
MACHINE LEARNING

Probably Approximately Correct (PAC) Learning

Alessandro Moschitti

Department of Information Engineering and Computer Science

University of Trento

Email: moschitti@disi.unitn.it



Introduction

- Objective: learning the concept of **medium-built people** from examples:
 - *Interesting features* are: Height and Weight.
 - The **training-set** of examples has a cardinality of m .
(m people for who we know if they are medium-built people size, their height and their size).
- Find m to learn this concept *well*.
- The adjective “well” can be expressed with probability error.



PAC Learning Definition (1)

- Let c be the function (i.e. a *concept*) we want to learn
- Let h be the learned concept and x an instance (e.g. a person)
- $error(h) = Prob [c(x) \neq h(x)]$
- It would be useful if we could find:
- $Pr(error(h) > \epsilon) < \delta$
- Given a target error ϵ , the probability to make a larger error is less δ

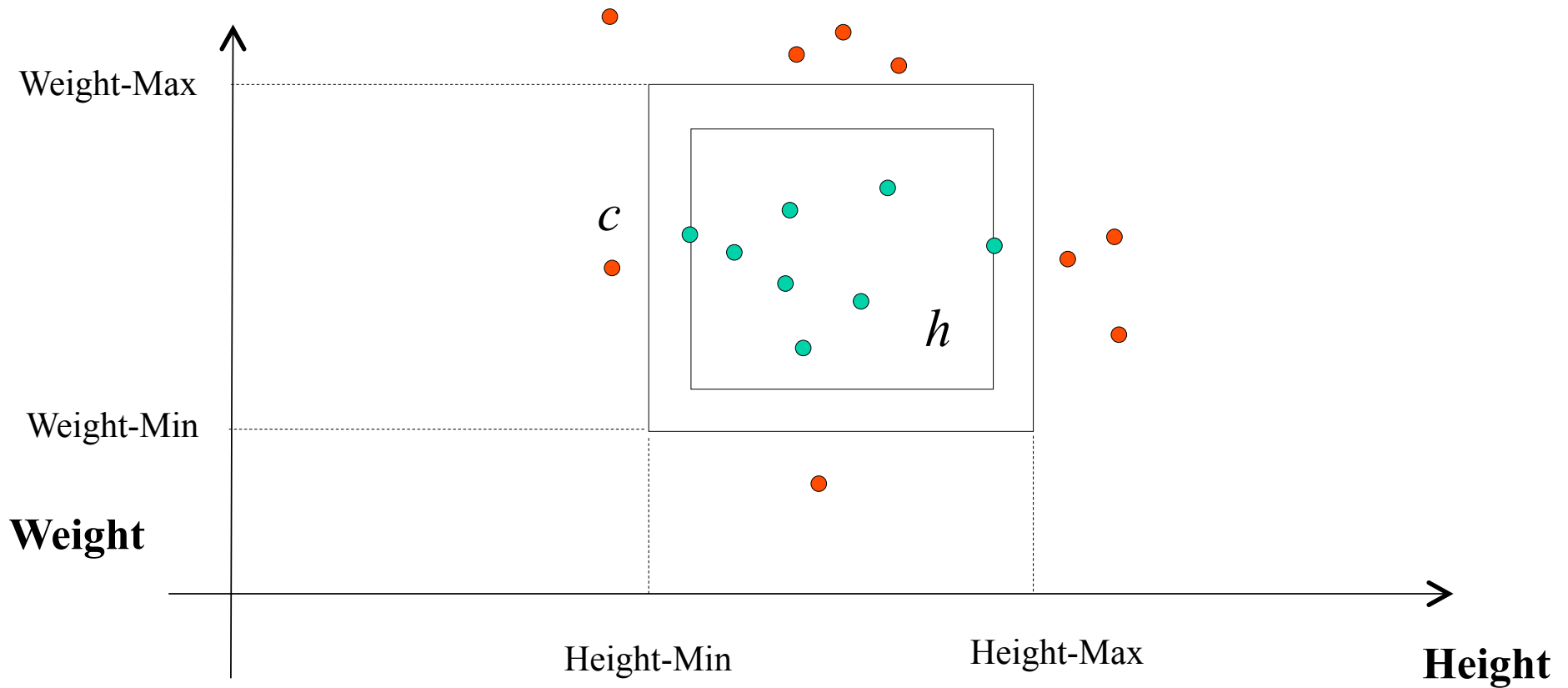


Definizione di PAC Learning (2)

- This methodology is called Probability Approximately Correct Learning
- The more ε and δ are smaller the better the learning is
- Problem:
 - Given ε and δ , determine the size m of the training-set.
 - Such size may be independent of the learning algorithm

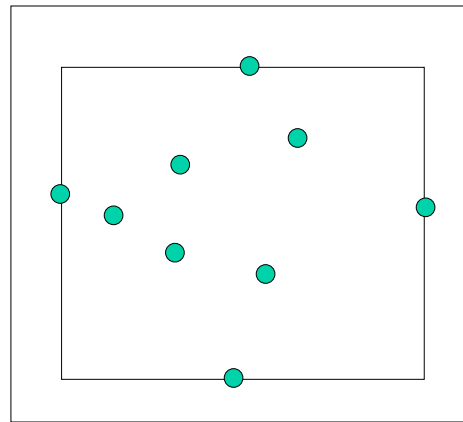


Graphical Representation of the target learning problem

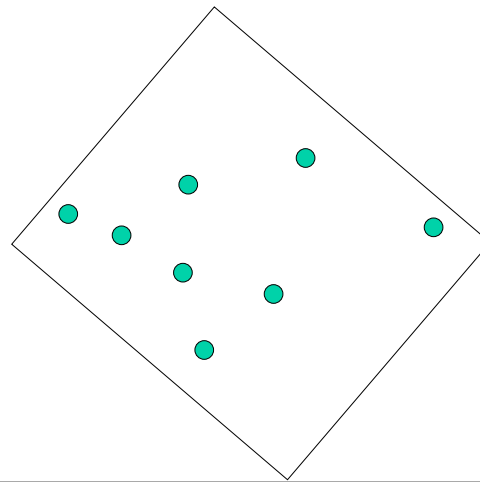


Learning Algorithm/Learning Function Class

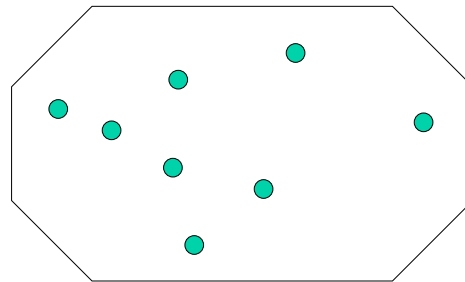
1. If no positive examples of the concept are available
⇒ the learned concept is NULL
2. Else the concept is the smallest rectangular (parallel to the axes) containing all positive examples



We don't consider other complex hypothesis

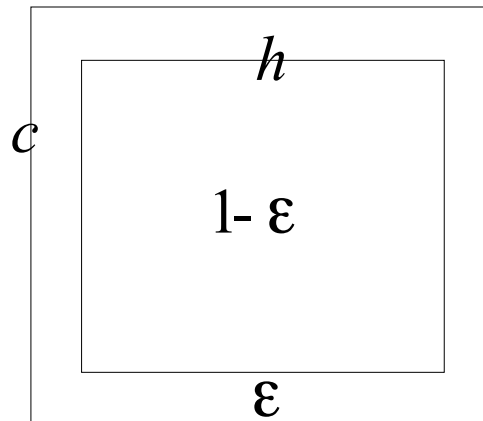


We don't consider other complex hypothesis



How good is our algorithm?

- An example x is misclassified if h falls between two rectangles.
 - Let ε be the measure of the area
- ⇒ The error probability (error) of h is ε
- With which assumption?



Proving PAC Learnability

- Given an error ϵ and a probability δ , how many training examples m are needed to learn the concept?
- We can find a bound to δ , *i.e.* the probability of learning a function h with an error $> \epsilon$.
- For this purpose: let us compute the probability of selecting a hypothesis h which correctly classifies m training example and shows an error greater than ϵ .
 - This is a *bad* function



Probability of Bad Hypothesis

- Given x , $P(h(x) \neq c(x)) < \epsilon$
 - since the error of bad function is greater than ϵ
- Given ϵ , m examples fall in the rectangle h with a probability $< (1-\epsilon)^m$
- The probability of choosing a bad hypothesis h is $(1-\epsilon)^m \cdot N$
 - where N is the number of hypothesis with an error $> \epsilon$.



Upper-bound Computation

- If we set a bound on the probability of bad hypothesis

$$N \cdot (1-\varepsilon)^m < \delta$$

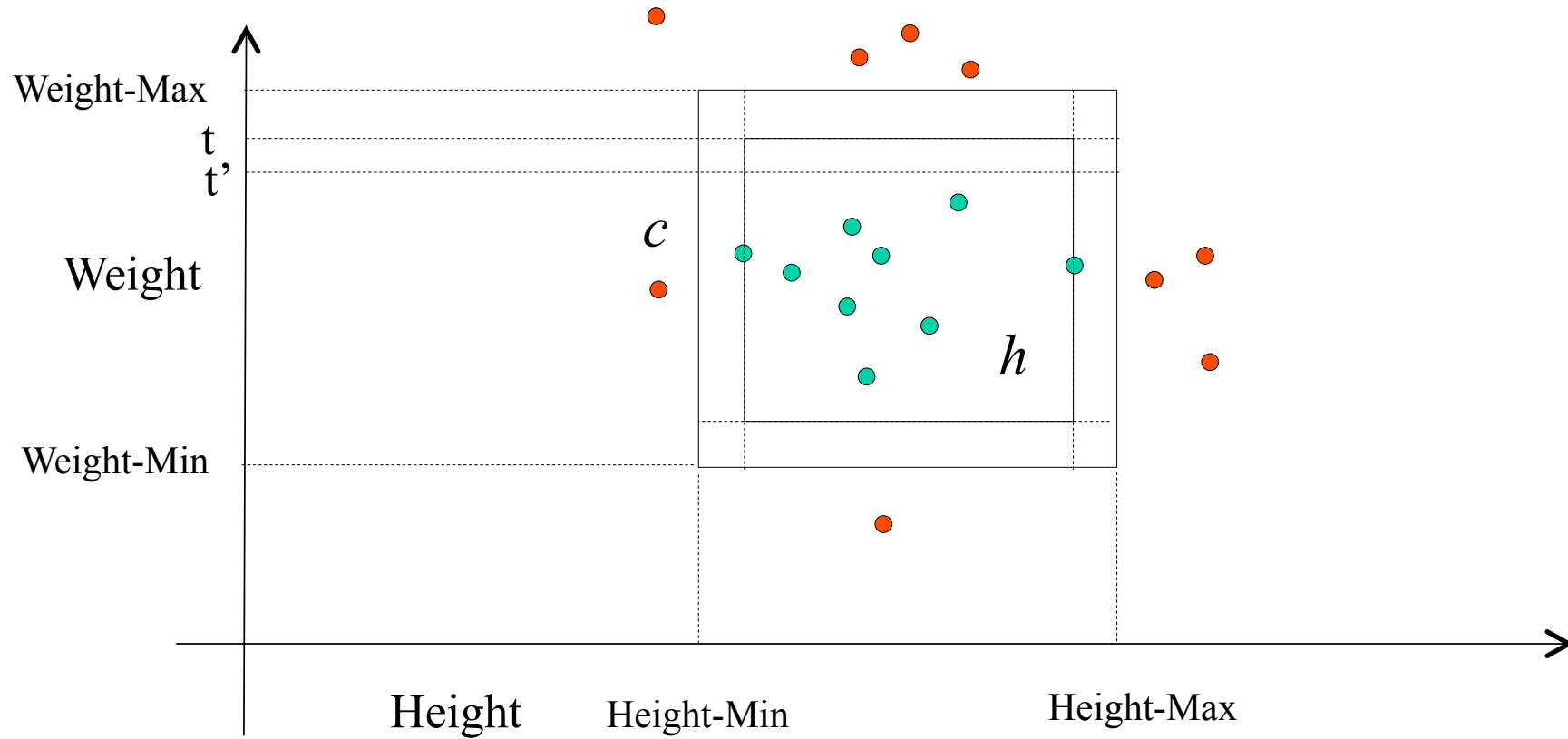
- we would be done but we don't know N

\Rightarrow we have to find a bound, independent of the number of bad hypothesis.

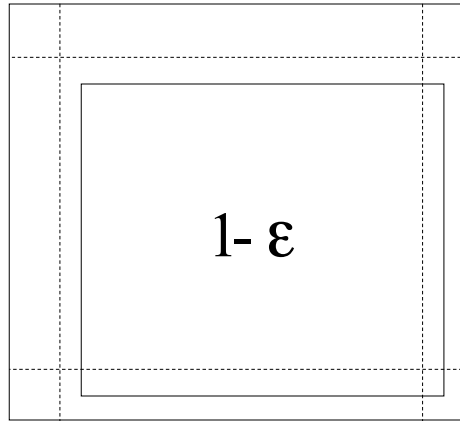
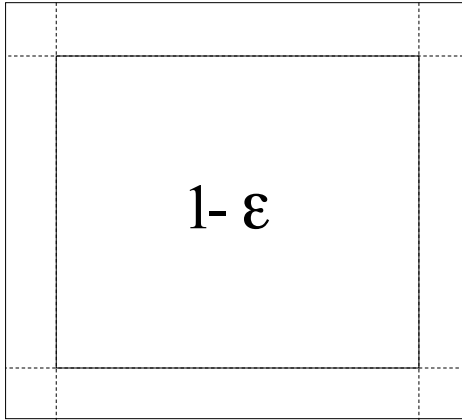
- Let us divide our rectangle in four strip of area $\varepsilon/4$



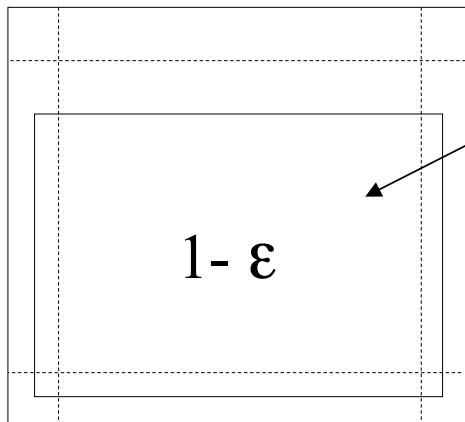
Initial Example



A bad hypothesis *cannot intersect more than 3 strips at a time*



Bad hypothesis with error $> \epsilon$ are contained in those having an error $= \epsilon$



To intersect 3 edges I can increase the rectangle length but I must decrease the height to have an area $\leq 1 - \epsilon$



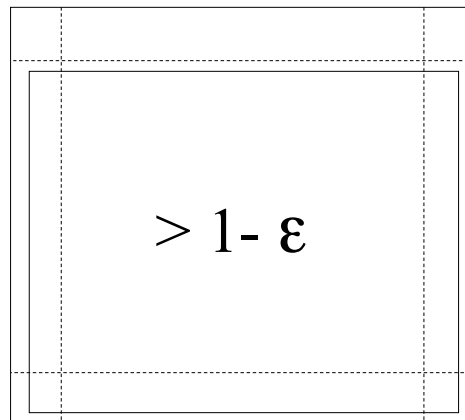
Upper-bound computation (2)

- A bad hypothesis has error $> \epsilon \Rightarrow$ it has an area $< 1 - \epsilon$
- A rectangle of area $< 1 - \epsilon$ cannot intersect 4 strips \Rightarrow if the examples fall into all the 4 strips they cannot be part of the same bad hypothesis.
- A necessary condition to have a bad hypothesis is that all the m examples are at least outside of one strip.
- In other words, when m examples are outside of one of the 4 strips we may have a bad hypothesis.
 \Rightarrow the probability of “*outside at least one of the strips*” $>$ probability of bad hypothesis.



Logic view

- Bad Hypothesis \Rightarrow examples out of at least one strip
 - (viceversa is not true)



- $A \Rightarrow B$
- $P(A) < P(B)$
- $P(\text{bad hyp.}) < P(\text{out of one strip})$



Upper-bound computation (3)

- $P(x \text{ out of the target strip}) = (1 - \varepsilon/4)$
 - $P(m \text{ points out of the target strip}) = (1 - \varepsilon/4)^m$
 - $P(m \text{ points out of at least one strip}) < 4 \cdot (1 - \varepsilon/4)^m$
- $\Rightarrow P(\text{error}(h) > \varepsilon) < 4 \cdot (1 - \varepsilon/4)^m$



Expliciting m

- Our upperbound must be lower than δ

- $4 \cdot (1 - \varepsilon/4)^m < \delta$

- $m > \ln(\delta/4) / \ln(1 - \varepsilon/4)$

- $-\ln(1-y) = y + y^2/2 + y^3/3 + \dots$

$\Rightarrow (1-y) < e^{(-y)}$

- From $m > \ln(\delta/4) / \ln(1 - \varepsilon/4)$

$\Rightarrow m > (4/\varepsilon) \cdot \ln(4/\delta)$



Numeric Examples

ε	δ	m
0.1	0.1	148
0.1	0.01	240
0.1	0.001	332

0.01	0.1	1476
0.01	0.01	2397
0.01	0.001	3318

0.001	0.1	14756
0.001	0.01	23966
0.001	0.001	33176



Formale PAC-Learning Definition

- Let f be the function we want to learn, $f: X \rightarrow I, f \in F$
- D is a probability distribution on X
 - used to draw training and test test
- $h \in H$,
 - h is the learned function and H the set of such function class
- m is the training-set size
- $error(h) = Prob [f(x) \neq h(x)]$
- F is a PAC learnable function class if there is a **learning algorithm** such that for each f , for all distribution D over X and for each $0 < \epsilon, \delta < 1$, **produces** $h : P(error(h) > \epsilon) < \delta$



Lower Bound on training-set size

- Let us reconsider the first bound that we found:
 - h is bad: $error(h) > \epsilon$
 - $P(f(x)=h(x))$ for m examples is lower than $(1 - \epsilon)^m$
 - Multiplying by the number of bad hypothesis we calculate the probability of selecting a bad hypothesis
 - $P(\text{bad hypothesis}) < N \cdot (1 - \epsilon)^m < \delta$
 - $P(\text{bad hypothesis}) < N \cdot (e^{-\epsilon})^m = N \cdot e^{-\epsilon m} < \delta$

$$\Rightarrow m > (1/\epsilon) (\ln(1/\delta) + \ln(N))$$

This is a general lower bound



Example

- Suppose we want to learn a boolean function in n variable
- The maximum number of different function are 2^{2^n}

$$\begin{aligned}\Rightarrow m &> (1/\varepsilon) (\ln(1/\delta) + \ln(2^{2^n})) = \\ &= (1/\varepsilon) (\ln(1/\delta) + 2^n \ln(2))\end{aligned}$$



Some Numbers

n		epsilon		delta		m
=====						
5		0.1		0.1		245
5		0.1		0.01		268
5		0.01		0.1		2450
5		0.01		0.01		2680

10		0.1		0.1		7123
10		0.1		0.01		7146
10		0.01		0.1		71230
10		0.01		0.01		71460
===== =						



References

- PAC-learning:
 - **MY SLIDES:** <http://disi.unitn.it/moschitti/teaching.html>
 - **MY BOOK:**
 - Artificial Intelligence: a modern approach (Second Edition) by Stuart Russell and Peter Norvig
 - <http://www.cis.temple.edu/~ingargio/cis587/readings/pac.html>
 - Machine Learning, Tom Mitchell, McGraw-Hill.

