

## Lower bounds on sample complexity

• Let  $H \subseteq Y^X$  be a hypothesis class

• Let  $d = VC(H)$

• We want show that

$$m(H, \epsilon, \delta) \gtrsim \frac{d}{\epsilon}$$

Theorem: Let  $H$  be a hypothesis class over a domain  $X$  such that

$d = VC(H) \geq 2$ . Let  $\epsilon \in (0, \frac{1}{8})$ .

Let  $A: (X \times Y)^* \rightarrow Y^X$  be any learning algorithm. There exists

distribution  $D$  over  $X$  and  $h^* \in H$

such that if  $S$  is an i.i.d.

sample from  $D$  labeled by  $h^*$

then  $\text{err}_{D, h^*}(A(S)) > \epsilon$  with probability at least  $\frac{1}{100}$  unless

$$m \geq \frac{d-1}{32\epsilon}.$$

Proof:

• Since  $VC(H) = d$ , there exists a shattered set of size  $d$ .

• Let  $C = \{x_0, x_1, x_2, \dots, x_{d-1}\}$

be a shattered set.

• Define  $D$  as follows:

$$\sim 1 - \delta\epsilon \quad \text{if } x = x_0$$

$$D(x) = \begin{cases} \frac{8\varepsilon}{d-1} & \text{if } x \in \{x_1, x_2, \dots, x_{d-1}\} \\ 0 & \text{if } x \notin C \end{cases}$$

•

$x_0$

•

$x_1$

•

$x_2$

...

•

$x_{d-1}$

• Let  $U$  be an i.i.d. sample from  $D$  of size  $m < \frac{d-1}{32\varepsilon}$ .

• Let  $Z$  be the number of samples

in  $\{x_1, x_2, \dots, x_{d-1}\}$

(That is, samples other than  $x_0$ .)

- $Z = \sum_{i=1}^m z_i$  where  $z_1, z_2, \dots, z_m$

are i.i.d.  $\Pr[z_i = 1] = \delta \epsilon = p$

- With "high" probability  $Z$  is "small".

We use Chernoff bound:

$$\Pr[Z \geq 2mp] \leq \left( \frac{e}{2^2} \right)^{mp}$$

- $2mp = 2m \cdot \delta \epsilon < 2 \cdot \frac{d-1}{32\epsilon} \cdot \delta \epsilon = \frac{d-1}{2}$

- $\left( \frac{e}{2^2} \right)^{mp} = \left( \frac{\sqrt{e}}{2} \right)^{2mp} < \left( \frac{\sqrt{e}}{2} \right)^{\frac{d-1}{2}} \leq \left( \frac{\sqrt{e}}{2} \right)^{\frac{1}{2}} \leq \frac{93}{100}$

- So with probability more than

$$\frac{7}{100}, \quad z < \frac{d-1}{2}$$

---

- The rest of the proof is similar to the no free lunch theorem.
- Since  $C$  is shattered, there exists  $2^{|C|}$  functions  $h_1, h_2, \dots, h_{2^{|C|}} \in H$  realizing all possible behaviors on  $C$ .
- Let  $S_i$  be  $U$  labeled by  $h_i$ .
- Let  $T = \{x \in C : x \notin U, x \neq x_0\}$
- $|T| \geq (d-1) - z$

• Condition on the event

$$|T| > \frac{d-1}{2}$$

$$\bullet \Pr \left[ |T| > \frac{d-1}{2} \right] \geq 1 - \frac{93}{100} = \frac{7}{100}$$

$$\frac{1}{2^{|d|}} \sum_{i=1}^{2^{|d|}} \mathbb{E} \left[ \text{err}_{D, h_i} (A(s_i)) \mid T > \frac{d-1}{2} \right]$$

$$\geq \frac{1}{2^{|d|}} \sum_{i=1}^{2^{|d|}} \mathbb{E} \left[ \frac{8\epsilon}{d-1} \sum_{x \in T} \mathbb{1}[A(s_i)(x) \neq h_i(x)] \mid T > \frac{d-1}{2} \right]$$

$$= \frac{8\epsilon}{d-1} \mathbb{E} \left[ \sum_{x \in T} \frac{1}{2^{|d|}} \sum_{i=1}^{2^{|d|}} \mathbb{1}[A(s_i)(x) \neq h_i(x)] \mid T \right] \mid T > \frac{d-1}{2}$$

$$\rightarrow \frac{8\epsilon}{d-1} \mathbb{E} \left[ |T| \mid T > \frac{d-1}{2} \right]$$

$$\leq d-1 \cdot 4 \left[ \frac{1}{2} \right] \cdot \left[ \frac{1}{4} \right]$$

$$> \frac{8\varepsilon}{d-1} \cdot \frac{d-1}{4}$$

$$= 2\varepsilon$$

So there exists  $h^* \in H$

$$\mathbb{P} \left[ \text{err}_{h^*, D} (A(S)) \mid |T| > \frac{d-1}{2} \right] \geq 2\varepsilon$$

$$\text{Let } \gamma = 1 - \text{err}_{h^*, D} (A(S))$$

$$\mathbb{P} \left[ \text{err}_{h^*, D} (A(S)) < \varepsilon \mid |T| > \frac{d-1}{2} \right]$$

$$P\left[ \left| \frac{1}{n} \sum_{i=1}^n X_i - \mu \right| \geq \frac{\epsilon}{2} \right]$$

$$= P\left[ Y \geq 1 - \epsilon \mid |T| > \frac{d-1}{2} \right]$$

$$\leq \frac{E\left[ Y \mid |T| > \frac{d-1}{2} \right]}{1 - \epsilon}$$

$$\leq \frac{1 - 2\epsilon}{1 - \epsilon}$$

$$\leq \frac{1 - 2 \cdot \frac{1}{8}}{1 - \frac{1}{8}}$$

$$\boxed{\epsilon \leq \frac{1}{8}}$$

$$= \frac{3/4}{7/8}$$



$$= \frac{1}{7}$$

$$\bullet \Pr \left[ \text{err}_{D, h^*}(A(s)) > \epsilon \mid |T| > \frac{d-1}{2} \right] \geq \frac{1}{7}$$

$$\bullet \Pr \left[ \text{err}_{D, h^*}(A(s)) > \epsilon \right]$$

$$\geq \Pr \left[ |T| > \frac{d-1}{2} \right] \cdot \Pr \left[ \text{err}_{D, h^*}(A(s)) \mid |T| > \frac{d-1}{2} \right]$$

$$\geq \frac{7}{100} \cdot \frac{1}{7}$$

$$= \frac{1}{100}$$

~~TCMA~~

14