

# Machine Learning

## Uniform Convergence

Fabio Vandin

November 6<sup>th</sup>, 2023

# When is an Hypothesis Class PAC Learnable?

Previously seen result: for binary classification with

- realizability assumption
- 0-1 loss

any finite hypothesis class is PAC learnable by ERM.

What about the more general PAC learning model we have seen last? Recall the (agnostic) PAC learnability for general loss:

$\subset X \times Y$

## Definition

A hypothesis class  $\mathcal{H}$  is *agnostic PAC learnable* with respect to a set  $Z$  and a loss function  $\ell : \mathcal{H} \times Z \rightarrow \mathbb{R}_+$  if there exist a function  $m_{\mathcal{H}} : (0, 1)^2 \rightarrow \mathbb{N}$  and a learning algorithm such that for every  $\delta, \varepsilon \in (0, 1)$ , for every distribution  $\mathcal{D}$  over  $Z$ , when running the learning algorithm on  $m \geq m_{\mathcal{H}}(\varepsilon, \delta)$  i.i.d. examples generated by  $\mathcal{D}$  the algorithm returns a hypothesis  $h$  such that, with probability  $\geq 1 - \delta$  (over the choice of the  $m$  training examples):

$$L_{\mathcal{D}}(h) \leq \min_{h' \in \mathcal{H}} L_{\mathcal{D}}(h') + \varepsilon$$

where  $L_{\mathcal{D}}(h) = \mathbb{E}_{z \sim \mathcal{D}}[\ell(h, z)]$

# Uniform Convergence and Learnability

**Uniform convergence:** the empirical risks (training error) of *all* members of  $\mathcal{H}$  are good approximations of their true risk (generalization error).

## Definition

A training set  $S$  is called  $\varepsilon$ -representative (w.r.t. domain  $Z$ , hypothesis class  $\mathcal{H}$ , loss function  $\ell$ , and distribution  $\mathcal{D}$ ) if

$$\forall h \in \mathcal{H}, |L_S(h) - L_{\mathcal{D}}(h)| \leq \varepsilon$$



$$L_S(h) - \varepsilon \leq L_{\mathcal{D}}(h) \leq L_S(h) + \varepsilon$$

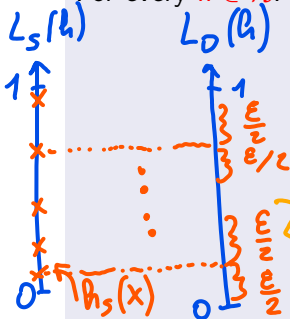
## Proposition

Assume that training set  $S$  is  $\frac{\varepsilon}{2}$ -representative (w.r.t. domain  $Z$ , hypothesis class  $\mathcal{H}$ , loss function  $\ell$ , and distribution  $\mathcal{D}$ ). Then, any output of  $\text{ERM}_{\mathcal{H}}(S)$  (i.e., any  $h_S \in \arg \min_{h \in \mathcal{H}} L_S(h)$ ) satisfies

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

## Proof.

For every  $h \in \mathcal{H}$ :



$$\begin{aligned}
 L_{\mathcal{D}}(h_S) &\leq L_S(h_S) + \frac{\varepsilon}{2} && \text{ } S \text{ is } \frac{\varepsilon}{2} \text{ representative} \\
 &\leq L_S(h) + \frac{\varepsilon}{2} && \text{ } h_S \text{ picked by ERM} \\
 &\leq L_{\mathcal{D}}(h) + \frac{\varepsilon}{2} + \frac{\varepsilon}{2} && \text{ } S \text{ is } \frac{\varepsilon}{2} \text{ representative} \\
 &= L_{\mathcal{D}}(h) + \varepsilon
 \end{aligned}$$

the ERM hypothesis is at most  $\varepsilon$  different from the optimal  $\square$

Uniform convergence depends on training set: when do we have uniform convergence?

### Definition

A hypothesis class  $\mathcal{H}$  has the *uniform convergence property* (w.r.t. a domain  $Z$  and a loss function  $\ell$ ) if there exists a function  $m_{\mathcal{H}}^{UC} : (0, 1)^2 \rightarrow \mathbb{N}$  such that for every  $\varepsilon, \delta \in (0, 1)$  and for every probability distribution  $\mathcal{D}$  over  $Z$ , if  $S$  is a sample of  $m \geq m_{\mathcal{H}}^{UC}(\varepsilon, \delta)$  i.i.d. examples drawn from  $\mathcal{D}$ , then with probability  $\geq 1 - \delta$ ,  $S$  is  $\varepsilon$ -representative.

### Proposition

If a class  $\mathcal{H}$  has the uniform convergence property with a function  $m_{\mathcal{H}}^{UC}$  then the class is agnostically PAC learnable with the sample complexity  $m_{\mathcal{H}}(\varepsilon, \delta) \leq m_{\mathcal{H}}^{UC}(\varepsilon/2, \delta)$ . Furthermore, in that case the  $\text{ERM}_{\mathcal{H}}$  paradigm is a successful agnostic PAC learner for  $\mathcal{H}$ .

What classes of hypotheses have uniform convergence?

# Finite Classes are Agnostic PAC Learnable

We prove that finite sets of hypotheses are agnostic PAC learnable under some restriction for the loss.

## Proposition

Let  $\mathcal{H}$  be a finite hypothesis class, let  $Z$  be a domain, and let  $\ell : \mathcal{H} \times Z \rightarrow [0, 1]$  be a loss function. Then:

- $\mathcal{H}$  enjoys the uniform convergence property with sample complexity

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

different

- $\mathcal{H}$  is agnostically PAC learnable using the ERM algorithm with sample complexity

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

Idea of the proof:

- 1 prove that uniform convergence holds for a finite hypothesis class
- 2 use previous result on uniform convergence and PAC learnability

Useful tool: Hoeffding's Inequality

### Hoeffding's Inequality

Let  $\theta_1, \dots, \theta_m$  be a sequence of i.i.d. random variables and assume that for all  $i$ ,  $\mathbb{E}[\theta_i] = \mu$  and  $\mathbb{P}[a \leq \theta_i \leq b] = 1$ . Then, for any  $\varepsilon > 0$

$$\mathbb{P} \left[ \left| \frac{1}{m} \sum_{i=1}^m \theta_i - \mu \right| > \varepsilon \right] \leq 2e^{-\frac{2m\varepsilon^2}{(b-a)^2}}$$

average of observations

expectation

## Proof (see also the book)

Fix  $\epsilon, \delta$  in  $(0,1)$  we need a sample size  $m$  such that, for any  $D$ , with probability  $\geq 1-\delta$  (over the choice of  $S=(z_1, z_2, \dots, z_m)$ ) we have for all  $h$  in  $H$

$$|L_S(h) - L_D(h)| \leq \epsilon$$

that is:  $\mathbb{P}(\{S \mid \exists h \in H, |L_S(h) - L_D(h)| > \epsilon\}) \leq \delta$

CONTINU A









# Bibliography

[UML] Chapter 4