

# Review of Lecture 6

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

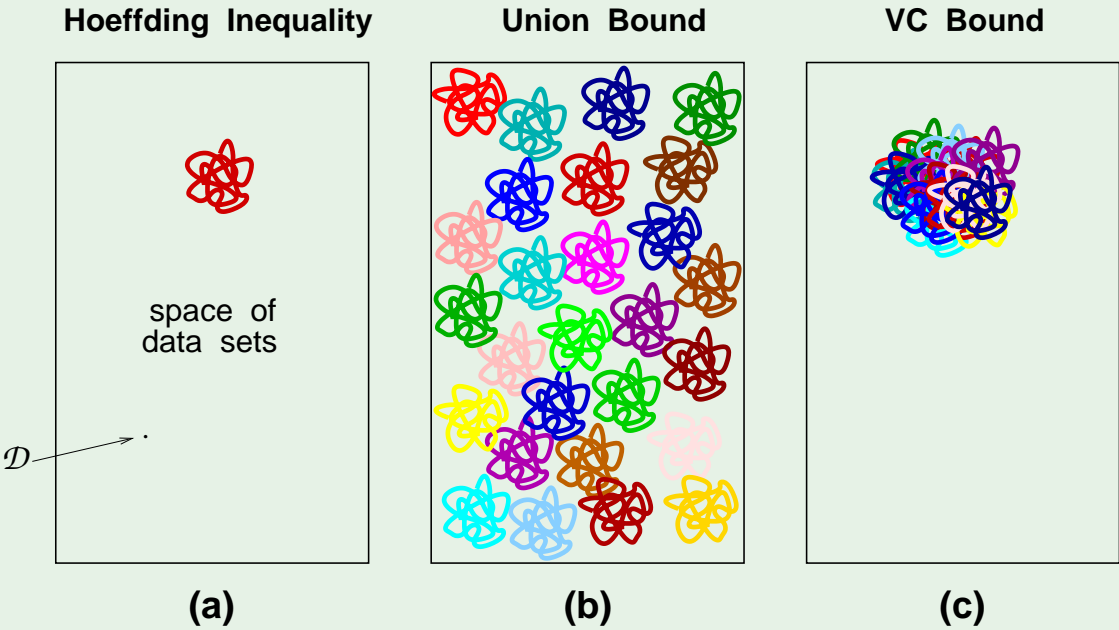
if  $\mathcal{H}$  has a break point  $k$

		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  $N^{k-1}$

- The VC Inequality



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

$\downarrow$   
 $\downarrow$

$\downarrow$   
 $\downarrow$

$\downarrow$   
 $\downarrow$

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

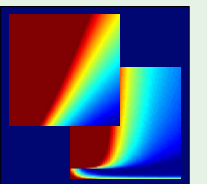
# Learning From Data

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”

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

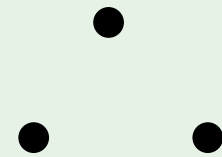
- $\mathcal{H}$  is positive rays:

$$d_{VC} = 1$$



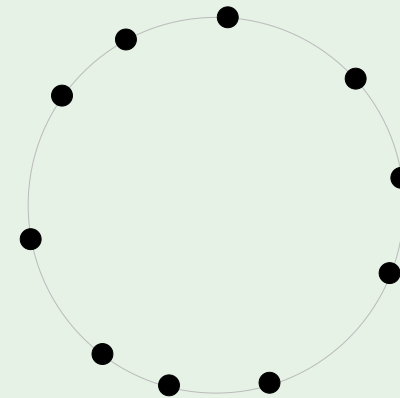
- $\mathcal{H}$  is 2D perceptrons:

$$d_{VC} = 3$$



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

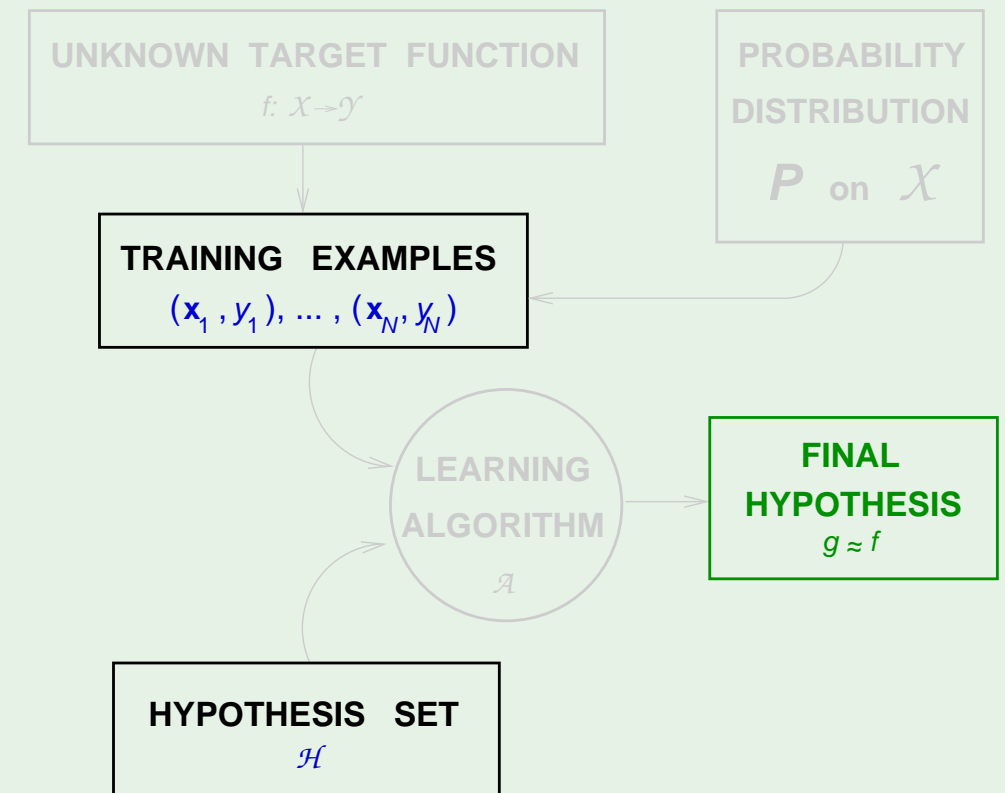
$$d_{VC} = \infty$$



# 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

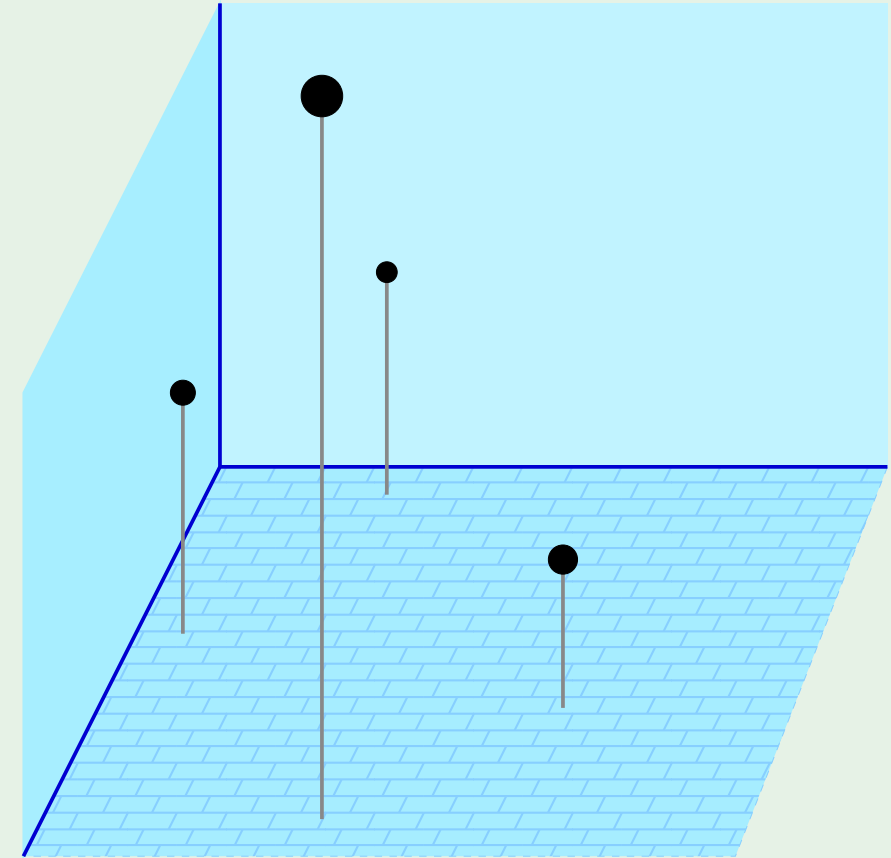
For  $d = 2$ ,  $d_{\text{VC}} = 3$

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

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

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

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

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

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

where not all the  $a_i$ 's are zeros

So?

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

Consider the following dichotomy:

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

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

No perceptron can implement such dichotomy!

Why?

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

If  $y_i = \text{sign}(\mathbf{w}^\top \mathbf{x}_i) = \text{sign}(a_i)$ , then  $a_i \mathbf{w}^\top \mathbf{x}_i > 0$

This forces

$$\mathbf{w}^\top \mathbf{x}_j = \sum_{i \neq j} a_i \mathbf{w}^\top \mathbf{x}_i > 0$$

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

## 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

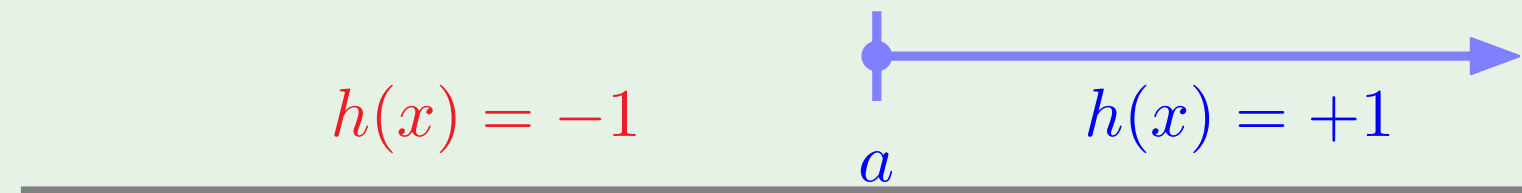
# of parameters: **analog** degrees of freedom

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

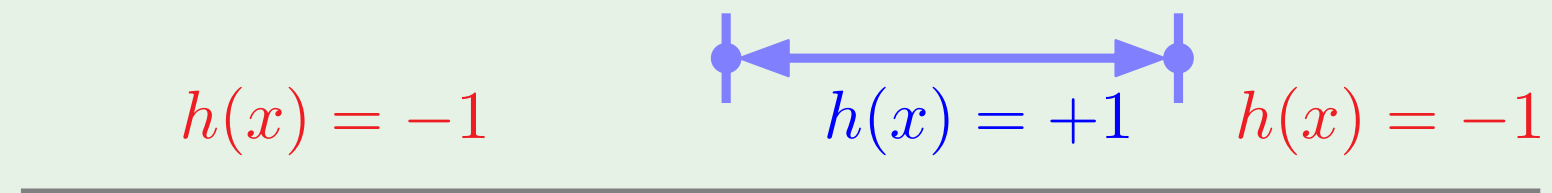


# The usual suspects

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

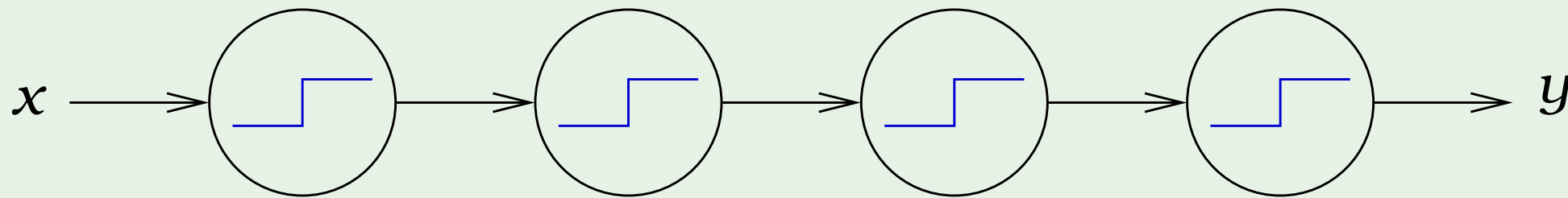


Positive intervals ( $d_{VC} = 2$ ):



## Not just parameters

Parameters may not contribute degrees of freedom:



$d_{VC}$  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

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

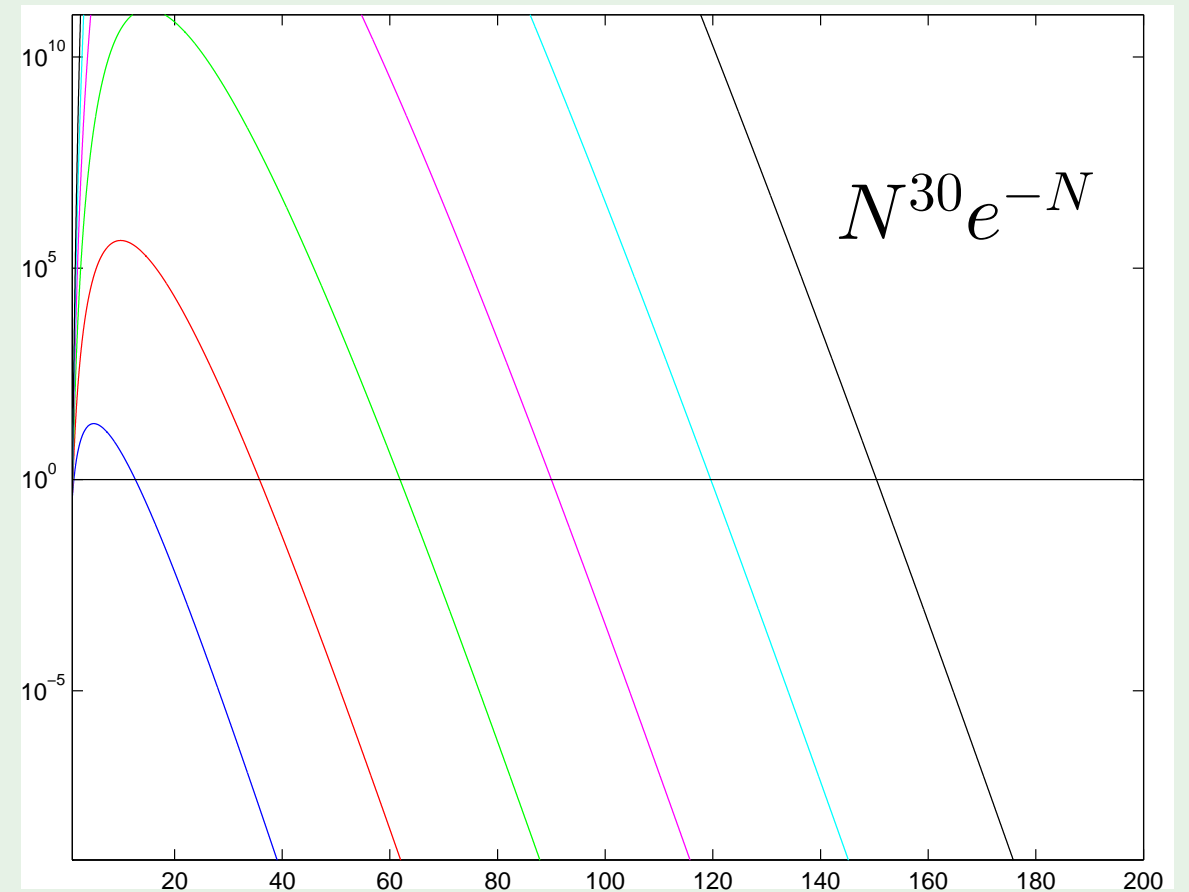
$$N^d e^{-N}$$

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

How does  $N$  change with  $d$ ?

Rule of thumb:

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



# 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  $\delta$ :

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

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



# Generalization bound

With probability  $\geq 1 - \delta$ ,  $E_{\text{out}} - E_{\text{in}} \leq \Omega$

$\Rightarrow$

With probability  $\geq 1 - \delta$ ,

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