

MACHINE LEARNING

giacomo.megla@inria.fr
othmane.marzouq@inria.fr

EVALUATION

- | (1) EXAM (40%)
- | (2) HOMEWORK (30%)
- | (3) QUESTION AT EVERY LECTURE (30%)

2020/2021

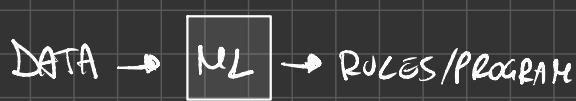


16/09

INTRODUCTION

MACHINE
LEARNING

- way of extract rules from data



ISSUES

- To large set of rules (Reu Ridgeon)
- To small set of rules (Reu Rat)

if too many reduces learning flexibility

INDUCTIVE
PRIOR

- A PRIORI KNOWLEDGE to prevent useless conclusion

INDUCTIVE
INFERENCE

- Ability to proceed from general examples to a broader generalization

KINDS OF LEARNING

SUPERVISED /
UNSUPERVISED

- In supervised learning output label are given ; Unsupervised instead finds unexpected correlation. (SPAM MAIL VS ANOMALY DETECTION)

- Labels are given but after a while

- Change the way the data is given from a teacher

- Data can be given all together or PASSO PASSO

- we'll see SUPERVISED BATCH LEARNING WITH PASSIVE TEACHER

①

- Statistic wants to check an HYPOTHESIS (smoke affects heart?)

- ML wants to find the HYPOTHESIS (what affects heart?)

- Statistic starts from the Prob. Distr. of data.

↳ ASYMPTOTICS: if you look at enough data you find a GAUSSIAN

- ML doesn't know the distribution

↳ FINITE SAMPLES: you get the distribution you have

STATISTIC
&
ML

②

STATISTICAL LEARNING FRAMEWORK

- X input space } (vector of features)
 - y output space }
 - S dataset of size m (#samples)
-

δ and D are unknown to the learner

- δ correct labelling function; D distribution over X
 - $A: S \subset (X \times Y)^m \rightarrow \{\text{functions: } X \rightarrow Y\}$
 - $\Rightarrow A(S) = h$
 - $h: X \rightarrow Y$ called HYPOTHESIS or PREDICTOR or CLASSIFIER
it can be seen as $A(S)$, so the output of the ML Algo. A given the DS. S
-

h is the output of an ML Algo.
(is a PREDICTION RULE)

EXPECTED LOSS

$$\bullet L_D(h) = \mathbb{E}_D(1_{h(x) \neq y}) \in [0, 1]$$

INDICATOR FUNCTION: 0 if prediction is ok it is 1 else is 0

$$\bullet L_S(h) = \frac{\sum_{i=1}^m 1_{h(x_i) \neq y_i}}{m}$$

\Rightarrow for Large number theory, EMPIRICAL LOSS will converge to EXPECTED LOSS.

Note

- we can write $y = \delta(x)$ so we can rewrite
- $L_{D,\delta}(h) = \mathbb{E}_D(1_{h(x) \neq \delta(x)})$; note also that our hope is to find h such that $L_{D,\delta}(h) < \epsilon$ with Prob $> 1 - \delta$ because in the best case we can have $L_{D,\delta}(h) = 0$ (that is obv. impossible)

APPROXIMATE WITH HIGH PROBABILITY

ERM - EMPIRICAL RISK MINIMIZATION

Def

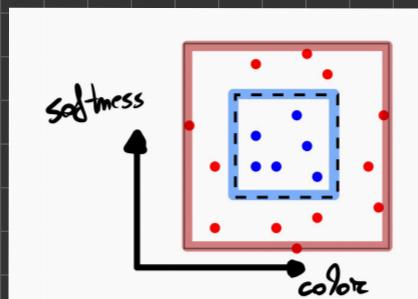
- ERM is the process of finding a PREDICTOR h which minimizes the EMPIRICAL LOSS $L_S(h)$

OVERFITTING

- When a predictor is excellent on the training set but very poor on the distribution



BRU is bad?



- Uniform D
- $\text{AREA}(\square) = 1$
- $\text{AREA}(\square) = 2$
- $g(x) \begin{cases} 1 & \text{if } x \in \square \\ 0 & \text{otherwise} \end{cases}$

Overfitting example

- Let $S = \square$ and define:

$$h_S(x) = \begin{cases} g, & \forall x \in S \\ 0 & \text{otherwise} \end{cases}$$

it means that $h_S(x)$ is perfect for the training set, so $L_S(h_S) = 0$; but, having $P = \frac{1}{2}$ of taking an element $x \notin S \Rightarrow L_D(h_S) = \frac{1}{2}$, so it's bad for the distribution D



- MEMORIZATION ALGO.

- Suppose: $h_M(x) = \begin{cases} 1 & \text{if } x = x_i \in S \\ \text{Toss A COIN} & \text{otherwise} \end{cases}$

- the P of taking an already known SAMPLE from D over the distribution D is ϕ , so you always toss a coin.

→ It perfectly explains the database, but is shit in general

Other
OVERFITTING
EXAMPLE

ERM WITH INDUCTIVE BIAS

⇒ ERM is not bad but learns in a too large set of HYPOTESIS, so it depends from the set of hypothesis you apply it.

Def

- The IDEA is to reduce the SEARCH SPACE of ERM;
- we define, before looking at the data, a set of predictor H called HYPOTHESIS CLASS, and we use ERM rule to find $h_{\text{ERM}_H} = \text{ERM}_H(S) \in \underset{h \in H}{\operatorname{argmin}} L_S(h)$.
- mean biasing a LEARNER to a finite set of predictor.

?

| HOW TO REDUCE THE SEARCH SPACE H ?

- - upperbound its size

~~REGULARITY ASSUMPTION~~

- There exist always $h^* \in H \mid L_{D,g}(h^*) = 0$, it means $h^* = g$

⇒ it implies that, if the set H is finite, then for every S taken from D $L_S(h^*) = 0$, so ERM_H will always find it.

ISSUE

- We are interested in $L_{D,g}(h_{\text{ERM}_H})$, so to be well done S should be enough representative of the underlying dist. D

↓
i.i.d.
ASSUMPTION

- All sample of S are INDEPENDENTLY and IDENTICALLY DISTRIBUTED according to D

↳ it implies that the bigger is S the more it is representative of D and g

NOTE

- $L_{D,g}(h_S)$ is a random variable because the choice of S is random, so can always happen that S would be not representative of D

CONFIDENCE
PARAMETER

- δ is the Prob. of getting NON-REPR. S
- $(1 - \delta)$

ACCURACY
PARAMETER

- ε is needed when to address the QUALITY of a Prediction

⇒ if $L_{D,g}(h_S) > \varepsilon$ FAIL; else if $L_{D,g} \leq \varepsilon$ we got an APPROX. CORRECT PREDICTION h_S

?

Can we upp. bound. the prob. of sampling S from D in a way that S won't be MISLEADING?

IDEA

We can find the right size $m = |S|$ such that ERM_H won't choose a $h_{\text{ERM}_H} = h_S$

Some Defs

- BAD HYPOTHESIS SET → $H_B := \{h \in H \mid L_{D,g}(h) > \varepsilon\}$
- POSSIBLY MISLEADING S →

$$Sp := \{S \mid \exists h_S \in H_B, L_S(h_S) = 0\}$$

BAD HYPOT. THAT LOOK GOOD ON S

- MISLEADING S →

$$S_u \in \{S \mid L_{D,g}(h_S) > \varepsilon\} = \{S \mid h_S \in H_B\} \subseteq Sp$$

IT HAPPENS ONLY IF $h_S \in H_B$

- we can rewrite $S_p = \bigcup_{h_s \in H_B} \{S \mid L_s(h_s) = 0\}$

Prob. of picking S_M

- we now calculate $\mathbb{P}_{D^m}(S_M)$

$$\begin{aligned}\mathbb{P}_{D^m}(S_M) &\leq \mathbb{P}_{D^m}(S_p) = \mathbb{P}_{D^m}\left(\bigcup_{h_s \in H_B} \{S \mid L_s(h_s) = 0\}\right) \leq \\ &\leq \sum_{h_s \in H_B} \mathbb{P}_{D^m}(\{S \mid L_s(h_s) = 0\}) \dots\end{aligned}$$

- Notice that $L_s(h_s) = 0 \iff \forall (x_i, y_i) \in S \quad h_s(x_i) = y_i \Rightarrow$

$$\Rightarrow \mathbb{P}(h_s(x_i) = y_i) = 1 - L_{D^m}(h_s) \stackrel{L_{D^m}(h_s) \geq \varepsilon \text{ (separating } S_M \text{!)}}{\leq} 1 - \varepsilon \Rightarrow$$

$$\Rightarrow \mathbb{P}(h_s(x_i) = y_i \quad \forall i=1, \dots, m) = \mathbb{P}_{D^m}(\{S \mid L_s(h_s) = 0\}) \leq (1 - \varepsilon)^m$$

$$\begin{aligned}\dots \sum_{h_s \in H_B} \mathbb{P}_{D^m}(\{S \mid L_s(h_s) = 0\}) &\leq \sum_{h_s \in H_B} (1 - \varepsilon)^m = \\ &= |H_B| (1 - \varepsilon)^m \stackrel{\text{over } H_B}{\leq} |H| e^{-\varepsilon m}\end{aligned}$$



- In conclusion $\mathbb{P}_{D^m}(S_M) \leq |H_B| e^{-\varepsilon m}$

BOUND

- we want $\mathbb{P}_{D^m}(S_M) \leq \delta$ so we should choose m carefully:

$$|H_B| e^{-\varepsilon m} \leq \delta \Rightarrow$$

$$m \geq \frac{\log(|H|/\delta)}{\varepsilon}$$

Conclusion

COROLLARY 2.3 Let \mathcal{H} be a finite hypothesis class. Let $\delta \in (0, 1)$ and $\varepsilon > 0$ and let m be an integer that satisfies

$$m \geq \frac{\log(|\mathcal{H}|/\delta)}{\epsilon}.$$

Then, for any labeling function, f , and for any distribution, \mathcal{D} , for which the realizability assumption holds (that is, for some $h \in \mathcal{H}$, $L_{(\mathcal{D},f)}(h) = 0$), with probability of at least $1 - \delta$ over the choice of an i.i.d. sample S of size m , we have that for every ERM hypothesis, h_S , it holds that

$$L_{(\mathcal{D},f)}(h_S) \leq \epsilon.$$

The preceding corollary tells us that for a sufficiently large m , the $\text{ERM}_{\mathcal{H}}$ rule over a finite hypothesis class will be *probably* (with confidence $1 - \delta$) *approximately* (up to an error of ϵ) correct.

21/09/21

PAC - LEARNABILITY

PAC LEARNABILITY

Def¹

- Given REGUZ. ASS on 0-1 loss funct $\mathbb{1}_{h \neq y}$) H is PAC-LEARNABLE if:

$$\exists A: (\mathcal{X} \times \mathcal{Y})^m \rightarrow H$$

$$\exists m_H: (0, 1)^2 \rightarrow N // \text{SAMPLE COMPLEXITY how many samples are required to guarantee an approx. correct solution?}$$

They are such that:

$\forall D$ over \mathcal{X} , $\forall \delta$, $\forall \delta \in (0, 1)$, $\forall \epsilon \in (0, 1)$ if we get S of m i.i.d. samples according to D from \mathcal{X} such that $|S| = m \geq m_H(\epsilon, \delta)$, then:

$$\mathbb{P}_{\delta, \epsilon}(A(S)) < \epsilon \text{ w.p. } (1 - \delta)$$

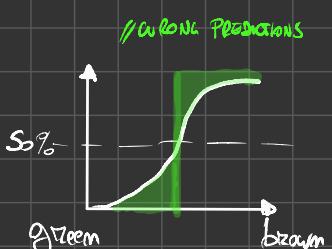
weak points of Def¹

- ① REGULARIZABILITY ASSUMPTION $\exists \delta$, $\delta \in H$
- ② 0-1 loss (only binary classification)

 \Rightarrow

- $\delta(x)$ is wasted and we'll use $P(y|x)$
- We need a THRESHOLD for $P(y|x)$, and put it to 50% is the choice that reduces $L_{DP}(h)$

$$h(x) = \begin{cases} 1 & P(1|x) = 0,5 \\ 0 & \end{cases}$$



MARGINAL TEST

CONDITIONAL LABELLING PROB

$$\cdot P(x)$$

$$\cdot P(y|x) = \frac{P(x,y)}{P(x)}$$

LOSS GENERALIZATION

some losses

- There exist different kinds of loss other than the 0-1 loss:

$$\cdot \vartheta(h, (x, y)) \in \mathbb{R}$$

$$\cdot \vartheta(h, (x, y)) = \mathbf{1}_{h(x) \neq y} \quad (\text{cause we want } \vartheta \geq 0)$$

$$\cdot \vartheta(h, (x, y)) = (h(x) - y)^2$$

- We can generalize them by $\mathcal{L}_D(h) = \mathbb{E}[\vartheta(h, (x, y))]$

REAL RISK EVALUATION

- Until now, with Realizability ass., we were sure that the best possible $\mathcal{L}_D(h)$ was 0 (because $\exists h^* \in H | h^* = g$); relaxing this assumption we need the to compare our loss to be the nearest possible to the BEST ONE we can achieve in H .

$$\mathcal{L}_D(h) \leq \min_{h' \in H} \mathcal{L}_D(h') + \varepsilon$$

AGNOSTIC PAC LEARNABILITY

Defn

- ~~(Given REAZ. ASS am 0-1 loss funct $\vartheta_{0,1}$)~~ H is PAC-LEARNABLE w.r.t. the loss ϑ :

$$- \exists A: (X \times \mathcal{Y})^m \rightarrow H$$

$$- \exists m_H: (0, 1)^2 \rightarrow \mathbb{N}$$

They are such that:

~~If D over $(X \times \mathcal{Y})$, $\forall \delta, \forall \epsilon \in (0, 1)$, $\forall S \subseteq \mathcal{X}$ if we get S of size $m \geq m_H(\epsilon, \delta)$ i.i.d. sampled according to D from $(X \times \mathcal{Y})$ such that~~

~~$|S| = m \geq m_H(\epsilon, \delta)$, then:~~

$$\underline{\mathcal{L}_{D,\vartheta}(A(S)) \leq \varepsilon \text{ w.p. } (1-\delta)}$$

$$\mathcal{L}_D(A(S)) \leq \min_h \mathcal{L}_D(h) + \varepsilon \text{ w.p. } \geq 1-\delta$$

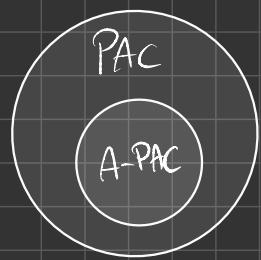
BEST POSSIBLE ERROR
OVER H .
W.R.T. REALIZABILITY WAS 0
BECAUSE OF $f \in H$

NB

• What is true?

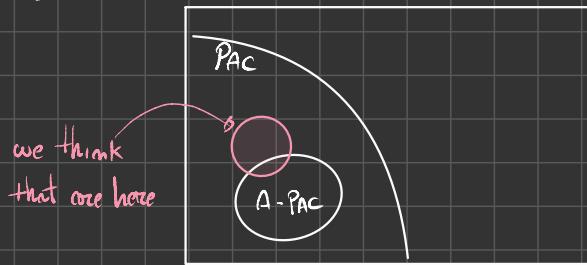
① H is A-PAC $\Rightarrow H$ is PAC

② H is PAC $\Rightarrow H$ is A-PAC



\Rightarrow ① : If H satisfies PAC LEARN. constraints, it's obvious that it also satisfies AGNOSTIC PAC LEARN. constraints.

Where are finite class of hypotheses? we know are PAC, but no one they say



ϵ -REPRESENTATIVE

ϵ -REPRESENTATIVE DATASET

- S is ϵ -REPRESENTATIVE if $|L_D(h) - L_S(h)| < \epsilon$ $\forall h$
- If you have S ϵ -repr, then ERM_h finds a good predictor

- Let $h_S = \text{ERM}_h(S)$

- Let S be ϵ/ϵ -REPRESENTATIVE

- we know by def that

$$L_D(h_S) \leq L_S(h_S) + \epsilon/\epsilon$$

$$\leq L_S(h) + \epsilon/\epsilon \quad \forall h$$

↑ by Def

$$\leq L_D(h) + \epsilon/\epsilon + \epsilon/\epsilon \quad \forall h$$

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

↓

By A-PAC LEARNABILITY

UNIFORM CONVERGENCE PROPERTY (UC)

- H has UC prop. w.r.t. the loss ℓ $\exists m_H^{\text{uc}}: (0,1)^2 \rightarrow N$ | HD
 $\forall \epsilon, \delta \in (0,1)$, if you draw a dataset S i.i.d from D with $|S| = m_H^{\text{uc}}(\epsilon, \delta)$ then S is ϵ -repr. w.p. $\geq 1 - \delta$

Note

- If H has UC prop. then S is ϵ -repr. with high prob., so $L_D(h_S) \leq \min_{h \in H} L_D(h) + \epsilon$,
 so H is AGNOSTIC PAC LEARNABLE, and $\text{ERM}_H(S)$ is AGNOSTIC PAC LEARNER for H

TH

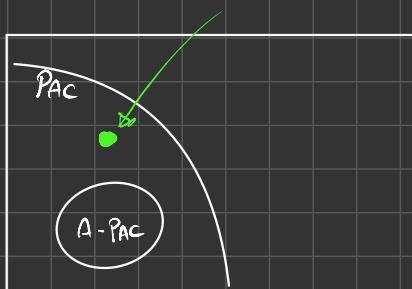
- If H is finite $\Rightarrow H$ has UC prop

Proof

Proceeding...

TH

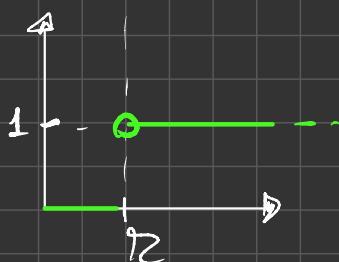
- There is at least an INFINITE PAC CLASS H



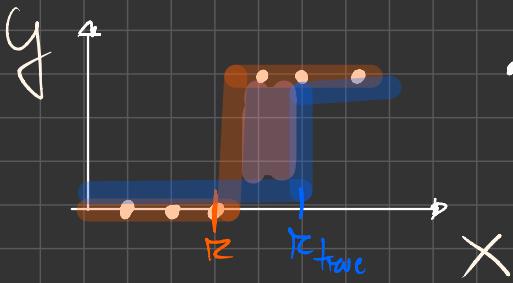
Proof

- Let H be a threshold function

$$H = \{h_r, r \in [0,1], h_r: [0,1] \rightarrow \{0,1\}, h_r(x) = \begin{cases} 1 & \forall x > r \\ 0 & \forall x \leq r \end{cases}\}$$



- Given the dataset

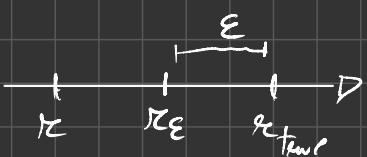


• ERM_H gives me h_r
where $r = \max\{x_i | g_i = 0\}$

- Suppose the δ , every time we observe a point in the pink region we have a loss of 1

$$L_D(h_r) = \mathbb{P}(X \in [r, r_{true}])$$

- Now we look a point in $[r, r_{true}]$ such that the error on it is ϵ



$$\cdot r_\epsilon : \mathbb{P}(X \in [r, r_{true}]) = \epsilon$$

$$\cdot \text{if } r < r_\epsilon \Rightarrow L_D(h_r) > \epsilon \quad \textcircled{1}$$

$$\cdot \text{if } r > r_\epsilon \Rightarrow L_D(h_r) < \epsilon \quad \textcircled{2}$$

$$\mathbb{P}(r < r_\epsilon) =$$

$$\Rightarrow \mathbb{P}(X \notin [r, r_{true}]) = 1 - \epsilon$$

$$= \mathbb{P}(X \notin [r, r_{true}] \text{ } \forall i=1, \dots, m) = \prod_{i=1}^m \mathbb{P}(X \notin [r, r_{true}]) =$$

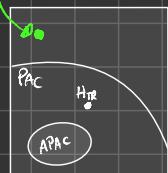
$$= \prod_{i=1}^m (1 - \epsilon) \Rightarrow (1 - \epsilon)^m \Rightarrow \text{we want } t \leq \delta$$

$$\Rightarrow \text{if we take } m \geq \frac{\log(\delta)}{\log(1 - \epsilon)} \text{ then } L_D(\text{ERM}_S) \leq \epsilon$$

w.p. $\geq 1 - \delta$ so H is PAC

28/9

- We will show that There is a class $H_{\{0,1\}^X} = \{\text{All binary partition over } X\} \mid H_{\{0,1\}^X} \notin \text{PAC}$



TH
[NO FREE LUNCH]

A

- $\exists D$ over X and a labeling function f , & A learning algorithm A for binary classification, those are s.t. by picking S of $m < \frac{|X|}{2}$ i.i.d. according to D you have

$$L_{D,f}(A(S)) \geq \frac{1}{8} \text{ w.p. } \geq \frac{1}{2} \quad \text{non-decidable}$$

PAC NEGATION

- H is not pac $\Leftrightarrow \forall A, \forall m_H : (0,1)^2 \rightarrow N, \exists D$ over X and $\exists f \in H$, $\exists \varepsilon_0, \delta_0 \in (0,1)$, $\nexists \exists m \geq m_H(\varepsilon_0, \delta_0) \mid |S|=m$ then every A learn with an error of

$$L_{D,f}(A(S)) \geq \varepsilon_0 \text{ w.p. } \geq 1 - \delta_0$$

→ FROM NO FREE LUNCH TH we can enforce the \neg PAC definition

- $\exists \delta$ instead of $\exists f \in H$
- $\varepsilon_0 = \frac{1}{3}$ and $\delta_0 = \frac{1}{4}$
- the requirement $\exists m \geq m_H(\varepsilon_0, \delta_0)$ that becomes $\nexists m$

Free lunch proof

B

- Pick $D \mid D(x_i) = \frac{1}{2}^m \text{ H.}$
- Let's proof that $\exists f \mid \mathbb{E}_{S \sim D^m} [L_{D,f}(A(S))] \geq \frac{1}{4}$
- Build a table of $|H|$ rows with many possible S on column

h_{z^m}			
:			
h_2	\oplus	\square	$L_{D,h_2}(A(S))$
h_1		\square	
S'	S'	$S^3 \dots$	

\oplus
rows of h_2 on S'

// pick g sources at at time t

// the average on each \square has avg error of $\frac{1}{4}$

// the set of all columns \square has ... $\frac{1}{4}$

// the all BS too in BFR $\frac{1}{4}$

// Avg of elements on raw \square $\frac{1}{4}$

⇒ I draw with Avg $> \frac{1}{4}$, it's the raw of f

↳ Proof if the whole set has Avg $\geq \frac{1}{4}$
there is at least one with Avg $> \frac{1}{4}$

- Now show that $A \Rightarrow B$ by showing $\neg A \Rightarrow \neg B$

$$(A) L_{D,+}(A(S)) \geq \frac{1}{8} \text{ w.p. } \geq \frac{1}{7}$$

$$(B) \exists \delta \mid \mathbb{E}_{S \sim D^m} [L_{D,+}(A(S))] \geq \frac{1}{6}$$

How can $L_{D,+}(A(S)) = 1$? happens w.p. $< \frac{1}{7}$ so

$$\mathbb{E}_{S \sim D^m} [L_{D,+}(A(S))] < \overbrace{1 \cdot \frac{1}{7} + \frac{1}{8} \cdot \frac{6}{7}}^{< \frac{1}{6}} = \frac{1}{6}$$

COROLLARY

[No Free Lunch]

Proof

- $|X| = +\infty$, $H_{\{0,1\}^X}$ is not PAC learnable w.r.t. 0-1 loss
- Assume H is PAC and choose $\varepsilon < \frac{1}{8}$ and $\delta < \frac{1}{7}$
- By PAC def: $L_D(A(S)) \leq \varepsilon$ w.p. $> 1 - \delta$
- By No-free lunch: $\exists D \mid L_D(A(S)) > \frac{1}{8} > \varepsilon$ w.p. $> \frac{1}{7} > \delta$ \perp

VC DIMENSION

- Finiteness of H is a sufficient but not necessary condition for it's learnability

why
VC-DIMENSION?

$$\begin{aligned} \text{PAC} &= A - \text{PAC} = VC = \{ \text{everything we can form via DPM} \} \\ &= VC - \dim(H) < +\infty \end{aligned}$$

SHATTERING

shattering example

- Given X , H over X and $A \subset X$ subset of X
- H shatters A if $\forall g: A \rightarrow \{0,1\} \exists h \in H | h(x) = g(x) \forall x \in A$



$$H = \{h_1, h_2, h_3, h_4\}$$

- $A = \{x_1\} \Rightarrow H$ shatters x_1

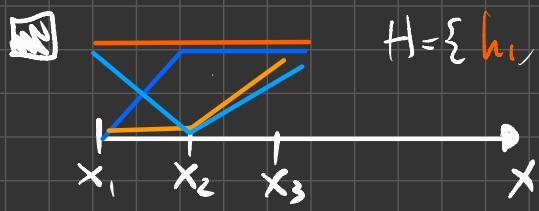
- $A = \{x_3\} \Rightarrow H$ DOESN'T (You can't get $x_3=1$ with $h=0$)

- $A = \{x_1, x_2\} \Rightarrow H$ Shatters A (we have 00, 01, 10, 11 in h)

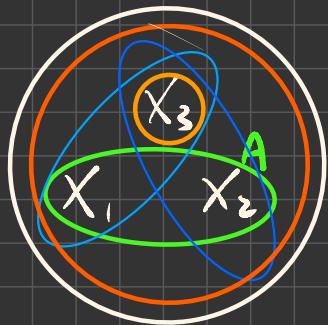


why Shattering

- Given a subset of point, we can always define H in a way that shatters them, so H can be seen as a SUBSET of X .
- $\Rightarrow H$ shatters A if $\forall g$ subset of $A \exists h \in H | h \cap g = 0$



$$H = \{h_1, h_2, h_3, h_4\}$$



- To get $\{\emptyset\}$ we get x_3 set cause $\{\emptyset\} \cap A = \{\emptyset\}$
- To get x_1 we get \circlearrowleft
- To get x_2 we get \circlearrowright

□

VC-dim

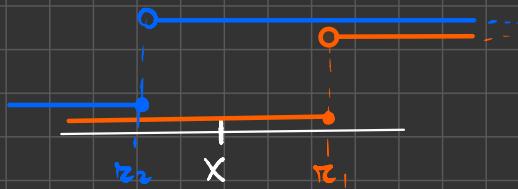
- $\text{VC-dim}(H) = \max \{ |A| \mid H \text{ shatters } A \}$



Find VC-dim of $H_{\text{Threshold}} = \{h_{x_i} \mid h_{x_i}(x) \begin{cases} 1 & \text{if } x > x_i \\ 0 & \text{otherwise} \end{cases}\}$

row 1 point

row 2 set



- we can pick him by taking the set $X \setminus \{x\}$ cause $A \cap X \setminus \{x\} = \{\emptyset\}$

$$A \cap X \setminus \{x\} = \{\emptyset\}$$

- Now we know $\text{VC-dim}(H) \geq 1$

- we can't do the same for 2 points because there is no way of picking only one of them

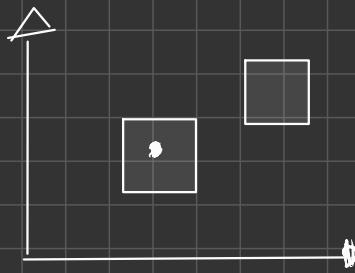
$$\Rightarrow \text{VC-dim}(M) = 1$$



Find VC-dim of $H_{rect} = \sum h_{a,b,c,d}$ s.t.

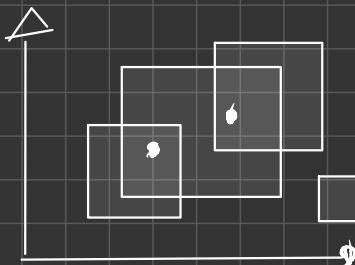
$$h_{a,b,c,d} \begin{cases} 1 & \text{if } 0 \leq a \leq b \text{ and } 0 \leq c \leq d \\ 0 & \text{otherwise} \end{cases}$$

1 point

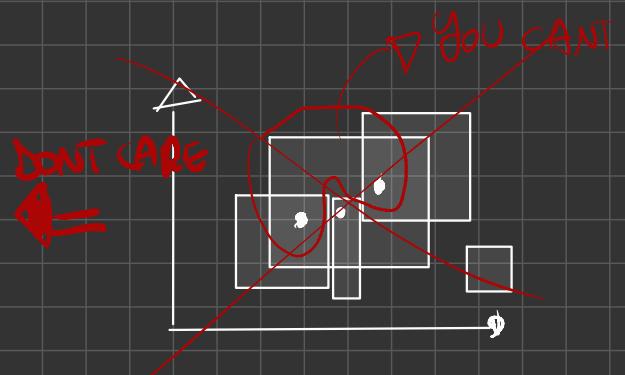
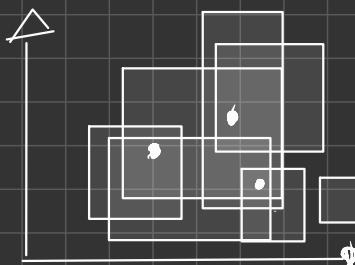


$$\Rightarrow \text{VC-dim}(H_{rect}) \geq 1$$

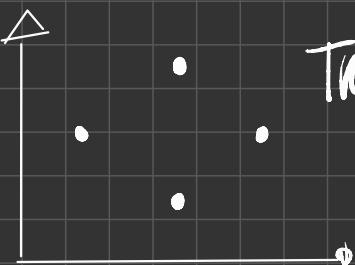
2 point



3 point



TRUST, WE CAN!



WE CANT SHATTER 6 POINTS !
WITHOUT THE 5TH !

\Rightarrow So you proved also for more than 5 points

$$\Rightarrow \text{VC-dim}(H) = 4$$

MIND

FUNDAMENTAL TH OF STATISTICAL LEARNING

TH FTSL

$\cdot X, 0-1 \text{ LOSS}, H$

(1) H IS PAC LEARNABLE

(2) ERM_H IS PAC LEARNER

(3) H IS AGNOSTIC PAC LEARNABLE

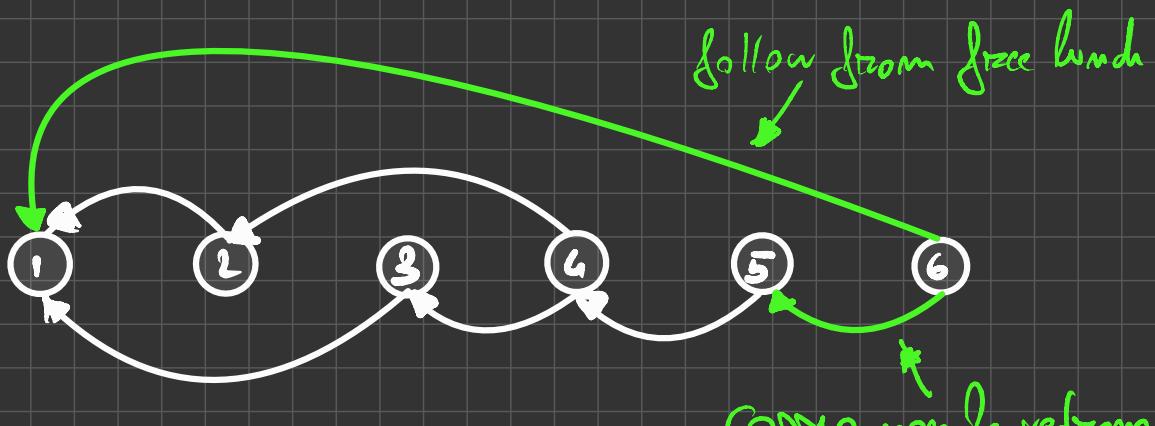
(4) ERM_H IS AGNOSTIC PAC LEARNER

(5) H HAS VC PROPERTY

(6) $\text{VC-dim}(H) < +\infty$

Proof

we don't know



(G-1) Proof

TH • H IS PAC $\Rightarrow \text{VC-dim}(H) < +\infty$

Proof.

NEGATE TH:

$\text{VC-dim}(H) > +\infty \Rightarrow H$ IS NOT PAC

.. ~ Vedi Alessio, io non
so la faccio per i sottive
Poncodio

5/10/21

TH QUANTITATIVE FTSL

// You can see how VCdim characterize the complexity of different class of learning

- Suppose $d = \text{VCdim}(H) < +\infty$, $\exists c_1, c_2 \in \mathbb{R}$ s.t.

(1) H is PAC-LEARNABLE with $m_H(\epsilon, \delta)$:

$$C_1 \frac{d + \log(1/\delta)}{\epsilon} \leq m_H(\epsilon, \delta) \leq C_2 \frac{d \log(1/\epsilon) + \log(1/\delta)}{\epsilon}$$

(2) H is APAC-LEARNABLE with $m_H(\epsilon, \delta)$:

$$C_1 \frac{d + \log(1/\delta)}{\epsilon^2} \leq m_H(\epsilon, \delta) \leq C_2 \frac{d + \log(1/\delta)}{\epsilon^2}$$

(3) H has UC-PROPERTY with $m_H^{uc}(\epsilon, \delta)$:

$$C_1 \cdot 1/\epsilon \leq m_H^{uc}(\epsilon, \delta) \leq C_2 \cdot 1/\epsilon$$

LINEAR PREDICTORS

- Linear predictors are a FAMILY OF HYPOTHESIS CLASS

CLASSES	ALGORITHMS
HALFSPLANES	LP
LINEAR REGRESSION PREDICTORS	PERCEPTRON
LOGISTIC REGRESSION PREDICTORS	LEAST SQUARE

AFFINE FUNCTIONS

- $L_d = \{h_{wb}: X \rightarrow \langle w, x \rangle + b \mid w \in \mathbb{R}^d \wedge b \in \mathbb{R}\}$

$$h_{wb}(x) = \langle w, x \rangle + b =$$

- We can incorporate the bias b in w

$$w' = (b, w_1, \dots, w_d) \in \mathbb{R}^{d+1}, x' = (1, x_1, \dots, x_d) \in \mathbb{R}^{d+1}$$

HALFSPACES

- Designed for BINARY CLASSIFICATION problems

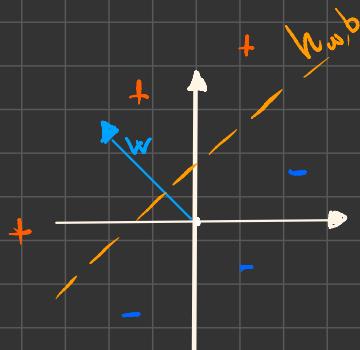
- $X = \mathbb{R}^d$, $Y = \{-1, +1\}$

$$HS_d = L_d \circ \text{SIGN} = \{x \mapsto \text{SIGN}(h_{w,b}(x)) \mid h_{w,b} \in L_d\}$$



$\exists x \cdot w / d = 2$

EACH HYPOTHESIS FORM AN
HYPERPLANE PERPENDICULAR TO w



SEPARABLE AND UNSEPARABLE CASES

- Separable mean that we can separate all Positive samples from Negative with a hyperplane.
 - The Non-Separable is NP-HARD
- \Rightarrow We have 2 method for IMPLEMENTING ERM FOR HALFSPACES IN SEPARABLE CASE (ie. LP and PERCEPTRON)

LINEAR PROGRAMMING FOR HS

- ERM predictor should have \emptyset error on the TS, so we looks for a vector w^* s.t.

$$\text{SIGN}(\langle w^*, x_i \rangle) = y_i \quad \forall i = 1, \dots, m \quad \Rightarrow$$

$$\Rightarrow y_i \langle w^*, x_i \rangle > 0 \quad \forall i = 1, \dots, m$$

- Let $\varphi = \min_{i=1-m} \{y_i \langle w^*, x_i \rangle\}$ and $\bar{w} = \frac{w^*}{\varphi}$ we have

$$y_i \langle \bar{w}, x_i \rangle = \frac{1}{\varphi} y_i \langle w^*, x_i \rangle \geq 1 \quad \forall i = 1, \dots, m$$

- So $\exists w \mid y_i \langle w, x_i \rangle \geq 1 \quad \forall i = 1 \dots m$

- This is our ERM predictor

• So the ERM L⁰ is:

$$\max_{w \in \mathbb{R}^d} \langle w, \mu \rangle \\ \text{s.t. } Aw \geq v$$

where

$$A_{m \times d} \mid A_{i,j} \cdot g_i \cdot x_{ij}$$

$$v = (1, -1) \in \mathbb{R}^m$$

$$\mu = (0, \dots, 0) \in \mathbb{R}^d // \text{because every } w \text{ that satisfies the constraints are equal candidates}$$

PERCEPTRON FOR HS

- ERM wants 0 error on the ts so $\text{SIGN}(\langle w^*, x_i \rangle) = y_i \quad \forall i=1 \dots m$
- At each step t if there is an $x_i \mid \text{SIGN}(\langle w^t, x_i \rangle) \neq y_i$ the algo. update w^t in such a way $w^{t+1} = w^t + g_i x_i$ to accomplish $y_i \langle w^{t+1}, x_i \rangle > 0$

Batch Perceptron

```

input: A training set  $(x_1, y_1), \dots, (x_m, y_m)$ 
initialize:  $w^{(1)} = (0, \dots, 0)$ 
for  $t = 1, 2, \dots$ 
    if ( $\exists i$  s.t.  $y_i \langle w^{(t)}, x_i \rangle \leq 0$ ) then
         $w^{(t+1)} = w^{(t)} + y_i x_i$ 
    else
        output  $w^{(t)}$ 

```

- At the end of the execution all samples of TS will be correctly classified

TH

- Given a SEPARABLE TS

$$\cdot \text{Let } B = \min \left\{ \|w\| \mid y_i \langle w, x_i \rangle \geq 1 \quad \forall i=1 \dots m \right\}$$

$$\cdot \text{Let } R = \max_i \{ \|x_i\| \}$$

$\Rightarrow (RB)^2$ iterations

- It stops with $y_i \langle w^t, x_i \rangle > 0 \quad \forall i=1 \dots m$

VCdim(HS)

- **HOMOGENEOUS HALFSPACE** : is an halfspace that does contain the \emptyset -vector
 \Rightarrow So the induced hyperplane pass through the ORIGIN.

TH

Proof

- $\text{VCdim}(\text{HS}) = d$ for the class of HOMOGENOUS HALFSPACES in \mathbb{R}^d
- Consider the set $(e_1, \dots, e_d) \mid e_{ij} = 0 \forall i \neq j \wedge e_{ii} = 1$
- this set is shattered by the Hom. HS class, by simply setting w as the labelling $w = (g_1, \dots, g_d) \Rightarrow \langle w, e_i \rangle = g_i \forall i$
- Let $x_1, \dots, x_d, x_{d+1} \in \mathbb{R}^d$ them must exist $a_1, \dots, a_{d+1} \in \mathbb{R}$ not all zeros s.t. $\sum_{i=1}^{d+1} a_i \cdot x_i = 0$ (LINEAR DEPENDENCE)
- Let $I := \{i \mid a_i > 0\}$, $J := \{j \mid a_j < 0\}$
- Assume both non empty:

$$\sum_{i \in I} a_i x_i - \sum_{j \in J} |a_j| x_j = 0$$

- Suppose x_1, \dots, x_{d+1} are shattered them

$$\exists w \mid \langle w, x_i \rangle > 0 \forall i \in I \wedge \langle w, x_j \rangle < 0 \forall j \in J$$

it implies

$$0 < \sum_{i \in I} a_i \langle x_i, w \rangle = \left\langle \sum_{i \in I} a_i x_i, w \right\rangle = \left\langle \sum_{j \in J} |a_j| x_j, w \right\rangle = \sum_{j \in J} |a_j| \langle x_j, w \rangle < 0,$$

- if I is empty

$$0 = \sum_{i \in I} a_i \langle x_i, w \rangle = \left\langle \sum_{i \in I} a_i x_i, w \right\rangle = \left\langle \sum_{j \in J} |a_j| x_j, w \right\rangle = \sum_{j \in J} |a_j| \langle x_j, w \rangle < 0,$$

- if J is empty

$$0 < \sum_{i \in I} a_i \langle x_i, w \rangle = \left\langle \sum_{i \in I} a_i x_i, w \right\rangle = \left\langle \sum_{j \in J} |a_j| x_j, w \right\rangle = \sum_{j \in J} |a_j| \langle x_j, w \rangle = 0,$$

TH • $\text{VC dim}(\text{HS}) = d+1$ for the class of non-homogeneous PAC spaces in \mathbb{R}^{d+1}

Proof

- As before consider the set $(0, e_1, \dots, e_d)$ is shattered. Then suppose $x_1, \dots, x_{d+1} \in \mathbb{R}^{d+1}$, we can reach a \perp as before.

ERROR DECOMPOSITION

$$L_D(A(S)) = \boxed{L_D(A(S)) - \min_{h \in H} L_D(h)} + \boxed{\min_{h \in H} L_D(h)}$$

- ESTIMATION ERROR • APPROXIMATION ERROR
 $(\leq \epsilon \text{ for PAC learnability})$

NOTE

- If $H \subset H'$ we have

$$\epsilon_{\text{Approx}}(H') \leq \epsilon_{\text{Approx}}(H)$$

$$\epsilon_{\text{Estimation}}(H) \leq \epsilon_{\text{Estimation}}(H')$$

- The bigger is H the easier the OVERFITTING PROBABILITY

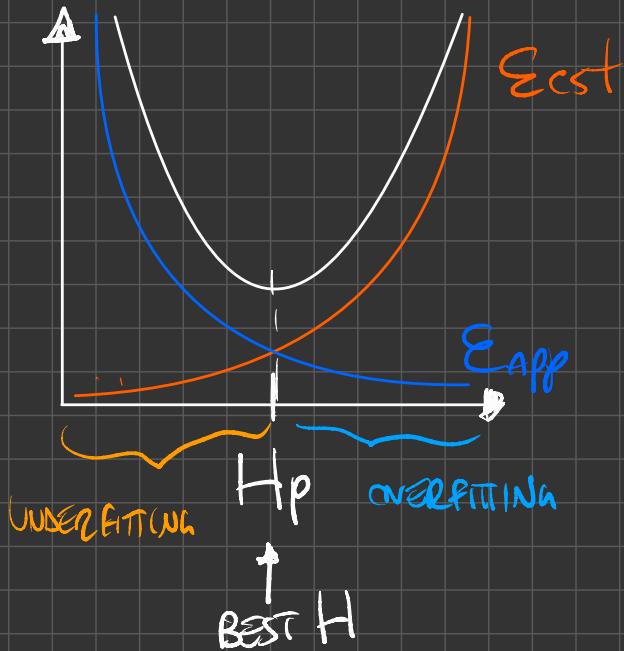
- from THE QUANT. FTSL :

$$m_H(\epsilon, d) \approx \frac{d_H + \vartheta_m(1/\delta)}{\epsilon^2}$$

- So given m :

$$\epsilon_{\text{estim}} \approx \sqrt{\frac{d_H + \vartheta_m(1/\delta)}{m}}$$

- Graphically:



- What is the best possible error?

$$L_D(A(S)) =$$

$$L_D(A(S)) = \underbrace{\min_{h \in H} L_D(h)}_{\text{ESTIMATION ERROR}} + \underbrace{\min_{h \in H} L_D(h) - \min_{h \in H} L_D(h)}_{\text{APPROXIMATION ERROR}} + \underbrace{\min_{h \in H} L_D(h)}_{\text{BIAS ERROR (IRREDUCIBLE)}}$$

WEAK LEARNABILITY, BOOSTING

- We look for a simpler LEARNABILITY definition leading to more efficient solution.
- Something NP-HARD for PAC LEARN. becomes POLYNOMIAL for WEAK LEARN. (ob course they fucked up)
- WL \equiv PAC L.
- SHAPIRE showed that a POLY. ALGO. to weakly LEARN H can be reduced in Polynomial time in a PAC LEARNER for H (che sign)

HOPES
CONCESSION

- The Poly. Reduction Algorithm is called **BOOSTING**

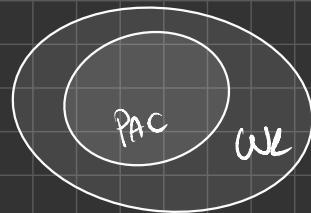
WEAK LEARNABILITY

- Assumed 0-1 loss and learnability
- H is δ -WL if $\exists A: (\mathcal{X} \times \mathcal{Y})^m \rightarrow H$ and $\exists m_H^{\text{WL}}: (0,1) \rightarrow N$ s.t. $\forall D \text{ over } \mathcal{X}, \forall \delta \in H, \forall \gamma \text{ s.t. } S \sim D^m \text{ i.i.d. } |S|=m \geq m_H^{\text{WL}}(\delta)$ then:

$$L_{D,\gamma}(A(S)) \leq \frac{1}{2} - \gamma \quad \text{w.p. } 1 - \delta$$

$\delta \in (0, \frac{1}{2})$ // we want the learner a bit better than a coin tosser

- Because WL require a stupid learner:



- But the truth is $\text{PAC-L} = \text{WL}$, we will show why

TH

Proof

- $\text{WL} = \text{PAC}$

① $\text{PAC-L} \Rightarrow \text{WL}$: $\exists \gamma \in (0, \frac{1}{2}) \mid H \text{ is } \gamma\text{-WL?}$

- $\text{if } H \in \text{PAC} \Rightarrow L_{D,\gamma}(A(S)) < \varepsilon \quad \text{w.p. } \geq (1-\delta)$

- Suppose $H \in \text{PAC}$ is also $\text{H} \in \text{WL}$:

- We can choose $\gamma = \frac{1}{n} \Rightarrow L_{D,\gamma}(A(S)) \leq \frac{1}{n} = \frac{1}{2} - \gamma$

② $\text{WL} \Rightarrow \text{PAC-L}$: proof $\neg \text{PAC} \Rightarrow \neg \text{WL}$

- Let H NOT PAC $\Rightarrow \text{VCd}(H) = +\infty$

For QNT. FCLS $\forall \varepsilon, d \Rightarrow m_H(\varepsilon) \leq \frac{\text{VCd}(H) + \log(\frac{1}{\delta})}{\varepsilon} \leq m_H(\varepsilon, d)$

- so $L_{D,\gamma}(A(S)) < \varepsilon$ in particular $\varepsilon = \frac{1}{2} - \gamma$

- Let $H_3 = \{h_{a,b} : \mathbb{R} \rightarrow \{-1, +1\} \mid a, b \in \mathbb{R} \cup \{-\infty, +\infty\}\}$, $h_{a,b}(x) = \begin{cases} +1 & \text{if } x \in [a, b] \\ -1 & \text{otherwise} \end{cases}$

• H_3 is W.L. by a learner that works over a smaller class $H_{DS} \subset H_3$

- H_{DS} = class of decision stamps in

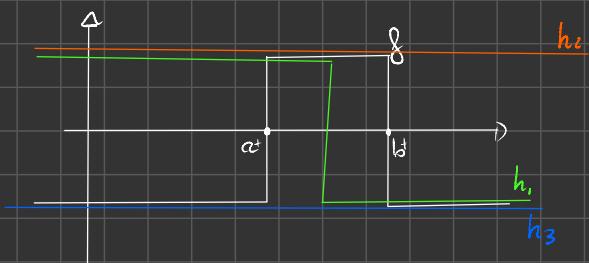
$$1-d = \{h_{a,b} \mid \mathbb{R} \rightarrow \{-1, +1\}, a \in \mathbb{R}, b \in \{-1, 1\}\}, h_{a,b}(x) = \begin{cases} b & \text{for } x \leq a \\ -b & \text{for } x > a \end{cases}$$

$\# D, g$

- $\exists h \in H_{DS} \mid L_{D,g}(h) \leq \frac{1}{3}$

• If $g \in H_3$ and $g \in H_{DS}$ then $L_{D,g}(h) = 0 < \frac{1}{3}$

- $\{g \in H_3 \setminus H_{DS}$



function does an error:

$$\begin{aligned} \cdot L_{D,g}(h_3) &= P(x \in [a, b]) \\ \cdot L_{D,g}(h_2) &= P(x > b_g) \\ \cdot L_{D,g}(h_1) &= P(x < a_g) \end{aligned}$$

*

$$\exists P(h) \leq \frac{1}{3}$$

- H_{DS}' is APAC with VCdim = 2

• We can take $m_H(\epsilon, \delta) = \frac{2 + \ln(1/\delta)}{\epsilon^2}$ to have

$$L_{D,g}(A(s)) \leq \min_{h \in H_{DS}'} L_{D,g}(h) + \epsilon \stackrel{*}{\leq} \frac{1}{3} + \epsilon$$

($\epsilon \in \mathcal{D}_n[\delta]$)

• We can say that this is a WL for H_3 just by choosing ϵ, δ such that $\frac{1}{3} + \epsilon = \frac{1}{2} - \gamma$ so

we choose $\epsilon = \frac{1}{12} = \gamma$

$$\left(\frac{1}{3} + \frac{1}{12} = \frac{1}{2} - \frac{1}{12} \right)$$

- ERM over H_{DS}' is a $\frac{1}{12}$ -Weak Learner for H_3

Why? \rightarrow Solving ERM $_{H_{DS}'}$ is EASIER than ERM $_{H_3}$; we can then transform it in a PAC-LEARNER for H_3 by BOOSTING ALGO.

CONCLUSION

BOOSTING ALGO

- ERM_{H_{DS}^d} CAN be implemented in polynomial time
 - $H_{DS}^d = \{h_{a,b}: \mathbb{R}^d \rightarrow \{-1, +1\} \mid a \in \mathbb{R}, b \in \{-1, +1\}\}$
 - $h_{a,b,i} = \begin{cases} b & \text{if } x_i < a \\ -b & \text{otherwise} \end{cases}$
 - Given $S = \{(x_i, y_i) \mid i=1, \dots, m\}$
 - $\text{ERM}_{H_{DS}^d} = \underset{h \in H_{DS}^d}{\text{argmin}} L_S(h) = \underset{a, b, i}{\text{argmin}} L_S(h_{a,b,i})$ for a we take just $m+1$ value
 - The complexity is:
- COMPLEXITY

$$(m-1) \cdot 2 \cdot d \cdot \underbrace{m}_{\text{useful } |b|} + \underbrace{d \cdot m \log m}_{\text{to compute } L_S(h_{a,b,i})}$$

 \Rightarrow we can update it instead of recalculating every step

 \Rightarrow ORDER acc d sets of point
- $$\Rightarrow md + md \log m$$

- QUA BOI, NON SO PERCHÉ MA:

$$\underset{\text{argmin}}{\sum_{i=1}^m D_i \delta(h_i(x_i, y_i))} \quad / D_i \text{ IS A DISTRIBUTION OF WEIGHTS}$$

ADA BOOST

- To learn over H it uses a δ -WL over H_B WL

$$\underset{h \in H_B}{\operatorname{argmin}} \sum_{i=1}^m D_i \vartheta(h_i(x_i, y_i))$$

$$D_i^{(1)} = \frac{1}{m} \quad \text{for } i = 1 - m$$

For ($t = 1 - T$)

$$h_t = \text{WL}(D^{(t)}, \gamma)$$

$$\varepsilon_t = \sum_{i=1}^m D_i^{(t)} \vartheta(h_i, (x_i, y_i))$$

$$\omega_t = 2 \ln \left(\frac{1}{\varepsilon_t} - 1 \right)$$

$$D^{(t+1)} = \frac{\exp(-\omega_t y_i h_t(x_i)) D_i^t}{\sum_{j=1}^m \exp(-\omega_t y_j h_t(x_j)) D_j^t}$$

// POSITIVE PAPAYA HAVE
POSITIVE NUMERATOR,
ELSE NEGATIVE
(ENFORCE THE WEIGHTS!)

RETURN $h = \text{SIGN} \left(\sum_{t=1}^T \omega_t h_t \right)$

TH

- If WL outputs h with $\sum D_i \vartheta(h_i, (x_i, y_i)) \leq \frac{1}{2} - \delta$ w.p. δ
then AdaBoost outputs h_{Ada} with $L_S(h_{\text{Ada}}) \leq e^{-2\delta^2 T}$ w.p. $(1 - \delta)^T$