

today: theory basics

theory basics

decision theory setup

estimate $h_w: X \rightarrow Y$

test loss



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

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 $(x^{(i)}, y^{(i)})_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} P$ \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

$\hat{L}_n(w) \xrightarrow[n \rightarrow \infty]{\text{a.s.}} L_P(w)$ (LLN) for each w (pointwise)

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

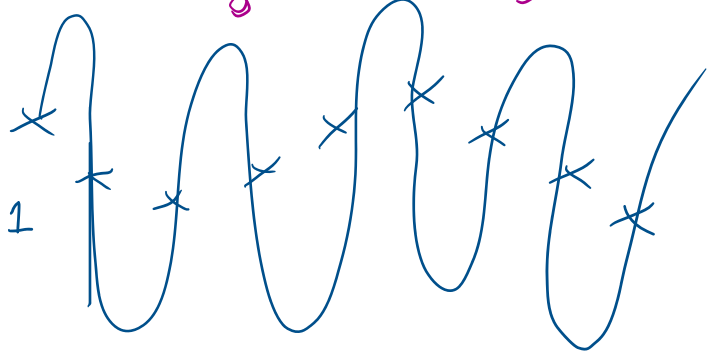
note: minimizing the training error gives no guarantee in general ∇
 $L_P(\hat{w}_n)$

\rightsquigarrow 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: study properties of learning alg $A: D_n \rightarrow W$

in particular, what can we say about $L_P(\underbrace{A(D_n)}_n)$

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

\hat{w}_n
 D_n is random

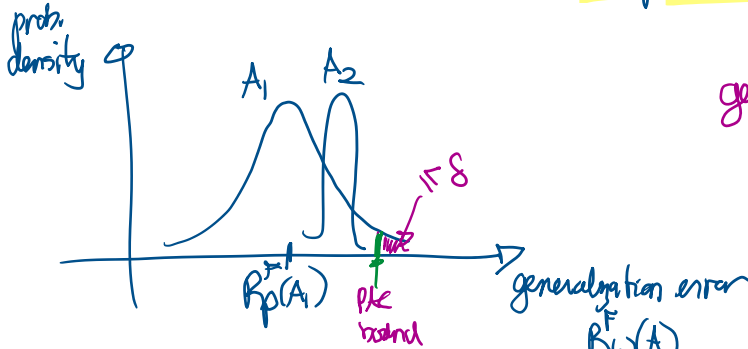
different approaches:

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

b) PAC framework "probably approximately correct" $P\{L_p(A(D_n)) > \text{some bound}\} \leq \delta$

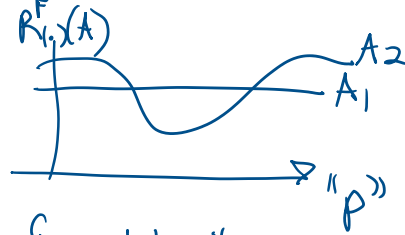
ie. $L_p(A(D_n)) \leq \epsilon$ with prob $\geq 1 - \delta$

generalization error bound



issue with $R_p^F \rightarrow$ depends on p

"risk profiles":



weighted frequentist risk $\mathbb{E}_{\Theta \sim \pi(\Theta)} [R_{p_\Theta}^F(A)]$

c) "Bayesian posterior risk"

$R^{\text{post}}(w | D_n) \triangleq \mathbb{E}_{\Theta \sim p(\Theta | D_n)} [L_p(w)]$
 Bayesian estimate $\hat{w}_n^{\text{Bayes}} = \underset{w}{\text{argmin}} R^{\text{post}}(w | D_n)$

- prior $p(\Theta)$ over dist.
- observation model $p(D_n | \Theta)$
- \Rightarrow posterior $p(\Theta | D_n)$

A Bayesian is optimal for weighted frequentist risk using $\pi(\Theta) = p(\Theta)$

No free lunch!

frequentist risk analysis learning alg. A

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

sample complexity of A with respect to \mathcal{P}

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

we have that $\sup_{p \in \mathcal{P}} [R_p^F(A; n) - L_p(h_p^*)] \leq \epsilon$

we have that $\sup_{P \in \mathcal{P}} [R_P^F(A; n) - L_P(h_P^*)] \leq \epsilon$
 "uniform result" \rightarrow $h_P^* = \arg \min_{h: X \rightarrow \mathcal{Y}} L_P(h)$

terminology: \bullet A is consistent for fixed dist. P

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

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

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

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

I) if X is finite; then the " Voting procedure" (assign the most frequent label to an output x)
 is uniformly and universally consistent

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

with (universal) sample complexity

$n(\mathcal{P}, \epsilon, \text{Averaging}) \leq \frac{|X|}{\epsilon^2}$

(free lunch!)
 \Downarrow

10h28

II) if X is infinite

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

for any n and any learning alg. A

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

i.e. \exists always a dist $P_{A, n}$ s.t. your A is worse than random prediction for $P_{A, n}$

NFL II: [Thm. 7.2 in Perceptron & al. 1996] — $\epsilon_n \leq \frac{1}{16}$

Let ϵ_n be any non-increasing seq. converging to 0

(could be arbitrarily slow)

for any A , then $\exists P_A$

s.t. $[R_{P_A}^F(A; n) - L_{P_A}(h_{P_A}^*)] \geq \epsilon_n \cdot n$

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

⊗ consequence: we need assumptions on \underline{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. $\sum_{w \in W} \pi(w) = 1$ $\pi(w) \geq 0$ tho
 $|w|_{\pi} = \text{"description length" of } w \triangleq \log_2 \frac{1}{\pi(w)}$

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

"Kraft's inequality"

Occam's bound

for any fixed P ; with prob $\geq 1-\delta$ over training $D_n \sim P^{\otimes n}$

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

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

⊗ minimizing bound \rightarrow universally consistent alg.

⊗ bound is useful only for dist. P st. $|w^*|_{\pi}$ is small
 $\hookrightarrow \arg \min_{w \in W} L_P(w)$

note: 0-1 loss assumption appears in constant of bound

proof: use 3 things

1) Chernoff bound (concentration inequality)

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

2) union bound

$$P\{\exists x \text{ st. prop}(x) \text{ is true}\} \leq \sum_{x \in X} P\{\text{prop}(x) \text{ is true}\}$$

\uparrow countable

3) "Kraft's ineq." $\sum_w 2^{-|w|_{\pi}} \leq 1$

we say that w is "naughty" if bound fails
 "bad" $L - \epsilon > \hat{L}_n$

$$\text{bad}(w) = \mathbb{1}\{L(w) > \hat{L}_n(w) + \frac{1}{\sqrt{2n}} \Omega_{\pi}(w; \delta)\}$$

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

$$P\{\text{bad}(w) = 1\} \leq \exp(-2n\epsilon^2) = \exp(-2n / (1 / (\ln 2) |w|_{\pi} + \ln \frac{1}{\delta}))$$

using Chernoff, $L_n(w) \approx L(w) - \epsilon$ with small prob

$$P\{\text{bad}(w)=1\} \leq \exp(-2n\epsilon_n(w)^2) = \exp(-2n \left(\frac{1}{2n} (\ln 2) \left(\frac{1}{\epsilon} + \ln \frac{1}{\delta} \right) \right)^2)$$

$$= \delta 2^{-|w|n}$$

using union bound

$$P\{\exists w: \text{bad}(w)\} \leq \sum_w P\{\text{bad}(w)\} \leq \sum_w \delta 2^{-|w|n} \stackrel{\text{Kraft inequality}}{\leq} \delta //$$

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

then $|w|n = \|w\|^2 + \text{cst.}$

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

e.g. • hinge loss
• log-loss

next: countable \rightarrow uncountable
"PAC-Bayes"