

Linear Algebra

L-p norm

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

e.g.

$$\begin{aligned} \|x\|_1 &= (|x_1| + |x_2| + \dots + |x_n|)^1 \\ &= |x|_1 = \max\{|x_1|, \dots, |x_n|\} \\ &\quad \text{if } |x|_1 = \max\{|x_1|, \dots, |x_n|\} \\ &\quad \Rightarrow \frac{\partial}{\partial x_i} (|x_1|^1 + |x_2|^1 + \dots + |x_n|^1) \\ &\quad \quad \text{for } x_i = \max\{|x_1|, \dots, |x_n|\} \\ &\quad \Rightarrow \frac{\partial}{\partial x_i} (|x_i|^1) \\ &\quad \Rightarrow \frac{\partial}{\partial x_i} (|x_i|) = |x_i| \\ &\quad \text{from } (1) \cdot (2) \\ &\therefore \|x\|_1 = \max\{|x_1| + \dots + |x_n|\} \end{aligned}$$

Frobenius Norm

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

e.g.

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

Trace

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

e.g.

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

Matrix

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

Eigen decomposition

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

Notes:

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

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

Singular Value Decomposition

All matrix fit SVD

$$A = UDV^T$$

more more more right singular vectors
U, V, E orthogonal
D = Diag (singular values)

More - Pseudo inverse

$$\begin{aligned} A &= X \\ &= X \hat{X}^{-1} \quad \text{minimize}_{\hat{X}} \|Ax - b\|_2 \\ &= X(X^T X)^{-1} X^T b \\ &= X(X^T X)^{-1} X^T b \end{aligned}$$

Matrix Derivatives

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

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

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

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

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

$$\frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial f}{\partial x_{11}} & \frac{\partial f}{\partial x_{12}} & \dots & \frac{\partial f}{\partial x_{1n}} \\ \frac{\partial f}{\partial x_{21}} & \frac{\partial f}{\partial x_{22}} & \dots & \frac{\partial f}{\partial x_{2n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial x_{n1}} & \frac{\partial f}{\partial x_{n2}} & \dots & \frac{\partial f}{\partial x_{nn}} \end{bmatrix}$$

$$\bullet f = \|x - y\|^2 \quad x, y \in \mathbb{R}^n$$

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

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

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

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

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

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

$$x = (x^T)^{-1} y^T$$

$$(x = X^T y)$$

via SVD

$$\text{minimize}_w \|x - y\|$$

$$= x - UZV^T$$

$$U^T U = I$$

$$V^T V = I$$

$$Z = \text{diag} \{ \sigma_1, \sigma_2, \dots, \sigma_n, 0, \dots \}$$

$$\therefore w = X^T y$$

$$= Y Z^T U^T$$

$$= \cancel{Y} Z^T \cancel{U^T}$$

$$= Z^T - U^T y$$

$$= Z^T - U^T y$$

$$\Rightarrow \cancel{Z^T} - \cancel{U^T} y = 0$$

$$\Rightarrow \cancel{Z^T} = \cancel{U^T} y$$

$$\Rightarrow Z^T Z = Z^T U^T y$$

$$\Rightarrow Z = Z^T U^T y$$

$$\Rightarrow Z = Z^T U^T y$$

$$\Rightarrow Z = V^T W^T$$

$$\Rightarrow V W^T = Z^T U^T y$$

$$\Rightarrow W = V^T U^T y$$

$$\Rightarrow W = V^T U^T y$$

$$\text{SVD calculation}$$

recall eigen decomposition

$$\text{given } A = \begin{bmatrix} 26 & 18 \\ 18 & 74 \end{bmatrix}$$

$$\det(A - \lambda I) = \det \begin{bmatrix} 26-\lambda & 18 \\ 18 & 74-\lambda \end{bmatrix}$$

$$\Rightarrow \lambda^2 - 100\lambda + 1600 = (2 - \lambda)(2 - 74)$$

$$\Delta = \lambda^2 - 100\lambda + 1600 = 200$$

$$A - \lambda I = \begin{bmatrix} 6 & 18 \\ 18 & 56 \end{bmatrix}$$

$$\lambda = 80$$

$$A - \lambda I = \begin{bmatrix} -54 & 18 \\ 18 & -6 \end{bmatrix}$$

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

$$= \begin{bmatrix} -3 & 1 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} 20 & 0 \\ 0 & 80 \end{bmatrix} \begin{bmatrix} -3 & 1 \\ 1 & 3 \end{bmatrix}$$

$$= \begin{bmatrix} -3 & 1 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} 20 & 0 \\ 0 & 80 \end{bmatrix} \begin{bmatrix} -3 & 1 \\ 1 & 3 \end{bmatrix}$$

$$= \begin{bmatrix} 50 & 0 \\ 0 & 180 \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{50} & 0 \\ 0 & \sqrt{180} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

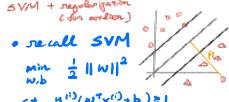
$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$= \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{6} \end{bmatrix}$$

$$$$



recall SVM
 $\min_{w,b} \frac{1}{2} \|w\|^2$
 $s.t. g^{(1)}(N^T x^i + b) \geq 1$
+ add L1 regularization
 $\min_{w,b} \frac{1}{2} \|w\|^2 + C \sum_i \xi_i$
 $s.t. g^{(1)}(N^T x^i + b) \geq 1 - \xi_i$
 $\xi_i \geq 0$
 $\Leftrightarrow \min_w \frac{1}{2} \|w\|^2 + C \sum_i \xi_i$
 $s.t. g^{(1)}(N^T x^i + b) \geq 1 - \xi_i$
 $(we can always do w/ \xi_i)$
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
- $f = \phi(x)$
import features
ambiguity
also allow for more complex boundaries
+ Kernel Trick
- $\langle a, b \rangle$ replaced by
 $K(a, b)$ for nonlinear 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
 $f'(x) \geq 0$
 $f''(x) \geq 0$
 $f''(x) > 0$
 $f''(x) \geq 0$
 $x \in \mathbb{R}$

FACT 1st order condition
 $f'(x) \geq f(x) + f'(x)(x-x)$
FACT 2nd order condition
 $f''(x) \geq f''(x) + 2f''(x)(x-x)$
 $f''(x) \geq 0$
 $f''(x) \geq 0 \quad \forall x \in \mathbb{R}^n$

e.g.
- $f(w) = g(x^T w + y)$
 $g(z) = \text{convex}$
how about $f(w)$?
proof:
 $\Rightarrow f$ is convex i.e.
 $f(\theta w + (1-\theta)w_0) \leq \theta f(w) + (1-\theta)f(w_0)$
 $\Rightarrow g(x^T [\theta w + (1-\theta)w_0] + y) \leq \theta g(x^T w + y) + (1-\theta)g(x^T w_0 + y)$
 $\Rightarrow g(\theta x^T w + (1-\theta)x^T w_0 + y + (1-\theta)y) \leq \theta g(x^T w + y) + (1-\theta)g(x^T w_0 + y)$
 $\leq \theta g(x^T w + y) + (1-\theta)(x^T w_0 + y)$
 $\leq \theta g(x^T w + y) + (1-\theta)g(x^T w_0 + y)$
 $\Rightarrow g(\theta w + (1-\theta)w_0) \leq \theta g(w) + (1-\theta)g(w_0)$
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:
 $\|f'(x) - f'(x_0)\|_2 \leq L \|x - x_0\|_2$

FACT $f'(x) \leq L$ if f is L-smooth

proof
mean value theorem, $\exists c \in [a,b]$

$$f(x) - f(a) = \frac{d}{dx} f(c)(x-a)$$

$$\therefore \exists c \in [a,b] s.t.$$

$$\sqrt{f'(c)(x-a)} = |f(x) - f(a)|$$

L-smooth property:

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

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

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

from Cauchy-Schwarz inequality

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

$$\therefore \|f'(x)\|_2 \leq L$$

$f'(x)$ is symmetric

$$\|f'(x)\|_2 \leq L$$

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

$$\therefore \|f'(x)\|_2 \leq L$$

Theory of Generalization

FACT

- Δ with distribution in training data & testing data
 \Rightarrow low training error
 $\qquad \qquad \qquad \#$
 \Rightarrow 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 (L on D)

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

Summary

if $E(h) = 0$

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

ie

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

Q: How do we make sure

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

throw this

FACT Hoeffding's inequality



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

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

→ we DO NOT know u

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

△ $V \rightarrow u$?
 perhaps

Hoeffding's inequality

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

note: u & u 的差距, 越大 P , 很小
 ϵ +? → $\epsilon 2e^{-2\epsilon^2 N}$ 遠小

△ statement $u = V$ is
 $\text{probably approximately correct}$
 $(PAC!)$

$V = u$ probably approximately correct

FACT

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

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

- valid for N
- $\epsilon > 0$
- independent from u
 $\text{create probability}$



in learning :

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

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

sample data error

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

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

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

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 , H is a hypothesis set

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

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

summary

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

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

$$\downarrow$$

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

$$\downarrow$$

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

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

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

NOTE More on Hoeffding's inequality

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

infared hypothesis
(假設現在
我選出的
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 .
infared 的 h on D
可能導致 "Bad"
即 $E_{in}(h) \neq E_{out}(h)$

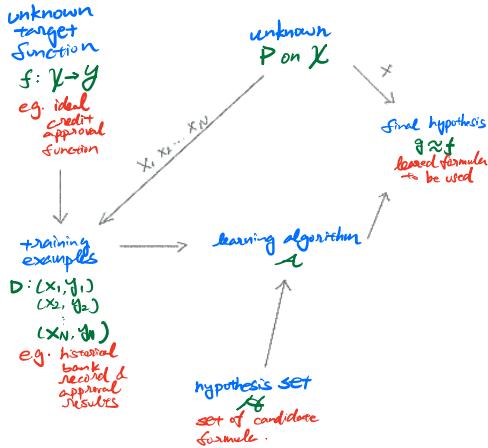
之所以要
Hoeffding's inequality
就是為了量化
 $P[\text{Bad}]$ 机率
多高?

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

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

- 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 $|\mathcal{H}|$: $P[\text{BAD}] \leq 2|\mathcal{H}| e^{-2\epsilon N}$
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.

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

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 $E_{out}(h_1) \approx E_{out}(h_2)$

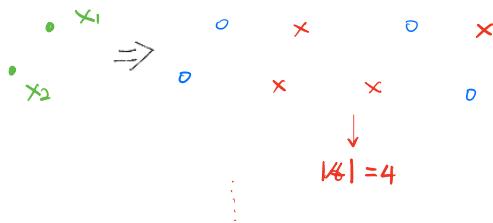
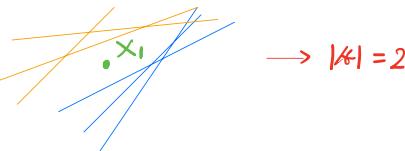
② for most D $E_{in}(h_1) \approx E_{out}(h_2)$

- should be instead of



- So: can we group similar kinds?

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



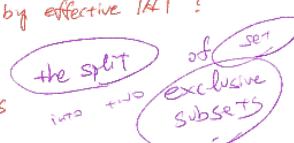
$N=3 \quad |\mathcal{H}|=8 \quad$ but if on same line, different

$N=4 \quad |\mathcal{H}|=14 \quad$ 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

