

# 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

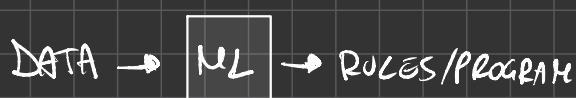


16/09

## 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

?

- From a set of predictors which one should we choose?  
Cause we do not have the D variable we cannot calculate the EXPECTED LOSS, so we use the EMPIRICAL LOSS, and take the one that minimize it. (There can be more than one)

## ERM

- EMPIRICAL RISK MINIMIZATION

$$h_{\text{ERM}} \in \arg \min L_S(h)$$

?

ERM is bad?

## MEMORIZATION ALGORITHM

- MEMORIZATION  $\rightarrow$  you learn all the couple  $\{x, y\}$ , so when in future you see a sample that is in the dataset, you know it's label, else you TOSS A COIN (or a DICE, or ...)

$$h_M(x) = \begin{cases} y_i & \text{if } x = x_i \in S \\ \text{Toss a coin} & \end{cases}$$

- We can proof that  $\exists D, \exists L_D, h_M = \frac{1}{2}$  w.p. 1  
 $\hookrightarrow$  Assume uniform prob. dist. D over X

$y(x)=1$	$y(x)=0$	// 2 feature space
----------	----------	--------------------

$\hookrightarrow$  The predictor is wrong  $\frac{1}{2}$  of the time, because you have infinite new possible sample, so the prob. of getting a sample already in the dataset is 0 (NOT IMPOSSIBLE BUT PROB=0), so the Prob. of taking a sample not in original dataset is 1, so the predictor always toss a coin! it explain perfectly the dataset but can't explain anything else.

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

?

When is good?

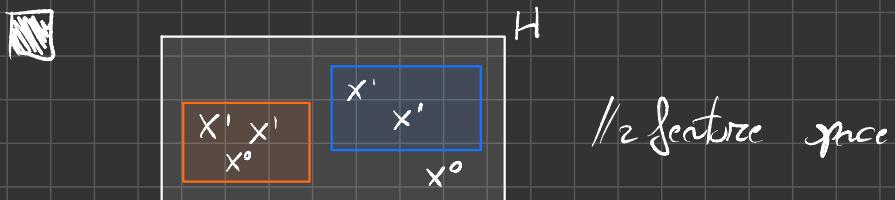
$\text{ERM}_H$

→ we should search the minimum in a finite set of hypotheses, so

$$h_{\text{ERM}_H} \in \underset{h \in H}{\operatorname{argmin}} L_S(h)$$

REALIZABILITY ASSUMPTION

- Realizability assumption impose  $f \in H$ , so if  $H$  is finite  $\text{ERM}_H$  can learn



- BLUE CH is better because it has positive sample inside and a bad outside



- A predictor is bad if

$$H_B = \{ h \in H \mid L_{D,+}(h) > \varepsilon \}$$

- A Dataset is MISLEADING if

$$S_M = \{ S \mid h_{\text{ERM}_H} \in H_B \} \subset S_p$$

POTENTIALLY  
MISLEADING  
DATASET

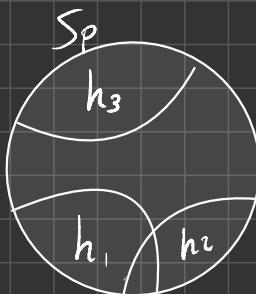
- $S_p = \{ S \mid \exists h \in H_B \text{ with } L_{D,+}(h) = 0 \}$

There is a finite number of  $h \in H_B \mid L_S(h) = 0$

$$\Pr_{D^m}(S_M) \leq \Pr_{D^m}(S_p) =$$

$$= \Pr_{D^m}\left(\bigcup_{h \in H_B} \{S \mid L_S(h) = 0\}\right) \leq$$

$$\leq \sum_{h \in H_B} \Pr_{D^m}(\{S \mid L_S(h) = 0\})$$



PROBABILITIES TO PICK DATASET

TRIVIAL •  $L_S(h) = 0 \iff \forall (x_i, y_i) \in S \quad h(x_i) = y_i$

•  $\text{Prob}(h(x_i) = y_i) \leq 1 - \varepsilon$  so

•  $\text{Prob}(h(x_i) \neq y_i, \forall i=1, \dots, m) \leq (1 - \varepsilon)^m$

$$\leq \sum_{h \in H_B} (1 - \varepsilon)^m = |H_B| (1 - \varepsilon)^m = |H| (1 - \varepsilon)^m$$

so a  
BAD ONE

we have  
proved

•  $\text{Prob}(\text{ERM}_H \text{ picks a predictor with } L_{D,t}(h) > \varepsilon) \leq |H| (1 - \varepsilon)^m$

↳ we want this quantity  $|H| (1 - \varepsilon)^m \leq \delta$

Can we choose how small should be?

↳ so if picking  $m$  large enough:  $m \geq \frac{\ln |H| / \delta}{\varepsilon}$

Conclusion • given finite  $H$ ,  $\text{ERM}_H$  can learn  $\forall \varepsilon, \delta, D, t$

$$\text{if } m \geq \frac{\ln |H| / \delta}{\varepsilon}$$

## PAC LEARNABILITY

Def

• A class  $H$  is probably approximately correct if  $\exists$  an algorithm  $A: (X \times Y)^m \rightarrow H$  (AND  $\exists$  also a function  $m_H: (0, 1)^2 \rightarrow \mathbb{N}$ ) such that  $\forall D$  over  $X$  and  $\forall \delta$  over  $Y$ ,  $\forall \varepsilon \in (0, 1)$ ,  $\forall S \subseteq (0, 1)$  whenever we pick a dataset of size  $|S|=m \geq m_H(\varepsilon, \delta)$  then we can guarantee that our algo. is such that

$$L_{D,t}(A(S)) \leq \varepsilon \text{ with } \Pr \geq (1 - \delta)$$

- A finite class  $H$  is pack-learnable with  $\text{ERM}_H$ .

$\delta$

$$m_H(\varepsilon, \delta) = \frac{\log |H|}{\varepsilon} / \delta$$



