

# Review of Lecture 6

- $m_{\mathcal{H}}(N)$  is polynomial

if  $\mathcal{H}$  has a break point  $k$

		1	2	3	4	5	6	..
	$k$							
	1	1	2	2	2	2	2	..
	2	1						
	3	1						
$N$	4	1						
	5	1						
	6	1						
	:	:						

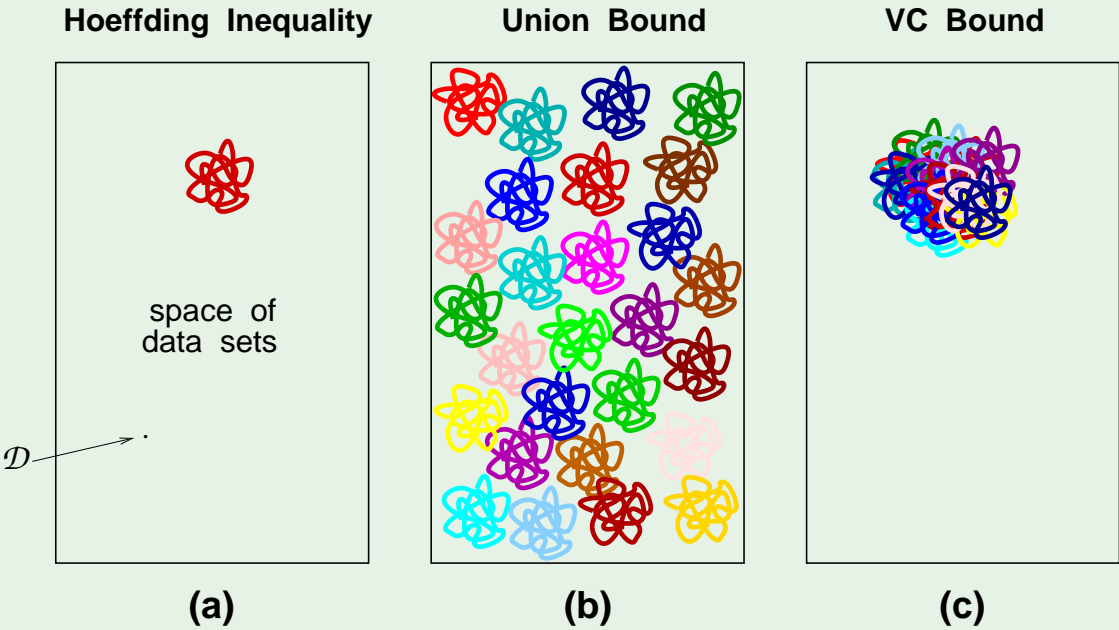
$$m_{\mathcal{H}}(N) \leq \sum_{i=0}^{k-1} \binom{N}{i}$$

maximum power is  $N^{k-1}$

Growth function characterizes the redundancy in the bad regions that we need to understand to be able to switch from Hoeffding to VC inequality.

# The VC Inequality

Can see the redundancy resulting from the fact that different hyp have, by and large, overlapping bad regions



$$\mathbb{P} [ |E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon ] \leq 2 M e^{-2 \epsilon^2 N}$$

$\downarrow$                        $\downarrow$                        $\downarrow$   
 $\downarrow$                        $\downarrow$                        $\downarrow$

$$\mathbb{P} [ |E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon ] \leq 4 m_{\mathcal{H}}(2N) e^{-\frac{1}{8} \epsilon^2 N}$$

VC inequality characterizes generalization of learning

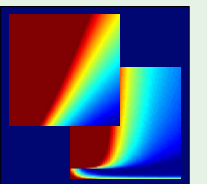
# Learning From Data

Yaser S. Abu-Mostafa  
*California Institute of Technology*

## Lecture 7: The VC Dimension



Sponsored by Caltech's Provost Office, E&AS Division, and IST • Tuesday, April 24, 2012



# Outline

- The definition
- VC dimension of perceptrons
- Interpreting the VC dimension
- Generalization bounds

# Definition of VC dimension

The VC dimension of a hypothesis set  $\mathcal{H}$ , denoted by  $d_{\text{VC}}(\mathcal{H})$ , is

the largest value of  $N$  for which  $m_{\mathcal{H}}(N) = 2^N$

“the most points  $\mathcal{H}$  can shatter”

Note that if  $d_{\text{VC}}=N$ , it does not say every  $N$  points can be shattered - only one set of  $N$  points that can be shattered is sufficient for the above statement (this has always been the case in our analysis).

$N \leq d_{\text{VC}}(\mathcal{H}) \implies \mathcal{H}$  can shatter  $N$  points

$k > d_{\text{VC}}(\mathcal{H}) \implies k$  is a break point for  $\mathcal{H}$

# The growth function

In terms of a break point  $k$ :

$$m_{\mathcal{H}}(N) \leq \sum_{i=0}^{k-1} \binom{N}{i}$$

In terms of the VC dimension  $d_{\text{VC}}$ :

$$m_{\mathcal{H}}(N) \leq \underbrace{\sum_{i=0}^{d_{\text{VC}}} \binom{N}{i}}_{\text{maximum power is } N^{d_{\text{VC}}}}$$

# Examples

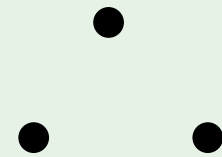
- $\mathcal{H}$  is positive rays:

$$d_{VC} = 1$$



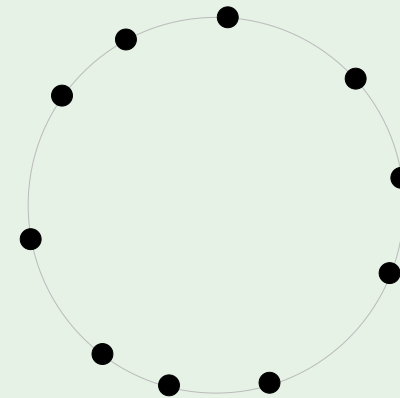
- $\mathcal{H}$  is 2D perceptrons:

$$d_{VC} = 3$$



- $\mathcal{H}$  is convex sets:

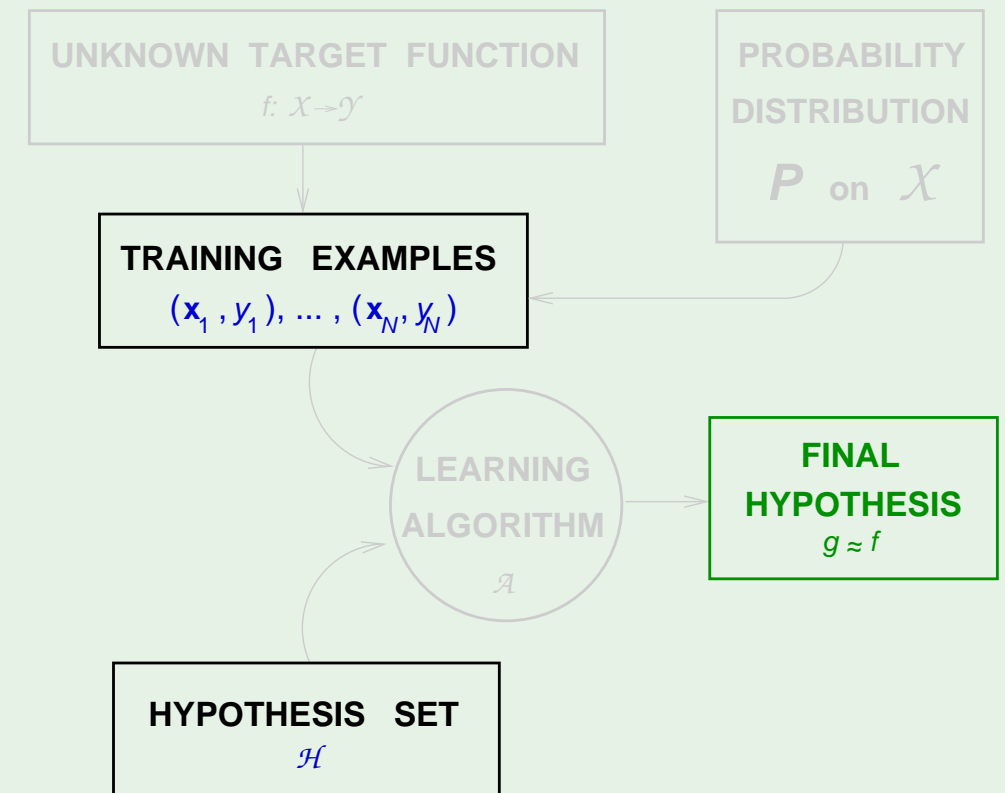
$$d_{VC} = \infty$$



# VC dimension and learning

$d_{\text{VC}}(\mathcal{H})$  is finite  $\implies g \in \mathcal{H}$  will generalize

- Independent of the **learning algorithm**
- Independent of the **input distribution**
- Independent of the **target function**



# VC dimension of perceptrons

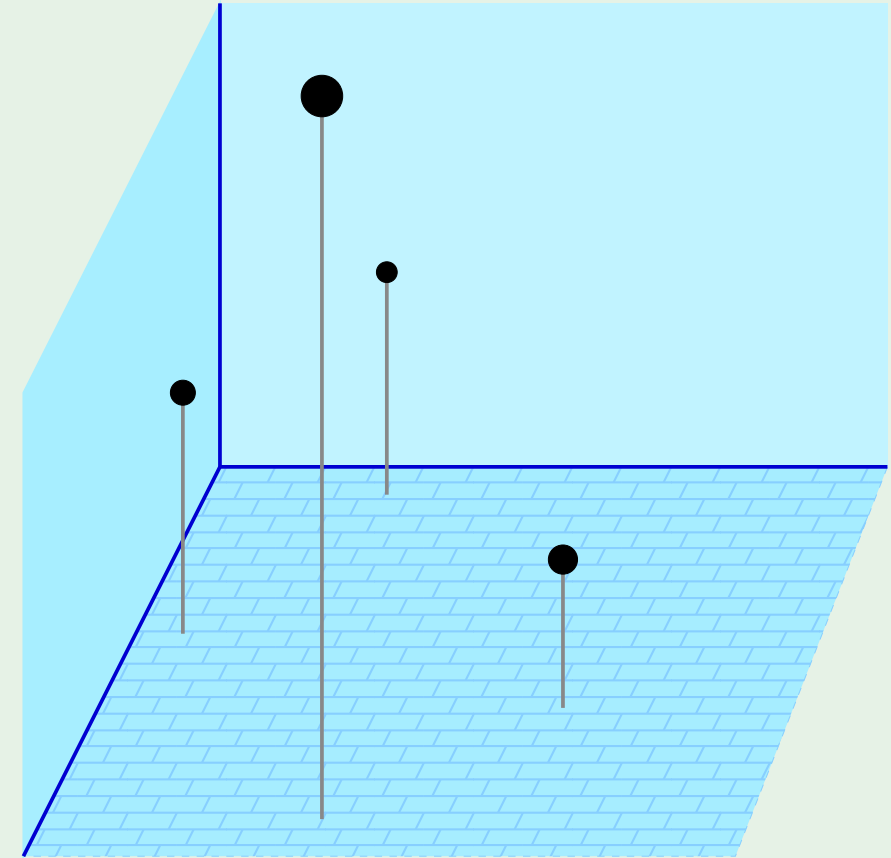
For  $d = 2$ ,  $d_{\text{VC}} = 3$

In general,  $d_{\text{VC}} = d + 1$

We will prove two directions:

$$d_{\text{VC}} \leq d + 1$$

$$d_{\text{VC}} \geq d + 1$$





Here is one direction

A set of  $N = d + 1$  points in  $\mathbb{R}^d$  shattered by the perceptron:

$$X = \begin{bmatrix} \text{---} \mathbf{x}_1^\top \text{---} \\ \text{---} \mathbf{x}_2^\top \text{---} \\ \text{---} \mathbf{x}_3^\top \text{---} \\ \vdots \\ \text{---} \mathbf{x}_{d+1}^\top \text{---} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 1 & 1 & 0 & \dots & 0 \\ 1 & 0 & 1 & & 0 \\ & \vdots & & \ddots & 0 \\ 1 & 0 & \dots & 0 & 1 \end{bmatrix}$$

$X$  is invertible

## Can we shatter this data set?

For any dichotomy that is picked from  $y$ , we want to show that we can find a perceptron that can realize this, and therefore we show that we can shatter the set.

For any  $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_{d+1} \end{bmatrix} = \begin{bmatrix} \pm 1 \\ \pm 1 \\ \vdots \\ \pm 1 \end{bmatrix}$ , can we find a vector  $\mathbf{w}$  satisfying

$$\text{sign}(\mathbf{X}\mathbf{w}) = \mathbf{y}$$

Easy! Just make  $\mathbf{X}\mathbf{w} = \mathbf{y}$

which means  $\mathbf{w} = \mathbf{X}^{-1}\mathbf{y}$

We can shatter these  $d + 1$  points

This implies what?

[a]  $d_{\text{VC}} = d + 1$

[b]  $d_{\text{VC}} \geq d + 1$  ✓

[c]  $d_{\text{VC}} \leq d + 1$

[d] No conclusion

Now, to show that  $d_{\text{vc}} \leq d + 1$

We need to show that:

- [a] There are  $d + 1$  points we cannot shatter
- [b] There are  $d + 2$  points we cannot shatter
- [c] We cannot shatter *any* set of  $d + 1$  points
- [d] We cannot shatter *any* set of  $d + 2$  points ✓

Take any  $d + 2$  points

For any  $d + 2$  points,

$$\mathbf{x}_1, \dots, \mathbf{x}_{d+1}, \mathbf{x}_{d+2}$$

More points than dimensions  $\implies$  we must have

$$\mathbf{x}_j = \sum_{i \neq j} a_i \mathbf{x}_i$$

if there are more vectors than dimensions, they must be linearly dependant

where not all the  $a_i$ 's are zeros

since  $x_0 = 1$  for each point due to the bias/threshold term, so not all  $a_i$ 's can be zero

So?

$$\mathbf{x}_j = \sum_{i \neq j} \mathbf{a}_i \mathbf{x}_i$$

Consider the following dichotomy:

$\mathbf{x}_i$ 's with non-zero  $\mathbf{a}_i$  get  $y_i = \text{sign}(\mathbf{a}_i)$

$\mathbf{x}_i$ 's with zero  $\mathbf{a}_i$  get  $\pm 1$ , we will not consider them

and  $\mathbf{x}_j$  gets  $y_j = -1$

No perceptron can implement such dichotomy!

# Why?

$$\mathbf{x}_j = \sum_{i \neq j} a_i \mathbf{x}_i \implies \mathbf{w}^\top \mathbf{x}_j = \sum_{i \neq j} a_i \mathbf{w}^\top \mathbf{x}_i \quad (1)$$

we asserted this dichotomy in the last slide

If  $y_i = \text{sign}(\mathbf{w}^\top \mathbf{x}_i) = \text{sign}(a_i)$ , then  $a_i \mathbf{w}^\top \mathbf{x}_i > 0$  (2) since  $a_i$  and  $\mathbf{w}^\top \mathbf{x}_i$  must have the same sign

This forces

$$\mathbf{w}^\top \mathbf{x}_j = \sum_{i \neq j} a_i \mathbf{w}^\top \mathbf{x}_i > 0 \quad \text{combining (1) and (2)}$$

Therefore,  $y_j = \text{sign}(\mathbf{w}^\top \mathbf{x}_j) = +1$

which contradicts the dichotomy - therefore we cannot shatter (for any set we choose) a set of  $d+2$  points with a  $d$ -dimensional perceptron

## Putting it together

We proved  $d_{\text{VC}} \leq d + 1$  and  $d_{\text{VC}} \geq d + 1$

$$d_{\text{VC}} = d + 1$$

What is  $d + 1$  in the perceptron?

It is the number of parameters  $w_0, w_1, \dots, w_d$



# Outline

- The definition
- VC dimension of perceptrons
- Interpreting the VC dimension
- Generalization bounds

# 1. Degrees of freedom

Parameters create degrees of freedom

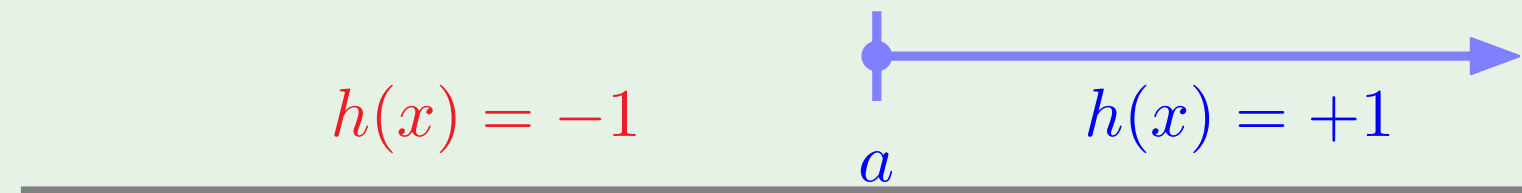
# of parameters: **analog** degrees of freedom

$d_{VC}$ : equivalent '**binary**' degrees of freedom

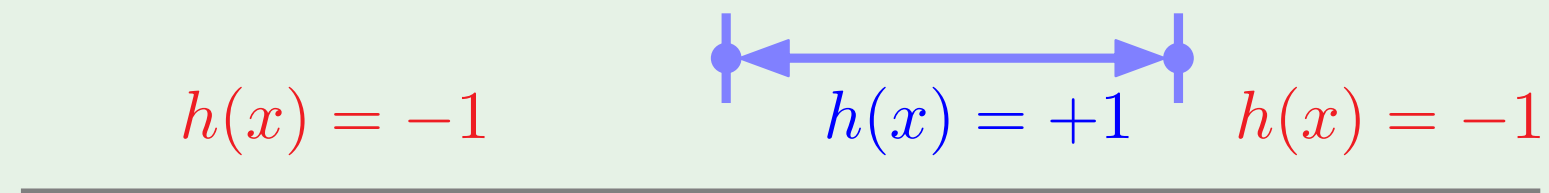


# The usual suspects

Positive rays ( $d_{VC} = 1$ ):

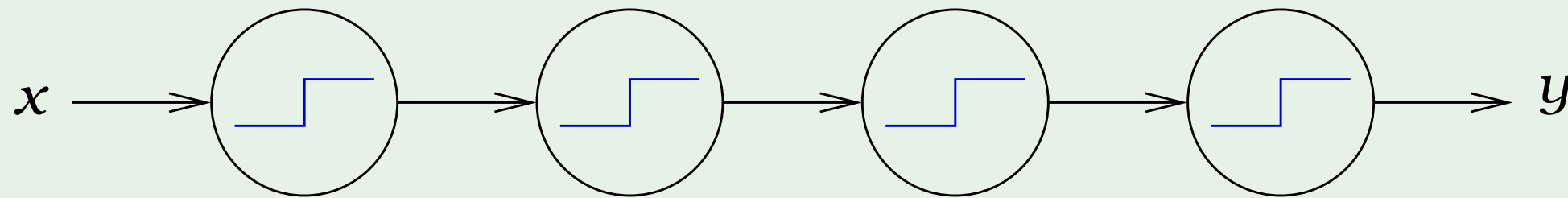


Positive intervals ( $d_{VC} = 2$ ):



# Not just parameters

Parameters may not contribute degrees of freedom:



perceptrons 2,3,4 are redundant, so only 2 effective parameters rather than 8

$d_{VC}$  measures the **effective** number of parameters

## 2. Number of data points needed

Two small quantities in the VC inequality:

$$\mathbb{P}[|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon] \leq \underbrace{4m_{\mathcal{H}}(2N)}_{\delta} e^{-\frac{1}{8}\epsilon^2 N}$$

If we want certain  $\epsilon$  and  $\delta$ , how does  $N$  depend on  $d_{\text{VC}}$ ?

Let us look at

$$N^d e^{-N}$$

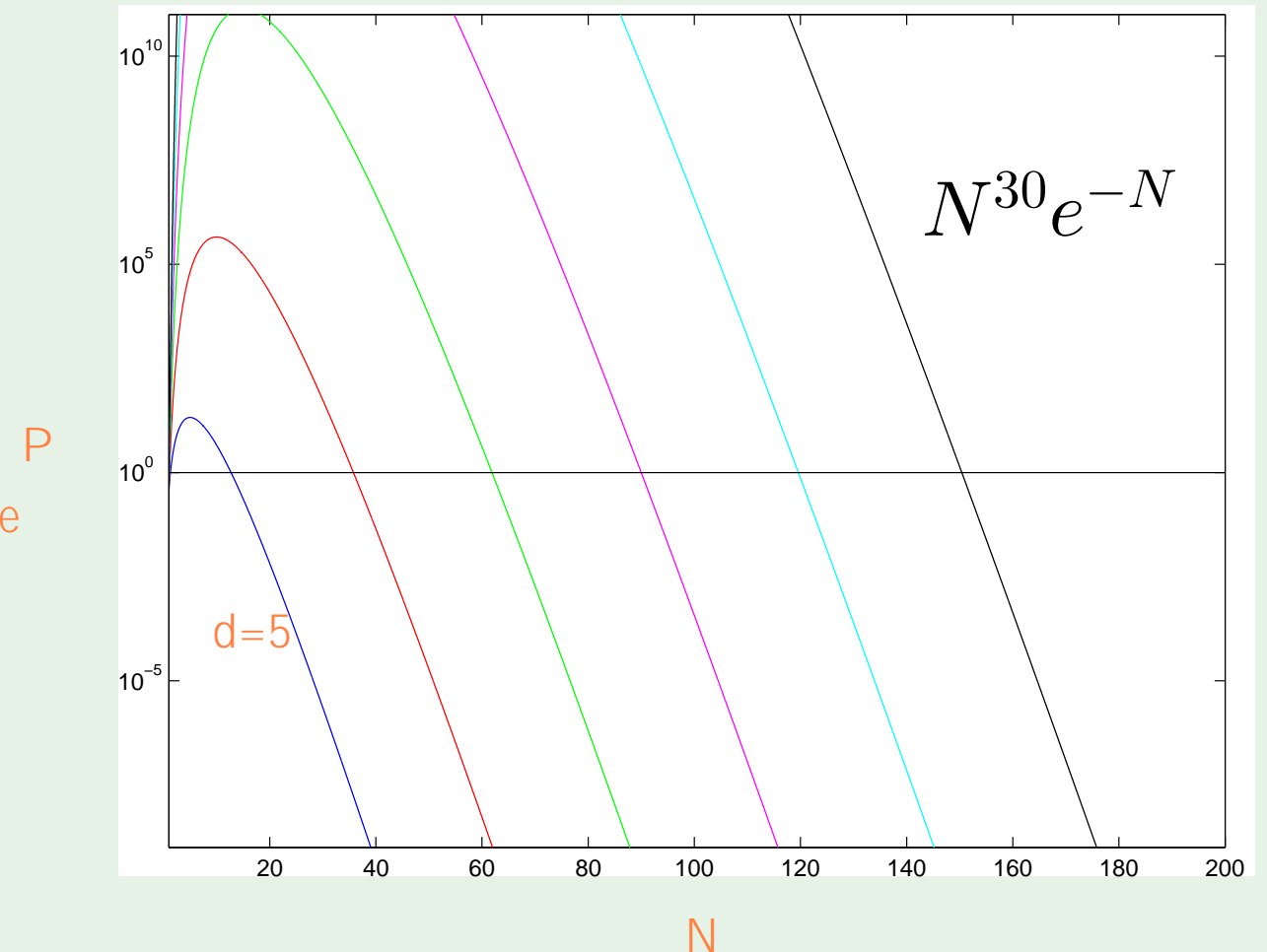
$$N^d e^{-N}$$

Fix  $N^d e^{-N} = \text{small value}$

How does  $N$  change with  $d$ ?

for a large range of delta and epsilon, and large range  
**Rule of thumb:** of practical applications

$$N \geq 10 d_{\text{VC}}$$



Practical observation: the actual quantity we are trying to bound follows the same monotonicity as the bound.- so in using bigger d<sub>VC</sub>, the quantities you get are bigger to achieve a certain level of performance and actually close to proportional. In spite of the fact we cannot get an exact value due to the bound, the relative aspect of the VC dimension holds - so a larger VC dimension requires more examples.

# Outline

- The definition
- VC dimension of perceptrons
- Interpreting the VC dimension
- Generalization bounds

# Rearranging things

Start from the VC inequality:

$$\mathbb{P}[|E_{\text{out}} - E_{\text{in}}| > \epsilon] \leq \underbrace{4m_{\mathcal{H}}(2N)e^{-\frac{1}{8}\epsilon^2 N}}_{\delta}$$

Get  $\epsilon$  in terms of  $\delta$ :

$$\delta = 4m_{\mathcal{H}}(2N)e^{-\frac{1}{8}\epsilon^2 N} \implies \epsilon = \underbrace{\sqrt{\frac{8}{N} \ln \frac{4m_{\mathcal{H}}(2N)}{\delta}}}_{\Omega}$$

so bigger dVC, larger omega  
and therefore larger epsilon, so  
worse generalization -  
conversely more data we have  
a smaller omega

good event

With probability  $\geq 1 - \delta$ ,  $|E_{\text{out}} - E_{\text{in}}| \leq \Omega(N, \mathcal{H}, \delta)$



# Generalization bound

generalization error

With probability  $\geq 1 - \delta$ ,

$$E_{\text{out}} - E_{\text{in}} \leq \Omega$$

$E_{\text{in}}$  is generally smaller than  $E_{\text{out}}$  as that is the quantity we are deliberately minimizing

$\Rightarrow$

With probability  $\geq 1 - \delta$ ,

$$E_{\text{out}} \leq E_{\text{in}} + \Omega$$