# MACHINE LEARNING

## Vapnik-Chervonenkis (VC) Dimension

Alessandro Moschitti

Department of Information Engineering and Computer Science
University of Trento
Email: moschitti@disi.unitn.it

# Computational Learning Theory

- The approach used in rectangular hypotheses is just one simple case:
  - Medium-built people
  - No general rule has been derived

- Is there any means to determine if a function is PAC learnable and derive the right bound?

- The answer is yes and it is based on the Vapnik-Chervonenkis dimension  (VC-dimension,  [Vapnik 95])
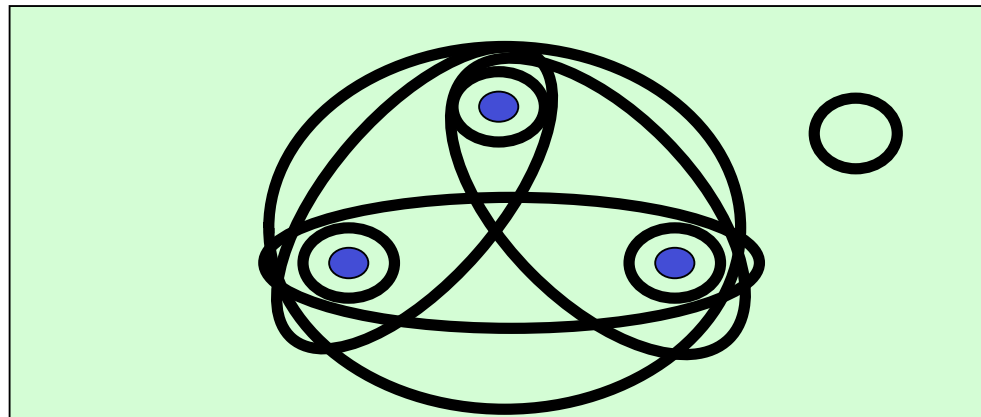
# VC-Dimension definition (1)

- Def.1: (*set shattering*): a subset S of instances of a set X is shattered by a collection of function $F$ if $\forall$ S'$\subseteq$ S there is a function $f \in F$ such data:

$$f(x) = \begin{cases} 1 & x \in S' \\ 0 & x \in S - S' \end{cases}$$

# VC-Dimension definition (2)

- Def. 2: the VC-dimension of a function set $F$ (VC-dim($F$)) is the cardinality of the largest dataset that can be shattered by $F$

- Observation: the type of the functions used for shattering data determines the VC-dim
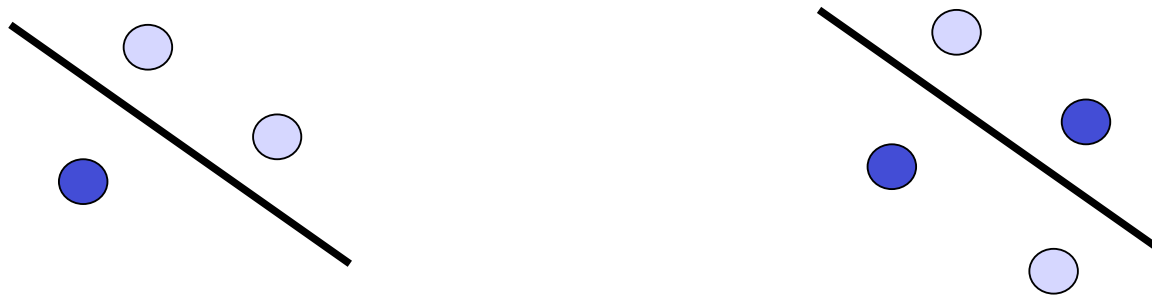
# VC-Dim of linear functions (hyperplane)

- In the plane (hyperplane = line):
  - VC (Hyperplanes) is at least 3
  - VC (Hyperplanes) < 4 since there is no set of 4 points, which can be shattered by a line.

$\Rightarrow$ VC(H)=3. In general, for a k-dimension space VC(H)=k+1

- NB: It is useless selecting a set of linearly independent points

# Upper Bound on Sample Complexity

**Theorem 2.9** *(upper bound on sample complexity, [Blumer et al., 1989])*
*Let $H$ and $F$ be two function classes such that $F \subseteq H$ and let $A$ an algorithm that derives a function $h \in H$ consistent with $m$ training examples. Then, $\exists c_0$ such that $\forall f \in F$, $\forall D$ distribution, $\forall \epsilon > 0$ and $\delta < 1$ if*

$$m > \frac{c_0}{\epsilon}\left( VC(H) \times ln\frac{1}{\epsilon} + \frac{1}{\delta} \right)$$

*then with a probability $1 - \delta$,*

$$error_D(h) \leq \epsilon,$$

*where VC(H) is the VC dimension of $H$ and $error_D(h)$ is the error of $h$ according to the data distribution $D$.*

# Lower Bound on Sample Complexity

**Theorem 2.10** *(lower bound on sample complexity, [Blumer et al., 1989])*
*To learn a concept class $F$ whose VC-dimension is $d$, any PAC algorithm requires $m = \Omega((d(H) + \ln(1/\delta))/\epsilon)$*

# Bound on the Classification error using VC-dimension

**Theorem 2.11** *(Vapnik and Chervonenkis, [Vapnik, 1995])*

*Let $H$ be a hypothesis space having VC dimension $d$. For any probability distribution $D$ on $X \times \{-1, 1\}$, with probability $1 - \delta$ over $m$ random examples $S$, any hypothesis $h \in H$ that is consistent with $S$ has error no more than*

$$error(h) \leq \epsilon(m, H, \delta) = \frac{2}{m}\Big(d \times ln\frac{2e \times m}{d} + ln\frac{2}{\delta}\Big),$$

*provided that $d \leq m$ and $m \geq 2/\epsilon$.*

# Example: Rectangles for learning medium-built person concept have VC-dim > 4

- We must choose 4-point set, which can be shattered in all possible ways

- Given such 4 points, we assign them the {+,-} labels, in all possible ways.

- For each labeling it must exist a rectangle which produces such assignment, i.e. such classification

# Example (cont'd)

- Our classifier: inside the rectangle positive and outside negative examples, respectively

- Given 4 points (linearly independent), we have the following assignments:

a) All points are "+" $\Rightarrow$ use a rectangle that includes them

b) All points are "-" $\Rightarrow$ use a empty rectangle

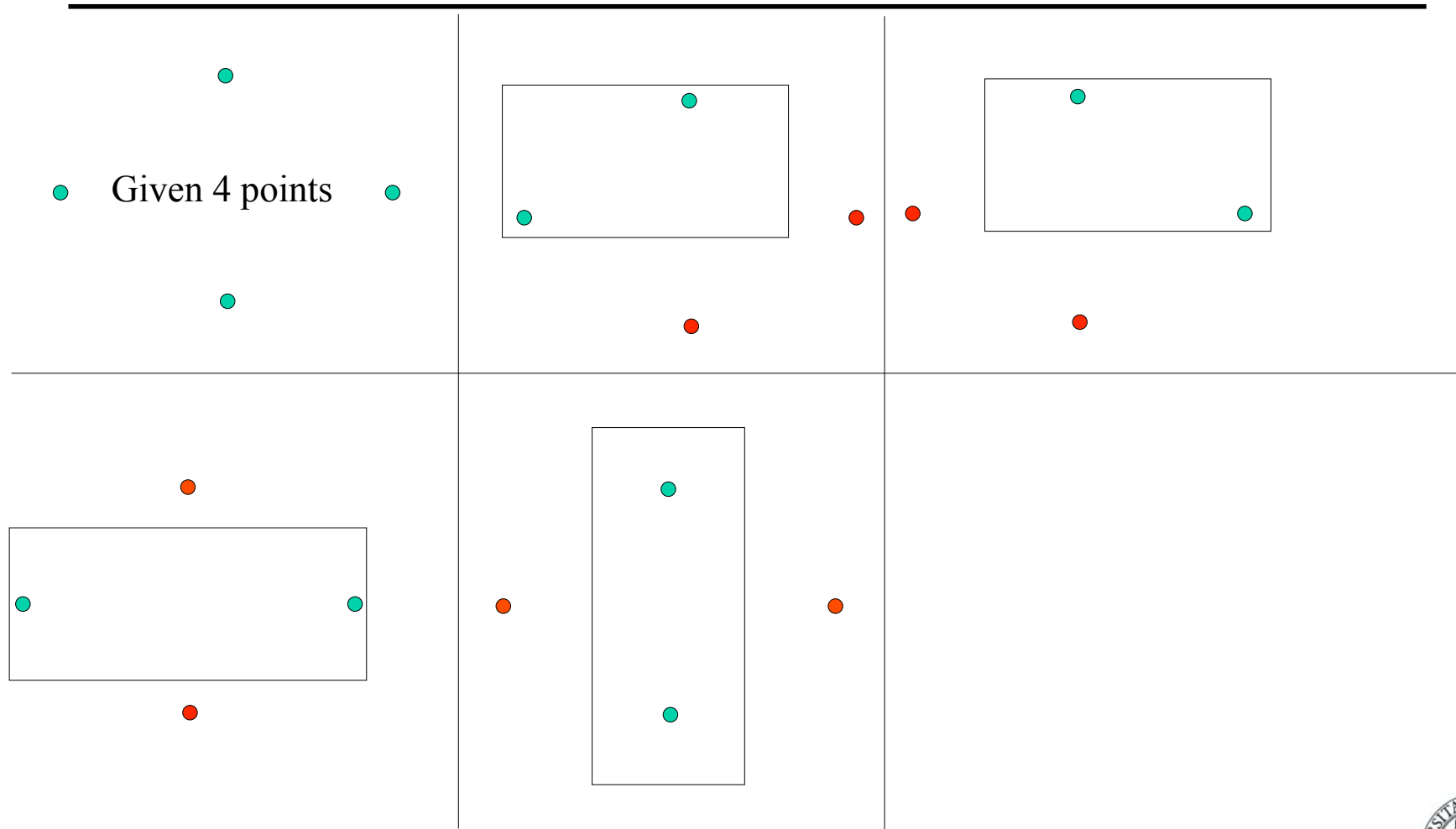c) 3 points "-" and 1 "+" $\Rightarrow$ use a rectangle centered on the "+" points

# Example (cont'd)

d) 3 points "+" and one "-" $\Rightarrow$ we can always find a rectangle which excludes the "-" points

e) 2 points "+" and 2 points "-" $\Rightarrow$ we can define a rectangle which includes the 2 "+" and excludes the 2 "-".

- To show d) and e) we should check all possibilities

# For example, to prove e)

Given 4 points

# *VC-dim* cannot be 5

- For any 5-point set, we can define a rectangle which has the most extern points as vertices

- If we assign to such vertices the "+" label and to the internal point the "-" label, there will not be any rectangle which reproduces such assigment

# Applying general lower bound to rectangles

**Theorem 2.10** *(lower bound on sample complexity, [Blumer et al., 1989])*
*To learn a concept class $F$ whose VC-dimension is $d$, any PAC algorithm requires $m = \Omega((d(H) + \ln(1/\delta))/\epsilon)$*

- $m = O((4 + \ln(1/\delta))/\epsilon))$

# Bound Comparison (lower bound)

- $m > (4/\varepsilon) \cdot ln(4/\delta)$   (ad hoc bound)

- $m = O((1/\varepsilon) \cdot (ln(1/\delta) + 4)) =$  (lower bound based on VC-dim)

- Does the ad hoc bound satisfy the general bound?

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

$\Leftrightarrow ln(4/\delta) > ln(1/\delta)/4 + 1 \Leftrightarrow ln(1/\delta) + ln(4) > ln(1/\delta)/4 + 1$

$\Leftrightarrow ln(4) > (-1 + 1/4)ln(1/\delta) + 1 \Leftarrow ln(4) > 1$

$\Leftrightarrow ln(4) > ln(e)$

# References

- VC-dimension:

  - **MY SLIDES: http://disi.unitn.it/moschitti/ teaching.html**

  - **MY BOOK:**

    - Automatic text categorization: from information retrieval to support vector learning
    - Roberto Basili and Alessandro Moschitti

# References

- *A tutorial on Support Vector Machines for Pattern Recognition*
  - **Downlodable from the web**

- *The Vapnik-Chervonenkis Dimension and the Learning Capability of Neural Nets*
  - **Downlodable from the web**

- Computational Learning Theory
  (Sally A Goldman Washington University St. Louis Missouri)
  - **Downlodable from the web**

- *AN INTRODUCTION TO SUPPORT VECTOR MACHINES*
  *(and other kernel-based learning methods)*
  N. Cristianini and J. Shawe-Taylor Cambridge University Press
  - **You can buy it also on line**

# Other Web References

- On the sample complexity of PAC learning half spaces against the uniform distribution, Philip M. Long.

- A General Lower Bound on the Number of Examples Needed for Learning, Andrzej Ehrenfeucht, David Haussler, Michael Kearns and Leslie Valiant

- BOUNDS ON THE NUMBER OF EXAMPLES NEEDED FOR LEARNING FUNCTIONS, Hans Ulrich Simon

- Learnability and the Vapnik-Chervonenkis Dimension, ANSELM BLUMER, ANDRZEJ EHRENFEUCHT, DAVID HAUSSLER AND MANFRED K. WARMUTH

- A Preliminary PAC Analysis of Theory Revision, Raymond J. Mooney

- The Upper Bounds of Sample Complexity, http://mathsci.kaist.ac.kr/~nipl/am621/lecturenotes.html

# Proposed Exercises

- Try to formulate the concept medium-built people with squares instead of rectangles and apply the content of the PAC learning lecture to this new class of functions.

- Could you build a better ad-hoc bound than the one we evaluated in class? (assume that the concept to learn is a square and not a rectangle)

# Propose Exercises

- Evaluate the VC-dimension (of course in a plane) for
  - squares
  - circles
  - equilateral triangles
  - Sketch the proof of VC < k but do not spend to much time in formalizing such proof.

- Compare the lower-bound to the sample complexity using squares (calculated with VC dimension) with your ad hoc bound derived from medium-built people (as we did it in class for rectangles).