

Given $f: X \rightarrow \{0, 1\}$, and a sample $\{(x_i, y_i)\}_{i=1}^n$

Define

Empirical Loss : $\hat{R}_n(f) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(f(x_i) \neq y_i)$

True Loss : $R(f) = \mathbb{E}_{P_{XY}}[\mathbb{1}(f(x) \neq y)]$

Definition A hypothesis class $\mathcal{H} \subseteq \mathcal{Y}^X$ is PAC-learnable if there is a function $m_{\mathcal{H}}: [0, 1]^2 \rightarrow \mathbb{N}$ and a learning algorithm A with the following property: For every $\epsilon, \delta \in (0, 1)$ and every choice of P_{XY} , when running the algorithm on $m \geq m_{\mathcal{H}}(\epsilon, \delta)$ iid examples from \mathcal{D} , the algorithm returns $h \in \mathcal{H}$ s.t. w/ probability $> 1 - \delta$

$$R(h) \leq \min_{h' \in \mathcal{H}} R(h') + \epsilon$$

\uparrow
returned by
 A

\uparrow best classifier
in \mathcal{H}

Q: What hypothesis classes \mathcal{H} are PAC-learnable?

• $|\mathcal{H}| < \infty$ - PAC learnable

↳ finite

• A: $\hat{h} = \min_{h' \in \mathcal{H}} \hat{R}_S(h')$: ERM

• HW 3 ex 8

$$R(\hat{h}) \leq \min_{h' \in \mathcal{H}} R(h') + \sqrt{\frac{\log |\mathcal{H}| / \delta}{m}}$$

size of the class

↑ # samples

Sketch

by Hoeffding

$\forall h \in \mathcal{H}$

$$\hat{R}_S(h) - R(h) \leq \sqrt{\frac{\log(1/\delta)}{m}}$$

$$\rightarrow \sqrt{\frac{\log(|\mathcal{H}|/\delta)}{m}} < \epsilon$$

$$\rightarrow m_{\mathcal{H}}(\epsilon, \delta) = \frac{\log(|\mathcal{H}|/\delta)}{\epsilon^2}$$

Q: Is $\mathcal{H} = \mathcal{Y}^X$, $\mathcal{Y} = \{0, 1\}$, all functions,
PAC-learnable?

→ NO

No Free Lunch Theorem

Consider $m \in \mathbb{N}$, X s.t. $|X| \geq 2m$,
any algorithm A which outputs

$A(S)$, given $S \subseteq X$, $|S| = m$. Then
there exists $f: X \rightarrow [0, 1]$ s.t.

w/ probability $\geq \frac{1}{7}$

$R_{P_f}(A(S)) \geq \frac{1}{8}$

↑ randomness from
sample

Corollary \mathcal{Y}^X is not PAC-learnable

Given $f: X \rightarrow \{0, 1\}$ let

$$P_f((x, y)) = \begin{cases} 1/|X| & y = f(x) \\ 0 & y \neq f(x) \end{cases}$$

Proof WLOG $|X| = 2^m$.

$$|Y^X| = 2^{2^m} \rightarrow \text{uniform dist on } Y^X$$

Will Show ~~*~~

$$\mathbb{E}_{f \sim Y^X} \mathbb{E}_{S \sim P_f^m} R(A(S)) \geq \frac{1}{4}$$

↑
sample of
size m
labeled by f

$$\rightarrow \exists f \in Y^X, \mathbb{E}_{S \sim P_f^m} R(A(S)) \geq \frac{1}{4}$$

↓
Markov's inequality

Prove ~~*~~.

$$\mathbb{E}_{f \sim Y^X} \mathbb{E}_{S \sim P_f^m} \mathbb{E}_{x \sim P_f} [\mathbb{1}\{A(S)(x) \neq f(x)\}]$$

$$= \mathbb{E}_S \mathbb{E}_{F \sim \mathcal{P}_F} \mathbb{E}_{x \sim \mathcal{P}_F} [\mathbb{1}(A(S)(x) \neq f(x))]]$$

$$= \mathbb{E}_S \mathbb{E}_F [\mathbb{1}(A(S)(x) \neq f(x) \mid x \in S)] \mathbb{P}(x \in S)$$

$$+ \mathbb{E}_S \mathbb{E}_F [\mathbb{1}(A(S)(x) \neq f(x) \mid x \notin S)] \mathbb{P}(x \notin S)$$

$$\geq \frac{1}{2} \quad \circ \quad \frac{1}{2}$$

$$= \frac{1}{4}$$

• $h \in H$, if H is PAC-learnable

$$\frac{1}{8} \leq R(h) = \underbrace{\min_{h' \in H} R(h')}_{\text{approximation error}} + \underbrace{R(h) - \min_{h' \in H} R(h')}_{\substack{\rightarrow \text{w/ enough samples} \\ < \epsilon \\ \uparrow \\ \text{estimation error}}}$$

• Finite Classes

$$R(h) \approx \underbrace{\min_{h' \in H} R(h')}_{\substack{\text{bias} \\ m \uparrow}}$$

$$\sqrt{\frac{\log(|H|/\delta)}{m}}$$

Complexity

\downarrow
0

What if $\mathcal{H} = \infty$?

Definition The shattering coefficient of \mathcal{H} is

$$S(\mathcal{H}, n) = \max_{x_1, \dots, x_n} |\{ (f(x_1), \dots, f(x_n)) : f \in \mathcal{H} \}|$$

↓
bounded above
by 2^n

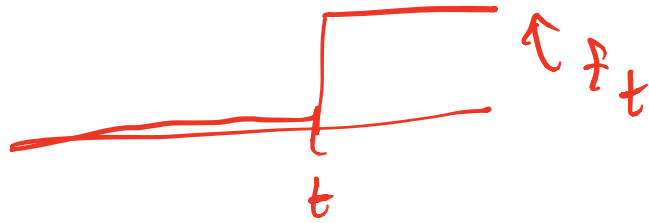
Definition The VC-dim of \mathcal{H} is the maximum $k \in \mathbb{N}$ s.t.

$$S(\mathcal{H}, k) = 2^k$$

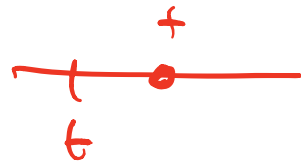
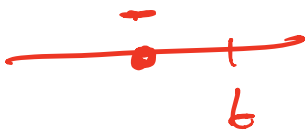
→ get all possible labelings of k points.

Ex 1

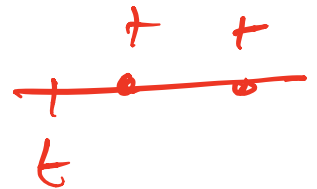
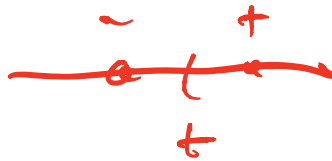
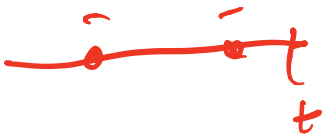
$$\mathcal{H}_1 = \{ F_t = \mathbb{1}\{x \geq t\}, x, t \in [0, 1] \}$$



• $n=1$ $S(\mathcal{H}_1, 1) = 2$



• $n=2$ $S(\mathcal{H}_1, 2) = 3$



• n points $S(\mathcal{H}_1, n) = n+1$



$$VC(\mathcal{H}_1) = 1$$

Ex 2

$$X = \mathbb{R}^2$$

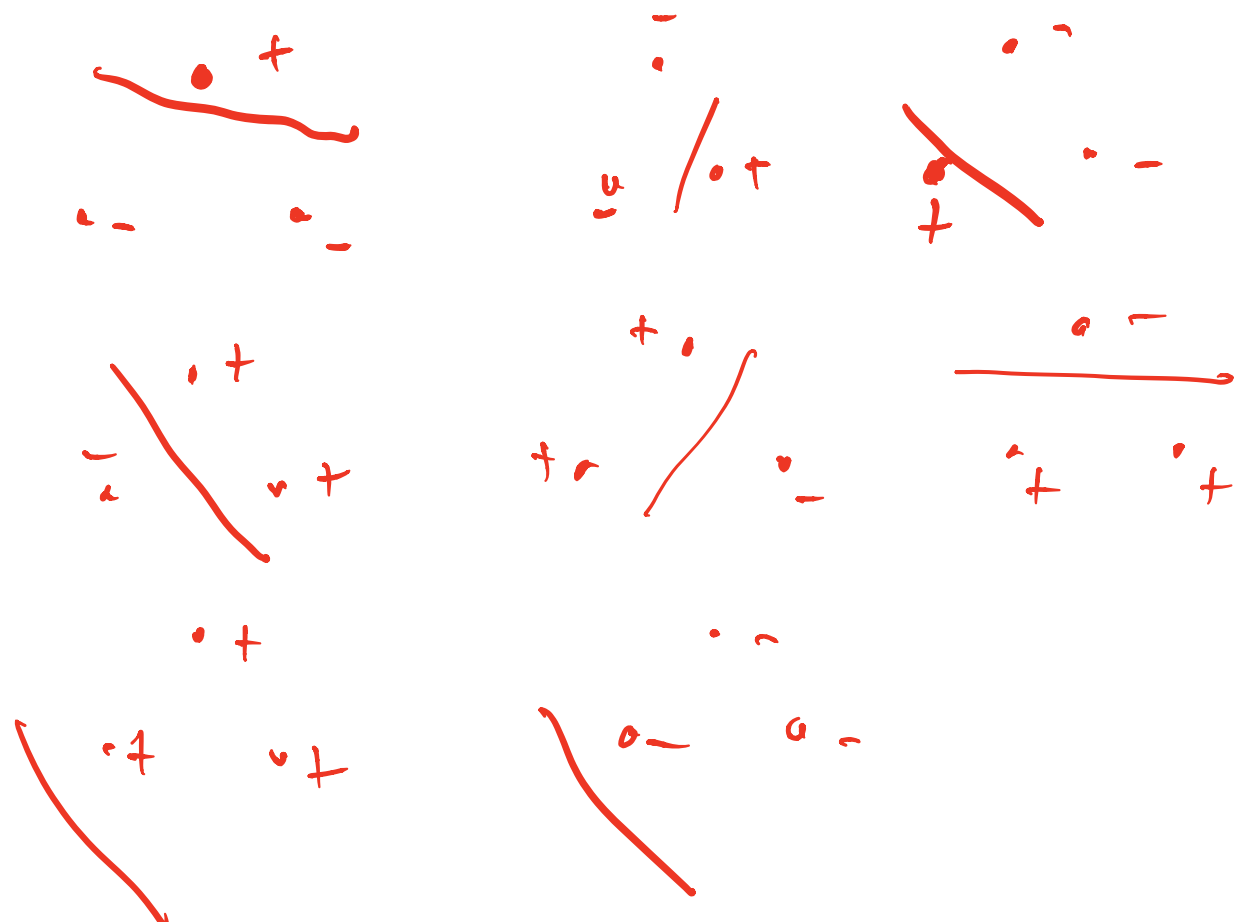
$\mathcal{H}_2 = \{ \text{hyperplane classifications} \}$

$n=1$

$S(\mathcal{H}_2, 1) = 2$

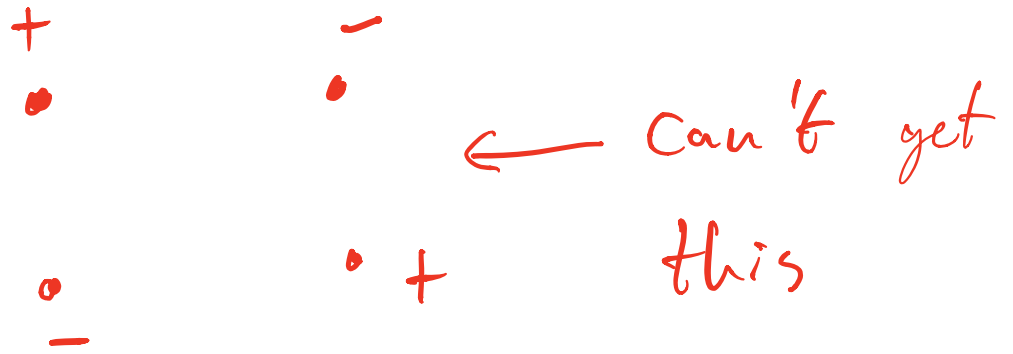


$n=3$



$S(\mathcal{H}_2, 3) = 8 = 2^3$

$n=4$
→



$$VC(H_2) = 3$$

→ In general

$H_d = \{\text{hyperplanes in } \mathbb{R}^d\}$

$$VC(H_d) = d + 1, \quad d \geq 2$$

Sauer's Lemma

$$S(H, n) \leq (n+1)^{VC(H)}$$

→ Effective # of classifiers on n points is same as # labelings





Theorem Vapnik-Chervonenkis

$$R(h) \leq \min_{h' \in H} R(h') + \sqrt{\frac{\log(n+1)^{VC(H)} / \epsilon}{n}}$$

$$= \min_{h' \in H} R(h') + \sqrt{\frac{VC(H) \log n / \epsilon}{n}}$$



→ hyperplanes in d dimensions

→ $\frac{d \log(1/\epsilon)}{\epsilon^2}$ samples

can learn a hyperplane

Fundamental Theorem of Learning

IF H has finite VC dim



PAC Learnable

