

Lecture 6 - statistical decision theory

Thursday, September 21, 2023 2:26 PM

today: Statistical decision theory
evaluation of estimators

statistical decision theory - formal setup

- 1) • a random observation $D \sim P$ unknown distribution which models the world / phenomenon
(often p_θ)

- 2) • action space A

- 3) • loss $L(P, a) =$ statistical loss of doing action $a \in A$ when the world is P } describe the goal / task
 ↳ often write $L(G, a)$ if we have a parametric model of world i.e. P is a pdf / pmf p_θ for some $\theta \in \Theta$

- $\delta: D \rightarrow A$ "decision rule"

examples: a) parameter estimation:

$$\mathcal{P} = \Theta \text{ for parametric family } P_\theta$$

δ is a parameter estimator from data

$$D = (X_1, \dots, X_n)$$

$$\text{typical loss } L(G, a) = \|G - a\|_2^2$$

"squared loss"

[usually $X_i \stackrel{iid}{\sim} p_\theta$]

but another loss is $k h(\rho_\theta || \rho_a)$

b) $A = \{0, 1\}$; this is hypothesis testing

δ describes a statistical test

$$\text{loss} \rightarrow \text{usually 0-1 loss } L(\delta, a) = \underbrace{\mathbb{1}_{\{\delta \neq a\}}}_{\mathbb{1}_{\{\delta \neq a\}^c} = \Theta \setminus \{\delta\}} \quad (\triangleq \mathbb{1}_{\{\delta \neq a\}})$$

$$\mathbb{1}_A(u) = \begin{cases} 1 & \text{if } u \in A \\ 0 & \text{o.w.} \end{cases}$$

\rightarrow need to minimize $\mathbb{E}_{\theta} L(\delta, a)$ over δ . $\Pr :=$...

$$g_{\theta}^{(k)} = \text{ReLU}_{\theta}$$

c) prediction in ML: learn a prediction fct. in supervised learning
(function estimation)

$$\text{here } D = (X_i, Y_i)_{i=1}^n$$

$X_i \in \mathcal{X}$ (input space)
 $Y_i \in \mathcal{Y}$ (output space)

$\mathcal{Y} = \{0, 1\} \rightarrow \text{classification}$

P_θ gives joint on (X, Y)

$\mathcal{Y} = \mathbb{R} \rightarrow \text{regression}$

$$\mathcal{P} = \mathcal{Y}^X \text{ (set of fct's from } X \text{ to } \mathcal{Y})$$

$$|\mathcal{P}| = |\mathcal{Y}|^{|\mathcal{X}|}$$

$$D \sim P \text{ where } P = P_1 \otimes P_2 \otimes \dots \otimes P_n \text{ (n times)}$$

$$P(x_1, \dots, x_n) = \prod_{i=1}^n p_i(x_i)$$

in ML

$$L(P_\theta, f) \triangleq \mathbb{E}_{P_\theta} [l(Y, f(X))]$$

(x, y) ~ P_θ

"prediction loss"

"generalization error"

"classification error"

in ML, is often call the "risk"

↳ e.g. classification
 $l(Y, f(X)) = \mathbb{1}\{Y \neq f(X)\}$

Sunior calls it "Vapnik risk"

to distinguish from frequentist risk

$$\mathbb{E}_D [L(P_\theta, S(D))]$$

* decision rule

$$f = g(D)$$

prediction fct/
classifier/
etc.

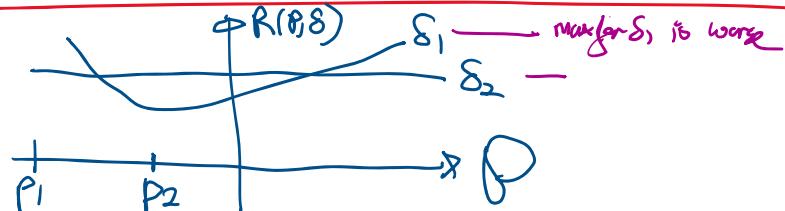
think
cross-validation e.g.

Comparing procedures?

S_1 vs. S_2

"risk profiles"

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



* transform to scalar

- "minimax" analysis : $\max_{D \in \mathcal{D}} R(P, S)$ "worst case"

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

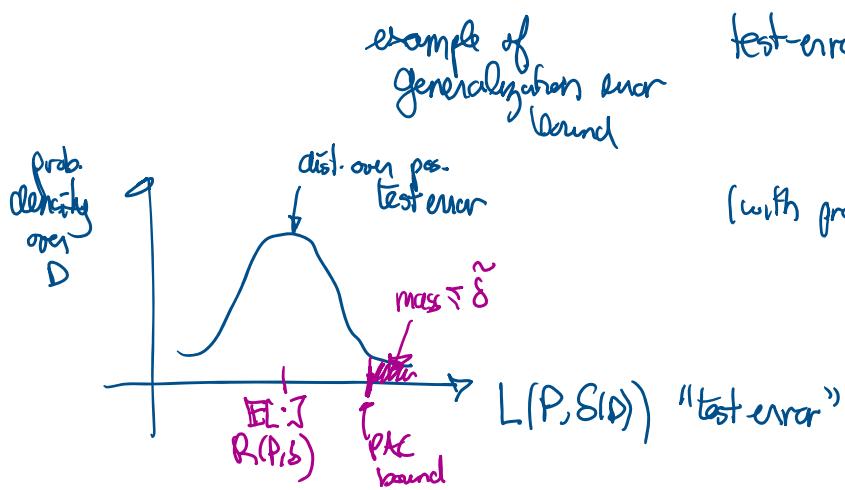
- weighted average $\int_{\mathcal{H}} R(f_\theta, \delta) \pi(\theta) d\theta$ (kind of a Bayesian feel)

|Sh3)

PAC theory vs frequentist risk:

in ML, usually they look at tail bound for dist. of $L(P, \delta(D))$ where D is random

PAC theory
↳ "probably approx. correct"
 $P\{L(P, \delta(D)) > \text{stuff}\} \leq \tilde{\delta}$



$$\text{test-error}(\hat{f}) \leq \text{train-error}(\hat{f}) + \frac{1}{\sqrt{n}} \underbrace{\sqrt{\text{Complexity}(\hat{f})}}_{\text{Complexity}(\hat{f})} + \log \frac{1}{\tilde{\delta}}$$

(with prob $\geq 1 - \tilde{\delta}$ on D)

Bayesian decision theory

→ condition on data D

Bayesian posterior risk $R_B(a|D) = \int_{\mathcal{H}} L(\theta, a) p(\theta|D) d\theta$

posterior over "possible worlds"

Bayesian optimal action: $S_{\text{Bayes}}(D)$

$$\triangleq \underset{a \in \mathcal{A}}{\operatorname{argmin}} R_B(a|D)$$

$$\propto p(\theta) p(D|\theta)$$

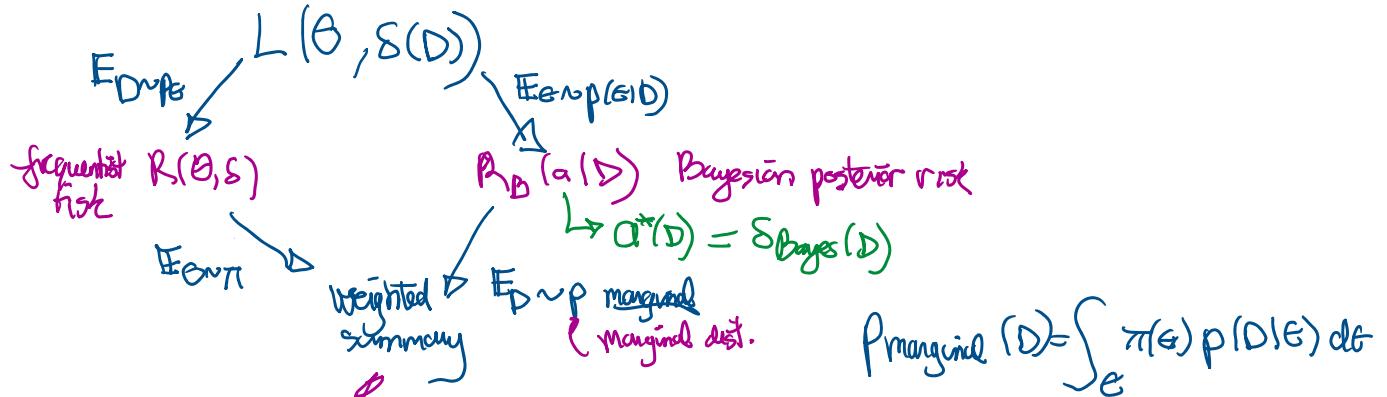
example: if $\mathcal{A} = \mathbb{R}$ ("estimation")

$$L(\theta, a) = \| \theta - a \|_2^2$$

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

but if we $L(\theta, a) = (\theta - a)^2$ (ID)

then $\hat{\theta}_{\text{Bayes}}(D) = \underline{\text{posterior median}}$



Bayesian procedure S_{Bayes} minimizes the weighted summary among all S 's
when using prior $\pi(\theta)$ as the weighting fact. i.e. $\pi(\theta) = p(\theta)$

Examples of estimators: $\delta: \mathcal{D} \rightarrow \mathbb{H}$

- 1) MLE
- 2) MAP
- 3) method of moments (MoM)

idea: find an injective mapping from \mathbb{H} to "moments" of r.v.

$$\mathbb{E}X, \mathbb{E}X^2, \dots$$

and then insert it from empirical moments to get $\hat{\theta}$

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

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

$$\begin{aligned}\mathbb{E}X &= \mu \\ \mathbb{E}X^2 &= \sigma^2 + \mu^2\end{aligned}$$

$$f(\mu, \sigma^2) \triangleq \begin{pmatrix} \mu \\ \sigma^2 + \mu^2 \end{pmatrix}$$

$$\begin{pmatrix} \hat{\mu} \\ \hat{\sigma}^2 \end{pmatrix} \triangleq f\left(\begin{pmatrix} \hat{\mathbb{E}}[x] \\ \hat{\mathbb{E}}[x^2] \end{pmatrix}\right) = \begin{pmatrix} \hat{\mathbb{E}}[x] \\ \hat{\mathbb{E}}[x^2] - (\hat{\mathbb{E}}[x])^2 \end{pmatrix}$$

(here, this estimator is same as MLE)
 [general property in exponential family → see later]

④ MoM is quite used for latent variable models

↳ ("spectral methods" e.g.)



4) prediction example: $\mathcal{F} = \{f: \mathcal{X} \rightarrow \mathcal{Y}\}$

\mathcal{X} = input space
 \mathcal{Y} = output space

Example of $S: \mathcal{D} \rightarrow \mathcal{F}$

is using empirical "risk" minimization (ERM)

↳ Vapnik risk i.e. generalization/test error

$$\text{i.e. } L(P, f) = \mathbb{E}_{(x,y) \sim P} [l(y, f(x))]$$

$$\text{replace this with } \hat{\mathbb{E}} [l(y, f(x))] = \frac{1}{n} \sum_{i=1}^n l(y_i, f(x_i))$$

$$\text{ERM} = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \hat{\mathbb{E}} [l(y, f(x))]$$

hypothesis class

James-Stein estimator:

estimator to estimate the mean of $N(\vec{\mu}, \sigma^2 I)$ ← d independent Gaussian variables
 δ_{JS} is biased, but much lower variance than MLE

$$\begin{aligned} \text{recall bias-variance decomposition: } R(G, \hat{\theta}) &= \mathbb{E}[\|\hat{\theta} - \theta\|_2^2] \\ &= \underbrace{\|\mathbb{E}\hat{\theta} - \theta\|_2^2}_{\text{bias}^2} + \underbrace{\mathbb{E}[\|\hat{\theta} - \mathbb{E}\hat{\theta}\|_2^2]}_{\text{variance}} \end{aligned}$$

δ_{JS} actually strictly dominates SMLE

for $d \geq 3$
 ↑ dimension of $\vec{\mu}$

i.e. $R(G, \delta_{JS}) \leq R(G, \text{SMLE}) + G$

for $d \geq 3$
↑ dimension of $\vec{\mu}$

i.e. $R(\theta, \delta_{JS}) \leq R(\theta, \delta_{MLE}) + G$
and $\exists \theta$ s.t. $R(\theta, \delta_{JS}) < R(\theta, \delta_{MLE})$

→ MLE is inadmissible in this case [note $n=1$ here]

(can interpret the δ_{JS} as an "empirical" Bayesia method) ↪ $d \geq 3$