

# Lineare Algebra

Lp norm

$$\|x\|_p = \left( \sum_{i=1}^n |x_i|^p \right)^{1/p}$$

$$\begin{aligned} \|x\|_1 &= (|x_1| + |x_2| + |x_3| + \dots + |x_n|)^1 \\ \|x\|_2 &= (\|x_1\|^2 + \|x_2\|^2 + \|x_3\|^2 + \dots + \|x_n\|^2)^{1/2} \\ \|x\|_\infty &= \max\{|x_1|, |x_2|, \dots, |x_n|\} \\ \Rightarrow \|x\|_\infty &= \max\{\|x_1\|^p, \|x_2\|^p, \dots, \|x_n\|^p\} \\ \Rightarrow \|x\|_\infty &= \max\{\|x_1\|, \|x_2\|, \dots, \|x_n\|\} \\ \Rightarrow \|x\|_\infty &= \max\{\|x_1\|^p, \|x_2\|^p, \dots, \|x_n\|^p\} \\ \Rightarrow \|x\|_\infty &= \max\{\|x_1\|, \|x_2\|, \dots, \|x_n\|\} \end{aligned}$$

Frobenius Norm

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2}$$

$$\begin{aligned} \text{e.g. } A &= \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix} \\ \Rightarrow \|A\|_F &= \sqrt{5^2 + 6^2 + 7^2 + 8^2} = \sqrt{174} \end{aligned}$$

$$\|A\|_F = \sqrt{\text{tr}(AA^T)} \quad (\text{should be } AA^T)$$

$$\text{e.g. } AA^T = \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix} \begin{bmatrix} 5 & 7 \\ 6 & 8 \end{bmatrix}$$

$$= \begin{bmatrix} 61 & 63 \\ 63 & 113 \end{bmatrix}$$

$$\text{tr}(AA^T) = 174$$

Trace

$$\text{tr}(A) = \sum_{i=1}^n A_{ii}$$

e.g.

$$\begin{aligned} \text{tr}\left(\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}\right) &= 15 \\ \text{tr}\left(\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}\right) &= 6 \end{aligned}$$

$$\begin{aligned} \text{tr}(A^T) &= \text{tr}(A) \\ \text{tr}(A+B) &= \text{tr}(A) + \text{tr}(B) \\ \text{tr}(ABC) &= \text{tr}(CAB) = \text{tr}(BCA) \end{aligned}$$

Matrix (<sup>here</sup> normally no extensional matrix)

Orthogonal, orthonormal  
columns norm = 1

$$AA^T = I$$

$$\text{e.g. } \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}$$

$$\begin{bmatrix} \cos \theta & \sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \begin{bmatrix} \cos \theta & \sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}$$

$$= \begin{bmatrix} \cos^2 \theta + \sin^2 \theta & \cos \theta \sin \theta + \sin \theta \cos \theta \\ \sin \theta \cos \theta + \cos \theta \sin \theta & \sin^2 \theta + \cos^2 \theta \end{math>$$

$$= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

Eigen decomposition

$$- Av = \lambda v \quad \text{right eigen vector} \quad (\text{usually unitary})$$

$$- v^T A = \lambda v^T \quad \text{left eigen vector}$$

- n independent eigenvectors  
{v<sup>(1)</sup>, ..., v<sup>(n)</sup>}

$$\{v<sup>(1)</sup>, ..., v<sup>(n)</sup>\}$$

$$- v = [v<sup>(1)</sup>, ..., v<sup>(n)</sup>]$$

$$\lambda = [\lambda_1, \lambda_2, \dots, \lambda_n]^T$$

$$- A = V \text{diag}(\lambda) V^{-1}$$

- if A is symmetric then

$$\exists A = Q \Lambda Q^T$$

scaling v<sup>(i)</sup> by λ<sub>i</sub>

Notes:

$$\begin{aligned} \Delta f(x) &= x^T A x \quad \text{s.t. } Ax = b \\ \text{if } x \in \{v^{(1)}, \dots, v^{(n)}\} \text{ then } f(x) &= \lambda_i \end{aligned}$$

$$\Delta \lambda \approx 0, A \in \text{PD}$$

$$\lambda \approx 0, A \in \text{PSD} \quad x^T A x \geq 0$$

$$\lambda \approx 0, A \in \text{ND}$$

$$\lambda \approx 0, A \in \text{NSD}$$

Singular Value Decomposition

- All matrix  $\not\in$  SVD

$$- A = UDV^T$$

$$\begin{aligned} \text{more more more} \\ \text{left } U, V \text{ right } V^T \\ \text{singular } D \in \text{Diag} \quad (\text{singular values}) \end{aligned}$$

Moore-Penrose Pseudoinverse

$$- A = X$$

$$- X = \arg \min_x \|Ax - b\|_2$$

$$- \|Ax - b\|^2 = 2A^T A x - 2A^T b$$

$$A^T = (A x)^T A^T$$

$$A^T = (A^T x)^T A$$

$$A^T = A x$$

$$x = (A^T A)^{-1} A^T b$$

$$x = A^T b$$

$$A^T = (A^T b)^T A$$

$$A^T = A^T b$$

$$A^T$$



# Theory of Generalization

## FACT

- $\Delta$  with distribution in training data & testing data
- $\Rightarrow$  low training error  
     $\Downarrow$   
    low testing error

## Def

$\Delta$  Training error:

- $E_{tr}(h) = \frac{1}{N} \sum_{n=1}^N e(h(x_n), f(x_n))$
- where  $x_1, \dots, x_N$  sampled from  $D$
- $h$  is determined by  $x_1, \dots, x_N$

$\Delta$  Testing error:

- $E_{te}(h) = \frac{1}{N} \sum_{n=1}^N e(h(x_n), f(x_n))$
- where  $x_1, \dots, x_N$  sampled from  $D$

$h$  is independent from  $x_1, \dots, x_N$

$\Delta$  Generalization error  $E(h)$

- G. error = Test error (<sup>expected performance</sup>) on  $D$
- $E(h) = E_{x \sim D}[e(h(x), f(x))] = E_{te}(h)$

$\Delta$  Summary

if  $E(h) = 0$

then  $E(h) \approx E_{tr}(h)$   $\rightarrow$  How?  $\Leftarrow$  but  
 $\Downarrow$

$E_{tr}(h) \approx 0 \rightarrow$  Training

Q: How do we make sure

$E(h) \approx E_{tr}(h)$

thru this

FACT Hoeffding's inequality



$\Delta P[\text{pick red ball}] = \mu$

$P[\text{pick green ball}] = 1 - \mu$

$\rightarrow$  we DO NOT know  $\mu$

$\Delta$  by pick ball's independently  
we get fraction of  $V$

$\Delta V \rightarrow \mu?$   
perhaps

$\Delta$  Hoeffding's inequality

$$P[|V - \mu| > \epsilon] \leq 2e^{-2\epsilon^2 N}$$

annotation:  $V$  &  $\mu$  的差距, 越大  $\epsilon$  大的  $P$ , 越小  
 $\epsilon$ ?  $\rightarrow$  越  $2e^{-2\epsilon^2 N}$  小

$\Delta$  statement  $\mu = V$  is  
probably approximately correct  
(PAC!)

$$V = \mu$$

probably approximately correct

## FACT

$$\Delta P[|V - \mu| > \epsilon] \leq 2e^{-2\epsilon^2 N}$$

$$P[|V - \mu| > \epsilon] \leq 2e^{-2\epsilon^2 N}$$

- valid for  $N$

-  $\epsilon > 0$

- independent from  $\mu$   
(real probability)



$\Delta$  in learning :

- given a function  $h$

- we randomly draw  $x_1, \dots, x_n$   $\Downarrow$  independent

- generalization error

$$E(h) = E_{x \sim D}[h(x) \neq f(x)] \Leftrightarrow \mu$$

sample data error

$$E_{tr}(h) = \frac{1}{N} \sum_{n=1}^N [h(x_n) \neq y_n] \Leftrightarrow V$$

$$P[|V - \mu| > \epsilon] \leq 2e^{-2\epsilon^2 N}$$

$$P[|V - \mu| > \epsilon] \leq 2e^{-2\epsilon^2 N}$$

known

## FACT

$\Delta$  for each  $h$  -  $h$  is a hypothesis

$$P[|E_{tr}(h) - E(h)| > \epsilon] \leq 2e^{-2\epsilon^2 N}$$

$\Delta$  for all  $h$ .  $H$  is a hypothesis set

$$P[|E_{tr}(h_1) - E(h_1)| > \epsilon], \quad P[|V - \mu| > \epsilon] \leq 2e^{-2\epsilon^2 N}$$

$$P[|E_{tr}(h_2) - E(h_2)| > \epsilon], \quad P[|E_{tr}(h_k) - E(h_k)| > \epsilon] \leq 2e^{-2\epsilon^2 N}$$

$$\vdots \quad P[|E_{tr}(h_{1k}) - E(h_{1k})| > \epsilon]$$

$$\leq P[\sup_{h \in H} |E_{tr}(h) - E(h)| > \epsilon] \leq P[\sup_{h \in H} |E_{tr}(h) - E(h)| > \epsilon] \leq 2^{|\mathcal{H}|} e^{-2\epsilon^2 N}$$

$$\leq \sum_{m=1}^{|\mathcal{H}|} P[|E_{tr}(h_m) - E(h_m)| > \epsilon] \leq 2^{|\mathcal{H}|} e^{-2\epsilon^2 N}$$

$$\text{from } P(\bigcup_{i=1}^{\mathcal{H}} A_i) \leq \sum_{i=1}^{\mathcal{H}} P(A_i)$$

$\Delta$  summary

$$P[|E_{tr}(h) - E(h)| > \epsilon] \leq P[\sup_{h \in H} |E_{tr}(h) - E(h)| > \epsilon] \leq 2^{|\mathcal{H}|} e^{-2\epsilon^2 N}$$

$$P[|V - \mu| > \epsilon] \leq 2e^{-2\epsilon^2 N}$$

$$\Downarrow \quad P[|E_{tr}(h) - E(h)| > \epsilon] \leq 2e^{-2\epsilon^2 N}$$

$$P[|G_{tr}(h) - E(h)| > \epsilon]$$

$$\leq P[\sup_{h \in H} |G_{tr}(h) - E(h)| > \epsilon]$$

$$\leq 2^{|\mathcal{H}|} e^{-2\epsilon^2 N}$$

NOTE More on Hoeffding's inequality

$$\begin{array}{ccccccc} D_1 & D_2 & \dots & D_N & P[\text{BAD } D \text{ for } h_1] \leq \dots \\ h_1 & \text{Bad} & & \text{Bad} & P[\text{BAD } D \text{ for } h_2] \leq \dots \\ h_2 & \text{Bad} & & \text{Bad} & P[\text{BAD } D \text{ for } h_m] \leq \dots \\ \vdots & & & & & & \\ h_m & \text{Bad} & & \text{Bad} & & & \end{array}$$

inferred hypothesis  
(假設現在我認定的模型 i.e., 不是學出來的)

對應到我手上  
的資料  $D_1, D_2, \dots, D_N$ .  
inferred 的  $h$  on  $D$   
可能導致 "Bad"  
亦即  $E_{\text{in}}(h) \neq E_{\text{out}}(h)$

→ 答案是假的:  
bound by  
 $2e^{-2\epsilon^2 N}$

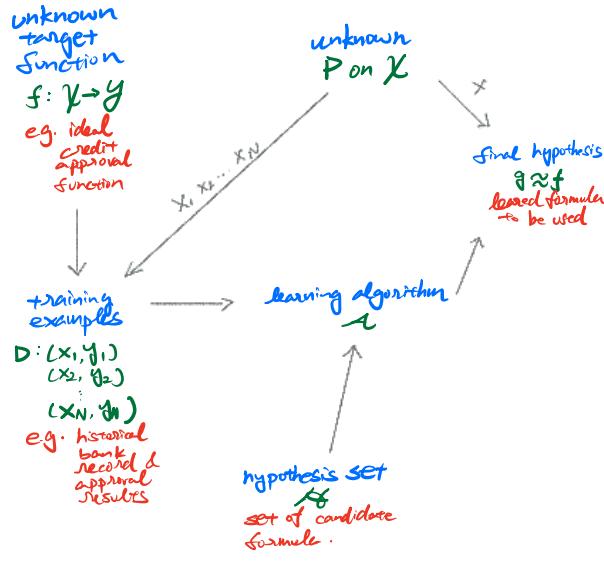
→ 之後以  $\epsilon$   
Hoeffding's ineq.  
就是為了量化  
 $P[\text{Bad}]$  機率  
多高?

∴  $P_D[\text{BAD } D]$

$$\begin{aligned} &= P[\text{BAD } D \text{ for } h_1 \text{ or } \text{BAD } D \text{ for } h_2 \dots \text{ or } \text{BAD } D \text{ for } h_m] \\ &\leq P[\text{BAD } D \text{ for } h_1] + P[\text{BAD } D \text{ for } h_2] + \dots + P[\text{BAD } D \text{ for } h_m] \\ &\quad \text{union bound} \\ &\leq 2M e^{-2\epsilon^2 N} = 2|\mathcal{H}| e^{-2\epsilon^2 N} \end{aligned}$$

- finite-bin version of Hoeffding
  - & hope...  $E_{\text{in}}(g) = E_{\text{out}}(g)$  is PAC.
- A will pick  $h_m$  w/ min.  $E_{\text{in}}(h_m)$  as  $g$

## △ Summary: statistical learning flow



& hope  $E_{\text{out}}(g) \approx E_{\text{in}}(g) \approx 0$

- for batch & supervised learning,  $g \approx f \Leftrightarrow E_{\text{out}}(g) \approx 0$  achieved through  $E_{\text{out}}(g) \approx E_{\text{in}}(g) \wedge E_{\text{in}}(g) \approx 0$
- ① can we make sure  $E_{\text{out}}(g) \approx E_{\text{in}}(g)$  G, Error
- ② can we make  $E_{\text{in}}(g)$  small enough Training

FACT  $|\mathcal{H}| = \infty$

△ Question: How do we deal with it?

- Small  $|\mathcal{H}|$ :  $P[\text{BAD}] \leq 2|\mathcal{H}| e^{-2\epsilon^2 N}$ 
  - small! great!
  - but  $|\mathcal{H}|$  too little
  - $E_{\text{in}}(g) \uparrow$
- large  $|\mathcal{H}|$ :  $E_{\text{in}}(g) \rightarrow 0$ 
  - small error! great!
  - but  $|\mathcal{H}|$  too large
  - $P[\text{BAD}] \uparrow$

註: 我們如何找  $g$  使得  $|\mathcal{H}|$  但能控制 the number!

FACT establish a finite quantity replace  $|\mathcal{H}|$

let  $|\mathcal{H}|$  replaced by  $M_H$

s.t.

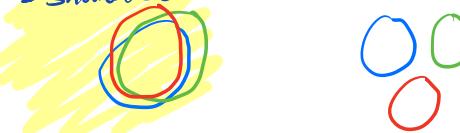
$$P[|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon] \leq 2 M_H e^{-2\epsilon^2 N}$$

FACT  $|\mathcal{H}|$  is over-estimated for BAD events

- BAD events  $B_m$ :  $|E_{\text{in}}(h_m) - E_{\text{out}}(h_m)| > \epsilon$
- overlapping for similar hypothesis  $h_1 \approx h_2$
- as  $\mathbb{E}_{\text{out}}(h_1) \approx \mathbb{E}_{\text{out}}(h_2)$

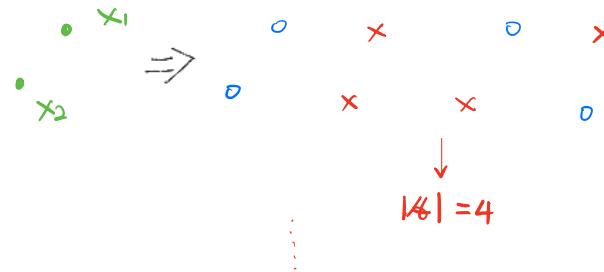
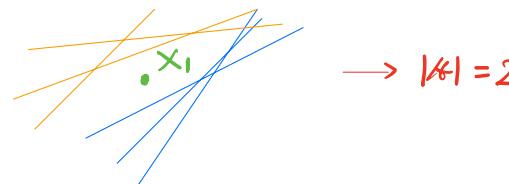
② for most  $D$   $E_{\text{in}}(h_1) \approx E_{\text{out}}(h_2)$

- should be instead of



→ So: can we group similar kinds?

e.g.  $H$  in  $\mathbb{R}^2$   $|\mathcal{H}| \rightarrow \infty$



$$\begin{array}{ll} N=3 & |\mathcal{H}|=8 \quad \text{but if on same line, different} \\ N=4 & |\mathcal{H}|=14 \quad \text{but if on same line, different} \end{array}$$

FACT observation: effective  $|\mathcal{H}| \leq 2^N$

perhaps can replace  $|\mathcal{H}|$  by effective  $|\mathcal{H}|$ ?  
need more rigorous proof

the split of set into two exclusive subsets

FACT Dichotomies: mini-hypotheses

△ limited hypothesis:  $H(x_1, x_2, \dots, x_N)$

△  $H(x_1, x_2, \dots, x_N)$ : depend on inputs  $(x_1, x_2, \dots, x_N)$

△ growth function:

remove dependence by taking max of all possible  $(x_1, x_2, \dots, x_N)$

$$M_H(N) = \max_{x_1, x_2, \dots, x_N \in X} |H(x_1, x_2, \dots, x_N)|$$

△ Finite, upper-bounded by  $2^N$

Q: How to calculate growth function

## FACT shattered

△ if  $m_H(N) = 2^N \Leftrightarrow$  exists  $N$  inputs that can be shattered.

△ e.g. convex set

## FACT summary of 4 growth function

- positive rays  $N+1$
- positive intervals  $C_2^{N+1} + 1$
- convex sets  $2^N$
- 2D perceptrons  $< 2^N$

polynomial good!

exponential bad!

## FACT Break point $\rightsquigarrow k$ paths,無法被 shattered

△ if no  $k$  inputs can be shattered by  $H$  call  $k$  a break point for  $H$

△  $m_H(k) < 2^k$  (in binary case)

△  $k+1, k+2, k+3 \dots$  are all break points

△ study minimum break point

e.g. linear case break point  $k=4$

note: 4↑無法被 shattered

## FACT conjecture:

△ no break point:  $m_H(N) = 2^N$

△ break point  $k$ :  $m_H(N) = O(N^{k-1})$

proof?

## FACT $m_H(N) \leq$ maximum possible $m_H(N)$ given $k$ $\leq \text{poly}(N)$

## FACT Bounding function 如果我有 break point $k$ upper bound is off?

• 從上面結論: 露出一絲亮光!

- 在 break point =  $k$

- 且  $k$  級度, 不能被 shattered (不是  $2^k$  分級) exists no  $2^k$  dichotomies!

• Now, for bounding function:

$$\max m_H(N) @ \text{break point } = k$$

## • $B(N, k)$

	PINK:					
	$\Delta B(4,3) = 2x + \beta$					
	1	2	3	4	5	6
1	1	2	2	2	2	2
2	1	3	4	4	4	4
3	1	4	7	8	8	8
4	1	5	11	15	16	16
5	1	6	16	26	31	32
6	1	7	22	42	57	63
:						

## • YELLOW

△  $k=1$ , max. dichotomies = 1

△  $k>N$ , max. dichotomies =  $2^N$

△  $k=N$ , max. dichotomies =  $2^{N-1}$

△  $k < N$ ? PINK

$$\bullet B(N, k) \leq \sum_{j=0}^{k-1} \binom{N}{j}$$

• Growth Function

$$\leq \text{Bounding Function} \quad \leq \sum_{j=0}^{k-1} \binom{N}{j} = \frac{N!}{(N-k)!k!} = \frac{N(N-1)\dots(N-k+1)}{k!}$$

$$\leq \frac{N^k}{(N-k)!} \leq \frac{N^k}{(N-k)^{N-k}} = \frac{N^k}{N^{N-k}}$$

$$\leq \frac{N^k}{N^k} = 1$$

只要有 break point  $k$  就是 1!

可被 polynomial 上限!

## PINK:

$$\Delta B(4,3) = 2x + \beta$$

Identify:

$\Delta B(4,3) = 2x + \beta$

if only consider  $x_1 x_2 x_3$

we have dichotomies  $\alpha + \beta$

$\therefore B(4,3) \rightarrow$  break pt = 3

$\therefore B(3,2) \rightarrow$  break pt = 2

i.e., cannot be shattered

i.e., cannot be shattered by  $2^3$

$\therefore \alpha + \beta \leq B(3,3)$

<math