

# Machine Learning

## (機器學習)

### Lecture 4: Theory of Generalization

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)



# Roadmap

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

## Lecture 4: Theory of Generalization

- Effective Number of Lines
- Effective Number of Hypotheses
- Break Point
- Definition of VC Dimension
- VC Dimension of Perceptrons
- Physical Intuition of VC Dimension
- Interpreting VC Dimension

Is  $M = \infty$  Feasible?

- input  $x \in [-1, +1] \subset \mathbb{R}^1$ , uniform iid
- target  $f(x) = \text{sign}(x)$ , taking  $\text{sign}(0) = +1$
- hypothesis set:  $h_a(x) = \text{sign}(x - a)$  for  $a \in [-1, 1]$   
**infinitely many  $a$**
- algorithm:  $g = h_{a^*}$  with  $a^* = \min_{y_n = +1} x_n$ ,  
assuming at least one  $y_n = 1$

- for  $\epsilon < 0.5$ ,  $E_{\text{out}}(g) > \epsilon$  if every  $y_n = +1$  satisfies  $x_n > 2\epsilon$

$$\mathbb{P} \left[ \underbrace{|E_{\text{in}}(g) - E_{\text{out}}(g)|}_0 > \epsilon \right] \leq \left( \frac{2 - 2\epsilon}{2} \right)^N$$

**BAD data** can happen rarely  
even for infinitely many hypotheses

Where Did **M** Come From?

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

- BAD events**  $\mathcal{B}_m$ :  $|E_{\text{in}}(h_m) - E_{\text{out}}(h_m)| > \epsilon$
  - to give  $\mathcal{A}$  freedom of choice: bound  $\mathbb{P}[\mathcal{B}_1 \text{ or } \mathcal{B}_2 \text{ or } \dots \mathcal{B}_M]$
  - worst case: all  $\mathcal{B}_m$  non-overlapping
- 把所有 hypothesis 發生 BAD 的機率加總

$$\mathbb{P}[\mathcal{B}_1 \text{ or } \mathcal{B}_2 \text{ or } \dots \mathcal{B}_M] \leq \frac{\mathbb{P}[\mathcal{B}_1] + \mathbb{P}[\mathcal{B}_2] + \dots + \mathbb{P}[\mathcal{B}_M]}{\text{union bound}}$$

↑ 都不 overlap.

union bound 有可能變成  $\infty$

where did union bound fail  
to consider for  $M = \infty$ ?

\*BAD.  $E_{\text{in}}$  and  $E_{\text{out}}$  差很多

# Where Did Union Bound Fail?

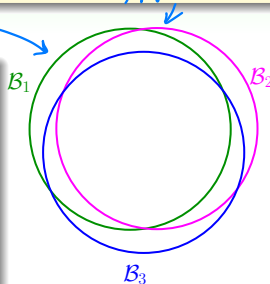
union bound  $\mathbb{P}[B_1] + \mathbb{P}[B_2] + \dots + \mathbb{P}[B_M]$  实际上会

但  $B_1 \sim B_M$  真的不会 overlap 吗?

- BAD events**  $B_m$ :  $|E_{\text{in}}(h_m) - E_{\text{out}}(h_m)| > \epsilon$

overlapping for similar hypotheses  $h_1 \approx h_2$   
(e.g. if  $a_1 \approx a_2$  in previous example)

- why? ①  $E_{\text{out}}(h_1) \approx E_{\text{out}}(h_2)$   
② for most  $\mathcal{D}$ ,  $E_{\text{in}}(h_1) = E_{\text{in}}(h_2)$
- union bound **over-estimating**



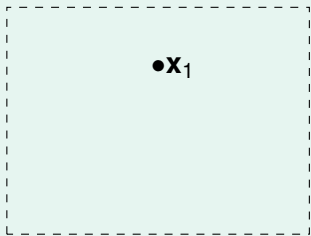
to account for overlap,  
can we group similar hypotheses by kind?

想办法找重叠的部分

# How Many Lines Are There? (1/2)

$$\mathcal{H} = \left\{ \text{all lines in } \mathbb{R}^2 \right\}$$

- how many lines?  $\infty$
- how many **kinds of** lines if viewed from one input vector  $\mathbf{x}_1$ ?



**2 kinds:**  $h_1\text{-like}(\mathbf{x}_1) = \circ$  or  $h_2\text{-like}(\mathbf{x}_1) = \times$

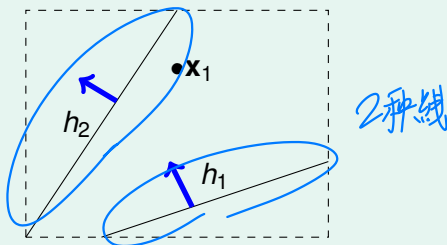
2种线

如果我们只看一个点：线不是过  $\mathbf{x}_1$  变成  $\circ, \times$

# How Many Lines Are There? (1/2)

$$\mathcal{H} = \left\{ \text{all lines in } \mathbb{R}^2 \right\}$$

- how many lines?  $\infty$
- how many **kinds of** lines if viewed from one input vector  $\mathbf{x}_1$ ?



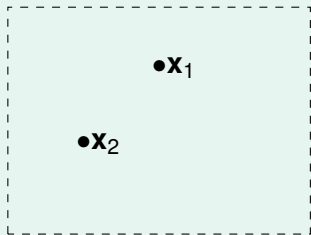
**2 kinds:**  $h_1\text{-like}(\mathbf{x}_1) = \circ$  or  $h_2\text{-like}(\mathbf{x}_1) = \times$

# How Many Lines Are There? (2/2)

$$\mathcal{H} = \left\{ \text{all lines in } \mathbb{R}^2 \right\}$$

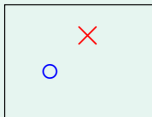
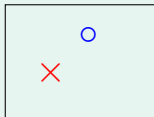
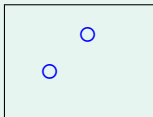
- how many **kinds of** lines if viewed from two inputs  $\mathbf{x}_1, \mathbf{x}_2$ ?

二个点



4种线

4:



:

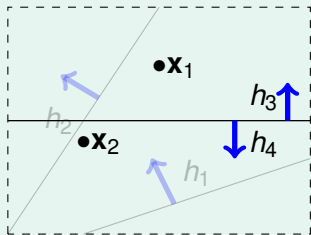
one input: 2; two inputs: 4; **three inputs?**



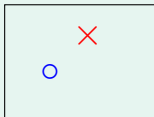
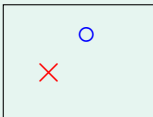
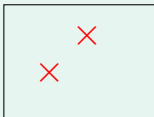
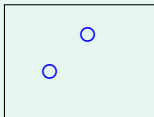
# How Many Lines Are There? (2/2)

$$\mathcal{H} = \left\{ \text{all lines in } \mathbb{R}^2 \right\}$$

- how many **kinds of** lines if viewed from two inputs  $\mathbf{x}_1, \mathbf{x}_2$ ?



4:



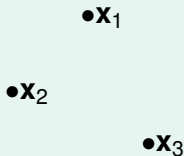
one input: 2; two inputs: 4; **three inputs?**

# How Many Kinds of Lines for Three Inputs? (1/2)

$$\mathcal{H} = \left\{ \text{all lines in } \mathbb{R}^2 \right\}$$

必須 linear separable

for three inputs  $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3$

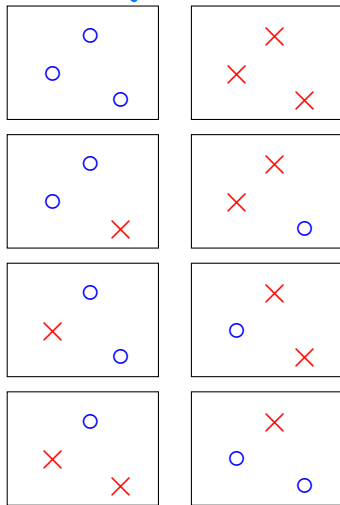


3點

always 8 for three inputs?

$2^n$

8:

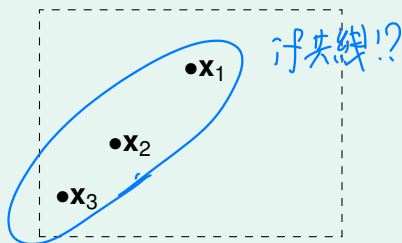


# How Many Kinds of Lines for Three Inputs? (2/2)

$$\mathcal{H} = \left\{ \text{all lines in } \mathbb{R}^2 \right\}$$

for **another** three inputs

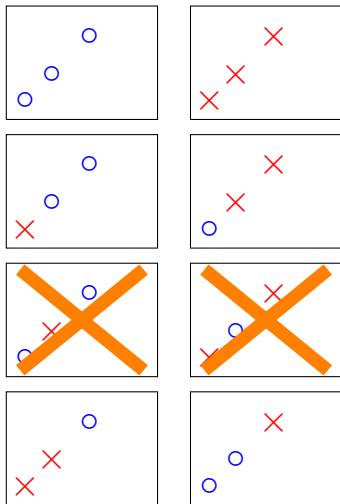
$\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3$



**'fewer than 8'** when degenerate  
(e.g. collinear or same inputs)

6:

不是 linearly  
separable

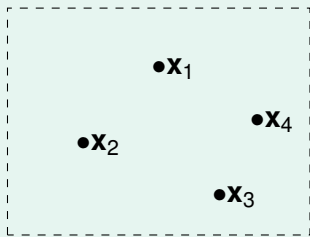


# How Many Kinds of Lines for Four Inputs?

$$\mathcal{H} = \left\{ \text{all lines in } \mathbb{R}^2 \right\}$$

任意  $x_1 \sim x_4$  的位置

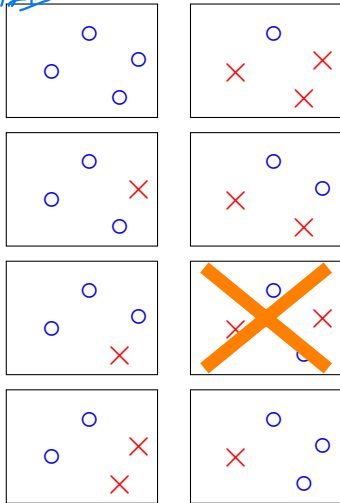
for four inputs  $x_1, x_2, x_3, x_4$



for any four inputs

at most 14

14:  $2 \times$



# Effective Number of Lines

maximum kinds of lines with respect to  $N$  inputs  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$

$\iff$  **effective number of lines**

- must be  $\leq 2^N$  (why?) ← 這是所有 0, 1 的可能性
- finite 'grouping' of infinitely-many lines  $\in \mathcal{H}$
- wish:

$$\begin{aligned} & \mathbb{P} [|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon] \\ & \leq 2 \cdot \underbrace{\text{effective}(N)}_{\text{取代原來的 } M} \cdot \exp(-2\epsilon^2 N) \end{aligned}$$

遠小於  $2^N$

## lines in 2D

$N$	effective( $N$ )
1	2
2	4
3	8
4	14 <span style="color: orange;">&lt; <math>2^N</math></span>

- if
- ① effective( $N$ ) can replace  $M$  and
  - ② effective( $N$ )  $\ll 2^N$  ← upper bound.

**learning possible with infinite lines :-)**

# Questions?

# Dichotomies: Mini-hypotheses

$$\mathcal{H} = \{\text{hypothesis } h: \mathcal{X} \rightarrow \{\times, \circ\}\}$$

- call

$$h(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N) = (\underbrace{h(\mathbf{x}_1), h(\mathbf{x}_2), \dots, h(\mathbf{x}_N)}_{N\text{个集分成 } 0 \text{ 或 } \times}) \in \{\times, \circ\}^N$$

a **dichotomy**: hypothesis 'limited' to the eyes of  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$

- $\mathcal{H}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$ : *这是 dichotomy*

*2分*  
all dichotomies 'implemented' by  $\mathcal{H}$  on  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$

*针对 N 个集定的集*

	<u>hypotheses <math>\mathcal{H}</math></u>	<u>dichotomies <math>\mathcal{H}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)</math></u>
e.g.	all lines in $\mathbb{R}^2$	$\{\circ\circ\circ\circ, \circ\circ\circ\times, \circ\circ\times\times, \dots\}$
size	possibly infinite	upper bounded by $2^N$

*efficient(N)*

$|\mathcal{H}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)|$ : candidate for **replacing  $M$**

# Growth Function

- $|\mathcal{H}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)|$  depend on inputs  $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$  *不希望這值與input怎麼變*
  - growth function: *有關*  
remove dependence by **taking max of all possible**  $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$  *選最大的 dichotomies set*
- $$m_{\mathcal{H}}(N) = \max_{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N \in \mathcal{X}} |\mathcal{H}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)|$$
- 取代 M*  
*稱 growth function*
- finite, upper-bounded by  $2^N$

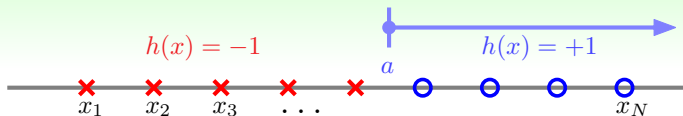
## lines in 2D

$N$	$m_{\mathcal{H}}(N)$
1	2 <i>共幾</i>
2	4
3	$\max(\dots, \underline{6}, 8)$ $= \underline{8}$
4	$14 < 2^N$

how to 'calculate' the growth function?



# Growth Function for Positive Rays



- $\mathcal{X} = \mathbb{R}$  (one dimensional) 一維數線上
- $\mathcal{H}$  contains  $h$ , where **each**  $h(x) = \text{sign}(x - a)$  **for threshold**  $a$  不同的 threshold 決定不同的 hypothesis
- 'positive half' of 1D perceptrons 1D perceptrons

one dichotomy for  $a \in$  each spot  $(x_n, x_{n+1})$ :

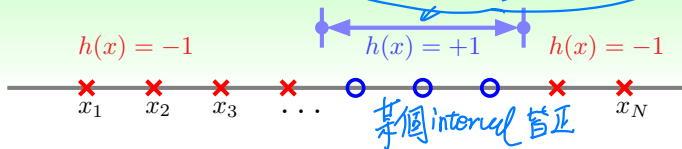
$$m_{\mathcal{H}}(N) = \underline{N+1}$$

← 植樹問題

$x_1$	$x_2$	$x_3$	$x_4$
o	o	o	o
x	o	o	o
x	x	o	o
x	x	x	o
x	x	x	x

$(N+1) \ll 2^N$  when  $N$  large!

# Growth Function for Positive Intervals



- $\mathcal{X} = \mathbb{R}$  (one dimensional)
  - $\mathcal{H}$  contains  $h$ , where **each**  $h(x) = +1$  iff  $x \in [l, r]$ , **-1 otherwise**
- Handwritten blue text: "for interval, others are negative" (for interval, others are negative)

one dichotomy for each 'interval kind'

$$m_{\mathcal{H}}(N) = \underbrace{(N+1)}_{\text{interval ends in } N+1 \text{ spots}} \underbrace{2}_{\text{任意2个作为interval}} + \underbrace{1}_{\text{all } \times}$$

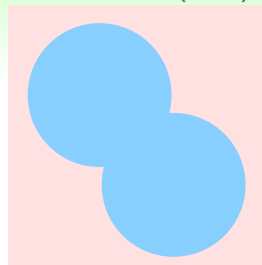
$$= \frac{1}{2}N^2 + \frac{1}{2}N + 1$$

Handwritten blue text: "N+1 个空位" (N+1 empty spaces), "任意2个作为interval" (any 2 points as interval), "all x" (all x).

$x_1$	$x_2$	$x_3$	$x_4$
○	×	×	×
○	○	×	×
○	○	○	×
○	○	○	○
×	○	×	×
×	○	○	×
×	○	○	○
×	×	○	×
×	×	○	○
×	×	×	○
×	×	×	×

$$\left(\frac{1}{2}N^2 + \frac{1}{2}N + 1\right) \ll \underline{2^N} \text{ when } N \text{ large!}$$

## Growth Function for Convex Sets (1/2)

Blue:  $\oplus$ Pink:  $\ominus$ 

convex region in blue

non-convex region

- $\mathcal{X} = \mathbb{R}^2$  (two dimensional)
- $\mathcal{H}$  contains  $h$ , where  $h(\mathbf{x}) = +1$  iff  $\mathbf{x}$  in a convex region,  $-1$  otherwise

what is  $m_{\mathcal{H}}(N)$ ?

## Growth Function for Convex Sets (2/2)

- one possible set of  $N$  inputs:  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$  on a big circle
- every dichotomy can be implemented** by  $\mathcal{H}$  using a convex region slightly extended from **contour of positive inputs**

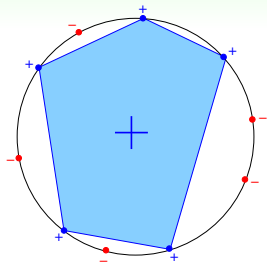
$$m_{\mathcal{H}}(N) = 2^N$$

這就是  $2^N$

- call those  $N$  inputs **'shattered'** by  $\mathcal{H}$

找到特定的  $N$  個點，使  $m_{\mathcal{H}}(N) = 2^N$ ，將此  $\mathcal{H}$  shattered.

$m_{\mathcal{H}}(N) = 2^N \iff$   
**exists**  $N$  inputs that can be shattered



bottom

# The Four Growth Functions

- positive rays:

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

- positive intervals:

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

- convex sets:

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

- 2D perceptrons: 至少  $N=4$  時,  $m_{\mathcal{H}}(4)=14$   
 $m_{\mathcal{H}}(N) < 2^N$  in some cases

what if  $m_{\mathcal{H}}(N)$  replaces  $M$ ?

$$\mathbb{P} \left[ |E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon \right] \stackrel{?}{\leq} 2 \cdot m_{\mathcal{H}}(N) \cdot \exp(-2\epsilon^2 N) \leftarrow \text{跟這項比}$$

polynomial: good; exponential: bad

↑ 增加的比 exp 慢

for 2D or general perceptrons,

$m_{\mathcal{H}}(N)$  **polynomial**?

# Break Point of $\mathcal{H}$

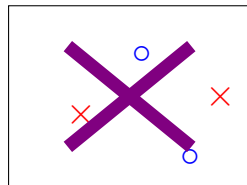
what do we know about 2D perceptrons now?

**three inputs: 'exists' shatter;**  
**four inputs, 'for all' no shatter**

if no  $k$  inputs can be shattered by  $\mathcal{H}$ ,  
 call  $k$  a **break point** for  $\mathcal{H}$

- $m_{\mathcal{H}}(k) < 2^k$  某處小於  $2^k$  之後後面也都小於
- $k + 1, k + 2, k + 3, \dots$  also break points!
- will study minimum break point  $k$

第一個做不到  $2^n$  的臭



2D perceptrons: **minimum break point at 4**

$$m_{\mathcal{H}}(4) = 14 < 2^4$$

# The Four Minimum Break Points

- positive rays:

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

minimum break point at 2  $m_{\mathcal{H}}(2) = 2 < 2^2 = 4$

- positive intervals:

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

minimum break point at 3  $m_{\mathcal{H}}(3) = \frac{9}{2} + \frac{1}{2} + 1 = 6 < 2^3$

- convex sets:

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

no break point

- 2D perceptrons:

$$m_{\mathcal{H}}(N) < 2^N \text{ in some cases}$$

minimum break point at 4  $m_{\mathcal{H}}(4) = 14$

theorem from combinatorics

(not going to prove in class):

- no break point:  $m_{\mathcal{H}}(N) = 2^N$  (sure!) By definition
- minimum break point  $k$ :

$m_{\mathcal{H}}(N) = O(N^{k-1})$   $\leftarrow$  猜想 grow function 与  $k$  有关

**Questions?**



BAD Bound for General  $\mathcal{H}$ 

want:

$$\mathbb{P}\left[\exists h \in \mathcal{H} \text{ s.t. } |E_{\text{in}}(h) - E_{\text{out}}(h)| > \epsilon\right] \leq 2 \underbrace{m_{\mathcal{H}}(N)}_{\text{取代數}} \cdot \exp\left(-2 \epsilon^2 N\right)$$

actually, when  $N$  large enough,

$$\mathbb{P}\left[\exists h \in \mathcal{H} \text{ s.t. } |E_{\text{in}}(h) - E_{\text{out}}(h)| > \epsilon\right] \leq 2 \cdot \underline{2} m_{\mathcal{H}}(\underline{2}N) \cdot \exp\left(-2 \cdot \underline{\frac{1}{16}} \epsilon^2 N\right)$$

不好證

called Vapnik-Chervonenkis (VC) Bound

# Interpretation of Vapnik-Chervonenkis (VC) Bound

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

break point

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

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

## VC Dimension

break point . 的正式名稱

the formal name of **maximum non-break** point  $d_{VC}$   
 = (minimum break point  $k - 1$ ) *最大的 non-break point*

## Definition

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

**largest**  $N$  for which  $m_{\mathcal{H}}(N) = 2^N$   
 (the **most** inputs  $\mathcal{H}$  that can shatter)

$$d_{VC} = \text{minimum } k - 1$$

這  $N$  个 point 可能被 shattered.

2D perception

$$m_{\mathcal{H}}(1) = 2$$

$$m_{\mathcal{H}}(2) = 4$$

$$m_{\mathcal{H}}(3) = 8 \leftarrow d_{VC}(\mathcal{H}) = 3$$

$$m_{\mathcal{H}}(4) = 14 \leftarrow k=4 \text{ (break point)}$$


$N \leq d_{VC} \implies \mathcal{H}$  can shatter some  $N$  inputs

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

比  $d_{VC}$  大一定是 break point

$$\text{if } N \geq 2, d_{VC} \geq 2, \underbrace{m_{\mathcal{H}}(N)}_{\text{grow function}} \leq \underbrace{N^{d_{VC}}}_{\text{upper bound.}}$$

## The Four VC Dimensions

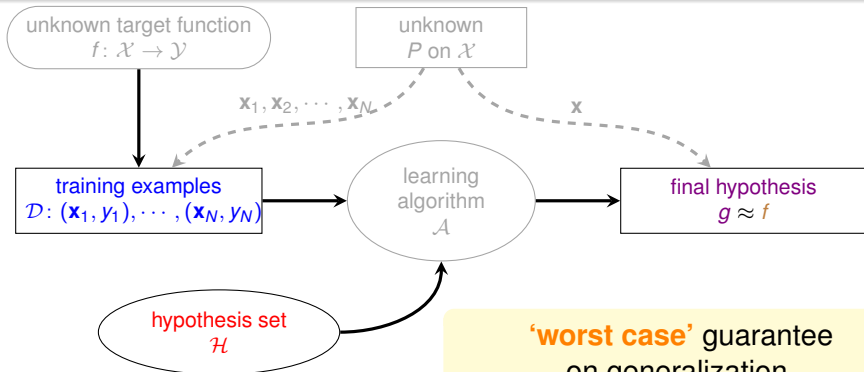
- positive rays:  $m_H(1)=2$ ,  $m_H(2)=3$  <sup>break point</sup>  $m_H(N) = N + 1$   
 $d_{VC} = 1$
- positive intervals:  $m_H(1)=2$ ,  $m_H(2)=4$   $m_H(N) = \frac{1}{2}N^2 + \frac{1}{2}N + 1$   
 $d_{VC} = 2$   $m_H(3)=7$  <sup>break point</sup>
- convex sets:  $m_H(N) = 2^N$   
 $d_{VC} = \infty$  
- 2D perceptrons:  $m_H(3)=8$   $m_H(4)=14$   $m_H(N) \leq N^3$  for  $N \geq 2$   
 $d_{VC} = 3$   <sup>$N^{d_{VC}}$</sup>

good: finite  $d_{VC}$

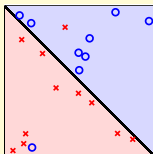
# VC Dimension and Learning

**finite  $d_{VC} \implies g$  'will' generalize ( $E_{out}(g) \approx E_{in}(g)$ )**

- regardless of learning algorithm  $\mathcal{A}$
  - regardless of input distribution  $P$
  - regardless of target function  $f$
- VC dimension 跟這些東西無关*



# From Noiseless VC to Noisy VC



real-world learning problems are often **noisy**

age	23 years
gender	female
annual salary	NTD 1,000,000
year in residence	1 year
year in job	0.5 year
current debt	200,000

credit? {no(-1), yes(+1)}

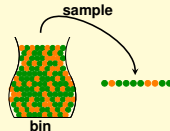
but more!

- **noise in  $\mathbf{x}$**  (covered by  $P(\mathbf{x})$ ): inaccurate customer information?
- **noise in  $y$**  (covered by  $P(y|\mathbf{x})$ ): good customer, 'misabeled' as bad?

does VC bound work under **noise**?

# Probabilistic Marbles

one key of VC bound: **marbles!**



## 'deterministic' marbles

- marble  $\mathbf{x} \sim P(\mathbf{x})$
- deterministic color  
 $\llbracket f(\mathbf{x}) \neq h(\mathbf{x}) \rrbracket$

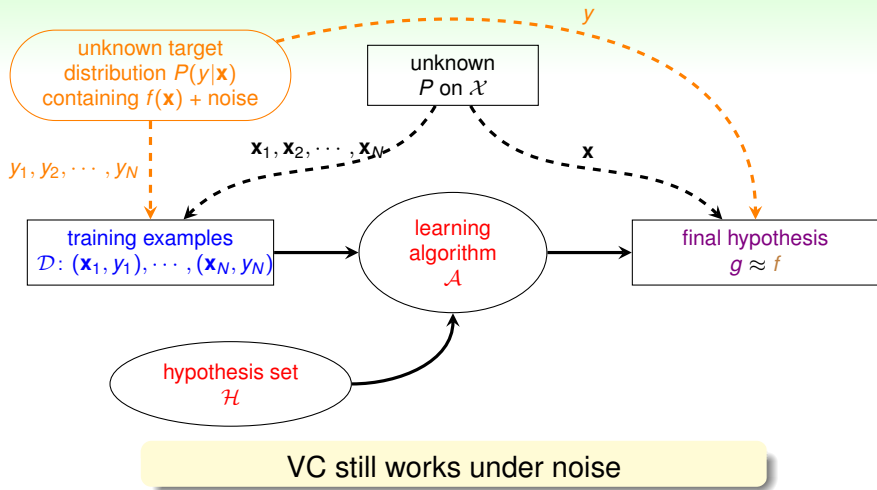
## 'probabilistic' (noisy) marbles

- marble  $\mathbf{x} \sim P(\mathbf{x})$
- probabilistic color  
 $\llbracket y \neq h(\mathbf{x}) \rrbracket$  with  $y \sim P(y|\mathbf{x})$

**same nature:** can estimate  $\mathbb{P}[\text{orange}]$  if  $\overset{i.i.d.}{\sim}$

VC holds for  $\underbrace{\mathbf{x} \overset{i.i.d.}{\sim} P(\mathbf{x}), y \overset{i.i.d.}{\sim} P(y|\mathbf{x})}_{(\mathbf{x}, y) \overset{i.i.d.}{\sim} P(\mathbf{x}, y)}$

# The New Learning Flow





**Questions?**

## 2D PLA Revisited

linearly separable  $\mathcal{D}$ with  $\mathbf{x}_n \sim P$  and  $y_n = f(\mathbf{x}_n)$ 

PLA can converge

 $\mathbb{P}[|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon] \leq \dots$  by  $d_{\text{VC}} = 3$  $T$  large $N$  large

$E_{\text{in}}(g) = 0$



$E_{\text{out}}(g) \approx E_{\text{in}}(g)$

$E_{\text{out}}(g) \approx 0 \text{ :-)}$

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

多維度

# VC Dimension of Perceptrons

- 1D perceptron (pos/neg rays):  $d_{VC} = 2$
- 2D perceptrons:  $d_{VC} = 3$   $m_{H(3)} = 8$ 
  - $d_{VC} \geq 3$ : 
  - $d_{VC} \leq 3$ : 
- $d$ -D perceptrons:  $d_{VC} \stackrel{?}{=} d + 1$  *guess*

*Proof*

two steps:

- $d_{VC} \geq d + 1$
  - $d_{VC} \leq d + 1$
- 證這 2 件事情*

# Extra Fun Time

What statement below shows that  $d_{VC} \geq d + 1$ ?

- ① There are some  $d + 1$  inputs we can shatter. ]  $d_{VC}$  至少大於  $d+1$
- ② We can shatter any set of  $d + 1$  inputs.
- ③ There are some  $d + 2$  inputs we cannot shatter. ] 但搞不好有一些可以
- ④ We cannot shatter any set of  $d + 2$  inputs. ⇒ break point  $\Rightarrow d_{VC} = d+1$

## Reference Answer: ①

$d_{VC}$  is the maximum that  $m_{\mathcal{H}}(N) = 2^N$ , and  $m_{\mathcal{H}}(N)$  is the most number of dichotomies of  $N$  inputs. So if we can find  $2^{d+1}$  dichotomies on *some*  $d + 1$  inputs,  $m_{\mathcal{H}}(d + 1) = 2^{d+1}$  and hence  $d_{VC} \geq d + 1$ .

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

There are **some**  $d + 1$  **inputs** we can shatter.

- some 'trivial' inputs:  $\leftarrow$  定義 - 一个特别 input 使它能够被 shattered

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

- visually in 2D:  $\begin{matrix} \bullet \\ \bullet \end{matrix} \Rightarrow \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} = X$

3 point: 可以被 shattered.

note: **X invertible!**

## Can We Shatter X?

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

to shatter ...

for any  $\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_{d+1} \end{bmatrix}$ , find  $\mathbf{w}$  such that

$\leftarrow$   $O$  or  $X$  的排列組合

結論, 若  $X$  invertible  
任何  $\mathbf{y}$ , 都能得出  $\mathbf{w}$

$$\boxed{\text{sign}(X\mathbf{w})} = \mathbf{y} \quad \leftarrow \quad \boxed{(X\mathbf{w}) = \mathbf{y}} \xleftrightarrow{X \text{ invertible!}} \mathbf{w} = \boxed{X^{-1}\mathbf{y}}$$

perception      更嚴格

'special'  $X$  can be shattered  $\Rightarrow d_{VC} \geq d + 1$

$X$  可逆

# Extra Fun Time

What statement below shows that  $d_{VC} \leq d + 1$ ?

- ① There are some  $d + 1$  inputs we can shatter.  $\leftarrow$  但可能有些不行
- ② We can shatter any set of  $d + 1$  inputs.  $\leftarrow$  那繼續找
- ③ There are some  $d + 2$  inputs we cannot shatter.  $\leftarrow$  但可能有些可以
- ④ We cannot shatter any set of  $d + 2$  inputs.  $d+2$  是 break point, 所以  $d_{VC}$  不可能大於  $d+1$

Reference Answer: ④

$d_{VC}$  is the maximum that  $m_{\mathcal{H}}(N) = 2^N$ , and  $m_{\mathcal{H}}(N)$  is the most number of dichotomies of  $N$  inputs. So if we cannot find  $2^{d+2}$  dichotomies on *any*  $d + 2$  inputs (i.e. break point),  $m_{\mathcal{H}}(d + 2) < 2^{d+2}$  and hence  $d_{VC} < d + 2$ . That is,  $d_{VC} \leq d + 1$ .

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

## A 2D Special Case

$$\begin{matrix} \bullet & \bullet \\ \bullet & \bullet \end{matrix} \quad X = \begin{bmatrix} -\mathbf{x}_1^T- \\ -\mathbf{x}_2^T- \\ -\mathbf{x}_3^T- \\ -\mathbf{x}_4^T- \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 1 \end{bmatrix} \quad \text{四行}$$

○ ?

× ○

○ ×

× ○

← 無法做出来

? cannot be ×

$$\mathbf{w}^T \mathbf{x}_4 = \underbrace{\mathbf{w}^T \mathbf{x}_2}_{\oplus} + \underbrace{\mathbf{w}^T \mathbf{x}_3}_{\oplus} - \underbrace{\mathbf{w}^T \mathbf{x}_1}_{\oplus} > 0 \quad \rightarrow \text{一定是 } \oplus \text{ 代表 } \mathbf{w}^T \mathbf{x}_4 \text{ 只能是 } \bigcirc, \text{ 不能是 } \times$$

這是因為  $\mathbf{w}^T \mathbf{x}_4$  是  $\mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_1$  的線性組合

linear dependence **restricts dichotomy**



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

## d-D General Case

$$X = \begin{bmatrix} \text{---} \mathbf{x}_1^T \text{---} \\ \text{---} \mathbf{x}_2^T \text{---} \\ \vdots \\ \text{---} \mathbf{x}_{d+1}^T \text{---} \\ \text{---} \mathbf{x}_{d+2}^T \text{---} \end{bmatrix}$$

$d+2$

$d+1$

more rows than columns:

linear dependence (some  $a_i$  non-zero)

$$\mathbf{x}_{d+2} = a_1 \mathbf{x}_1 + a_2 \mathbf{x}_2 + \dots + a_{d+1} \mathbf{x}_{d+1}$$

$d+2 \times d+1$  rank(X) 必小於  $d+2$ .

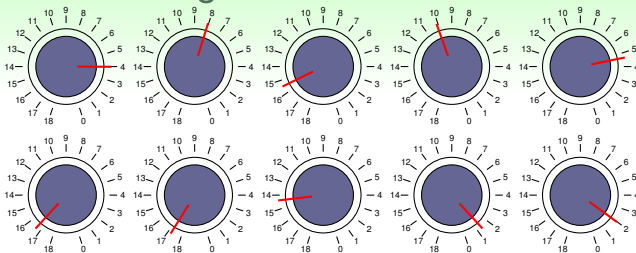
- can you generate  $(\text{sign}(a_1), \text{sign}(a_2), \dots, \text{sign}(a_{d+1}), \times)$ ? if so, what  $\mathbf{w}$ ?

$$\begin{aligned} \mathbf{w}^T \mathbf{x}_{d+2} &= a_1 \underbrace{\mathbf{w}^T \mathbf{x}_1}_0 + a_2 \underbrace{\mathbf{w}^T \mathbf{x}_2}_{\times} + \dots + a_{d+1} \underbrace{\mathbf{w}^T \mathbf{x}_{d+1}}_{\times} \\ &> 0 \text{ (contradiction!)} \end{aligned}$$

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

# Questions?

# 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**

這是 hypothesis set 的強度，  
它能 shatter 到  $N = d_{VC}$

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

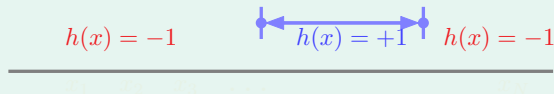
## Two Old Friends

Positive Rays ( $d_{VC} = 1$ ) ← 很弱的 hypothesis



free parameters:  $a$

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



free parameters:  $\ell, r$

practical rule of thumb:

$d_{VC} \approx \# \text{free parameters}$  (but not always, e.g.,  
mystery about deep learning models)

$M$  and  $d_{VC}$ 

copied from Lecture 3 :-)

- 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!,  $E_{in} \approx E_{out}$   
 $\mathbb{P}[\mathbf{BAD}] \leq 2 \cdot M \cdot \exp(\dots)$
- ② No!, too few choices

↓ 取代  $M$

small  $d_{VC}$ 

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

large  $M$ 

- ① No!,  $E_{in} \neq E_{out}$   
 $\mathbb{P}[\mathbf{BAD}] \leq 2 \cdot M \cdot \exp(\dots)$
- ② Yes!, many choices  $E_{in} \downarrow$

large  $d_{VC}$ 

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

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

# Questions?

# VC Bound Rephrase: Penalty for Model Complexity

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

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

## Rephrase

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

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

$$\frac{\delta}{4(2N)^{d_{VC}}} = \exp\left(-\frac{1}{8}\epsilon^2 N\right)$$

$$\ln\left(\frac{4(2N)^{d_{VC}}}{\delta}\right) = \frac{1}{8}\epsilon^2 N$$

$$\sqrt{\frac{8}{N} \ln\left(\frac{4(2N)^{d_{VC}}}{\delta}\right)} = \epsilon$$

# VC Bound Rephrase: Penalty for Model Complexity

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

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

## Rephrase

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

generalization error

gen. error  $|E_{in}(g) - E_{out}(g)| \leq \sqrt{\frac{8}{N} \ln \left( \frac{4(2N)^{d_{VC}}}{\delta} \right)}$

$$E_{in}(g) - \sqrt{\frac{8}{N} \ln \left( \frac{4(2N)^{d_{VC}}}{\delta} \right)} \leq E_{out}(g) \leq E_{in}(g) + \sqrt{\frac{8}{N} \ln \left( \frac{4(2N)^{d_{VC}}}{\delta} \right)}$$

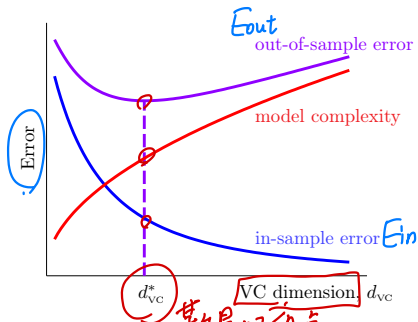
$\underbrace{\sqrt{\dots}}_{\Omega(N, \mathcal{H}, \delta)}$  : penalty for **model complexity**  
 $\nwarrow$   $E_{out}$  最壞的情況



# THE VC Message

with a high probability,

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



- $d_{\text{VC}} \uparrow$ :  $E_{\text{in}} \downarrow$  but  $\Omega \uparrow$ .
- $d_{\text{VC}} \downarrow$ :  $\Omega \downarrow$  but  $E_{\text{in}} \uparrow$ .
- best  $d_{\text{VC}}^*$  in the middle

最好的点

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  ~~$N \geq 2$~~ ,  $d_{VC} \geq 2$

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

$P[\text{BAD}] = 10\%$

2D perception hypothesis set.

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

# of data	$N$	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  $N \approx 10,000 d_{VC}$  in theory

100,000 e.g. 2D perception 要 30K points.

← BAD 發生的機率很小

practical rule of thumb:

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

這代表 VC bound 很鬆

Why?

- Hoeffding for unknown  $E_{\text{out}}$  *grow function ← 任何資料*
- $m_{\mathcal{H}}(N)$  instead of  $|\mathcal{H}(\mathbf{x}_1, \dots, \mathbf{x}_N)|$
- $N^{d_{\text{VC}}}$  instead of  $m_{\mathcal{H}}(N)$
- union bound on worst cases

any distribution, any target

'any' data

'any'  $\mathcal{H}$  of same  $d_{\text{VC}}$

any choice made by  $\mathcal{A}$

— **but hardly better, and 'similarly loose for all models'**

\* 不要一味追求 hypothesis set complexity

**philosophical message** of VC bound  
important for improving ML

# Questions?

# Summary

## 1 When Can Machines Learn?

### Lecture 3: Feasibility of Learning

## 2 Why Can Machines Learn?

### Lecture 4: Theory of Generalization

- Effective Number of Lines
  - Effective Number of Hypotheses
  - Break Point
  - Definition of VC Dimension
  - VC Dimension of Perceptrons
  - Physical Intuition of VC Dimension
  - Interpreting VC Dimension
- **next: beyond VC theory, please :-)**