

# 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

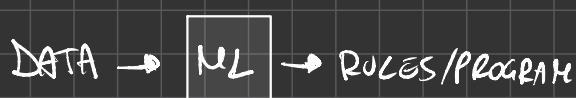


16/09

## 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)
- 

$\delta$  and  $D$  are unknown to the learner

- $\delta$  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

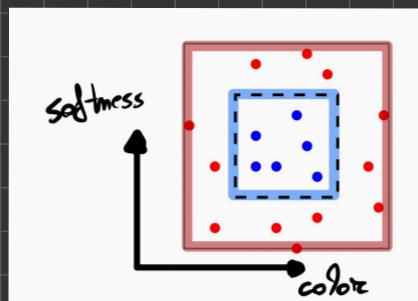
- ERM is the process of finding a PREDICTOR  $h$  which minimizes the EMPIRICAL LOSS  $L_S(h)$

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  $\phi$ , 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  $H$ ?

o

- upperbound its size

~~REGULARIZABILITY 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  $\text{ERM}_H$  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

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

ACCURACY  
PARAMETER

- $\varepsilon$  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  $h_{ERM_H} = h_S$

Some Defs

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

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

BAD HYPOT. THAT LOOK GOOD ON  $S$

- MISLEADING  $S \rightarrow$

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

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  $\epsilon$ ) correct.

21/09/21

# PAC - LEARNABILITY

## PAC LEARNABILITY

Def<sup>1</sup>

- 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:

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

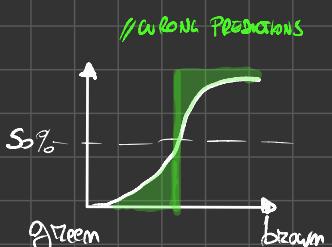
weak points  
of Def<sup>1</sup>

- 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:

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

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

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

- 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  $\vartheta_{(x,y)}$ )~~  $H$  is PAC-LEARNABLE w.r.t. the loss  $\vartheta$ :

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

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

They are such that:

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

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

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

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

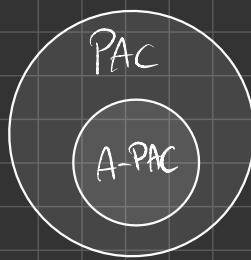
BEST POSSIBLE ERROR  
OVER  $H$ .  
w.r.t. REALIZABILITY WAS 0  
BECAUSE OF  $\exists h^* \in H$

NB

• What is true?

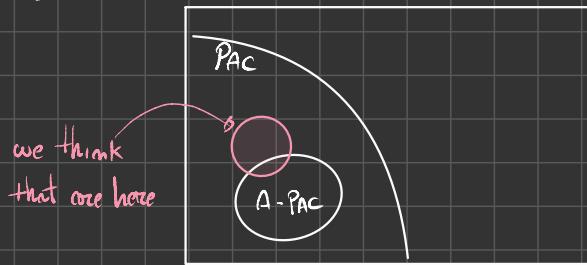
①  $H$  is A-PAC  $\Rightarrow H$  is PAC

②  $H$  is PAC  $\Rightarrow H$  is A-PAC



$\Rightarrow$  ① : If  $H$  satisfies PAC LEARN. constraints, it's obvious that it also satisfies A-agnostic PAC LEARN. constraints.

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



## $\epsilon$ -REPRESENTATIVE

### $\epsilon$ -REPRESENTATIVE DATASET

- $S$  is  $\epsilon$ -REPRESENTATIVE if  $|L_D(h) - L_S(h)| < \epsilon$   $\forall h$
- If you have  $S$   $\epsilon$ -repr, then  $\text{ERM}_h$  finds a good predictor

• Let  $h_S = \text{ERM}_h(S)$

• Let  $S$  be  $\epsilon/\epsilon$ -REPRESENTATIVE

• we know by def that

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

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

$\uparrow$  by def

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

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

$\downarrow$

By A-PAC LEARNABILITY

## UNIFORM CONVERGENCE PROPERTY (UC)

- $H$  has UC prop. w.r.t. the loss  $\ell$   $\exists m_H^{\text{uc}}: (0,1)^2 \rightarrow N$  | 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  $\epsilon$ -repr. w.p.  $\geq 1 - \delta$

Note

- If  $H$  has UC prop. then  $S$  is  $\epsilon$ -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

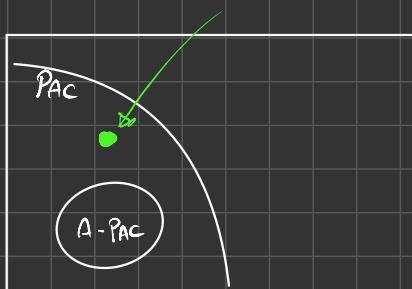
- If  $H$  is finite  $\Rightarrow H$  has UC prop

Proof

Proceeding...

TH

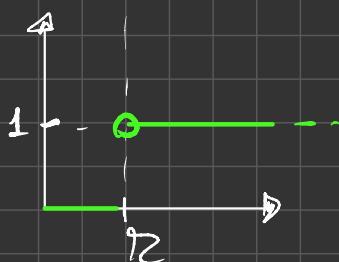
- There is at least an INFINITE PAC CLASS  $H$



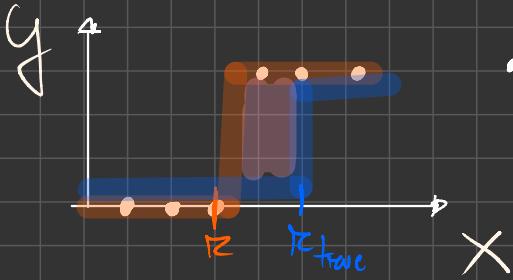
Proof

- 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

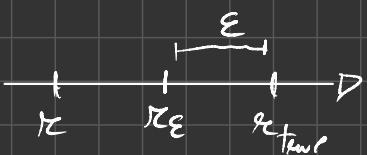


• ERM<sub>H</sub> gives me  $h_r$   
where  $r = \max\{x_i | g_i = 0\}$

- Suppose the  $\delta$ , 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  $\epsilon$



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

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

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

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

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

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

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

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

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

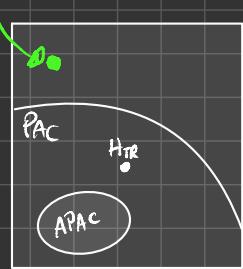
28/9 TH

Proof

• There is a class  $H \notin \text{PAC}$

•  $H_{\{\{0,1\}^X\}} = \{\text{All binary partition over } X\}$  // think of it as ALL BINARY NUMBER

• There is not a universal Learner (FREE LUNCH TH)



TH  
[NO FREE LUNCH]

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

$$L_{D,+}(A(S)) \geq \frac{1}{8} \text{ w.p. } \geq \frac{1}{2} \quad \text{NON REVERSIBLE}$$

A

COROLLARY  
[NO FREE LUNCH]

• If  $|X| = +\infty$ ,  $H_{\{\{0,1\}^X\}}$  is not PAC learnable w.r.t. 0-1 loss

PAC NEGATION

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

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

→ From NO FREE LUNCH TH we can enforce the  $\neg$ -PAC definition.

$\cdot \exists \delta$  instead of  $\exists \delta_{\text{f}}$

$\cdot \varepsilon_0 = \frac{1}{3}$  and  $\delta_0 = \frac{1}{4}$

$\cdot$  the requirement  $\exists m \geq m_H(\varepsilon_0, \delta_0)$  that becomes  $m$

free lunch proof

B

• Pick  $D | D(x_i) = \frac{1}{2^m} + i$

• Let's proof that  $\exists \delta | \mathbb{E}_{S \sim D^m} [L_{D,+}(A(S))] \geq \frac{1}{4}$

• Build a table of  $|H|$  rows with many possible  $S$  on columns

$h_{2^m}$	$\vdots$	$\vdots$	$L_{D_{h_2},+}(A(S))$
$h_2$	$\circlearrowleft$		
$h_1$			
$S'$	$S'$	$S^3 \dots$	

$\overset{\text{es}}{\text{rows of }} h_2 \text{ on } S'$

// pick 4 samples at a time

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

// the set of  $j$ th column  $\square$  has  $\frac{1}{4}$

// the all  $S$  too have BPP  $\frac{1}{4}$

// AVH of elements on row  $\square$   $\frac{1}{4}$

$\Rightarrow \exists \text{ raw with } \text{Avg} > \frac{1}{4}$ ,  $H$ 's the raw of  $\mathcal{F}$

$\hookrightarrow$  Prob if the whole Set has  $\text{Avg} \geq \frac{1}{4}$   
there is at least one with  $\text{Avg} > \frac{1}{4}$

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

$$(A) L_{D,+}(A(S)) \geq \frac{1}{8} \text{ w.p. } \geq \frac{1}{2}$$

$$(B) \exists \mathcal{F} \mid \mathbb{E}_{S \sim D^m} [L_{D,+}(A(S))] \geq \frac{1}{4}$$

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

$$\mathbb{E}_{S \sim D^m} [L_{D,+}(A(S))] < \overbrace{1 \cdot \frac{1}{2} + \frac{1}{8} \cdot \frac{6}{7}}^{< \frac{1}{4}} = \frac{1}{4} \quad \blacksquare$$

## VC DIMENSION

- why VC-DIMENSION?
- $\text{PAC} = A - \text{PAC} \equiv \text{VC} = \{ \text{everything we can form via ERM} \} = \text{VC-dim}(H) < +\infty$

SHATTERING

- Given  $X$ ,  $H$  over  $X$  and  $A \subset X$ :

- $H$  shatters  $A$  if  $\forall g: A \rightarrow \{0,1\} \exists h \in H \mid h(x) = g(x) \forall x \in A$



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



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

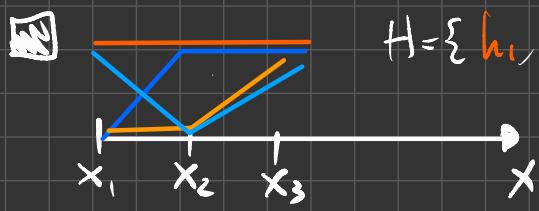
- $A = \{x_3\} \Rightarrow H$  DOESN'T (you can't get  $x_3=1$  with  $h \in H$ )

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

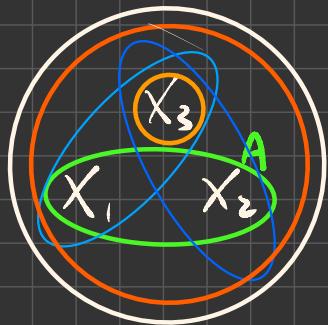


why Shattering

- 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  $\forall$  subset of  $A \exists h \in H \mid h \cap A = 0$



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



- To get  $\{\emptyset\}$  we get  $x_3$  set cause  $\{x_3\} \cap A = \{\emptyset\}$
- 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  $H_{\text{Threshold}} = \{h_{\tau} \mid h_{\tau}(x) \begin{cases} 1 & \text{if } x > \tau \\ 0 & \text{otherwise} \end{cases}\}$

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

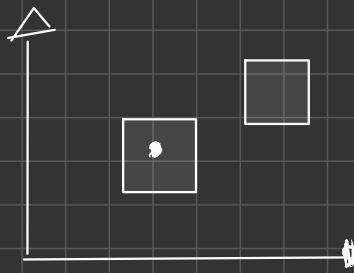
$$\Rightarrow \text{VC-dim}(M) = 1$$



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

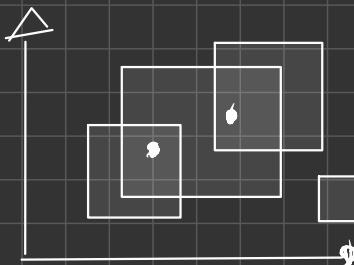
$$h_{a,b,c,d} \begin{cases} 1 & \text{if } 0 \leq a \leq b \text{ and } 0 \leq c \leq d \\ 0 & \text{otherwise} \end{cases}$$

1 point

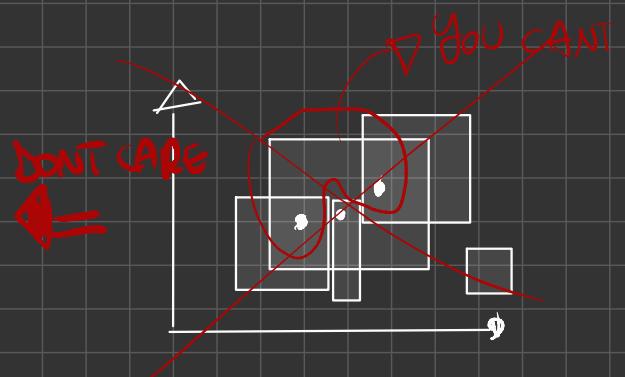
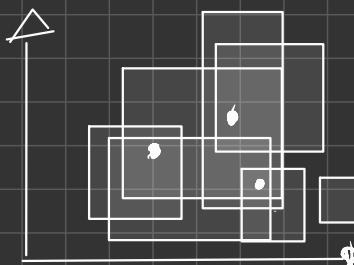


$$\Rightarrow \text{VC-dim}(H_{rect}) \geq 1$$

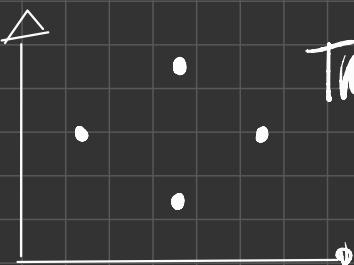
2 point



3 point



TRUST, WE CAN!



WE CANT SHATTER 6 POINTS !  
WITHOUT THE 5<sup>TH</sup> !

$\Rightarrow$  So you proved also for more than 5 points

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

MIND

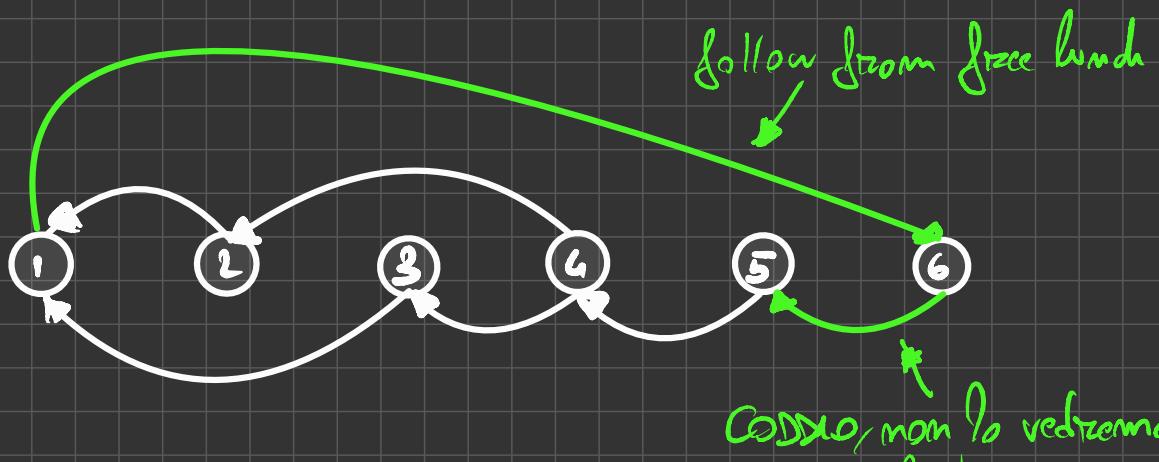
# FUNDAMENTAL TH OF STATISTICAL LEARNING

TH

- $X$ , 0-1 LOSS,  $H$
- (1)  $H$  IS PAC LEARNABLE
- (2)  $\text{ERM}_H$  IS PAC LEARNER
- (3)  $H$  IS AGNOSTIC PAC LEARNABLE
- (4)  $\text{ERM}_H$  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