

Linear Algebra

L-p norm

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

e.g.

$$\begin{aligned} \|x\|_1 &= (|x_1| + |x_2| + \dots + |x_n|) \\ &\vdots \\ n \times 1 &= (|x_1|^p + |x_2|^p + \dots + |x_n|^p)^{1/p} \\ \vdots \\ \|x\|_\infty &= \max\{|x_1|, \dots, |x_n|\} \\ &\vdots \\ \text{e.g. } \|x\|_2 &= \sqrt{\sum_{i=1}^n |x_i|^2} \\ \Rightarrow \|x\|_2 &= \sqrt{\sum_{i=1}^n (x_i^2 + y_i^2)} \\ \text{e.g. } \|x\|_1 &= \max\{|x_1|, \dots, |x_n|\} \\ \text{e.g. } \|x\|_0 &= \max\{1, \dots, 1\} \\ \text{e.g. } \|x\|_1 &= \max\{|x_1| + \dots + |x_n|\} \end{aligned}$$

Frobenius Norm

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

e.g.

$$\begin{aligned} A &= \begin{bmatrix} 2 & 6 \\ 7 & 8 \end{bmatrix} \\ \Rightarrow \|A\|_F &= \sqrt{2^2 + 6^2 + 7^2 + 8^2} = \sqrt{174} \\ \bullet \|A\|_F &= \sqrt{\text{tr}(AA^T)} \quad (\text{should be } AA^T) \\ \text{e.g. } &AA^T = \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix} \begin{bmatrix} 5 & 7 \\ 6 & 8 \end{bmatrix} \\ &= \begin{bmatrix} 61 & 63 \\ 63 & 68 \end{bmatrix} \\ \text{tr}(AA^T) &= 174 \end{aligned}$$

Trace

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

e.g.

$$\begin{aligned} \text{tr}\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix} &= 15 \\ \text{tr}\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix} &= 6 \\ \bullet \text{tr}(A^T) &= \text{tr}(A) \\ \bullet (A+B) &= \text{tr}(A) + \text{tr}(B) \\ \bullet (AB) &= \text{tr}(CAB) = \text{tr}(BCA) \end{aligned}$$

Matrix

$$\begin{aligned} \text{orthogonal} &\quad (\text{possibly no orthonormal}) \\ \bullet \text{Orthogonal, orthonormal} & \downarrow \\ AA^T &= I \quad (\text{columns norm = 1}) \\ \text{e.g. } & \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \\ & \begin{bmatrix} \cos \theta & 0 & 0 \\ 0 & \cos \theta & 0 \\ 0 & 0 & \cos \theta \end{bmatrix} \\ & \begin{bmatrix} \cos \theta & \cos \theta & \cos \theta \\ \cos \theta & \cos \theta & \cos \theta \\ \cos \theta & \cos \theta & \cos \theta \end{bmatrix} \\ & = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \end{aligned}$$

Eigen decomposition

$$\begin{aligned} \bullet AV &= \lambda V \quad (\text{left eigen}) \\ &\quad \text{eigenvectors (usually unit)} \\ \bullet V^T A = V^T \lambda V \quad (\text{left eigen}) \\ \bullet n \text{ independent eigenvectors} \\ &\quad \{V^{(1)}, \dots, V^{(n)}\} \\ &\quad \{v_1, \dots, v_n\} \\ \bullet V = [V^{(1)}, \dots, V^{(n)}] \\ &\quad \lambda = [\lambda_1, \lambda_2, \dots, \lambda_n]^T \\ \bullet A = V \text{ diag}(\lambda) V^{-1} \\ \bullet \text{if } A \text{ is symmetric then} \\ \quad \exists A = Q \Lambda Q^T \\ &\quad (\text{coding } \lambda^{(1)} \text{ by } \lambda_1) \end{aligned}$$

Notes:

$$\begin{aligned} \bullet f(x) &= x^T A x + b^T x + 1 \\ &\quad \text{if } x \in \{V^{(1)}, \dots, V^{(n)}\} \\ \text{then } f(x) &= \lambda_1 \end{aligned}$$

$$\begin{aligned} \bullet \lambda \neq 0, \quad A \in \text{PD} \\ \lambda \neq 0, \quad A \in \text{PSD} \quad x^T A x \geq 0 \\ \lambda \neq 0, \quad A \in \text{ND} \\ \lambda \neq 0, \quad A \in \text{NSD} \end{aligned}$$

Singular Value Decomposition

$$\begin{aligned} \bullet \text{All matrix fit SVD} \\ \bullet A = UDV^T \\ &\quad \text{more more more right singular vectors} \\ &\quad \text{U, V, } \in \text{orthogonal} \\ &\quad D \in \text{Diag} \quad (\text{singular values here}) \\ \bullet \text{More - Pseudo inverse} \\ \bullet A = X \\ \bullet X = \arg \min_X \|Ax - b\|_2 \\ \bullet \|Ax - b\|_2^2 = 2A^T Ax - 2A^T b \\ \bullet A = (A^T A)^{-1} A^T \end{aligned}$$

Matrix Derivatives

$$\begin{aligned} \bullet f \in \mathbb{R} \\ \bullet x \in \mathbb{R}^n \end{aligned}$$

$$\frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix} \quad \frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix}$$

$$Df(x) = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} & \dots & \frac{\partial f}{\partial x_n} \\ \frac{\partial f}{\partial x_2} & \frac{\partial f}{\partial x_1} & \dots & \frac{\partial f}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial x_n} & \frac{\partial f}{\partial x_n} & \dots & \frac{\partial f}{\partial x_1} \end{bmatrix}$$

$$J = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \dots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

$$\bullet f \in \mathbb{R}^m \\ \bullet x \in \mathbb{R}^n$$

$$\frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \dots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

$$Df(x) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \dots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

$$f = \|x - y\|^2 \quad Wx - y \in \mathbb{R}^m$$

$$\frac{\partial f}{\partial x} = \frac{\partial}{\partial x} ((Wx - y)^T (Wx - y))$$

$$= \frac{\partial}{\partial x} (W^T W x^2 - 2W^T y x + y^T y)$$

$$= 2W^T W x - 2W^T y$$

$$= 2x^T (Wx - y)$$

$$\text{argmin}_x \|x - y\|^2$$

$$\Rightarrow \text{let } 2x^T (Wx - y) = 0$$

$$\Rightarrow \cancel{2x^T W} x = \cancel{2x^T y}$$

$$W = (X^T X)^{-1} X^T y$$

$$(W = X^T + y)$$

via SVD

$$\bullet \text{minimize}_w \|Wx - y\|$$

$$\bullet x = UZV^T$$

$$\begin{aligned} U^T U &= I \\ V^T V &= I \\ Z^T Z &= \text{diag}[v_1, v_2, \dots, v_n, 0, \dots, 0] \end{aligned}$$

$$\bullet W = X^T + y$$

$$= VZ^T + U^T + \cancel{y} \quad \text{[cancel } U^T \text{ and } Z^T \text{]} \quad \text{[cancel } U^T \text{ and } Z^T \text{]}$$

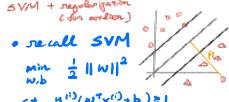
$$= \|VZ^T + U^T\|$$

$$\min_w \|Wx - y\| \equiv \|UZ^T + U^T\|$$

$$\min_w \|UZ^T + U^T\| \equiv \|U\| \|Z^T + U^T\|$$

$$\min_w \|UZ^T + U^T\| \equiv \|U\| \|Z^T\|$$

$$\min_w \|UZ^T + U^T\| \equiv \|U\| \|Z\|$$



recall SVM
 $\min_{w,b} \frac{1}{2} \|w\|^2$
 $s.t. g^{(1)}(N^T x^i + b) \geq 1$
+ add L1 regularization
 $\min_{w,b} \frac{1}{2} \|w\|^2 + C \sum_i \delta_i$
 $s.t. g^{(1)}(N^T x^i + b) \geq 1 - \delta_i$
 $\delta_i \geq 0$
 $\Leftrightarrow \min_w \frac{1}{2} \|w\|^2 + C \sum_i \delta_i$
 $s.t. g^{(1)}(N^T x^i + b) \geq 1 - \delta_i$
 $(we can always do w/ \delta_i)$
become $\|b\|_1 + \frac{1}{2} \|w\|_2^2$

Non-linear SVM
+ some class is linear separable
+ project pth dimension transformation to a space that is linear separable
- $\tilde{x} = \phi(x)$
import features
ambiguity
also allow for more complex decision boundary
- $\langle a, b \rangle$ replaced by
 $K(a, b)$ for non-linear transformation

Model Class Classification
+ 1 v. All
+ build k-class classifier
- 1st class
- ...
- k-th class
+ decision rule for k-th class
+ 1 v. 1
+ L(L-1) binary classifiers
+ S-t
- 1st, 2nd class
- ...
- L-th class
+ multi-class loss function
+ softmax
+ cross entropy
- $-\log(\frac{e^{x_i}}{\sum_j e^{x_j}})$
+ optimization problem
 $\min_w \frac{1}{n} \sum_i -\log(\frac{e^{x_i^T w}}{\sum_j e^{x_j^T w}}) + \lambda \|w\|_2^2$

Optimization
+ $f(x_1 + (1-\theta)x_2) \leq \theta f(x_1) + (1-\theta)f(x_2)$
Convexity
+ 1st order
+ convex function
FACT
 $f(y) \geq f(x) + \nabla f(x)^T (y-x)$
 $f(x) : R^n \rightarrow R$
 $\nabla f(x) \in R^{n \times n}$
 $x \in I\!O, \mathbb{R}^n$

FACT 1st order condition
 $f'(x) \geq f(x) + \nabla f(x)^T (y-x)$
FACT 2nd order condition
 $\nabla^2 f(x) \geq 0$
 $\nabla^2 f(x) \geq 0 \quad \forall x \in \mathbb{R}^n$

e.g.
- $f(w) = g(x^T w + y)$
 $g(z) = \text{convex}$
how about $f(w)$?
proof:
 $\Rightarrow f$ is convex iff.
 $f(\theta w + (1-\theta)w_0) \leq \theta f(w) + (1-\theta)f(w_0)$
 $\Rightarrow g(x^T [(\theta w + (1-\theta)w_0)] + y) \leq \theta g(x^T w + y) + (1-\theta)g(x^T w_0 + y)$
 $\Rightarrow g(\theta x^T w + (1-\theta)x^T w_0 + y + (1-\theta)y) \leq \theta g(x^T w + y) + (1-\theta)g(x^T w_0 + y)$
 $\leq \theta g(x^T w + y) + (1-\theta)(x^T w_0 + y)$
 $\leq \theta g(x^T w + y) + (1-\theta)g(x^T w_0 + y)$
 $\Rightarrow g(\theta x^T w + (1-\theta)x^T w_0) \leq \theta g(x^T w) + (1-\theta)g(x^T w_0)$
 $\Rightarrow g(\theta w + (1-\theta)w_0) \leq \theta g(w) + (1-\theta)g(w_0)$
as $g(z)$ is convex
 $f(w)$ is convex

Lipschitz Continuity & Smoothness

- L is a constant
- L-Lipschitz continuous:
 $\|f(x) - f(x_0)\|_2 \leq L \|x - x_0\|_2$
- L-Lipschitz smooth:
 $\|\nabla f(x) - \nabla f(x_0)\|_2 \leq L \|x - x_0\|_2$

FACT $\|\nabla f(x)\|_2 \leq L$ if f is L-smooth

proof
mean value theorem, $\exists c \in I\!O, x$

$$f(x) = f(x_0) + \frac{\partial f(x)}{\partial x}(x-x_0)$$

$$\therefore \exists \theta \in I\!O, x \in x_0$$

$$\nabla f(x) - \nabla f(x_0) = \nabla f(x) - \nabla f(x_0)$$

L-smooth property:

$$\|\nabla f(x) - \nabla f(x_0)\|_2 \leq L \|x - x_0\|_2$$

$$\therefore \|\nabla f(x)(x-y)\|_2 \leq L \|x-y\|_2$$

from Cauchy-Schwarz inequality

$$\|\nabla f(x)(x-y)\|_2 \leq \|\nabla f(x)\|_2 \|x-y\|_2$$

$$\therefore \|\nabla f(x)(x-y)\|_2 \leq L \|x-y\|_2$$

$\nabla^2 f(x)$ is symmetric

$$\|\nabla^2 f(x)\|_2 \leq L$$

$$\Rightarrow \|\nabla f(x)\|_2 \leq L$$

FACT $\|\nabla f(x)\|_2 = \sqrt{f(x)^T \nabla^2 f(x) f(x)}$

$$\Delta f(x) = x + z(y-x) \quad \therefore f(z) = f(x)$$

$$g(x) = f(z(x))$$

$$\therefore g(x) = f(x) + \nabla f(x)^T (y-x)$$

+ Fundamental theorem

$$g(x) - g(x_0) = \int_{x_0}^x \nabla g(t) dt$$

$$\therefore g(x) - g(x_0) = \int_{x_0}^x \nabla g(t) dt - \int_{x_0}^x \nabla f(t) dt$$

$$\therefore \|\nabla g(x) - \nabla f(x)\|_2 \leq \int_{x_0}^x \|\nabla g(t) - \nabla f(t)\|_2 dt$$

$$= \int_{x_0}^x \|\nabla^2 f(t)\|_2 dt \leq L \int_{x_0}^x dt = L \|x-x_0\|_2$$

$$\therefore \|\nabla g(x) - \nabla f(x)\|_2 \leq L \|x-x_0\|_2$$

$$= \|\nabla f(x) - \nabla f(x_0)\|_2 \leq L \|x-x_0\|_2$$

$$\therefore \|\nabla f(x)\|_2 \leq L \|x-x_0\|_2$$

$$= \frac{1}{2} \|\nabla^2 f(x)\|_2 \|x-x_0\|_2^2$$

$$\Delta \|\nabla f(x)\|_2 \leq \frac{1}{2} \|\nabla^2 f(x)\|_2 \|x-x_0\|_2^2$$

FACT gradient descent converges

If $\alpha < \frac{1}{L}$ $\alpha < \frac{1}{L}$

proof:
 $\Delta f(x) = f(x) - f(x_0)$

$$= f(x) - f(x) - \nabla f(x)^T (x-x_0)$$

$$= -\nabla f(x)^T (x-x_0)$$

$$\leq \frac{1}{2} \|x-x_0\|_2^2$$

$$\Delta \|\nabla f(x)\|_2 \leq \frac{1}{2} \|x-x_0\|_2^2$$

$$\therefore f(x) - f(x_0) \leq \frac{1}{2} \|x-x_0\|_2^2$$

$$\therefore f(x) - f(x_0) \leq \frac{1}{2} \|\nabla f(x)\|_2^2$$

Theory of Generalization

FACT

- with distribution in training data & testing data
⇒ low training error
 ||
 low testing error

Def

Training error:

$$E_{\text{tr}}(h) = \frac{1}{N} \sum_{n=1}^N e(h(x_n), f(x_n))$$

where x_1, \dots, x_N sampled from D

h is determined by x_1, \dots, x_N

Testing error:

$$E_{\text{te}}(h) = \frac{1}{N} \sum_{n=1}^N e(h(x_n), f(x_n))$$

where x_1, \dots, x_N sampled from D

h is independent from x_1, \dots, x_N

Generalization error $E(h)$

- G.error = Test error (on D)

$$- E(h) = E_{\text{tr}}[e(h(x), f(x))] = E_{\text{te}}(h)$$

Summary

if $E(h) = 0$

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

else

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

Q: How do we make sure

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

throw this

FACT Hoeffding's inequality



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

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

→ we DO NOT know μ

△ by pick ball's independently
we get fraction of V

△ $V \rightarrow \mu$?
perhaps

Hoeffding's inequality

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

note: V & μ 的差距, 越大 P , 很小
多大? → $2e^{-2\epsilon^2 N}$ 遠小

△ 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 μ
(create probability)



in learning :

- given a function h
- we randomly draw x_1, \dots, x_n
- generalization error

$$E(h) = E_{\text{tr}}[e(h(x), f(x))] \Leftrightarrow \mu$$

sample data error

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

known

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

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

FACT

△ for each h - h is a hypothesis

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

△ for all h , \mathcal{H} is a hypothesis set

$$\begin{aligned} & P[|E_{\text{tr}}(h_1) - E(h_1)| > \epsilon], \quad P[|V - \mu| > \epsilon] \leq 2e^{-2\epsilon^2 N} \\ & P[|E_{\text{tr}}(h_2) - E(h_2)| > \epsilon], \quad \downarrow \\ & \vdots \quad P[|E_{\text{tr}}(h_{10}) - E(h_{10})| > \epsilon] \leq 2e^{-2\epsilon^2 N} \\ & \leq P[\sup_{h \in \mathcal{H}} |E_{\text{tr}}(h) - E(h)| > \epsilon] \leq P[\sup_{h \in \mathcal{H}} |E_{\text{tr}}(h) - E(h)| > \epsilon] \leq 2^{|\mathcal{H}|} e^{-2\epsilon^2 N} \\ & \leq \sum_{m=1}^{|\mathcal{H}|} P[|E_{\text{tr}}(h_m) - E(h_m)| > \epsilon] \leq 2^{|\mathcal{H}|} e^{-2\epsilon^2 N} \end{aligned}$$

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

summary

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

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

$$\downarrow$$

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

$$\downarrow$$

$$P[|G(h) - E(h)| > \epsilon]$$

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

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

NOTE More on Hoeffding's inequality

$$\begin{array}{cccc} D_1 & D_2 & \dots & D_N \\ h_1 & \text{Bad} & & \text{Bad} \\ h_2 & \text{Bad} & \dots & \text{Bad} \\ \vdots & & & \vdots \\ h_m & \text{Bad} & & \text{Bad} \end{array}$$

$P_D[\text{BAD } D \text{ for } h_1] \leq \dots$
 $P_D[\text{BAD } D \text{ for } h_2] \leq \dots$
 \dots
 $P_D[\text{BAD } D \text{ for } h_m] \leq \dots$

↓
 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"
 即即 $E_{in}(h) \neq E_{out}(h)$

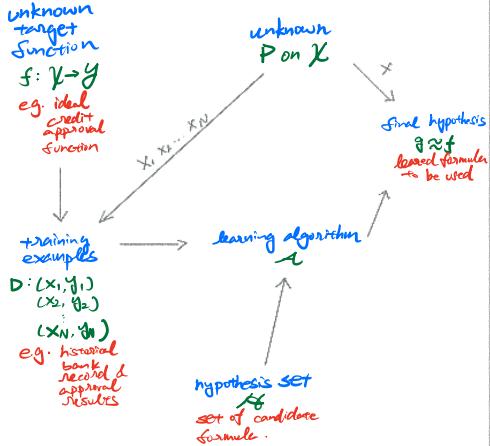
→ 該點要
 Hoeffding's inequality
 就是為了量化
 $P[\text{Bad}]$ 机率
 多高?

→ 答案是依循:
 bounded by
 $2e^{-2\epsilon^2 N}$
 → 不過剛剛是 1 個 h 啦...
 我 also 在找 h ...
 upper bound 是啥?

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

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

△ Summary: statistical learning flow



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

- for batch & supervised learning, $g \approx f \Leftrightarrow E_{out}(g) \approx 0$ achieved through $E_{out}(g) \approx E_{in}(g) \& E_{in}(g) \approx 0$
- ① can we make sure $E_{out}(g) \approx E_{in}(g)$ G. $E_{in}(g)$
- ② can we make $E_{in}(g)$ 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
 $E_{in}(g) \uparrow$
 - large $|\mathcal{H}|$:
 $E_{in}(g) \rightarrow 0$
small error! great!
 but $|\mathcal{H}|$ too large
 $P[\text{BAD}] \uparrow$
- 註: 我們如何找 finite $|\mathcal{H}|$? but can control the number!

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

let $|\mathcal{H}|$ replaced by M
 s.t.

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

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

- BAD events $B_m : |E_{in}(h_m) - E_{out}(h_m)| > \epsilon$
- overlapping for similar hypothesis $h_1 \approx h_2$
- as ① $E_{out}(h_1) \approx E_{out}(h_2)$

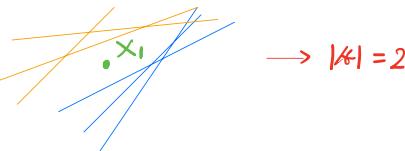
② for most D $E_{in}(h_1) \approx E_{out}(h_2)$

- should be instead of

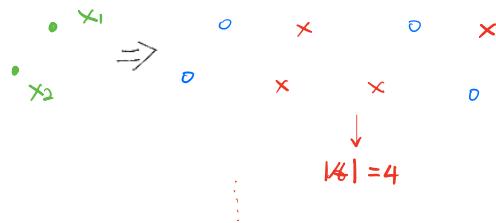


- So: can we group similar kinds?

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



$$\rightarrow |\mathcal{H}| = 2$$



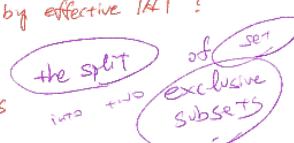
$$N=3 \quad |\mathcal{H}|=8 \quad \text{but if on same line, different}$$

$$N=4 \quad |\mathcal{H}|=14 \quad \text{but if on same line, different}$$

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

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

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 $\rightsquigarrow k$ 開始無法被 shattered

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

△ $m_H(k) < 2^k$ (in wrong 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 在哪?

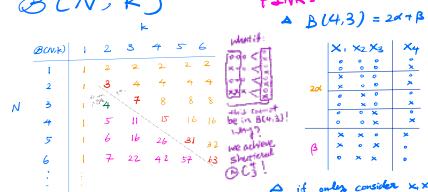
○ 從上面結論:

- 在 break point $= k$ \rightsquigarrow 跑出一條直線!
- 但 k 級度, 不能被 shattered (不是 2^k 組合)
 ○ 有 2^k dichotomies!

○ Now, for bounding function:

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

• $B(N, k)$



YELLOW

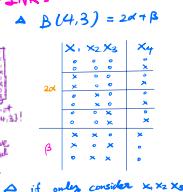
- △ $k=1$, max. dichotomies = 1
- △ $k=2$, max. dichotomies = 2^2
- △ $k=N$, max. dichotomies = 2^{N-1}
- △ $k=N$? PINK

• $B(N, k) \leq \frac{k!}{(k-1)!} \binom{N}{k}$

$$\begin{aligned} &\triangle \text{Growth Function} \quad \binom{N}{k} \leq \binom{N}{k-1} \cdot \binom{N}{k} = \binom{N}{k-1} \\ &\triangle \text{Bounding Function} \quad \approx \frac{1}{2} N! \\ &\triangle \frac{1}{2} N! \leq \frac{1}{2} \binom{N}{k} \end{aligned}$$

只要 k 有 break point
有異志!
可被 polynomial 上界住!

PINK:



△ if only consider $x_1 x_2 x_3$

we have dichotomies $2^3 = 8$

△ $B(4,3) = 2^3 = 8$

△ $B(4,3) \rightarrow$ break pt = 3

△ $B(3,3) \rightarrow$ break pt = 3

△ cannot be shattered

△ we cannot do $2^3 = 8$

△ $a+b \leq B(3,3)$

△ $a+b \leq B(3,3)$