

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

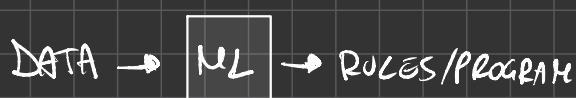


16/09

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 STATISTICAL 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

- ML doesn't know the distribution

↳ FINITE SAMPLES: you get the distribution you have

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

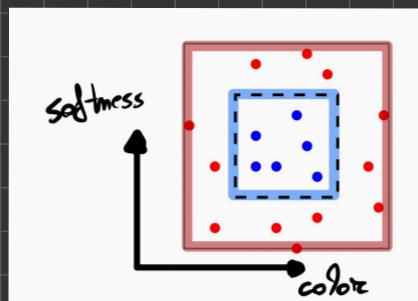
- ERM is the process of finding a PREDICTOR h which minimizes the EMPIRICAL LOSS $L_S(h)$

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

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

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

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



- 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 REGZ. 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:

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

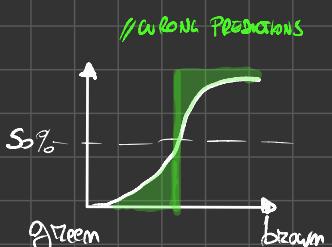
weak points of Def¹

- 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:~~

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

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

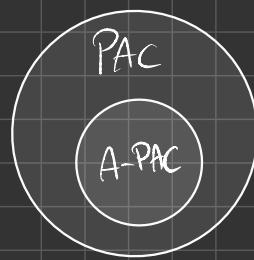
BEST POSSIBLE ERROR
OVER H .
W.R.T. REALIZABILITY WAS 0
BECAUSE OF $f \in H$

NB

• What is true?

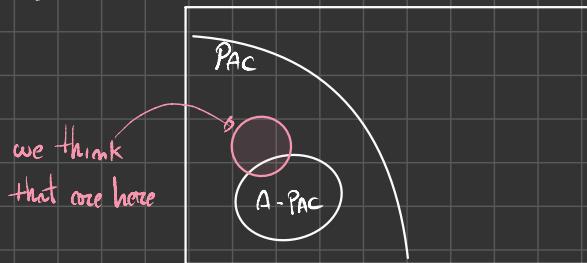
① H is A-PAC $\Rightarrow H$ is PAC

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

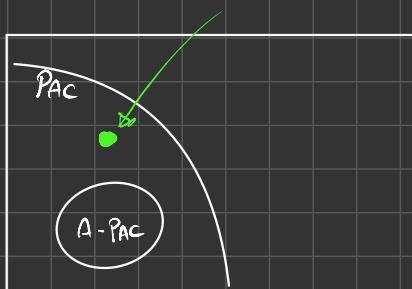
- If H is finite $\Rightarrow H$ has UC prop

Proof

Proceeding...

TH

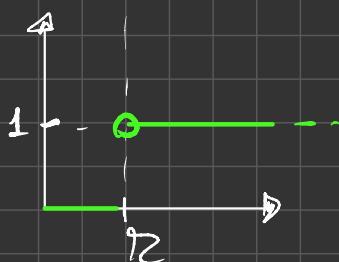
- There is at least an INFINITE PAC CLASS H



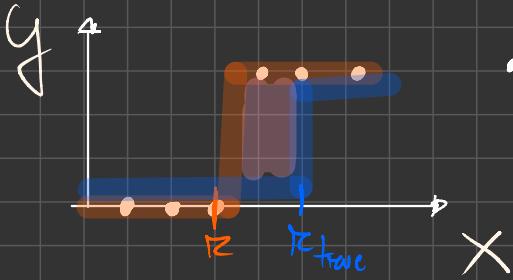
Proof

- 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

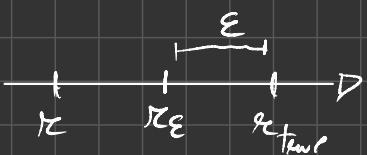


• ERM_H gives me h_R
where $\hat{r} = \max\{x \mid g_i = 0\}$

- suppose the \hat{r} , every time we observe a point in the pink region we have a loss of 1

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

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



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

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

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

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

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

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

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

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

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

