

# LECTURE 26

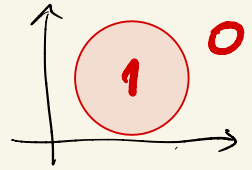
Examples continued:

5. Circles in  $\mathbb{R}^2$ :

$$\mathcal{H} = \left\{ \mathbb{1}_{\{(x(1)-a_1)^2 + (x(2)-a_2)^2 \leq r^2\}} : a_1, a_2, r \in \mathbb{R} \right\}$$

$vc(\mathcal{H}) = 3$

HW

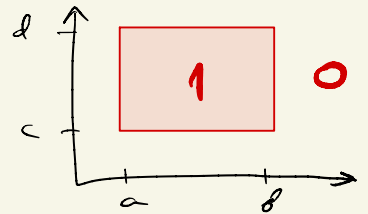


6. Axis-aligned rectangles in  $\mathbb{R}^2$ :

$$\mathcal{H} = \left\{ \mathbb{1}_{[a,b] \times [c,d]} : a < b, c < d \right\}$$

$vc(\mathcal{H}) = 4$

HW



Generalizing Example 4 from last class:

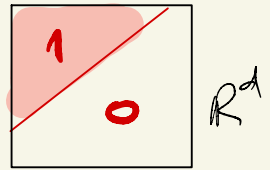
7. Half-spaces in  $\mathbb{R}^d$ :

$$\mathcal{H} = \left\{ \mathbb{1}_{\{\langle w, x \rangle + b \geq 0\}} : w \in \mathbb{R}^d, b \in \mathbb{R} \right\}$$

$vc(\mathcal{H}) = d+1$

↖ "Linear classifier" (svm)

HW



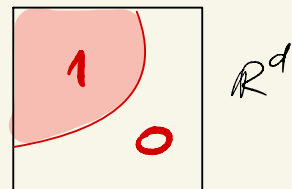
More generally:

8. Polynomial surfaces in  $\mathbb{R}^d$ :  $\mathcal{H} = \left\{ \mathbb{1}_{\{p(x) \geq 0\}} : \deg(p) \leq r \right\}$

$$vc(\mathcal{H}) = \binom{d+r}{r}$$

[Anthony 1995]

"Polynomial classifier"

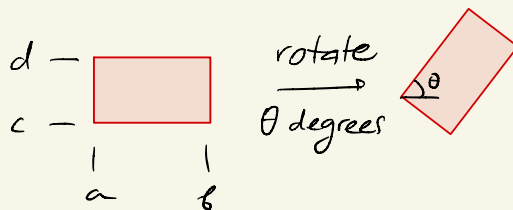


Remark. In all examples above,  $vc(\mathcal{H}) = \# \text{ parameters}$  that describe a function in  $\mathcal{H}$ .

- This is not true in general. For rectangles in  $\mathbb{R}^2$  that are not necessarily axis-aligned as in Ex. 5,

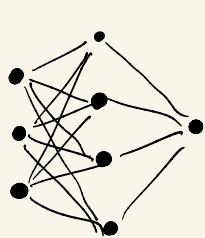
$$vc(\mathcal{H}) = 7$$

while the # parameters is 5:  
 $a, b, c, d, \theta$ .



- But heuristically, and "approximately", this is often true:

8.  $\mathcal{H} = \{ \text{functions a given neural network can compute} \}$ ,



if activation function = ,

then  $vc(\mathcal{H}) \leq C \underbrace{W}_{\# \text{ connections } (= \# \text{ weights, parameters})} \log W$  (Cover 69; Baum-Haussler 89)

- $\exists$  networks for which this bound is tight [Maas 94]

- If activation function is piecewise linear (e.g. ReLU),

then  $vc(\mathcal{H}) \leq C \underbrace{L}_{\# \text{ layers}} \log W$  [Bartlett-Karvev-Liaw-Mehradian 2017]

and  $\exists$  examples showing tightness.

Lem [Pajor 85]  $\forall$  finite class of Boolean functions  $\mathcal{H}$  on  $X$ ,  
 $|\mathcal{H}| \leq \#(\text{subsets of } X \text{ shattered by } \mathcal{H})$

Convention:  $\emptyset$  is shattered by  $\forall$  nonempty  $\mathcal{H}$ .

Proof WLOG,  $X = \{1, \dots, n\}$ . Denote by  $\text{sh}(\mathcal{H})$  the family of all subsets of  $X$  shattered by  $\mathcal{H}$ . To prove

$$|\mathcal{H}| \leq |\text{sh}(\mathcal{H})|,$$

- partition  $\mathcal{H}$  according to the value at point  $n$ , i.e.

$$\mathcal{H} = \mathcal{H}_0 \sqcup \mathcal{H}_1$$

where  $\mathcal{H}_0 = \{h \in \mathcal{H} : h(n) = 0\}$  and  $\mathcal{H}_1 = \{h \in \mathcal{H} : h(n) = 1\}$ .

- $\forall$  subset  $\{i_1, \dots, i_d\} \subset X$  shattered by  $\mathcal{H}_0$  or  $\mathcal{H}_1$  is also shattered by  $\mathcal{H}$ . Thus

$$|\text{sh}(\mathcal{H})| \geq |\text{sh}(\mathcal{H}_0)| + |\text{sh}(\mathcal{H}_1)| \quad (*)$$

domain =  $\{1, \dots, n-1\}$

- Iterate: partition  $\mathcal{H}_0$  and  $\mathcal{H}_1$  according to the value  $h(n-1)$ :

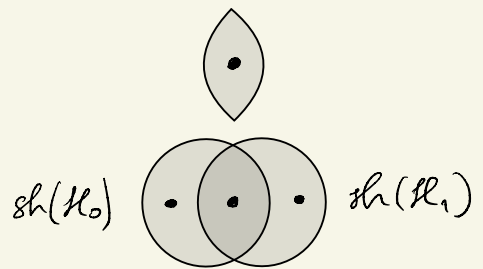
$$\geq |\text{sh}(\mathcal{H}_{00})| + |\text{sh}(\mathcal{H}_{01})| + |\text{sh}(\mathcal{H}_{10})| + |\text{sh}(\mathcal{H}_{11})| \geq$$

... down to single-point classes, each of which shatters one set  $\emptyset \dots \geq |\mathcal{H}|$

• **MISTAKE**: we double counted in (\*)

the sets that are shattered

by both  $\mathcal{H}_0$  and  $\mathcal{H}_1$



• **FIX**: suppose  $\{i_1, \dots, i_d\}$  is shattered by both  $\mathcal{H}_0$  and  $\mathcal{H}_1 \Rightarrow$

$\forall$  label assignment  $y_1, \dots, y_d \in \{0, 1\}$

$\exists h \in \mathcal{H}_0: h(i_1) = y_{i_1}, \dots, h(i_d) = y_{i_d}, h(n) = 0$

$\exists g \in \mathcal{H}_1: g(i_1) = y_{i_1}, \dots, g(i_d) = y_{i_d}, g(n) = 1$

$\Rightarrow \underbrace{\{i_1, \dots, i_d, n\}}_{\uparrow}$  is shattered by  $\mathcal{H} = \mathcal{H}_0 \cup \mathcal{H}_1$  (choose either  $h$  or  $g$ )

This set is NOT shattered by either  $\mathcal{H}_0$  or  $\mathcal{H}_1 \Rightarrow$  it was NOT counted before

•  $\Rightarrow \forall$  set that we double counted, we find a set we never counted

$\Rightarrow$  (\*) is true. Proceed as before. QED



- By def of vc dimension,  $\forall$  subset shattered by  $\mathcal{H}$  has cardinality  $\leq \text{vc}(\mathcal{H}) =: d$ . So Pajor's lemma yields  $|\mathcal{H}| \leq \#(\text{subsets of } \{1, \dots, n\} \text{ with card. } \leq d) \leq \sum_{k=0}^d \binom{n}{k}$ .



Cor (Sauer-Shelah lemma) let  $\mathcal{H}$  be a class of Boolean functions on an  $n$ -point domain. Then

$$|\mathcal{H}| \leq \sum_{k=0}^d \binom{n}{k} \quad \text{where } d = \text{vc}(\mathcal{H})$$

### Examples (skip?)

(a) Integer intervals:  $\mathcal{H} = \left\{ \begin{array}{c} \text{---} \text{---} \text{---} \text{---} \text{---} \text{---} \text{---} \text{---} \text{---} \text{---} \text{---} \\ | \quad | \quad | \quad | \quad | \quad | \quad | \quad | \quad | \quad | \\ 1 \quad 2 \quad \quad a \quad \quad b \quad \quad \quad n \end{array} : 1 \leq a \leq b \leq n \right\}$

$\text{vc}(\mathcal{H}) = 2$  (as in the previous lecture for real intervals)

$$|\mathcal{H}| = 1 + n + \binom{n}{2} = \sum_{k=0}^2 \binom{n}{k} \Rightarrow \text{Pajor lemma is sharp}$$

↑ zero function
↖  $a=b$ 
↖ # pairs  $(a < b)$

(b)  $\mathcal{H} = \{\text{all functions on an } n\text{-point domain supported by } \leq d \text{ pts}\}$   
 $\text{vc}(\mathcal{H}) = d$  **(HW?)** and  $|\mathcal{H}| = \sum_{k=0}^d \binom{n}{k} \Rightarrow \text{sharp again!}$

Remarks ①  $\sum_{k=1}^d \binom{n}{k} \leq \left(\frac{en}{d}\right)^d$  (HW 3, Problem 3)  $\xrightarrow{\text{Pajor}}$

HW 13 ↑

$$d \leq \log |\mathcal{H}| \leq d \log \left(\frac{en}{d}\right), \quad \text{where } d = \text{vc}(\mathcal{H}).$$

② Heuristically,  $\log |\mathcal{H}| = \# \text{bits to specify a function in } \mathcal{H}$   
 $d = \text{vc}(\mathcal{H}) \sim \# \text{parameters that describe functions in } \mathcal{H}$

$\Rightarrow \log |\mathcal{H}| \asymp d$ , is expected.