

Analysis of finite hypothesis class

Setup

$h_S \in \text{ERM}_H(S)$, assume $|H| < \infty$,

assume realizability ($\exists h^* \in H, L_{D,f}(h^*) = 0$)

and have binary classification w/ 0-1 loss

$$\text{so } L_{D,f}(h) = \mathbb{P}_{x \sim D} (h(x) \neq f(x))$$

Goal is PAC learning: $\forall \epsilon, \delta \in (0, 1)$

$$\text{prove } L_{D,f}(h_S) \leq \epsilon \quad \text{w.p.} \geq 1 - \delta$$

$$\text{aka } D^m(\{S : L_{D,f}(h_S) > \epsilon\}) < \delta$$

Analysis:

let H_B = "bad hypotheses"

$$= \{h \in H : L_{D,f}(h) > \epsilon\} \subseteq H$$

$$\text{so } |H_B| \leq |H|$$

let M = Mistake samples

$$= \{S : \exists h \in H_B \text{ s.t. } \underbrace{h \in \text{ERM}_H(S)}\}$$

equiv., $\hat{L}_S(h) = 0$
by realizability

Our $h_S \in \text{ERM}_H(S)$

so if $S \notin M$ we know h_S is good.

How big is M ?

$$M = \bigcup_{h \in H_B} \{S : \hat{L}_S(h) = 0\}$$

$$\text{so } D^m(M) \leq \sum_{h \in H_B} D^m(\{S : \hat{L}_S(h) = 0\})$$

via union bound

P.1

(aside: union bound: $\mathbb{D}(A \cup B) \leq \mathbb{D}(A) + \mathbb{D}(B)$
 \uparrow "or"

Directly useful for bounding failures.

For successes, use $\mathbb{D}(A) = 1 - \mathbb{D}(A^c)$

and De Morgan's laws: $(\cup A_i)^c = \cap A_i^c$
 $(\cap A_i)^c = \cup A_i^c$

Each term in sum: $\mathbb{D}^m(\{S: \hat{L}_S(h) = 0\})$, fix $h \in \mathcal{H}_B$

Since $\hat{L}_S(h) = 0$ means $h(x_i) = \underbrace{f(x_i)}_{\text{true label}} \forall i \in [m]$

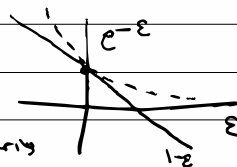
$$\begin{aligned} \mathbb{D}^m(\dots) &= \mathbb{D}^m(\{S: h(x_i) = f(x_i) \forall i \in [m]\}) \\ &= \prod_{i=1}^m \mathbb{D}(\{x: h(x) = f(x)\}) \quad \text{by independence} \\ &= \mathbb{D}(\{x: h(x) = f(x)\})^m \quad \text{since identically distr.} \\ &= (1 - \mathbb{D}(\{x: h(x) \neq f(x)\}))^m \quad \text{take complement} \\ &= (1 - \underbrace{L_{D,f}(h)}_{> \varepsilon \text{ since } h \in \mathcal{H}_B})^m \quad \text{since 0-1 loss} \end{aligned}$$

$$\leq (1 - \varepsilon)^m$$

$$\leq e^{-\varepsilon m} \quad \text{since } (1 - \varepsilon) \leq e^{-\varepsilon}$$

proof via Taylor Series
or derivative + convexity

nicer looking



Going back to union bound,

$$\mathbb{D}^m(M) \leq \sum_{h \in \mathcal{H}_B} e^{-\varepsilon m} = |\mathcal{H}_B| \cdot e^{-\varepsilon m} \leq \underbrace{|\mathcal{H}|}_{\text{"J"}} \cdot e^{-\varepsilon m}$$

want $\delta = |\mathcal{H}| \cdot e^{-\epsilon m}$, or, $-\epsilon m = \log(\delta/|\mathcal{H}|)$

So...

Corollary 2.3 (binary class, 0-1 loss, realizable, $|\mathcal{H}| < \infty$)
 $\forall \delta \in (0, 1) \forall \epsilon > 0$, if $m \in \mathbb{Z}$ satisfies

$$m \geq \frac{1}{\epsilon} \log(|\mathcal{H}|/\delta) \text{ then } \forall f, \forall D$$

(as long as realizable), if S has m iid samples
then $\forall h_S \in \text{ERM}_{\mathcal{H}}(S)$, w. prob $\geq 1 - \delta$,

$$\underbrace{L_{D,f}(h_S)}_{\text{"approximate"}} \leq \underbrace{\epsilon}_{\text{"probably"}}$$

m is our "sample complexity":

smaller δ
and/or
smaller $\epsilon \quad \Rightarrow \quad \text{need more samples}$

Later we'll drop realizability assumption and
allow for "agnostic" case, and also more
general loss functions: Corollary 4.6

Then $m \geq \frac{1}{\epsilon^2} \log\left(\frac{2|\mathcal{H}|}{\delta}\right)$
 \uparrow yuck!

So, from training examples, we can learn arbitrarily
well (given enough data)

... if our choice of \mathcal{H} (inductive bias)
was good enough to include a good
classifier.

more generally, (§3), define

Def A hypothesis class \mathcal{H} is "PAC learnable" if

\exists function $m_{\mathcal{H}} : (0,1)^2 \rightarrow \mathbb{N}$ "sample complexity"
and some learning algo "A" (eg. ERM $_{\mathcal{H}}$)
such that
 $\forall D$ over X , $\forall f: X \rightarrow \{0,1\}$ \rightarrow still doing binary classif.
if (\mathcal{H}, D, f) is realizable, then

if $m \geq m_{\mathcal{H}}(\epsilon, \delta)$ iid samples are used in algo A,
A returns a classifier/hypothesis h s.t. w.p.
 $\geq 1 - \delta$, $L_{D,f}(h) \leq \epsilon$.

Aside Since S is a random variable (r.v.),

so is $X_m = L_{D,f}(A(S))$.

We're saying $0 \leq X_m \leq \epsilon$ w.p. $\geq 1 - \delta$ if $m \geq m(\epsilon, \delta)$

How does this relate to usual notions of convergence of r.v.?

You can show that existence of a finite $m(\epsilon, \delta)$ is

iff $\lim_{m \rightarrow \infty} \mathbb{E} \sup X = 0$, i.e., "L' convergence"

(see exer. 4.1)

$X_m \xrightarrow{L^p} X$ means

$$\lim_{m \rightarrow \infty} \mathbb{E} |X_m - X|^p = 0$$

(aside still)

and conv. in measure / probability :

$$X_m \xrightarrow{P} X \text{ means } \forall \varepsilon > 0$$

$$\lim_{m \rightarrow \infty} P(\{\omega : |X_m(\omega) - X(\omega)| > \varepsilon\}) = 0$$

and almost sure (a.s.) convergence:

$$P(\{\omega : X_m(\omega) \rightarrow X(\omega)\}) = 1 \quad \text{i.e. ptwise a.e.}$$

Agnostic PAC learning means relaxing assumption of an oracle $f: X \rightarrow Y$

Now, let D be a ^{joint} distribution over $Z := X \times Y$

and D_x its marginal wrt X ,

$D(Z/x)$ its conditional.

Before, in (pure) PAC learning, we're assuming

$D(Z/x)$ collapses to a deterministic function $f(x)$

\hat{L}_g is unchanged but (for 0-1 loss)

$$\text{turn } L_{D,f}(h) = P_{x \sim D}(h(x) \neq f(x))$$

$$\text{into } L_D(h) = P_{(x,y) \sim D}(h(x) \neq y) \quad \left. \vphantom{L_D(h)} \right\} \text{AGNOSTIC CASE}$$

Now, it's typically too much to ask for $L_D(h) \leq \varepsilon$

...

△ One algo for all distributions

Def \mathcal{H} is agnostic PAC learnable if $\exists m_H: (0,1)^2 \rightarrow \mathbb{N}$,
 \exists algo A , s.t. $\forall \epsilon, \delta \in (0,1)$, \forall distr. \mathcal{D} over $X \times Y$,
 A s.t. if S has $m \geq m_H(\epsilon, \delta)$ iid samples, then
w.p. $\geq 1 - \delta$
$$L_D(h) \leq \epsilon + \underbrace{\min_{h' \in \mathcal{H}} L_D(h')}_{\text{new}}$$

where $h = A(S)$

Note: Mohri puts it this way: agnostic = inconsistent + stochastic
consistent: means $\min_{h' \in \mathcal{H}} L(h') = 0$
deterministic: means $D(y(x)) = f(x)$

Remark An alternative comparison would be the Bayes Optimal Predictor, $f_D(x) = \begin{cases} 1 & \text{if } P(y=1/x) \geq 1/2 \\ 0 & \text{else} \end{cases}$ uses knowledge of \mathcal{D} !

Optimal, but

- ① you have to know \mathcal{D}
- ① it changes when \mathcal{D} changes
- ② it's asking too much

Loss function $\hat{L}_S(h) = \frac{1}{n} \sum_{i=1}^n l(h, z_i)$
Allow general risk/loss $L_D(h) = \mathbb{E}_{z \sim \mathcal{D}} l(h, z)$

Ex: 0-1 loss for binary classif,

$$l_{0-1}(h, (x, y)) = \begin{cases} 0 & h(x) = y \\ 1 & \text{else} \end{cases} = \mathbb{I}_{h(x) \neq y}(h, z)$$

Ex regression, l_p losses

$$l_p(h, (x, y)) = (h(x) - y)^p, \quad p \geq 1 \text{ for convex}$$

$p=2$ most common $p \in \mathbb{Q}$

Brief History:

a lot of the framework due to

Vapnik + Chervonenkis '71

PAC due to Valiant '84

Vladimir Vapnik (1936-) from Samarkand (Uzb.) PhD '64

'71 VC theory

'90 move to US, AT&T, develop SVMs

'98 "Stat. Learning Theory" book 60,000+ citations

'02 NEC

'03 Columbia

'14 Facebook AI Research, '16 Venere (abs

Alexey Chervonenkis (1938-2014)

'71 VC theory

Kept researching, stayed in USSR/Russia

Leslie Valiant (1949-) from Budapest

Prof. at Harvard since '82

Many works in Theoretical Computer Science (TCS)

two sons on faculty at Stanford + Brown