

(1)

## Uniform Convergence -

The goal is to prove that finite  $\mathcal{H}$  are agnostic PAC learnable by the ERM algo.

$$A(S) \in \arg\min L_S(h).$$

The general idea is to show that if  $|H| < +\infty$  and  $|S|$  is large enough (how large is to be determined) then  $L_S(h)$  is sufficiently close to  $L_{\mathcal{D}}(h)$  uniformly in  $h \in H$ .

This will allow to check the condition:

$$\Pr_S \left\{ L_{\mathcal{D}}(A(S)) \leq \min_{h \in H} L_{\mathcal{D}}(h) + \epsilon \right\} \geq 1 - \delta.$$

for  $|S|$  large enough  $|S| = m > M_{\mathcal{D}}(\epsilon, \delta)$  for some  $M_{\mathcal{D}}(\epsilon, \delta)$ .

(2)

Definition: a training set  $S$  is  $\epsilon$ -representative

$$\text{if } \forall h \in \mathcal{H} : |L_S(h) - L_{\mathcal{D}}(h)| \leq \epsilon .$$

In other words  $L_S(h)$  is uniformly close to  $L_{\mathcal{D}}(h)$

for all  $h \in \mathcal{H}$ . [“closeness” does not depend on  $h$  ].

Lemma: Let  $S$  be  $\frac{\epsilon}{2}$ -representative. Let

$$h_S = \text{ERM}_{\mathcal{H}}(S) \equiv \underset{h \in \mathcal{H}}{\arg \min} L_S(h) \text{ be an}$$

empirical risk minimizer. Then

$$L_{\mathcal{D}}(h_S) \leq \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \epsilon .$$

Proof:  $L_{\mathcal{D}}(h_S) \leq L_S(h_S) + \frac{\epsilon}{2}$ : by  $\frac{\epsilon}{2}$ -repr def of  $h_S$

$$L_{\mathcal{D}}(h) - \epsilon \leq L_S(h) \leq L_{\mathcal{D}}(h) + \epsilon \quad \leq L_S(h) + \frac{\epsilon}{2} , \text{ by } h_S \text{ minimizer}$$

$$\leq L_{\mathcal{D}}(h) + \frac{\epsilon}{2} + \frac{\epsilon}{2} : \text{by } \frac{\epsilon}{2} \text{-repr def of } h_S$$

$$= L_{\mathcal{D}}(h) + \epsilon .$$

(3)

Lemma: Let  $|H| < +\infty$  a finite hyp. class.

and let  $\ell : \mathcal{H} \times \mathbb{Z} \rightarrow \{0, 1\}$  the 0-1

loss function (for classification say). Then  $\ell$  has the

uniform convergence property:

$$\left\{ \begin{array}{l} \Pr_S \{ \text{there is } h \in \mathcal{H} : |L_S(h) - L_\Delta(h)| \leq \epsilon \} \geq 1 - \delta \\ \Pr_S \{ S \text{ is } \epsilon\text{-representable} \} \geq 1 - \delta \end{array} \right.$$

for training sets  $S$  with size  $m \geq m_{\mathcal{H}}^{uc}(\epsilon, \delta)$

where,

$$m_{\mathcal{H}}^{uc}(\epsilon, \delta) = \left\lceil \frac{\log \left( \frac{2|\mathcal{H}|}{\delta} \right)}{2\epsilon^2} \right\rceil$$

Note: The uniform convergence property just means that the empirical risk  $L_S(h)$  is uniformly close to the true risk  $L_\Delta(h)$  with high probability w.r.t  $S$  if  $|S|$  is large enough.

(4)

## Proof of Lemma.

We show the statement for the converse event:

$$\Pr \{ \exists h \in \mathcal{H} : |L_S(h) - L_{\mathcal{D}}(h)| > \epsilon \} \leq \delta$$

$$\Pr \{ \exists h \in \mathcal{H} : |L_S(h) - L_{\mathcal{D}}(h)| > \epsilon \} = \Pr \left\{ \bigcup_{h \in \mathcal{H}} |L_S(h) - L_{\mathcal{D}}(h)| > \epsilon \right\}$$

$$\leq \sum_{\text{union bound}} \sum_{h \in \mathcal{H}} \Pr \{ |L_S(h) - L_{\mathcal{D}}(h)| > \epsilon \}$$

$$= \sum_{h \in \mathcal{H}} \Pr \left\{ \underbrace{\left| \frac{1}{m} \sum_{i=1}^m l(h, z_i) - \underbrace{\mathbb{E}_{\mathcal{D}} l(h, z)} \right|}_{\text{empirical mean}} > \epsilon \right\}$$

true mean

$$\leq \sum_{h \in \mathcal{H}} 2 \exp(-2m\epsilon^2) \quad \begin{array}{l} \text{by Hoeffding's} \\ \text{inequality for} \\ z_i \stackrel{iid}{\sim} \mathcal{D} \end{array}$$

$$= |\mathcal{H}| 2 e^{-2m\epsilon^2}.$$

Imposing  $|\mathcal{H}| 2 e^{-2m\epsilon^2} \leq \delta$  implies that

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



## Hoeffding inequality : Reminder or see exercise.

Let  $\delta_1, \dots, \delta_m$  be iid with  $E(\delta_i) = \mu$

and  $a \leq \delta_i \leq b$ . (RV have bounded support)

$$\Pr \left\{ \left| \frac{1}{m} \sum_{i=1}^m \delta_i - \mu \right| > \epsilon \right\} \leq 2 e^{-2m \frac{\epsilon^2}{(b-a)^2}}$$

Note: for us  $|a-b|=1$  for 0-1 loss

since  $\delta_i = \ell(h, z_i) \in \{0, 1\}$ .

but proof works also for any bounded loss fn.

(6)

Now we arrive at the main result of this chapter:

Corollary: If  $\mathcal{H}$  is finite then it is

agnostic PAC learnable with the ERM rule with

sample complexity  $m_{\text{err}}^{\text{UC}}(\epsilon, \delta) = \left\lceil \frac{1}{2\epsilon^2} \log \frac{|\mathcal{H}|}{2\delta} \right\rceil$ .

Proof: for  $m \geq m_{\text{err}}^{\text{UC}}(\epsilon, \delta)$  we have that

$$\Pr \left\{ S \text{ is } \frac{\epsilon}{2}\text{-repr} \right\} \geq 1 - \delta$$

We also proved that:

$$S \text{ is } \frac{\epsilon}{2}\text{-repr} \implies L_{\mathcal{D}}(h_S) \leq \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \epsilon$$

$$\text{so } \Pr \left\{ S \text{ is } \frac{\epsilon}{2}\text{-repr} \right\} \leq \Pr \left\{ L_{\mathcal{D}}(h_S) \leq \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \epsilon \right\}$$

Thus we obtain

$$\Pr \left\{ L_{\mathcal{D}}(h_S) \leq \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \epsilon \right\} \geq 1 - \delta$$