

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  $p_\theta$ )

- 2) action space  $A$

- 3) loss  $L(P, a)$  = statistical loss of doing action  $a \in A$  } describe the goal/task  
when world is  $P$

↳ often write  $L(\theta, a)$  if we have a parametric model of world  
ie.  $P$  has a pdf/pmf  $p_\theta$  for some  $\theta \in \Theta$

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

examples: a) parameter estimation

$A = \Theta$  for some parametric family  $P_\Theta$

$S$  is a parameter estimator from data  $D = (x_1, x_2, \dots, x_n)$   
[usually  $x_i \stackrel{iid}{\sim} p_\theta$ ]

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

"squared loss"

but another loss is  $k L(p_\theta \| p_a)$

b)  $A = \{\theta, \bar{A}\}$ ; this is hypothesis testing

$S$  describes a statistical test

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

loss → usually 0-1 loss  $L(\theta, a) = I_{\{\theta \neq a\}}$  ( $\triangleq I_{\{\theta \neq a\}}$ )

$$I_{\{\theta \neq a\}} = \Theta \setminus \{a\}$$

→ related to  $\text{fun I} \in \text{fun II}$  aware

$\{\theta\}^C$  "whole.."

$$T_{\Sigma} = \{ \cdot \} \cap \Sigma$$

→ related to type I & type II errors

$$\begin{matrix} \{ \cdot \}^C \\ A^C = B \setminus A \end{matrix}$$

"whole world"

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

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

$X_i \in \mathcal{X}$	(input space)	$\mathcal{Y} = \{0, 1\} \rightarrow$ classification
$Y_i \in \mathcal{Y}$	(output space)	$\mathcal{Y} = \mathbb{R} \rightarrow$ regression

Let  $p_g$  be joint on  $(X, Y)$

$\mathcal{F} = \mathcal{Y}^{\mathcal{X}}$  (set of functions from  $\mathcal{X}$  to  $\mathcal{Y}$ )

$D \sim P$  where  $P = p_0 \otimes p_1 \otimes p_2 \cdots \otimes p_n$

$$p(\tilde{x}_1, \dots, \tilde{x}_n) = \prod_{i=1}^n p_i(\tilde{x}_i)$$

in ML

$$L(p_g, f) \triangleq \mathbb{E}_{(x,y) \sim p_g} [l(Y, f(X))]$$

[ "prediction loss" ]

"generalization error"

"classification error"

in ML, is often called the "risk"

Suvor calls it the "Vapnik risk"

to distinguish it from frequentist risk  $\mathbb{E}_D [L(p_g, S(D))]$

↳ classification

$$l(Y, f(X)) = \mathbb{1}\{Y \neq f(X)\}$$

0-1 error

↗ decision rule

$$\hat{f} = S(D)$$

"learning alg"

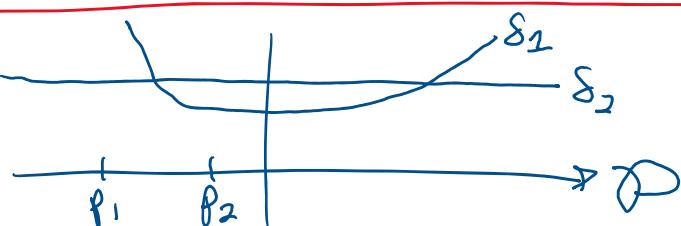
prediction fct. classifier/regressor/etc..

comparing procedures?

$\delta_1$  vs.  $\delta_2$

"risk profiles"

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



\* transform to scalar

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

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

13h34

PAC theory vs. frequentist risk:

in ML, usually they look at test bound for dist. on  $L(P, S(D))$  where  $D$  is random

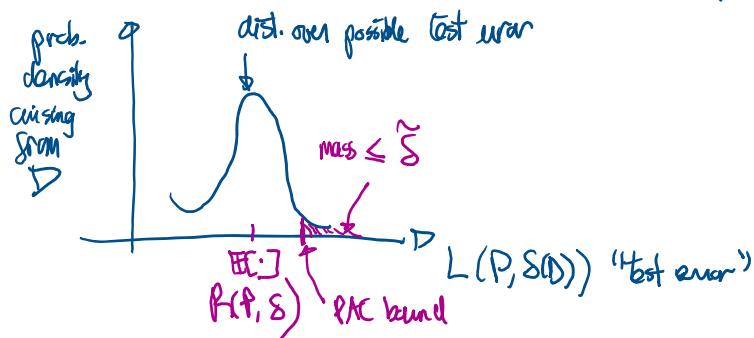
PAC theory  
↳ "probably approx. correct"

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

example of generalization error bound

$$\text{test-err}(f) \leq \text{train-err}(f) + \frac{1}{\sqrt{n}} \sqrt{\text{complexity}(f) + \log \frac{2}{\delta}}$$

(with prob  $\geq 1 - \tilde{\delta}$ )  
↳ randomness of  $D$



Bayesian decision theory

→ condition on data  $D$

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

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

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

posterior over "possible worlds"  
 $\propto p(\epsilon) p(D | \epsilon)$

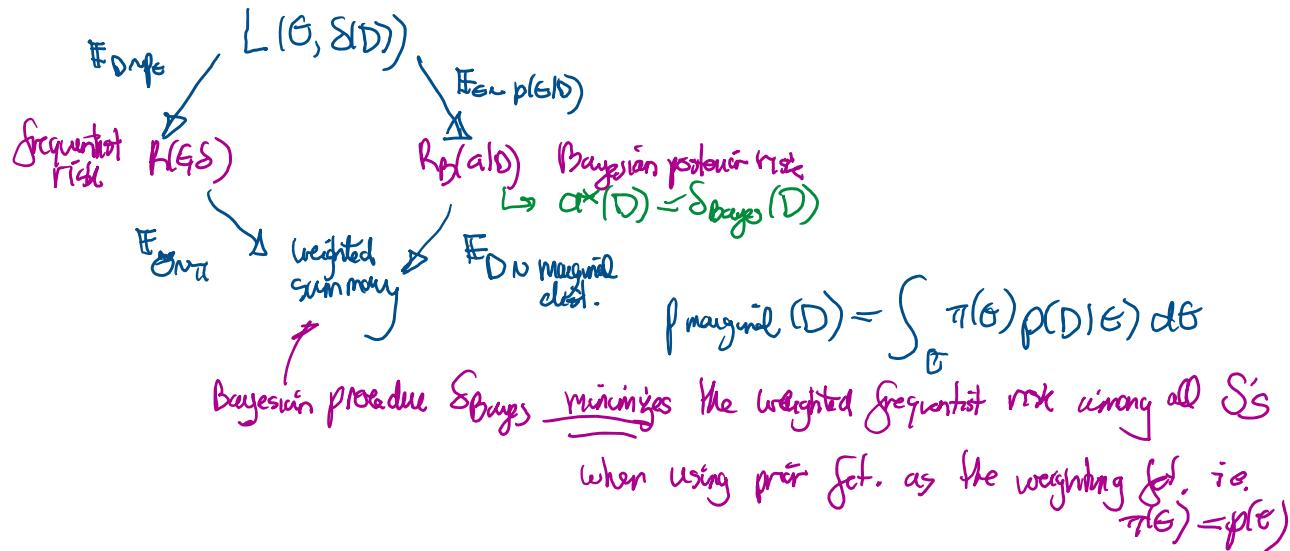
Example: if  $\mathcal{A} = \mathbb{R}$  "estimation"

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

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

but if we use  $L(G, a) = (G - a)^2$  (1D)

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



Examples of estimators:  $S: D \rightarrow \Theta$

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

Idea: find an injective mapping from  $\Theta$  to "moments" of RV.

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

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

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

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

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

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

(here, this estimator gives same answer as MLE)

[General property in exponential family  $\rightarrow$  see notes]

$\textcircled{*}$  MoM is quite useful for latent variable models

$\hookrightarrow$  ("spectral methods" e.g.)



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

$X \leftarrow$  input space

$\mathcal{Y} \leftarrow$  output space

example of  $\delta: D \rightarrow A$

is using empirical "risk" minimization (ERM)

$\hookrightarrow$  Vapnik risk i.e. test error

$$\text{i.e. } L(\rho, f) = \#_{(x_i, y_i) \in D} [l(y_i, f(x_i))]$$

$$\text{replace this with } \hookrightarrow E[l(Y, f(X))] = \sum_{i=1}^n l(Y_i, f(X_i))$$

$$\hat{f}_{\text{ERM}} = \underset{f \in \mathcal{H}}{\operatorname{argmin}} E[l(Y, f(X))]$$

$\hookleftarrow$  hypothesis class

James-Stein estimator:

estimator to estimate the mean of  $N(\bar{\mu}, \sigma^2 I)$   $\leftarrow$  d independent Gaussian variables

$\delta_{JS}$  is biased, but much lower variance than MLE

$X_i \sim N(\mu_i, \sigma^2)$

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

$\delta_{JS}$  actually strictly dominates SMLE

for  $d \geq 3$   
dimension of  $\mu$

$$\text{i.e. } R(\theta, \delta_{JS}) \leq R(\theta, \delta_{MLE}) + \epsilon$$

and for st.  $R(\theta, \delta_{JS}) < R(\theta, \delta_{MLE})$

$\rightarrow$  MLE is inadmissible in this case [note  $n=1$  here  $\hookleftarrow d \geq 3$ ]

From statement the L is not an "unbiased" method

(can interpret the  $S_{JS}$  as an "empirical" Bayesian method)

here  
 $d \geq 3$