

LEARNING VIA UNIFORM CONVERGENCE

Earlier we showed finite H class is PAC LEARNABLE

What about AGNOSTIC PAC LEARNABILITY?

Here we will show via the tool of uniform convergence.

We define below condition on training set for which all members of H have their risks close to true risk more or less.

(ϵ -representative sample) A training set S is called ϵ -representative (w.r.t. $\mathcal{Z}, H, \ell, \mathcal{D}$) if

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



LEMMA: If S is $\frac{\epsilon}{2}$ -representative sample.

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

Output of ERM Algorithm

i.e. $h_S = \underset{h \in H}{\operatorname{argmin}} L_S^{\text{ERM}}(h)$

Proof: $\left. \begin{array}{l} \text{(\frac{\epsilon}{2}\text{-representative} \\ \text{sample})} \end{array} \right\}$

$$L_D(h_S) \leq L_S(h_S) + \frac{\epsilon}{2}$$

ERM
Minimizer
= h_S

$$\leq L_S(h_S) + \frac{\epsilon}{2} \leq L_D(h) + \epsilon$$

$\left. \begin{array}{l} \text{(\frac{\epsilon}{2}\text{-representative sample})} \end{array} \right\}$

$$\therefore L_D(h_S) \leq L_D(h) + \epsilon \quad (\text{any } h \in H)$$

$$\Rightarrow L_D(h_S) \leq \min_{h \in H} L_D(h) + \epsilon$$

NOTE: Thus in order to infer Agnostic PAC learnability, it suffices to show that i.i.d. samples are $\frac{\epsilon}{2}$ -representative w.r.t. H, Z, ℓ, D with high probability ($\geq 1 - \delta$)

formalize this for hypothesis class

(UNIFORM CONVERGENCE PROPERTY [UCP])

A hypothesis class H has uniform convergence property (w.r.t. Z, ℓ)
 $\forall \epsilon, \delta \in (0, 1)$ and $\forall D$
 $\exists m_{H, UC}^{\epsilon, \delta} \in \mathbb{N}$ s.t. w.p. $\geq 1 - \delta$, a sample sequence S of m

$(m \geq m_{H, UC}^{\epsilon, \delta})$ i.i.d. samples from D is ϵ -representative (w.r.t. H, Z, ℓ, D)

Consequence of UCP

If a class H has UCP with a function m_H^{UC} then the class is agnostic PAC learnable with

- ① Sample Complexity $m_H(\epsilon, \delta) \leq m_H^{UC}(\frac{\epsilon}{2}, \delta)$
- ② A is ERM_H paradigm

Next Big Question: Are finite hypothesis classes H agnostic PAC learnable?

(Equivalently: Does UCP hold for finite hypothesis class H ?)

CLAIM: Yes, UCP holds for finite hypothesis class H .

PROOF: We have to find a sample size m such that
(fix $\epsilon, \delta \in (0, 1)$)

$$\mathbb{P}_{x, y \sim D, S \sim D^m} \left(\left\{ S : \forall h \in H, |L_S(h) - L_D(h)| > \epsilon \right\} \right) < \delta$$

$$\leq \sum_{h \in H} \mathbb{P}_{S \sim D^m} \left(\left\{ S : |L_S(h) - L_D(h)| > \epsilon \right\} \right)$$

$$L_S(h) = \frac{1}{m} \sum_{i=1}^m \ell(h, z_i)$$

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

(Hoeffding's Inequality) Let $\theta_1, \theta_2, \dots, \theta_m$ be a sequence of i.i.d. r.v.s and assume that $\forall i, \mathbb{E}[\theta_i] = \mu$ and $\mathbb{P}[a \leq \theta_i \leq b] = 1$. Then, $\forall \varepsilon > 0$

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

In our problem (w.l.o.g. ^{after normalizing} we can assume $\ell \in [0, 1]$)

$$\Rightarrow a = 0, b = 1, \mu = L_D(h), \theta_i = \ell(h, z_i)$$

$$\therefore \mathbb{P}_{S \sim \mathcal{D}^m} \left(\{S : |L_S(h) - L_D(h)| > \varepsilon\} \right)$$

$$\leq 2 \exp\left(-\frac{2m\varepsilon^2}{1}\right)$$

$$\Rightarrow \mathbb{P}_{S \sim \mathcal{D}^m} \left(\{S : \forall h \in \mathcal{H}, |L_S(h) - L_D(h)| > \varepsilon\} \right)$$

$$\leq 2|\mathcal{H}| \exp(-2m\varepsilon^2)$$

$$\leq \delta \quad \text{if } m \geq \frac{\log(2|\mathcal{H}|/\delta)}{2\varepsilon^2}$$

$\Rightarrow H$ (finite hypothesis class) w.r.t to domain Z and $\ell(\cdot \in [0,1])$ has UCP with

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

$\Rightarrow H$ is agnostic PAC learnable and ERM_H is the learning algorithm with sample complexity

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

Note: A general H with d parameters can be represented (and approximated) after quantizing each parameter into 64 bits as a finite class with 2^{64d} members.

It is then Agnostic PAC learnable with sample complexity $\leq \left\lceil \frac{128d + 2\log 2/\delta}{\epsilon^2} \right\rceil$

(drawback depends on d)
and machine representation