

today: theory basics

theory basics

decision theory setup

estimate $h_w: X \rightarrow Y$

task loss

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

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

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

suppose $(x^{(i)}, y^{(i)})_{i=1}^n \stackrel{\text{iid.}}{\sim} P \rightarrow$ we could look at
 $\hat{L}_n(w) = \frac{1}{n} \sum_{i=1}^n l(y^{(i)}, h_w(x^{(i)}))$
 $\triangleq D_n$ training data

from statistics/prob. theory

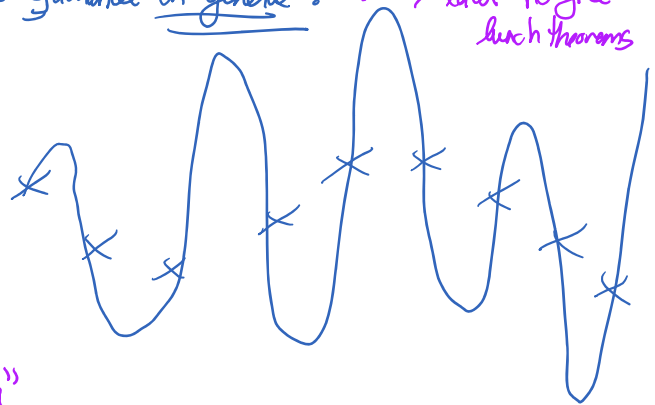
$\hat{L}_n(w) \xrightarrow[n \rightarrow \infty]{\text{a.s.}} L_P(w)$ for each fixed w (pt. wise)
 (LLN)

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 in general \rightarrow later no free lunch theorems

e.g. polynomial regression

for n points, can get zero training error with polynomial of deg. $n-1$



\Rightarrow "overfitting"

in learning theory \rightarrow study properties of learning alg.

in particular, what can we say about $L_P(\hat{w}_n)$

different approaches

a) "frequentist risk"

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

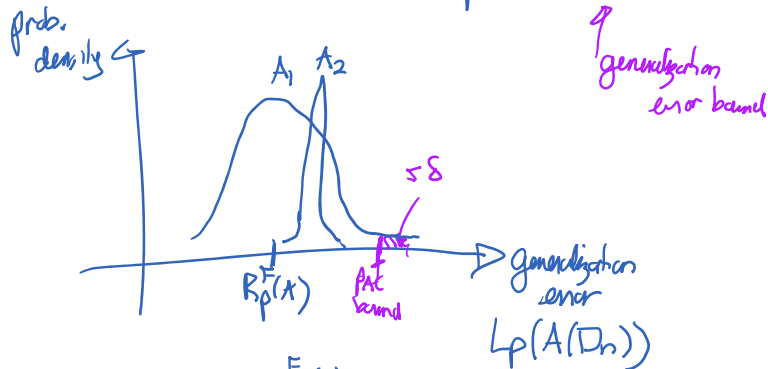
D_n is random

b) PAC framework
"probably approximately correct"

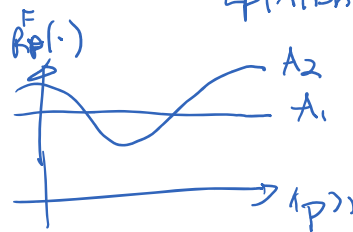
$$P \{ L_p(A(D_n)) > \text{some bound} \} \leq \delta$$

("tail bound")

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



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



weighted frequentist risk

$$\mathbb{E}_{\theta \sim \pi(\theta)} [R_{p, \theta}^F(A)]$$

c) "Bayesian posterior risk"

$$R^{\text{post}}(w | D_n) \triangleq \mathbb{E}_{\theta \sim p(\theta | D_n)} [L_{p, \theta}(w)]$$

prior $p(\theta)$ over \mathcal{H}

Bayesian estimate $\hat{w}_n^{\text{Bayes}} = \underset{w}{\text{argmin}} R^{\text{post}}(w | D_n)$

• observation model $p(D_n | \theta)$

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

A^{Bayesian} is optimal for weighted freq. risk using $\pi(\theta) = p(\theta)$

14/20

No free lunch P

frequentist risk analysis learning algorithm A

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

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

is the smallest $n(P, A, \epsilon) \Rightarrow \forall n \geq n(P, A, \epsilon)$

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

"uniform result"

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

terminology: \bullet A is consistent for 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 $Y = \{-1, +1\}$

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 on $X \times Y$

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

II) if X is infinite

no free lunch theorem (for binary with ℓ the 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 s.t. your alg. A is worse than random prediction (?)

NFLT II:

[thm. 7.2 in Devroye *et al.* 1996] $\epsilon_n \leq \frac{1}{16}$

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

then $\exists P$ s.t. $R_P^F(A; n) - L_P(h_P^*) \geq \epsilon_n \forall n$

(could be arbitrarily slowly)

eg. $\int \log \log \log \dots \log (1/n)$

⚠️ Consequence: we need assumptions on \mathcal{F} to say anything

Occam's generalization error bound

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

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

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

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

Occam's bound

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

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

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

complexity
 measure

⊛ bound is useful only for dist. P s.t. $|w|_{\pi}$ is small
 $\hookrightarrow \arg \min_{w \in W} L_P(w)$

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

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

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

note: 0-1 loss ass. appears in constants of Chernoff 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

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

3) "Kraft's inequality"

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

we say w is naughty if bound fails

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

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

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

$$\begin{aligned} P\{\text{bad}(w)\} &\leq \exp(-2n \epsilon_n(w)^2) = \exp\left(-2n \frac{1}{2n} \left((\ln 2) |w|_1 + \ln \frac{1}{2} \right)\right) \\ &= \delta 2^{-|w|_1} \end{aligned}$$

using union bound

$$P\{\exists w; \text{bad}(w)\} \leq \sum_w P\{\text{bad}(w)\} \leq \sum_w \delta 2^{-|w|_1} \stackrel{\text{Kraft}}{\leq} \delta$$

Sumoate loss

NP hard to minimize $\hat{L}_n(w)$; replace with $\hat{J}_n(w)$ which is "sumoate"
eg. hinge loss
log-loss