

today: statistical decision theory  
properties of estimator

Statistical decision theory

formal setup:

- a random observation  $D \sim P$  (often  $P_\theta$ )   
 *unknown distribution which models the world/phenomenon*
- action space  $\mathcal{A}$
- loss  $L(P, a) =$  <sup>statistical</sup> loss of doing action  $a \in \mathcal{A}$  when the world is  $P$  } describe the goal/task
- ↳ often write  $L(\theta, a)$  if we have a parametric model of world i.e.  $P$  had pdf/pdf  $P_\theta$  for some  $\theta \in \Theta$
- $\delta: \mathcal{D} \rightarrow \mathcal{A}$  "decision rule"   
  $\uparrow$   $\Omega_D$

examples: a) parameter estimation:

$\mathcal{A} = \Theta$  for a parametric family  $P_\theta$

$\delta$  is a parameter estimator from data

$\mathcal{D} = (X_1, \dots, X_n)$    
 *typically*

typical loss  $L(\theta, a) = \|\theta - a\|_2^2$    
  $a \in \mathcal{A} = \Theta$    
 "Squared loss"

[usually,  $X_i \stackrel{iid}{\sim} P_\theta$ ]   
  $\theta$  unknown

$\hat{\theta} = \delta(D)$

but other losses are used, e.g.  $KL(P_\theta || P_a)$

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

$\delta$  describes a statistical test

loss  $\rightarrow$  usually 0-1 loss  $L(\theta, a) = \mathbb{1}\{\theta \neq a\}$

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

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

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

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

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

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

$P_\theta$  gives joint  $(X, Y)$

$D \sim P$  where  $P = P_{X_1} \otimes P_{X_2} \dots \otimes P_{X_n}$

$P_0$  gives joint  $(X, Y)$

risk  $\rightarrow$  (input space)  $\rightarrow$  risk regression

$D \sim P$  where  $P = P_0 \otimes P_1 \dots \otimes P_n$   
n times

$\mathcal{H} = \mathcal{Z}^X$  (set of fct.'s from  $X$  to  $\mathcal{Z}$ )

in machine learning

$L(P_0, f) \triangleq \mathbb{E}_{P_0} [l(Y, f(X))]$

"generalization error"

"classification error"

in ML, is often called the "risk"

$(X, Y) \sim P_0$  prediction loss

e.g. classification  $l(Y, f(X)) = \mathbb{1}\{Y \neq f(X)\}$   
0-1 error

Simon calls it "Bayesian risk" to distinguish it from frequentist risk

\* decision rule

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

prediction fct. / classifier / etc.

"learning algorithm"

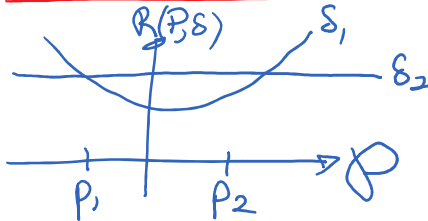
frequentist risk  $\mathbb{E}_D [L(P_0, \delta(D))]$

comparing procedure?

$\delta_1$  vs.  $\delta_2$

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

"risk profiles"



\* transform to scalar

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

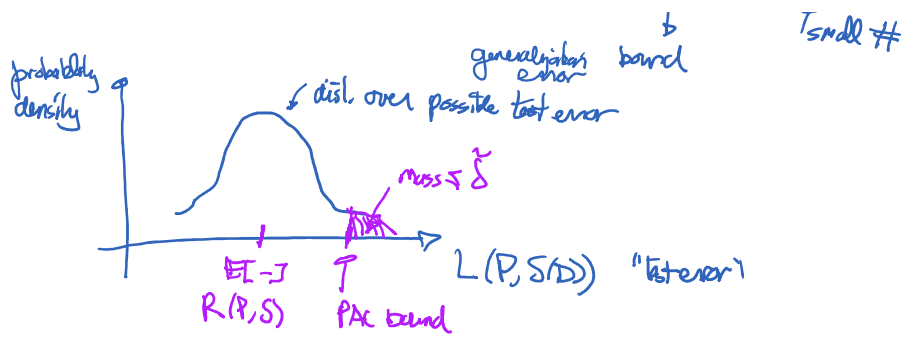
• weighted average  $\int_{\mathcal{H}} R(P, \delta) \pi(\delta) d\delta$  (kind of a Bayesian seed)  
weighting

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  $\rightarrow$  "probably approx. correct"

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



example of error bound  $\text{test error}(\hat{\xi}) \leq \text{train error}(\hat{\xi}) + \sqrt{\frac{\text{complexity}(\hat{\xi})}{n} + \log(\frac{1}{\delta})}$

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

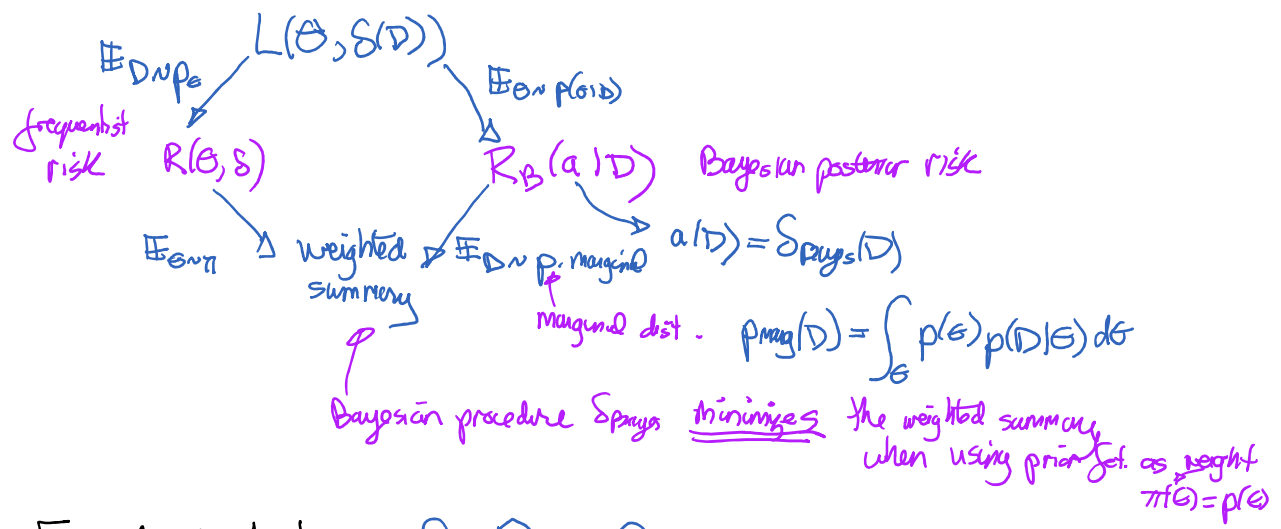
example: if  $\mathcal{A} = \Theta$  ("estimation")

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

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

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

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



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

$$\pi(\theta) = p(\theta)$$

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

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

idea: find an injective mapping from  $\mathcal{G}$  to "moments" of RVs

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

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

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

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

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

$$\mathbb{E}X = \mu$$

$$\mathbb{E}X^2 = \sigma^2 + \mu^2$$

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

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

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

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

