

VC dimension : exercises

October 15, 2017

Notation :

We will consider functions $f : \chi \mapsto \{0, 1\}$.

If F is a class of such functions and x_1, \dots, x_n is a family of n points in χ , we define the set $N_F(x_1, \dots, x_n)$ as the set of all images of this family of points by the functions in F :

$$N_F(x_1, \dots, x_n) = \{(f(x_1), \dots, f(x_n)), f \in F\}$$

We define then the shattering coefficient of F with respect to n points sets in χ , denoted $S(F, n)$, as :

$$S(F, n) = \max |N_F(x_1, \dots, x_n)|$$

where the maximum is taken over all possible sets $(x_1, \dots, x_n) \in \chi^n$.

Finally, we define the VC dimension of F as :

$$\text{VC}(F) = \max \{n \geq 1, S(F, n) = 2^n\}$$

Exercises :

Determine the VC dimension of the next sets of functions where $\chi = [0, 1]$:

- $F = \{f : \chi \mapsto \{0, 1\}, f(x) = 1_{x < t}, t \in [0, 1]\}$
- $F' = \{f : \chi \mapsto \{0, 1\}, f(x) = 1_{x < t} \text{ or } f(x) = 1 - 1_{x < t}, t \in [0, 1]\}$
- $F = \{f : \chi \mapsto \{0, 1\}, f(x) = 1_{t_1 \leq x < t_2}, t_1 < t_2 \in [0, 1]\}$
- $F' = \{f : \chi \mapsto \{0, 1\}, f(x) = 1_{t_1 \leq x < t_2} \text{ or } f(x) = 1 - 1_{t_1 \leq x < t_2}, t_1 < t_2 \in [0, 1]\}$
- $F_k = \{f : \chi \mapsto \{0, 1\}, f(x) = \sum_{i=0}^k 1_{t_{2i} \leq x < t_{2i+1}}, \text{ for } 0 \leq t_0 < \dots < t_{2k+1} \leq 1\} \text{ for any } k \geq 1$

Note here that for any F , F' is essentially the same set of functions, the only difference being that it allows to label the points indifferently 1 against 0, or 0 against 1. This apparently harmless technical enhancement is actually not totally insignificant as the VC dimension of F and F' are different.

Solutions

- Obviously any set of one point can be shattered, so $VC(F) \geq 1$. Moreover, if you take two points x_1 and x_2 (assume $x_1 < x_2$ without loss of generality) then if x_1 is labeled 1 and x_2 labeled 0, the set cannot be shattered by any function in F . Therefore $VC(F) = 1$.
- Now, if you take two points x_1 and x_2 (assume $x_1 < x_2$ without loss of generality) all possible labeling of the points is reachable by putting $x_1 < t < x_2$, $t < x_1$ or $t > x_2$. So $VC(F') \geq 2$. If you take three points x_1 , x_2 and x_3 (assume $x_1 < x_2 < x_3$ without loss of generality), then for example there is no way that you can label x_1 and x_3 with the value 1, and x_2 with the value 0. So $VC(F') = 2$.
- If you take two points x_1 and x_2 (assume $x_1 < x_2$ without loss of generality) all possible labeling of the points is reachable by putting $t_1 < x_1 < t_2 < x_1$, $x_1 < t_1 < x_2 < t_2$, $t_1 < x_1 < x_2 < t_2$ or $x_1 < t_1 < t_2 < x_2$. So $VC(F) \geq 2$. If you take three points x_1 , x_2 and x_3 (assume $x_1 < x_2 < x_3$ without loss of generality) then there is no way that you can label x_1 and x_3 with the value 1, and x_2 with the value 0. Thus $VC(F) = 2$.
- With F' you can label x_1 and x_3 with the value 1, and x_2 with the value 0. This was the only labeling that was impossible with the previous F , therefore $VC(F') \geq 3$. With four points x_1 , x_2 , x_3 and x_4 (assumed increasing as always), you cannot label x_1 and x_3 with the value 1 and x_2 and x_4 with the value 0 for example. So $VC(F') = 3$.
- It's clear that the "worst" labeling you can encounter is when the labels are alternating $(0, 1, 0, 1, \dots)$. Here "worst" means that if you can do this one you can do any other labeling. Now in this kind of configuration, if you have $k + 1$ labels 1 (and therefore $2(k + 1)$ points in your set) it's clear that you can label all of them by putting one of the $k + 1$ "doors" of your function over each one of the $k + 1$ labels 1 of your set of points (the set F_k being the set of all functions with $k + 1$ doors). So $VC(F_k) \geq 2(k + 1)$. Moreover if you have $2(k + 1) + 1$ points, you can create a configuration of alternating labels with $k + 2$ labels 1, by starting and ending by 1 $(1, 0, 1, \dots, 0, 1)$. this last configuration is unreachable with $k + 1$ doors. Therefore $VC(F_k) = 2(k + 1)$.

Svolgimento Dettagliato. $X = [0, 1]$

③

Schemi Generale procedura

$\text{VCdim}(F) = \text{VC}(F) = d$ se esiste un sottoinsieme di X , diciamo Y , tale che $|F_Y| = 2^d = 2^{|Y|}$

dove $F_Y = \{(\mathbb{1}_Y(x), -, \mathbb{1}_Y(x)) ; f \in F\}$

e $\forall Z \subset X$ con $|Z| = d+1$ si ha $|F_Z| < 2^{d+1} = 2^{|Z|}$

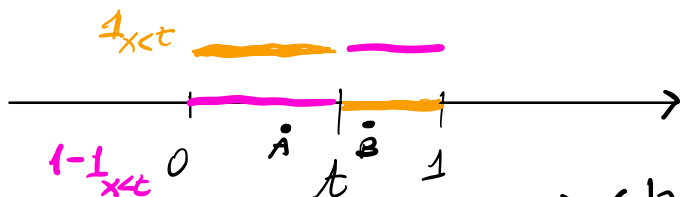
① $F = \{f: X \rightarrow \{0, 1\} : f(x) = \mathbb{1}_{x < t}, t \in [0, 1]\}$

Per un qualsiasi $x \in [0, 1)$ si ha che $f(x) = 1$ per $t = 1$. D'altro canto, prendendo $t \leq x$ si ha $f(x) = 0$ quindi $\text{VCdim}(F) \geq 1$.

D'altro canto \forall coppia di punti distinti $x < y$ in $[0, 1]$ si ha che la funzione $(x, y) \rightarrow (0, 1)$ non è realizzabile da f per alcuna $f \in F$.

Quindi $\text{VCdim}(F) = 1$

② $F' = \{f: X \rightarrow \{0, 1\} : f(x) = \mathbb{1}_{x < t} \text{ o } f(x) = 1 - \mathbb{1}_{x < t}, t \in [0, 1]\}$

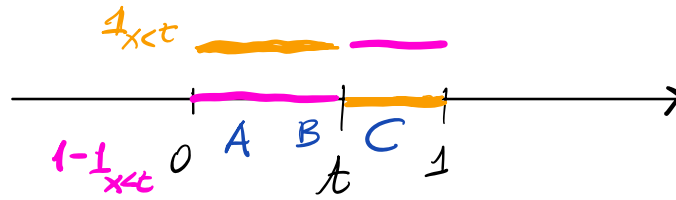


Chiusamente $\text{VCdim}(F') \geq 1$ (vedi sopra e considera che $F' \supset F \Rightarrow \text{VCdim}(F') \geq \text{VCdim}(F)$).

Considerando A, B come sopra, allora è chiaro che muovendo t possiamo ottenere $(0, 0), (1, 0), (0, 1), (1, 1)$ quindi $\text{VCdim}(F') \geq 2$

→ Segue

Seque



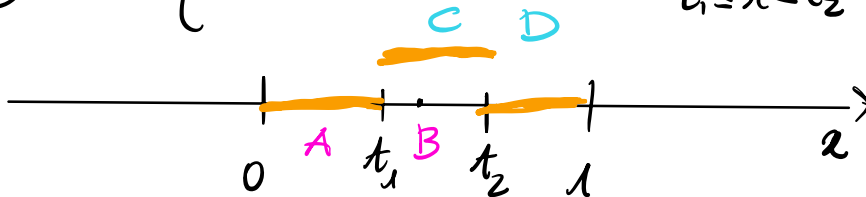
④

Quali terni di $\{0,1\}^3$ possiamo realizzare muovendo t ? Oss. le uniche configurazioni possibili

Sono $t < A < B < C$ $A < t < B < C$
 $A < B < t < C$ $A < B < C < t$

per cui si possono produrre tutte le stringhe binarie, tranne $(0,1,0)$ e $(1,0,1)$
 quindi la dimensione $\text{Vdim}(F') = 2$.

⑤ $F = \{f: X \rightarrow \{0,1\}; f(x)_{t_1 \leq x < t_2, t_1, t_2 \in [0,1]}\}$



Chiaramente $\text{Vdim}(F) \geq 1$ e anche due punti sono shattered

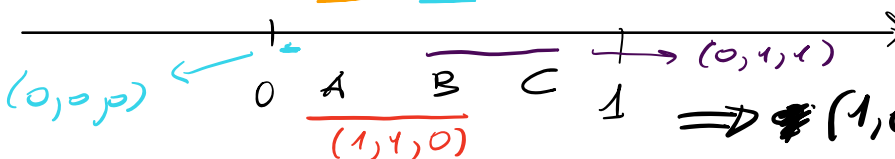
$A < B \leq t_1 < t_2 \rightarrow (1,1)$

$t_1 < t_2 < A < B \rightarrow (0,0)$

$A < t_1 < B < t_2 \Rightarrow (0,1)$

$t_1 < A < t_2 < B \Rightarrow (1,0)$

$(1,1,1) \leftarrow \text{---} \Rightarrow (1,0,0), (0,1,0), (0,0,1)$

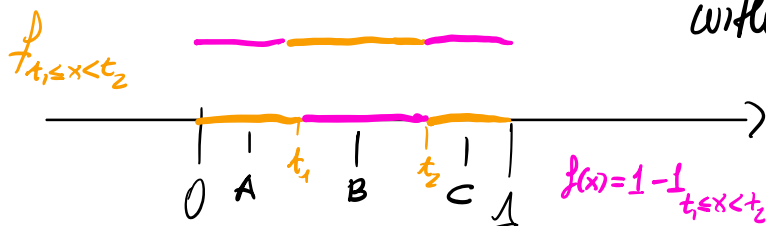


$(0,0,0) \leftarrow \text{---} \Rightarrow (1,0,1)$ non è possibile

→ quindi $\text{Vcdim}(F) = 2$

(5)

① $F' = \{f: X \rightarrow \{0,1\}; f(x) = 1_{t_1 \leq x < t_2} \text{ o } f(x) = 1 - 1_{t_1 \leq x < t_2}$
 with $t_1 < t_2 \in [0,1]$

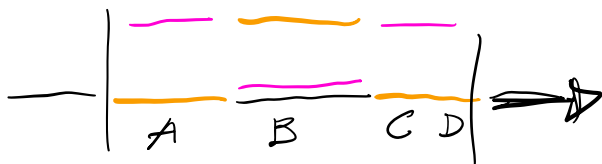


$F' \supset F \Rightarrow \text{Vcdim}(F') \geq 2$

Controlliamo 3 punti $A < B < C$, considerando che $F' \supset F$
~~con~~ guardiamo solo il caso $(1,0,1)$ che prima non
 era possibile: si vede che $f(x) = 1 - 1_{t_1 \leq x < t_2}$
 con $A < t_1 < B < t_2 < C$ produce $(1,0,1)$

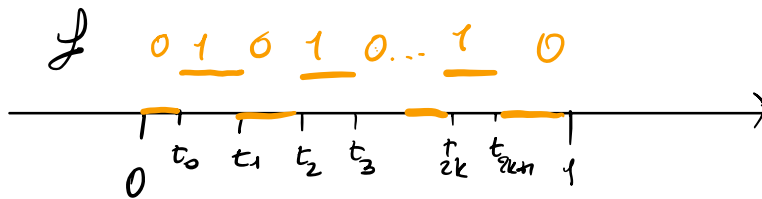
quindi $\text{Vcdim}(F') \geq 3$

Consideriamo 4 punti e vediamo che non c'è modo
 di produrre ad esempio $(1,0,1,0)$



Quindi da $\text{Vcdim}(F') = 3$

② $F_k = \{f: X \rightarrow \{0,1\} : f(x) = \sum_{i=0}^k 1_{t_{2i} \leq x < t_{2i+1}}, 0 \leq t_0 < \dots < t_{2k+1} \leq 1\}$
 - con $k \geq 1$.



Le sequenze alternate sono il problema.

...
 $(1, 0, 1, 0, \dots, 1, 0) \Rightarrow 2(k+1)$
 (stesso con $(0, 1, 0, 1, \dots, 0, 1)$ sempre $2(k+1)$)

Se prendo $2(k+1) + 1$ ottengo che
 $(1, 0, 1, 0, \dots, 1, 0, 1)$ non si può produrre, perché ha
 $k+2$ volte uguali a "1" il che non è possibile

Quindi $|K_{du}(F_k)| = 2(k+1)$

Fine