

# Machine Learning Foundations (機器學習基石)



Lecture 7: The VC Dimension

Hsuan-Tien Lin (林軒田)

htlin@csie.ntu.edu.tw

Department of Computer Science  
& Information Engineering

National Taiwan University  
(國立台灣大學資訊工程系)



The VC Dimension

## Roadmap

- ① When Can Machines Learn?
- ② Why Can Machines Learn?

### Lecture 6: Theory of Generalization

$E_{\text{out}} \approx E_{\text{in}}$  possible  
if  $m_{\mathcal{H}}(N)$  breaks somewhere and  $N$  large enough

### Lecture 7: The VC Dimension

- Definition of VC Dimension
- VC Dimension of Perceptrons
- Physical Intuition of VC Dimension
- Interpreting VC Dimension

- ③ How Can Machines Learn?
- ④ How Can Machines Learn Better?

## Recap: More on Growth Function

$$m_{\mathcal{H}}(N) \text{ of break point } k \leq B(N, k) = \underbrace{\sum_{i=0}^{k-1} \binom{N}{i}}_{\text{highest term } N^{k-1}}$$

	$k$				
$B(N, k)$	1	2	3	4	5
$N$	1	1	2	2	2
	2	1	3	4	4
	3	1	4	7	8
	4	1	5	11	15
	5	1	6	16	26
6	1	7	22	42	57

	$k$				
$N^{k-1}$	1	2	3	4	5
1	1	1	1	1	1
2	1	2	4	8	16
3	1	3	9	27	81
4	1	4	16	64	256
5	1	5	25	125	625
6	1	6	36	216	1296

provably & loosely, for  $N \geq 2, k \geq 3$ ,

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

## Recap: More on Vapnik-Chervonenkis (VC) Bound

For any  $g = \mathcal{A}(\mathcal{D}) \in \mathcal{H}$  and 'statistical' large  $\mathcal{D}$ , for  $N \geq 2, k \geq 3$

$$\begin{aligned} & \mathbb{P}_{\mathcal{D}}[|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon] \\ & \leq \mathbb{P}_{\mathcal{D}}[\exists h \in \mathcal{H} \text{ s.t. } |E_{\text{in}}(h) - E_{\text{out}}(h)| > \epsilon] \\ & \leq 4m_{\mathcal{H}}(2N) \exp\left(-\frac{1}{8}\epsilon^2 N\right) \\ & \text{if } k \text{ exists} \\ & \leq 4(2N)^{k-1} \exp\left(-\frac{1}{8}\epsilon^2 N\right) \end{aligned}$$

- if ①  $m_{\mathcal{H}}(N)$  breaks at  $k$  (good  $\mathcal{H}$ )
- ②  $N$  large enough (good  $\mathcal{D}$ )
- ⇒ probably generalized ' $E_{\text{out}} \approx E_{\text{in}}$ ', and
- if ③  $\mathcal{A}$  picks a  $g$  with small  $E_{\text{in}}$  (good  $\mathcal{A}$ )
- ⇒ probably learned! (-:) good luck)

## VC Dimension

the formal name of **maximum non-break point**

### Definition

VC dimension of  $\mathcal{H}$ , denoted  $d_{\text{VC}}(\mathcal{H})$  is

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

- the **most** inputs  $\mathcal{H}$  that can shatter
- $d_{\text{VC}} = \text{minimum } k' - 1$

$$\begin{aligned} N \leq d_{\text{VC}} &\implies \mathcal{H} \text{ can shatter some } N \text{ inputs} \\ k > d_{\text{VC}} &\implies k \text{ is a break point for } \mathcal{H} \end{aligned}$$

if  $N \geq 2, d_{\text{VC}} \geq 2, m_{\mathcal{H}}(N) \leq N^{d_{\text{VC}}}$

## The Four VC Dimensions

- positive rays:

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

$$m_{\mathcal{H}}(N) = N + 1$$



- positive intervals:

$$d_{\text{VC}} = 2$$

$$m_{\mathcal{H}}(N) = \frac{1}{2}N^2 + \frac{1}{2}N + 1$$



- convex sets:

$$d_{\text{VC}} = \infty$$

$$m_{\mathcal{H}}(N) = 2^N$$



- 2D perceptrons:

$$d_{\text{VC}} = 3$$

$$m_{\mathcal{H}}(N) \leq N^3 \text{ for } N \geq 2$$

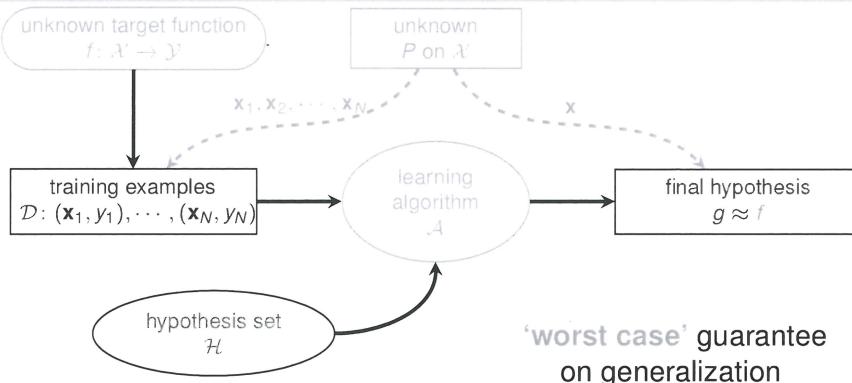


good: finite  $d_{\text{VC}}$

## VC Dimension and Learning

**finite  $d_{VC} = g$  ‘will’ generalize ( $E_{out}(g) \approx E_{in}(g)$ )**

- regardless of learning algorithm  $\mathcal{A}$
- regardless of input distribution  $P$
- regardless of target function  $f$



## Fun Time

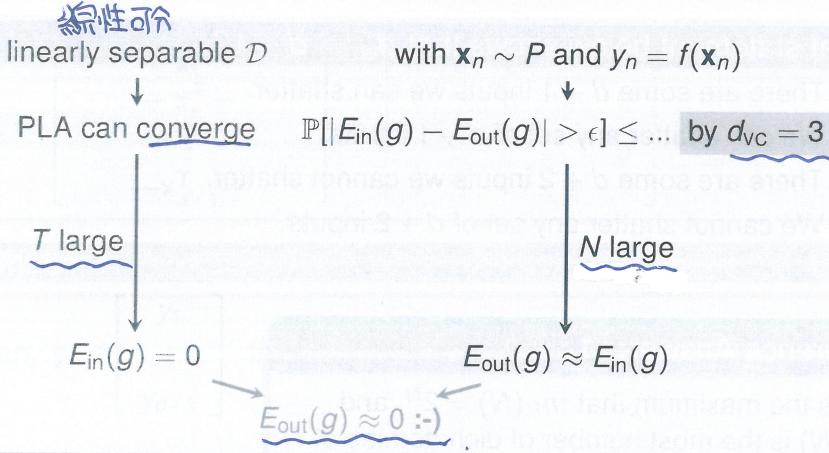
If there is a set of  $N$  inputs that cannot be shattered by  $\mathcal{H}$ . Based only on this information, what can we conclude about  $d_{VC}(\mathcal{H})$ ?

- ①  $d_{VC}(\mathcal{H}) > N$
- ②  $d_{VC}(\mathcal{H}) = N$
- ③  $d_{VC}(\mathcal{H}) < N$
- ④ no conclusion can be made

### Reference Answer: ④

It is possible that there is another set of  $N$  inputs that can be shattered, which means  $d_{VC} \geq N$ . It is also possible that no set of  $N$  input can be shattered, which means  $d_{VC} < N$ . Neither cases can be ruled out by one non-shattering set.

## 2D PLA Revisited



general PLA for  $\mathbf{x}$  with more than 2 features?

## VC Dimension of Perceptrons

- 1D perceptron (pos/neg rays):  $d_{\text{VC}} = 2$
- 2D perceptrons:  $d_{\text{VC}} = 3$ 
  - $d_{\text{VC}} \geq 3$ :
  - $d_{\text{VC}} \leq 3$ :
- $d$ -D perceptrons:  $d_{\text{VC}} \stackrel{?}{=} d + 1$

two steps:

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

$$d_{VC} \leq d + 1 \quad (1/2)$$

## A 2D Special Case

$$\begin{array}{ccc} \vdots & \vdots & X = \left[ \begin{array}{c} -\mathbf{x}_1^T- \\ -\mathbf{x}_2^T- \\ -\mathbf{x}_3^T- \\ -\mathbf{x}_4^T- \end{array} \right] = \left[ \begin{array}{ccc} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 1 \end{array} \right] \end{array}$$

○ ?

× ○

? cannot be ×

$$\mathbf{w}^T \mathbf{x}_4 = \underbrace{\mathbf{w}^T \mathbf{x}_2}_{\circ} + \underbrace{\mathbf{w}^T \mathbf{x}_3}_{\circ} - \underbrace{\mathbf{w}^T \mathbf{x}_1}_{\times} > 0$$

linear dependence restricts dichotomy

$$d_{VC} \leq d + 1 \quad (2/2)$$

## d-D General Case

$$X = \left[ \begin{array}{c} -\mathbf{x}_1^T- \\ -\mathbf{x}_2^T- \\ \vdots \\ -\mathbf{x}_{d+1}^T- \\ -\mathbf{x}_{d+2}^T- \end{array} \right]$$

more rows than columns:  
linear dependence (some  $a_i$  non-zero)  
 $\mathbf{x}_{d+2} = a_1 \mathbf{x}_1 + a_2 \mathbf{x}_2 + \dots + a_{d+1} \mathbf{x}_{d+1}$

- can you generate  $(\text{sign}(a_1), \text{sign}(a_2), \dots, \text{sign}(a_{d+1}), \times)$ ? if so, what  $\mathbf{w}$ ?

$$\begin{aligned} \mathbf{w}^T \mathbf{x}_{d+2} &= a_1 \underbrace{\mathbf{w}^T \mathbf{x}_1}_{\circ} + a_2 \underbrace{\mathbf{w}^T \mathbf{x}_2}_{\times} + \dots + a_{d+1} \underbrace{\mathbf{w}^T \mathbf{x}_{d+1}}_{\times} \\ &> 0 \text{(contradiction!)} \end{aligned}$$

'general' X no-shatter  $\implies d_{VC} \leq d + 1$

## Fun Time

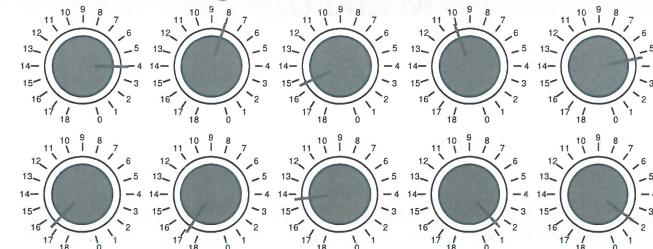
Based on the proof above, what is  $d_{VC}$  of 1126-D perceptrons?

- ① 1024
- ② 1126
- ③ 1127
- ④ 6211

Reference Answer: (3)

Well, too much fun for this section! :-)

## Degrees of Freedom



(modified from the work of Hugues Vermeiren on <http://www.texample.net>)

- hypothesis parameters  $\mathbf{w} = (w_0, w_1, \dots, w_d)$ : **creates degrees of freedom**
- hypothesis quantity  $M = |\mathcal{H}|$ : **'analog' degrees of freedom**
- hypothesis 'power'  $d_{VC} = d + 1$ : **effective 'binary' degrees of freedom**

$d_{VC}(\mathcal{H})$ : powerfulness of  $\mathcal{H}$

## Two Old Friends

### Positive Rays ( $d_{VC} = 1$ )

A horizontal line segment with arrows at both ends. A tick mark labeled 'a' is on the left side. The left part of the line is labeled  $h(x) = -1$  and the right part is labeled  $h(x) = +1$ .

free parameters:  $a$

### Positive Intervals ( $d_{VC} = 2$ )

A horizontal line segment with arrows at both ends. There are two tick marks labeled 'l' and 'r'. The interval between them is labeled  $h(x) = +1$ . The regions outside this interval are labeled  $h(x) = -1$ .

free parameters:  $\ell, r$

practical rule of thumb:

$d_{VC} \approx \# \text{free parameters}$  (but not always)

## $M$ and $d_{VC}$

copied from Lecture 5 :-)

- ① can we make sure that  $E_{out}(g)$  is close enough to  $E_{in}(g)$ ?
- ② can we make  $E_{in}(g)$  small enough?

### small $M$

- ① Yes!,  $\mathbb{P}[\text{BAD}] \leq 2 \cdot M \cdot \exp(\dots)$
- ② No!, too few choices

### large $M$

- ① No!,  $\mathbb{P}[\text{BAD}] \leq 2 \cdot M \cdot \exp(\dots)$
- ② Yes!, many choices

### small $d_{VC}$

- ① Yes!,  $\mathbb{P}[\text{BAD}] \leq 4 \cdot (2N)^{d_{VC}} \cdot \exp(\dots)$
- ② No!, too limited power

### large $d_{VC}$

- ① No!,  $\mathbb{P}[\text{BAD}] \leq 4 \cdot (2N)^{d_{VC}} \cdot \exp(\dots)$
- ② Yes!, lots of power

using the right  $d_{VC}$  (or  $\mathcal{H}$ ) is important

## Fun Time

Origin-crossing Hyperplanes are essentially perceptrons with  $w_0$  fixed at 0. Make a guess about the  $d_{VC}$  of origin-crossing hyperplanes in  $\mathbb{R}^d$ .

- ① 1
- ②  $d$
- ③  $d + 1$
- ④  $\infty$

## Reference Answer: ②

The proof is almost the same as proving the  $d_{VC}$  for usual perceptrons, but it is the **intuition** ( $d_{VC} \approx \# \text{free parameters}$ ) that you shall use to answer this quiz.

## VC Bound Rephrase: Penalty for Model Complexity

For any  $g = \mathcal{A}(\mathcal{D}) \in \mathcal{H}$  and ‘statistical’ large  $\mathcal{D}$ , for  $N \geq 2$ ,  $d_{VC} \geq 2$

$$\mathbb{P}_{\mathcal{D}} \left[ \underbrace{|E_{in}(g) - E_{out}(g)| > \epsilon}_{\text{BAD}} \right] \leq \underbrace{4(2N)^{d_{VC}} \exp\left(-\frac{1}{8}\epsilon^2 N\right)}_{\delta}$$

## Rephrase

..., with probability  $\geq 1 - \delta$ , **GOOD**:  $|E_{in}(g) - E_{out}(g)| \leq \epsilon$

$$\begin{aligned} \text{set } \delta &= 4(2N)^{d_{VC}} \exp\left(-\frac{1}{8}\epsilon^2 N\right) \\ \frac{\delta}{4(2N)^{d_{VC}}} &= \exp\left(-\frac{1}{8}\epsilon^2 N\right) \\ \ln\left(\frac{4(2N)^{d_{VC}}}{\delta}\right) &= \frac{1}{8}\epsilon^2 N \\ \sqrt{\frac{8}{N} \ln\left(\frac{4(2N)^{d_{VC}}}{\delta}\right)} &= \epsilon \end{aligned}$$

## VC Bound Rephrase: Penalty for Model Complexity

For any  $g = \mathcal{A}(\mathcal{D}) \in \mathcal{H}$  and 'statistical' large  $\mathcal{D}$ , for  $N \geq 2$ ,  $d_{\text{vc}} \geq 2$

$$\mathbb{P}_{\mathcal{D}} \left[ \underbrace{|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon}_{\text{BAD}} \right] \leq \underbrace{4(2N)^{d_{\text{vc}}} \exp\left(-\frac{1}{8}\epsilon^2 N\right)}_{\delta}$$

### Rephrase

..., with probability  $\geq 1 - \delta$ , **GOOD!**

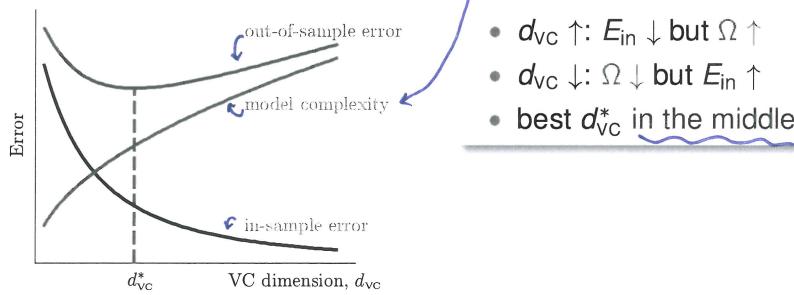
$$\begin{aligned} \text{gen. error } |E_{\text{in}}(g) - E_{\text{out}}(g)| &\leq \sqrt{\frac{8}{N} \ln\left(\frac{4(2N)^{d_{\text{vc}}}}{\delta}\right)} \\ E_{\text{in}}(g) - \sqrt{\frac{8}{N} \ln\left(\frac{4(2N)^{d_{\text{vc}}}}{\delta}\right)} &\leq E_{\text{out}}(g) \leq E_{\text{in}}(g) + \sqrt{\frac{8}{N} \ln\left(\frac{4(2N)^{d_{\text{vc}}}}{\delta}\right)} \end{aligned}$$

$\underbrace{\sqrt{\dots}}_{\Omega(N, \mathcal{H}, \delta)}$  : penalty for model complexity

## THE VC Message

with a high probability,

$$E_{\text{out}}(g) \leq E_{\text{in}}(g) + \underbrace{\sqrt{\frac{8}{N} \ln\left(\frac{4(2N)^{d_{\text{vc}}}}{\delta}\right)}}_{\Omega(N, \mathcal{H}, \delta)}$$



powerful  $\mathcal{H}$  not always good!

## VC Bound Rephrase: Sample Complexity

For any  $g = \mathcal{A}(\mathcal{D}) \in \mathcal{H}$  and ‘statistical’ large  $\mathcal{D}$ , for  $N \geq 2, d_{\text{VC}} \geq 2$

$$\mathbb{P}_{\mathcal{D}} \left[ \underbrace{|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon}_{\text{BAD}} \right] \leq \underbrace{4(2N)^{d_{\text{VC}}} \exp\left(-\frac{1}{8}\epsilon^2 N\right)}_{\delta}$$

given specs  $\epsilon = 0.1, \delta = 0.1, d_{\text{VC}} = 3$ , want  $4(2N)^{d_{\text{VC}}} \exp\left(-\frac{1}{8}\epsilon^2 N\right) \leq \delta$

$N$	bound
100	$2.82 \times 10^7$
1,000	$9.17 \times 10^9$
10,000	$1.19 \times 10^8$
100,000	$1.65 \times 10^{-38}$
29,300	$9.99 \times 10^{-2}$

sample complexity:  
need  $N \approx 10,000d_{\text{VC}}$  in theory

practical rule of thumb:

$N \approx 10d_{\text{VC}}$  often enough!

## Looseness of VC Bound

$$\mathbb{P}_{\mathcal{D}} \left[ |E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon \right] \leq 4(2N)^{d_{\text{VC}}} \exp\left(-\frac{1}{8}\epsilon^2 N\right)$$

theory:  $N \approx 10,000d_{\text{VC}}$ ; practice:  $N \approx 10d_{\text{VC}}$

### Why?

- Hoeffding for unknown  $E_{\text{out}}$  any distribution, any target
- $m_{\mathcal{H}}(N)$  instead of  $|\mathcal{H}(\mathbf{x}_1, \dots, \mathbf{x}_N)|$  ‘any’ data
- $N^{d_{\text{VC}}}$  instead of  $m_{\mathcal{H}}(N)$  ‘any’  $\mathcal{H}$  of same  $d_{\text{VC}}$
- union bound on worst cases any choice made by  $\mathcal{A}$

—but hardly better, and ‘similarly loose for all models’

philosophical message of VC bound  
important for improving ML

## Fun Time

Consider the VC Bound below. How can we decrease the probability of getting **BAD** data?

$$\mathbb{P}_{\mathcal{D}} \left[ |E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon \right] \leq 4(2N)^{d_{\text{VC}}} \exp \left( -\frac{1}{8}\epsilon^2 N \right)$$

- ① decrease model complexity  $d_{\text{VC}}$
- ② increase data size  $N$  a lot
- ③ increase generalization error tolerance  $\epsilon$
- ④ all of the above

Reference Answer: ④

Congratulations on being  
Master of VC bound! :-)

## Summary

- ① When Can Machines Learn?
- ② Why Can Machines Learn?

Lecture 6: Theory of Generalization

Lecture 7: The VC Dimension

- Definition of VC Dimension  
**maximum non-break point**
- VC Dimension of Perceptrons  
 $d_{\text{VC}}(\mathcal{H}) = d + 1$
- Physical Intuition of VC Dimension  
 $d_{\text{VC}} \approx \# \text{free parameters}$
- Interpreting VC Dimension  
**loosely: model complexity & sample complexity**

- next: more than noiseless binary classification?

- ③ How Can Machines Learn?
- ④ How Can Machines Learn Better?