

MAC 0460 / 5832

Introduction to Machine Learning

10 – Is learning feasible ?

• VC dimension •

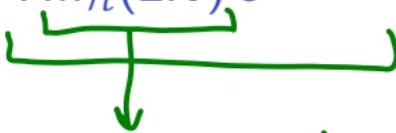
IME/USP (17/05/2021)

Hoeffding inequality

$$P\left(\left|E_{in}(g) - E_{out}(g)\right| > \epsilon\right) \leq 2\underbrace{M e^{-2\epsilon^2 N}}$$

VC inequality

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


Growth function

Hypothesis space: \mathcal{H}_{II}

Growth-function: $m_{\mathcal{H}}(N)$ (counts max number of dichotomies)

Break point: k is a break point for \mathcal{H} if there is no dataset of size k for which \mathcal{H} generates all 2^k dichotomies

$m_{\mathcal{H}}(N)$ is polynomial if there is a break-point

The bound $4 m_{\mathcal{H}}(2N) e^{-\frac{1}{8}\epsilon^2 N}$ in the VC inequality tends to zero as N increases (The negative exponential starts to dominate the polynomial at some point)



VC dimension $d_{\text{vc}}(\mathcal{H})$:

The largest number of points that can be shattered by \mathcal{H}
(The largest value of N for which $m_{\mathcal{H}}(N) = 2^N$)

Break point:

k is a break point for \mathcal{H} if there is no dataset of size k shattered by \mathcal{H}

If k is a break point for \mathcal{H} , then $d_{\text{vc}}(\mathcal{H}) < k$

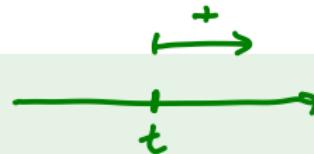
$d_{\text{vc}}(\mathcal{H}) + 1$ is a *break-point* for \mathcal{H}

[]
↑ ↑

Examples

- \mathcal{H} is positive rays:

$$d_{VC} = 1$$



- \mathcal{H} is 2D perceptrons:

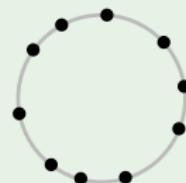
$$d_{VC} = 3$$

- $K=2$
• •

- \mathcal{H} is convex sets:

$$d_{VC} = \infty$$

- •
• • $K=4$ is break-point



Example: VC dimension of the perceptron

Let d be the dimension of input data ($x = (x_1, x_2, \dots, x_d)$)

For perceptrons, $d_{vc} = d + 1$

To prove it, it is enough to show that

- { (a) $d_{vc} \geq d + 1$, and ↗
- (b) $d_{vc} \leq d + 1$ ↗ .

2D
 $K=4$ is break-point
 $K=3 \Rightarrow 2^3$ partitions
 $d_{vc} \leq 3$ ($d=2$)



What do we need to do to prove (a) $d_{vc} \geq d + 1$?



What do we need to do to prove (a) $d_{vc} \geq d + 1$?

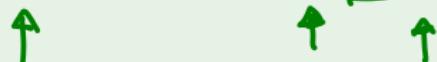
A. We need to show that there is a set of $d + 1$ points that can be shattered by the perceptron

How? Carefully choose $d + 1$ points, assign arbitrary labels in $\{-1, +1\}$ for each of them, and then show that there is a hypothesis that agrees with the labels

Here is one direction

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

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

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

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



What do we need to do to prove (b) $\underline{d_{VC} \leq d + 1}$?

- 1) Mostrar que $K = d + 2$ é break-point \leftarrow
- 2) provar que algum conjunto com $d + 2$
~~é~~ pode ser shattord .



What do we need to do to prove (b) $d_{vc} \leq d + 1$?

A. We need to show that no set of $d + 2$ points can be shattered
by the perceptron

How? Take any set of $d + 2$ points and show that it is always possible to build a dichotomy that can not be generated by any of the hypotheses

Take any $d+2$ points

For any $d+2$ points,

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

More points than dimensions \implies we must have

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


where not all the a_i 's are zeros



So?

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

Consider the following dichotomy:

$$\underbrace{\mathbf{x}_i \text{'s with non-zero } \mathbf{a}_i}_{\leftarrow} \text{ get } \mathbf{y}_i = \text{sign}(\mathbf{a}_i) \quad \leftarrow$$

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

No perceptron can implement such dichotomy!

Why?

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

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

This forces

$$\mathbf{w}^\top \mathbf{x}_j = \underbrace{\sum_{i \neq j} a_i}_{\text{---}} \mathbf{w}^\top \mathbf{x}_i > 0 \quad \leftarrow$$

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

Putting it together

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

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

What is $d + 1$ in the perceptron?

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

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_{VC} :

$$d_{VC} < k$$

$$m_{\mathcal{H}}(N) \leq \sum_{i=0}^{d_{VC}} \binom{N}{i}$$

maximum power is $N^{d_{VC}}$

$k = d_{VC} + 1$ is a break-point

$$m_{\mathcal{H}}(N) \leq \sum_{i=0}^{d_{VC}} \binom{N}{i} \leq N^{d_{VC}} + 1$$



(Proof by induction; Problem 2.5 of the book)

VC inequality

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

$$\underline{m_{\mathcal{H}}(2N)} \leq (2N)^{d_{VC}} + 1$$

VC bound

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

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

0.05

If $P(|a - b| > \epsilon) \leq \delta$, then with probability $1 - \delta$ we have $|a - b| \leq \epsilon$, i.e.,
 $b - a \leq \epsilon \leq a - b$

Letting $a = E_{in}$ and $b = E_{out}$, with probability $1 - \delta$ we have:

$$E_{out} \leq E_{in} + \sqrt{\frac{8}{N} \ln \frac{4m_{\mathcal{H}}(2N)}{\delta}}$$

Rearranging things

bad event

Start from the VC inequality:

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

$\delta = 0.05$

Get ϵ in terms of δ :

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

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

Generalization bound

With probability $\geq 1 - \delta$,



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



$$\delta \rightarrow 5\%$$

With probability $\geq 1 - \delta$,

$$1 - \delta = 95\% \longrightarrow$$

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



Remarks

- Generalization error:

- in this course it refers to $|E_{in}(h) - E_{out}(h)|$ ←
- However, often it is used to name $E_{out}(h)$

- VC generalization bound: with probability $1 - \delta$ ($\delta > 0$)

$$E_{out} \leq E_{in} + \underbrace{\sqrt{\frac{8}{N} \ln \frac{4m_{\mathcal{H}}(2N)}{\delta}}}_{\Omega}$$

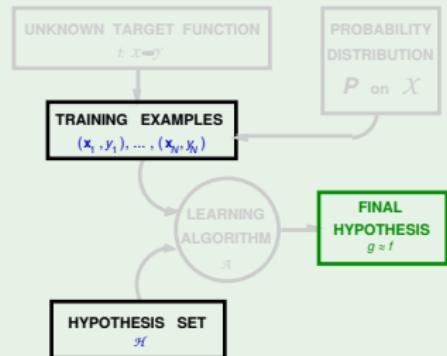
Some discussions

- Intuitive meaning of VC dimension
- Sample complexity

VC dimension and learning

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

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



1. Degrees of freedom

Parameters create degrees of freedom

of parameters: **analog** degrees of freedom

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



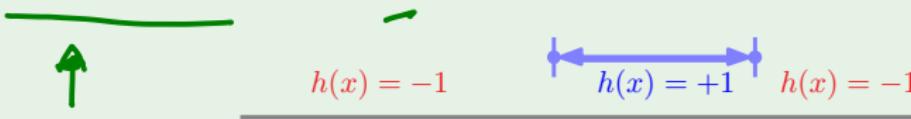
The usual suspects



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

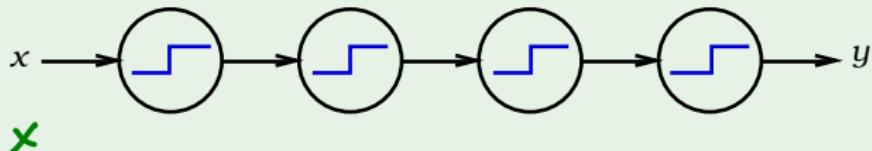


Positive intervals ($d_{VC} = 2$):



Not just parameters

Parameters may not contribute degrees of freedom:



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

Sample complexity

If d_{VC} is finite, learning generalizes

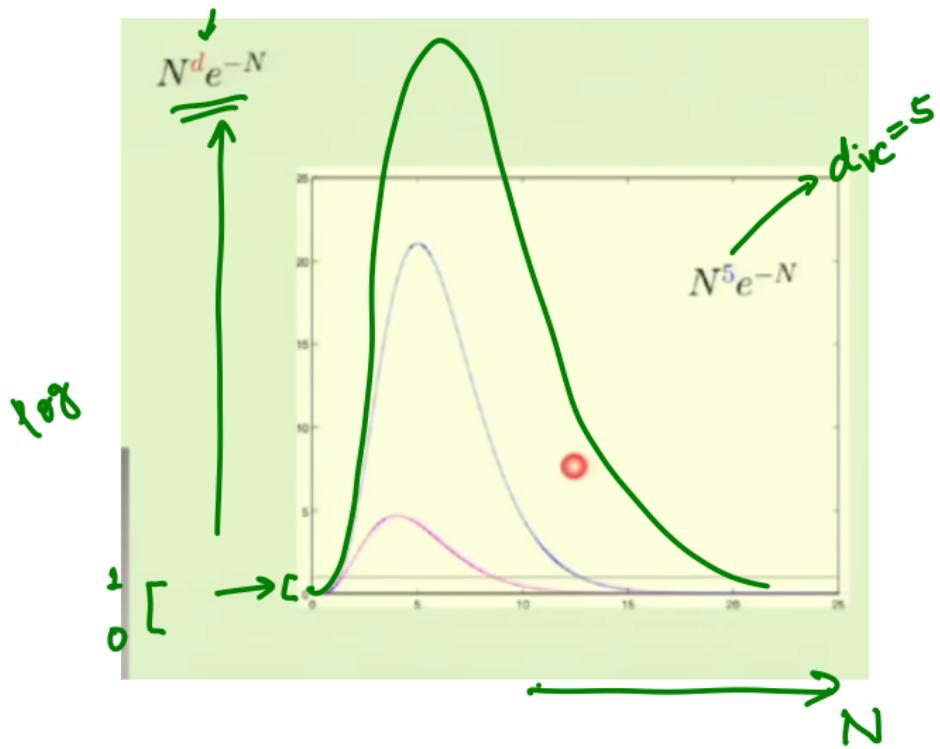
How many examples do we need ?

Let us examine the behavior of a rough approximation for the bound:

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

$$N^{d_{VC}} e^{-N}$$

(Recall that $m_{\mathcal{H}}(N) \leq N^{d_{VC}} + 1$)

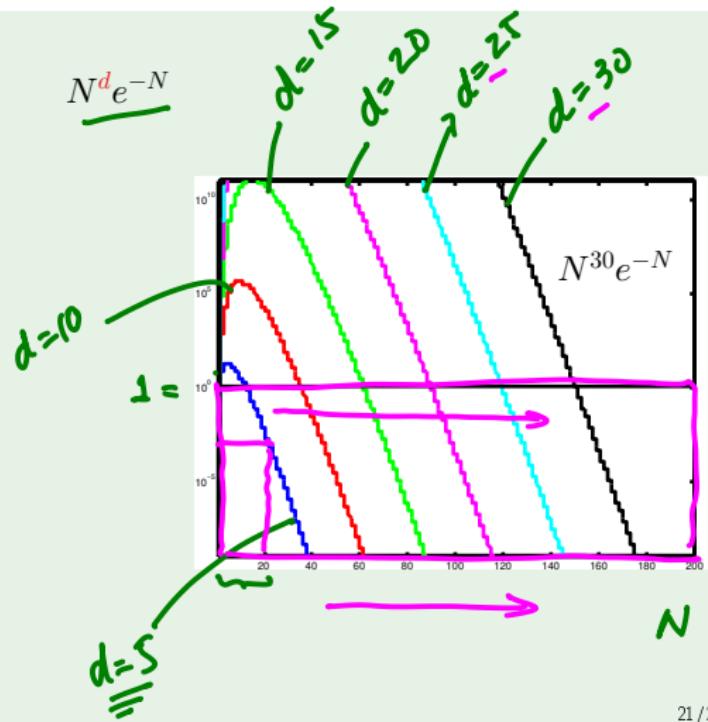


Fix $N^d e^{-N}$ = small value

How does N change with d ?

Rule of thumb:

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



1. Dichotomies are the key for the definition of VC dimension

2. The VC dimension replaces M (size of \mathcal{H}) in the Hoeffding inequality bound

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

$(m_{\mathcal{H}}(2N) \leq (2N)^{d_{VC}} + 1)$

3. VC dimension is related to the expressiveness of \mathcal{H}

$$E_{out} \leq E_{in} + \sqrt{\frac{8}{N} \ln \frac{4m_{\mathcal{H}}(2N)}{\delta}}$$

Ω

d_{VC}	E_{in}	Ω
small	large	small
large	small	large

