

Machine Learning

VC-Dimension

Fabio Vandin

December 11th, 2023

PAC Learning

Question: which hypothesis classes \mathcal{H} are PAC learnable?

Up to now: if $|\mathcal{H}| < +\infty \Rightarrow \mathcal{H}$ is PAC learnable.

What about \mathcal{H} : $|\mathcal{H}| = +\infty$? Not PAC learnable?

We focus on:

- *binary classification:* $\mathcal{Y} = \{0, 1\}$
- 0-1 loss

but similar results apply to other learning tasks and losses.

Restrictions

Definition (Restriction of \mathcal{H} to \mathcal{C})

Let \mathcal{H} be a class of functions from \mathcal{X} to $\{0, 1\}$ and let $\mathcal{C} = \{c_1, \dots, c_m\} \subset \mathcal{X}$. The restriction $\mathcal{H}_{\mathcal{C}}$ of \mathcal{H} to \mathcal{C} is:

$$\mathcal{H}_{\mathcal{C}} = \{[h(c_1), \dots, h(c_m)] : h \in \mathcal{H}\}$$

where we represent each function from \mathcal{C} to $\{0, 1\}$ as a vector in $\{0, 1\}^{|\mathcal{C}|}$.

Note: $\mathcal{H}_{\mathcal{C}}$ is the set of functions from \mathcal{C} to $\{0, 1\}$ that can be derived from \mathcal{H} .

VC-dimension and Shattering

Definition (Shattering)

Given $C \subset \mathcal{X}$, \mathcal{H} shatters C if \mathcal{H}_C contains all $2^{|C|}$ functions from C to $\{0, 1\}$. $|C|=m \Rightarrow 2^{|C|}=2^m=|\mathcal{H}_C|$

Definition (VC-dimension)

The *VC-dimension* $VCdim(\mathcal{H})$ of a hypothesis class \mathcal{H} , is the maximal size of a set $C \subset \mathcal{X}$ that can be shattered by \mathcal{H} .

Notes: *ogni predizione possibile*

- VC = Vapnik-Chervonenkis, that introduced it in 1971
- if \mathcal{H} can shatter sets of arbitrarily large size then we say that $VCdim(\mathcal{H}) = +\infty$; $\nexists |C|=m \Rightarrow 2^m \text{ insiemi possibili} \Rightarrow 2^m \leq |\mathcal{H}|$
- if $|\mathcal{H}| < +\infty \Rightarrow VCdim(\mathcal{H}) \leq \log_2 |\mathcal{H}|$

Intuition: the VC-dimension measures the complexity of \mathcal{H} (\approx how large a dataset that is perfectly classified using the functions in \mathcal{H} can be)

Example

\uparrow
 X

	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9
h_1	0	0	1	0	0	0	1	0	0
h_2	0	1	0	0	0	1	0	0	0
h_3	1	0	0	0	1	1	0	0	0
h_4	0	0	0	1	1	0	0	0	1
h_5	0	0	1	0	0	0	0	1	0
h_6	0	1	0	0	0	0	1	0	0
h_7	1	0	0	0	0	1	0	0	0
h_8	0	0	0	0	0	0	0	0	0

$H \Leftarrow$

VC dimension?

Scoprire $C \subset X$ + grande | C frantumato da $H \Rightarrow$ tutte le funzioni possibili = tutti i vettori possibili

$$\text{VC dim}(H) \geq 1 \Rightarrow C = \{x_3\} \Rightarrow H_C = \{[0], [1]\}$$

simile per $\geq 2 (\{x_5, x_6\}) \Rightarrow$ non può essere 3: serve almeno un numero ≤ 1 , non c'è

Note

To show that $\text{VCdim}(\mathcal{H}) = d$ we need to show that:

- ① $\text{VCdim}(\mathcal{H}) \geq d$
- ② $\text{VCdim}(\mathcal{H}) \leq d$

that translates to

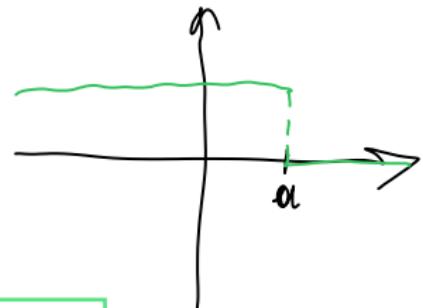
- ① there exists a set C of size d which is shattered by \mathcal{H}
- ② every set of size $d + 1$ is not shattered by \mathcal{H}

Question: why don't we need to consider sets of size $> d + 1$?

se potessi avere C di grandezza $d+1$ da sottrarre, ma non ce l'ho

Example: Threshold Functions

$$\mathcal{H} = \{h_a : a \in \mathbb{R}\}$$



where $h_a : \mathbb{R} \rightarrow \{0, 1\}$ is

$$h_a(x) = \mathbb{1}[x < a] = \begin{cases} 1 & \text{if } x < a \\ 0 & \text{if } x \geq a \end{cases}$$

$$|\mathcal{H}| = +\infty$$

VC-dimension?

$$\text{VCdim}(\mathcal{H}) = 1 ? \Rightarrow$$

For $x = c$:

- $h_{a_1}(c) = 0 \Rightarrow a_1 < c$
- $h_{a_2}(c) = 1 \Rightarrow a_2 > c$

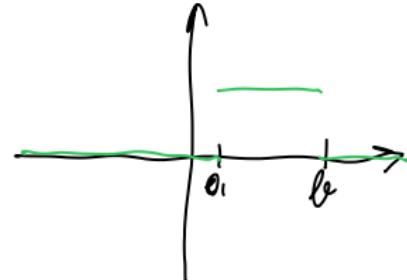
VCdim(\mathcal{H}) = 2 ?

impossible

$$\text{VCdim}(\mathcal{H}) \leq 1 \Rightarrow = 1$$

Example: Intervals

$$\mathcal{H} = \{h_{a,b} : a, b \in \mathbb{R}, a < b\}$$

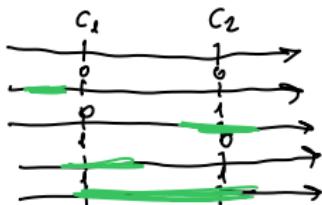


where $h_{a,b} : \mathbb{R} \rightarrow \{0, 1\}$ is

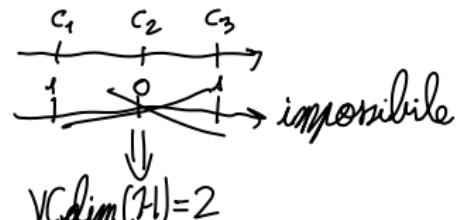
$$h_{a,b}(x) = \mathbb{1}[x \in (a, b)] = \begin{cases} 1 & \text{if } a < x < b \\ 0 & \text{otherwise} \end{cases}$$

VC-dimension?

$$\text{VC}_{\text{dim}}(\mathcal{H}) \geq 2$$



$$\text{VC}_{\text{dim}}(\mathcal{H}) \leq 2:$$

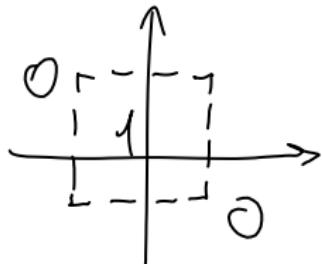


Example: Axis Aligned Rectangles

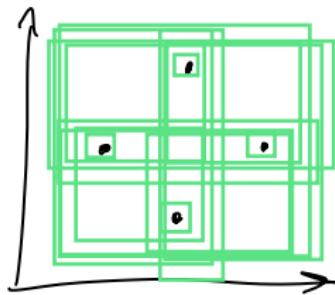
$$\mathcal{H} = \{ h_{(a_1, a_2, b_1, b_2)} : a_1, a_2, b_1, b_2 \in \mathbb{R}, a_1 \leq a_2, b_1 \leq b_2 \}$$

$$h_{(a_1, a_2, b_1, b_2)}(x_1, x_2) = \begin{cases} 1 & \text{if } a_1 \leq x_1 < a_2, b_1 \leq x_2 \leq b_2 \\ 0 & \text{otherwise} \end{cases}$$

VC-dimension?



$$VC_{dim}(\mathcal{H}) \geq 4$$



$$VC_{dim}(\mathcal{H}) \leq 4$$

\downarrow $i \leftarrow i + \alpha dx$ impossible

$i = 1 \leftarrow 1 + \alpha dx$

$i = 0 \leftarrow 0 + \alpha dx$

$$\begin{matrix} q_1 \\ q_2 \end{matrix}$$

$i \leftarrow i + \alpha dx$

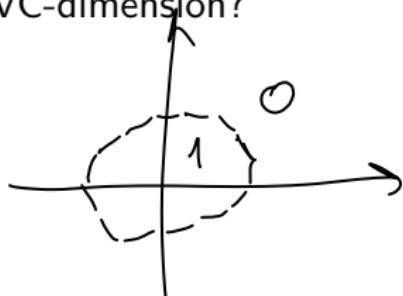
Example: Convex Sets

Model set \mathcal{H} such that for $h_s \in \mathcal{H}, h_s: \mathbb{R}^2 \rightarrow \{0, 1\}$ with

$$h_s(x) = \begin{cases} 1 & \text{if } x \in S \\ 0 & \text{otherwise} \end{cases}$$

where S is a convex subset of \mathbb{R}^2

VC-dimension?

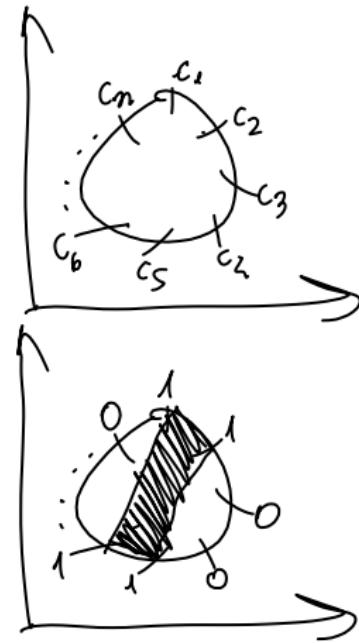


Prendiamo valore arbitrario di $n \in \mathbb{N}$ (

dim. di insieme da fronteggiare)

Poi, labeling arbitrario

ipotesi corrispondente a insieme convesso
dato da (x_i, y_i) con $y_i = 1$, da labeling
desiderato $\Rightarrow H$ può fronteggiare set di
 n punti $\forall n \Rightarrow V_{\text{dim}}(H) = +\infty$



in \mathbb{R}^2



soltuzione: 3
i) $VC_{dim} \geq 3 \Rightarrow$



possere coprire ogni
funzione possibile

Exercise

Consider the classification problem with $\mathcal{X} = \mathbb{R}^2$, $\mathbb{Y} = \{0, 1\}$.

Consider the hypothesis class $\mathcal{H} = \{h_{(\mathbf{c}, a)}, \mathbf{c} \in \mathbb{R}^2, a \in \mathbb{R}\}$ with

$$h_{(\mathbf{c}, a)}(\mathbf{x}) = \begin{cases} 1 & \text{if } \|\mathbf{x} - \mathbf{c}\| \leq a \\ 0 & \text{otherwise} \end{cases}$$

Find the VC-dimension of \mathcal{H} .

i) $VC_{dim} \leq 3 \Rightarrow$ 3 intarsio \Rightarrow 3 vari:

- 3 intarsio formo triangolo e 3° è dentro

- // // // // fiori 1° 0° i 0° \Rightarrow imp.

$\therefore \Rightarrow$ diag. + lunga \Rightarrow imp.

- stessa linea 1° 0° 0° 1° \Rightarrow imp.

The Fundamental Theorems of Statistical Learning

Theorem

Let \mathcal{H} be a hypothesis class of functions from a domain \mathcal{X} to $\{0, 1\}$ and consider the 0-1 loss function. Assume that $V\text{Cdim}(\mathcal{H}) = d < +\infty$. Then there are absolute constants C_1, C_2 such that

- \mathcal{H} has the uniform convergence property with sample complexity

$$C_1 \frac{d + \log(1/\delta)}{\varepsilon^2} \leq m_{\mathcal{H}}^{UC}(\varepsilon, \delta) \leq C_2 \frac{d + \log(1/\delta)}{\varepsilon^2}$$

- \mathcal{H} is agnostic PAC learnable with sample complexity

$$C_1 \frac{d + \log(1/\delta)}{\varepsilon^2} \leq m_{\mathcal{H}}(\varepsilon, \delta) \leq C_2 \frac{d + \log(1/\delta)}{\varepsilon^2}$$

Equivalently:

Theorem

Let \mathcal{H} be an hypothesis class with VC-dimension

$VCdim(\mathcal{H}) < +\infty$. Then, with probability $\geq 1 - \delta$ (over $S \sim \mathcal{D}^m$) we have:

$$\forall h \in \mathcal{H}, L_{\mathcal{D}}(h) \leq L_S(h) + C \sqrt{\frac{VCdim(\mathcal{H}) + \log(1/\delta)}{2m}}$$

where C is a universal constant.

Note: finding $h \in \mathcal{H}$ that minimizes the upper bound (above) to $L_{\mathcal{D}}(h) \Rightarrow$ ERM rule

Theorem

Let \mathcal{H} be a class with $VCdim(\mathcal{H}) = +\infty$. Then \mathcal{H} is not PAC learnable.

Notes:

- the VC-dimension characterizes PAC learnable hypothesis classes

Exercise

Let

$$\mathcal{H}_d = \{h_{\mathbf{w}}(\mathbf{x}) : h_{\mathbf{w}}(\mathbf{x}) = \text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle)\}$$

where $\mathcal{X} = \mathbb{R}^d$.

Prove that $VCdim(\mathcal{H}_d) = d$.

An Interesting Example...

Note: in previous examples the VC-dimension is equivalent to the number of parameters that define the model... but it is not always the case!

Function of one parameter: $f_\theta(x) = \sin^2 [2^{8x} \arcsin \sqrt{\theta}]$

VC-dimension of $f_\theta(x)$ is infinite!

In fact, $f_\theta(x)$ can approximate any function $\mathbb{R} \rightarrow \mathbb{R}$ by changing the value of θ !

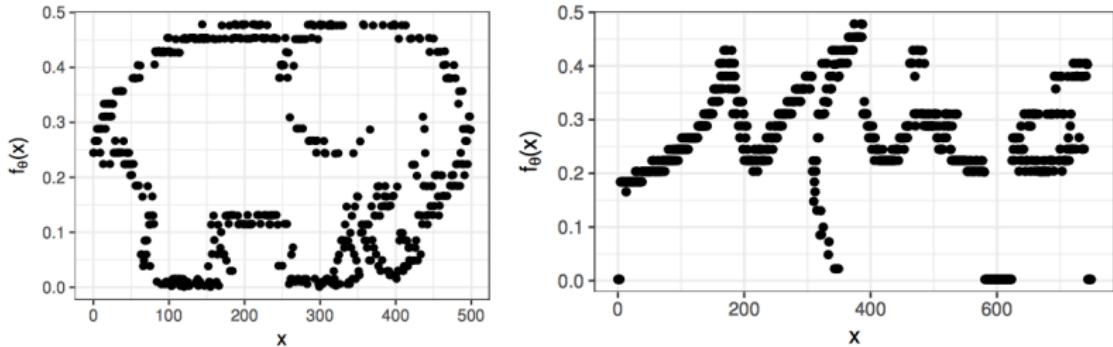


FIG. 1: A scatter plot of f_θ for $\theta = 0.2446847266734745458227540656\cdots$ plotted at integer x values, showing that a single parameter can fit an elephant (left). The same model run with parameter $\theta = 0.0024265418055000401935387620\cdots$ showing a fit of a scatter plot to Joan Miró's signature (right). Both use $r = 8$ and require hundreds to thousands of digits of precision in θ .

[“One parameter is always enough”, Piantadosi, 2018]

Bibliography

[UML] Chapter 6