

Lecture 4 - scribbles

Thursday, January 17, 2019 13:25

today: theory basics

theory basics

decision theory setup

estimate: $h_w: X \rightarrow \mathcal{Y}$

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

task loss

ultimate goal is to find $w^* = \underset{w \in W}{\operatorname{arg\,min}} L_p(w)$

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

suppose $\underbrace{(x^{(i)}, y^{(i)})}_{\triangleq D_n} \stackrel{\text{i.i.d.}}{\sim} P$ training data

→ we could look at

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

from statistics / prob. theory

$$\hat{L}_n(w) \xrightarrow[n \rightarrow \infty]{a.s.} L_p(w) \quad \forall w$$

→ this is weaker

$$\text{than } \sup_w |\hat{L}_n(w) - L_p(w)| \xrightarrow{a.s.} 0$$

$$\left[\begin{array}{c} X_n \xrightarrow{a.s.} X \\ \dots \end{array} \right]$$

$$\text{learning alg.: } \hat{w}_n = A(D_n)$$

algorithm

$$P\{w \in \Omega : X_n(w) \xrightarrow{n \rightarrow \infty} X(w)\} = 1$$

note: minimizing training error gives no guarantee in general

e.g. polynomial regression

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



in Learning Theory, study properties of learning alg

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

different approaches:

a) "frequentist risk"

$$R_{P,n}^F(A) \triangleq \mathbb{E}_{D_n \sim P^{(n)}} [L_p(A(\hat{w}_n))]$$

\hat{w}_n is random

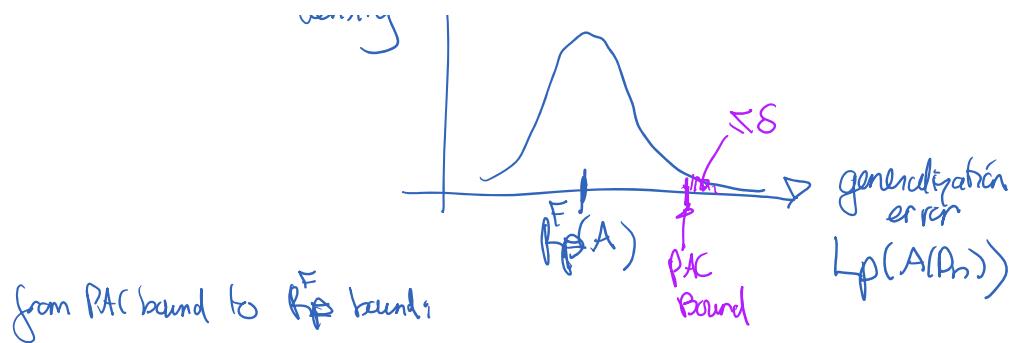
b) PAC framework
"probably approximately correct"

$$P\{L_p(A(D_n)) > \text{some bound}\} \leq \delta$$

i.e. $L_p(A(D_n)) \leq \text{"some bound}$ with prob. $\geq 1 - \delta$

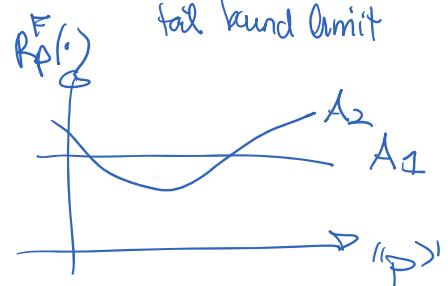
↳ generalization error bound





$$\mathbb{E}X = \mathbb{E}_{\{B \sim \mathcal{B}\}} X \mathbb{I}_{\{X \leq B\}} + \mathbb{E}_{\{B \sim \mathcal{B}\}} X \mathbb{I}_{\{X > B\}}$$

tail bound limit



weighted frequentist risk
 $\mathbb{E}_{D \sim p(\cdot|P)} [R_p^F(A)]$

c) "Bayesian posterior risk"

$$R_{\text{post}}^{\text{post}}(w | D_h) \triangleq \mathbb{E}_{G \sim p(G | D_h)} [L_{p_G}(w)]$$

Bayesian estimate $\hat{w}_n = \underset{w}{\operatorname{argmin}} R_{\text{post}}^{\text{post}}(w | D_h)$

- prior $p(\theta)$ over distributions
- observation model $p(D_h | \theta)$
- posterior $\Rightarrow p(\theta | D_h)$

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

14h20

No free lunch

frequentist risk analysis of learning algorithm A

Let \mathcal{P} be a set of distribution $X \times \mathcal{S}$

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 $\sup_{P \in \mathcal{P}} [R_p^F(A; n) - L_p(h_p^*)] < \epsilon$

"uniform result" \uparrow
 $h_p^* = \operatorname{arg\min}_{h: X \rightarrow \mathcal{S}} L_p(h)$

terminology:

- A is consistent for dist. P

if $\lim_{M \rightarrow \infty} R_p^F(A; M) - L_p(h_p^*) = 0$

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

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

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

I) if X is finite, then the "voting procedure" (assign the most frequent label to an input x) is uniformly and universally consistent

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

with (universal) sample complexity $n(\mathcal{P}, \epsilon, \text{Averaging}) \leq \frac{|X|}{\epsilon^2}$ (free lunch? !!)

II) If X is infinite

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

for any n and any learning alg. A

$$\text{then } \sup_{\substack{\mathcal{P} \text{ all} \\ \text{distrn}}} [R_p^F(A; n) - L_p(h_p^*)] \geq \frac{1}{2}$$

i.e. \exists always a dist. P s.t. your alg. A is worse than random prediction (?)

NFLT II:

[Thm. 7.2 in Devroye et al. 1996]

Let ε_n be any non-increasing sequence converging to zero (could be arbitrarily slowly e.g. $\frac{1}{\log(\log(\dots(n)))}$)

$$\text{then } \exists P \text{ s.t. } R_p^F(A; n) - L_p(h_p^*) \geq \varepsilon_n \quad \forall n$$

⊗⊗ consequence: we need assumptions on P to say anything

Occam's generalization error bound

- binary class, and 0-1 loss
- consider W to be a countable set

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

$|w|_\pi$ = "description length" of w

$$\triangleq \log_2 \frac{1}{\pi(\omega)}$$

Occam's bound

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

$$\forall \omega \in W \quad L_p(\omega) \leq \hat{L}_p(\omega) + \frac{1}{\sqrt{2n}} D_{\pi}(\omega; \delta)$$

$$\text{where } D_{\pi}(\omega; \delta) \triangleq \sqrt{(2n\delta)} \frac{\|\omega\|_{\pi} + \ln \frac{1}{\delta}}{8}$$

Complexity measure

(*) bound is only useful for distribution P s.t. $\|\omega^*\|_{\pi}$ is small

$$\omega^* = \operatorname{argmin}_{\omega \in W} L_p(\omega)$$

$$\|\omega\|_{\pi} = \log \frac{1}{\pi(\omega)}$$

$$\text{if } \pi(\omega) \propto \exp(-\|\omega\|^2)$$

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

note: O-1 loss assumption appears in constants of Chernoff bound

proof: use 3 things

1) Chernoff bound.
(concentration inequality)

2) union bound

3) "Kraft's inequality"
.. "bad". ..

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

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

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

$$\dots \Rightarrow L - \varepsilon > \hat{L}_n$$

we say w is "bad" if bound fails

$$\text{bad}(w) = \mathbb{1}\{L(w) > \hat{L}_n(w) + \underbrace{\Omega_n(w; \delta)}_{\varepsilon_n(w)}\}$$

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

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

$$\Pr\{\text{bad}(w)\} \leq \exp(-2n\varepsilon_n(w)^2) = \exp\left(-2n\frac{1}{n}\left((\ln 2)(|w|_n + \ln \frac{1}{\delta})\right)\right) = \delta 2^{-|w|_n}$$

using union bound

$$\Pr\{\exists w : \text{bad}(w)\} \leq \sum_w \Pr\{\text{bad}(w)\} \leq \sum_w S 2^{-|w|_n} \leq S$$

//

Surrogate loss:

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

e.g. hinge loss
log-loss
...