

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

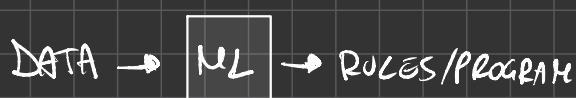


16/09

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

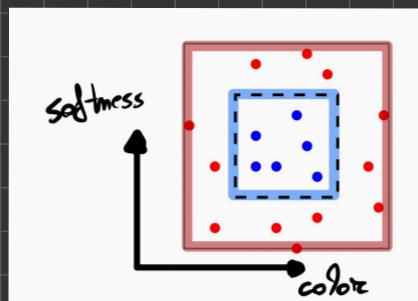
- ERM is the process of finding a PREDICTOR h which minimizes the EMPIRICAL LOSS $L_S(h)$

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:

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

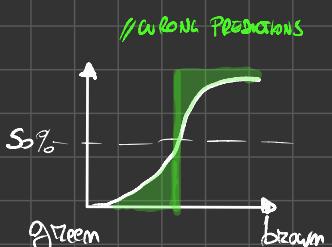
weak points
of Def¹

- 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:

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

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

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

- 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 \iff
 w.r.t. the loss ϑ

They are such that:

~~If D over (X, 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, y) such that

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

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

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

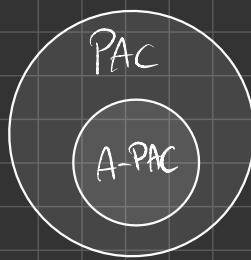
BEST POSSIBLE ERROR
OVER H .
W.R.T. REALIZABILITY WAS 0
BECAUSE OF $f \in H$

NB

• What is true?

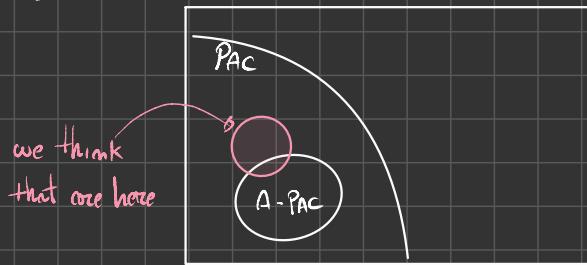
① H is A-PAC $\Rightarrow H$ is PAC

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

$$\begin{aligned} L_D(h_S) &\leq L_S(h_S) + \epsilon/\epsilon \\ &\leq L_S(h) + \epsilon/\epsilon \quad \forall h \end{aligned}$$

$$\begin{aligned} &\stackrel{\text{↑ by Def}}{\leq} L_D(h) + \epsilon/\epsilon + \epsilon/\epsilon \quad \forall h \end{aligned}$$

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

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

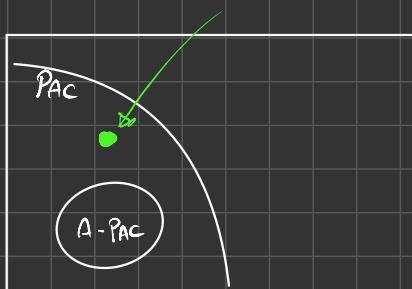
- If H is finite $\Rightarrow H$ has UC prop

Proof

Proceeding...

TH

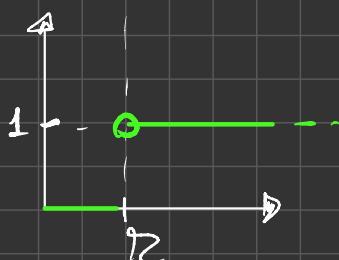
- There is at least an INFINITE PAC CLASS H



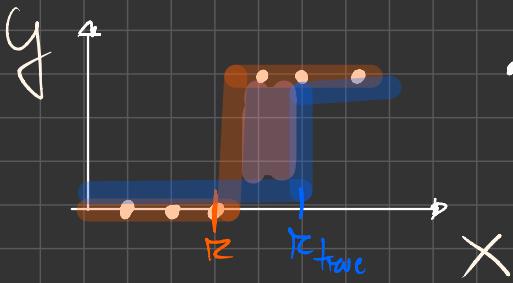
Proof

- 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

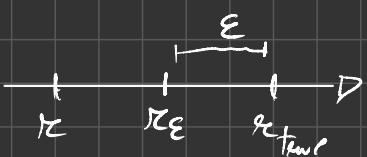


• ERM_H gives me h_r
where $r = \max\{x_i | g_i = 0\}$

- 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



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

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

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

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

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

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

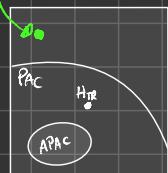
$$= \prod_{i=1}^m (1 - \epsilon) \Rightarrow (1 - \epsilon)^m \Rightarrow \text{we want } t \leq \delta$$

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

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
h_1		\square
S'	S'	$L_{D,h_1}(A(S))$
S^3	\dots	

\oplus
rows of h_2 on S'

// pick g sources at at time

// 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

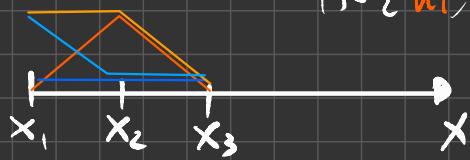
why
VC-DIMENSION?

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

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

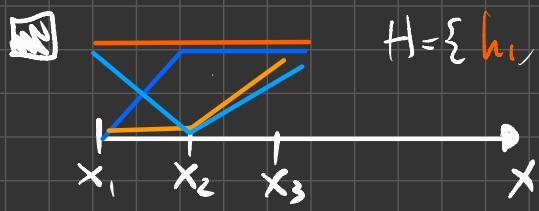
- $A = \{x_3\} \Rightarrow H$ DOESN'T (You can't get $x_3=1$ with $h(x)=0$)

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

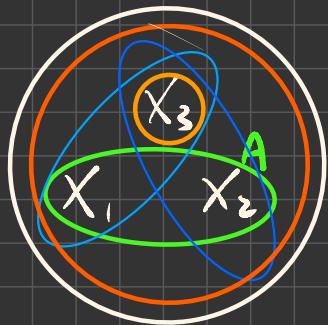


why Shattering

- 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

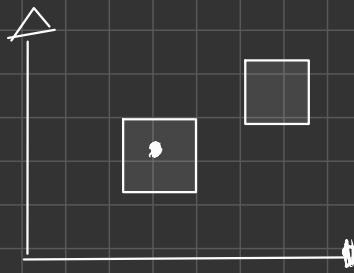
$$\Rightarrow \text{VC-dim}(M) = 1$$



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

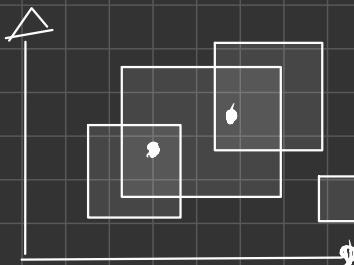
$$h_{a,b,c,d} \begin{cases} 1 & \text{if } 0 \leq a \leq b \text{ and } 0 \leq c \leq d \\ 0 & \text{otherwise} \end{cases}$$

1 point

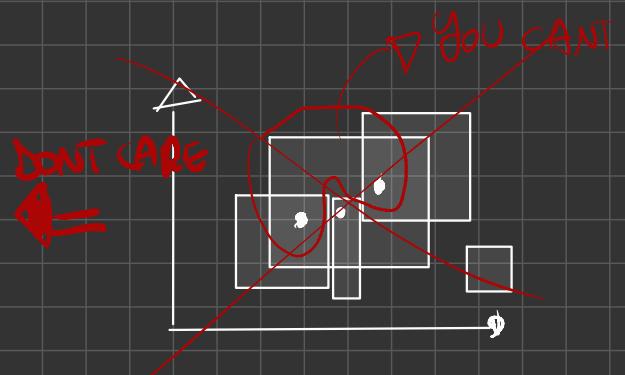
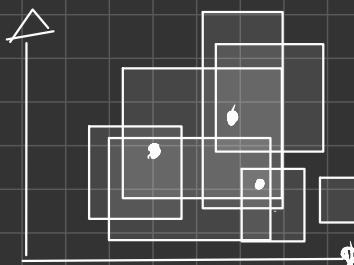


$$\Rightarrow \text{VC-dim}(H_{rect}) \geq 1$$

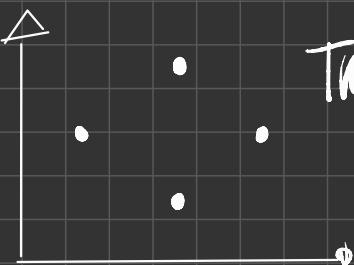
2 point



3 point



TRUST, WE CAN!



we CAN'T 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

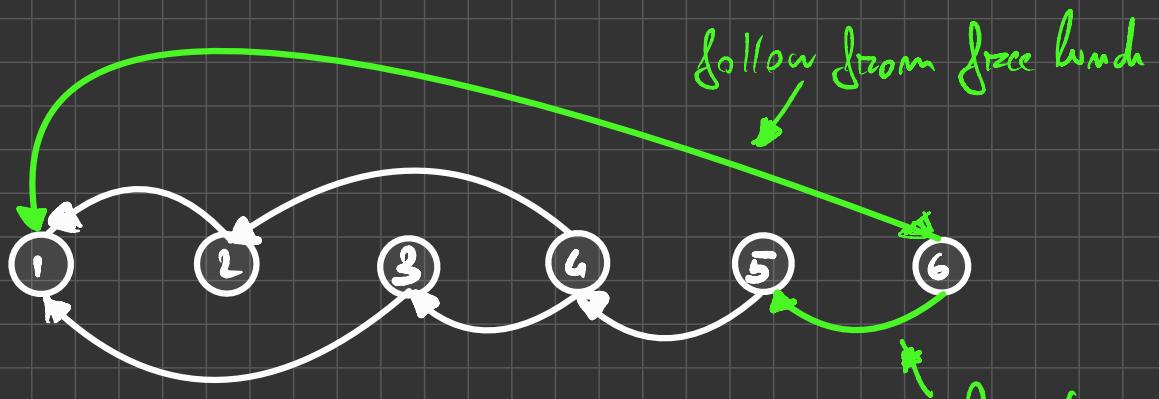
(4) ERM_H IS AGNOSTIC PAC LEARNER

(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

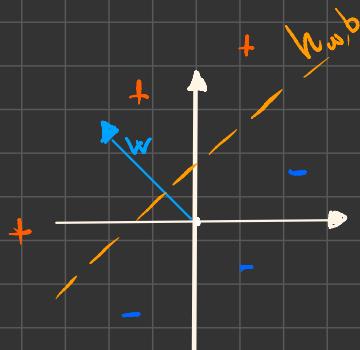
- $X = \mathbb{R}^d$, $Y = \{-1, +1\}$

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



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

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 • $\text{VC}_{\text{dim}}(\mathcal{H}_S) = d+1$ 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

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

- ESTIMATION ERROR ($\leq \epsilon$ for PAC learnability)
- APPROXIMATION ERROR

APPROX. ERR • Depends only by $\text{VC}(H)$, enlarging $|H|$ decreases ϵ_{app} (with $\delta \rightarrow 0 \Rightarrow \epsilon_{\text{app}} = 0$)

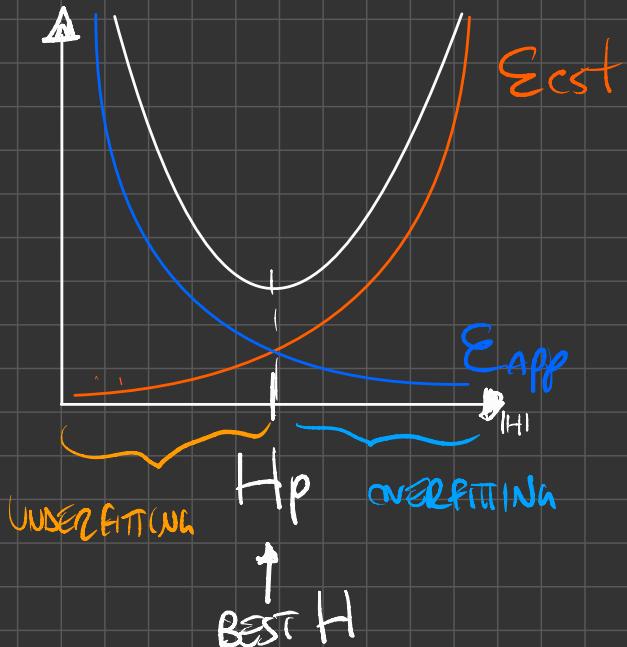
ESTIMATION ERR • Depends on $\text{VC}_d(H)$ and $|S|=m$, from The QUANT. FTSL:

$$m_H(\epsilon, d) \approx \frac{d+1 \ln(1/\delta)}{\epsilon^2} \Rightarrow \epsilon_{\text{estim}} \approx \sqrt{\frac{d+1 \ln(1/\delta)}{m}}$$

Note

- If $H \subset H'$ we have $\left\{ \begin{array}{l} \epsilon_{\text{approx}}(H') \leq \epsilon_{\text{approx}}(H) \\ \epsilon_{\text{estimation}}(H) \leq \epsilon_{\text{estimation}}(H') \end{array} \right.$

BIAS/COMPLEXITY TRADE-OFF



- What is the best possible error?

$$L_D(A(s)) =$$

Note that h is not held!

$$L_D(A(s)) - \min_{h \in H} L_D(h) + \min_{h \in H} L_D(h) - \min_h L_D(h) + \min_h L_D(h)$$

WEAK LEARNABILITY , BOOSTING

- We look for a simpler LEARNABILITY definition leading to more efficient solution.

SHARPER HOPE

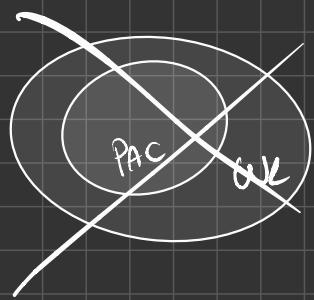
CONCLUSION

- Something NP-HARD for PAC LEARN. becomes POLYNOMIAL for WEAK LEARN. (of course they fucked up)
 - $\text{WL} = \text{PAC L.}$
 - SHAPIRE showed that a POLY ALGO. to weakly LEARN H can be reduced in Polynomial time in a PAC LEARNER for H (the sign)
 - The POLY. REDUCTION ALGORITHM is called BOOSTING
 - Boosting Allow the learner to have smooth control over BIAS/COMP. TRADEOFF
 - Think to a WEAK-LEARNER as an algo. which outputs h_{weak} performing just a bit better than a random guess.
 - Boosting is an AGGREGATOR of h_{weak} hypothesis to approximate predictors for hard to learn classes
 - ADABoost output a m h_{ADA} that is a LINEAR COMB. of h_{weak} . hypts

WEAK LEADS.

- Assuming $\delta \in H$ and 0-1 loss, H is γ -WC if $\exists A: (X, y) \xrightarrow{\text{m}} H$ and
 $\exists m_H^{\text{wc}}: (0, 1) \rightarrow \mathbb{N}$ | $\forall D$ over X , $H \not\models H$, $\forall S \subset D^m$ iid $|S| = m > m_H^{\text{wc}}(\delta)$
 theorem: $L_{D, \delta}(A(S)) \leq \frac{1}{2} - \delta$ w.p. $1 - \delta$ ($\delta \in (0, \frac{1}{2}]$)

• Because $\text{WL} \subsetneq \text{PAC-L}$
 It seems $\text{WL} \subset \text{PAC-L}$ but $\text{PAC-L} = \text{WL}$



TU

• $\text{WL} = \text{PAC-L}$

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

• If $H \in \text{PAC} \Rightarrow L_{D,+}(A(S)) < \varepsilon \text{ w.o.p. } \geq (1 - d)$

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

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

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 quasi-fcls $\forall \varepsilon, d \Rightarrow m_H(\varepsilon, d) \leq \frac{\text{VCd}(H) + \log(\frac{1}{\delta})}{\varepsilon} \leq m_H(\varepsilon, d)$

• So $L_D(A(S)) > \varepsilon$ in particular $\varepsilon = \frac{1}{2} - \gamma$



- In the example we show that computing $\text{ERM}_{H_{\text{DS}}}$ (crazier than ERM_{H_3}) is a weak learner for H_3 even if $H_{\text{DS}} \subset H_3$.
- Means that every sample selected from H_{DS} have an error w.t.f. of at most $1 - \gamma$ (obv H_{DS} must be polynomially learned by $\text{ERM}_{H_{\text{DS}}}$)

$$H_3 = \{h_{ab} : \mathbb{R} \rightarrow \{-1, +1\} \mid a, b \in \mathbb{R} \cup \{-\infty, +\infty\}\}$$

$$h_{ab}(x) = \begin{cases} +1 & \text{if } x \in [a, b] \\ -1 & \text{else} \end{cases}$$

$$\text{VCdim}(H_3) = 2$$

$$H_{\text{DS}}^1 = \{h_{a,b} \mid \mathbb{R} \rightarrow \{-1, +1\}, a \in \mathbb{R}, b \in \{-1, +1\}\}$$

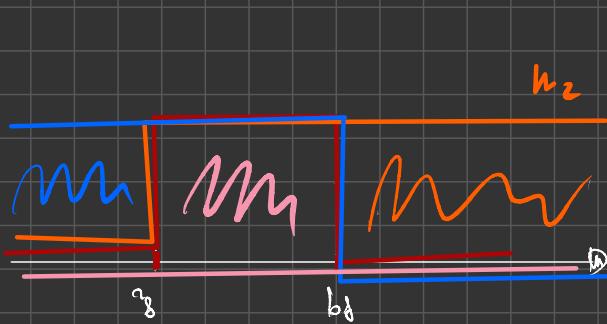
$$h_{a,b}(x) = \begin{cases} b & \text{for } x \leq a \\ -b & \text{for } x > a \end{cases}$$

$$\text{VCdim}(H_{\text{DS}}) = 2$$

- Notice that $H_D, \delta \Rightarrow \exists h \in H_{D\delta} \mid L_{D\delta}(h) \leq \frac{1}{3}$
[Also for $\delta \in H_3$]

Take $\delta \in H_3 \setminus H_{D\delta}$ can divide R in 3 regions; (δ doesn't in $H_{D\delta} \Rightarrow L_{D\delta}(h) = 0$)

- Obviously there is $h \in H_3$ that fall in 1 of the 3 regions:



- $h, h_1, h_2, h_3 \in H_{D\delta}$
- $\delta \in H_3$

$$L_{D\delta}(h_1) = \frac{1}{3}$$

$$L_{D\delta}(h_2) = \frac{1}{3}$$

$$L_{D\delta}(h_3) = \frac{1}{3}$$

- $\exists h \mid L_{D\delta}(h) \leq \frac{1}{3}$

- It implies H_D PAC because:

$$L_{D\delta}(A(s)) \leq \min_{h \in H_{D\delta}} L_{D\delta}(h) + \varepsilon \leq \frac{1}{3} + \varepsilon$$

- To prove that $H_{D\delta}$ is WCL for H_3 we must find $\varepsilon, \gamma \mid \frac{1}{3} + \varepsilon - \gamma = 1 - \gamma$ so

$$\text{for } \varepsilon = \frac{1}{12} - \gamma \text{ we have } \frac{1}{3} + \frac{1}{12} - \frac{1}{2} - \frac{1}{12}$$

$\Rightarrow H_{D\delta}$ is $\frac{1}{12}$ -WEAK LEARNER for H_3

Note

- Now with BOOSTING ALGO. we can transform it to a PAC LEARNER for H_3 .

BOOSTING ALGO

H_{DS}^d means that
x_i has only one
feature

- ERM _{H_{DS}^d} CAN be implemented in polynomial time

$$H_{DS}^d = \{h_{a,b,i} : \mathbb{R}^d \rightarrow \{-1, +1\} \mid a \in \mathbb{R}, b \in \{-1, +1\}\}$$

$$h_{a,b,i} = \begin{cases} b & \text{if } x_i < a \\ -b & \end{cases}$$

- Given $S = \{(x_i, y_i) \mid i=1, \dots, m\}$

$$\text{ERM}_{H_{DS}^d} = \underset{h \in H_{DS}^d}{\operatorname{argmin}} L_S(h) = \underset{a, b, i}{\operatorname{argmin}} L_S(h_{a,b,i})$$

- The complexity is:

COMPLEXITY

$$(m+1) \cdot 2 \cdot d \cdot \underbrace{m}_{(2)} \cdot \underbrace{|b|}_{(1)} \cdot \underbrace{d \cdot m \log m}_{(3)}$$

to compute $L_S(h_{a,b,i})$

order acc d sets of point

we can update it instead of recalculating every step

$$\Rightarrow md + md \log m$$

- QUA BOI, NON SO PERCHÉ MA:

$$\text{ERM}_{H_{DS}^d} = \underset{h}{\operatorname{argmin}} \sum_{i=1}^m D_i \vartheta(h_i; (x_i, y_i))$$

$\vartheta(h_i; (x_i, y_i)) \leq 1/(m+1)$

Probability of guess y_i

Ponré

$$L_D(h) = \sum_{i=1}^m D_i \vartheta(h_i; (x_i, y_i))$$

$$\text{so if } D = \{\frac{1}{m}, \dots, \frac{1}{m}\} \Rightarrow L_D(h) = L_S(h)$$

ADA BOOST

To learn over H_m it uses a WL over H_B WL

$$\text{ERM}_{H_B} = \underset{h \in H_B}{\operatorname{arg\,min}} \sum_{i=1}^m D_i \vartheta(h_i(x_i, y_i))$$

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

Get $i.i.d.$ samples from S

$$L_D(h_t) - \varepsilon_t \leq \frac{1}{2}\delta$$

Assigning a weight to h_t inversely proportional to ε_t

Decrease D_i for i that are well guessed and increase D_i for others
 Also will focus on problematic samples of S most

$$h_t = \text{WL}(D^{(+)}, S)$$

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

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

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

// Positive papaya have positive numerator, else negative (enforce the weights)

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

TH

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

Note

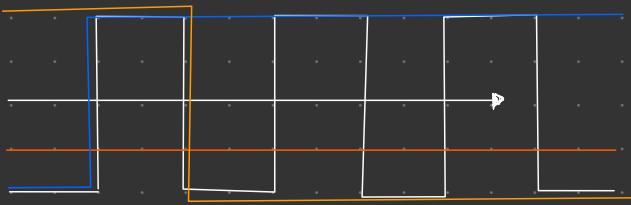
For $T \rightarrow \infty$ then $L_S(h) = 0 \rightarrow$ overfitting!

19/10

• By applying AdaBoost to $B = DS'$ how much can $h_{\text{ADA}}(x) = \text{sign}(\sum_{t=1}^T w_t h_t(x))$ become complex?

• Define weak hypothesis to sum up:

$$\begin{array}{l|l} h_1(x) = -1 & w_1 = \frac{1}{2} \\ h_2(x) = \text{sigm.}(x - x_1) & w_2 = 1 \\ h_3(x) = -\text{sigm.}(x - x_2) & w_3 = 1 \end{array}$$



• Summing those: $h_{\text{ADA}} = w_1 h_1 + w_2 h_2 + w_3 h_3$

• Notice that also if B is very poor () we can scale $h \in B$ by weights w and summing them to obtain more complex functions (h_{ADA})

• $\text{VCdim}(H_{\text{ADA}(B)}) \leq \tilde{\mathcal{O}}(T \cdot \text{VCdim}(B))$

Given m what is \rightarrow
the smallest ϵ I can
get

$$\begin{aligned} \forall \epsilon, \delta \quad \exists m \geq m_H^{\text{VC}}(\epsilon, \delta) \Rightarrow |L_S(A(S)) - L_D(A(S))| \leq \epsilon \quad w.p. \geq 1 - \delta \Rightarrow m = m_H^{\text{VC}}(\epsilon, \delta) \approx \frac{\Omega_m(1/\epsilon) + \text{VCdim}(H)}{\epsilon^2} \Rightarrow \\ \Rightarrow \epsilon = \sqrt{\frac{\text{VCdim}(H) + \Omega_m(1/\delta)}{m}} \Rightarrow |L_S(A(S)) - L_D(A(S))| \leq \sqrt{\frac{\text{VCdim}(H_{\text{ADA}}) + \Omega_m(1/\delta)}{m}} \end{aligned}$$

CONVEXITY

- LEARNABLE?
 $f_{\text{SUS}} \text{ AND } \text{VCdim}(H) < \infty$
- EFFICIENTLY LEARNABLE?
(1) BOOSTING
(2) CONVEXITY

CONVEX SET

CONVEX $\&$

- $X \subset \mathbb{R}^n$ is convex $\Leftrightarrow d x + (1-d)y \in X \quad \forall x, y \in X \quad \forall d \in [0, 1]$
- $f: X \rightarrow \mathbb{R}$ convex $\Leftrightarrow f(dx + (1-d)y) \leq d f(x) + (1-d)f(y) \quad \forall x, y \in X, d \in [0, 1], \forall X \subset \mathbb{R}^n | X \text{ convex set}$
- We can also say that f is convex if EPIGRAPH(f) is convex
- The nice property is LOCAL MIN = GLOBAL MIN



- $f: \mathbb{R} \rightarrow \mathbb{R}$, if $f''(x) \geq 0 \quad \forall x \Rightarrow f$ convex
- generalizing $f: \mathbb{R}^m \rightarrow \mathbb{R}$, $f' \rightarrow \nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_m} \end{bmatrix}, H_f = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_m} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_m \partial x_1} & \dots & \frac{\partial^2 f}{\partial x_m^2} \end{bmatrix}$
- $(H_f)_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j}$
- H_f is Positive 'SEMI-DEFINITE' $\forall x \quad x^T H_f x \geq 0$ and $\forall y \quad y^T H_f y \geq 0$

Some props.

- (1) f, g convex $\Rightarrow m \times (f, g)$ is convex
- (2) f, g convex AND g monotonically decreasing $\Rightarrow g(f(x))$ is convex
- (3) f linear, g convex $\Rightarrow g(f(x))$ convex
- (4) f, g convex $\Rightarrow df + \beta dg$ is convex $\forall \alpha, \beta \geq 0$

Convex (LEARNING) PROBLEM

$$E_{\text{RM}} = \arg \min_{h \in H} L_S(h) = \frac{1}{m} \sum_{i=1}^m \ell(h_i(x_i, y_i))$$

$$1) H = \{h_w : w \in \mathbb{R}^n\}$$

$$2) \ell(h_w(x, y)) \text{ convex in } w$$

CONVEX
OPT. PROB.

- By (g) we can define $H = \mathbb{E}_{\text{law}}[w]$, we $\mathbb{R}^d \setminus \{0\}$ so all classes with a real parameter (like H_{softmax} ...)
- So in $H = \mathbb{E}_{\text{law}}[w]$ we can look for a set of parameters $w \in \mathbb{R}^d$ instead of h : // W ARE PARAMS OF LINEAR FUNCTIONS

$$\text{ERH}_H = h \in \underset{h \in H}{\operatorname{argmin}} L_s(h) \equiv w \in \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} L_s(w)$$
- We want $L_s(w)$ to be convex in w

Note f is convex in w
 f given

$$f = w^2 + \sqrt{xy}$$

You don't care about other variables

Conclusion

- If f convex in w and $w \in \mathbb{W}$ (\mathbb{W} convex set) then ERH is a CONVEX OPT. PROBLEM

LINEAR REGRESSION (EXAMPLE)

- You can use a LINEAR CLASSIFIER for REGRESSION problem (guess a number)

$$y = \langle w, x \rangle = w \cdot x = w^T x = \sum_{i=1}^m w_i x_i$$

- The loss we use is LEAST SQUARE

$$l(h, (x, y)) = (h(x) - y)^2$$

$$l(h_w, (x, y)) = (\langle w, x \rangle - y)^2 \quad \text{// CONVEX}$$

$$l(h_{w,b}, (x, y)) = (\underbrace{\langle w, x \rangle + b - y}_{{h}_{w,b}})^2 \quad \text{// CONVEX}$$

- So given $L_s(h_{w,b}) = \frac{1}{m} \sum_{i=1}^m (\langle w, x_i \rangle + b - y_i)^2$

Prop (3)
 $\Rightarrow L_s$ is convex

- If the derivative (the gradient ∇) in higher dimensions

$$f'(x) = 0 \Rightarrow \nabla f = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$$

- Then derive w.r.t. every variable

$$\frac{\partial L_s}{\partial b} = \frac{1}{m} \sum_{i=1}^m 2(\langle w, x_i \rangle + b - y_i) = 0$$

$$\frac{\partial L_s}{\partial w_3} = \frac{1}{m} \sum_{i=1}^m 2(\langle w, x_i \rangle + b - y_i) x_{i3} = 0 \quad \text{if } i=1, \dots, m \quad \text{Attribute 3 of } x_i$$

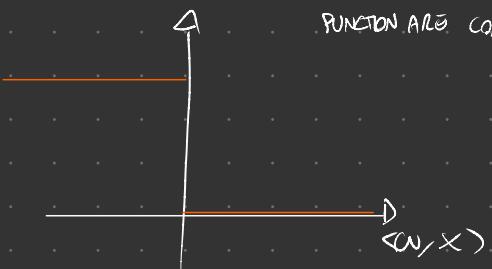
BINARY CLASSIFICATION (EXAMPLE)

$$l_{0,1}(h, (x_i, y_i)) = \mathbb{I}_{h(x_i) = y_i}$$

$$h_w(x) = \operatorname{sgn}(\langle w, x \rangle)$$

$$l_{0,1}(h_w, (x, y)) = \mathbb{I}_{y \neq \operatorname{sgn}(\langle w, x \rangle)} \quad \text{// NOT CONVEX}$$

NOTE NO PIECEWISE
 PUNCTUAL ARE CONVEX



W /
 MAX
 SPACES

SURROGATE LOSS
FOR 0-1 LOSS

- Let's find a convex function ϑ that works the same as ϑ_{0-1}

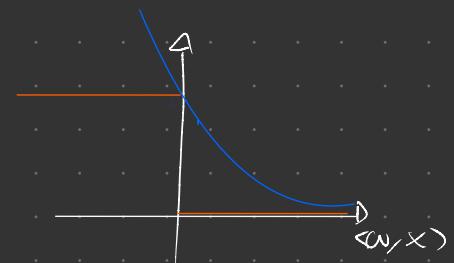
$$\vartheta(h_{\omega}, (x, y)) \geq \vartheta_{0-1}(h_{\omega}, (a, b))$$

- We can for example take ϑ as:

$$(e^{-x})^2 = e^{-x} \geq 0 \Leftrightarrow \log_2(1 + \exp(-\langle \omega, x \rangle \cdot y))$$

$$\Rightarrow \vartheta_{\text{LOG}}(1 + \exp(-\langle \omega, x \rangle \cdot y)) \quad \text{Called LOGISTIC LOSS,}$$

that is a SURROGATE FUNCTION for 0-1 loss



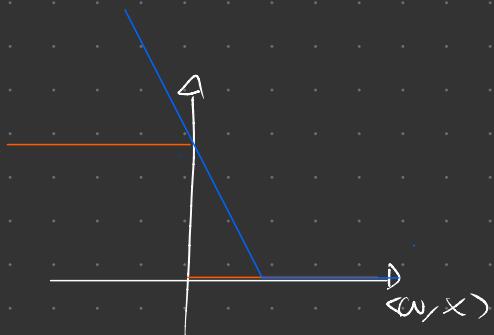
① LOGISTIC LOSS

- Minimizing Log. Loss means solving LOG. REG. PROBLEM

② HINGE LOSS

- Another surrogate loss for 0-1 loss is the blue function.

$$\max(0, -\langle \omega, x \rangle \cdot y + 1) = \text{HINGE LOSS}$$



$$\mathcal{L}_D^{0-1}(h_s) =$$

$$\mathcal{L}_D^{0-1}(h_s) - \mathcal{L}_D^{0-1}(h_s)$$

$$< 0$$

ESTIMATION ERROR

$\mathcal{L}_D^{0-1}(h_s) - \min_{h \in \mathcal{H}} \mathcal{L}_D^{0-1}(h)$

OPT. ERROR

$$\min_{h \in \mathcal{H}} \mathcal{L}_D^{0-1}(h) - \min_{h \in \mathcal{H}} \mathcal{L}_D^{0-1}(h)$$

$$\min_{h \in \mathcal{H}} \mathcal{L}_D^{0-1}(h)$$

APP. ERROR

$$\text{So } \mathcal{L}_D^{0-1}(h_s) \leq \mathcal{E}_{\text{est}} + \mathcal{E}_{\text{opt}} + \mathcal{E}_{\text{app}} \quad (+ \text{BLAS se te va!})$$

26/10

MODEL SELECTION

LINEAR CLASSIFIER

- $h_{\omega, b} = \langle \omega, x \rangle + b = \sum_{i=1}^m \omega_i x_i + b$

DISLINEAR
IN ω AND b
(NOT x)

ALSO

$$h_{\omega, b}(x) = \omega_1 x_1 + \omega_2 x_2 + \omega_3 x_3^2 + b$$

IS LINEAR IN ω BUT QUADRATIC IN x

- How to choose H ?



1) STRUCTURAL RISK MINIMIZATION (SRM)

- H is not Uniform Learnable (NUL) if

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

$$\exists m_H: (0,1)^2 \times H \rightarrow \mathbb{N} \quad // \text{NON UNIFORM w.r.t. } h \text{ depends on } h$$

Then $\forall D, \forall \epsilon, \delta \in (0,1), \forall \delta \in (0,1), \forall h \in H$, sampling $|S|=m$ iid sample according to D , if $m \geq M_H(\epsilon, \delta, h)$ then

$$L_D(A(S)) \leq L_D(h) + \epsilon \quad \text{w.p. } 1-\delta$$

COR.



- H is NUL $\Leftrightarrow H = \bigcup_{m=1}^{\infty} H_m \mid \text{VC}(H_m)$ is finite

- Given the class of all polynomials

$$H = \{ \text{SIGN}(P(x)) \mid P(x) \text{ is Polynomial} \}$$

$$\cdot \text{VC}_{\text{distr}}(H) = +\infty$$

$\Rightarrow H$ can be rewritten as $H = \bigcup_{m=1}^{\infty} H_m \mid \text{VC}_{\text{distr}}(H_m)$ finite

where H_m :

$$H_m = \{ \text{SIGN}(P(x)) \mid P(x) \text{ Polynomial with degree at most } m \}$$

$\Rightarrow H$ is NUL

CONCLUSION

- If we have a NUL-learner then we inductively solve the model selection problem whom $H = \bigcup_m H_m \mid \text{VC}_{\text{distr}}(H_m) < \infty \Rightarrow$ We need an algo. that learns in a more uniform way (No ERH!)

TM

- If $H = \bigcup_m H_m \mid H_m$ has VC prop. $\Rightarrow H$ is NUL

SRM
(NVL LEARNER)

- Given $H = \bigcup_m H_m$ with $h \in H$
- $m(h) = \min \{m \mid h \in H_m\} // \text{im } H_{\text{polynomial}}: h(x) = \text{sign}(x^2 - c) \text{ has } m(h) = 3$

Define $\varepsilon_m: N \times (0,1) \rightarrow (0,1)$

$$\varepsilon_m(m, \delta) = \min \left\{ \varepsilon \mid \underbrace{m_{H_m}^{\text{VC}}(\varepsilon, \delta)}_{m_{H_m}^{\text{VC}}(\varepsilon, \delta) \approx C \frac{\text{VCd}(H_m) - \delta_m(1/\delta)}{\varepsilon^2}} \leq m \right\} \Rightarrow$$

$$\Rightarrow \min \left\{ \varepsilon \mid \varepsilon \geq \sqrt{\frac{\text{VCd}(H_m) - \delta_m(1/\delta)}{m}} \right\} = \varepsilon_m(m, \delta)$$

- Introduce a weight for each H_m
- $H_m \rightarrow w_m \geq 0$
- $\sum_m w_m \leq 1$

IN PRACTICE IF
 $\text{VC}(H_m) \leq \text{VC}(H_{m+1})$
we give weights
 $w_m \geq w_{m+1}$

- We know that $|L_D(h) - L_S(h)| \leq \varepsilon_m(m, \delta) \leq \varepsilon_m(m, w_m \delta)$
- $L_D(h) \leq L_S(h) + \min_{m \in \text{hett}} \varepsilon_m(m, w_m \delta) \Rightarrow \text{SRM WANT TO MINIMIZE THIS QUANTITY}$

$$\Rightarrow \min_{h \in H} L_S(h) + \boxed{\sqrt{C \frac{d_{m(h)} + \delta_m(1/\delta w_{m(h)})}{m}}}$$

SRM AGO.
MINIMIZZA
QUESTA
MELDA

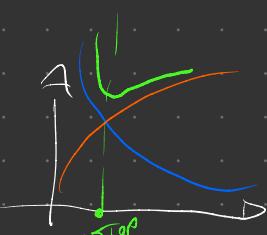
QUESTO TERMINETI
IMPEDISCE DI PRENDERE
CLASSI TROPPO COMPLESSE

Apply SRM
as an algo

- Apply $\text{ERM}_{H_m} \nmid m \Rightarrow$ we have reduced SRM to ERM, so if ERM_{H_m} is Polynomial $\nmid n$, then SFR is too.

- We take the smallest δ_{ram}

$$\min_{h \in H} L_S(h) + \sqrt{C \frac{d_{m(h)} + \delta_m(1/\delta w_{m(h)})}{m}}$$



We know that the quantity decrease until it starts to increase, when it does we stop.

ASSIGN WEIGHTS

- $h \rightarrow |h| = \# \text{ bits to code } h$

① way

$$H = \bigcup_m \{h_m\} \quad w_m = \frac{1}{2^{|h_m|}}$$

$$\cdot \mathcal{E}_m(m, \delta) = \min_m \{ \varepsilon \mid m_{H_m}^{uc}(\varepsilon, \delta) \leq m \}$$

$$m_{H_m}^{uc} = \frac{\log |H_m| + \delta m / \delta}{\varepsilon^2} = \frac{O + \delta m / \delta}{\varepsilon^2} \Rightarrow \mathcal{E}_m(m, \delta) = \sqrt{\frac{\delta m / \delta}{m}}$$

composed by one
 summation $H_m = \{h_m\}$

② way

$$\Rightarrow \min_h L_S(h) + \sqrt{\frac{\delta m / \delta}{m}} = \min_h L_S(h) + \sqrt{\frac{\delta m |h| + \delta m / \delta}{m}}$$

$$\Rightarrow w_m \propto \frac{1}{m^2}$$

\uparrow
 Proportional

2) VALIDATION/TEST SET

- Split S into : T_S , $testS$, Validation S

\downarrow TO FIND AN HYPOTHESIS \downarrow TEST HOW GOOD IS YOUR MODEL \downarrow TO SOLVE MODEL SELECTION PROBLEM

- Given $h = h_{\arg \min_{h'} L_{TS}(h')}$ we check if

$$L_b(h) \approx L_{test}(h) = \underbrace{\frac{1}{m_{test}} \sum_{(x, y) \in testS} g(h, (x, y))}_{\text{ERROR ON TESTS}}$$

- Then for h we look at

$$|L_{test}(h) - L_b(h)| \leq \mathcal{E}_{test} \sqrt{\frac{\delta m_2(H) + \delta m / \delta}{2m_{test}}}$$

- If we haven't found already h we are in model selection problem H_1, H_2, \dots, H_m so we find $ERM_{H_i} \quad H_i = 1, \dots, m$

- You have to choose from: $\{h_1, h_2, \dots, h_m\}$
- we must compute $L_D(h_1), \dots, L_D(h_m)$ but not having to recompute L over validations

$$L_{\text{vals}}(h_1), \dots, L_{\text{vals}}(h_m)$$

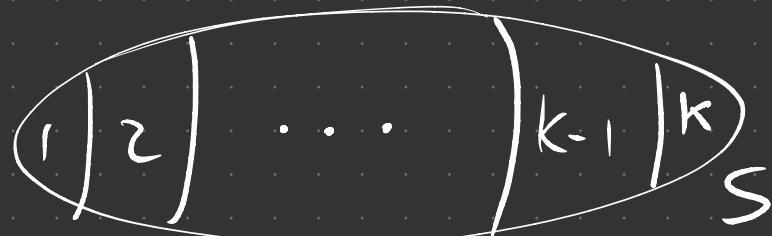
- So $h \in \arg \min_{i=1 \dots m} L_{\text{vals}}(h_i)$

$$\frac{1}{m_{\text{vals}}} \sum_{(x, y) \in \text{vals}} l(h_i(x, y))$$

- We lastly check

$$|L_{\text{vals}}(h_i) - L_D(h_i)| \leq \sqrt{\frac{f_m(2H/\delta)}{2m_{\text{vals}}}}$$

K-FOLD VALIDATION



Also

For $i = 1, \dots, m$ // m classes

FOR $j = 1, \dots, k$

$$h_{i,j} = \underset{H_i}{\text{ERM}}(S_j)$$

$$l_{i,j} = L_{S_j}(h_{i,j})$$

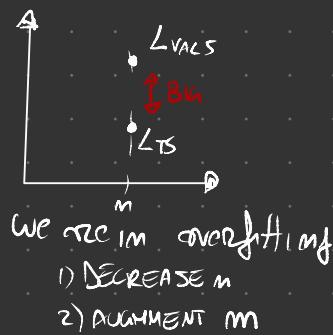
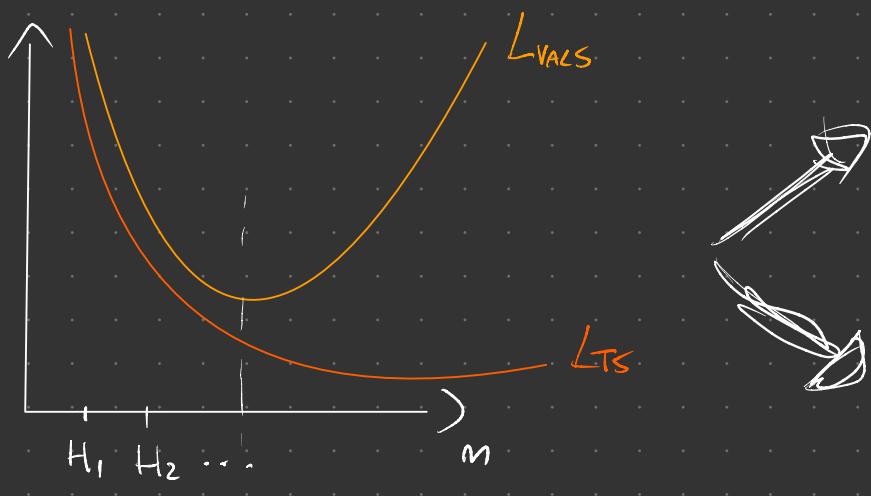
$$\bar{l}_i = \frac{1}{k} \sum_{j=1}^k l_{i,j}$$

$$i^* = \arg \min_{i=1 \dots m} \bar{l}_i$$

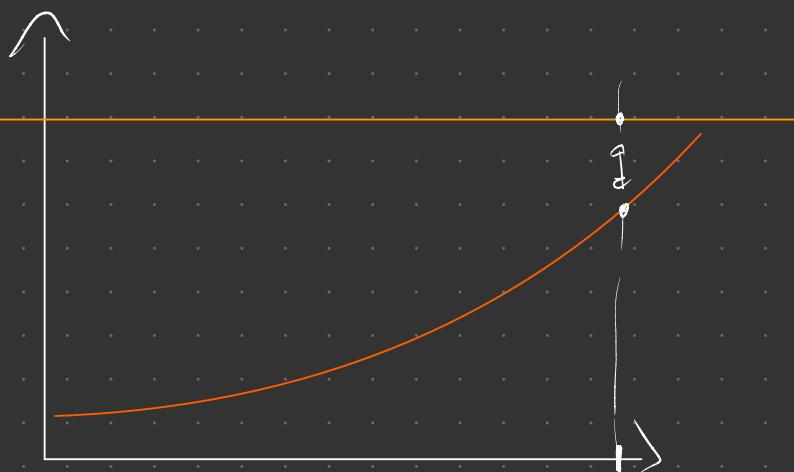
$$h = \underset{H_{i^*}}{\text{ERM}}(S)$$

VITO
PERKELE
SAATANA





if L_{TRAIN} BIG \rightarrow UNDERFITTING
 if L_{TRAIN} SMALL \rightarrow OK!



OVERFITTING
 // WE ARE NOT LEARNING ANYTHING, WE MUST SIMPLIFY THE MODEL (PICK SIMPLER CLASS)

