# Non Uniform Learnability and learning by Structural Risk Minimization.

Recap till now we have seen :

## Notion of PAC learnability :

$\mathcal{H}$ is agnostic PAC learnable if we can find $A$ & $m_{\mathcal{H}}(\epsilon, \delta)$ s.t

$\forall (\epsilon, \delta) \in (0,1)^2$ and $\forall \mathcal{D}$ on $\mathcal{X} \times \mathcal{Y}$ for $m \geq m_{\mathcal{H}}(\epsilon, \delta)$

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

## Uniform convergence property :

for $m \geq m_{\mathcal{H}}(\epsilon, \delta)$

$$\Pr_{S \sim \mathcal{D}^m} \left\{ \sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)| \leq \epsilon \right\} \geq 1 - \delta$$

implies PAC learnability through Empirical Risk Minimization $A^{ERM}(S) = \arg\min_h L_S(h)$ .

finally we saw that unif conv prop holds
iff VC dim $(\mathcal{H})$ is finite and also

$$m_{ve}(\varepsilon, \delta) \underset{\sim}{\approx} \frac{VC\,dim(\mathcal{H}) + \log 1/\delta}{\varepsilon^2} \quad .$$

⊨

There are natural relaxed notions of
learnability for infinite classes and associated
criteria as well as algorithms.

Here we explore one such notion:

Notion of non-uniform learnability and
associated Structurel Risk Minimization algo.

Main idea of non-uniform learnability is to relax fact that $m_{\mathcal{H}}$ depends only on $\epsilon, \delta$. We may make it $h$-dependent: indeed some hypothesis sets might require more or less samples to be learned.

Definition. $\mathcal{H}$ is non-uniform learnable if $\exists A(S)$ and $m_{\mathcal{H}}^{NUL}(\epsilon, \delta, h)$ such that $\forall (\epsilon, \delta) \in (0,1)^2$ and $\forall \mathcal{D}$ we have:

$$\mathbb{P}_{S \sim \mathcal{D}^m} \left\{ L_{\mathcal{D}}(A(S)) \leq \min_{\substack{h \ s.t \\ m \geq m_{\mathcal{H}}^{NUL}(\epsilon, \delta, h)}} L_{\mathcal{D}}(h) + \epsilon \right\} \geq 1 - \delta$$

Meaning: Fix $m = |S|$. Compare $L_{\mathcal{D}}(A(S))$ with best hypothesis in a subset of $\mathcal{H}$ namely such that $m_{\mathcal{H}}^{NUL}(\epsilon, \delta, h) \leq m$. (instead of comparing $L_{\mathcal{D}}(A(S))$ with best hypothesis in all of $\mathcal{H}$).

Obviously we have:

$$L_{\mathcal{D}}(A(S)) \leq \min_{\substack{h \text{ s.t} \\ m \geq m_{\mathcal{H}}^{NUL}(\varepsilon, \delta, h)}} L_{\mathcal{D}}(h) + \varepsilon \implies L_{\mathcal{D}}(A(S)) \leq \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \varepsilon$$

Thus    Non-Unif Learnability $\implies$ Agnostic PAC Learnability

## Main Theorem:

An hypothesis class $\mathcal{H}$ is non-uniformly learnable if and only if it is a countable union $\mathcal{H} = \bigcup_{m \in \mathbb{N}} \mathcal{H}_m$ of agnostic PAC learnable classes $\mathcal{H}_m$.

- Proof of $\mathcal{H}$ NUL $\implies$ $\mathcal{H} = \bigcup \mathcal{H}_m$ with $\mathcal{H}_m$ learnable is easy and we do it first.

- Proof of converse $\mathcal{H} = \bigcup \mathcal{H}_m$, $\mathcal{H}_m$ learnable $\implies$ $\mathcal{H}$ NUL is more involved and requires notion of Structural Risk Min.

## Proof of first direction:

$$\left[ \mathcal{H} \text{ NUL} \implies \mathcal{H} = \bigcup_m \mathcal{H}_m \text{ with } \mathcal{H}_m \text{ learnable} \right]$$

Assume $\mathcal{H}$ is NUL with some Algorithm $A(S)$.

Let $\mathcal{H}_m = \left\{ h \in \mathcal{H} \text{ s.t } m_{\mathcal{H}}^{\text{NUL}}(\varepsilon = \tfrac{1}{8}, \delta = \tfrac{1}{7}, h) \leq m \right\}$

Obviously $\bigcup_m \mathcal{H}_m = \mathcal{H}$.

By def of NUL:

$$\mathbb{P}_{S \sim \mathcal{D}^m} \left\{ L_{\mathcal{D}}(A(S)) \leq \min_{h \in \mathcal{H}_m} L_{\mathcal{D}}(h) + \tfrac{1}{8} \right\} \geq \tfrac{6}{7}$$

Take $\mathcal{D}$ realizable s.t $\min_{h \in \mathcal{H}_m} L_{\mathcal{D}}(h) = 0$. Then for

those $\mathcal{D}$: $\mathbb{P}_{S \sim \mathcal{D}^m} \left\{ L_{\mathcal{D}}(A(S)) \leq \tfrac{1}{8} \right\} \geq \tfrac{6}{7}$

$$\implies \mathbb{P}_{S \sim \mathcal{D}^m} \left\{ L_{\mathcal{D}}(A(S)) \geq \tfrac{1}{8} \right\} \leq \tfrac{1}{7}$$

$\implies VCdim(\mathcal{H}_m) < +\infty$ (finite). Indeed

otherwise the No Free lunch thm would hold and

$\mathbb{P}_{S \sim \mathcal{D}^m} \left\{ L_{\mathcal{D}}(A(S)) \geq \tfrac{1}{8} \right\} \geq \tfrac{1}{2}$. $\boxed{\implies \mathcal{H}_m \text{ PAC learnable}}$

Proof of converse direction.

$$\left[ \mathcal{H} = \underset{m}{\cup} \mathcal{H}_m \;,\; \mathcal{H}_m \text{ agnostic PAC learnable} \Rightarrow \mathcal{H} \text{ NUL} \right]$$

This requires a lemma (proved later on):

Lemma:

If $\mathcal{H} = \underset{m}{\cup} \mathcal{H}_m$ where $\mathcal{H}_m$ has unif convergence property then $\mathcal{H}$ is non-uniformly learnable.

Now assume $\mathcal{H} = \underset{m}{\cup} \mathcal{H}_m$, $\mathcal{H}_m$ agnostic PAC learnable. Then by Fund Thm of learning $\mathcal{H}_m$ must have the unif conv property. Thus by the lemma $\mathcal{H}$ is non-unif. learnable.

The rest of this lecture is devoted to the tools allowing to prove the lemma.

## Structural Risk Minimization.

Some intuitions first:

* $\mathcal{H} = \underset{m}{\cup} \mathcal{H}_m$ and we imagine that we have some prior confidence associated to each $\mathcal{H}_m$. So we associate some "weight", $0 \leq W_m \leq 1$, $\sum_m w_m = 1$, to each $\mathcal{H}_m$. The higher the weight the more we believe that the algorithm should belong to $\mathcal{H}_m$.

* in the lemma each class $\mathcal{H}_m$ has the unif conv property for some sample complexity $m_{\mathcal{H}_m}^{UC}(\epsilon, \delta)$.

* Given $m = |S| =$ sample size and confidence param $\delta$ we will consider the best "slack"

$$\epsilon_m(m, \delta) = \min\left(\epsilon \; ; \; m \geq m_{\mathcal{H}_m}^{UC}(\epsilon, \delta)\right)$$

$$\left(\text{Recall } m_{\mathcal{H}_m}^{UC}(\epsilon, \delta) \approx \frac{vcdim \, \mathcal{H}_m + \log 1/\delta}{\epsilon^2}\right).$$

## Theorem:

Let $\mathcal{H} = \bigcup_m \mathcal{H}_m$ with $\mathcal{H}_m$ satisfying UC property

with some $m^{UC}_{\mathcal{H}_m}(\epsilon, \delta)$.

Let $\mathcal{E}_m(m, \delta) = \min(\epsilon : m \geq m^{UC}_{\mathcal{H}_m}(\epsilon, \delta))$

Let $0 < w_m \leq 1$, $\sum_m w_m = 1$ <u>any</u> set of "weights"

associated to $\mathcal{H}_m$ (set $w_0 \geq w_1 \geq w_2 \geq \cdots$).

Then

$$\mathbb{P}_{S \sim \mathcal{D}^m}\left\{ \forall m, \sup_{h \in \mathcal{H}_m} \left| L_{\mathcal{D}}(h) - L_S(h) \right| \leq \mathcal{E}_m(m, w_m \delta) \right\}$$

$$\geq 1 - \delta.$$

## Proof of Theorem.

By UC for each $\mathcal{H}_m$ we have for $m \geq m_{\mathcal{H}_m}^{UC}(\varepsilon, \delta)$

$$\mathbb{P}\left\{ \sup_{h \in \mathcal{H}_m} |L_{\mathcal{D}}(h) - L_S(h)| \leq \varepsilon \right\} \geq 1-\delta .$$

Here $m$ and $\delta$ are fixed. So choose the best possible $\varepsilon$ i.e:

$$\varepsilon \to \varepsilon_m(m, \delta) = \min\left( \varepsilon : m_{\mathcal{H}_m}^{UC}(\varepsilon, \delta) \leq m \right).$$

$$\mathbb{P}\left\{ \sup_{h \in \mathcal{H}_m} |L_{\mathcal{D}}(h) - L_S(h)| \leq \varepsilon_m(m, \delta) \right\} \geq 1-\delta$$

Apply this to $\delta \to w_m \delta$ for each $m \in \mathbb{N}$.

$$\mathbb{P}_{S \sim \mathcal{D}^m}\left\{ \sup_{h \in \mathcal{H}_m} |L_{\mathcal{D}}(h) - L_S(h)| \leq \varepsilon_m(m, w_m \delta) \right\} \geq 1 - w_m \delta$$

$$\mathbb{P}_{S \sim \mathcal{D}^m}\left\{ \sup_{h \in \mathcal{H}_m} |L_{\mathcal{D}}(h) - L_S(h)| \geq \varepsilon_m(m, w_m \delta) \right\} \leq w_m \delta$$

this is valid for all $m \in \mathbb{N}$.

Now we sum these inequalities over $m \in \mathbb{N}$:

$$\sum_{m \in \mathbb{N}} \mathbb{P}_{S \sim \mathfrak{D}^m} \left\{ \sup_{h \in \mathcal{H}_m} \left| L_{\mathfrak{D}}(h) - L_S(h) \right| \geqslant \varepsilon_m(m, w_m \delta) \right\} \leq \delta$$

By the union bound: $\mathbb{P} \left\{ \exists m : \ldots \right\} \leqslant \sum_{m \in \mathbb{N}} \mathbb{P} \left\{ \ldots \right\}$

Taking then converse probabilities:

$$\mathbb{P} \left\{ \forall m : \sup_{h \in \mathcal{H}_m} \left| L_{\mathfrak{D}}(h) - L_S(h) \right| \leq \varepsilon_m(m, w_m \delta) \right\} \geqslant 1 - \delta$$

Let us discuss some implications of this thm.

This will allow us to introduce the <u>notion of</u>

<u>structural risk</u> .

Thm says that with high prob we have

$$\forall m : \quad \sup_{h \in \mathcal{H}_m} |L_{\mathcal{D}}(h) - L_S(h)| \leq \varepsilon_m (m, w_m \delta)$$

This means also that

$$\forall m, \quad \forall h \in \mathcal{H}_m : \quad |L_{\mathcal{D}}(h) - L_S(h)| \leq \varepsilon_m (m, w_m \delta)$$

Now we can fix $h \in \mathcal{H}$. We have

$$|L_{\mathcal{D}}(h) - L_S(h)| \leq \varepsilon_m (m, w_m \delta); \quad \forall m \text{ s.t } h \in \mathcal{H}_m .$$

$$\Rightarrow \quad |L_{\mathcal{D}}(h) - L_S(h)| \leq \min_{m : h \in \mathcal{H}_m} \varepsilon_m (m, w_m \delta)$$

$$\leq \varepsilon_{m(h)} (m, w_{m(h)} \delta)$$

we choose $m(h) = \min (m : h \in \mathcal{H}_m)$ $\xrightarrow{}$ (First $\mathcal{H}_m$ where h occurs)

Thus the theorem states that with high probability we have :

$$L_{\mathscr{D}}(h) \leq L_S(h) + \varepsilon_{m(h)}(m, w_{m(h)}, \delta)$$

This suggests that to make $L_{\mathscr{D}}(h)$ small a good idea is to use the

Structural risk minimization rule or also:

$$A^{SRM}(S) = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \left( L_S(h) + \varepsilon_{m(h)}(m, w_{m(h)}, \delta) \right)$$

This means we penalize the empirical risk with an additional term $\varepsilon_{m(h)}(m, w_{m(h)}, \delta)$. This term is constructed from the decomposite $\mathcal{H} = \bigcup_m \mathcal{H}_m$ with $m_{\mathcal{H}_m}^{UC}(\varepsilon, \delta)$ and $\varepsilon_m(m, \delta) = \min(\varepsilon ; m \geq m_{\mathcal{H}_m}^{UC}(\varepsilon, \delta))$ $m(h) = \min(m ; h \in \mathcal{H}_m)$ and a set of weights $0 \leq w_m \leq 1$.

Recall that this construction was motivated by the lemma:

**Lemma**  Let $\mathcal{H} = \underset{m}{\cup} \mathcal{H}_m$ with each $\mathcal{H}_m$ satisfying the UC property. Then $\mathcal{H}$ is non uniformly learnable with

$$m_{\mathcal{H}}^{NUL}(\varepsilon, \delta, h) \leq m_{\mathcal{H}_{m(h)}}^{UC}\left(\frac{\varepsilon}{2}, w_{m(h)}\delta\right).$$

**Proof:**

By definition of $A^{SRM}(S)$ we have

$$L_{\mathcal{D}}(A^{SRM}(S)) = \underset{h \in \mathcal{H}}{\min}\left(L_S(h) + \varepsilon_{m(h)}(m, w_{m(h)}\delta)\right)$$

$$\leq L_S(h) + \varepsilon_{m(h)}(m, w_{m(h)}\delta)$$

Take $m \geq m_{\mathcal{H}_{m(h)}}^{UC}\left(\frac{\varepsilon}{2}, w_{m(h)}\delta\right).$

From the above definitions it can then be seen that

$$\varepsilon_{m(h)}(m, w_{m(h)}, \delta) \leq \frac{\varepsilon}{2},$$

Then for $m \geq m^{UC}_{\mathcal{H}_{m}(h)}(\frac{\varepsilon}{2}, w_{m(h)}, \delta)$ we have

$$L_{\mathcal{D}}(A^{SRM}(S)) \leq L_S(h) + \frac{\varepsilon}{2}$$

Now $\mathcal{H}_{m(h)}$ (recall $h \in \mathcal{H}_{m(h)}$ by def of $m(h)$) satisfies UC property. So with prob at least $1 - \delta$ we have $L_S(h) \leq L_{\mathcal{D}}(h) + \frac{\varepsilon}{2}$.

Therefore we conclude that with prob $\geq 1 - \delta$

$$L_{\mathcal{D}}(A^{SRM}(S)) \leq L_{\mathcal{D}}(h) + \varepsilon$$

for $h \in \mathcal{H}_{m(h)}$ and $m \geq m^{UC}_{\mathcal{H}_{m(h)}}(\frac{\varepsilon}{2}, w_{m(h)}, \delta)$.

Now take any $m_{\mathcal{H}}^{NUL}(\varepsilon, \delta, h) = m_{\mathcal{H}_{m(h)}}^{UC}\left(\frac{\varepsilon}{2}, w_{m(h)} \delta\right)$

The above statement is equivalent to

$$\mathbb{P}\left( L_{\mathcal{D}}\left( A^{SRM}(S) \right) \leq \min_{\substack{h \text{ such that} \\ m_{\mathcal{H}}^{NUL}(\varepsilon, \delta, h) \leq m}} L_{\mathcal{D}}(h) + \varepsilon \right)$$

$$\geq 1 - \delta$$

This is the def of Non Unif learnability.