

Lecture 6 - scribbles

Friday, September 21, 2018 13:39

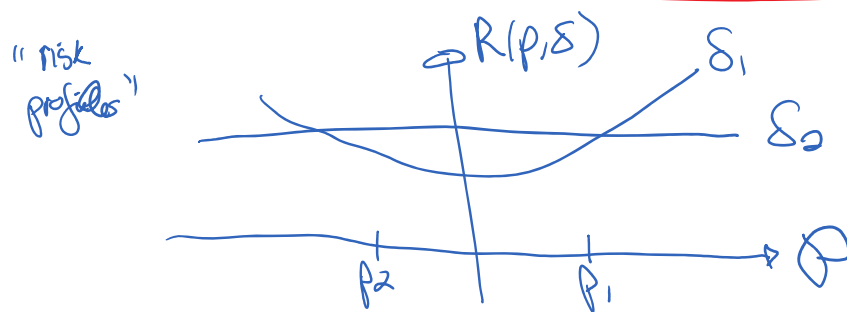
today: finish decision theory
gen/disc. learning

stat decision theory ctd. :

given previous framework, how can we compare stat procedures? e.g. δ_1 vs. δ_2

first property

$$\text{(frequentist) risk } R(p, \delta) \triangleq \mathbb{E}_{D \sim p} [L(p, \delta(D))]$$



$$\delta_i: D \rightarrow \mathcal{A}$$

* transform to scalar:

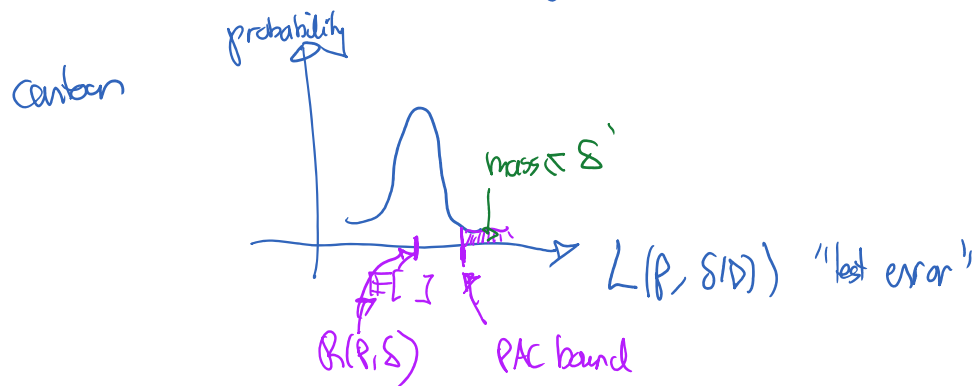
• "minimax" analysis: $\max_{p \in \mathcal{P}} R(p, \delta)$ ↖ "worst case"

• weighted average: $\int_{\mathcal{P}} R(p, \delta) \pi(p) d\mathcal{P}$
↖ frequentist

• weighted average : $\int R(\theta, \delta) \pi(\theta) d\theta$
 (H)

* alternatively: in ML theory, is PAC theory

→ look at tail bounds for dist. of $L(P, \delta(D))$ [D is random]



$\hookrightarrow P\{L(P, \delta(D)) \geq \underbrace{\text{stuff}}_{\text{bound}}\} \leq \underbrace{\delta'}_{\text{small number}}$
 "with high proba" statement

Bayesian decision theory:

→ condition on data D

Bayesian posterior risk

$$R_B(a|D) = \int_{\Theta} L(\theta, a) \underbrace{p(\theta|D)}_{\text{posterior or } p(\theta)p(D|\theta)} d\theta$$

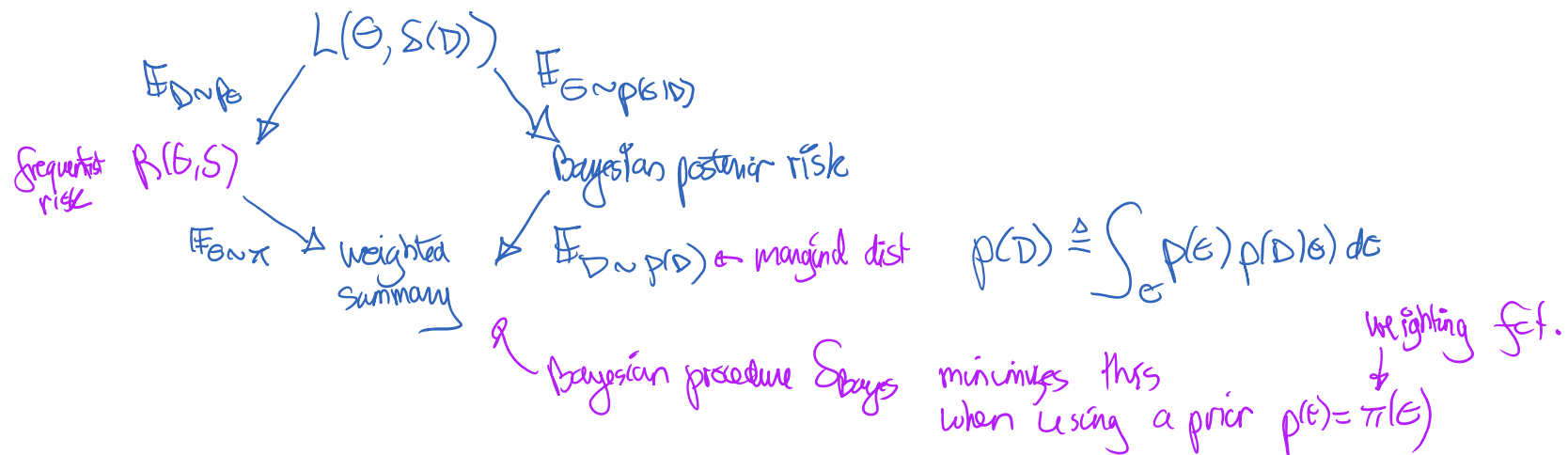
Bayesian optimal action:

$$\delta_{\text{Bayes}}(D) \triangleq \arg\min_{a \in \mathcal{A}} R_B(a|D)$$

example: if $f_\theta = \theta$ ("estimation")

$$L(\theta, a) = \|\theta - a\|^2$$

then (exercise) $S_{\text{Bayes}}(D) = \mathbb{E}[\theta | D]$ (posterior mean)



Examples of estimators: $S: \mathcal{D} \rightarrow \Theta$

1) • MLE

2) • MAP

3) • method of moments

idea: find an injective mapping from Θ to "moments" $\begin{matrix} \mathbb{E} X \\ \mathbb{E} X^2 \end{matrix}$

and subjective on "possible moments"

and then invert it from empirical moments

$$\begin{aligned} \hat{\mathbb{E}}[x] &\triangleq \frac{1}{n} \sum_{i=1}^n x_i \\ \hat{\mathbb{E}}[x^2] & \\ \dots \end{aligned}$$

example: for Gaussian $X \sim N(\mu, \sigma^2)$

$$\mathbb{E}[x] = \mu$$

$$\mathbb{E}[x^2] = \sigma^2 + \mu^2$$

$$f(\mu, \sigma^2) \triangleq \begin{pmatrix} \mathbb{E}[x] \\ \mathbb{E}[x^2] \end{pmatrix}$$

$$\begin{pmatrix} \hat{\mu} \\ \hat{\sigma}^2 \end{pmatrix} \triangleq f^{-1} \begin{pmatrix} \hat{\mathbb{E}}[x] \\ \hat{\mathbb{E}}[x^2] \end{pmatrix}$$

(here, this estimator is same as MLE)
[property of exponential family]

⊗ this is useful for latent variable models (e.g. mixture of Gaussians)
("spectral methods" e.g.)

4) in the context of prediction $\mathcal{A} = \{f: \mathcal{X} \rightarrow \mathcal{Y}\}$ \mathcal{X} input space
 \mathcal{Y} output "

example of $\mathcal{S}: \mathcal{D} \rightarrow \mathcal{A}$

is using empirical "risk" minimization (ERM)

↳ "various risk" i.e. generalization error

recall: $L(p, f) = \mathbb{E}_{(x, y) \sim p} [l(y, f(x))]$

replace with

$$\mathbb{E} [\ell(y, f(x))] = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i))$$

$$\hat{f}_{ERM} = \underset{f \in \mathcal{H}}{\operatorname{argmin}} \mathbb{E}[\ell(Y, f(x))]$$

← hypothesis class

14h34

bias-variance decomposition:

id setting $\mathcal{D} = \mathcal{P}_G^{\otimes n}$

$$D_n = (X_i)_{i=1}^n$$

$$X_i \text{ iid } p_0$$

rotation; $\hat{S}_n = S_n(D_n)$

- ↳ highlight dependence on n

Study $R(E, S_n)$ as a function of n

in particular, would take $R(G, S_n) \xrightarrow{n \rightarrow \infty} R(G, S_\infty)$

in particular, would like $R(\theta, \delta_n) \rightarrow R(\theta, \delta_\theta)$
 "consistency" \uparrow
 $\arg\min_{\delta} R(\theta, \delta)$

* for estimation, typical loss: squared loss $L(\theta, \delta_n(D)) = \|\theta - \delta_n(D)\|_2^2$

standard statistical consistency $\hat{\theta}_n \xrightarrow{P} \theta$
 \uparrow "in probability"

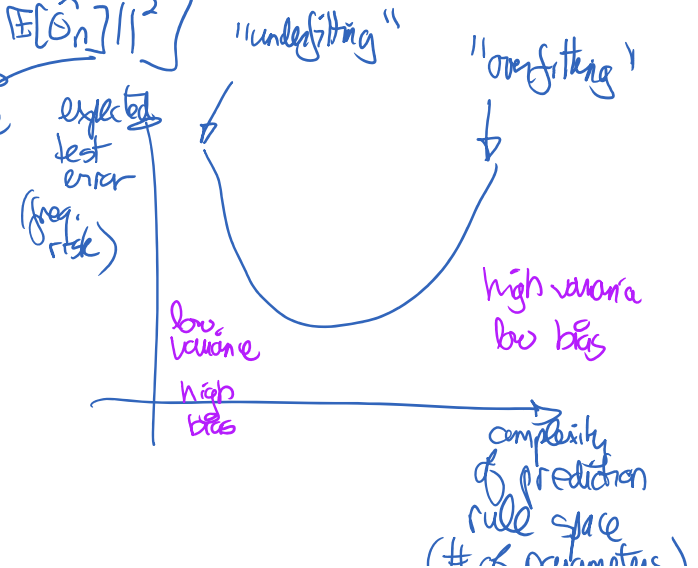
i.e. $\forall \epsilon > 0, P\{\|\hat{\theta}_n - \theta\| > \epsilon\} \xrightarrow{n \rightarrow \infty} 0$

\uparrow randomness is from \underline{D}_n

$$\hat{\theta}_n = \delta_n(D_n)$$

* last time, $R(\theta, \delta_n) = \mathbb{E}_{D_n} [\|\theta - \hat{\theta}_n\|_2^2]$

$$= \underbrace{\|\theta - \mathbb{E}[\hat{\theta}_n]\|_2^2}_{\text{bias}^2} + \underbrace{\mathbb{E}[\|\hat{\theta}_n - \mathbb{E}[\hat{\theta}_n]\|_2^2]}_{\text{variance}}$$



James-Stein estimator:

for estimating the mean of $N(\vec{\mu}, \sigma^2 I)$

SJS is biased, but lower variance than MLE

SJS is admissible \Rightarrow $d \geq 2$

OJS is bias, but lower variance than MLE

δ_{JS} actually strictly dominates δ_{MLE} for $d \geq 3$

of reduction
rule space
(# of parameters)

$$\text{i.e. } R(\theta, \delta_{JS}) \leq R(\theta, \delta_{MLE}) \quad \forall \theta \text{ dimension } (p)$$

$$\text{and } \exists \theta \text{ s.t. } R(\theta, \delta_{JS}) < R(\theta, \delta_{MLE})$$

\Rightarrow MLE is sometimes inadmissible

in assignment, "consistency" mean $R(\theta, \delta_n) \xrightarrow{n \rightarrow \infty} 0$ ($E \|\theta - \hat{\theta}_n\|^2 \xrightarrow{n \rightarrow \infty} 0$)
"convergence in L_2 "

by bias-variance decomposition:

[it turns out that L_2 -convergence \Rightarrow convergence in probability
i.e. $\hat{\theta}_n \xrightarrow{P} \theta$]

$$\left. \begin{array}{l} \text{bias}(\delta_n) \xrightarrow{n \rightarrow \infty} 0 \\ \text{variance}(\delta_n) \xrightarrow{n \rightarrow \infty} 0 \end{array} \right\} \Rightarrow R(\theta, \delta_n) \rightarrow 0 \Rightarrow \text{consistency} \quad \hat{\theta}_n \xrightarrow{P} \theta$$

properties (asymptotic) of MLE:

under regularity conditions on $\Theta \ni p(x; \theta)$

a) $\hat{\theta}_n \xrightarrow{P} \theta$ "consistent" "in distribution"

b) CLT: $\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{d} N(0, \underbrace{I(\theta)^{-1}}_{\text{information matrix}})$
central limit theorem (iid data)

theorem (iid data)

information matrix

c) asymptotically optimal

(Cramer-Rao lower bound)

ie. it has minimal asymptotic variance
among all "reasonable" estimators
↳ consistent

d) invariance: MLE is preserved under reparameterization

suppose have bijection $f: \Theta \rightarrow \Phi$

$$\text{then } f(\hat{\theta}) = \hat{f}(\hat{\theta})$$

* if not a bijection, can generalize MLE with "profile likelihood"

suppose $g: \Theta \rightarrow \Lambda$ profile likelihood $L(\eta) \triangleq \max_{\theta: \eta=g(\theta)} p(\text{data}; \theta)$

$$\text{define } \hat{\eta}_{MLE} \triangleq \arg \max_{\eta \in g(\Theta)} L(\eta)$$

then we have

$$\hat{\eta}_{MLE} = g(\hat{\theta}_{MLE})$$

"plug in" estimator

example: $\hat{(\sigma^2)} = (\hat{\sigma})^2$

$$\hat{\sin \sigma^2} = \sin \hat{\sigma}^2$$