

Learning Infinite Classes I.

①

In this lecture and the next one we explore PAC learning when the hypothesis class has infinite cardinality. We will ultimately prove the "fundamental theorem of learning" next time and see that one can learn infinite hyp classes as long as the "VC dimension" is finite. The notion of VC dimension is introduced in next lecture but this one prepares the ground -

‡

We start with a recap of what has been done till now :

- Model of PAC learning and definition
- Uniform convergence property.
- Learning finite classes with ERM.
- No free lunch theorem.

Notations:

X domain set $x \in X$.

Y label set $y \in Y$, $Z = \underbrace{X \times Y}_{\text{sample set}}$.

unknown sample distr:

$$D(x, y) = D(x) D(y|x).$$

$$S = \{z_1, \dots, z_m\} = \{x_1, y_1, x_2, y_2, \dots, x_m, y_m\}$$

$$|S| = m$$

$H \ni h$ $h: X \rightarrow Y$ "hypothesis"
 $x \mapsto h(x) = y$

$$\text{loss fun } l(h, z) = \begin{cases} \mathbb{1}(y \neq h(x)) \\ (y - h(x))^2 \\ \dots \end{cases}$$

$$\text{True loss / risk } L_D(h) = \mathbb{E}_{z \sim D} (l(h, z))$$

$$\text{Empirical loss } L_S(h) = \frac{1}{m} \sum_{i=1}^m l(h, z_i).$$

Definition. \mathcal{H} is agnostic PAC learnable w.r.t \mathcal{L} and \mathcal{Z} , if given any $\epsilon > 0, \delta > 0$, $\exists m_{\mathcal{L}, \epsilon}(\epsilon, \delta) \in \mathbb{N}$ and an algorithm or rule $A(S) : \mathcal{X} \rightarrow \mathcal{Y}$ such that $\forall \mathcal{D}$ if $m \geq m_{\mathcal{L}, \epsilon}(\epsilon, \delta)$ we have

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

Definition: \mathcal{H} is simply PAC learnable if above holds for just realizable distributions, that is $\forall \mathcal{D}$ such that $\exists f$ with $L_{\mathcal{D}}(f) = 0$.

(in other words $\min L_{\mathcal{D}} = 0$ and $\mathbb{P} \{ L_{\mathcal{D}}(A(S)) \leq \epsilon \} \geq 1 - \delta$.)

Definition: We say \mathcal{H} satisfies the uniform convergence property if for $m > m_{\mathcal{L}, \epsilon}^{\text{UC}}(\epsilon, \delta)$:

$$\mathbb{P} \left\{ \forall h \in \mathcal{H} : |L_S(h) - L_{\mathcal{D}}(h)| \leq \epsilon \right\} \geq 1 - \delta$$

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

⑨

Theorem. Let $|\mathcal{H}| < +\infty$ (finite \mathcal{H} .)
and let loss l be bounded (say $0 \leq l \leq 1$).

Then \mathcal{H} has the uniform convergence property

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

Corollary: A finite \mathcal{H} is agnostic PAC

learnable with $A(S) = \text{ERM}(S)$ with

sample complexity $m_{\mathcal{H}}^{\text{UC}}(\epsilon, \delta)$, i.e. $|S| = m_{\mathcal{H}}^{\text{UC}}$.

Remark $\text{ERM}(S) = \arg \min L_S(h)$.

(5)

Theorem : No Free Lunch. (labeling problems)

Let A be a fixed algo or rule. Let S be sample set with $m \leq \frac{|\mathcal{X}|}{2}$ (so for infinite \mathcal{X} this is true for any integer $m \in \mathbb{N}$). Then

$\exists \mathcal{D}$ such that

a) we are in realizable case : $\exists f$ with $L_{\mathcal{D}}(f) = 0$

b) $\prod_{S \sim \mathcal{D}^m} \{ L_{\mathcal{D}}(A(S)) \geq \frac{1}{8} \} \geq \frac{1}{2}$.

Proof idea : we considered $C \subset \mathcal{X}$ with $|C| = 2m$ and all possible 2^{2m} labeling sets $C \rightarrow \{0, 1\}$

$f_1, f_2, \dots, f_{2^{2m}}$ and $\mathcal{D}_i(x, y) = \underbrace{\mathcal{D}(x)}_{\frac{1}{|C|}} \underbrace{\mathcal{D}(y|x)}_{\mathcal{D}_{y, f(x)}}$

and showed that a) and b) hold for at least one such set.

①

Corollary: Take X to be infinite and let

$\mathcal{H} : X \rightarrow \{0, 1\}$ be all possible labeling fcts.

This hypothesis class is not learnable.

Proof: We proceed by contradiction. Assume \mathcal{H}

is PAC learnable. Fix $\epsilon > 0, \delta > 0$. There must

exist A and $m_{\mathcal{H}}(\epsilon, \delta)$ s.t. $\forall \mathcal{D}$ with realizable property $L_{\mathcal{D}}(f) = 0$ (some f) we have

$$\mathbb{P} \left\{ L_{\mathcal{D}}(A(S)) \leq \epsilon \right\} \geq 1 - \delta$$

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

But by the No Free Lunch Theorem when $|X| > 2m$

(which is the case here for any $m \in \mathbb{N}$) there exists \mathcal{D}

$$\text{with } \mathbb{P} \left\{ L_{\mathcal{D}}(A(S)) \geq \frac{1}{8} \right\} \geq \frac{1}{7}.$$

Thus

$$\mathbb{P} \left\{ L_{\mathcal{D}}(A(S)) \geq \frac{1}{8} \right\} + \mathbb{P} \left\{ L_{\mathcal{D}}(A(S)) \leq \frac{1}{8} \right\} \geq \frac{1}{7} + 1 - \delta > 1 \text{ for } \delta < \frac{1}{7}$$

\Rightarrow contradiction. 

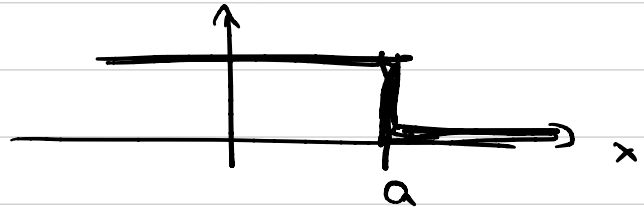
Learnability of Infinite classes.

Now we come to infinite classes. By the previous corollary we see that if an hypothesis class \mathcal{H} is too big it is not PAC learnable.

However we have seen in the exercise sessions that threshold functions:

$$\mathcal{H} = \{ h_a ; h_a(x) = \mathbb{1}(x \leq a), a \in \mathbb{R} \}$$

are PAC learnable.



$$\text{we have } m_{\mathcal{H}}(\epsilon, \delta) \leq \left\lceil \frac{\log 2/\delta}{\epsilon} \right\rceil .$$

We will ultimately show that if some "effective size of \mathcal{H} " is finite \mathcal{H} is agnostic PAC learnable with ERM .

Definition, Growth rate of \mathcal{H} .

The growth rate of an hypothesis class is by definition:

$$\tau_{\mathcal{H}}(m) = \max_{C \subset \mathcal{X}: |C|=m} |\mathcal{H}_C|$$

where \mathcal{H}_C is the restriction of \mathcal{H} to the set C ;
~~#~~

We restrict the functions $h: \mathcal{X} \rightarrow \mathcal{Y}$ to the arguments in C . So we consider $h: C \rightarrow \mathcal{Y}$. This forms set \mathcal{H}_C .

For $\mathcal{Y} = \{0, 1\}$ and $\mathcal{H} =$ labeling functions this amounts to consider only those labelings of points $C = \{c_1, c_2, \dots, c_m\}$ for $|C|=m$.

Example.

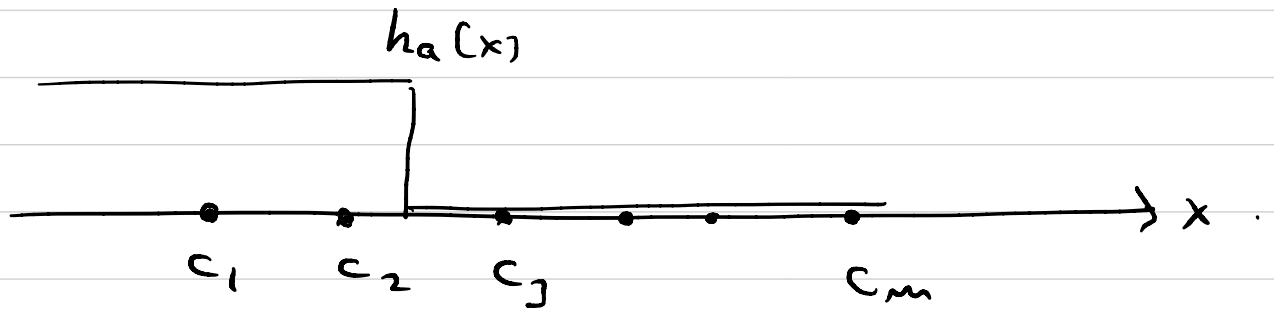
- $\mathcal{H} =$ all possible labeling fcts

$\Rightarrow \mathcal{H}_C =$ all possible labeling fcts of $C_1 \dots C_m$

and $|\mathcal{H}_C| = 2^m$ so $\tau_{\mathcal{H}}(m) = 2^m$

The growth rate is exponential.

- Threshold fcts $h_a(x) = \mathbb{1}(x \leq a)$.



possible labelings are

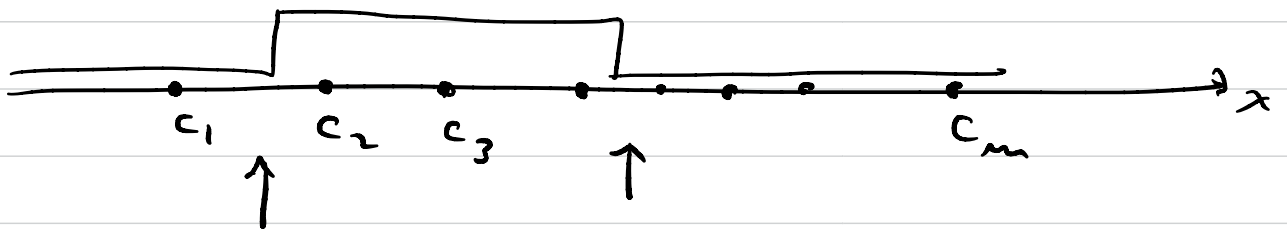
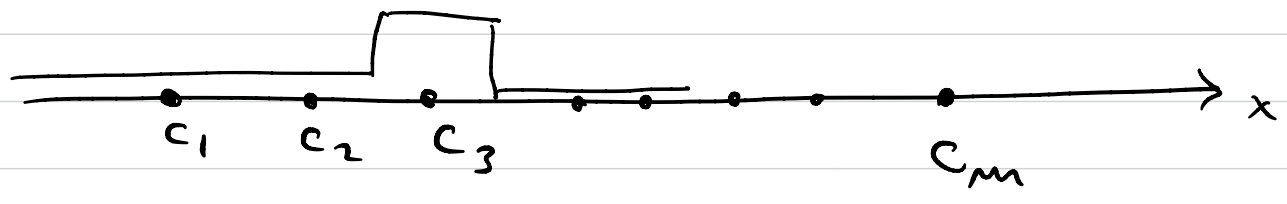
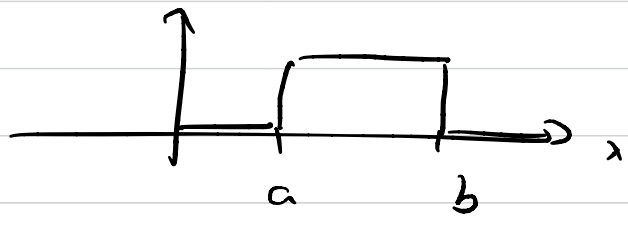
$$\begin{array}{l} 0 \ 0 \ \dots \ 0 \\ 1 \ 0 \ \dots \ 0 \\ 1 \ 1 \ \dots \ 0 \\ 1 \ 1 \ 1 \ \dots \ 0 \\ 1 \ 1 \ 1 \ \dots \ 1 \end{array}$$

$$\left. \begin{array}{l} \\ \\ \\ \\ \end{array} \right\} |\mathcal{H}_C| = m+1.$$

$\Rightarrow \tau_{\mathcal{H}}(m) = m+1$

linear growth rate.

• Rectangle $h_{a,b}(x) = \mathbb{1}(a \leq x \leq b)$



~ m choices for one threshold ~ m choices for other threshold.

$\Rightarrow \tau_{\text{TE}}(m) \sim m^2$

Squared growth rate.

Final lesson Growth rate can be exponential or polynomial etc...

We will see ultimately that pol growth rate means a class is learnable.

(11)

Main Theorem: fix $\delta > 0$. for any \mathcal{D} we

have :

$$\mathbb{P}_{S \sim \mathcal{D}^m} \left\{ \sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)| \leq \frac{4 + \log \tau_{\mathcal{D}}(2m)}{f \sqrt{2m}} \right\} \geq 1 - \delta$$

Corollary. Let $\tau_{\mathcal{D}}(m)$ has polynomial growth rate, specifically $\tau_{\mathcal{D}}(m) \leq \left(\frac{em}{d}\right)^d$

then \mathcal{H} has the uniform convergence property and is therefore PAC learnable by the ERM rule.

Remark : Next time we give a combinatorial characterization of \mathcal{H} with pol growth rate and will see that $d = VCdim(\mathcal{H})$.

Sketch of proof of corollary:

For large enough m and some constant C :

$$m \geq C \frac{d + \log(1/\delta)}{\epsilon^2}$$

$$\Rightarrow \frac{\sqrt{4 + d \log\left(\frac{em}{d}\right)}}{\delta \sqrt{2m}} \leq \epsilon$$

Thus for $m \geq C \frac{d + \log 1/\delta}{\epsilon^2}$ we have

$$\mathbb{P} \left\{ \sup_{h \in \mathcal{H}} |L_{\infty}(h) - L_S(h)| \leq \frac{4 + \sqrt{d \log \frac{em}{d}}}{\delta \sqrt{2m}} \right\}$$

$$\leq \mathbb{P} \left\{ \sup_{h \in \mathcal{H}} |L_{\infty}(h) - L_S(h)| \leq \epsilon \right\}$$

Thus $\mathbb{P} \left\{ \sup_{h \in \mathcal{H}} |L_{\infty}(h) - L_S(h)| \leq \epsilon \right\} \geq 1 - \delta$

which means emp. cov. prop. with $m_{\text{de}}(\epsilon, \delta) = C \frac{d + \log 1/\delta}{\epsilon^2}$.

Proof of Main Theorem.

We will show that

$$\mathbb{E}_{S \sim \mathcal{D}^m} \left[\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)| \right] \leq \frac{4 + \sqrt{\log \frac{1}{\delta} (2m)}}{\delta \sqrt{2m}}$$

Then by Markov's inequality:

$$\mathbb{P} \left[\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)| \geq \frac{4 + \sqrt{\log \frac{1}{\delta} (2m)}}{\delta \sqrt{2m}} \right]$$

$$\leq \frac{\mathbb{E} \left[\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)| \right]}{\frac{4 + \sqrt{\log \frac{1}{\delta} (2m)}}{\delta \sqrt{2m}}}$$

$$\leq \delta.$$

$$\mathbb{E}_{S \sim \mathcal{D}^m} \left\{ \sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)| \right\} = (*).$$

$$\mathbb{E}_{S' \sim \mathcal{D}^m} (L_{S'}(h)) \quad \text{by definition.}$$

$$\Rightarrow (*) = \mathbb{E}_{S \sim \mathcal{D}^m} \left[\sup_{h \in \mathcal{H}} \left| \mathbb{E}_{S' \sim \mathcal{D}^m} (L_{S'}(h) - L_S(h)) \right| \right]$$

$$\leq \mathbb{E}_{S \sim \mathcal{D}^m} \sup_{h \in \mathcal{H}} \mathbb{E}_{S' \sim \mathcal{D}^m} |L_{S'}(h) - L_S(h)|$$

$$\leq \mathbb{E}_{S \sim \mathcal{D}^m} \sup_{h \in \mathcal{H}} |L_{S'}(h) - L_S(h)|.$$

$$\left| \frac{1}{m} \sum_{i=1}^m (\ell(h, z'_i) - \ell(h, z_i)) \right|$$

where we pick $2m$ samples according to \mathcal{D}

$z_1 \quad z_2 \quad z_3 \quad \dots \quad z_m$

$z'_1 \quad z'_2 \quad z'_3 \quad \dots \quad z'_m$

This equivalent to

$$z_1, z_2', z_3, \dots, z_m$$

$$z_1', z_2, z_3', \dots, z_m'$$

In general it is equivalent to exchange z_i & z_i' for each i . This means that the r.h.s of inequality is equal to

$$\mathbb{E}_{\substack{S \sim \mathcal{D}^m \\ S' \sim \mathcal{D}^m}} \sup_{h \in \mathcal{H}} \left| \frac{1}{m} \sum_{i=1}^m \sigma_i (l(h, z_i) - l(h, z_i')) \right|$$

for any $(\sigma_1, \dots, \sigma_m) \in \{-1, +1\}^m$.

We can further average (*) over $\sigma_1, \dots, \sigma_m$ uniformly and obtain

$$(*) \leq \mathbb{E}_{\sigma_1, \dots, \sigma_m} \mathbb{E}_{\substack{S \sim \mathcal{D}^m \\ S' \sim \mathcal{D}^m}} \sup_{h \in \mathcal{H}} \left| \frac{1}{m} \sum_{i=1}^m \sigma_i (l(h, z_i) - l(h, z_i')) \right|$$

$$(*) \leq \mathbb{E}_{\substack{S \sim \mathcal{D}^m \\ S' \sim \mathcal{D}^m}} \mathbb{E}_{\sigma_1, \dots, \sigma_m} \sup_{h \in \mathcal{H}} \left| \frac{1}{m} \sum_{i=1}^m \sigma_i (\ell(h, z_i) - \ell(h, z'_i)) \right|$$

- for S and S' given "fixed" in this expression h see only

$$\text{the samples } \left\{ \begin{array}{l} x_1, \dots, x_m \\ x'_1, \dots, x'_m \end{array} \right\} = C$$

Here $C \subset \mathcal{X}$ and $|C| = 2m$.

So even if \mathcal{H} is infinite in fact

this expression depends only on \mathcal{H}_C .

\Rightarrow

$$\begin{aligned} & \mathbb{E}_{\sigma_1, \dots, \sigma_m} \sup_{h \in \mathcal{H}} \left| \frac{1}{m} \sum_{i=1}^m \sigma_i (\ell(h, z_i) - \ell(h, z'_i)) \right| \\ &= \mathbb{E}_{\sigma_1, \dots, \sigma_m} \sup_{h \in \mathcal{H}_C} \underbrace{\left| \frac{1}{m} \sum_{i=1}^m \sigma_i (\ell(h, z_i) - \ell(h, z'_i)) \right|}_{\mathcal{D}_h} \end{aligned}$$

$$\mathbb{E}_{\sigma_1, \dots, \sigma_m}(\vartheta_h) = 0 \quad \text{since } \sigma_i \sim \pm 1 \text{ uniformly.}$$

By Hoeffding's inequality:

$$P_{\sigma_1, \dots, \sigma_m}(|\vartheta_h - \mathbb{E}_{\sigma_1, \dots, \sigma_m}(\vartheta_h)| > \rho) \leq 2e^{-\frac{2m\rho^2}{4}}$$

$$\text{Remark } -1 \leq \vartheta_h \leq +1 \quad \uparrow \quad 2e^{-\frac{2m\rho^2}{(b-a)^2}}$$

By the union bound:

$$P_{\sigma_1, \dots, \sigma_m}(\sup_{h \in \mathcal{H}_c} |\vartheta_h| > \rho) = P\left\{ \bigcup_{h \in \mathcal{H}_c} (|\vartheta_h| > \rho) \right\}$$

$$\leq |\mathcal{H}_c| 2e^{-\frac{m\rho^2}{2}}$$

This is finite here!

An analysis calculation shown later implies

from here

$$\mathbb{E}_{\sigma_1, \dots, \sigma_m} \left[\sup_{h \in \mathcal{H}_c} |\vartheta_h| \right] \leq \frac{4 + \sqrt{\log |\mathcal{H}_c|}}{\sqrt{2m}}$$

Since this bound is even independent of S & S'
we find (see bottom of page 15).

$$G^* \leq \frac{4 + \log |\mathcal{H}_c|}{\sqrt{2m}} \leq \frac{4 + \sqrt{\log \tau_{\mathcal{H}_c}(2m)}}{\sqrt{2m}}$$

(Recall $\tau_{\mathcal{H}_c}(m) = \max_{|C|=m} |\mathcal{H}_c|$ and here $|C|=2m$)



Lemma (needed on page 17)

$$\mathbb{P}(X > p) \leq a e^{-bp^2} \Rightarrow \mathbb{E}(X) \leq \frac{2 + \sqrt{\log a}}{\sqrt{b}}$$

Proof, let $x \geq 0$

$$\begin{aligned} \mathbb{E}(X) &= \int_0^{+\infty} dp \, p \, \mathbb{P}(p) = - \int_0^{+\infty} dp \, p \, \frac{d}{dp} \underbrace{\int_0^{+\infty} p(x) dx}_{\mathbb{P}(x \geq p)} \\ &= \int_0^{+\infty} dp \, \mathbb{P}(x \geq p) \\ &= \int_0^{\sqrt{\frac{\log a}{b}}} dp \, \mathbb{P}(x \geq p) + \int_{\sqrt{\frac{\log a}{b}}}^{+\infty} dp \, \mathbb{P}(x \geq p) \\ &\leq \sqrt{\frac{\log a}{b}} + \int_{\sqrt{\frac{\log a}{b}}}^{+\infty} dp \, a e^{-bp^2} \\ &\quad \underbrace{\int_0^{+\infty} dp \, a e^{-b\left(p + \sqrt{\frac{\log a}{b}}\right)^2}}_{\leq e^{-bp^2} \text{ for } p > 0} \end{aligned}$$

check: $a e^{-b \left(p + \sqrt{\frac{bja}{b}}\right)^2} \leq e^{-bp^2}$?

$\Leftrightarrow \log a - b \left(p^2 + 2p \sqrt{\frac{bja}{b}} + \frac{bja}{b}\right) \leq -bp^2$

$\Leftrightarrow -2pb \sqrt{\frac{bja}{b}} \leq 0 \quad \checkmark$

So we have

$E(X) \leq \sqrt{\frac{bja}{b}} + \underbrace{\int_0^{+\infty} dp e^{-bp^2}}_{\sqrt{\frac{\pi}{b}} \leq \frac{2}{\sqrt{b}}}$

