=1}^m \min \left(\frac{|B_j|}{n}, \| \theta_{B_j} \|^2 \right) \leq \sum_{j=1}^m \min \left(\frac{3^j}{n}, \| \theta_{B_j} \|^2 \right)$$

$$\sum_{j=1}^{m^*} \frac{3^j}{n} + \sum_{j>m^*} \| \theta_{\beta_j} \|^2$$

Choose m^* to be the smallest int. s.t. $3^{m^*} \geq n^{\frac{1}{2d+1}}$

$$\Rightarrow 3^{m^*-1} \leq n^{\frac{1}{2d+1}} \leq 3^{m^*}$$

$$\sum_{j=1}^{m^*} \frac{3^j}{n} = \frac{1}{n} \sum_{j=1}^{m^*} 3^j \lesssim \frac{3^{m^*}}{n} \times n^{-\frac{2d}{2d+1}}$$

$$\sum_{j>m^*} \| \theta_{\beta_j} \|^2 \leq \sum_{j=3^{m^*-1}+1}^{\infty} \| \theta_j \|^2 \leq n^{-\frac{2d}{2d+1}}$$

$$\Rightarrow \exists c > 0 \text{ s.t. } E_\theta \| \hat{\theta} - \theta \|^2 \leq c n^{-\frac{2d}{2d+1}} \quad \forall d > 0 \quad \text{Johnstone}$$

Q. How to prove lower bound?

$$\inf_{\hat{\theta} \in \mathcal{H}_0(\mathbb{R})} \sup_{\theta \in \mathcal{H}} E_\theta \| \hat{\theta} - \theta \|^2 \geq C' n^{-\frac{2d}{2d+1}} \text{ for some } C' > 0$$

In order to prove lower bound for nonparametric estimation problem, we needs to reduce the problem to the hypothesis testing problem

two-point test (simple vs. simple)

$$H_0: X \sim P \quad H_1: X \sim Q$$

testing procedure $\phi: X \mapsto \phi(x) \in [0,1]$

$$\text{Type I error } P\phi = \mathbb{E}_{X \sim P} \phi(X)$$

$$\text{Type II error } Q(1-\phi)$$

$$\text{testing error } P\phi + Q(1-\phi)$$

$$\text{optimal testing error } \inf_{\phi} (P\phi + Q(1-\phi))$$

def. total variance distance

$$TV(P, Q) = \sup_B |P(B) - Q(B)|$$

$$\text{Theorem: } TV(P, Q) = P(p(x) > q(x)) - Q(p(x) > q(x))$$

$$= \frac{1}{2} \int |p - q| = \underbrace{\int \min(p, q)}_{\text{total variation distance}}$$

$$(F) A = \{p(x) > q(x)\}$$

$$TV(P, Q) \geq P(A_0) - Q(A_0)$$

$$VB = |P(B) - Q(B)| = |\int_B (P-Q)|$$

$$= \left| \int_{B \cap A_0} (P-Q) + \int_{B \cap A_0^c} (Q-P) \right|$$

$$\leq \max \left(\int_{B \cap A_0} (P-Q), \int_{B \cap A_0^c} (Q-P) \right)$$

$$\leq \max \left(\int_{A_0} (P-Q), \int_{A_0^c} (Q-P) \right)$$

$$= \max \left(P(A_0) - Q(A_0), Q(A_0^c) - P(A_0^c) \right)$$

$$= P(A_0) - Q(A_0)$$

Take sup over B.

$$\frac{1}{2} \int |P-Q| = \frac{1}{2} \left(\int_{P > Q} P-Q + \int_{Q > P} Q-P \right)$$

$$= \frac{1}{2} (P(A_0) - Q(A_0)) + \frac{1}{2} (Q(A_0^c) - P(A_0^c))$$

$$= P(A_0) - Q(A_0)$$

$$= TV(P, Q)$$

$$\int_{P > Q} (P-Q) = \int_{P \geq Q} Q + \int_{P \leq Q} P$$

$$= Q(A_0) + P(A_0^c)$$

$$= Q(A_0) + 1 - P(A_0)$$

$$= 1 - TV(P, Q)$$

Theorem (Neyman-Pearson lemma)

$$\inf_{\phi} (P\phi + Q(1-\phi)) = 1 - TV(P, Q) = \int \min(P, Q)$$

the optimal test ϕ is $\phi(x) = 1/P(x) < Q(x)$ ($\leftrightarrow MLE$)

Proof $\forall \phi$

$$P\phi + Q(1-\phi) = \int P\phi + \int Q(1-\phi)$$

$$= \int P\phi + Q(1-\phi)$$

$$\geq \int \min(P, Q)$$

$$\Rightarrow \inf_{\phi} (P\phi + Q(1-\phi))$$

$$\inf_{\phi} (P\phi + Q(1-\phi)) \leq P(P(X \leq g(x)) + Q(P(X > g(x)))$$

$$= 1 - TV(P, Q) = \int m_n(P, Q)$$

$$H_0: X_1, \dots, X_n \stackrel{\text{IID}}{\sim} P$$

$$H_1: X_1, \dots, X_n \stackrel{\text{IID}}{\sim} Q$$

$$\inf_{\phi} [P^n\phi + Q^n(1-\phi)]$$

$$= 1 - TV(P^n, Q^n)$$

$$= \int m_n(\pi P(X), \pi Q(X))$$

Q. What is the rate of convergence as $n \rightarrow \infty$?

Chernoff information

$$\text{Def. } C(P, Q) = - \min_{t \in (0, 1)} \log \int P^{1-t} Q^t$$

Theorem Under some mild conditions

$$\inf_{\phi} [P^n\phi + Q^n(1-\phi)] = \exp(-((1+\epsilon(1))n)C(P, Q))$$

Carrying up: high dimensional linear regression

Lecture 6

$$H_0: X_1, \dots, X_n \stackrel{\text{iid}}{\sim} P \quad \inf_{\phi} (P^n \phi + Q^n (1-\phi)) = 1 - TV(P, Q)$$

$$H_1: X_1, \dots, X_n \stackrel{\text{iid}}{\sim} Q$$

Theorem. Under mild assumptions, as $n \rightarrow \infty$

$$\inf_{\phi} (P^n \phi + Q^n (1-\phi)) = \exp(-(-1 + o(1)) n C(P, Q))$$

$$\text{where } C(P, Q) = -\min_{t \in (0, 1)} \log \int p^{1-t} q^t$$

proof: Upper bound

$$\inf_{\phi} (P^n \phi + Q^n (1-\phi)) = P^n \left(\prod_{i=1}^n q(X_i) > \prod_{i=1}^n p(X_i) \right) + Q^n \left(\prod_{i=1}^n q(X_i) \leq \prod_{i=1}^n p(X_i) \right)$$

$$P^n \left(\prod_{i=1}^n q(X_i) > \prod_{i=1}^n p(X_i) \right) = P^n \left(\left(\prod_{i=1}^n \frac{q(X_i)}{p(X_i)} \right)^t > 1 \right) \quad \forall t > 0$$

$$\leq P^n \left(\prod_{i=1}^n \frac{q(X_i)}{p(X_i)} \right)^t = \int \prod_{i=1}^n p(X_i) \left(\prod_{i=1}^n \frac{q(X_i)}{p(X_i)} \right)^t$$

$$= \left(\int p^{1-t} q^t \right)^n = \exp(n \underbrace{\log \int p^{1-t} q^t}_{o(1)})$$

$$\Rightarrow \inf_{\phi} (P^n \phi + Q^n (1-\phi)) \leq \exp(-n C(P, Q))$$

$$\text{Lower bound } \tilde{F} = \max_{t \in [0, 1]} \log \int p^{1-t} q^t$$

$$\inf_{\phi} (P^n \phi + Q^n (1-\phi)) \geq P^n \left(\prod_{i=1}^n q(X_i) > \prod_{i=1}^n p(X_i) \right)$$

$$\geq P^n \left(0 < \sum_{i=1}^n \log \frac{q(X_i)}{p(X_i)} < L \right)$$

$$= \int \left\{ 0 < \sum_{i=1}^n \log \frac{q(X_i)}{p(X_i)} < L \right\} \prod_{i=1}^n p(X_i)$$

$$= \int \frac{\prod_{i=1}^n p(X_i) \left(\prod_{i=1}^n \frac{q(X_i)}{p(X_i)} \right)^{\tilde{t}}}{\int \prod_{i=1}^n p(X_i) \left(\prod_{i=1}^n \frac{q(X_i)}{p(X_i)} \right)^{\tilde{t}}} p^{1-\tilde{t}} q^{\tilde{t}}$$

$$\left\{ 0 < \sum_{i=1}^n \log \frac{q(X_i)}{p(X_i)} < L \right\}$$

$$= \left(\int p^{1-\tilde{t}} q^{\tilde{t}} \right)^n \int \frac{1}{\prod_{i=1}^n \left(\frac{q(X_i)}{p(X_i)} \right)^{\tilde{t}}} \left(\prod_{i=1}^n \frac{p(X_i)}{q(X_i)} \right)^{\tilde{t}} \frac{q(X_i)}{\int p^{1-\tilde{t}} q^{\tilde{t}}}$$

$$\left\{ 0 < \sum_{i=1}^n \log \frac{q(X_i)}{p(X_i)} < L \right\}$$

$$\text{density } \tilde{p}$$

$$\frac{\prod_{i=1}^n \frac{p(X_i)}{q(X_i)}^{\tilde{t}}}{\int p^{1-\tilde{t}} q^{\tilde{t}}}$$

$$\geq e^{-\tilde{t}L}$$

$$\geq \left(\int p^{1-\tilde{t}} q^{\tilde{t}} \right)^n e^{-\tilde{t}L} \int \left\{ 0 < \sum_{i=1}^n \log \frac{q(X_i)}{p(X_i)} < L \right\} \prod_{i=1}^n \tilde{p}(X_i)$$

$$= \exp(-n C(P, Q)) \underline{e^{-\tilde{t}L} \tilde{p}^n \left(0 < \sum_{i=1}^n \log \frac{q(X_i)}{p(X_i)} < L \right)}$$

$$\tilde{t} = \underset{t}{\operatorname{argmax}} \log \int p^{1-t} q^t \text{ satisfies}$$

$$\frac{d}{dt} \left. \log \int p^{1-t} q^t \right|_{t=\tilde{t}} = 0$$

~~~~~ since enough  
↑ ↓  
 $\int p \frac{\frac{d}{dt} q^t}{q^t}$

$$= \frac{\int p \frac{\frac{d}{dt} q^t}{q^t} (q^t)}{\int p^{1-t} q^t}$$

$$= \int \tilde{p} \log \frac{q}{p} = 0 \quad (\text{set } t = \tilde{t})$$

$$\Rightarrow E_{\tilde{p}^n} \sum_{k=1}^n \log \frac{q(x_k)}{p(x_k)} = 0$$

$$\text{Set } L = \sqrt{\operatorname{Var}_{\tilde{p}^n} \left( \sum_{k=1}^n \log \frac{q(x_k)}{p(x_k)} \right)} = \sqrt{n \tilde{p} \left( \log \frac{q}{p} \right)^2}$$

$$\inf_p (P^n p + Q^n (1-p)) \geq \exp \left( -n C(p, q) - \tilde{t} \sqrt{n \tilde{p} \left( \log \frac{q}{p} \right)^2} \right) \tilde{p}^n \left( 0 < \sum_{k=1}^n \log \frac{q(x_k)}{p(x_k)} \middle/ \sqrt{n \tilde{p} \left( \log \frac{q}{p} \right)^2} < 1 \right)$$

**Assumption ①**  $\sqrt{\tilde{p} \left( \log \frac{q}{p} \right)^2} = o(\sqrt{n} C(p, q))$

**② CLT for**  $\sum_{k=1}^n \log \frac{q(x_k)}{p(x_k)}$  under  $\tilde{p}^n$

$$\left( \Rightarrow \tilde{p}^n \left( 0 < \frac{\sum_{k=1}^n \log \frac{q(x_k)}{p(x_k)}}{\sqrt{n \tilde{p} \left( \log \frac{q}{p} \right)^2}} < 1 \right) \rightarrow P(0 < N(0, 1) < 1) \right)$$

Under ① & ②

$$\inf_p (P^n p + Q^n (1-p)) \geq \exp \left( - (1 + o(1)) n C(p, q) \right)$$

Le Cam two point method

$$\inf_{\hat{\theta}} \sup_{\theta \in \Theta} E_{\theta} (\hat{\theta} - \theta)^2 \geq \inf_{\hat{\theta}} \sup_{\theta \in \{\theta_1, \theta_2\}} E (\hat{\theta} - \theta)^2 \quad (\theta_1, \theta_2 \in \Theta)$$

$$\geq \frac{(\theta_1 - \theta_2)^2}{4} \int_{\emptyset}^{} \min(P_{\theta_1}, P_{\theta_2})$$

**inf**  $(P_{\theta_1} q + P_{\theta_2} (1-q))$

pf)  $\inf_{\hat{\theta}} \sup_{\theta \in \{\theta_1, \theta_2\}} E_{\theta} (\hat{\theta} - \theta)^2$

$$\geq \inf_{\hat{\theta}} \frac{1}{2} E_{\theta_1} (\hat{\theta} - \theta_1)^2 + \frac{1}{2} E_{\theta_2} (\hat{\theta} - \theta_2)^2$$

$$= \inf_{\hat{\theta}} \frac{1}{2} \int (\hat{\theta} - \theta_1)^2 P_{\theta_1} + (\hat{\theta} - \theta_2)^2 P_{\theta_2}$$

$$\geq \inf_{\hat{\theta}} \frac{1}{2} \int [\hat{\theta} - \theta_1]^2 + [\hat{\theta} - \theta_2]^2 \min(P_{\theta_1}, P_{\theta_2})$$

$$(x+y)^2 \leq 2x^2 + 2y^2$$

$$\geq \frac{1}{2} \int \frac{(\theta_1 - \theta)^2}{2} m_{\theta}(\theta_1, \theta_2)$$

$$= \frac{(\theta_1 - \theta)^2}{4} \int \min(\theta_1, \theta_2)$$

$X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\theta, 1)$

$$\inf_{\hat{\theta}} \sup_{\theta \in \mathbb{R}} E_{\theta} (\hat{\theta} - \theta)^2 \geq \inf_{\hat{\theta}} \sup_{\theta \in [0, \theta_1]} E_{\theta} (\hat{\theta} - \theta)^2$$

$H_0: X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\theta_1, 1)$

$$H_1: X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\theta_2, 1) \quad \inf_{\phi} [P_{\theta_1}^n \phi + P_{\theta_2}^n (-\phi)] = \exp(-(-1 + \infty) n C(N(\theta_1), N(\theta_2, 1)))$$

$$C(N(\theta_1, 1), N(\theta_2, 1)) \propto (\theta_1 - \theta_2)^2$$

$$\text{Choose } \theta_1 = 0, \theta_2 = \frac{1}{m}$$

$$\inf_{\hat{\theta}} \sup_{\theta \in [0, \theta_1]} E_{\theta} (\hat{\theta} - \theta)^2 \geq \frac{1}{m} \inf_{\phi} (P_0^n \phi + P_{\frac{1}{m}}^n (-\phi))$$

$$\inf_{\phi} (P_0^n \phi + P_{\frac{1}{m}}^n (-\phi)) \geq P_0^n (P_0^n(x) < P_{\frac{1}{m}}^n(x))$$

$$= P_0^n \left( \frac{\prod_{i=1}^n e^{-\frac{1}{2}(X_i - \frac{1}{m})^2}}{e^{-\frac{1}{2} \sum X_i^2}} > 1 \right)$$

$$= P_0^n \left( \sum_{i=1}^n (X_i - \frac{1}{m})^2 - \bar{X}_i^2 < 0 \right)$$

$$= P_0^n \left( \frac{1}{m} \sum X_i > \frac{1}{2} \right)$$

$$= P(N(0, 1) > \frac{1}{2})$$

$$\Rightarrow \inf_{\hat{\theta}} \sup_{\theta \in \mathbb{R}} E_{\theta} (\hat{\theta} - \theta)^2 \geq \frac{P(N(0, 1) > \frac{1}{2})}{4} = \frac{1}{\pi}$$

Gaussian sequence  $X_j = \theta_j + \frac{1}{\sqrt{n}} Z_j \quad j=1, \dots, n$   
 $Z_j \sim N(0, 1)$

$$X_j \sim N(\theta_j, \frac{1}{n})$$

$$\theta \in \Theta_R = \left\{ \theta: \sum_j j^{2d} \theta_j^2 \leq R^2 \right\}$$

$$\text{Thm. } \inf_{\hat{\theta}} \sup_{\theta \in \Theta_R} E_{\theta} \| \hat{\theta} - \theta \|^2 \geq C n^{-\frac{2d}{2d+1}}$$

$$\text{Pf) } \Theta \subset \Theta_R$$

$$\Theta_0 = \left\{ \theta: \theta_j = \begin{cases} 0 & j \leq k \\ \frac{1}{j} & j > k \end{cases} \right\} \quad |\Theta_0| = 2^k \quad \text{hardest possible problem: "lower bound"}$$

$$\forall \theta \in \Theta_0 \quad \sum_j j^{2d} \theta_j^2 = \sum_{j=k}^k j^{2d} \theta_k^2 \leq \frac{1}{n} \sum_{j=1}^k \theta_k^2 \leq \frac{k^{2d+1}}{n} \leq R^2$$

↑ integration      ↓ word

Choose  $k \propto n^{\frac{1}{2d+1}}$  Choose  $\pi$  to be uniform prior on  $\Theta_0$

$$\inf_{\hat{\theta}} \sup_{\theta \in \Theta_0} E_\theta \|\hat{\theta} - \theta\|^2 \geq \inf_{\hat{\theta}} \sup_{\theta \in \Theta} E_\theta \|\hat{\theta} - \theta\|^2$$

$$\geq \inf_{\hat{\theta}} \int \pi(\theta) E_\theta \|\hat{\theta} - \theta\|^2 d\theta \quad n^{\frac{1}{2d+1}}$$

$$\geq \inf_{\hat{\theta}} \int \pi(\theta) E_\theta \sum_{j=1}^k (\hat{\theta}_j - \theta_j)^2 d\theta \sim \frac{k}{n}$$

$$= \inf_{\hat{\theta}} \sum_{j=1}^k \text{ave}_{\theta \in \Theta} E_\theta (\hat{\theta}_j - \theta_j)^2 \quad \theta = (\theta_j, \theta_{-j})$$

$$= \inf_{\hat{\theta}} \sum_{j=1}^k \text{ave}_{\theta_j} \text{ave}_{\theta_{-j}} E_{(\theta_j, \theta_{-j})} (\hat{\theta}_j - \theta_j)^2$$

$$= \inf_{\hat{\theta}} \sum_{j=1}^k \text{ave}_{\theta_j} \left[ \frac{1}{2} E_{(\theta_j = 0, \theta_{-j})} (\hat{\theta}_j - 0)^2 + \frac{1}{2} E_{(\theta_j = \frac{1}{n}, \theta_{-j})} (\hat{\theta}_j - \frac{1}{n})^2 \right]$$

$$\geq \sum_{j=1}^k \text{ave}_{\theta_j} \inf_{\hat{\theta}_j} \left[ \frac{1}{2} E_{(\theta_j = 0, \theta_{-j})} (\hat{\theta}_j - 0)^2 + \frac{1}{2} E_{(\theta_j = \frac{1}{n}, \theta_{-j})} (\hat{\theta}_j - \frac{1}{n})^2 \right]$$

$$\geq \sum_{j=1}^k \text{ave}_{\theta_j} C \frac{1}{n} = C \frac{k}{n} \times n^{-\frac{2}{2d+1}} \quad \text{previously derived}$$

## Lecture 9

### ① empirical Bayes

$$X|\theta \sim p(x|\theta) \quad \theta|\eta \sim \pi(\theta|\eta)$$

$$x|\eta \sim p(x|\eta) = \int p(x|\theta)\pi(\theta|\eta) d\theta$$

"marginal"

$$\text{MLE } \hat{\eta} = \underset{\eta}{\operatorname{argmax}} \ p(x|\eta)$$

### ② model with latent variables

$$X|\theta, \eta \sim p(x|\theta, \eta) \quad \theta \sim \pi(\theta)$$

↑  
model parameter  
latent variable

(like missing data, has prior)

$$x|\eta \sim p(x|\eta) = \int p(x|\theta, \eta) \pi(\theta) d\theta$$

$$\text{MLE } \hat{\eta} = \underset{\eta}{\operatorname{argmax}} \ p(x|\eta)$$

*Bayesian frequentist*

③ In general,  $x|\theta, \alpha \sim p(x|\theta, \alpha)$

$$\theta|\beta \sim \pi(\theta|\beta)$$

$$p(x|\alpha, \beta) = \int p(x|\theta, \alpha) \pi(\theta|\beta) d\theta$$

### example. Gaussian mixture model

$$x_i | z_i, \mu, \sigma^2 \sim N(\mu_{z_i}, \sigma^2 I_p) \quad z_i: \text{latent variable}$$

clustering label  $z_i \in \{1 \dots k\}$

special case  $k=2$   $x_i \sim \begin{cases} N(\mu_1, \sigma^2 I_p) & z_i=1 \\ N(\mu_2, \sigma^2 I_p) & z_i=2 \end{cases}$

$$\mu_1, \mu_2, \dots, \mu_n \in \mathbb{R}^p$$

$$z_i \stackrel{iid}{\sim} \text{Unif}(\{1 \dots k\})$$

$$p(x_1, \dots, x_n | \mu, \sigma^2) = \frac{1}{K} \sum_{i=1}^K \prod_{j=1}^n p(x_j | z_i, \mu, \sigma^2)$$

$$\text{MMLE } \max_{\mu, \sigma^2} \frac{1}{K} \sum_{i=1}^K \prod_{j=1}^n p(x_j | z_i, \mu, \sigma^2)$$

### • EM algorithm

Some basics Kullback-Leibler divergence (relative entropy)

$$D(P||Q) = \int p \log \frac{p}{q} d\eta \quad (P \ll Q, \text{ if not } \infty) \quad P, Q \ll \pi$$

Lemma  $D(P||Q) \geq 0$

Why MLE works?

$$\begin{aligned}
 X_1, \dots, X_n &\stackrel{\text{iid}}{\sim} P_{\theta^*} \quad \text{MLE} \quad \hat{\theta} = \underset{\theta}{\operatorname{argmax}} \frac{1}{n} \sum_{i=1}^n \log P_{\theta}(X_i) \\
 &= \underset{\theta}{\operatorname{argmax}} \frac{1}{n} \sum_i \log \frac{P_{\theta}(X_i)}{P_{\theta^*}(X_i)} \\
 &= \underset{\theta}{\operatorname{argmax}} \left( \frac{1}{n} \sum_i \log \frac{P_{\theta}(X_i)}{P_{\theta^*}(X_i)} - \frac{1}{n} \sum_i \log P_{\theta^*}(X_i) \right) \\
 &= \underset{\theta}{\operatorname{argmin}} \frac{1}{n} \sum_i \log \frac{P_{\theta}(X_i)}{P_{\theta^*}(X_i)}
 \end{aligned}$$

$$\text{LLN: } \frac{1}{n} \sum_i \log \frac{P_{\theta}(X_i)}{P_{\theta^*}(X_i)} \rightarrow \mathbb{E}_{\theta^*} \log \frac{P_{\theta^*}}{P_{\theta}} = \int P_{\theta^*} \log \frac{P_{\theta^*}}{P_{\theta}} = D(P_{\theta^*} || P_{\theta})$$

$$\theta^* = \underset{\theta}{\operatorname{argmin}} D(P_{\theta^*} || P_{\theta})$$

EM algorithm

$$\max_{\mu, \sigma^2} \log \sum_z P(x|z, \mu, \sigma^2)$$

$$\max_{\eta} \log \int p(x|\theta, \eta) \pi(\theta) d\theta$$

find a lower bound:  $\log \int p(x|\theta, \eta) \pi(\theta) d\theta$

$$= \log \int p(x|\theta, \eta) \frac{\pi(\theta)}{g(\theta)} g(\theta) d\theta$$

$$\text{Jensen} \geq \int g(\theta) \log \left( p(x|\theta, \eta) \frac{\pi(\theta)}{g(\theta)} \right) d\theta \quad (\text{replace log sum with sum log})$$

$$= \int g(\theta) \log p(x|\theta, \eta) d\theta - \int g(\theta) \log \frac{g(\theta)}{\pi(\theta)} d\theta$$

$$= \int g(\theta) \log p(x|\theta, \eta) d\theta - D(g \parallel \pi)$$

$$= F(\eta, g)$$

$$\text{EM algorithm} \quad \eta^{t+1} = \underset{\eta}{\operatorname{argmax}} F(\eta, g^t) \quad \text{M-step}$$

$$g^{t+1} = \underset{g}{\operatorname{argmax}} F(\eta^{t+1}, g) \quad \text{E-step}$$

$$\text{Theorem (Neal & Hinton)}: \max_{\eta} \log \int p(x|\theta, \eta) \pi(\theta) d\theta$$

$$= \max_{\eta, g} F(\eta, g)$$

$$\begin{aligned}
 \text{pf)} \quad F(\eta, g) &= \int q(\theta) \log p(x|\theta, \eta) d\theta - D(g||\pi) \\
 &= \int q(\theta) \log \frac{p(x|\theta, \eta) \pi(\theta)}{q(\theta)} \\
 &= \int q(\theta) \log \frac{\frac{p(x|\theta, \eta) \pi(\theta)}{\int p(x|\theta, \eta) \pi(\theta) d\theta}}{q(\theta)} d\theta \\
 &\quad + \underline{\int q(\theta) \log \int p(x|\theta, \eta) \pi(\theta) d\theta} \\
 &\quad \log \int p(x|\theta, \eta) \pi(\theta) d\theta
 \end{aligned}$$

$$= -D(g || \frac{p(x|\theta, \eta) \pi(\theta)}{\int p(x|\theta, \eta) \pi(\theta)}) + \log \int p(x|\theta, \eta) \pi(\theta) d\theta. \quad \text{Done}$$

$$g^t(x) = \frac{p(x|\theta, \eta) \pi(\theta)}{\int p(x|\theta, \eta) \pi(\theta)} : \text{posterior expectation}$$

$$\text{EM: } \eta^{th} = \arg \max_{\eta} F(\eta, g^t)$$

$$g^{th} = \arg \max_g F(\eta^{th}, g)$$

$$\begin{cases} \eta^{th} = \arg \max_{\eta} \int g^t(\theta) \log p(x|\theta, \eta) d\theta & \text{M-step} \\ g^{th} = \frac{p(x|\theta, \eta^{th}) \pi(\theta)}{\int p(x|\theta, \eta^{th}) \pi(\theta) d\theta} & \text{E-step} \end{cases}$$

Gaussian mixture

$$\begin{cases} x_i | z_i, M \sim N(\mu_{z_i}, \Sigma_p) \\ z_i \sim \text{Unif}([1, \dots, k]) \end{cases} \quad \text{find: } \max_M \sum_z p(x|z, M)$$

$$\begin{aligned}
 \text{E-step} \quad g(z_1 \dots z_n) &\propto \prod_{i=1}^n p(x_i | z_i, M) \quad \checkmark \\
 \rightarrow g(z_1 \dots z_n) &= \prod g(z_i)
 \end{aligned}$$

$$\begin{aligned}
 g(z_i = j) &= \frac{p(z_i=j, \mu_j)}{\sum_{i=1}^k p(z_i=i, \mu_i)} \\
 p(z_i=j, \mu_j) &\sim N(\mu_{z_i}, \Sigma_p) \\
 g_{ij} &= \frac{p(x_i | z_i=j, \mu_j)}{p(x_i | z_i=1, \mu_1) + \dots + p(x_i | z_i=k, \mu_k)} \\
 &\sim N(\mu_j, \Sigma_p)
 \end{aligned}$$

$$\begin{aligned}
 e^{-\frac{1}{2} \|x_i - \mu_j\|^2} \\
 = \frac{e^{-\frac{1}{2} \|x_i - \mu_1\|^2} + \dots + e^{-\frac{1}{2} \|x_i - \mu_k\|^2}}{e^{-\frac{1}{2} \|x_i - \mu_1\|^2} + \dots + e^{-\frac{1}{2} \|x_i - \mu_k\|^2}}
 \end{aligned}$$

M-step

$$\log \prod_{i=1}^n p(x_i | z_i, \mu)$$

$$= \sum_{i=1}^n \log p(x_i | z_i, \mu)$$

objective  $\mathbb{E}_{z \sim q} \left( \sum_{i=1}^n \log p(x_i | z_i, \mu) \right)$

$$= \sum_{i=1}^n \mathbb{E}_{z \sim q} \log p(x_i | z_i, \mu)$$

$$= \sum_{i=1}^n \sum_{j=1}^k g_{ij} \log p(x_i | z_i = j, \mu)$$

$$= \sum_i \sum_j g_{ij} \frac{\log p(x_i | z_i = j, \mu)}{e^{-\frac{1}{2} \|x_i - \mu_j\|^2}}$$

$$= \sum_{i=1}^n \sum_{j=1}^k g_{ij} \left( -\frac{1}{2} \|x_i - \mu_j\|^2 \right) + C$$

$$- \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^k g_{ij} \|x_i - \mu_j\|^2$$

M-step  $\Leftrightarrow \min_{\mu_1, \dots, \mu_k} \sum_{i=1}^n \sum_{j=1}^k g_{ij} \|x_i - \mu_j\|^2$  ↗ weighted least square problem

$$= \sum_{j=1}^k \min_{\mu_j} \sum_{i=1}^n g_{ij} \|x_i - \mu_j\|^2$$

$$\Rightarrow \mu_j = \frac{\sum_{i=1}^n g_{ij} x_i}{\sum_{i=1}^n g_{ij}}$$

EM-algorithm for clustering

$$g_{ij}^{t+1} = \frac{e^{-\frac{1}{2} \|x_i - \mu_j^t\|^2}}{\sum_{j=1}^k e^{-\frac{1}{2} \|x_i - \mu_j^t\|^2}}$$

$$\mu_j^{t+1} = \frac{\sum_{i=1}^n g_{ij}^{t+1} x_i}{\sum_{i=1}^n g_{ij}^{t+1}}$$

\* if  $g$  is concentrated on some  $j$ 'th cluster \* M-convex optimization

$\mu_j \approx$  sample mean

so we may run several times

post-processing  $\hat{z}_i = \arg \max_{j \in \{1, \dots, k\}} g_{ij}^T$ .

## Lecture 8

EM algorithm (Dempster, Laird, & Rubin)

$$x | \theta, \eta \sim p(x|\theta, \eta) \quad \theta \sim \pi \quad \hat{\eta} = \arg \max_{\eta} \log \int p(x|\theta, \eta) \pi(\theta) d\theta$$

$$\max_{\eta} \log \int p(x|\theta, \eta) \pi(\theta) d\theta = \max_{\eta, \theta} F(\eta, \theta)$$

$$F(\eta, \theta) = \int q(\theta) \log p(x|\theta, \eta) d\theta - D(q||\pi)$$

Gaussian mixture model

$$x_i | z_i, \mu \sim N(\mu_{z_i}, \Sigma_{z_i}) \quad z_i \sim U_{\text{ref}}(k)$$

$$\begin{cases} q_{ij}^{t+1} = \frac{e^{-\frac{1}{2} \|x_i - \mu_j^t\|^2}}{\sum_{j=1}^k e^{-\frac{1}{2} \|x_i - \mu_j^t\|^2}} \\ \mu_j^{t+1} = \frac{\sum_{i=1}^n q_{ij}^{t+1} x_i}{\sum_{i=1}^n q_{ij}^{t+1}} \end{cases}$$

$$x | \theta, \eta \sim p(x|\theta, \eta) \quad \max_{\theta, \eta} \log p(x|\theta, \eta)$$

$$\eta^{t+1} = \arg \max_{\eta} \log p(x|\theta^t, \eta)$$

$$\theta^{t+1} = \arg \max_{\theta} \log p(x|\theta, \eta^{t+1})$$

Gaussian mixture Model

$$\begin{aligned} \sum_{i=1}^n \log p(x_i | \mu_{z_i}) &= -\frac{1}{2} \sum_{i=1}^n \|x_i - \mu_{z_i}\|^2 + C \\ &= -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^k 1_{\{z_i=j\}} \|x_i - \mu_j\|^2 + C \\ &= -\frac{1}{2} \sum_{j=1}^k \sum_{i=1}^n 1_{\{z_i=j\}} \|x_i - \mu_j\|^2 + C \end{aligned}$$

$$\max_z \sum_{i=1}^n \log p(x_i | \mu_{z_i}) \Leftrightarrow \min_{z_i} \sum_{j=1}^k 1_{\{z_i=j\}} \|x_i - \mu_j\|^2 \quad \text{for } i=1 \dots n$$

$$z_i = \arg \min_{j \in \{1 \dots k\}} \|x_i - \mu_j\|^2$$

$$\max_{\mu} \sum_{i=1}^n \log p(x_i | \mu_{z_i}) \Leftrightarrow \min_{\mu_j} \sum_{i=1}^n 1_{\{z_i=j\}} \|x_i - \mu_j\|^2 \quad \text{for } j=1 \dots k$$

$$\mu_j = \frac{\sum_{i=1}^n 1_{\{z_i=j\}} x_i}{\sum_{i=1}^n 1_{\{z_i=j\}}}$$

$$\Rightarrow \begin{cases} z_i^{t+1} = \arg \min_{j \in \{1 \dots k\}} \|x_i - \mu_j^t\|^2 \\ \mu_j^{t+1} = \frac{\sum_{i=1}^n 1_{\{z_i^{t+1}=j\}} x_i}{\sum_{i=1}^n 1_{\{z_i^{t+1}=j\}}} \end{cases}$$

↑ thresholded version of E-step

k-means algorithm

"M-step"

"EM algorithm is a soft version of k-means algorithm"

k-means: purely frequentist

EM: hybrid philosophy

Which is better?

|         | clustering | parameter estimation |
|---------|------------|----------------------|
| EM      | ✓          | ✓                    |
| k-means | ✓          | inconsistency        |

but can do "likelihood" likelihood test

As long as they don't get stuck at local optimum, they give us maximum a posteriori cluster

Special case  $k=2$ ,  $\mu_1 = -\mu_2 = \mu$

$$x_i | z_i, \mu \sim \begin{cases} N(\mu, I_p) & z_i=1 \\ N(\mu, I_p) & z_i=2 \end{cases}$$

you use other cluster's data  
unless  $\mu \rightarrow \infty$ , error = constant

$$l(\hat{z}, z) = \frac{1}{n} \sum_{i=1}^n 1_{\{\hat{z}_i \neq z_i\}} + \frac{n}{2} \sum_{i=1}^n \log \left( \frac{\|x_i - \mu_1\|^2}{\|x_i - \mu_2\|^2} \right)$$

all equivalent permutations

$\inf_{\hat{z}} \sup_{z \in \{1,2\}^n} l(\hat{z}, z) = \exp \left( -(1+o(1)) \frac{\|\mu\|^2}{2} \right)$

$$P(X | \theta, \eta)$$

Neyman - Scott problem

$$\begin{cases} X_i \sim N(\mu_i, \sigma^2) \\ Y_i \sim N(\mu_i, \sigma^2) \quad i=1 \dots n \end{cases}$$

goal : estimate  $\sigma^2$

method ① (like k-means)  $\max_{\mu_1, \dots, \mu_n} P(X, Y | \mu_1, \dots, \mu_n, \sigma^2)$

$$\hat{\sigma}_{MLE}^2 \rightarrow c \quad (\text{homework})$$

method ② REML  $X_i - Y_i \sim N(0, 2\sigma^2) \quad i=1 \dots n$

$$\hat{\sigma}^2 = \frac{1}{2n} \sum_{i=1}^n (X_i - Y_i)^2 \xrightarrow{LLN} \sigma^2$$

homework: equivalent to putting prior on  $\mu_1, \dots, \mu_n$  (EM-algorithm)

Lessini put prior on local parameter (if you're a frequentist)

EM algorithm

$$X | \theta, \eta \sim P(X | \theta, \eta) \quad \text{hybrid}$$

# Bayes frequentist

More generally:

$$x|\theta, \eta \sim p(x|\theta, \eta) \quad \theta \sim \pi_c \quad (\theta|c \sim \pi(\cdot|c))$$

$\hookrightarrow$  a hyper parameter

$$\max_{\eta, c} \log \int p(x|\theta, \eta) \pi_c(\theta) d\theta$$

$\hookrightarrow$  empirical Bayes principle

$$\begin{aligned} \log \int p(x|\theta, \eta) \pi_c(\theta) d\theta &\geq \int g(\theta) \log p(x|\theta, \eta) d\theta - D(g||\pi_c) \\ &= F(\eta, c, g) \end{aligned}$$

$$\left\{ \begin{array}{l} \text{E-Step} \quad q^{t+1} = \underset{g}{\operatorname{argmax}} F(\eta^t, c^t, g) \\ \text{M-Step} \quad (\eta^{t+1}, c^{t+1}) = \underset{\eta, c}{\operatorname{argmax}} F(\eta, c, q^{t+1}) \end{array} \right.$$

① only  $\eta$   $x|\eta \sim p(x|\eta)$

$$\max_{\eta} \log p(x|\eta) \quad \text{MLE}$$

② only  $\theta$   $x|\theta \sim p(x|\theta) \quad \theta \sim \pi$

$$\max_g \left[ \int g(\theta) \log p(x|\theta) d\theta - D(g||\pi) \right]$$

some info abt data      penalty

$$\Leftrightarrow q(\theta) = \frac{p(x|\theta) \pi(\theta)}{\int p(x|\theta) \pi(\theta) d\theta} \quad \text{Bayes}$$

③ only  $\theta$  (unknown  $c$ ) empirical Bayes

Variational Bayes  $x|\theta \sim p(x|\theta) \quad \theta \sim \pi(\theta) \quad \text{Computing posterior!}$

$$\max_{g \in S} \left[ \int g(\theta) \log p(x|\theta) d\theta - D(g||\pi) \right]$$

$\hookrightarrow$  Variational class

$$\Leftrightarrow \min_{g \in S} D(g||\pi(\cdot|x))$$

mean-field class

$$S = \left\{ q(\theta) = \prod_j q_j(\theta_j) : q_j \text{ is any distribution} \right\}$$

$$\text{e.g. } x|\theta, \eta \sim p(x|\theta, \eta) \quad \theta, \eta \sim \pi$$

$$S = \left\{ q(\theta, \eta) = f(\theta) g(\eta) : f, g \right\}$$

$$V\theta : \min_{g \in S} D(g||\pi(\cdot|x))$$

$$\Leftrightarrow \max_{f,g} \left[ \int f(\theta) g(\eta) \log p(x|\theta, \eta) - D(f \otimes g \| \pi) \right]$$

$$\Leftrightarrow \min_{f,g} \left[ \int f(\theta) g(\eta) \log \frac{1}{p(x|\theta, \eta)} d\theta d\eta + D(f \otimes g \| \pi) \right]$$

$$= \min_{f,g} F(f, g)$$

CAVI (coordinate ascent variational inference)

$$\begin{cases} f^{t+1} = \arg \min_f F(f, g^t) \\ g^{t+1} = \arg \min_g F(f^{t+1}, g) \end{cases}$$

$$\begin{aligned} F(f, g) &= \iint f(\theta) g(\eta) \log \frac{1}{p(x|\theta, \eta)} d\theta d\eta + \iint f(\theta) g(\eta) \log \frac{f(\theta) \pi(\eta)}{\pi(\theta, \eta)} d\theta d\eta \\ &= \iint f(\theta) g(\eta) \log \frac{1}{p(x|\theta, \eta)} d\theta d\eta + \int f(\theta) \log f(\theta) d\theta + \underbrace{\int g(\eta) \log g(\eta) d\eta}_{\text{constant}} \\ &\quad - \int f(\theta) g(\eta) \log \pi(\theta, \eta) d\theta d\eta \end{aligned}$$

$$= \iint f(\theta) \left[ g(\eta) \log \frac{1}{p(x|\theta, \eta) \pi(\theta, \eta)} d\eta \right] d\theta + C + \int f(\theta) \log f(\theta) d\theta$$

$$= \int f(\theta) \log e^{\int g(\eta) \log \frac{1}{p(x|\theta, \eta) \pi(\theta, \eta)} d\eta} d\theta + \int f(\theta) \log f(\theta) d\theta + C$$

$$= \int f(\theta) \log \frac{f(\theta)}{e^{\int g(\eta) \log p(x|\theta, \eta) \pi(\theta, \eta) d\eta}} d\theta + C$$

Solution is  $f(\theta) \propto e^{\int g(\eta) \log p(x|\theta, \eta) \pi(\theta, \eta) d\eta}$

## Lecture 9

Variational Bayes

$$\min_{q \in S} \int q(\theta) \log \frac{1}{P(x|\theta)} d\theta + D(q||\pi)$$

$$\Leftrightarrow \min_{q \in S} D(q||\pi(\cdot|x))$$

$$S = \left\{ q(\theta, \eta) = f(\theta) g(\eta); f, g \right\}$$

(Jensen)

$$F(f, g) = \int f(\theta) g(\eta) \log \frac{1}{P(x|\theta)} d\theta d\eta + D(f \otimes g || \pi)$$

$$= \int f(\theta) g(\eta) \log \frac{1}{P(x|\theta, \eta)} d\theta d\eta + \iint f(\theta) g(\eta) \log \frac{f(\theta) g(\eta)}{\pi(\theta, \eta)} d\theta d\eta$$

$$= \iint f(\theta) g(\eta) \log \frac{1}{P(x|\theta, \eta)} d\theta d\eta + \int f(\theta) \log f(\theta) d\theta + \underbrace{\int g(\eta) \log g(\eta) d\eta - \int f(\theta) g(\eta) \log \pi(\theta, \eta) d\theta d\eta}_{\text{constant}}$$

$$= \int f(\theta) \log f(\theta) d\theta + \iint f(\theta) g(\eta) \log \frac{1}{P(x|\theta, \eta) \pi(\theta, \eta)} d\eta d\theta + \text{const.}$$

$$= \int f(\theta) \log \frac{f(\theta)}{\exp \int g(\eta) P(x|\theta, \eta) \pi(\theta, \eta) d\eta} d\theta + \text{const}$$

minimizer is  $f(\theta) \propto \exp \left( \int g(\eta) P(x|\theta, \eta) \pi(\theta, \eta) d\eta \right)$

$$\text{CAVI} \quad \left\{ \begin{array}{l} f^{t+1}(\theta) \propto \exp \left( \mathbb{E}_{\eta \sim g^t} [\log P(x|\theta, \eta) \pi(\theta, \eta)] \right) \\ g^{t+1}(\eta) \propto \exp \left( \mathbb{E}_{\theta \sim f^{t+1}} [\log P(x|\theta, \eta) \pi(\theta, \eta)] \right) \end{array} \right.$$

$\xrightarrow{\text{switch to E-step}}$

$$\Leftrightarrow f^{t+1}(\theta) \propto \exp \left( \mathbb{E}_{\eta \sim g^t} [\log \pi(\theta, \eta | x)] \right)$$

$$g^{t+1}(\eta) \propto \exp \left( \mathbb{E}_{\theta \sim f^{t+1}} [\log \pi(\theta, \eta | x)] \right)$$

Gaussian mixture

$$x_i | z_i, M \sim N(M_{z_i}, I_p)$$

$$z_i \sim \text{Unif}(\{1, \dots, k\})$$

$$M_j \sim N(0, \tau^2 I_p) \quad j=1, \dots, k$$

$$\pi(M, z | x) \propto \prod_{i=1}^n P(x_i | M_{z_i}) \prod_{j=1}^k P(M_j) \propto \exp \left( -\frac{1}{2} \sum_{i=1}^n \|x_i - M_{z_i}\|^2 - \frac{1}{2\tau^2} \sum_{j=1}^k \|M_j\|^2 \right)$$

$$\text{find } q(\mu, z) = f(\mu) g(z)$$

$$\log \pi(\mu, \sigma^2 | x) = -\frac{1}{2} \sum_i \sum_j \mathbb{1}_{\{z_i=j\}} \|x_i - \mu_j\|^2 - \frac{1}{2\sigma^2} \sum_{j=1}^k \|\mu_j\|^2 + C$$

$$g(z_1, \dots, z_n) \propto \exp \left( -\frac{1}{2} \sum_{j=1}^k \sum_{i=1}^n \mathbb{1}_{\{z_i=j\}} E_{M_{ij}} \|x_i - \mu_j\|^2 - \frac{1}{2\sigma^2} \sum_{j=1}^k (E_{M_{jj}} \|\mu_j\|^2) \right)$$

$$\propto \prod_{i=1}^n \exp \left( -\frac{1}{2} \sum_{j=1}^k \mathbb{1}_{\{z_i=j\}} E_{M_{ij}} \|x_i - \mu_j\|^2 \right)$$

$$g(z_1, \dots, z_n) = \prod_{i=1}^n g(z_i)$$

$$g_{i,j} = g(z_i=j) \propto \exp \left( -\frac{1}{2} E_{M_{jj}} \|x_i - \mu_j\|^2 \right) \quad \text{"Like E-step"}$$

$$f(\mu_1, \dots, \mu_k) \propto \exp \left( -\frac{1}{2} \sum_{j=1}^k \sum_{i=1}^n \mathbb{1}_{\{z_i=j\}} \|x_i - \mu_j\|^2 - \frac{1}{2\sigma^2} \sum_{j=1}^k \|\mu_j\|^2 \right)$$

$$\propto \prod_{j=1}^k \exp \left( -\frac{1}{2} \sum_{i=1}^n \mathbb{1}_{\{z_i=j\}} \|x_i - \mu_j\|^2 - \frac{1}{2\sigma^2} \|\mu_j\|^2 \right)$$

$$f(\mu_1, \dots, \mu_k) = \prod_{j=1}^k f(\mu_j)$$

$$f(\mu_j) \propto \exp \left( -\frac{1}{2} \sum_{i=1}^n \mathbb{1}_{\{z_i=j\}} \|x_i - \mu_j\|^2 - \frac{1}{2\sigma^2} \|\mu_j\|^2 \right)$$

$$\sum_{i=1}^n \mathbb{1}_{\{z_i=j\}} \|x_i - \mu_j\|^2 + \frac{1}{2\sigma^2} \|\mu_j\|^2$$

$$= \underbrace{\|\mu_j\|}_{} - \underbrace{\|\mu_j\|^2}_{} + C$$

$$\beta_j = \frac{\sum_{i=1}^n \mathbb{1}_{\{z_i=j\}} x_i}{\sum_{i=1}^n \mathbb{1}_{\{z_i=j\}} + \frac{1}{2\sigma^2}}$$

$$K_j^{-2} = \sum_{i=1}^n \mathbb{1}_{\{z_i=j\}} + \frac{1}{2\sigma^2}$$

$$f(\mu_j) = N(\beta_j, K_j^{-2} I_p)$$

$$\begin{cases} \beta_j = \frac{\sum_{i=1}^n \mathbb{1}_{\{z_i=j\}} x_i}{\sum_{i=1}^n \mathbb{1}_{\{z_i=j\}} + \frac{1}{2\sigma^2}} \\ K_j^{-2} = \frac{1}{\sum_{i=1}^n \mathbb{1}_{\{z_i=j\}} + \frac{1}{2\sigma^2}} \end{cases}$$

$$E_{M_{ij}} \|x_i - \mu_j\|^2 = E_{\mu_j \sim N(\beta_j, K_j^{-2} I_p)} \|x_i - \mu_j\|^2$$

$$= \|x_i - \beta_j\|^2 + K_j^{-2} P$$

CAVI for Gaussian mixture

$$g(z_i=j) = \int_{\mathbb{R}^p} d\mu e^{-\frac{1}{2} \|x_i - \beta_j\|^2 - \frac{1}{2} K_j^{-2} P}$$

$$(g_{ij}^{t+1}) = \frac{e^{-\frac{1}{2} \|x_i - \beta_j\|^2 - \frac{1}{2} K_j^{-2} P}}{\sum_{l=1}^k e^{-\frac{1}{2} \|x_i - \beta_l\|^2 - \frac{1}{2} K_l^{-2} P}}$$

$$\begin{cases} \beta_j^{t+1} = \frac{\sum_{i=1}^n \mathbb{1}_{\{z_i=j\}} x_i}{\sum_{i=1}^n \mathbb{1}_{\{z_i=j\}} + \frac{1}{2\sigma^2}} \\ K_j^{t+2} = \frac{1}{\sum_{i=1}^n \mathbb{1}_{\{z_i=j\}} + \frac{1}{2\sigma^2}} \end{cases}$$

$$\mu_j \sim N(\beta_j^{t+1}, K_j^{t+2} I_p)$$

$$\beta_j^{t+1} = \frac{\sum_{i=1}^n g_{ij}^{t+1} x_i}{\sum_{i=1}^n g_{ij}^{t+1} + 1/2\sigma^2}$$

$$k_{\theta}^{(t+1)} = \frac{1}{\sum_{j=1}^n \frac{q_{\theta}^{(t+1)}(x_j)}{p(x_j)} + \frac{\lambda}{n}}$$

As good as EM algorithm. Usually difficult analytically (more than 1 step)

$$\pi(\theta_1, \dots, \theta_d | X) = f(\theta_1, \dots, \theta_d) g(\theta_4, \dots, \theta_d) \quad \text{Want as small groups as possible (approximation error)}$$

$$= \prod_{j=1}^d f_j(\theta_j)$$

Similar to Gibbs Sampling (exact)

Some extensions

$$\max_{q \in \mathcal{S}} \max_{\eta} \int q(\theta) \log p(x|\theta, \eta) d\theta - D(q||\pi)$$

prior

"Variational EM algorithm"

$$\max_{q \in \mathcal{S}} \max_{\eta} \int q(\theta) \log p(x|\theta, \eta) d\theta - D(q||\pi_q)$$

Variational empirical Bayes

$$\max_{q \in \mathcal{S}} \max_{\eta, \tau} \int q(\theta) \log p(x|\theta, \eta) d\theta - D(q||\pi_{\eta, \tau})$$

• Bayesian asymptotics

review asymptotics of MLE

$$X_1, \dots, X_n \stackrel{iid}{\sim} P_{\theta} \quad \text{MLE} \quad \hat{\theta} = \arg \max_{\theta} L_n(\theta)$$

$$L_n(\theta) = \sum_{i=1}^n \log P_{\theta}(x_i)$$

$$\sqrt{n}(\hat{\theta} - \theta^*) \sim N(0, I_{\theta^*}^{-1})$$

$$\text{notation} \quad l_{\theta}(x) = \log P_{\theta}(x) \quad \dot{l}_{\theta}(x) = \frac{d}{d\theta} \log P_{\theta}(x)$$

$$\ddot{l}_{\theta}(x) = \frac{d^2}{d\theta^2} \log P_{\theta}(x)$$

Under regularity conditions

$$\begin{cases} \mathbb{E}_{\theta^*} \dot{l}_{\theta^*}(x) = 0 \\ -\mathbb{E}_{\theta^*} \ddot{l}_{\theta^*}(x) = \mathbb{E}_{\theta^*} [\dot{l}_{\theta^*}(x)]^2 = I_{\theta^*} \end{cases}$$

$$L_n(\theta) - L_n(\theta^*) = (\theta - \theta^*) \sum_{i=1}^n \dot{l}_{\theta^*}(x_i) + \frac{1}{2} (\theta - \theta^*)^2 \sum_{i=1}^n \ddot{l}_{\theta^*}(x_i) + \text{remainder}$$

(TE around  $\theta^*$ )

$$\theta = \theta^* + \frac{1}{\sqrt{n}}$$

$$L_n(\theta^* + \frac{1}{\sqrt{n}}) - L_n(\theta^*) = \lambda \frac{1}{n} \sum_{i=1}^n \hat{l}_{\theta^*}(x_i) - \frac{1}{2} h^2 \frac{1}{n} \sum_{i=1}^n (-\hat{l}_{\theta^*}(x_i)) + \text{remainder}$$

$$\text{CLT: } \frac{1}{n} \sum_{i=1}^n \hat{l}_{\theta^*}(x_i) \rightarrow N(0, I_{\theta^*})$$

$$\text{LLN: } \frac{1}{n} \sum_{i=1}^n (-\hat{l}_{\theta^*}(x_i)) \rightarrow I_{\theta^*}$$

$$\Rightarrow L_n(\theta^* + \frac{1}{\sqrt{n}}) - L_n(\theta^*) = h \Delta_n - \frac{1}{2} h^2 I_{\theta^*} + o_p(1)$$

$$\Delta_n \sim N(0, I_{\theta^*})$$

$$\text{MLE } \hat{\theta} = \underset{\theta}{\operatorname{argmax}} L_n(\theta)$$

$$h = \sqrt{n}(\hat{\theta} - \theta^*) = \underset{h}{\operatorname{argmax}} L_n(\theta^* + \frac{h}{\sqrt{n}})$$

$$\underset{\text{"uniform"}}{\approx} \underset{h}{\operatorname{argmax}} \left( h \Delta_n - \frac{1}{2} h^2 I_{\theta^*} \right)$$

$$= I_{\theta^*}^{-1} \Delta_n \sim N(0, I_{\theta^*}^{-1})$$

## Lecture 11

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

$$L_n(\theta) = \sum_{i=1}^n \log P_\theta(X_i) = \sum_{i=1}^n l_\theta(x_i)$$

$$L_n(\theta + \frac{h}{n}) - L_n(\theta^*) = h \Delta_n - \frac{1}{2} h^2 I_{\theta^*} + o_p(1)$$

$$\Delta_n = \frac{1}{n} \sum_{i=1}^n l'_{\theta^*}(x_i) \sim N(0, I_{\theta^*}^{-1})$$

$$I_{\theta^*} = -E_{\theta^*} l'_{\theta^*}(x) = E_{\theta^*} (l'_{\theta^*}(x))^2$$

$$\hat{\theta} = \operatorname{argmax}_{\theta} L_n(\theta) \quad \sqrt{n}(\hat{\theta} - \theta^*) = I_{\theta^*}^{-1} \Delta_n + o_p(1) \sim N(0, I_{\theta^*}^{-1})$$

$\Rightarrow$  frequentist POV

$$X_1, \dots, X_n | \theta \stackrel{\text{iid}}{\sim} P_\theta \quad \theta \sim \pi$$

$$P(\theta | X_1, \dots, X_n) \propto \pi(\theta) \prod_{i=1}^n P(X_i | \theta)$$

Choose  $\hat{C}$  such that  $\pi(\theta \in \hat{C} | X_1, \dots, X_n) = 0.95$

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

Is  $P_{\theta^*}(\theta^* \in \hat{C}) \approx 0.95$ ?

Theorem (Bernstein-Von-Mises)

Under regularity conditions  $\hat{\theta} = \operatorname{argmax}_{\theta} L_n(\theta)$

$$\sqrt{n}(\hat{\theta} - \theta^*) | X_1, \dots, X_n \xrightarrow{P_{\theta^*}} N(0, I_{\theta^*}^{-1})$$

$\Rightarrow$  frequentist confidence interval  $\propto$  Bayesian credible interval

\* only for parametric model (high-dim, non-parametric Model: not true)

Comparison  $\sqrt{n}(\hat{\theta} - \theta) | \theta = \theta^* \xrightarrow{P_{\theta^*}} N(0, I_{\theta^*}^{-1})$

BVM  $\sqrt{n}(\hat{\theta} - \theta) | X_1, \dots, X_n \xrightarrow{\text{frequentist POV}} N(0, I_{\theta^*}^{-1})$

$$\lim_{n \rightarrow \infty} P_{\theta^*}(|\sqrt{n}(\hat{\theta} - \theta) - z| > \varepsilon | X_1, \dots, X_n) \rightarrow 0$$

Lemma If  $\int e^{tx} dP_n(x) \rightarrow \int e^{tx} dP(x)$  in probability  $\forall t \in \mathbb{R}$

random probability (the posterior dist.)

then  $\sup_{t \in \mathbb{R}} |\mathbb{P}_n((-\infty, t]) - \mathbb{P}((-\infty, t])| \rightarrow 0$  in prob  $\lim_{n \rightarrow \infty} \mathbb{P}_{\theta^*} \left( \left| \int_{-\infty}^t \sqrt{n}(\hat{\theta} - \theta) d\pi(\theta | X_1, \dots, X_n) - \int_{-\infty}^t N(0, I_{\theta^*}^{-1}) d\theta \right| > \varepsilon \right) \rightarrow 0 \quad \forall t$

Kolmogorov-Smirnov distance

$$\begin{aligned}
 & \text{SIS } \int e^{t\sqrt{n}(\hat{\theta} - \theta)} d\pi(\theta | x_1, \dots, x_n) \xrightarrow{\text{prob}} e^{\frac{1}{2} I_{\theta^*} t^2} \\
 &= \frac{\int e^{t\sqrt{n}(\hat{\theta} - \theta) + L_n(\theta) - L_n(\theta^*) - \log \pi(\theta)}}{\int e^{L_n(\theta) - L_n(\theta^*) + \log \pi(\theta)} d\theta} \\
 &= \frac{\int e^{t\sqrt{n}(\hat{\theta} - \theta) + L_n(\theta) - L_n(\theta^*) + \log \pi(\theta) + t\sqrt{n}(\theta^* - \hat{\theta})} d\theta}{\int e^{L_n(\theta) - L_n(\theta^*) + \log \pi(\theta)} d\theta}
 \end{aligned}$$

$$\begin{aligned}
 \text{COV} \quad \frac{\int e^{th + L_n(\theta^* + \frac{h}{n}) - L_n(\theta^*) + \log \pi(\theta^* + \frac{h}{n}) + t\sqrt{n}(\theta^* - \hat{\theta})} dh}{\int e^{L_n(\theta^* + \frac{h}{n}) - L_n(\theta^*) + \log \pi(\theta^* + \frac{h}{n})} dh}
 \end{aligned}$$

$$\begin{aligned}
 & \cancel{\int e^{th + h\Delta_n - \frac{1}{2} h^2 I_{\theta^*} + \log \pi(\theta^* + \frac{h}{n}) - t\sqrt{I_{\theta^*}} \Delta_n} dh} \\
 \downarrow \approx & \frac{\int e^{th + h\Delta_n - \frac{1}{2} h^2 I_{\theta^*} + \cancel{\log \pi(\theta^* + \frac{h}{n})}} dh}{\int e^{h\Delta_n - \frac{1}{2} h^2 I_{\theta^*}} dh}
 \end{aligned}$$

$$\begin{cases} \sqrt{n}(\hat{\theta} - \theta^*) = I_{\theta^*}^{-1} \Delta_n + o_p(1) \\ (L_n(\theta^* + \frac{h}{n}) - L_n(\theta^*))n = h\Delta_n - \frac{1}{2} h^2 I_{\theta^*} + o_p(1) \end{cases} ?$$

Why prior is "washed out" as we get more data

since  $\pi$  is cancelled

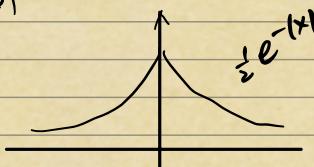
$$= e^{-t I_{\theta^*}^{-1} \Delta_n} \mathbb{E}_{h \sim N(I_{\theta^*}^{-1} \Delta_n, I_{\theta^*}^{-1})} e^{th}$$

$$= e^{-t I_{\theta^*}^{-1} \Delta_n} e^{+I_{\theta^*}^{-1} \Delta_n + \frac{1}{2} t^2 I_{\theta^*}}$$

$$= e^{\frac{1}{2} t^2 I_{\theta^*}} \quad \text{MGF of } N(0, I_{\theta^*})$$

Second derivative may not exist

$$\text{e.g. } P_\theta(x) = \frac{1}{2} e^{-|x-\theta|}$$

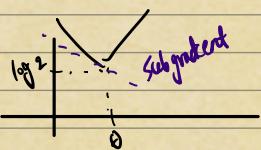


Laplace density

$$\text{MLE } \hat{\theta} = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^n |x_i - \theta|$$

Sample median

$$-\hat{I}_\theta(x) = -\hat{I}_\theta P_\theta(x) = |x - \theta| + \log 2$$



$$-\hat{I}_\theta(x) = \begin{cases} -1 & \theta < x \\ [0, 1] & \theta = x \\ 1 & \theta > x \end{cases}$$

$$-\hat{I}_\theta(x) = \begin{cases} 0 & \theta < x \\ \infty? & \theta = x \\ 0 & \theta > x \end{cases}$$

$$I_{\theta^*} = -E_{\theta^*} \hat{I}_\theta(x) = E_{\theta^*} (\hat{I}_{\theta^*}(x))^2$$

" "  
?  
|

not well-defined

Le Cam & Hajek

LAN (local asymptotic normality)

$$\nabla \left( L_n(\theta + \frac{h}{n}) - L_n(\theta) \right) = h \Delta_{n,\theta} - \frac{h^2}{2} I_\theta + o(h) \quad (\text{previously seen})$$

for some  $\Delta_{n,\theta}$  and  $I_\theta$  such that  $\Delta_{n,\theta} \sim N(0, I_\theta)$

DQM (differentiability under quadratic mean)

$$\exists \hat{I}_\theta \text{ st. } \int \left( \sqrt{P_{\theta+h}} - \sqrt{P_\theta} - \frac{1}{2} h \hat{I}_\theta \sqrt{P_\theta} \right)^2 = o(h^2) \text{ as } h \rightarrow 0$$

if it exists

$$\text{e.g. can take } \hat{I}_\theta = \frac{\partial}{\partial \theta} \log P_\theta(x)$$

$$\frac{1}{2} \hat{I}_\theta \sqrt{P_\theta} = \frac{\partial}{\partial \theta} \sqrt{P_\theta} \quad (\text{Slightly stronger than the first derivative})$$

$$\text{for } P_\theta = \frac{1}{2} x^{(2-\theta)} \text{ take } \hat{I}_\theta(x) = \begin{cases} -1 & \theta < x \\ 0 & \theta = x \\ 1 & \theta > x \end{cases}$$

Can check DQM is satisfied

Theorem If DQM, then

①  $E_\theta \hat{I}_\theta(x) = 0 \rightarrow$  call this  $\text{is score}$

②  $I_\theta = E_\theta (\hat{I}_\theta(x))^2 < \infty$  (call this as Fisher information)

③ LAN holds with  $\Delta_{n,\theta} = \frac{1}{n} \sum_{i=1}^n \hat{I}_{\theta^*}(x_i)$

$$I_0 = \mathbb{E}_\theta (\hat{J}_\theta(\omega))^2$$

enough to cover Laplace  
Without this condition we can prove that  
we can create an alternative  
estimator that only needs LAN  
for asymptotic normality

$$\sqrt{n}(\hat{\theta}_{MLE} - \theta^*) = \bar{J}_{\theta^*}^{-1} \Delta_{n,\theta^*} + O_p(1)$$

$$\sim N(0, \bar{J}_{\theta^*}^{-1})$$

Theorem (Le Cam) If

① LAN

②  $\mathbb{P}(|\sqrt{n}(\theta - \theta^*)| > M_n | X_1, \dots, X_n) = o_p(1)$  for some  $M_n \rightarrow \infty$  "  $\sqrt{n}$  - consistent "

$\Leftrightarrow |\theta - \theta^*| = O_p(\frac{1}{\sqrt{n}})$  under posterior

then BVM holds

$$\boxed{\begin{aligned} X &= Q(\mathbf{u}) \\ \Leftrightarrow P(X > M_n) &\rightarrow 0 \quad \forall M_n \rightarrow \infty \end{aligned}}$$

$$TV(P_{\sqrt{n}(\theta - \theta^*) | X_1, \dots, X_n}, N(\bar{J}_{\theta^*}^{-1} \Delta_{n,\theta^*}, \bar{J}_{\theta^*}^{-1})) \xrightarrow{\bar{P}_{\theta^*}} 0$$

$$( \sup_{t \in \mathbb{R}} |P_n((-\infty, t]) - P((-\infty, t])| \rightarrow 0 \quad \text{vs. } \sup_{A} |P_n(A) - P(A)|, \text{ doesn't assume MLE} )$$

Stronger

Corollary: If additionally  $\sqrt{n}(\hat{\theta}_{MLE} - \theta^*) = \bar{J}_{\theta^*}^{-1} \Delta_{n,\theta^*} + o_p(1)$  then

$$\hat{\theta} = \theta^* + \frac{1}{\sqrt{n}} \bar{J}_{\theta^*}^{-1} \Delta_n + o_p(\frac{1}{\sqrt{n}})$$

$\nwarrow$  can be posterior mean/mode  
MLE, etc

$$TV(P_{\sqrt{n}(\theta - \hat{\theta}_{MLE}) | X_1, \dots, X_n}, N(0, \bar{J}_{\theta^*}^{-1})) \xrightarrow{\bar{P}_{\theta^*}} 0$$

$$X_1, \dots, X_n \stackrel{iid}{\sim} \text{Unif}([0, \theta]) \in \text{LAE}$$

$$|\hat{\theta}_{MLE} - \theta^*| = O_p(\frac{1}{\sqrt{n}}) \quad (\text{local asymptotic exponentiality})$$

$$\sqrt{n}(\hat{\theta}_{MLE} - \theta^*) \sim \bar{E}_p \dots$$

LAMN: financial mathematics

## Lecture 12

$(P_\theta : \theta \in \Theta)$

$$\text{LAN} : L_n(\theta) = \sum_{i=1}^n \log P_\theta(x_i)$$

$$L_n(\theta + \frac{h}{m}) - L_n(\theta) = h \Delta_{n,\theta} - \frac{1}{2} h^2 I_\theta + o_p(1)$$

for some  $\Delta_{n,\theta}$  and  $I_\theta$  s.t.  $\Delta_{n,\theta} \sim N(0, I_\theta)$

Thm (BvM, Le Cam) If ① LAN

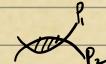
②  $\pi(|\sqrt{n}(\theta - \theta^*)| > M_n | x_1, \dots, x_n) = o_p(1)$  for some  $M_n \rightarrow \infty$

③  $\pi$  is smooth at  $\theta^*$

Then  $TV(P_{\pi(\theta-\theta^*)}(x_1, \dots, x_n), N(I_{\theta^*}^{-1} \Delta_{n,\theta}, I_{\theta^*}^{-1})) \xrightarrow{\theta^*} 0$

Proof) Let  $h = \sqrt{n}(\theta - \theta^*)$ ,  $C_n = \{|\sqrt{n}(\theta - \theta^*)| \leq M_n\} = \{|h| \leq M_n\}$

$$\pi^{C_n}(B|x) = \frac{\pi(B \cap C_n|x)}{\pi(C_n|x)}$$



$$TV(\pi^{C_n}(\cdot|x), \pi(\cdot|x)) = \sup_B |\pi^{C_n}(B|x) - \pi(B|x)|$$

$$= \sup_B \left| \frac{\pi(B \cap C_n|x) - \pi(B|x)\pi(C_n|x)}{\pi(C_n|x)} \right|$$

$$\leq \sup_B \left| \frac{\pi(B \cap C_n|x) - \pi(B|x)}{\pi(C_n|x)} \right| + \sup_B \left| \frac{\pi(B|x) - \pi(B|x)\pi(C_n|x)}{\pi(C_n|x)} \right|$$

$$\leq \left| \frac{\pi(C_n^c|x)}{\pi(C_n|x)} \right| \xrightarrow{\theta^*} 0 \quad \text{by ②}$$

also have  $TV(N(I_{\theta^*}^{-1} \Delta_{n,\theta^*}, I_{\theta^*}^{-1}), N^{C_n}(I_{\theta^*}^{-1} \Delta_{n,\theta^*}, I_{\theta^*}^{-1})) \xrightarrow{\theta^*} 0$

STS to bound

$$TV(P_{h|x}^{C_n}, N^{C_n}(I_{\theta^*}^{-1} \Delta_{n,\theta^*}, I_{\theta^*}^{-1}))$$

density

$$\frac{1_{C_n}(h) \prod_{i=1}^n P_{\theta^* + \frac{h}{m}}(x_i) \pi(\theta^* + \frac{h}{m})}{\int_{C_n} \prod_{i=1}^n P_{\theta^* + \frac{h}{m}}(x_i) \pi(\theta^* + \theta/m) d\theta}$$

density  $N^{C_n}(I_{\theta^*}^{-1} \Delta_{n,\theta^*}, I_{\theta^*}^{-1})(h)$

$$TV(P, Q) = \frac{1}{2} \int |P - Q| = \frac{1}{2} \int_{P > Q} (P - Q) + \frac{1}{2} \int_{Q > P} (Q - P)$$

$$= \int (\bar{P} - g)_+ \quad X_+ = \max(X, 0)$$

$$= \int (1 - \frac{g}{\bar{P}})_+ d\bar{P}$$

$TV(\quad)$

$$= \int \left( - \frac{\int_{C_n} dN^{C_n}(I_{\theta^*}^{-1} \Delta_{n,\theta^*}, I_{\theta^*}^{-1})(h) \int_{C_n} \prod_{j=1}^n P_{\theta^* + \frac{j}{m}}(x_j) \pi(\theta^* + \frac{j}{m}) dj}{\int_{C_n} 1_{C_n}(h) \prod_{j=1}^n P_{\theta^* + \frac{j}{m}}(x_j) \pi(\theta^* + \frac{j}{m})} \right) dP_{hx}^{C_n}$$

$$= \int \left( - \int_{C_n} \frac{dN^{C_n}(I_{\theta^*}^{-1} \Delta_{n,\theta^*}, I_{\theta^*}^{-1})(h) \int_{C_n} \prod_{j=1}^n P_{\theta^* + \frac{j}{m}}(x_j) \pi(\theta^* + \frac{j}{m}) dj}{1_{C_n}(h) \prod_{j=1}^n P_{\theta^* + \frac{j}{m}}(x_j) \pi(\theta^* + \frac{j}{m})} \right)_+ dP_{hx}^{C_n}$$

outside of support

$$= \int \left( - \int_{C_n} \frac{\prod_{j=1}^n P_{\theta^* + \frac{j}{m}}(x_j) dN^{C_n}(I_{\theta^*}^{-1} \Delta_{n,\theta^*}, I_{\theta^*}^{-1})(h) \pi(\theta^* + \frac{j}{m})}{\prod_{j=1}^n P_{\theta^* + \frac{j}{m}}(x_j) dN^{C_n}(I_{\theta^*}^{-1} \Delta_{n,\theta^*}, I_{\theta^*}^{-1})(j) \pi(\theta^* + \frac{j}{m})} dN^{C_n}(I_{\theta^*}^{-1} \Delta_{n,\theta^*}, I_{\theta^*}^{-1})(j) \right)_+ dP_{hx}^{C_n}$$

$$\leq \int \left( - \int_{C_n} \frac{\prod_{j=1}^n P_{\theta^* + \frac{j}{m}}(x_j) dN^{C_n}(I_{\theta^*}^{-1} \Delta_{n,\theta^*}, I_{\theta^*}^{-1})(h) \pi(\theta^* + \frac{j}{m})}{\prod_{j=1}^n P_{\theta^* + \frac{j}{m}}(x_j) dN^{C_n}(I_{\theta^*}^{-1} \Delta_{n,\theta^*}, I_{\theta^*}^{-1})(j) \pi(\theta^* + \frac{j}{m})} \right)_+ dN^{C_n}(I_{\theta^*}^{-1} \Delta_{n,\theta^*}, I_{\theta^*}^{-1})(j) dP_{hx}^{C_n}$$

Jensen

$\exists j \in M_n$  more stable than  $m_n$  ?

$$\frac{\pi(\theta^* + \frac{j}{m})}{\pi(\theta^* + \frac{m}{m})} \rightarrow 1 \quad (\text{prior effect washed out})$$

$\rightarrow 0$  by DCT

(there is subtlety,  
the measures depend on  
a. ready  
assign)

$$= \frac{e^{\ln(\theta^* + \frac{m}{m}) - \ln(\theta^*)}}{e^{\ln(\theta^* + \frac{m}{m}) - \ln(\theta^*)}} \frac{dN(I_{\theta^*}^{-1} \Delta_{n,\theta^*}, I_{\theta^*}^{-1})(h)}{dN(I_{\theta^*}^{-1} \Delta_{n,\theta^*}, I_{\theta^*}^{-1})(j)} \rightarrow 1 \quad \text{by } \Phi$$

Other Bayesian asymptotic results

① Least-informative prior

$$D(\pi(\theta|x) \| \pi(\theta)) \quad p(x|\theta) \Rightarrow \text{proportional. (why)}$$

$$m(x) = \int p(x|\theta) \pi(\theta) d\theta$$

$$\max_T \int_D D(\pi(\theta|x) \| \pi(\theta)) m(x) dx$$

$$\frac{\pi(\theta|x)}{\pi(\theta)} = \frac{p(x|\theta)}{m(x)}$$

$$\begin{aligned} & \int D(\pi(\theta|x) \| \pi(\theta)) m(x) dx \\ &= \iint \pi(\theta|x) \log \frac{\pi(\theta|x)}{\pi(\theta)} d\theta m(x) dx \\ &= \iint \pi(\theta) p(x|\theta) \log \frac{p(x|\theta)}{m(x)} d\theta dx \\ &= \iint D(p(x|\theta) \| m(x) \pi(\theta)) dx d\theta \end{aligned}$$

$$x_1 \dots x_n | \theta \stackrel{iid}{\sim} p(x|\theta)$$

$$\frac{m(x)}{p(x|\theta)} = \frac{\int_{\tilde{\theta}} \pi(\tilde{\theta}) \pi(p(x|\tilde{\theta})) d\tilde{\theta}}{\int_{\theta} p(x|\theta)} \quad (\theta, \tilde{\theta} \text{ are in neighborhood of } \theta^*)$$

$$= \int \pi(\theta) e^{L_n(\tilde{\theta}) - L_n(\theta)} d\tilde{\theta}$$

$$\approx \pi(\theta) \int e^{L_n(\tilde{\theta}) - L_n(\theta)} d\tilde{\theta}$$

$$\approx \pi(\theta) \int e^{\ln(\tilde{\theta} - \theta) \Delta_{n,\theta} - \frac{1}{2} n(\tilde{\theta} - \theta)^2 I_\theta} d\tilde{\theta}$$

$$h = \sqrt{n}(\tilde{\theta} - \theta) \quad h \Delta_{n,\theta} - \frac{1}{2} h^2 I_\theta = -\frac{1}{2} I_\theta (h^2 - 2I_\theta^{-1} h \Delta_{n,\theta})$$

$$= -\frac{1}{2} I_\theta (h - I_\theta^{-1} \Delta_{n,\theta})^2 + \frac{\Delta_{n,\theta}^2}{2 I_\theta}$$

$$= \pi(\theta) e^{\frac{1}{2} \frac{\Delta_{n,\theta}^2}{I_\theta}} \int e^{-\frac{1}{2} I_\theta (\ln(\tilde{\theta} - \theta) - I_\theta^{-1} \Delta_{n,\theta})^2} d\tilde{\theta}$$

$$= \pi(\theta) e^{\frac{1}{2} \frac{\Delta_{n,\theta}^2}{I_\theta}} \sqrt{\frac{2\pi}{n I_\theta}}$$

$$\text{In summary, } \frac{m(x)}{p(x|\theta)} \approx \pi(\theta) e^{\frac{1}{2} \frac{\Delta_{n,\theta}^2}{I_\theta}} \sqrt{\frac{2\pi}{n I_\theta}}$$

$$\int D(\pi(\theta|x) \| \pi(\theta)) m(x) dx$$

$$= \iint \underbrace{\pi(\theta)}_{\approx \mathbb{E}_{x,\theta}} \underbrace{p(x|\theta)}_{\approx \mathbb{E}_{x,\theta}} \underbrace{\log \frac{p(x|\theta)}{m(x)}}_{\approx -\frac{1}{2} \frac{\Delta_{n,\theta}^2}{I_\theta}} dx d\theta$$

$$\approx \mathbb{E}_{x,\theta} \left( -\log \pi(\theta) e^{\frac{1}{2} \frac{\Delta_{n,\theta}^2}{I_\theta}} \sqrt{\frac{2\pi}{n I_\theta}} \right)$$

$$= \mathbb{E}_{x,\theta} \left[ \log \left( \frac{1}{\pi(\theta)} e^{-\frac{1}{2} \frac{\Delta_{n,\theta}^2}{I_\theta}} \sqrt{\frac{n I_\theta}{2\pi}} \right) \right]$$

$$= \mathbb{E}_{x,\theta} \left( -\frac{1}{2} \frac{\Delta_{n,\theta}^2}{I_\theta} + \frac{1}{2} \log \frac{n}{2\pi} + \log \frac{\sqrt{I_\theta}}{\pi(\theta)} \right)$$

$$= \mathbb{E}_{\theta} \left[ \frac{1}{2} \log \frac{\pi(\theta)}{2\pi e} + \int \pi(\theta) \log \frac{\pi(\theta)}{\pi(\theta)} d\theta \right] \quad \theta \in \mathbb{R}$$

Multivariate case  $\theta \in \mathbb{R}^p$

Clarke - Barron expansion

$$\begin{aligned} & \int D(\pi(\theta|x) \| \pi(\theta)) m(x) dx \\ &= \frac{p}{2} \log \frac{\pi}{2\pi e} + \int \pi(\theta) \log \frac{\sqrt{\det(I_\theta)}}{\pi(\theta)} d\theta + o(1) \end{aligned}$$

The maximizer of C-B expansion is  $\pi(\theta) \propto \sqrt{\det(I_\theta)}$

Called Jeffrey's prior e.g.  $N(\theta, I)$

"Jeffrey's prior" is Lebesgue measure  $\curvearrowleft$  not requires

② Bayesian model selection

$$x|\theta_m, m \quad m \in \mathcal{M} \quad \theta_m \in \mathbb{R}^{d_m}$$

$$\theta_m \sim \pi_m$$

$$\int p(x|\theta_m, m) \pi_m(\theta_m) d\theta_m = P(x|m)$$

$$\hat{m} = \underset{m \in \mathcal{M}}{\operatorname{argmax}} p(x|m)$$

BIC

$X | \theta_m, m \quad m \in M \quad \theta_m \sim \pi_m \quad \pi_m \propto p_{\text{true}} \text{ on } \mathbb{R}^{d_m}$

$$\int p(x | \theta_m, m) \pi_m(\theta_m) d\theta_m = p(x|m)$$

$$\hat{m} = \arg \max_{m \in M} p(x|m)$$

Laplace approximation

$$\text{e.g. } e^{nX_1} + e^{nX_2} + \dots + e^{nX_k} = (1 + o(1)) e^{nX_{\max}}$$

$$X_{\max} = \max_{1 \leq i \leq k} X_i \quad \text{WLOG assume } X_{\max} = X_1$$

$$\sum_{i=1}^k e^{nX_i} = e^{nX_1} \left( 1 + e^{n(X_2 - X_1)} + \dots + e^{n(X_k - X_1)} \right) = e^{nX_1} (1 + o(1))$$

all negative

$$\int e^{nf(\theta)} d\theta \quad \hat{\theta} = \arg \max_{\theta} f(\theta) \Rightarrow f'(\hat{\theta}) = 0$$

$$= e^{nf(\hat{\theta})} \int e^{n(f(\theta) - f(\hat{\theta}))} d\theta$$

$$\approx e^{nf(\hat{\theta})} \int e^{n(\theta - \hat{\theta}) \frac{f'(\hat{\theta})}{=0} + \frac{1}{2} n(\theta - \hat{\theta})^2 f''(\hat{\theta})} d\theta$$

Laplace approximation

$$= e^{nf(\hat{\theta})} \int e^{\frac{1}{2} n(\theta - \hat{\theta})^2 f''(\hat{\theta})} d\theta \quad f''(\hat{\theta}) \neq 0 \text{ (assumption)}$$

$$= e^{nf(\hat{\theta})} \sqrt{\frac{2\pi}{-nf''(\hat{\theta})}}$$

Multivariate case  $\theta \in \mathbb{R}^d$

$$\int e^{nf(\theta)} d\theta \approx e^{nf(\hat{\theta})} \left( \frac{2\pi}{n} \right)^{\frac{d}{2}} \int \frac{1}{\det(-\nabla^2 f(\hat{\theta}))}$$

$$\int \pi(\theta) \prod_{i=1}^n p(x_i | \theta) d\theta = \int e^{nf(\theta)} d\theta$$

$$f(\theta) = \frac{1}{n} \sum_{i=1}^n \log p(x_i | \theta) + \frac{1}{n} \log \pi(\theta)$$

$$\hat{\theta} = \arg \max_{\theta} f(\theta) = \hat{\theta}_{\text{MAP}}$$

Maximum a posteriori estimation

$$\nabla f(\theta) = \frac{1}{n} \sum_{i=1}^n \nabla^2 (\log p(x_i | \theta) + \frac{1}{n} \nabla^2 \log \pi(\theta))$$

$$= -I_{\theta} + O_p(1)$$

$$\Rightarrow \nabla^2 f(\hat{\theta}) = -I_{\hat{\theta}} + \text{O}(1) \quad \text{not rigorous; need to verify this goes to zero}$$

apply Laplace approximation

$$\Rightarrow \int \pi(\theta) \prod_{i=1}^n p(x_i|\theta) d\theta \approx \pi(\hat{\theta}) \prod_{i=1}^n p(x_i|\hat{\theta}) \left(\frac{2\pi}{n}\right)^{\frac{1}{2}} \sqrt{\frac{1}{\det(I_{\hat{\theta}})}}$$

$$(\text{using } \int e^{nf(\theta)} d\theta \propto e^{f(\hat{\theta})} \xrightarrow{n \rightarrow \infty} \sqrt{\frac{1}{\det(-\nabla^2 f(\hat{\theta}))}})$$

Choose  $\pi$  to be Jeffreys prior

$$\left\{ \begin{array}{l} \text{① } \pi(\hat{\theta}) \sqrt{\frac{1}{\det(I_{\hat{\theta}})}} = \text{const} \\ \text{② } \hat{\theta} \cdot \hat{\theta}_{\text{MAP}} \approx \hat{\theta}_{\text{MLE}} \quad (\text{the two estimators are close in second order}) \end{array} \right.$$

$$\int \pi(\theta) \prod_{i=1}^n p(x_i|\theta) d\theta = C \prod_{i=1}^n p(x_i|\hat{\theta}_{\text{MLE}}) \left(\frac{2\pi}{n}\right)^{\frac{1}{2}}$$

$$\Rightarrow \max_{m \in M} p(X|m) \text{ is approximately } \max_{m \in M} \sum_{i=1}^n \log p(x_i|\hat{\theta}_m) - \frac{d_m}{2} \log \left(\frac{n}{2\pi}\right)$$

MLE for model m

BIC

$$\max_{m \in M} \left[ \sum_{i=1}^n \log p(x_i|\hat{\theta}_m) - \frac{d_m}{2} \log n \right]$$

$$\text{regression } y \sim N(X\beta, \sigma^2 I_n) \quad X \in \mathbb{R}^{n \times p} \quad \beta \in \mathbb{R}^p$$

$$p(y|\beta) \propto e^{-\frac{1}{2\sigma^2} \|y - X\beta\|^2}$$

more conservative: for large n, get sparser model

$$\text{BIC: } \min_{\beta} \|y - X\beta\|^2 + \frac{\sigma^2}{2} \log n$$

$\hookrightarrow$  dimension of our model

$$\text{AIC: } \min_{\beta} \|y - X\beta\|^2 + 2\sigma^2 p$$

SURE (Stein's unbiased risk estimate)

$$y \sim N(\mu, \sigma^2 I_n), \mu \in \mathbb{R}^n$$

$$\hat{\mu} = \hat{\mu}(y)$$

$$\epsilon z = y - \mu \sim N(0, \sigma^2 I_n), z \sim N(0, I_n)$$

$$\begin{aligned}
E \|\hat{\mu} - \mu\|^2 &= E \|\hat{\mu} - y + y - \mu\|^2 \\
&= E (\|\hat{\mu} - y\|^2 + \|y - \mu\|^2 + 2 \langle \hat{\mu} - y, y - \mu \rangle) \\
&= E \|\hat{\mu} - y\|^2 + n\sigma^2 + 2 E \langle \hat{\mu}, y - \mu \rangle - 2 E \langle \mu + \sigma z, \sigma z \rangle \\
&= E \|\hat{\mu} - y\|^2 - n\sigma^2 + 2 E \langle \hat{\mu}, \sigma z \rangle \\
&= E \|\hat{\mu} - y\|^2 - n\sigma^2 + 2\sigma E \sum_{i=1}^n \frac{d\hat{\mu}_i}{dy_i} \\
&= E \|\hat{\mu} - y\|^2 - n\sigma^2 + 2\sigma E \sum_{i=1}^n \frac{d\hat{\mu}_i}{dy_i} \cdot \frac{dy_i}{dz_i} \\
&= E \|\hat{\mu} - y\|^2 - n\sigma^2 + 2\sigma^2 E \sum_{i=1}^n \frac{d\hat{\mu}_i}{dz_i}
\end{aligned}$$

$$\text{SURE}(\hat{\mu}) = \underbrace{\|\hat{\mu} - y\|^2 - n\sigma^2}_{\text{fitness}} + \underbrace{2\sigma^2 \sum_{i=1}^n \frac{d\hat{\mu}_i}{dz_i}}_{\text{complexity}}, \text{ df of } \hat{\mu}$$

Theorem  $E \text{SURE}(\hat{\mu}) = E \|\hat{\mu} - \mu\|^2$  \* true data generating process

Linear model  $y \sim N(X\beta, \sigma^2 I_n)$  (Working Model)

$$\begin{aligned}
\hat{\mu} &= X\hat{\beta} = Hy & H &= X(X^T X)^{-1} X^T \text{ if } X^T X \text{ is invertible} \\
&\xrightarrow{\text{LSE}} & &= P_x \quad (\text{in general})
\end{aligned}$$

$$\begin{aligned}
\text{df}(\hat{\mu}) &= \sum_{i=1}^n \frac{d\hat{\mu}_i}{dz_i} = \sum_{i=1}^n H_{ii} = \text{Tr}(H) = p \text{ if } X^T X \text{ invertible} \\
&= \text{rank}(X) \quad (\text{in general})
\end{aligned}$$

$$\text{SURE}(X\hat{\beta}) = \|X\hat{\beta} - y\|^2 - n\sigma^2 + 2\sigma^2 p$$

$$\text{AIC} \quad \min_{\beta} \|y - X\beta\|^2 + 2\sigma^2 p$$

All-subset selection  $y \sim N(X\beta, \sigma^2 I_n)$   
full model

assume  $\beta$  is sparse:  $\beta_j = 0$  for many  $j$ 's

can we find the best  $X_S \beta_S$  for  $S \subseteq [p] = \{1, \dots, p\}$

$$\ell_0 \text{ penalty} \quad \|\beta\|_0 = \sum_{j=1}^p \mathbb{1}_{\{\beta_j \neq 0\}}$$

$$\min \left( \|y - X\beta\|^2 + \frac{\lambda}{n} \|\beta\|_0 \right)$$

→ penalty

$$= \min_{1 \leq s \leq p} \min_{\{\beta \in \mathbb{R}^p \mid \|s\|_1 \leq s\}} \min_{\beta_s} \|y - X_s \beta_s\|^2 + \gamma s$$

$$\text{AIC: } \pi = 2s^2$$

$$\text{BIC: } \pi = s^2 \log n$$

$$\inf_{\hat{\beta}} \sup_{\beta \in \mathbb{R}^p} E \|x(\hat{\beta} - \beta)\|^2 = s^2 p$$

$$\inf_{\hat{\beta}} \sup_{\|\beta\|_1 \leq s} E \|x(\hat{\beta} - \beta)\|^2 \times s^2 \log \frac{ep}{s}$$

group study

(Binge & Mourtz)

$$s \log \left( \frac{ep}{s} \right) \times \log \left( \frac{p}{s} \right)$$

cost to pay  
for not knowing the location of the nonzero  $\beta_j$ 's.

achieved by

$$\min_{\beta} \left[ \|y - X\beta\|^2 + C s^2 \|\beta\|_1 \log \left( \frac{ep}{\|\beta\|_1} \right) \right]$$

Some Constant

BIC seems to be a bit better. Usually  $n \gg p$   
(in an extremely high dimensional model)

(hard to solve)

Convex relaxation  $\ell_0 \rightarrow \ell_1$

$$\text{Lasso } \min_{\beta} \left[ \|y - X\beta\|^2 + \gamma \|\beta\|_1 \right] \quad \|\beta\|_1 = \sum_{j=1}^p |\beta_j|$$

↓ Statistical POV, equivalent. Optimization sense, same

Tibshirani ( $\min \|y - X\beta\|^2$  s.t.  $\|\beta\|_1 \leq p$ )

earlier works

|                |   |
|----------------|---|
| Donoho & Chen  | } |
| Fan & Friedman |   |

How to solve Lasso?

special case  $X = I_p$   $n = p$

$$\min_{\beta} \left[ \|y - \beta\|^2 + \gamma \|\beta\|_1 \right]$$

$$= \min_{\beta} \sum_{j=1}^p [(y_j - \beta_j)^2 + \gamma |\beta_j|]$$

$$= \sum_{j=1}^p \min_{\beta_j} [ (y_j - \beta_j)^2 + \gamma |\beta_j| ]$$

one-dimensional problem

$$\min_b [(y-b)^2 + \pi |b|] = \min_b f(b)$$

$$f'(b) = \begin{cases} 2(b-y) + \pi & b > 0 \\ 2(b-y) - \pi & b < 0 \end{cases}$$

$$2(b-y) + \pi = 0 \Leftrightarrow b = y - \frac{\pi}{2}$$

$$2(b-y) - \pi = 0 \Leftrightarrow b = y + \frac{\pi}{2}$$

$$\text{Case 1 } y > \frac{\pi}{2} \quad f'(y - \frac{\pi}{2}) = 0 \Rightarrow \hat{b} = y - \frac{\pi}{2}$$

$$\text{Case 2 } y < -\frac{\pi}{2} \quad f'(y + \frac{\pi}{2}) = 0 \Rightarrow \hat{b} = y + \frac{\pi}{2}$$

$$\text{Case 3 } y \in [-\frac{\pi}{2}, \frac{\pi}{2}]$$

$$\text{When } b > 0 \Rightarrow b - y \geq -\frac{\pi}{2} \Rightarrow 2(b-y) \geq -\pi$$

$$\Rightarrow f'(b) \geq 0$$

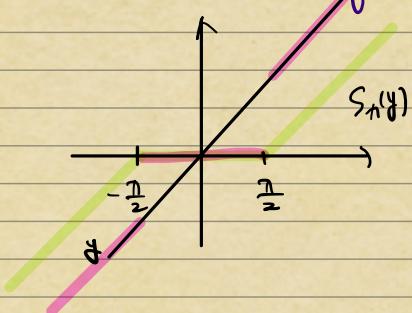
Similarly when  $b < 0 \Rightarrow f'(b) \leq 0$

$$\Rightarrow \hat{b} = 0$$

soft thresholding function

$$\arg \min_b [(y-b)^2 + \pi |b|] = S_\pi(y)$$

$$= \begin{cases} y - \frac{\pi}{2} & y > \frac{\pi}{2} \\ 0 & y \in [-\frac{\pi}{2}, \frac{\pi}{2}] \\ y + \frac{\pi}{2} & y < -\frac{\pi}{2} \end{cases}$$



general case: convex optimization

$$\min_x f(x) \quad f \text{ is convex}$$

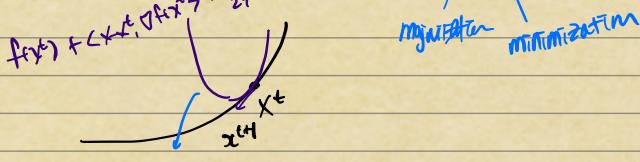
gradient descent

$$x^{t+1} = x^t - \eta \nabla f(x^t) \quad \text{step size}$$

$$= \arg \min_x [f(x^t) + \langle x - x^t, \nabla f(x^t) \rangle + \frac{1}{2\eta} \|x - x^t\|^2]$$

$$x^{t+1} = x^t - \frac{1}{\eta} \nabla f(x^t)$$

MM algorithm



majorization

$$\min_x (f(x) + g(x))$$

at  $x^t$ , an upper bound of  $f(x) + g(x)$  is  $f(x^t) + \langle x - x^t, \nabla f(x^t) \rangle + \frac{1}{2\eta} \|x - x^t\|^2 + g(x)$

$$x^{t+1} = \arg \min_x \left[ f(x^t) + \langle x - x^t, \nabla f(x^t) \rangle + \frac{1}{2\eta} \|x - x^t\|^2 + g(x) \right]$$

$$\text{Lasso: } \min_{\beta} \left[ \|y - X\beta\|^2 + \gamma \|\beta\|_1 \right]$$

$$\begin{cases} f(\beta) = \|y - X\beta\|^2 & \nabla f(\beta) = 2X^T(X\beta - y) \\ g(\beta) = \gamma \|\beta\|_1 & \end{cases}$$

$$\begin{aligned} \beta^{t+1} &= \arg \min_{\beta} \left[ \|\beta - [\beta^t - 2\gamma X^T(X\beta^t - y)]\|^2 + 2\gamma \|\beta\|_1 \right] \\ &= S_{2\gamma} (\beta^t - 2\gamma X^T(X\beta^t - y)) \Rightarrow \text{explain why Lasso solution is sparse} \end{aligned}$$

iterative soft thresholding algorithm (Back & Teboulle)

### risk analysis of Lasso

$$y \sim N(X\beta, \sigma^2 I)$$

$$\hat{\beta} = \arg \min_{\beta} \left[ \|y - X\beta\|^2 + \gamma \|\beta\|_1 \right]$$

basic inequality

$$\|y - X\hat{\beta}\|^2 + \gamma \|\hat{\beta}\|_1 \leq \|y - X\beta\|^2 + \gamma \|\beta\|_1$$

$\uparrow$   
true  $\beta$

$$\begin{aligned} \|y - X\hat{\beta}\|^2 &= \|y - X\beta + X\beta - X\hat{\beta}\|^2 \\ &= \|y - X\beta\|^2 + \|X(\hat{\beta} - \beta)\|^2 + 2 \langle y - X\beta, X(\hat{\beta} - \beta) \rangle \end{aligned}$$

$$\Rightarrow \|X(\hat{\beta} - \beta)\|^2 + 2 \langle y - X\beta, X(\hat{\beta} - \beta) \rangle + \gamma \|\hat{\beta}\|_1 \leq \gamma \|\beta\|_1$$

$$\Rightarrow \|X(\hat{\beta} - \beta)\|^2 \leq \underbrace{2 \langle y - X\beta, X(\hat{\beta} - \beta) \rangle + \gamma \|\beta\|_1}_{2(y - X\beta)^T X(\hat{\beta} - \beta)} - \gamma \|\hat{\beta}\|_1$$

$$\leq 2 \|x^T(y - x\beta)\|_\infty \|\hat{\beta} - \beta\|_1 \quad (\text{Hölder})$$

If  $\pi \geq 2 \|x^T(y - x\beta)\|_\infty$

$$\text{then } \|x(\hat{\beta} - \beta)\|^2 \leq \pi \|\hat{\beta} - \beta\|_1 + \pi \|\beta\|_1 - \pi \|\hat{\beta}\|_1 \\ \leq 2\pi \|\beta\|_1$$

Lemma If  $\pi \geq 2 \|x^T(y - x\beta)\|_\infty$  then  $\|x(\hat{\beta} - \beta)\|^2 \leq 2\pi \|\beta\|_1$

$$y \sim N(x\beta, \sigma^2 I_n) \quad g = x\beta + \sigma z \quad z \sim N(0, I_n)$$

$$2 \|x^T(y - x\beta)\|_\infty = 2\sigma \|x^T z\|_\infty$$

$$X = \begin{pmatrix} x_1 & \dots & x_p \end{pmatrix}$$

$$\|x^T z\|_\infty = \max_{1 \leq j \leq p} |x_j^T z|$$

$$= \max_{1 \leq j \leq p} \|x_j\| \left| \left( \frac{x_j}{\|x_j\|} \right)^T z \right|$$

$$\leq \max_{1 \leq j \leq p} \|x_j\| \max_{1 \leq j \leq p} \left| \left( \frac{x_j}{\|x_j\|} \right)^T z \right|$$

$$\downarrow \qquad \qquad \qquad \sim N(0, 1)$$

$$\text{assume } \max_{1 \leq j \leq p} \|x_j\| \leq L\sqrt{n}$$

constant

$$\leq L\sqrt{n} \max_{1 \leq j \leq p} \left| \left( \frac{x_j}{\|x_j\|} \right)^T z \right|$$

$$\downarrow$$

Claim: for  $w_1, w_2, \dots, w_p \sim N(0, 1)$

(not necessarily independent)

$$\text{then } \max_{1 \leq j \leq p} |w_j| \lesssim \sqrt{\log p} \quad \text{W.h.p.}$$

$$2 \|x^T(y - x\beta)\|_\infty = 2\sigma \|x^T z\|_\infty$$

$$\leq 2\sigma L\sqrt{n} \max_{1 \leq j \leq p} \left| \frac{|x_j^T z|}{\|x_j\|} \right|$$

$$\lesssim \sigma \sqrt{n \log p} \quad (\text{W.h.p.})$$

Choose  $\pi = C\sigma \sqrt{n \log p}$  then  $2 \|x^T(y - x\beta)\|_\infty \leq \pi \quad \text{W.h.p.}$

$\Rightarrow$  by the lemma  $\| \mathbf{X}(\hat{\beta} - \beta) \|_2^2 \lesssim \sigma \|\beta\|_1 \sqrt{n \log p}$

$$\Leftrightarrow \frac{1}{n} \| \mathbf{X}(\hat{\beta} - \beta) \|_2^2 \lesssim \sigma \|\beta\|_1 \sqrt{\frac{n \log p}{n}}$$

Lec 15

$$y = X\beta + \epsilon \quad X \in \mathbb{R}^{n \times p} \quad \epsilon \sim N(0, I) \quad \beta \in \mathbb{R}^p$$

$$\text{Lasso} \quad \hat{\beta} = \arg \min_{\beta} \left[ \|y - X\beta\|^2 + \gamma \|\beta\|_1 \right]$$

$$\frac{1}{n} \|X(\hat{\beta} - \beta)\|^2 \lesssim \sigma \|\beta\|_1 \sqrt{\frac{n}{n}} \quad \gamma = C \sigma \sqrt{\log p}$$

$$\text{W.h.p.} \quad \Rightarrow 2\|X^T(y - X\beta)\|_{\infty} = 2\sigma \|X^T z\|_{\infty} \leq \gamma$$

Claim: for  $w_1, \dots, w_p \sim N(0, 1)$

$$\max_{1 \leq j \leq p} |w_j| \leq \sqrt{\log p} \quad \text{W.h.p.}$$

$$\text{pf: } \mathbb{P}(|w_i| > t) = \mathbb{P}(w_i > t) + \mathbb{P}(w_i < -t)$$

$$= 2\mathbb{P}(w_i > t)$$

$$= 2\mathbb{P}(e^{pw_i} > e^{pt}) \quad p > 0$$

$$\leq 2e^{-pt} + e^{pt}$$

$$= 2e^{-pt + \frac{1}{2}p^2}$$

$$(p=t) \quad = 2e^{-\frac{1}{2}t^2}$$

$$\mathbb{P}\left(\max_{1 \leq j \leq p} |w_j| > t\right) = \mathbb{P}\left(\bigcup_{j=1}^p E_j\right) \quad E_j = \{|w_j| > t\}$$

$$\leq \sum_{j=1}^p \mathbb{P}(E_j)$$

$$\leq 2p e^{-\frac{1}{2}t^2}$$

$$= 2^{-\frac{1}{2}t^2 + (\log 2p)}$$

$$\text{Choose } t^2 = 4 \log(2p)$$

$$\mathbb{P}\left(\max_{1 \leq j \leq p} |w_j| > 2\sqrt{\log(2p)}\right) \leq \frac{1}{2p}$$

$$\Leftarrow \max_{1 \leq j \leq p} |w_j| \leq 2\sqrt{\log(2p)} \quad \text{W.P. at least } 1 - \frac{1}{2p}$$

$$\text{Binge \& Massart} \quad \inf_{\beta} \sup_{\|\beta\|_1 \leq S} \mathbb{E} \|X\hat{\beta} - X\beta\|^2 \propto \frac{C^2 S \log \frac{Cp}{S}}{n}$$

$$\inf_{\beta} \sup_{\|\beta\|_1 \leq S} \mathbb{E} \frac{1}{n} \|X(\hat{\beta} - \beta)\|^2 \propto C \sqrt{\frac{\log p}{n}} \quad \text{so for L, optimal}$$

Candes & Tao "Dantzig Selector"

$$\min \|\beta\|_1$$

$$\text{s.t. } \|X^T(y - X\beta)\|_\infty \leq \gamma \quad (\text{previously } 2\|X^T(y - X\beta)\|_\infty = 2\|X^Tz\|_\infty \leq \gamma)$$

$$\frac{1}{n} \|X(\hat{\beta} - \beta)\|^2 \leq \frac{\sigma^2 \text{supp}}{n R^2} \text{ with high probability}$$

*nonnegative can go to 0 for bad design matrix*

Bickel, Ritov & Tsybakov

$$\{\beta \in \mathbb{R}^p : \|\beta\|_0 \leq S\}$$

$$\|y - X\hat{\beta}\|^2 + \gamma \|\beta\|_1 \leq \|y - X\beta\|^2 + \gamma \|\beta\|_1$$

$$\Rightarrow \|X(\hat{\beta} - \beta)\|^2 \leq 2\sigma \|X^T z\|_\infty \|\hat{\beta} - \beta\|_1 + \gamma \|\beta\|_1 - \gamma \|\hat{\beta}\|_1$$

Last time, we regularize

$$\text{choose } \gamma \geq 4\sigma \|X^T z\|_\infty \quad (\gamma = C\sigma \sqrt{n(p)} )$$

$$\leq \frac{\gamma}{2} \|\hat{\beta} - \beta\|_1 + \gamma \|\beta\|_1 - \gamma \|\hat{\beta}\|_1$$

$$\leq \frac{\gamma}{2} \|\hat{\beta}\|_1 + \frac{\gamma}{2} \|\beta\|_1 + \gamma \|\beta\|_1 - \gamma \|\hat{\beta}\|_1$$

$$\leq \frac{3}{2} \gamma \|\beta\|_1$$

If  $\gamma \geq 4\sigma \|X^T z\|_\infty$ , then  $\|X(\hat{\beta} - \beta)\|^2 \leq \frac{3\gamma}{2} \|\beta\|_1$ .

↓

deduce the same  $\ell_1$ -norm bound (last time)

$$\text{notation } \Delta = \hat{\beta} - \beta, \quad S = \text{supp}(\beta) = \{j \in [p] : \beta_j \neq 0\}$$

$$\forall v \in \mathbb{R}^p \quad \|V_S\|_1 = \sum_{j \in S} |V_j| \quad \|V_{S^c}\|_1 = \sum_{j \in S^c} |V_j|$$

$$\|v\|_1 = \|V_S\|_1 + \|V_{S^c}\|_1$$

$$\|X(\hat{\beta} - \beta)\|^2 \leq \frac{\gamma}{2} \|\Delta\|_1 + \gamma \|\beta\|_1 - \gamma \|\beta + \Delta\|_1$$

$$= \frac{\gamma}{2} \|\Delta_S\|_1 + \frac{\gamma}{2} \|\Delta_{S^c}\|_1 + \gamma \|\beta_S\|_1 - \gamma \|\beta_S + \Delta_S\|_1 - \underline{\gamma \|\Delta_{S^c}\|_1}$$

$$\leq -\frac{\gamma}{2} \|\Delta_{S^c}\|_1 + \frac{\gamma}{2} \|\Delta_S\|_1 + \gamma \|\Delta_S\|_1$$

$$= -\frac{\gamma}{2} \|\Delta_{S^c}\|_1 + \frac{3\gamma}{2} \|\Delta_S\|_1$$

To summarize,  $\|X(\hat{\beta} - \beta)\|^2 \leq -\frac{\gamma}{2} \|\Delta_{S^c}\|_1 + \frac{3\gamma}{2} \|\Delta_S\|_1$

**Lemma** If  $\pi \geq 4 - \epsilon \|X^T z\|_\infty$ , then

$$\left\{ \begin{array}{l} \|X(\hat{\beta} - \beta)\|^2 \leq \frac{3\pi}{2} \|\Delta_S\|_1 \\ \|\Delta_S\|_1 \leq 3 \|\Delta_S\|_1 \quad (\text{cone condition}) \end{array} \right.$$

$$\|X(\hat{\beta} - \beta)\|^2 \leq \frac{3\pi}{2} \|\Delta_S\|_1 \leq \frac{3\pi}{2} \sqrt{s} \|\Delta\|_2$$

↓  
Cauchy-Schwarz

$$\frac{1}{n} \|X\Delta\|^2 = \Delta^T \underbrace{\left( \frac{1}{n} X^T X \right) \Delta}_{\text{sample covariance}} \geq \pi_{\min} \left( \frac{1}{n} X^T X \right) \|\Delta\|^2$$

$$\Rightarrow \|\Delta\| \leq \frac{\frac{1}{n} \|X\Delta\|}{\sqrt{\pi_{\min} \left( \frac{1}{n} X^T X \right)}}$$

$$\|X\Delta\|^2 \leq \frac{3\pi}{2} \sqrt{s} \|\Delta\| \leq \frac{3\pi}{2} \sqrt{s} \frac{\frac{1}{n} \|X\Delta\|}{\sqrt{\pi_{\min} \left( \frac{1}{n} X^T X \right)}}$$

$$\|X\Delta\| \leq \frac{3\pi}{2\sqrt{n}} \frac{\sqrt{s}}{\sqrt{\pi_{\min} \left( \frac{1}{n} X^T X \right)}} \quad \text{Choose } \pi = C \sqrt{n} \log p$$

$$\rightarrow \frac{1}{n} \|X\Delta\|^2 \leq \frac{\frac{6^2 s \log p}{n \pi_{\min} \left( \frac{1}{n} X^T X \right)}}{\frac{4^2 s \log \frac{C \sqrt{n} \log p}{n}}{0}}$$

need  $n \geq p$  so that  $\pi_{\min} \left( \frac{1}{n} X^T X \right) > 0$

$$\pi_{\min} \left( \frac{1}{n} X^T X \right) = \lim_{\Delta \rightarrow 0} \frac{\Delta^T \left( \frac{1}{n} X^T X \right) \Delta}{\|\Delta\|^2}$$

restricted eigenvalue

$$\text{define } K^2 = \min_{\substack{\Delta \neq 0 \\ \|\Delta_S\|_1 \leq 3 \|\Delta_S\|_1}} \frac{\Delta^T \left( \frac{1}{n} X^T X \right) \Delta}{\|\Delta\|^2}$$

With the same analysis

$$\frac{1}{n} \|X(\hat{\beta} - \beta)\|^2 \lesssim \frac{6^2 s \log p}{n K^2}$$

Theorem:  $X_F \stackrel{\text{iid}}{\sim} N(0, I)$ , if  $\frac{s \log p}{n}$  is small then

"each column is orthogonal"

$$K^2 \gtrsim 1 \quad \text{w.h.p.}$$

Thm (Lasso error bound)

Chose  $\pi = C\sigma \sqrt{n/p}$  then w.h.p.

hold for arbitrary design matrix

we have

$$\frac{1}{n} \|X(\hat{\beta} - \beta)\|^2 \lesssim \min \left( \frac{\sigma^2 \|\beta\|_0 p}{n K^2}, \|\beta\|_1 \sqrt{\frac{\sigma^2 p}{n}} \right)$$

$$\|\hat{\beta} - \beta\|^2 \lesssim \min \left( \frac{\sigma^2 \|\beta\|_0 p}{n K^2}, \|\beta\|_1 \sqrt{\frac{\sigma^2 p}{n}} \right)$$

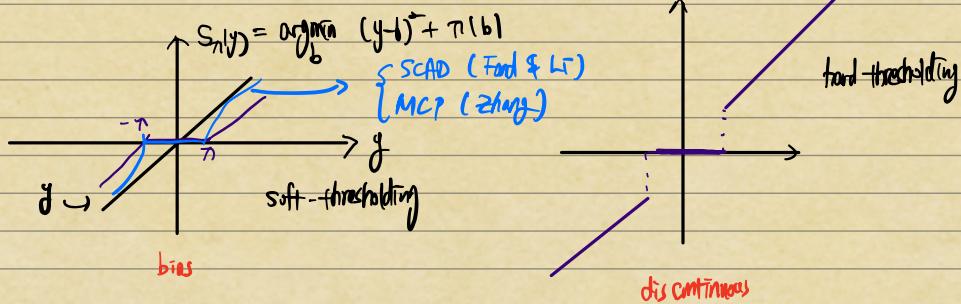
Zhang, Wainwright & Jordan  
 $K$  cannot be completely removed by any polynomial time algorithm

(due to  $\frac{1}{n} \|X\Delta\|^2 \geq \pi_{\min}(\frac{1}{n} X^T X) \|\Delta\|^2$  before)

extensions & improvements

① Lasso penalty

② non-convex penalty



for SCAD or MCP

$$\frac{1}{n} \|X(\hat{\beta} - \beta)\|^2 \lesssim \frac{\sigma^2 (S_1 + S_2 \sqrt{p})}{n K^2}$$

$$S_1 = \left| \left\{ j \in S : |\beta_j| > C \sqrt{\frac{\log p}{n}} \right\} \right| \quad S_2 = \left| \left\{ j \in S, |\beta_j| \leq C \sqrt{\frac{\log p}{n}} \right\} \right|$$

## Lec 16

$$y \sim N(X\beta, \sigma^2 I_n) \quad \hat{\beta} = \underset{\beta}{\operatorname{argmin}} \left( \|y - X\beta\|^2 + \tau \|\beta\|_1 \right)$$

Thm. If  $\max_{1 \leq j \leq p} |x_j| \lesssim \sqrt{n}$  choose  $\tau = C(6) n \log p$

$\uparrow$   
but know (3)

for some large  $C > 0$ , we have

$$\frac{1}{n} \|X(\hat{\beta} - \beta)\|^2 \lesssim \min \left( \frac{\sigma^2 \|\beta\|_0 / \log p}{n K^2}, \|\beta\|_1 \sqrt{\frac{\sigma^2 \log p}{n}} \right)$$

where  $K = \min_{\Delta \neq 0} \frac{\frac{1}{n} \|X\Delta\|}{\|\Delta\|}$   
 $\|\Delta\|_0 \leq 3 \|\Delta\|_1$

extension & improvements

①  $\ell_0$  penalty

② non convex penalty

③ SLOPE (Bogdan, van den Berg, Sabatti, Su, & Candès)

$$\text{achieves } \frac{1}{n} \|X(\hat{\beta} - \beta)\|^2 \lesssim \frac{\sigma^2 s \log(\frac{ep}{\delta})}{n K^2}$$

how to achieve?  $\min \left( \|y - X\beta\|^2 + \sum_{j=1}^p \tau_j |\beta_j|_{(p)} \right)$

$$|\beta|_{(1)} \geq |\beta|_{(2)} \geq \dots \geq |\beta|_{(p)}$$

④ scaled Lasso / square-root Lasso

(Sun & Zhang) (Belloni, Chernozhukov & Wang)

$$F(\beta, \sigma) = \frac{\|y - X\beta\|^2}{\sigma} + n\sigma + \tau \|\beta\|_1$$

jointly convex

$$\text{choose } \tau = C \sqrt{n \log p}$$

$$\left\{ \begin{array}{l} \beta^{t+1} = \underset{\beta}{\operatorname{argmin}} F(\beta, \sigma^t) = \underset{\beta}{\operatorname{argmin}} \left( \|y - X\beta\|^2 + \tau \sigma^t \|\beta\|_1 \right) \\ \sigma^{t+1} = \underset{\sigma}{\operatorname{argmin}} F(\beta^{t+1}, \sigma) = \sqrt{\frac{1}{n} \|y - X\beta^{t+1}\|^2} \end{array} \right.$$

$$\left\{ \begin{array}{l} \beta^{t+1} = \underset{\beta}{\operatorname{argmin}} F(\beta, \sigma^t) = \underset{\beta}{\operatorname{argmin}} \left( \|y - X\beta\|^2 + \tau \sigma^t \|\beta\|_1 \right) \\ \sigma^{t+1} = \underset{\sigma}{\operatorname{argmin}} F(\beta^{t+1}, \sigma) = \sqrt{\frac{1}{n} \|y - X\beta^{t+1}\|^2} \end{array} \right.$$

$$\min_{\beta} F(\beta, \sigma) = G(\beta)$$

square-root lasso

$$\min_{\beta} F(\beta) \Leftrightarrow \min_{\beta} \left( \|y - X\beta\|^2 + \rho \|\beta\|_1 \right)$$

## ⑤ Bayesian methods (Bayesian sparse linear regression)

$$\hat{\beta}_{\text{LASSO}} = \arg \max_{\beta} \left[ -\frac{1}{2n} \|y - X\beta\|^2 - \eta \|\beta\|_1 \right]$$

25:00

a useless posterior distribution

Gao, van der Vaart & Zhou

$$y | \beta \sim N(X\beta, \sigma^2 I_n)$$

$\beta \sim \pi$  is specified by the following sampling process

①  $s \sim \pi(s)$  supported on  $\{1, \dots, p\}$

$$\pi(s) \propto e^{-C s \text{tr}(s)} \frac{\Gamma(s)}{\Gamma(s_0)}$$

② Given  $s$ . Sample  $S \subseteq \{1, \dots, p\}$

uniformly over all subset sets.  $|S|=s$

$$\text{③ Given } (s, S), \beta = \begin{pmatrix} \beta_S \\ \beta_{S^c} \end{pmatrix}, \beta_{S^c} = 0$$

$$\beta_S \sim e^{-\eta \|X_S \beta_S\|} \quad (\text{elliptical Laplace distribution})$$

$$\begin{aligned} x_1, \dots, x_n &\sim N(\mu, 1) \\ \mu &\sim N(0, \tau^2) \quad \mu = E[X_1, \dots, X_n] \\ \rightarrow \sup_{\mu \in \mathbb{R}} E[\hat{\mu} - \mu] &= \infty \end{aligned}$$

Theorem  $\pi(\|\beta - \beta^*\|) \geq C \sigma^2 s^* \log\left(\frac{ep}{s^*}\right) \xrightarrow{\text{prob}} 0$

truth  
under posterior  $\|\beta^*\|_0$

density of elliptical Laplace is

$$\frac{1}{2} \frac{\det(X_S^T X_S)^{-\frac{1}{2}}}{\det(X^T X)^{-\frac{1}{2}}} \frac{\Gamma(\frac{s}{2})}{\Gamma(s)} e^{-\eta \|X_S \beta_S\|}$$

any prior that is compatible?

spike-and-slab prior (George & McCulloch)

①  $\eta \sim \text{Beta}(a, b)$  conjugate

②  $\gamma_1, \dots, \gamma_p \mid \eta \stackrel{\text{iid}}{\sim} \text{Bernoulli}(\eta)$  very small (likely 0, but compatible reason)

③  $\beta_j \mid \gamma_j \sim N(0, \sigma^2 V_j)$ ,  $\gamma_j = 0$  spike

$N(0, \sigma^2 V_j)$ ,  $\gamma_j = 1$  slab  
very large (heavy tail)

can do MCMC

but can be computed by EM (Rocková & George)

Kim & Stephen's method  $\rightarrow$  empirical Bayes

$$y | \beta, \sigma^2 \sim N(X\beta, \sigma^2 I_n)$$

$$\beta_k | \sigma^2, g \stackrel{\text{IID}}{\sim} g_c \quad g_c \stackrel{def}{=} \frac{1}{\sigma} g(\frac{\cdot}{\sigma})$$

$$g = \sum_{k=1}^K \frac{\pi_k}{\pi_k} N(0, \sigma_k^2)$$

$$\sum_{k=1}^K \pi_k = 1$$

$$f(x) = \prod_{i=1}^n \frac{\pi_i}{\pi_0} \cdot \sigma_i^2 \cdot \dots \cdot \sigma_K^2$$

Hypothesis test (multiple test & high dimension)

$$X \sim N(\theta, I_p)$$

$$\beta_0 = N(\theta, I_p)$$

$$H_0: \theta = 0 \quad H_1: \|\theta\| \geq \varepsilon$$

$$\text{minimax risk} \quad \inf_{\phi} \left[ P_0 \phi + \sup_{\theta \in \Sigma} P_\theta (\perp \phi) \right] = R(\varepsilon)$$

minimax separation rate

def.  $\varepsilon^*$  is minimax separation rate if

$$\textcircled{1} \quad \forall f \in C([0,1]), \exists C > 0 \text{ s.t. } R(C\varepsilon^*) < f \quad (\text{upper bound})$$

$$\textcircled{2} \quad \exists c_1, c_2 > 0 \text{ such that } R(G\varepsilon^*) > c_2 \quad (\text{lower bound})$$

Theorem The minimax separation rate is  $\underline{\varepsilon^* = p^{\frac{1}{2}}}$

1:14:00

rte-optimal testing procedure is

$$\phi = \mathbb{I}\left\{ \|x\|^2 > p + C\sqrt{p} \right\}$$

↑ expectation  
order of sd

(also an exact optimal)

$$"P+P"$$

estimation:  $\frac{p}{n}$  should be harder than  $(p^{\frac{1}{2}})^2 = p^{\frac{1}{2}}$

# Lec 19

$$X \sim N(\theta, I_p) \quad H_0: \theta = 0 \quad H_1: \|\theta\| > \varepsilon$$

minimum separation  $\varepsilon^* \sim p^{\frac{1}{2}}$

$$H_0: \theta = 0 \quad H_1: \theta \neq 0 \quad \theta \text{ is sparse}$$

Ingestor formulation

$$H_0: X_1, \dots, X_p \stackrel{iid}{\sim} N(0, 1)$$

$$H_1: X_1, \dots, X_p \stackrel{iid}{\sim} ((1-\varepsilon)N(0, 1) + \varepsilon N(\mu, 1)) \quad (\varepsilon \rightarrow 0 \text{ vs } p \rightarrow \infty)$$

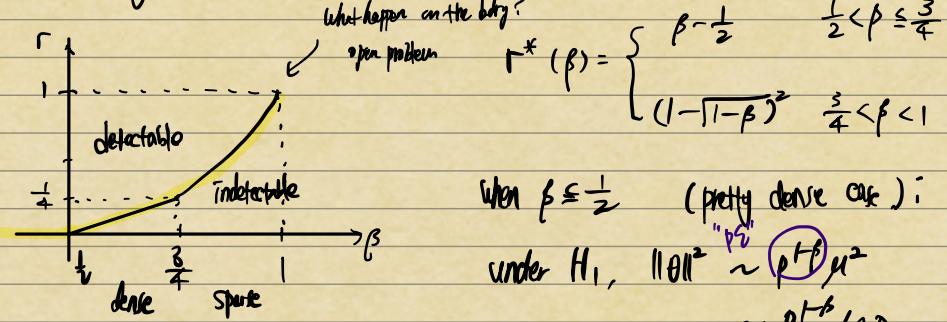
$\mu$  or  $\varepsilon$  small  $\rightarrow$  hard

$$\text{Thrunum (Ingestor)} \quad \text{Set} \quad \Sigma = p^{-\beta} \quad \mu = \sqrt{2r \log p}$$

If  $r > r^*(\beta)$  can have a consistent test

If  $r < r^*(\beta)$  does not exist consistent test

phase diagram



when  $\beta \leq \frac{1}{2}$  (pretty dense case): at least  $\sqrt{p}$  non-zero coefficients

$$\text{under } H_1, \quad \|\theta\|^2 \sim p^{1-\beta} \mu^2 \sim p^{1-\beta} \log p \geq \sqrt{p} \log p \gg \sqrt{p}$$

$$W_1, \dots, W_p \stackrel{iid}{\sim} N(0, 1)$$

$$\max_{1 \leq j \leq p} W_j = (1 + o_p(n)) \sqrt{\log p}$$

as long as  $r > 0$ , detectable using chi square test

$\beta < 1$  necessary