

today: finish decision theory
more properties of estimators

PAC theory vs. frequentist risk

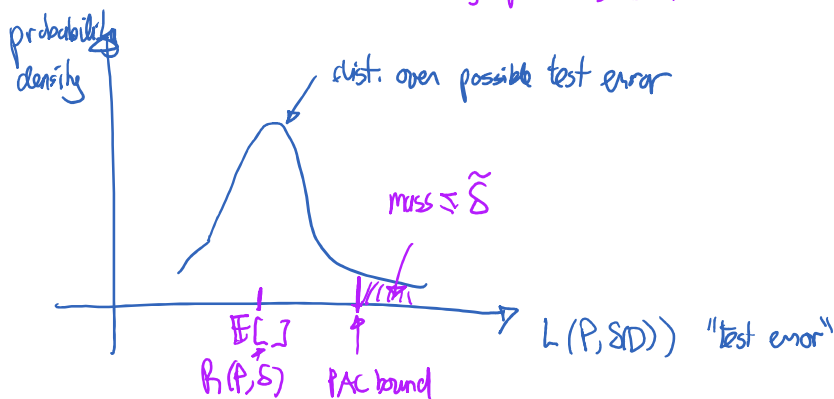
in ML, usually they look tail bounds for dist. $L(P, S(D))$ where D is random

PAC theory
↳ "probably approx. correct"

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

↓
bound

↑ small #
"with high prob." statement



Bayesian decision theory:

→ condition on data D

Bayesian posterior risk

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

Ⓡ

posterior over "possible worlds"
or $p(\theta) p(D|\theta)$

Bayesian optimal action: $S_{\text{Bayes}}(D) \triangleq \arg\min_{a \in A} R_B(a|D)$

example: if $A = \Theta$ ("estimation")

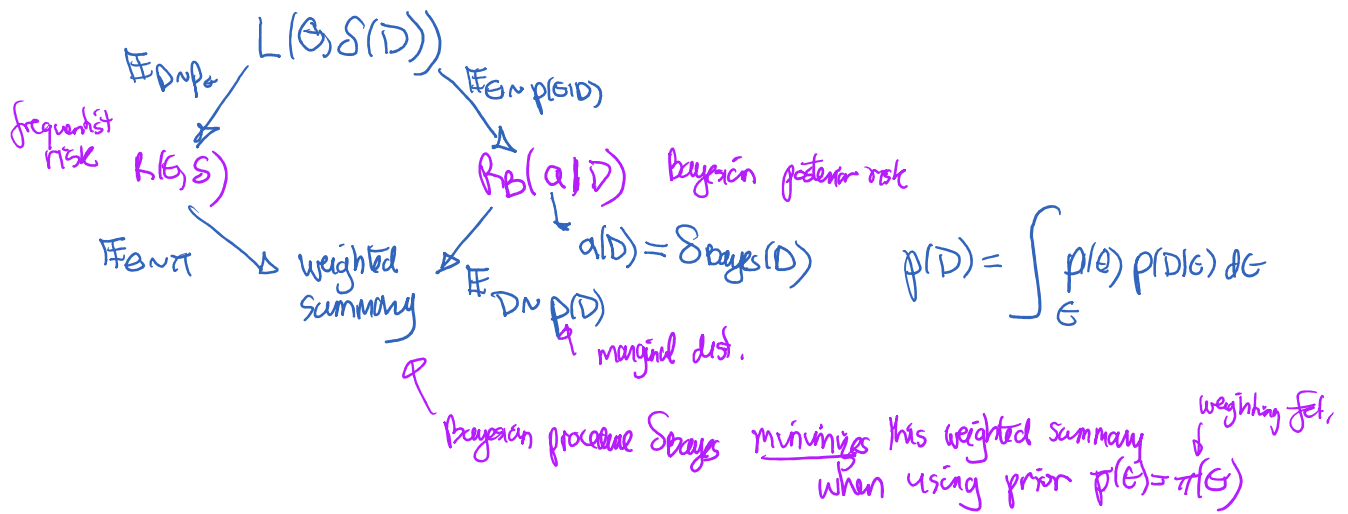
$$L(\theta, a) = \|\theta - a\|_2^2$$

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

but if $L(\theta, a) = |\theta - a|$ (4D)

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

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



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

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

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

$$E[X], E[X^2], \text{ etc. }$$

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

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

$$\hat{E}[X^2]$$

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

$$E[X] = \mu$$

$$E[X^2] = \sigma^2 + \mu^2$$

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

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

\mathbb{R}^L

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

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



e.g. mixture of Gaussian

15h10

4) in the context of prediction $\mathcal{F} = \{f: \mathcal{X} \rightarrow \mathcal{Y}\}$
 \mathcal{X} ← 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

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

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

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

↑ hypothesis class

James-Stein estimator:

estimator to estimate the mean of $N(\vec{\mu}, \sigma^2 I)$ ← d independent Gaussian variables with means μ_i & variance σ^2

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

recall the bias-variance decomposition for squared loss

$$R(\theta, \hat{\theta}) = \mathbb{E} \|\hat{\theta} - \theta\|^2 = \underbrace{\mathbb{E} \|\hat{\theta} - \theta\|^2}_{\text{bias}} + \underbrace{\mathbb{E} \|\hat{\theta} - \mathbb{E} \hat{\theta}\|^2}_{\text{variance}}$$

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

$$\text{i.e. } R(\theta, \mathcal{S}_{JS}) \leq R(\theta, \mathcal{S}_{MLE}) \quad \forall \theta$$

and $\exists \theta$ s.t. $R(\theta, \mathcal{S}_{JS}) < R(\theta, \mathcal{S}_{MLE})$

→ MLE is inadmissible in this case

(can interpret the SJS as an "empirical" Bayesian method)

Properties (asymptotic) of MLE:

under ^{jointly} regularity conditions on $\Theta \ni p(x; \theta)$ $\hat{\theta}_n = \underset{\theta \in \Theta}{\text{argmax}} \left(\sum_{i=1}^n \log p(x_i; \theta) \right)$

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

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

c) asymptotically optimal (Cramer-Rao lower bound)

ie. it has minimal asymptotic variance among all "reasonable" estimators
 \hookrightarrow consistent
 \dots

d) uniqueness: MLE is preserved under reparameterization

suppose have a bijection $f: \Theta \rightarrow \Theta'$
 $\hat{\theta}' = f(\hat{\theta})$

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

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

suppose $g: \Theta \rightarrow \mathcal{L}$

profile likelihood $\hat{L}(m) \stackrel{\text{def}}{=} \max_{\theta: m=g(\theta)} p(\text{data}; \theta)$

define $\hat{m}_{MLE} \stackrel{\text{def}}{=} \underset{m \in \mathcal{L}}{\text{argmax}} L(m)$

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

'plug-in' estimator