

today: theory basics

theory basics

decision theory setup

estimate  $h_w: X \rightarrow Y$

generalization error =  $L_P(w) \triangleq \mathbb{E}_{(x,y) \sim P} [l(y, h_w(x))]$

task obs

ultimate goal is to find  $w^* = \operatorname{argmin}_{w \in W} L_P(w)$

problem: do not know  $P$  ("true" distribution on  $(x,y)$ )

suppose  $(x^{(i)}, y^{(i)})_{i=1}^n \stackrel{\text{iid.}}{\sim} P \rightsquigarrow$  we could look at  
 $\triangleq D_n$  training dataset

$$\hat{L}_n(w) = \frac{1}{n} \sum_{i=1}^n l(y^{(i)}, h_w(x^{(i)}))$$

from statistics / prob. theory.

$$\hat{L}_n(w) \xrightarrow[n \rightarrow \infty]{\text{a.s.}} L_P(w) \quad \text{for each fixed } w \text{ (point wise)}$$

(LLN)

this is weaker than  $\sup_w |\hat{L}_n(w) - L_P(w)| \xrightarrow[n \rightarrow \infty]{\text{a.s.}} 0$

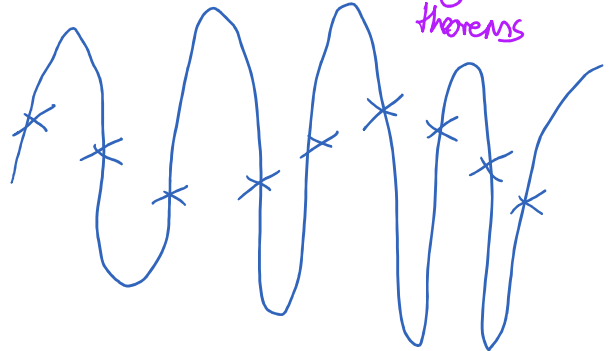
note: minimizing the training error gives no guarantee on  $L_P(\hat{w}_n)$  in general!

later no free lunch theorems

e.g. polynomial regression

for  $n$  points, can get zero training error with poly. of degree  $n-1$

$\Rightarrow$  "overfitting"



in learning theory:  $\rightarrow$  study properties of learning algo.

$\hat{w}_n$

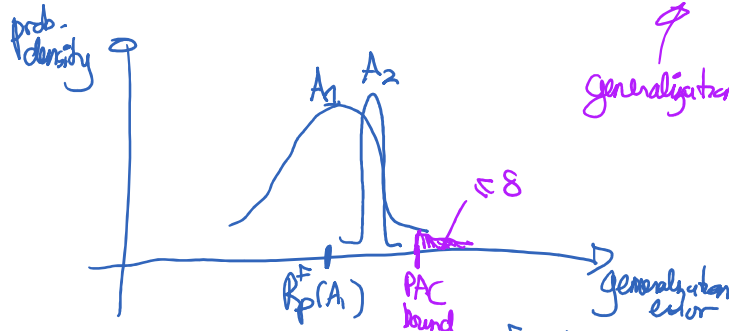
in particular, what can we say about  $L_p(A(D_n))$

different approaches:

a) "frequentist risk"  $R_{P, D_n}^F(A) \triangleq \mathbb{E}_{D_n \sim P^{\otimes n}} [L_p(A(D_n))]$

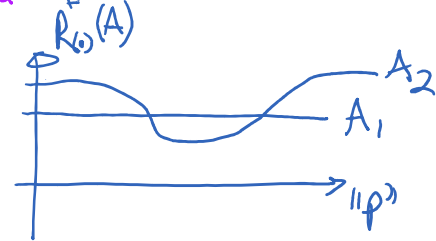
$D_n$  is random

b) PAC framework "probably approximately correct"  $P\{L_p(A(D_n)) > \text{some bound}\} \leq \delta$  ("tail bound")  
 i.e.  $L_p(A(D_n)) \leq \text{" "}$  with prob.  $\geq 1 - \delta$



generalization error bound

Issue with  $R_p^F \rightarrow$  depends on  $P$   
 "risk profiles":



weighted frequentist risk  $\mathbb{E}_{G \sim \pi(G)} [R_G^F(A)]$

c) "Bayesian posterior risk"

$R^{post}(w | D_n) \triangleq \mathbb{E}_{G \sim P(G|D_n)} [L_p(G|w)]$

• prior  $p(G)$  over  $\mathcal{H}$

Bayesian estimate  $w_n^{Bayes} = \underset{w}{\operatorname{argmin}} R^{post}(w | D_n)$

• observation model  $p(D_n | G)$

$\Rightarrow$  posterior  $p(G | D_n)$

$A^{Bayesian}$  is optimal for weighted frequentist risk using  $\pi(G) = p(G)$

15h25

No free lunch!

Frequentist risk analysis learning alg.  $A$

let  $\mathcal{D}$  be a set of distributions on  $X \times Y$

Sample complexity of  $A$  with respect to  $\mathcal{P}$

is the smallest  $n(\mathcal{P}, A, \epsilon) \geq 1$ ,  $\forall n \geq n(\mathcal{P}, A, \epsilon)$

we have  $\sup_{P \in \mathcal{P}} [R_P^E(A; n) - L_P(h_P^*)] \leq \epsilon$

"uniform result"

$h_P^* = \operatorname{argmin}_{h: X \rightarrow Y} L_P(h)$

terminology:  $A$  is consistent for dist.  $P$

if  $\lim_{n \rightarrow \infty} R_P^E(A; n) - L_P(h_P) = 0$

$A$  is uniformly consistent for a family  $\mathcal{P}$

if  $\lim_{n \rightarrow \infty} \left[ \sup_{P \in \mathcal{P}} [R_P^E(A; n) - L_P(h_P)] \right] = 0$

Binary classification  $Y = \{-1, +1\}$

I) if  $X$  is finite; then the "voting procedure" (assign the most frequent label to an input  $x$ )

is uniformly and universally consistent

$\hookrightarrow$  i.e.  $\mathcal{P}$  is all distributions  $X \times Y$

with (universal) sample complexity

$n(\mathcal{P}, \epsilon, A_{\text{voting}}) \leq \frac{|X|}{\epsilon^2}$  (free lunch?)

II) if  $X$  is infinite

no free lunch theorem (for binary with the 0-1 loss)

for any  $n$  and any learning alg.  $A$

then  $\sup_{\substack{P \text{ all} \\ \text{dist}}} [R_P^E(A; n) - L_P(h_P^*)] \geq \frac{1}{2}$

i.e.  $\exists$  always a dist  $P_A$  s.t. your  $A$  is worse than random prediction? for  $P_A$

NC, IT:

for PA

NFL II:

[thm. 7.2 in Devroye & al. 1996]

$\epsilon_1 \leq \frac{1}{16}$

let  $\epsilon_n$  be any non-increasing seq. converging to 0  
for any A,

(could be arbitrarily slowly)

then  $\exists PA$  st.  $[R_{PA}^F(A; n) - L_{PA}(h_{PA}^*)] \rightarrow \epsilon_n$   $\forall n$

e.g.  $\frac{1}{\lg(\lg(\lg(\dots(n))))}$

consequence: we need assumptions on  $\mathcal{P}$  to say anything useful

Occam's generalization error bound

- binary class. & 0/1 loss
- consider  $W$  to be a countable set

let's define a prior over  $W$ :  $\pi(w)$  i.e.  $\sum_{w \in W} \pi(w) = 1$   $\pi(w) \geq 0 \forall w$

$|w|_{\pi} = \text{"description length" of } w \triangleq \log_2 \frac{1}{\pi(w)}$

$\sum_w 2^{-|w|_{\pi}} \leq 1$   
"Kraft's inequality"

Occam's bound

for any fixed  $P_s$  with prob.  $\rightarrow$  1-S over training set  $D_n \sim P^{\otimes n}$

$\forall w \in W \quad L_P(w) \leq \hat{L}_P(w) + \frac{1}{\sqrt{2n}} \Omega_{\pi}(w; S)$

where  $\Omega_{\pi}(w; S) \triangleq \sqrt{(\ln 2) |w|_{\pi} + \ln \frac{1}{\delta}}$   
complexity measure

\* bound is useful only for dist  $P$  s.t.  $|w|_{\pi}$  is small

$\hookrightarrow$  argmin  $L_P(w)$   
 $w \in W$

$|w|_{\pi} = \log_2 \frac{1}{\pi(w)}$

if  $\pi(w) \propto \exp(-|w|^2)$

then  $|w|_{\pi} = |w|^2 + \text{const.}$

proof: use 3 things

1) Chernoff bound

note: 0/1 loss avg. appears in const. of Chernoff bound

$D \subseteq D_n$ :  $\hat{L}_n(w) \leq L(w) - \epsilon \} \leq \text{const} \cdot 2^{-n \epsilon^2} \forall \epsilon > 0$

proofs: use > unions

1) Chernoff bound  
(concentration inequality)

$$P\{\exists P_n : \hat{L}_n(w) \leq L(w) - \epsilon\} \leq \exp(-2n\epsilon^2) \quad \forall \epsilon > 0$$

appears on const. of Chernoff bound

2) union bound

$$P\{\exists x \text{ s.t. prop}(x) \text{ is true}\} \leq \sum_x P\{\text{prop}(x) \text{ is true}\}$$

3) "Kraft's inequ."

$$\sum_w 2^{-|w|/\pi} \leq 1$$

we say that  $w$  is naughty "bad" if bound fails

$$L - \epsilon > \hat{L}_n$$

$$\text{bad}(w) = \mathbb{1}\left\{L(w) > \hat{L}_n(w) + \frac{\epsilon_n(w)}{\sqrt{2n}}\right\}$$

using Chernoff,  $\hat{L}_n(w) \leq L(w) - \epsilon_n(w)$  with small prob.

$$\begin{aligned} P\{\text{bad}(w)\} &\leq \exp(-2n\epsilon_n(w)^2) = \exp\left(-2n \frac{1}{2n} ((\ln 2) w/\pi + (\ln 2))\right) \\ &= \delta 2^{-|w|/\pi} \end{aligned}$$

using union bound

$$P\{\exists w \text{ s.t. bad}(w)\} \leq \sum_w P\{\text{bad}(w)\} \leq \sum_w \delta 2^{-|w|/\pi} \stackrel{\text{Kraft}}{\leq} \delta$$

### Surrogate Loss

NP hard to minimize  $\hat{L}_n(w)$ ; replace with  $\hat{J}_n(w)$  which is "surrogate"

- eg. hinge loss
- log-loss

next: countable  $\rightarrow$  uncountable  
"PAC-Bayes"