

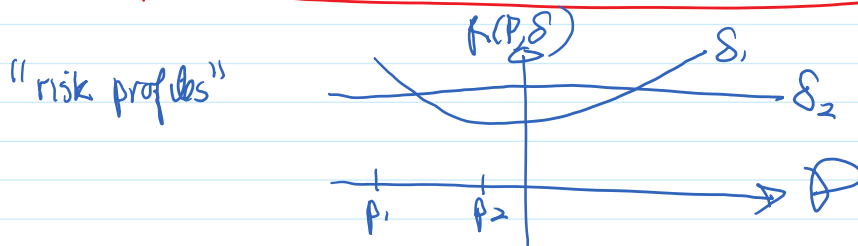
today: : statistical decision theory and properties of estimators

comparing procedures?

δ_1 vs. δ_2

(frequentist) risk $R(P, \delta) \equiv \mathbb{E}_{D \sim P} [L(P, \delta(D))]$

(think cross-validation)



* transform to scalars

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

• weighted average $\int_{\mathcal{D}} R(\theta, \delta) \pi(\theta) d\theta$ (kind of a Bayesian feel)
 (weighting)

PAC theory vs. frequentist risk:

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

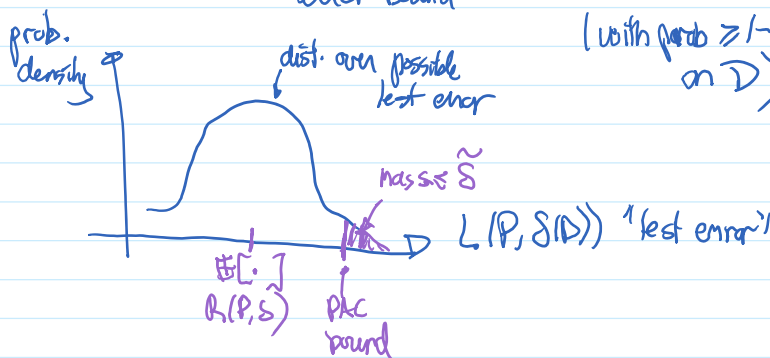
PAC theory
 \hookrightarrow "probably approx. correct"

$$P\{L(P, \delta(D)) \geq \text{stuff}\} \leq \tilde{\delta}$$

example of generalization error bound

$$\text{test error}(\hat{f}) \leq \text{train error}(\hat{f}) + \frac{1}{\sqrt{n}} \sqrt{\text{complexity}(\hat{f}) + \log\left(\frac{1}{\tilde{\delta}}\right)}$$

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



Bayesian decision theory

\rightarrow condition on data D

Bayesian posterior risk

$$R_D(a|D) = \int_{\mathcal{D}} L(\theta, a) p(\theta|D) d\theta$$

posterior over "possible worlds" $\propto N(\theta) \pi(D|\theta)$

Bayesian optimal action: $S_{\text{Bayes}}(D) \stackrel{\text{argmin}}{a \in \mathcal{A}} R_{\theta}(a|D)$ (4) posterior over "possible worlds" $\propto p(\theta) p(D|\theta)$

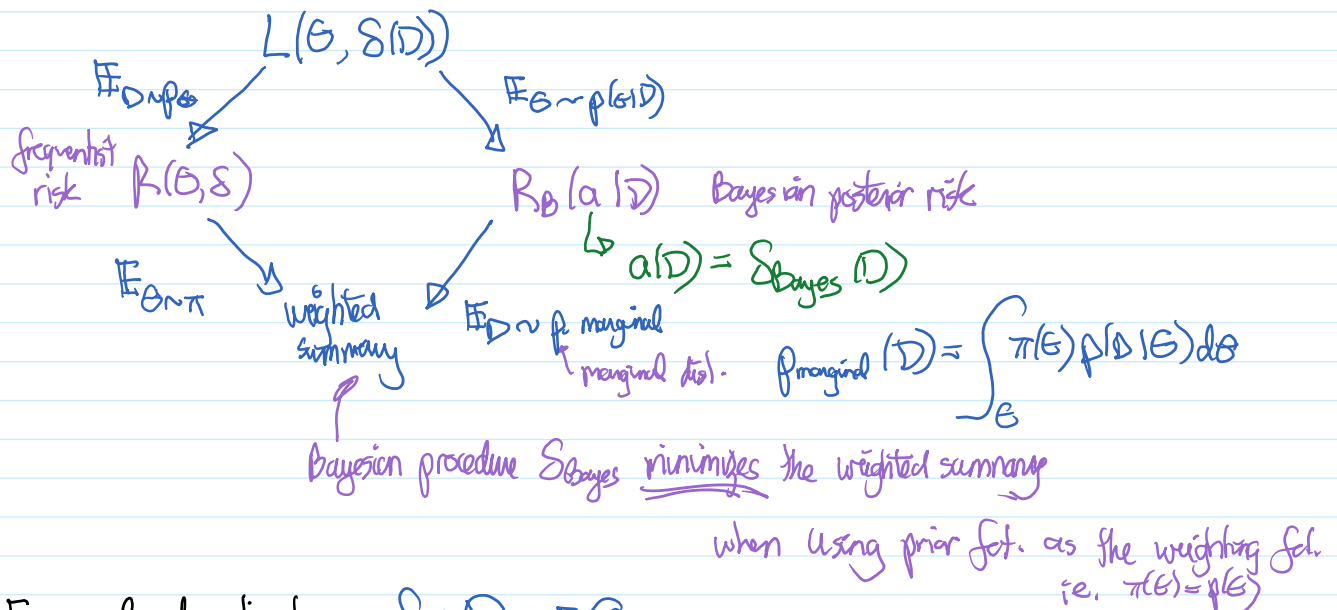
example: if $\mathcal{A} = \mathbb{R}$ ("estimation")

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

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

but if use $L(\theta, a) = |\theta - a|$ (1D)

then $S_{\text{Bayes}}(D) = \underline{\text{posterior median}}$



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

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

idea: find an injective mapping from \mathcal{B} to "moments" of R.V.

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

and then invert it from empirical moments to get θ

$$\hat{\mathbb{E}}[X] \triangleq \frac{1}{n} \sum_{i=1}^n x_i$$

$$\hat{\mathbb{E}}[X^2] \triangleq \frac{1}{n} \sum_{i=1}^n x_i^2 \dots$$

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

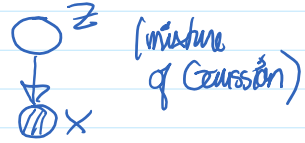
$$\begin{aligned} \mathbb{E}X &= \mu \\ \mathbb{E}X^2 &= \sigma^2 + \mu^2 \end{aligned}$$

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

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

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

⊛ MoM is quite used for latent variable models
↳ ("spectral methods" e.g.)



17h30

4) prediction example $\mathcal{F} = \{f: X \rightarrow \mathcal{Y}\}$ $X \leftarrow$ input space
 \mathcal{Y} ← output "

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

is using empirical "risk" minimization (ERM)

↳ Vapnik risk i.e. generalization error / test error

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

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

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

↳ hypothesis class

James - Stein estimator:

estimator to estimate the mean of $N(\vec{\mu}, \sigma^2 I)$ ← d independent Gaussian variables $X_i \sim N(\mu_i, \sigma^2)$

\mathcal{S}_{JS} is biased, but much lower variance than MLE

recall bias-variance decomposition $R(\mathcal{S}, \hat{\mathcal{S}}) = \mathbb{E} \|\mathcal{S} - \hat{\mathcal{S}}\|_2^2$

$$= \underbrace{\mathbb{E} \|\mathcal{S} - \mathbb{E}(\hat{\mathcal{S}})\|_2^2}_{\text{bias}^2} + \underbrace{\mathbb{E} \|\mathbb{E}(\hat{\mathcal{S}}) - \mathbb{E}(\mathcal{S})\|_2^2}_{\text{variance}}$$

\mathcal{S}_{JS} actually strictly dominates \mathcal{S}_{MLE}
for $d \geq 3$
↳ dimension of $\vec{\mu}$

ie. $R(\hat{\theta}, S_{JS}) \leq R(\hat{\theta}, S_{MLE}) \forall \theta$

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

→ MLE is inadmissible in this case [note $n=1$ here] $\leftarrow d \geq 0$

[can interpret the S_{JS} as an 'empirical' Bayesian method]

(asymptotic properties of MLE):

under suitable regularity conditions on $\Theta \ni \theta, p(x; \theta)$ $\hat{\theta}_n = \operatorname{argmax}_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \log p(x_i; \theta)$

a) $\hat{\theta}_n \xrightarrow{P} \theta$ 'consistent'

b) CLT (central limit theorem) $\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{dist.} N(0, I(\theta)^{-1/2})^n \sim P_{\theta}^{\otimes n}$
 information matrix

c) asymptotically optimal (Cramer-Rao lower bound)

ie. it has minimal asymptotic scaled variance among all 'reasonable' estimators

d) invariance: MLE is preserved under reparameterization

suppose have a bijection $f: \Theta \rightarrow \Theta'$

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

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

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

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

suppose $g: \Theta \rightarrow \Lambda$

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

define $\hat{\eta}_{MLE} \triangleq \operatorname{argmax}_{\eta \in \Lambda} L(\eta)$

then we have $\boxed{\hat{\eta}_{MLE} = g(\hat{\theta}_{MLE})}$

'plug in estimator'

$N(\mu, \sigma^2)$
 e.g. $g(\mu) = \mu^2$

prediction want to learn prediction fct. $h: X \rightarrow \mathcal{Y}$ $\mathcal{Y} = \{0, 1\} \rightarrow$ binary classification

prediction

want to learn prediction fct. $h: X \rightarrow \mathcal{Y}$

$\mathcal{Y} = \{0, 1\} \rightarrow$ binary classification
 $\mathcal{Y} = \{0, \dots, k-1\} \rightarrow$ multiclass classif.



$x \in \mathbb{R}^d$

"prediction model" model over X

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

$p(x, y) = p(y|x) p(x)$

$= \underbrace{p(x|y)}_{\text{"class conditional"}} \underbrace{p(y)}_{\text{prior over class}}$

"generative perspective" (in context of classification) \rightarrow model $p(x)$ as well

"conditional perspective"

\rightarrow only model $p(y|x)$

"more discriminative"

\leftarrow traditionally called "discriminative"

generative	conditional	"fully discriminative"
model $p_{\theta}(x, y)$ MLE	model $p_{\theta}(y x)$ max conditional likelihood	model $h_{\theta}: X \rightarrow \mathcal{Y}$ (not nec. derived from $p(y x)$) reg. ERM; etc. $\frac{1}{n} \sum_{i=1}^n \ell(y^{(i)}, h_{\theta}(x^{(i)}))$
more assumptions \rightarrow less robust for prediction		less assumptions more robust

$\hat{h}(x) \hat{=} \arg \min_{\tilde{y} \in \mathcal{Y}} \sum_y p_{\theta}^*(y|x) \ell(y, \tilde{y})$

if $\ell(y, \tilde{y}) = \mathbb{1}\{y \neq \tilde{y}\}$ (0-1 loss) then $\hat{h}(x) = \arg \max_{\tilde{y} \in \mathcal{Y}} p_{\theta}(\tilde{y}|x)$