

VC classes \Rightarrow unif concentration from emp. error
to generalization error

Q1:
$$\mathbb{E}_{S \sim D^m} \sup_{f \in \mathcal{F}} \mathbb{E}_S f(z) - \mathbb{E}_{S'} f(z), \quad (*)$$

 $S' \sim D^m$

can we directly bound this using Massart's
Lemma?

Q2: upper bounding (*) by

$$\mathbb{E}_{S \sim D^m} \sup_{f \in \mathcal{F}} \mathbb{E}_S f(z) + \mathbb{E}_{S' \sim D^m} \sup_{f \in \mathcal{F}} (-\mathbb{E}_{S'} f(z))$$

(w/o introducing Rademacher r.v.'s.)

will the pf still go through?

Lower bounds for statistical learning

$O\left(\frac{d}{\epsilon^2}\right)$ training examples \Rightarrow ERM has ^{excess} error $\leq \epsilon$.

what if we only have $O(d)$ examples?
 $o(d)$

(A, etc)

discuss learnability as a property of hypothesis

class only.

Def: (H, \mathcal{P}) , family of distn
 is said to be agnostic PAC learnable, if
 there exists an alg A , and a sample complexity $m = m(\epsilon, \delta)$
 $\rightarrow \mathbb{N}$, s.t. ~~for any distn $D \in \mathcal{P}$~~ for any $\epsilon > 0, \delta > 0, \exists$
 $m \geq m(\epsilon, \delta)$, then w.p. $(1 - \delta)$, over the draw of m
 training examples iid from D ,

$$\text{err}(\hat{h}_{\text{ACS}}, D) - \min_{h \in H} \text{err}(h, D) \leq \epsilon$$

Def: H is said to be (realizable) PAC learnable, if
 there exists an alg A , and a sample complexity $m = m(\epsilon, \delta)$
 $\rightarrow \mathbb{N}$, s.t. for any distn D , ^{realizable by H} for any $\epsilon > 0, \delta > 0, \exists$
 $m \geq m(\epsilon, \delta)$, then w.p. $(1 - \delta)$, over the draw of m
 training examples iid from D ,

$$\text{err}(\hat{h}_{\text{ACS}}, D) - \min_{h \in H} \text{err}(h, D) \leq \epsilon$$

Finite VC dim

PAC learnable



(\Rightarrow) Finite VC dim \Rightarrow uniform convergence
 \Rightarrow ERM sample complexity $O\left(\frac{d}{\epsilon^2}\right)$
 \Rightarrow \mathcal{H} is PAC learnable.

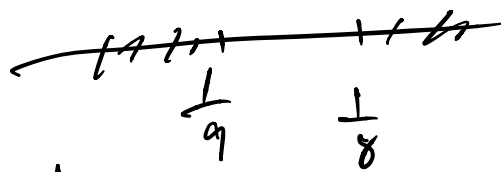
(\Leftarrow) see next.

Thm: \mathcal{H} : $VC(\mathcal{H}) \geq d$. If $m = \# \text{ training examples} \leq \frac{d}{2}$.
 then, \forall alg A , $\exists D$. realizable by \mathcal{H} ,

$$\mathbb{E}_{S \sim D^m} \underbrace{\text{err}(A(S), D)}_X \geq \frac{1}{4}. \quad (*)$$

Remark 1: $(*)$ also implies that

$$\mathbb{P}_{S \sim D^m} \left(\text{err}(A(S), D) > \frac{1}{8} \right) \geq \frac{1}{8} \quad (\Delta)$$



showing that A does not $(\epsilon = \frac{1}{8}, \delta = \frac{1}{9})$ -PAC
 learn \mathcal{H} with $m \leq \frac{d}{2}$ examples.

(If A $(\frac{1}{8}, \frac{1}{9})$ -learns \mathcal{H} w/ $m \leq \frac{d}{2}$ examples, then

$$\mathbb{P}_{S \sim D^m} \left(\text{err}(A(S), D) > \frac{1}{8} \right) \leq \frac{1}{9}, \text{ which contradicts w/ } \Delta$$

$(*) \Rightarrow (\Delta)$

Fact: $X \in [0,1]$, $\mathbb{E} X \geq \frac{1}{4}$. then $P(X > \frac{1}{8}) \geq \frac{1}{8}$

Pf:

$$= \mathbb{E} \left[\underbrace{X I(X \leq \frac{1}{8})}_{\leq \frac{1}{8}} \right] + \mathbb{E} \left[\underbrace{X I(X \in (\frac{1}{8}, 1])}_{\geq \frac{1}{8}} \right]$$

$$\leq \frac{1}{8} + P(X > \frac{1}{8})$$

$$\Rightarrow P(X > \frac{1}{8}) \geq \frac{1}{8}$$

Remark 2: $V(\mathcal{H}) = \infty \Rightarrow \mathcal{H}$ is not PAC learnable
 (d=2m in the above thm)

b/c: $V(\mathcal{H}) = \infty \Rightarrow \forall m, \forall \text{Alg}, \exists \text{distr } \mathcal{D} \text{ realizable by } \mathcal{H},$

$$P_{S \sim \mathcal{D}^m} (\text{err}(A(S), \mathcal{D}) > \frac{1}{8}) > \frac{1}{9}$$

$$f(\frac{1}{8}, \frac{1}{9})$$

max
 \mathcal{D} : realizable by \mathcal{H} $\mathbb{E}_{S \sim \mathcal{D}^m} \text{err}(A(S), \mathcal{D}) \geq \frac{1}{4}$

Pf of thm

$\forall \text{ alg } A, \exists \mathcal{D} \text{ realizable by } \mathcal{H},$

$$\mathbb{E}_{S \sim \mathcal{D}^m} \text{err}(A(S), \mathcal{D}) \geq \frac{1}{4} \quad (*)$$

$$\min_A \max_{\mathcal{D}: \text{realizable by } \mathcal{H}} \mathbb{E}_{S \sim \mathcal{D}^m} \text{err}(A(S), \mathcal{D}) \geq \frac{1}{4} \quad (\square)$$

(Minimax lower bound)

define a family of distributions $\mathcal{P} = \{D_b : b \in \{\pm 1\}^d\}$

want to show

$$\min_A \mathbb{E}_{b \sim U(\pm 1)^d} \mathbb{E}_{S \sim D^m} \text{err}(A(S), D) \geq \frac{1}{4}$$

(which implies (D))

First: find a set of examples z_1, \dots, z_d shattered by \mathcal{H} .

define $D_b : \mathcal{P}(x = z_i, y = b_i) = \frac{1}{d} \quad \forall i=1 \dots d$

\uparrow \hat{h}	-1	$+1$	$+1$	$d=3$	$b = (+1, +1, -1)$
	$+1$	-1	-1		$b_3 = -1$
	0	0	0		
	z_1	z_2	z_3		z_d

all D_b 's are realizable by \mathcal{H} .

shorthand: $\hat{h} = A(S)$

$$\forall A, \mathbb{E}_{b, S} \text{err}(\hat{h}, D_b) \geq \frac{1}{4}$$

$$\text{given } h, \text{err}(h, D_b) = \sum_{i=1}^d \frac{1}{d} I(h(z_i) \neq b_i)$$

$$\text{want to show } \sum_{i=1}^d \mathbb{E}_{b, S} I(\hat{h}(z_i) \neq b_i) \geq \frac{d}{4}$$

We are going to show:

$$P_{b, S} (\hat{h}(z_1) \neq b_1) \geq \frac{1}{4} \quad (\text{can show this for other } z_i \neq S_x)$$

(#)

Conditioned on $z_1 \notin S_x$ it's independent of b_1 .

b_1, \dots, b_d

$$S_x = \{z_2, z_3, z_5, z_7\}$$

$$S_x = \{z_1, z_2, z_5, z_7\}$$

$$S_x = \{x_1, \dots, x_m\} \stackrel{iid}{\sim} \text{unif}(\{z_1, \dots, z_d\})$$

the unlabeled training examples.

$$S = \{(x_1, y_1), \dots, (x_m, y_m)\}$$

\hat{h}

$$I(\hat{h}(z_1) \neq b_1)$$

(expectation?)

$$(\#) \geq P_{b, S_x} (\hat{h}(z_1) \neq b_1, z_1 \notin S_x)$$

$$= P_{b, S_x} (\hat{h}(z_1) \neq b_1 \mid z_1 \notin S_x) \cdot P_{S_x} (z_1 \notin S_x)$$

(1)

(2)

$$\textcircled{2} \quad P(z_1 \in S_X)$$

$$\approx P(z_1 \in \bigcup_i \{x_i\})$$

$$\leq \sum_{i=1}^m P(z_1 = x_i) = \frac{m}{d} \leq \frac{1}{2}$$

$$\textcircled{2} \approx \frac{1}{2}.$$

\textcircled{1} Conditioned on $z_1 \notin S_X$,

$\hat{h}(z_1)$ independent of b_1

$$P(\hat{h}(z_1) \neq b_1 \mid z_1 \notin S_X) = \frac{1}{2}.$$

$$\textcircled{\#} \approx \frac{1}{2} < \frac{1}{2} \approx \frac{1}{4}. \quad \star$$

Def: \mathcal{H} is said to satisfy the uniform convergence property if there exists a fn $f_{\mathcal{H}}: (0,1)^2 \rightarrow \mathbb{N}$, such that for any distn D , for any $\epsilon, \delta > 0$, if $m \geq f_{\mathcal{H}}(\epsilon, \delta)$, then, w.p. $1 - \delta$, over the draw of m iid training examples from D ,

$$| \text{err}(h, S) - \text{err}(h, D) | \leq \epsilon,$$

for all $h \in \mathcal{H}$.

Thm: (the fundamental thm of statistical learning)

The follow statements are equivalent:

1. \mathcal{H} satisfies the uniform convergence property
2. \mathcal{H} is agnostic PAC learnable w/ ERM
3. \mathcal{H} is agnostic PAC learnable
4. \mathcal{H} is PAC learnable
5. \mathcal{H} has finite VC dimension

S : observations in nature
 \mathcal{H} : scientific theory

\mathcal{H} has infinite VC dim

"falsifiable"

pf: $1 \Rightarrow 2$. Sample size $\Rightarrow f_{\epsilon}(\epsilon/2, \delta)$

$2 \Rightarrow 3$ trivial

$3 \Rightarrow 4$ (seen before)

$4 \Rightarrow 5$ (just proved)

$5 \Rightarrow 1$ (last class)

Ocean's Razor