

• Project Version 1 due today

• Theory HW 1 due Thursday

Lecture 13 Recap

Concentration of measure

① $\theta \geq 0, \mathbb{E}\theta = \mu \rightarrow \mathbb{P}(\theta > \mu + \epsilon) \leq \frac{\mu}{\epsilon}$

② $\mathbb{E}\theta = \mu, \text{Var}\theta = \sigma^2 \rightarrow \mathbb{P}(|\theta - \mu| > \epsilon) \leq \frac{\sigma^2}{\epsilon^2}$
 $\mathbb{P}(|\frac{1}{m} \sum_{i=1}^m \theta_i - \mu| > \epsilon) \leq \frac{\sigma^2}{m\epsilon^2}$

③ $\mathbb{E}\theta = \mu, a \leq \theta \leq b \rightarrow \mathbb{P}(|\frac{1}{m} \sum_{i=1}^m \theta_i - \mu| > \epsilon) \leq 2 \cdot e^{-\frac{2m\epsilon^2}{(b-a)^2}}$

④ $\mathbb{E}\theta_i = p_i, \theta_i \in \{0, 1\} \rightarrow \mathbb{P}(\frac{1}{m} \sum_{i=1}^m \theta_i > (1+\delta)\mu) \leq e^{-\frac{\delta^2}{2+\frac{1}{3}\delta} \cdot \mu m}$
 $\mathbb{E} \frac{1}{m} \sum_{i=1}^m \theta_i = \frac{1}{m} \sum_{i=1}^m p_i = \mu$

Recap.

Lecture 10. (ch 2 & 3)

- PAC learning ERM \leftarrow realizable $f \in \mathcal{H}, y \in \{0, 1\}, |\mathcal{H}| < \infty$
- technique: union bound

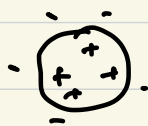
$$m_{\mathcal{H}} = \frac{\log |\mathcal{H}|}{\epsilon}$$

Lecture 11. (ch 4)

- Agnostic PAC learning ERM \leftarrow non-realizable, any bounded loss, $|\mathcal{H}| < \infty$
- technique: Hoeffding's ineq.

$$m_{\mathcal{H}} = \frac{\log |\mathcal{H}|}{\epsilon^2}$$

Q3.3



$$\mathcal{H} = \{h_r: r \geq 0\}, h_r = \mathbb{I}_{\{x: |x| \leq r\}}, |\mathcal{H}| = \infty$$

$$m_{\mathcal{H}} = \frac{\log \frac{1}{\epsilon}}{\epsilon} \rightarrow \text{How?}$$

Chapter 5. No free lunch theorem.

The reason Q.3.3 has low sample complexity is that we had a prior knowledge.

- Realizable: $\left\{ \begin{array}{l} \text{ground truth } D, f \text{ has a circular boundary} \\ H \text{ is a collection of circular decision boundaries.} \end{array} \right.$

fundamental Question:

- is such prior knowledge necessary for PAC learning? **YES.**

Learning Problem \Rightarrow Goal: minimize $L_D(h)$

- fix X, Y
- Nature chooses D_{xy}
- Learner observes $S = \{(x_i, y_i)\}_{i=1}^m \sim D$
- Applies Algo over H .
- Outputs $h \in H$.

Q. Can there be a Universal Algorithm that is PAC learnable for all D ?

Q. Is Prior Knowledge necessary for PAC learning?

- lec 10: realizable w.r.t. $H: |H| < \infty$
- lec 11: $\min_{h \in H} L_D(h)$ is small, $|H| < \infty$
- Q 3.3: realizable w.r.t. parametric family H

Q. How much prior knowledge is good to assume?

a lot : learner knows $D \rightarrow h^*$: Bayes optimal Predictor.



learner knows a class $\left\{ \begin{array}{l} \text{realizable} \\ H = \bigoplus \\ |H| < \infty \end{array} \right.$

No : learner knows nothing \rightarrow [No free lunch theorem]

No-free-lunch theorem

Let Algo be any learning algorithm

$Y = \{0, 1\}$

any X s.t. $|X| \geq 2 \cdot m$

Then exists a distribution D over $X \times Y$ and H s.t.

① $\exists f: X \rightarrow \{0, 1\}$ with $L_D(f) = 0$

② $\mathbb{P}_S (L_D(\text{Algo}(S)) \geq \frac{1}{8}) \geq \frac{1}{7}$

Negative result \rightarrow Constructive proof (Sketch)

since $|X| \geq 2m$, let $C \subseteq X$ s.t. $|C| = 2m$

we observe random m samples in X .

Let H be all possible ways to label C s.t. $|H| = 2^{2m}$
 $= \{h_0, h_1, \dots, h_{2^{2m}-1}\}$

\leftarrow * Very high degrees of freedom in assigning y 's.

* Unstructured

$m=2, |H|=2^{2 \cdot 2} = 16$

$|C|=2m$

$X = \{x_1, x_2, x_3, x_4, \dots\}$
 $H = \{(y_1, y_2, y_3, y_4)\}$
 $\left. \begin{array}{l} 0000, 0001, \\ 0010, 0011, \\ \vdots \end{array} \right\}$

for any algorithm that outputs $h \in \mathcal{H}$,

$\exists D$ s.t. one cannot do better than random guess on $\frac{1}{2}$ of C not observed.

$$\rightarrow \mathbb{E}_S [\underbrace{L_D(h_S)}_{\text{positive R.V.}}] \geq \underbrace{\frac{1}{4}}_{\frac{1}{2} \text{ of samples have error } \frac{1}{2}}$$

\rightarrow This implies the claim by tighter Markov inequality.

$$0 \leq \theta \leq 1.$$

Corollary. If $|\mathcal{X}| = \infty$ & \mathcal{H} is all possible functions from \mathcal{X} to $\{0, 1\}$ then \mathcal{H} is not PAC learnable.

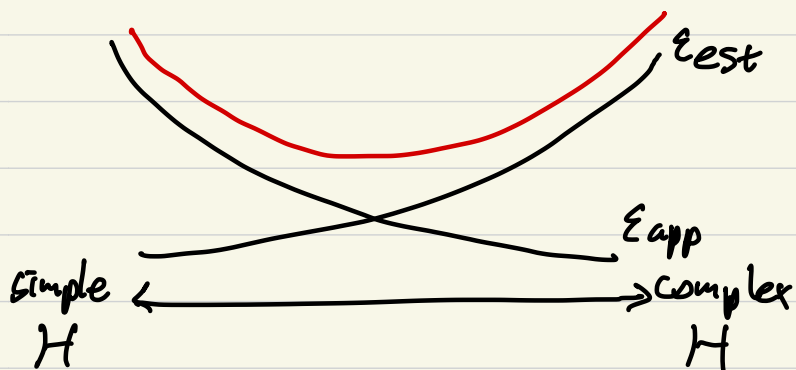
Q. How do we choose the right hypothesis class?

$$L_D(h_S) = \underbrace{\left(L_D(h_S) - \min_{h' \in \mathcal{H}} L_D(h') \right)}_{\epsilon_{\text{est}} = \text{estimation error}} + \underbrace{\min_{h' \in \mathcal{H}} L_D(h')}_{\epsilon_{\text{app}} = \text{approximation error}}$$

\uparrow
 ERM

If $m \uparrow$: $\epsilon_{\text{est}} \downarrow$ no change
 If complexity \uparrow : $\epsilon_{\text{est}} \uparrow$ $\epsilon_{\text{app}} \downarrow$
 \mathcal{H}

Bias-Complexity tradeoff



Chapter 6. VC-dimensions

fundamental Question

- what class \mathcal{H} is PAC learnable?

[Vladimir Vapnik & Alexey Chervonenkis, 1970]

Definition. [Restriction of \mathcal{H} to $C \subseteq \mathcal{X}$]

\mathcal{H} is class of functions: $\mathcal{X} \rightarrow \{0, 1\}$, and

$C = (c_1, c_2, \dots, c_m) \subseteq \mathcal{X}$, then

the restriction of \mathcal{H} to C is

$$\mathcal{H}_C \triangleq \left\{ \tilde{h} : C \rightarrow \{0, 1\} \mid h \in \mathcal{H} \right\}$$
$$c_i \mapsto \tilde{h}(c_i)$$

\tilde{h} can equivalently be represented by an m -dim vector

$$\tilde{h} = (\tilde{h}(c_1), \tilde{h}(c_2), \dots, \tilde{h}(c_m))$$

Definition. [Shattering]

A hypothesis class \mathcal{H} **shatters** a set $C \subseteq \mathcal{X}$

if the restriction of \mathcal{H} to C is the set of all functions from C to $\{0, 1\}$, that is $|\mathcal{H}| = 2^{|C|}$.

example $\mathcal{H} = \{h_a : \mathbb{R} \rightarrow \mathbb{R} = \mathbb{I}\{x \leq a\}\}$ is class of all threshold functions.

① $C = \{c_1\}$

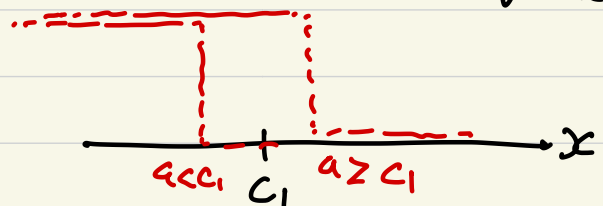
$$\mathcal{H}_C = \{0, 1\}$$

$$h_a(c_1) = 0$$

$$\text{for } a < c_1$$

$$h_a(c_1) = 1$$

$$\text{for } a \geq c_1$$



$\therefore \mathcal{H}$ shatters $C = \{c_1\}$.

$$\textcircled{2} \quad C = \{c_1, c_2\}$$

$$\mathcal{H} = \{(0,0), (1,0), (1,1)\}$$

$(h_1(c_1), h_2(c_2))$

$$\uparrow$$

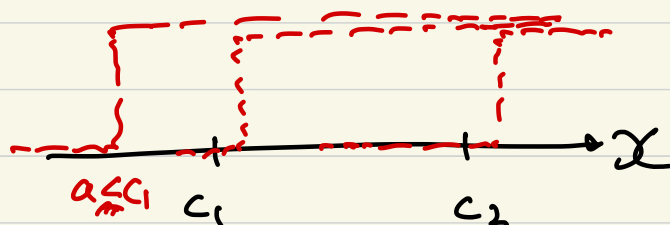
$$c_2 \leq a$$

$$\uparrow$$

$$c_2 \leq a < c_1$$

$$\uparrow$$

$$a < c_1$$



$(1,1)$ is not possible

$\therefore \mathcal{H}$ does not shatter $C = \{c_1, c_2\}$.

Corollary [No-free-lunch theorem restated]

for some $\mathcal{H}, \mathcal{X}, \mathcal{Y} = \{0,1\}$,

if $\exists C \subseteq \mathcal{X}, |C|=2m$, and \mathcal{H} **shatters** C ,

then for any learning algorithm, $\exists D$ & $h \in \mathcal{H}$ s.t.

$$L_D(h) = 0 \quad \text{but} \quad \mathbb{P}(L_D(\text{Algo}(S)) \geq \frac{1}{8}) \geq \frac{1}{7}$$

Definition [VC-dimension]

The **VC-dimension** of \mathcal{H} , defined $VC_{dim}(\mathcal{H})$,

is the maximal size of $C \subseteq \mathcal{X}$ that can be **shattered** by \mathcal{H} .

\uparrow No set of size $VC_{dim}(\mathcal{H}) + 1$ can be shattered.

* Need to show that $\textcircled{1} \exists C$ s.t. $|C|=d$ shattered by \mathcal{H}
 $\textcircled{2} \forall C$ s.t. $|C|=d+1$ not shattered.

Example [Threshold: $\mathcal{H} = \{ \mathbb{I}(x \leq a) \mid a \in \mathbb{R} \}$

$\exists C = \{c_1\}$ can be shattered

$\forall C = \{c_1, c_2\}$ cannot be shattered.

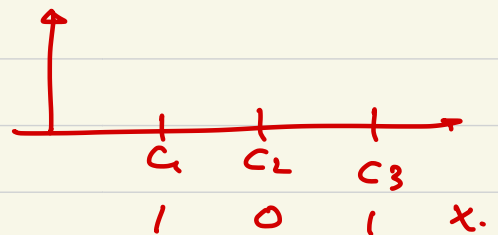
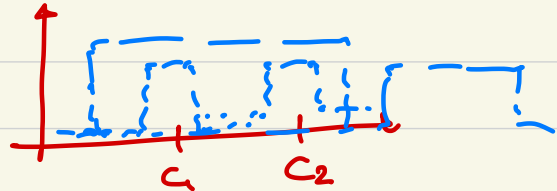
$\rightarrow VCdim(\mathcal{H}) = 1.$

Example [Intervals: $\mathcal{H} = \{ \mathbb{I}(a \leq x \leq b) \mid a, b \in \mathbb{R} \}$

$\exists C = \{c_1, c_2\}$ can be shattered

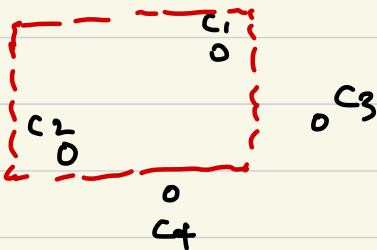
$\forall C = \{c_1, c_2, c_3\}$ cannot be shattered.

$\rightarrow VCdim(\mathcal{H}) = 2.$

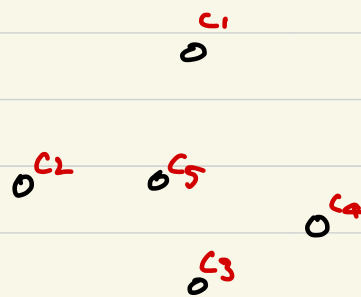


Example [axis aligned rectangles: $\mathcal{H} = \{ \mathbb{I} \{ \begin{matrix} a_1 \leq x_1 \leq b_1 \\ a_2 \leq x_2 \leq b_2 \end{matrix} \} \mid \begin{matrix} a_1 < b_1 \\ a_2 < b_2 \end{matrix} \}$

$VCdim(\mathcal{H}) = 4$



any subset can be in a rectangle



however you place 5 points,
1 of them has to be "inside".
and cannot be a unique \odot .

Example [finite class $|H| < \infty$]

$$VCdim \leq \log_2 |H|$$

Shattering requires at least $2^{|C|}$ hypotheses,

$$2^{|C|} \leq |H|$$

$$|C| \leq \log_2 |H|$$

$$VCdim \leq \log_2 |H|$$

← Tight in the worst case
can be smaller like interval.

Example [number of parameters]

- above examples have # of parameters match $VCdim$.

- not always true: $H = \{h_a: a \in \mathbb{R}\}$, $VCdim = \infty$

$$\begin{array}{c} \text{"} \\ [0.5 \sin(ax)] \end{array}$$

[Theory HW 2]
[Problem 4]