

CSE4088 Introduction to Machine Learning

The VC dimension

Slides are adopted from lecture notes of Yaser Abu-Mostafa

Outline

- The definition of VC dimension
- 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”

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

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

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_{VC} :

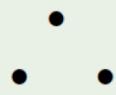
$$m_{\mathcal{H}}(N) \leq \underbrace{\sum_{i=0}^{d_{\text{VC}}} \binom{N}{i}}_{\text{maximum power is } N^{d_{\text{VC}}}}$$

Examples

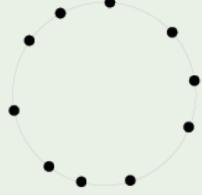
- \mathcal{H} is positive rays:

$$d_{VC} = 1$$


- \mathcal{H} is 2D perceptrons:

$$d_{VC} = 3$$


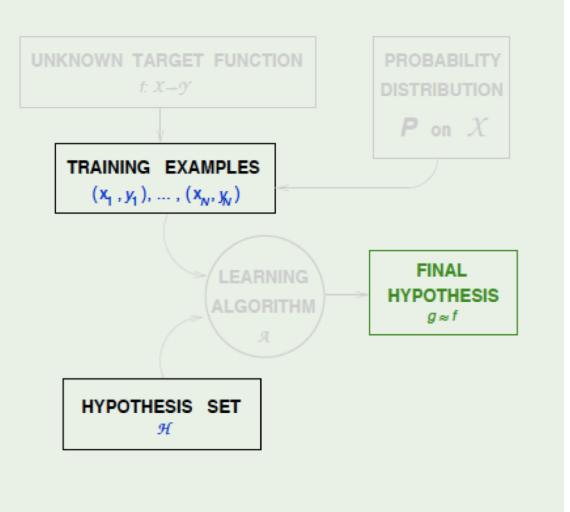
- \mathcal{H} is convex sets:

$$d_{VC} = \infty$$


VC dimension and learning

$d_{VC}(\mathcal{H})$ 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

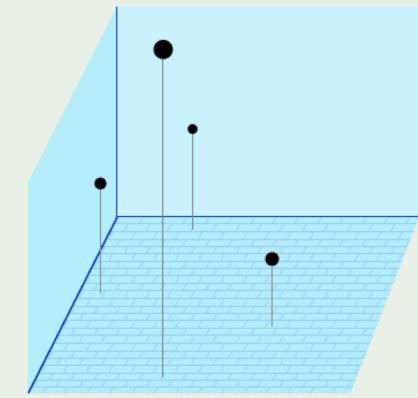
For $d = 2$, $d_{\text{VC}} = 3$

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:

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

\mathbf{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] $d_{VC} = d + 1$
- [b] $d_{VC} \geq d + 1$ ✓
- [c] $d_{VC} \leq d + 1$
- [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

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

where not all the \mathbf{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 $y_j = -1$

This is a dichotomy

No perceptron can implement such dichotomy!

Why?

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

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

This forces

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

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

Putting it together

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

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

What is $d + 1$ in the perceptron?

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

Outline

- The definition of VC dimension
- VC dimension of perceptrons
- Interpreting the VC dimension
 - 1. What does the VC dimension signify?
 - 2. Is VC dimension useful in practice?
- Generalization bounds

1. Degrees of freedom

Parameters create degrees of freedom

of parameters: **analog** degrees of freedom

d_{VC} : equivalent ‘**binary**’ degrees of freedom

- VC dimension represents how many different dichotomies we can get.
- VC dimension abstracts how expressive the model is.

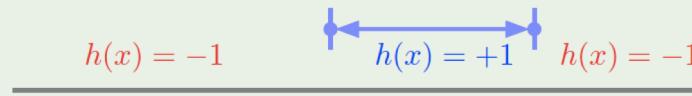


The usual suspects

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

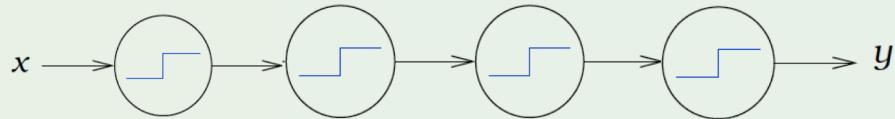


Positive intervals ($d_{VC} = 2$):



Not just parameters

Parameters may not contribute degrees of freedom:



- 1-D perceptron has 2 parameters: the weight and the threshold.
- If we cascade 4 perceptrons, we have $2 \times 4 = 8$ parameters
- Does it have 8 degrees of freedom?
- NO!: It is equivalent to a single perceptron.

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

2. Number of data points needed

Two small quantities in the VC inequality:

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

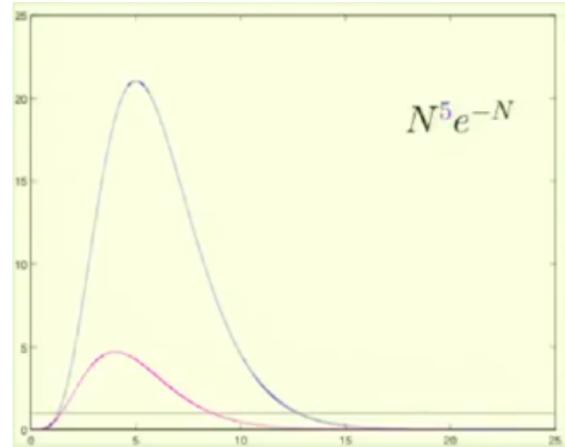
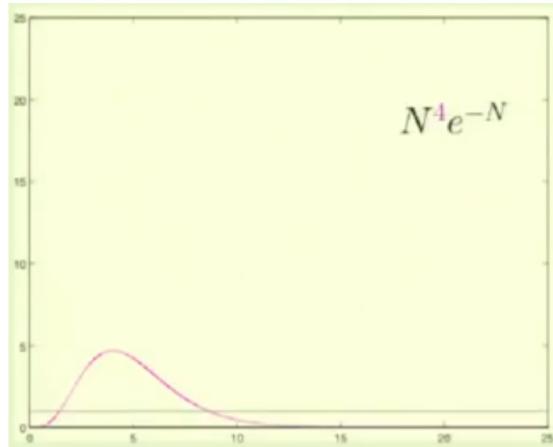
If we want certain ϵ and δ , how does N depend on d_{VC} ?

E.g. If we want E_{in} and E_{out} to be within 10% of each other 95% of the time, how many examples do we need?
($\epsilon = 0.1, \delta = 0.05$)

Let us look at

$$N^{\frac{d}{2}} e^{-N}$$

$$N^{\textcolor{red}{d}} e^{-N}$$



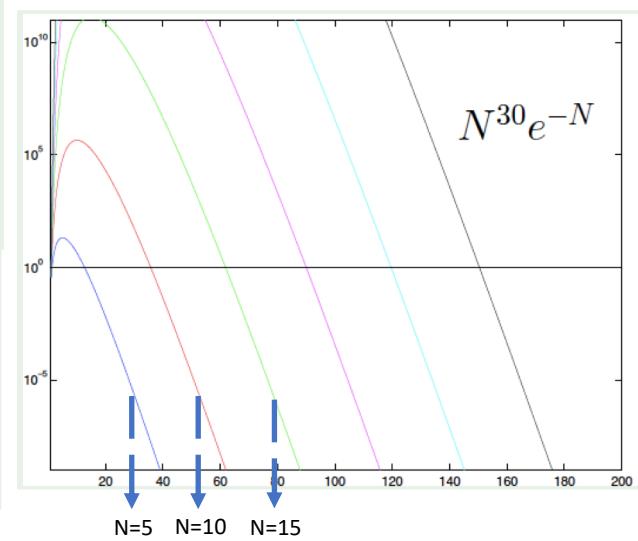
$$N^{\textcolor{red}{d}} e^{-N}$$

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

How does N change with $\textcolor{red}{d}$?

Rule of thumb:

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



Outline

- The definition of VC dimension
- 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 ϵ in terms of δ :

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

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

Generalization bound

With probability $\geq 1 - \delta$,

$$\underbrace{E_{\text{out}} - E_{\text{in}}}_{\substack{\text{Generalization error} \\ \implies}} \leq \Omega$$

With probability $\geq 1 - \delta$,

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