

LR { curve fitting  
classification , LR also applied on non-linear thing

## 6 Data

• Diff from Algo ML is data-driven process

• In ML, express everything as **points**.

• 人  $\rightarrow$  (height, weight, age)  
 ◦  $28 \times 28$  pixels  $\rightarrow$  4 dim

• Simple model Model  $\leftarrow$  solve params



• fit the line  $ax+b \Rightarrow \begin{cases} 7 = 4a+b \\ 2 = a+b \end{cases}$

• Input Data  $A = \begin{bmatrix} 4 & 1 \\ 1 & 1 \end{bmatrix}$  target  $\vec{b} = \begin{bmatrix} 7 \\ 2 \end{bmatrix}$ , Params  $\vec{x} = \begin{bmatrix} a \\ b \end{bmatrix}$

•  $A\vec{x} = \vec{b} \Rightarrow \vec{x} = A^{-1}\vec{b}$

• 解法① Gauss-Jordan  $\begin{bmatrix} 4 & 1 & | & 7 \\ 1 & 1 & | & 2 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & | & \frac{3}{3} \\ 0 & 0 & | & \frac{3}{3} \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & | & 1 \\ 0 & 1 & | & 1 \end{bmatrix}$

Note: For  $A_{N \times M}$ , Time =  $O(N^2 \times M)$

• 解法② LU decomposition

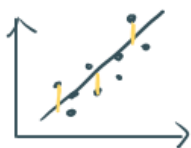
$A\vec{x} = \vec{b} \Rightarrow LU\vec{x} = \vec{b} \Rightarrow LY = \vec{b}$  解  $Y \Rightarrow U\vec{x} = Y$  解  $\vec{x}$

$\begin{bmatrix} 4 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 4 & \frac{1}{4} \\ 0 & \frac{3}{4} \end{bmatrix} \vec{x} = \vec{b} \Rightarrow \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} 7 \\ 2 \end{bmatrix}$  得  $\begin{bmatrix} y_1 = 7 \\ y_2 = \frac{1}{4} \end{bmatrix} \Rightarrow \begin{bmatrix} 4 & 1 \\ 0 & \frac{3}{4} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 7 \\ \frac{1}{4} \end{bmatrix}$  得  $\begin{bmatrix} x_1 = \frac{5}{3} \\ x_2 = \frac{1}{3} \end{bmatrix}$

Note: leading submatrices 的 det 都  $\neq 0$  才可 LU 分解

③ LL decomposition (cholesky)

## 6 LSE: least square error



• Concept: Minimum Loss (cost, distance, likelihood)

• 公式:  $\min \sum [f(x_i) - y_i]^2 = \min \|A\vec{x} - \vec{b}\|^2$

◦ 課本為  $\min_w \{y(x_n, w) - t_n\}^2$  Note: 对 poly basis func,  $y(x, w) = w_0 + w_1 x^1 + \dots + w_M x^M$

◦ 平方比  $\| \cdot \|$  容易微分 故以计算 error Note: error 有最小值:  $\text{error} \geq 0 \therefore \text{is impossible}$

## § LSE by matrix calculus

• concept :  $\|A\vec{x} - \vec{b}\|^2$

$$= \left\| \begin{bmatrix} x_0 & 1 \\ x_1 & 1 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} - \begin{bmatrix} y_0 \\ y_1 \end{bmatrix} \right\|^2$$

$$= (ax_0 + b - y_0)^2 + (ax_1 + b - y_1)^2$$

• Error

$$\|A\vec{x} - \vec{b}\|^2$$

$$L = (A\vec{x} - \vec{b})^T (A\vec{x} - \vec{b})$$

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

$$= \underbrace{\vec{x}^T A^T A \vec{x}}_{\text{Scalar}} - \underbrace{\vec{x}^T A^T \vec{b}}_{\text{Scalar}} - \underbrace{\vec{b}^T A \vec{x}}_{\text{Scalar}} + \vec{b}^T \vec{b}$$

$$= \vec{x}^T A^T A \vec{x} - 2\vec{x}^T A^T \vec{b} + \vec{b}^T \vec{b}$$

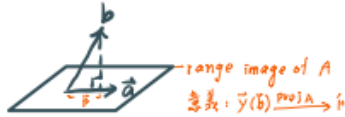
$$\frac{dL}{d\vec{x}} = 2A^T A \vec{x} - 2A^T \vec{b} = 0$$

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

Gram matrix, semi-positive definite ( $A^T A \geq 0$ ,  $\therefore$  not always invertible)  
 $\therefore$  不一定可解

• Geometry

Orthonormal projection 正交投影



$$a^T e = a^T (b - ax)$$

$$A^T (b - Ax) = 0$$

$$\Rightarrow a^T b = a^T a x$$

$$\Rightarrow A^T b = A^T A x$$

$$\Rightarrow \vec{x} = \frac{a^T b}{a^T a}$$

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

註:  $\frac{d(\vec{x}^T A^T A \vec{x})}{d\vec{x}}$  推导

1. 令  $A^T A$  为 symmetric matrix  $G = \begin{bmatrix} g_{11} & g_{12} & \dots & g_{1n} \\ g_{21} & g_{22} & \dots & g_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ g_{n1} & g_{n2} & \dots & g_{nn} \end{bmatrix}$

2.  $\frac{d(\vec{x}^T G \vec{x})}{d\vec{x}} = \begin{bmatrix} \frac{\partial}{\partial x_1} \\ \frac{\partial}{\partial x_2} \\ \vdots \\ \frac{\partial}{\partial x_n} \end{bmatrix} [x_1 \ x_2 \ \dots \ x_n] \begin{bmatrix} g_{11} & \dots & g_{1n} \\ g_{21} & \dots & g_{2n} \\ \vdots & \vdots & \vdots \\ g_{n1} & \dots & g_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$

$$= \begin{bmatrix} \frac{\partial}{\partial x_1} \\ \frac{\partial}{\partial x_2} \\ \vdots \\ \frac{\partial}{\partial x_n} \end{bmatrix} [x_1 (g_{11}x_1 + \dots + g_{1n}x_n) + x_2 (g_{21}x_1 + \dots + g_{2n}x_n) + \dots + x_n (g_{n1}x_1 + \dots + g_{nn}x_n)]$$

$$= \begin{bmatrix} 2g_{11}x_1 + (\sum_{i=1}^n g_{1i}x_i - g_{11}x_1) + (\sum_{i=1}^n g_{1i}x_i - g_{11}x_1) \\ \vdots \\ 2g_{nn}x_n + (\sum_{i=1}^n g_{ni}x_i - g_{nn}x_n) + (\sum_{i=1}^n g_{ni}x_i - g_{nn}x_n) \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^n g_{1i}x_i \\ \vdots \\ \sum_{i=1}^n g_{ni}x_i \end{bmatrix} + \begin{bmatrix} \sum_{i=1}^n g_{1i}x_i \\ \vdots \\ \sum_{i=1}^n g_{ni}x_i \end{bmatrix}$$

$$= G\vec{x} + G^T\vec{x} = 2G\vec{x} \quad (\because G = G^T)$$

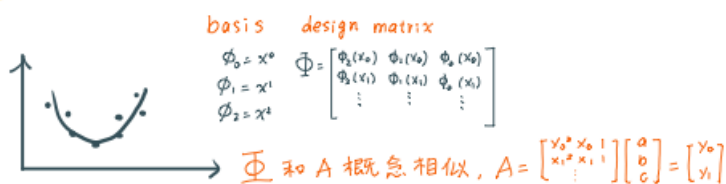
註:  $\frac{d\vec{x}^T A^T \vec{b}}{d\vec{x}}$  推导

$$= \begin{bmatrix} \frac{\partial}{\partial x_1} \\ \vdots \\ \frac{\partial}{\partial x_n} \end{bmatrix} [x_1 \ \dots \ x_n] \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nn} \end{bmatrix} \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix}$$

$$= \begin{bmatrix} \frac{\partial}{\partial x_1} \\ \vdots \\ \frac{\partial}{\partial x_n} \end{bmatrix} [x_1 (\sum_{i=1}^n a_{1i} b_i) + x_2 (\sum_{i=1}^n a_{2i} b_i) + \dots + x_n (\sum_{i=1}^n a_{ni} b_i)]$$

$$= [\sum_{i=1}^n a_{1i} b_i + \dots + \sum_{i=1}^n a_{ni} b_i] = A^T \vec{b}$$

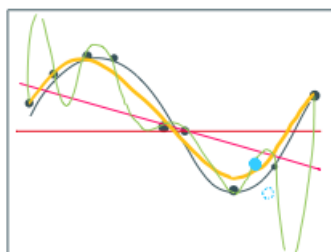
## LSE in non-linear case



- Concept:  $\|\Phi x - b\|^2$
- Error:  $\|\Phi x - b\|^2$ , similar to  $A$ . We can get  $x = (\Phi^T \Phi)^{-1} \Phi^T b$
- Problem1: Singularity  $\Rightarrow$  not invertible if singular
- Problem2: Too many bases  $\Rightarrow$  overfitting

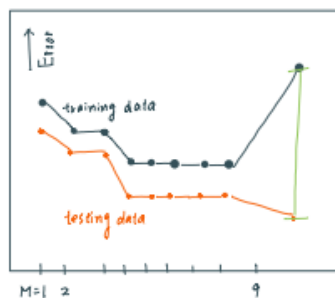
## Overfitting

- 決定 basic function 步驟
  1.  $W$  (係數)
  2. 改為 (in poly-func case)
  3. 決定元 basic function  $\Rightarrow$  對 Error  $E(W)$ , usually quadratic 微分  $\Rightarrow E(W)=0$  唯一解 = LSE
- Overfitting 示範



$M=3$   
 $M=0$   
 $M=1$   $M=0, 1$ : underfitting, 不足以描述 data

$M=9$   $M=9$ : overfitting, 太詳細, 以致再取 sine 上某時誤差變大  
 $M$  (basis function) 太多, 課本 suggest data 量的  $\frac{1}{5} \sim \frac{1}{10}$



Error 沒有趨小  $\searrow$ : overfit may occur

training data fitted perfectly, but won't fit well on other data

- 原因:
- ①  $M$  太大
  - ② data 太少, (若增加上圖  $\odot$ , 倒數 no.2 轉折可能不那么大)
  - ③ 雜訊 (上圖  $\bullet$ )

	$M=0$	$M=1$	$M=6$	$M=9$
$w_0$	0.19	0.82	0.31	0.35
$w_1$		1.27	7.99	232.37
$w_2$			-25.43	-3321.83
$w_3$			17.37	48568.31
$w_4$				-231639.30
$\vdots$				$\vdots$

overfit 特徵: function ||係數|| 很大

# 6 rLSE : regularization if overfitting

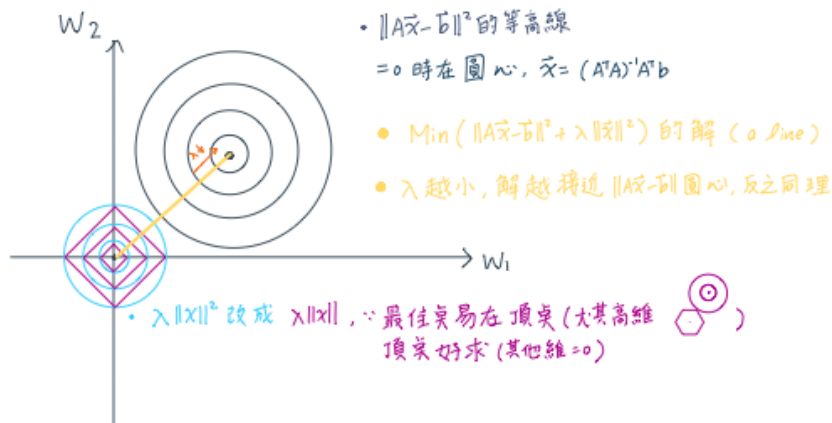
- Concept: 加  $\lambda$  penalty 懲罰太大的  $W$

- $E(w) : \sum_{n=1}^N [y(x_n, w) - t_n]^2 + \lambda \|w\|^2$

- 求 minimum  $E(w)$  by 微分:

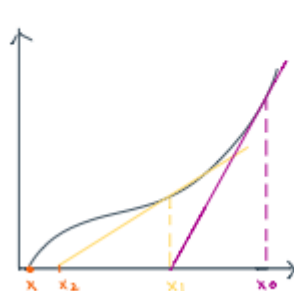
$$\begin{aligned}
 & \frac{d}{d\vec{x}} (\|A\vec{x} - \vec{b}\|^2 + \lambda \|\vec{x}\|^2) \\
 &= \frac{d}{d\vec{x}} [(A\vec{x} - \vec{b})^T (A\vec{x} - \vec{b}) + \lambda \vec{x}^T \vec{x}] \quad \text{Note: } (\vec{x}^T A^T - \vec{b}^T) (A\vec{x} - \vec{b}) \\
 &= \frac{d}{d\vec{x}} (\vec{x}^T A^T A \vec{x} - 2\vec{x}^T A^T \vec{b} + \vec{b}^T \vec{b} + \lambda \vec{x}^T \vec{x}) \quad \begin{aligned} &= \vec{x}^T A^T A \vec{x} - \vec{x}^T A^T \vec{b} - \vec{b}^T A \vec{x} + \vec{b}^T \vec{b} \\ &= \vec{x}^T A^T A \vec{x} - 2\vec{x}^T A^T \vec{b} + \vec{b}^T \vec{b} \quad \because (A\vec{x})^T \vec{b} = \vec{b}^T (A\vec{x}) \text{ 為純量} \end{aligned} \\
 &\Rightarrow 2A^T A \vec{x} - 2A^T \vec{b} + 2\lambda \vec{x} = 0 \\
 &\Rightarrow A^T A \vec{x} - A^T \vec{b} + \lambda \vec{x} = 0 \\
 &\Rightarrow (A^T A + \lambda I) \vec{x} = A^T \vec{b} \\
 &\Rightarrow \vec{x} = (A^T A + \lambda I)^{-1} A^T \vec{b} \quad \Rightarrow \text{for } A^T A \geq 0, A^T A + \lambda I \text{ 必} > 0 \text{ 為正定} \Rightarrow \text{必可逆}
 \end{aligned}$$

- 求 minimum  $E(w)$  by 圖:



# Newton's method : root finding $\Rightarrow$ optimization

## ● 观点①



Step 1 : choose  $x_0$ . find  $x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}$

Step 2 : continue from  $x_1$ .  $x_2 = x_1 - \frac{f(x_1)}{f'(x_1)}$

Step n : stop at  $x_n = x_{n-1} - \frac{f(x_{n-1})}{f'(x_{n-1})}$  till  $(x_n - x_{n-1})$  足够小  
或  
 $f'(x_n) = 0$

Ex :  $x^3 + 2x - 3 = (x+3)(x-1)$

取  $x_0 = 2$ ,  $x_1 = 2 - \frac{f(2)}{f'(2)} = 1.167$

$\Rightarrow x_2 = 1.167 - \frac{f(1.167)}{f'(1.167)} = 1.0$

$\Rightarrow \dots$

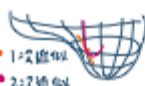
$\Rightarrow x_N$  收敛至 1.00

## ● 观点② : Taylor expansion

$f(x) \cong f(x_0) + \frac{1}{1!} f'(x_0)(x-x_0) + \frac{1}{2!} f''(x_0)(x-x_0)^2 + \dots$

$= \sum_{n=0}^{\infty} \frac{f^{(n)}(x_0)}{n!} (x-x_0)^n$

$= g(x)$  such that  $\begin{cases} g(x_0) = f(x_0) \\ g'(x_0) = f'(x_0) \\ \vdots \\ g^{(n)}(x_0) = f^{(n)}(x_0) \end{cases}$



1: 2 维图  
2: 3 维图

$(x_0=0) f(x) = f(x_0) + \frac{f'(x_0)}{1!} x + \frac{f''(x_0)}{2!} x^2 + \dots$  Maclaurin series

$(x_0=0) e^x = 1 + x + \frac{x^2}{2!} + \dots = \sum_{n=0}^{\infty} \frac{x^n}{n!}$

$(x_0=0) \sin x = \frac{\sin 0}{1!} x + \frac{\cos 0}{2!} x^2 + \frac{-\sin 0}{3!} x^3 + \frac{-\cos 0}{4!} x^4 + \frac{\sin 0}{5!} x^5 + \dots = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \dots$

$(x_0=0) \cos x = \frac{\cos 0}{1!} x + \frac{-\sin 0}{2!} x^2 + \frac{-\cos 0}{3!} x^3 + \frac{\sin 0}{4!} x^4 + \frac{\cos 0}{5!} x^5 + \dots = 1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \frac{x^6}{6!} + \dots$

$(x_0=0) e^{ix} = 1 + i \frac{e^{i\pi}}{1!} x + (-1) \frac{e^{i\pi}}{2!} x^2 + (-i) \frac{e^{i\pi}}{3!} x^3 + (1) \frac{e^{i\pi}}{4!} x^4 + \dots = i \sin x + \cos x$

$\Rightarrow$  for  $x=\pi$ ,  $e^{i\pi} + 1 = 0$

• Newton's Method in optimization

$g(x)$ : 把  $x_0$  的  $f(x)$  展開成 2 次近似

$$f(x) \cong f(x_0) + f'(x_0)\Delta x + \frac{f''(x_0)}{2!}\Delta x^2 = g(x)$$

註: Hessian matrix  
(Hessian func of  $f(x)$ )

$$\begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \dots \\ \frac{\partial^2 f}{\partial x_1 \partial x_2} & \frac{\partial^2 f}{\partial x_2^2} & \dots \\ \vdots & \vdots & \ddots \end{bmatrix}$$

目標:  $f'(x)=0$ , 從  $x_0$  出發, 找  $f'(x_0+\Delta x)=0$

$$\Rightarrow g'(x) = f'(x_0) + f''(x_0)\Delta x = 0$$

$$\Rightarrow \Delta x = \frac{-f'(x_0)}{f''(x_0)}$$

$$\begin{aligned} \therefore x_{n+1} &= x_n + \frac{-f'(x_n)}{f''(x_n)} \\ &= x_n - \frac{H^{-1}(f(x_n)) \nabla f(x_n)}{f''(x_n)} \Rightarrow \text{gradient} \end{aligned}$$

$\Rightarrow$  直取  $x_{i+1}$  至  $(x_i - x_{i-1})$  夠小

• Newton's method 驗證 LSE

$$\begin{aligned} \|A\bar{x} - \bar{b}\|^2 &= (A\bar{x} - \bar{b})^T (A\bar{x} - \bar{b}) \\ &= \bar{x}^T A^T A \bar{x} - 2\bar{x}^T A^T \bar{b} + \bar{b}^T \bar{b} \end{aligned}$$

$$\text{微分} \Rightarrow 2A^T A \bar{x} - 2A^T \bar{b} = \nabla f(\bar{x})$$

$$\text{再微} \Rightarrow 2A^T A = Hf(\bar{x})$$

$$\begin{aligned} \therefore \bar{x}_1 &= \bar{x}_0 - H^{-1}(\bar{x}_0) \nabla f(\bar{x}_0) \\ &= \bar{x}_0 - (2A^T A)^{-1} (2A^T A \bar{x}_0 - 2A^T \bar{b}) \\ &= \bar{x}_0 - \frac{1}{2} (A^T A)^{-1} (2A^T A \bar{x}_0 - 2A^T \bar{b}) \\ &= (A^T A)^{-1} A^T \bar{b} \end{aligned}$$

• 缺點

① Newton's method may be trapped in local

② Hessian takes  $O(n^3)$ , slow but takes less iteration

similar to  $\begin{cases} \text{gradient descent} \\ \text{conjugate gradient descent} \end{cases}$

③ Hessian 可能 invertible, 解法  $\begin{cases} \text{pseudo inverse} \\ \text{quasi-Newton approach} \end{cases}$

## L03-1 : probability

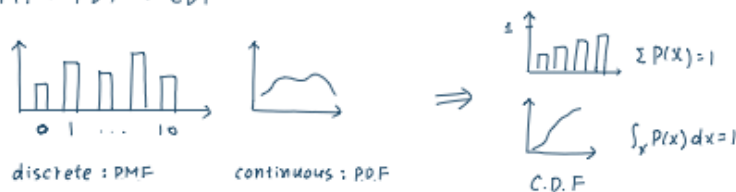
### Probability

Used for approximation since we know too few about the parameters that will affect the results

#### • 名詞 (例: 擲硬幣)

- trial: toss 2 coins
- outcome: HH, HT, TH, TT
- event: set of outcomes {HH, HT, TH, TT}
- sample space: all outcomes (U)
- random variable: mapping function  
 ex:  $\begin{matrix} HH \rightarrow 0 \\ HT \rightarrow 1 \\ TH \rightarrow 2 \\ TT \rightarrow 3 \end{matrix} \Rightarrow P(X=H) = P(X=1) \xrightarrow[\text{of trials}]{\text{no times}} P(X=1) = \frac{1}{2}$

#### • PMF, PDF, CDF



### Conditional Joint probability

- Conditional probability  $P(A|B) = \frac{P(A \cdot B)}{P(B)} = \frac{P(A) P(B|A)}{P(B)}$  (sample space shrink)
- Joint probability  $P(A \cdot B) = P(A) P(B|A) = P(A \cdot B | U)$ , U 為字集






#### • Conditional problem 三門問題

會換/不換 Y/N  
 得車/羊 C/G

Key: sample space changed  
 (shrink as ② is opened)  
 主持人必選羊 (非公正)

選①後知②, 換/不換?

#### 解法 1

			換	不換
Pick		host	T	F
host	Pick	host	F	T
host		Pick	T	F

$\therefore$  換而中獎機率 =  $\frac{2}{3}$   
 不換中獎機率 =  $\frac{1}{3}$

#### 解法 2

選	開	換	結果	P
車	$\frac{1}{3}$	羊	$\begin{cases} 0 \frac{1}{2} \\ 1 \frac{1}{2} \end{cases}$	$\begin{cases} \frac{1}{6} \\ \frac{1}{6} \end{cases}$
羊	$\frac{2}{3}$	羊	$\begin{cases} 0 \frac{1}{2} \\ 1 \frac{1}{2} \end{cases}$	$\begin{cases} \frac{2}{6} \\ \frac{2}{6} \end{cases}$

$\therefore P(\text{車|換}) = \frac{\frac{2}{6}}{\frac{3}{6}} = \frac{2}{3}$   
 $P(\text{車|不換}) = \frac{\frac{1}{6}}{\frac{3}{6}} = \frac{1}{3}$

✱ 考試可能出 4 門 2 車等

## L03-2 : Bayes Theorem

### Bayes's Theorem

Posterior = likelihood  $\frac{P(A)}{P(B)}$  Prior, A is parameter  

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)} \Rightarrow P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)}, D = \text{event}$$
marginal

Ex: A = Vote for 馬 P(A) = 10%  
 B = Vote for 洪 P(B) = 50%  

$$P(B|A) = 90\% \Rightarrow P(A|B)$$

Sum rule:  $P(A) = P(A, B) + P(A \text{ not } B) = \sum_v P(A, v)$  marginalize

Ex: A = lung cancer  
 B = smoke  

$$P(B|A) \text{ 好求 } \because A \text{ 人少}$$
  

$$\Rightarrow P(A|B) \text{ 吸烟得肺癌机率} = \frac{0.1 \times 0.9}{0.5} = 18\%$$

product rule:  $P(A_0, A_1, A_2, A_3) = P(A_0, A_1, A_2 | A_3) P(A_3)$   

$$= P(A_0, A_1 | A_2, A_3) P(A_2 | A_1) P(A_1) P(A_0)$$
  

$$= P(A_0 | A_1, A_2, A_3) P(A_1 | A_2, A_3) P(A_2 | A_1) P(A_1) P(A_0)$$

### Bayesian v.s. Frequentist (見 Los note)

Bayesian	Frequentist
Prior knowledge	observation
$P(B) = \frac{1}{2}$ though BB...B	$\therefore BB...B \Rightarrow$ guess next is B

### Distribution

Describe

location 1st moment  
 mean  $E(x) = \sum P(x_i) x_i$   
 median  
 mode (most frequent)

Dispersion - Variance 2nd moment  

$$\text{Var}(x) = E((x - \mu)^2) = \frac{1}{n} \sum (x_i - \mu)^2 = \frac{1}{n} \sum (x_i^2 - 2\mu x_i + \mu^2) = \frac{1}{n} \sum x_i^2 - \frac{1}{n} 2\mu \sum x_i + \frac{1}{n} \sum \mu^2 = \sum x_i^2 - \mu^2 = E(x^2) - E(x)^2$$

Skewness 3rd moment  

$$E((x - \mu)^3)$$

Kurtosis (peakness) 4th moment  

$$E((x - \mu)^4)$$



# Naïve Bayes classifier

- MLE: 找  $\theta = P$  得到 maximum likelihood  $P(D|\theta) \Rightarrow L(\theta|D)$

例: 对 Bayesian 全  $\theta$  的 distribution



trial 1:

$$P(\theta=0|H) = \frac{P(H|\theta=0)P(\theta)}{P(H)} = 0$$

$$P(\theta=0.3|H) = \frac{P(H|\theta=0.3)P(\theta)}{P(H)} = \frac{0.3 \cdot 0.1}{0.5} = 0.06 \xrightarrow{HH} 0.032 \dots \xrightarrow{H^{10}} 1.03E-06$$

$$P(\theta=0.5|H) = \frac{P(H|\theta=0.5)P(\theta)}{P(H)} = \frac{0.5 \cdot 0.1}{0.5} = 0.1 \xrightarrow{HH} 0.618 \dots \xrightarrow{H^{10}} 0.03$$

$$P(\theta=0.7|H) = \frac{0.7 \cdot 0.1}{0.5} = 0.14 \xrightarrow{HH} 0.173 \dots \xrightarrow{H^{10}} 0.027$$

$$P(\theta=1|H) = \frac{1 \cdot 0.05}{0.5} = 0.1 \xrightarrow{HH} 0.177 \dots \xrightarrow{H^{10}} 0.977$$

∴ H 次取 ↑, MLE 發生點從  $\theta=0.5 \rightarrow \theta=1$  (趨 frequentist)

例: 对 Frequentist

$$P(HH|P=1) = L(P=1|HH) = \textcircled{1} \text{ MLE}$$

$$P(HH|P=0.7) = 0.49$$

$$P(HH|P=0.5) = 0.25$$

$$P(HH|P=0.1) = 0.01$$

## Bayes 例題

$$\text{技① } P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)}$$

data 不夠假設

$$\text{技② conditional independence} = \frac{P(d_1|\theta)P(d_2|\theta)P(d_3|\theta)P(d_4|\theta)}{P(D)}$$

$$P(\text{play} = \text{yes} \mid \begin{matrix} \text{outlook} = \text{sunny} \\ \text{Temp} = \text{cool} \\ \text{Humidity} = \text{high} \\ \text{Wind} = \text{strong} \end{matrix})$$

$$\text{技③ } \frac{P(O=\text{sunny}, T=\text{cool}, H=\text{high}, W=\text{strong}|\text{yes})P(\text{yes})}{P(O=\text{sunny}, T=\text{cool}, H=\text{high}, W=\text{strong})}$$

∴ 最後全 normalize

$$= \frac{P(\text{sunny}|\text{yes})P(\text{cool}|\text{yes})P(\text{high}|\text{yes})P(\text{strong}|\text{yes}) \cdot \frac{9}{14}}{P(\text{sunny}, \text{cool}, \text{high}, \text{strong})}$$

$$= \frac{\frac{2}{9} \cdot \frac{3}{9} \cdot \frac{3}{9} \cdot \frac{3}{9} \cdot (\frac{9}{14})}{P(\text{sunny}, \text{cool}, \text{high}, \text{strong})}$$

$$= \frac{\frac{2}{9} \cdot \frac{3}{9} \cdot \frac{3}{9} \cdot \frac{3}{9} \cdot (\frac{9}{14})}{P(\text{sunny}, \text{cool}, \text{high}, \text{strong})}$$

註: conditional independence

$$P(A, B) = P(A)P(B)$$

$$P(A, B|U) = P(A|U)P(B|U)$$

$$P(A_1, A_2, A_3, A_4)$$

$$= P(A_1|A_2, A_3, A_4)P(A_2|A_3, A_4)P(A_3|A_4)P(A_4)$$

If  $A_2, A_3, A_4$  are independent

$$= P(A_1|A_4)P(A_2|A_4)P(A_3|A_4)$$

$$\textcircled{1} \text{ } P(A_1, A_2|A_3, A_4) = P(A_1|A_3, A_4)P(A_2|A_3, A_4)$$

$$\Rightarrow \frac{P(A_1, A_2, A_3, A_4)}{P(A_3, A_4)} = \frac{P(A_1, A_2, A_4)}{P(A_3, A_4)} \cdot \frac{P(A_3, A_4)}{P(A_3, A_4)}$$

$$\Rightarrow \frac{P(A_1, A_2, A_3, A_4)}{P(A_3, A_4)} = \frac{P(A_1, A_2, A_4)}{P(A_3, A_4)} \cdot P(A_3|A_4) \quad \therefore P(A_1|A_2, A_3, A_4) = P(A_1|A_3, A_4) = P(A_1|A_4)$$

$$\textcircled{2} \text{ } P(A_2, A_3|A_4) = P(A_2|A_4)P(A_3|A_4)$$

$$\Rightarrow \frac{P(A_2, A_3, A_4)}{P(A_4)} = \frac{P(A_2, A_4)}{P(A_4)} \cdot \frac{P(A_3, A_4)}{P(A_4)}$$

$$\Rightarrow \frac{P(A_2, A_3, A_4)}{P(A_3, A_4)} = \frac{P(A_2, A_4)}{P(A_4)} \quad \therefore P(A_2|A_3, A_4) = P(A_2|A_4)$$

$$\therefore \begin{cases} P(\text{play} = \text{Yes} | D) = \frac{\frac{2}{9} \cdot \frac{3}{9} \cdot \frac{3}{9} \cdot \frac{3}{9} \cdot (\frac{9}{14})}{P(D)} = \frac{1}{14} \cdot \frac{1}{P(D)} \cdot \frac{2}{3^3} \\ P(\text{play} = \text{No} | D) = \frac{\frac{3}{5} \cdot \frac{1}{3} \cdot \frac{4}{5} \cdot \frac{3}{5} \cdot (\frac{5}{14})}{P(D)} = \frac{1}{14} \cdot \frac{1}{P(D)} \cdot \frac{36}{5^3} \end{cases}$$

$$\Rightarrow \frac{P(\text{play} = \text{yes} | D)}{P(\text{play} = \text{no} | D)} = \frac{2 \cdot 5^3}{3^3 \cdot 36} = \frac{125}{486}$$

$$\Rightarrow \begin{cases} P(\text{play} = \text{yes} | D) = \frac{125}{611} \approx 20\% \\ P(\text{play} = \text{no} | D) = \frac{486}{611} \approx 80\% \end{cases}$$

# Information theory

- Entropy 熵 : randomness 亂度, to describe information

o unit of randomness / uncertainty :  $-\log_2 P$ , which shows how many bits to describe

$$\text{Ex: } \begin{cases} -\log \frac{1}{4} = 2 \text{ (take 2 bits)} \\ -\log \frac{1}{1024} = 10 : \text{rare event gain more info} \end{cases}$$

o Def: **Exp** (info) in events

$$\Delta \text{ review } \begin{cases} E(x) = \sum p(x) \cdot x \\ E(f(x)) = \sum p(x) \cdot f(x) \end{cases}$$

$$\Delta H(x) = -\sum p(x) \log p(x) = -E_{x \sim p} [\log_2 p]$$

△ Example : flip 2 fair coins

mapping : random variable  $X$

$$\begin{array}{ll} HH \rightarrow 2 & H(X) = -\sum_{\text{uniform}} \frac{1}{4} \cdot \log_2 \frac{1}{4} = 2 \text{ (得 2 bit 資訊量)} \\ HT \rightarrow 1 \\ TH \rightarrow 1 \\ TT \rightarrow 0 \end{array}$$

$$H(X) = -\left[\frac{1}{16} \log_2 \frac{1}{16} + 2 \cdot \left(\frac{1}{16} \log_2 \frac{1}{16}\right) + \frac{1}{16} \log_2 \frac{1}{16}\right] \text{ If } p(H) = \frac{3}{4}$$

△ Note :  $\begin{cases} \text{max(entropy)} = \text{base} \\ \text{max(entropy)} \text{ 發生在 uniform (最 predict, 越亂} \rightarrow \text{含更多 info)} \end{cases}$   
 $-\sum 1 \log 1 = 0$ , 若 always happen 則無 uncertainty

- Conditional entropy

$$H(Y|X) = -\sum p(x,y) \log \frac{p(x,y)}{p(x)} = (H(x,y) - H(x))$$

$$\text{pf: } H(Y|X) = -\sum_i p(x_i) H(Y|X=x_i)$$

$$= -\sum_i p(x_i) \sum_j p(y_j|x_i) \log p(y_j|x_i)$$

$$= -\sum_i \sum_j p(x_i, y_j) \log p(y_j|x_i)$$

$$= -\sum_i \sum_j p(x_i, y_j) \log \frac{p(x_i, y_j)}{p(x_i)}$$

$$= \sum_i \sum_j p(x_i, y_j) \log \frac{p(x_i)}{p(x_i, y_j)}$$

- Joint entropy

$$H(x,y) = H(x) + H(y|x)$$

Info needed by  $H(x)$  to get  $H(x,y)$

$$\text{pf: } H(x,y) = -\sum_{x \in X} \sum_{y \in Y} p(x,y) \log p(x,y)$$

$$= -\sum_{x \in X} \sum_{y \in Y} p(x,y) \log p(x) p(y|x)$$

$$= -\sum_{x \in X} \sum_{y \in Y} p(x,y) \log p(x) - \sum_{x \in X} \sum_{y \in Y} p(x,y) \log p(y|x)$$

$$= -\sum_{x \in X} p(x) \log p(x) - \sum_{x \in X} \sum_{y \in Y} p(x,y) \log p(y|x)$$

$$= H(x) + H(y|x)$$

- Relative entropy (KL divergence)

o usage : **Distance** between 2 distribution (random var)

$$o KL(M \parallel N) = H_N(M) - H(M)$$

$$= [-\sum_i M(x_i) \log N(x_i)] - [-\sum_i M(x_i) \log M(x_i)]$$

$$= -\sum_i M(x_i) \log \frac{N(x_i)}{M(x_i)}$$

$$= \sum_i M(x_i) \log \frac{M(x_i)}{N(x_i)}$$

$$\therefore KL(P \parallel Q) = \sum_{x \in X} P(x) \log \frac{P(x)}{Q(x)} = E_{x \sim P} [\log \frac{P(x)}{Q(x)}]$$

但不好用  $\because KL(P \parallel Q) \neq KL(Q \parallel P)$

$$\Rightarrow \text{改用 cross entropy } H(P, Q) = E_{x \sim P} [-\log Q(x)]$$

$$= -P \log Q - (1-P) \log (1-Q)$$

$$= KL(P \parallel Q) + H(P)$$

- Mutual Information  mutual entropy

o usage : to show the independence

$\therefore$  If  $X, Y$  independent,  $I(x,y) = 0$

o 公式 :  $I(x,y) = H(x) - H(x|y) \vee H(y) - H(y|x)$

$$\text{pf: } I(x,y) = \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p(x)p(y)}$$

$$= \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p(x)} - \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p(y)}$$

$$= H(x) - H(x|y) \vee H(y) - H(y|x)$$

o Problem : The scale is unpredictable

## L04-2 : Maximum entropy

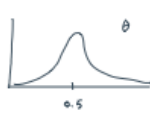
### Maximum Entropy Principle

◦ frequentist v.s. bayesian

$$L(H, H|\theta=0.5) = \frac{1}{2} \cdot 0$$

$$L(H, H|\theta=0.9) = \frac{0.9}{1.0}$$

$$L(H, H|\theta=1) = 1$$



↳ frequentist 看到会猜 coin always H | bayesian 有 prior ∴ 猜 coin 还是 fair

◦ uniform distribution → max(entropy)

$$Pf: H(x) = -\int_a^b p(x) \log p(x) dx, p(x) \geq 0 \text{ \& } \int_a^b p(x) dx = 1$$

$$\Rightarrow \frac{\delta L}{\delta p(x)} \left[ -\int_a^b p(x) \log p(x) dx + \lambda \left( \int_a^b p(x) dx - 1 \right) \right] \text{ 是 lagrange multiplier}$$

$$\Rightarrow \frac{\delta L}{\delta p(x)} \left[ -\log p(x) - 1 + \lambda \right] = 0$$

$$\Rightarrow [-\log p(x) - 1 + \lambda] = 0$$

$$\Rightarrow p(x) = e^{\lambda-1}$$

$$\because \int_a^b p(x) dx = 1 \Rightarrow p(x) \Big|_a^b = 1 \Rightarrow p(x) \frac{1}{b-a} = 1$$

$$\therefore p(x) = e^{\lambda-1} = \frac{1}{b-a}$$

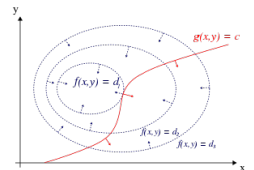
註: lagrange multiplier: f 在 constraint g 取極值

$$Ex: f = 4xy, \text{ constraint } x > 0, y > 0, \frac{x^2}{4} + \frac{y^2}{16} = 1$$

$$L = f - \lambda g, g = \frac{x^2}{4} + \frac{y^2}{16}$$

$$\nabla L = (4y, 4x) - \lambda \left( \frac{2x}{4}, \frac{2y}{16} \right)$$

$$\Rightarrow \begin{cases} 4y = \lambda \frac{x}{2} \\ 4x = \lambda \frac{y}{4} \end{cases} \Rightarrow \frac{y}{x} = \frac{\lambda}{8} \Rightarrow \frac{y}{x} = \frac{4}{3} \text{ 代 } \lambda g \Rightarrow y = \pm 2\sqrt{3} \text{ (取正)}$$



### Maximum Entropy (Given $\mu$ )

$$\text{Given expectation: } \int_0^\infty x p(x) dx = \mu \text{ 求: } \max = \int_0^\infty p(x) \ln p(x) dx$$

$$\hat{L} = \int_0^\infty p(x) \ln p(x) dx + \lambda_0 \left( \int_0^\infty p(x) dx - 1 \right) + \lambda_1 \left( \int_0^\infty x p(x) dx - \mu \right)$$

$$\frac{\delta L}{\delta p(x)} = -(\ln p(x) + 1) + \lambda_0 + \lambda_1 x = 0$$

$$\Rightarrow \ln p(x) = \lambda_0 + \lambda_1 x - 1$$

$$\Rightarrow p(x) = e^{\lambda_0 + \lambda_1 x - 1} \quad (a_1)$$

$$\text{代 } \lambda C_1 \Rightarrow \int_0^\infty e^{\lambda_0 + \lambda_1 x - 1} dx = 1$$

$$\Rightarrow \frac{1}{\lambda_1} e^{\lambda_0 + \lambda_1 x - 1} \Big|_0^\infty = 1$$

$$\Rightarrow \frac{e^{\lambda_0 - 1}}{\lambda_1} (e^{\lambda_1 \infty} - 1) = 1 \therefore \lambda_1 \leq 0 \text{ 使 } e^{\lambda_1 \infty} = 0$$

$$\Rightarrow \frac{e^{\lambda_0 - 1}}{\lambda_1} = -1 \Rightarrow \lambda_1 = -e^{\lambda_0 - 1} \quad (a_2)$$

$$\text{代 } \lambda C_2: \int_0^\infty x e^{\lambda_0 + \lambda_1 x - 1} dx = \mu \text{ Note: } \int f g' = f g - \int f' g$$

$$\Rightarrow e^{\lambda_0 - 1} \left[ x \frac{e^{\lambda_1 x}}{\lambda_1} \Big|_0^\infty - \int_0^\infty \frac{e^{\lambda_1 x}}{\lambda_1} dx \right] = \mu$$

$$\Rightarrow \frac{1}{-\lambda_1} \int_0^\infty e^{\lambda_0 + \lambda_1 x - 1} dx = \frac{1}{-\lambda_1} = \mu$$

$$\Rightarrow \lambda_1 = -\frac{1}{\mu} \quad (a_3)$$

$$\text{From } (a_1)(a_2)(a_3): p(x) = e^{\lambda_0 - 1} \cdot e^{\lambda_1 x}$$

$$= -\lambda_1 \cdot e^{-\frac{1}{\mu} x}$$

$$= \frac{1}{\mu} e^{-\frac{1}{\mu} x} \text{ exponential distribution}$$

### Maximum Entropy (Given $\mu, \sigma^2$ )




$$\hat{L} = \int p(x) \ln p(x) dx + \lambda_0 \left( \int p(x) dx - 1 \right) + \lambda_1 \left( \int x p(x) dx - \mu \right) + \lambda_2 \left( \int (x - \mu)^2 p(x) dx - \sigma^2 \right)$$

$$\frac{\delta L}{\delta p(x)} = -(\ln p(x) + 1) + \lambda_0 + \lambda_1 x + \lambda_2 (x - \mu)^2 = 0$$

$$\Rightarrow p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \text{ Gaussian distribution}$$

## L05-1 : Bernoulli distribution

### ● 整理

- uniform distr. 
- exponential distr. 
- Gaussian distr. 

## 🔗 Bernoulli Distribution

- $P(X=0,1|\theta) = \theta^X (1-\theta)^{1-X}$



$$E(X) = 1 \cdot \theta + 0 \cdot (1-\theta) = \theta \quad \text{Location}$$

$$\text{Var}(X) = E(X^2) - E(X)^2$$

$$= \theta - \theta^2$$

$$= \theta(1-\theta) \quad \text{Dispersion}$$

### ● Bernoulli distribution

- Given  $[0,0,1,\dots] = [x_1, x_2, \dots, x_N] = D$

$$P(D|\theta) = \prod_{i=1}^N \theta^{x_i} (1-\theta)^{1-x_i}$$

✧ 若已知結果  $D$ ，求參數  $\theta$  讓  $D$  發生機率最大 = MLE

$$\text{即找 } \arg \max_{\theta} \prod_{i=1}^N \theta^{x_i} (1-\theta)^{1-x_i} \Rightarrow \arg \max_{\theta} [\sum x_i \log \theta + \sum (1-x_i) \log (1-\theta)]$$

$$\text{令 } L = \sum x_i \log \theta + \sum (1-x_i) \log (1-\theta)$$

$$\frac{dL}{d\theta} = \sum_{i=1}^N x_i \frac{1}{\theta} - \sum_{i=1}^N (1-x_i) \frac{1}{1-\theta} = 0 \quad \text{Note: } \frac{d}{dx} \log x = \frac{1}{x}$$

$$\Rightarrow \sum_{i=1}^N x_i \frac{1}{\theta} = \sum_{i=1}^N \frac{1}{1-\theta} - \sum_{i=1}^N x_i \frac{1}{1-\theta}$$

$$\Rightarrow \sum_{i=1}^N x_i \left( \frac{1}{\theta} + \frac{1}{1-\theta} \right) = \sum_{i=1}^N \frac{1}{1-\theta}$$

$$\Rightarrow \left( \frac{1}{\theta(1-\theta)} \right) \sum_{i=1}^N x_i = \frac{N}{1-\theta}$$

$$\Rightarrow \theta = \frac{\sum_{i=1}^N x_i}{N} \Rightarrow \text{MLE 發生在 sample mean}$$

∴ 若 10 次有 5 次中，應猜  $\theta$  為  $\frac{5}{10}$

## Binomial distribution

- Multiple Bernoulli trial ( $N$ )  $\Rightarrow$  Binomial

$$P(X=m | N, \theta) = \binom{N}{m} \theta^m (1-\theta)^{N-m}, \quad X \text{ is \# of success in } N \text{ trials}$$

$$E(X) = N\theta$$

$$\text{Var}(X) = N\theta(1-\theta)$$

$$\text{MLE} = \theta = \frac{\sum m}{\sum N}$$

## conjugate prior

(Maximum Likelihood Estimation)

Frequentist based on MLE - dataset, Bayesian based on prior ex:  $P(\theta=0.5)$  = 大机率

Frequentist 忽略其他可能, Bayesian 要看 prior  $P(\theta)$  好壞  $P(\theta=1)$  = 小机率

$\therefore$  折衷找  $P(\theta|x)$ , posterior 為:

$$P(\theta|x) = \frac{\text{Binomial} \cdot \text{Beta}}{P(x)} = \frac{P(x|\theta)P(\theta)}{P(x)}$$

- conjugate prior - posterior

o 來由: 算 posterior 太累,  $\therefore$  找 distribution 使 prior - posterior in the same form  
此 distribution is called conjugate

## Beta distribution

- Beta function  $B(a, b)$

$$= \int_0^1 x^{a-1} (1-x)^{b-1} \frac{1}{B(a, b)} dx$$

$$= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \int_0^1 x^{a-1} (1-x)^{b-1} dx$$

$$E(X) = \frac{a}{a+b} \quad \text{Var}(X) = \frac{ab}{(a+b)^2(a+b+1)}$$

Pf:

$$\begin{aligned} E(X) &= \int_0^1 x \cdot x^{a-1} (1-x)^{b-1} \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} dx \\ &= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \int_0^1 x^{a+1-1} (1-x)^{b-1} dx \\ &= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \cdot \frac{\Gamma(a+1)\Gamma(b)}{\Gamma(a+b+1)} \end{aligned}$$

$$= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \cdot \frac{a\Gamma(a)\Gamma(b)}{(a+b)\Gamma(a+b)} = \frac{a}{a+b}$$

$$\begin{aligned} E(X^2) &= \int_0^1 x^2 \cdot x^{a-1} (1-x)^{b-1} \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} dx \\ &= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \int_0^1 x^{a+2-1} (1-x)^{b-1} dx \\ &= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \cdot \frac{\Gamma(a+2)\Gamma(b)}{\Gamma(a+b+2)} \\ &= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \cdot \frac{(a+1)a\Gamma(a)\Gamma(b)}{(a+b+1)(a+b)\Gamma(a+b)} = \frac{(a+1)a}{(a+b+1)(a+b)} \end{aligned}$$

$$\text{Var}(X) = E(X^2) - E(X)^2$$

$$= \frac{a}{a+b} \left( \frac{a+1}{a+b+1} - \frac{a}{a+b} \right) = \frac{a}{a+b} \left( \frac{b}{(a+b+1)(a+b)} \right)$$

$$= \frac{ab}{(a+b)^2(a+b+1)}$$

註: gamma function  $\Gamma(x) = \int_0^\infty t^{x-1} e^{-t} dt$

性質 ①  $\Gamma(x) = (x-1)\Gamma(x-1)$

$$\begin{aligned} \text{Pf: } \Gamma(x) &= \int_0^\infty t^{x-1} e^{-t} dt \\ &= -t^{x-1} e^{-t} \Big|_0^\infty + \int_0^\infty (x-1)t^{x-2} e^{-t} dt \\ &= (x-1) \int_0^\infty t^{x-2} e^{-t} dt \\ &= (x-1)\Gamma(x-1) \end{aligned}$$

因此推得

$$\Gamma(1) = \int_0^\infty e^{-t} dt = -e^{-t} \Big|_0^\infty = 0 - (-1) = 1$$

$$\Gamma(2) = \int_0^\infty t e^{-t} dt = -te^{-t} \Big|_0^\infty + \int_0^\infty 1 \cdot e^{-t} dt = 0 + 1 = \Gamma(1)$$

$$\Gamma(3) = \int_0^\infty t^2 e^{-t} dt = -t^2 e^{-t} \Big|_0^\infty + \int_0^\infty 2t e^{-t} dt = 0 + 2\Gamma(2)$$

$$\Gamma(x) = \begin{cases} 1, & x=1, 2 \\ (x-1)\Gamma(x-1) = (x-1)!, & \text{otherwise} \end{cases}$$

$$\text{性質 ② } \int_0^1 x^{a-1} (1-x)^{b-1} dx = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$$

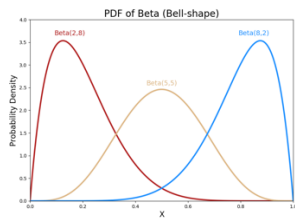
$$\begin{aligned} \text{Pf: } \Gamma(a)\Gamma(b) &= \int_0^\infty x^{a-1} e^{-x} dx \int_0^\infty y^{b-1} e^{-y} dy \\ &= \int_0^\infty \int_0^\infty e^{-(x+y)} x^{a-1} y^{b-1} dx dy \\ &\stackrel{\text{令 } x=uv, y=u(1-v)}{\Rightarrow} \int_0^\infty \int_0^1 e^{-u} u^{a-1} v^{a-1} (1-v)^{b-1} du dv \\ &= \int_0^\infty e^{-u} u^{a-1} v^{a-1} (1-v)^{b-1} du dv \\ &= \int_0^\infty v^{a-1} (1-v)^{b-1} dv \cdot \int_0^\infty e^{-u} u^{a-1} du \\ &= \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)} \end{aligned}$$

$$\therefore \beta(a+b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$$

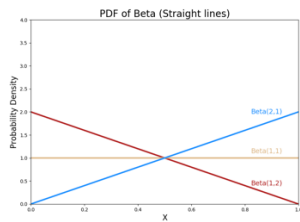
## L05 : 補充

The PDF of Beta distribution can be U-shaped with asymptotic ends, bell-shaped, strictly increasing/decreasing or even straight lines. As you change  $\alpha$  or  $\beta$ , the shape of the distribution changes.

### a. Bell-shape

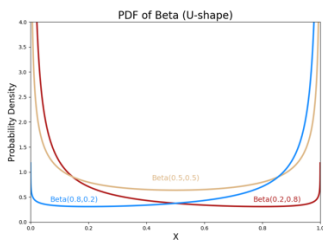


### b. Straight Lines



### c. U-shape

When  $\alpha < 1$ ,  $\beta < 1$ , the PDF of the Beta is U-shaped.



# Beta Binomial conjugate

likelihood:  $P(x=m | N, p) = \binom{N}{m} p^m (1-p)^{N-m}$

prior:  $B(\theta | a, b) = \theta^{a-1} (1-\theta)^{b-1} \frac{1}{\beta(a, b)}$ , 成功  $a$  次 trial  $a+b$  次

posterior:  $\frac{\text{likelihood} \cdot \text{prior}}{\text{marginal}} = B(p | a+m, b+N-m) = \frac{p^{a+m-1} (1-p)^{b+N-m-1} \frac{1}{\beta(a, b)}}{\int_0^1 \theta^{a-1} (1-\theta)^{b-1} \frac{1}{\beta(a, b)} d\theta}$

Pf: 
$$\frac{\binom{N}{m} p^m (1-p)^{N-m} \frac{1}{\beta(a, b)}}{\int_0^1 \binom{N}{m} \theta^m (1-\theta)^{N-m} \theta^{a-1} (1-\theta)^{b-1} \frac{1}{\beta(a, b)} d\theta}$$

$$= \frac{p^{m+a-1} (1-p)^{N-m+b-1}}{\int_0^1 \theta^{m+a-1} (1-\theta)^{N-m+b-1} d\theta} = \frac{\Gamma(m+a) \Gamma(N-m+b)}{\Gamma(N+a+b)}$$

$$= p^{m+a-1} (1-p)^{N-m+b-1} \frac{1}{\beta(a+m, b+N-m)}$$

$$= B(p, a+m, b+N-m)$$

註: review posterior



应用: online (sequential) learning

△ prior  $\xrightleftharpoons{\text{calculate}}$  posterior

△ 例

5 success 6 fail  $\rightarrow$  +1 success 2 fail  $\rightarrow$  +100 success



△ 优点

得到 distribution 而非 point estimation

一可得 mean, var

二可从 mode 得 most likely  $\theta$

## Frequentist v.s. Bayesian

Predict based on Data

算 Likelihood function

$\Rightarrow$  MLE to the  $\hat{\theta}_{MLE}$

Predict Based on prior  $\Rightarrow$  資料少也能利用 prior (knowledge) 預測

利用 posterior distribution

$\Rightarrow$  MLE to the  $\hat{\theta}_{MAP}$

$$\hat{\theta}_{MAP} = \frac{\text{likelihood} \cdot \text{prior}}{\text{marginal}} = \frac{\binom{N}{m} p^m (1-p)^{N-m} \frac{1}{\beta(a, b)} p^{a-1} (1-p)^{b-1}}{\int_0^1 \theta^{a-1} (1-\theta)^{b-1} \frac{1}{\beta(a, b)} d\theta}$$

$$= B(p, a+m, b+N-m) \Rightarrow \text{distribution}$$

$\Rightarrow$  算 MLE of  $\hat{\theta}_{MAP}$

$$\frac{dB(p, a, b)}{dp} = 0$$

$$\Rightarrow \frac{1}{\beta(a, b)} \frac{d[p^{a-1} (1-p)^{b-1}]}{dp} = \frac{1}{\beta(a, b)} [p^{a-2} (1-p)^{b-1} - p^{a-1} (1-p)^{b-2}] = 0$$

$$\Rightarrow p^{a-2} (1-p)^{b-2} (1-2p) = 0$$

$$\Rightarrow p = \frac{1}{2} (?)$$

## L05-4 : Multinomial

### ● Multi nomial

o Dice 100 times for example :



$$X = [m_1, m_2, \dots, m_6] = [30, 40, \dots]$$

$$M(x) = \binom{N}{m_1, m_2, \dots, m_6} \prod_i p_i^{m_i}$$

$$= \frac{N!}{m_1! m_2! \dots m_6!}$$

o Dirichlet distribution

likelihood	Model params	Conjugate Prior distribution	Posterior hyperparams	Interpretation of hyperparams	Posterior prediction
Bernoulli	$p$ (probability)	Beta	$\alpha, \beta$	$\alpha + \sum x_i, \beta + n - \sum x_i$	
Binomial	$p$	Beta	$\alpha, \beta$	$\alpha + \sum x_i, \beta + \sum N_i - \sum x_i$	
Multi nomial	$p$	Dirichlet			



# Gaussian Distribution

## Univariate / Multivariate Gaussin



## Gaussian Integral : $\int_{-\infty}^{\infty} e^{-x^2} dx = \sqrt{\pi}$

Pf:  $\int_{-\infty}^{\infty} e^{-x^2} dx$

$$= \sqrt{\left(\int_{-\infty}^{\infty} e^{-x^2} dx\right)^2}$$

$$= \sqrt{\left(\int_{-\infty}^{\infty} e^{-x^2} dx\right) \left(\int_{-\infty}^{\infty} e^{-y^2} dy\right)}$$

$$= \sqrt{\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{-(x^2+y^2)} dx dy}$$

$$= \sqrt{\int_0^{2\pi} \int_0^{\infty} e^{-r^2} (r dr d\theta)}$$

$$= \sqrt{2\pi \int_0^{\infty} r e^{-r^2} dr} = \sqrt{2\pi \left(-\frac{1}{2} e^{-r^2}\right)_0^{\infty}}$$

$$= \sqrt{\pi}$$

註①: polar coordinate 極座標



$$\text{面積} = \int_0^{2\pi} \int_0^{\infty} r dr d\theta = 2\pi \left(\frac{r^2}{2}\right)_0^{\infty} = 4\pi$$

註②: Jacobian determinant



$$J = \det \begin{pmatrix} \frac{\partial x}{\partial r} & \frac{\partial x}{\partial \theta} \\ \frac{\partial y}{\partial r} & \frac{\partial y}{\partial \theta} \end{pmatrix}$$

$$= \det \begin{pmatrix} \cos \theta & -r \sin \theta \\ \sin \theta & r \cos \theta \end{pmatrix} = r$$

原理:

① 令  $B = [\vec{e}_1, \vec{e}_2]$  的區域為  $\{x\vec{e}_1 + y\vec{e}_2 \mid 0 \leq x, y \leq 1\}$

經線性轉換  $T$  使  $A\vec{e}_i = \vec{a}_i$  故  $T(x\vec{e}_1 + y\vec{e}_2) = A(x\vec{e}_1 + y\vec{e}_2) = x\vec{a}_1 + y\vec{a}_2$

故面積  $= \det \begin{pmatrix} \vec{a}_1 & \vec{a}_2 \end{pmatrix} = \det \begin{pmatrix} A\vec{e}_1 & A\vec{e}_2 \end{pmatrix} = \det((AB)^T) = \det A \det B$

②  $\begin{bmatrix} du \\ dv \end{bmatrix} = \begin{bmatrix} \frac{\partial u}{\partial x} & \frac{\partial u}{\partial y} \\ \frac{\partial v}{\partial x} & \frac{\partial v}{\partial y} \end{bmatrix} \begin{bmatrix} dx \\ dy \end{bmatrix}$  展開  $\det \begin{pmatrix} \frac{\partial u}{\partial x} & \frac{\partial u}{\partial y} \\ \frac{\partial v}{\partial x} & \frac{\partial v}{\partial y} \end{pmatrix} (du dv) = \det(J(u,v)) du dv$

## Gaussian distribution 推導

### △ Step 1 求 $p(x)$

$$p(x, y) = p(x)p(y) = g(x)$$

$$\Rightarrow \frac{dg(x)}{dx} = \frac{d(\ln p(x))}{dx}$$

$$\Rightarrow 0 = \frac{dp(x)}{dx} p(y) + \frac{dp(y)}{dy} p(x), \begin{cases} x = r \cos \theta \Rightarrow \frac{dx}{dr} = -r \sin \theta = -y \\ y = r \sin \theta \Rightarrow \frac{dy}{dr} = r \cos \theta = x \end{cases}$$

$$\Rightarrow 0 = \frac{dp(x)}{dx} \frac{-y}{p(y)} + \frac{dp(y)}{dy} \frac{x}{p(x)}$$

$$\Rightarrow x p(x) \frac{dx}{dp(x)} = y p(y) \frac{dy}{dp(y)} = C \quad \because \text{若 } x \neq y \text{ 等式成立 則必為常數}$$

$$\Rightarrow \begin{cases} \frac{1}{C} x = \frac{dp(x)}{p(x)} \frac{1}{dx} \\ \frac{1}{C} y = \frac{dp(y)}{p(y)} \frac{1}{dy} \end{cases} \Rightarrow \begin{cases} \int \frac{1}{C} x dx = \int \frac{1}{p(x)} dp(x) \frac{1}{dx} dx = \ln p(x) \\ \int \frac{1}{C} y dy = \int \frac{1}{p(y)} dp(y) \frac{1}{dy} dy = \ln p(y) \end{cases}$$

$$\Rightarrow \begin{cases} \frac{2x^2}{C} = \ln p(x) \\ \frac{2y^2}{C} = \ln p(y) \end{cases} \Rightarrow \begin{cases} p(x) = e^{\frac{2}{C} x^2} = e^{-kx^2} \\ p(y) = e^{\frac{2}{C} y^2} = e^{-ky^2} \end{cases}$$

$$p(x) = \frac{1}{\sqrt{\pi}} e^{-kx^2}, \quad k = \frac{1}{2\sigma^2} \frac{\mu \cdot \sigma^2}{\text{given}} \frac{1}{\int_{-\infty}^{\infty} e^{-\frac{x^2}{2\sigma^2}} dx} \frac{1}{\frac{\sqrt{\pi}}{2\sigma}} \frac{1}{\int_{-\infty}^{\infty} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx} = x \sim N(\mu, \sigma^2)$$

### △ Step 2: Normalize $p(x) \rightarrow p(x)' = A p(x)$

$$\text{已知 } A \int_{-\infty}^{\infty} e^{-kx^2} dx = 1, \quad A \int_{-\infty}^{\infty} e^{-ky^2} dy = 1$$

$$\Rightarrow \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{-k(x^2+y^2)} dx dy = \frac{1}{A^2} = \int_0^{2\pi} \int_0^{\infty} e^{-kr^2} r dr d\theta$$

$$\Rightarrow 2\pi \left( \frac{e^{-kr^2}}{-2k} \right)_0^{\infty} = 2\pi \left( 0 + \frac{1}{2k} \right) = \frac{\pi}{k} = \frac{1}{A^2}$$

$$\Rightarrow A = \sqrt{\frac{k}{\pi}} \Rightarrow p(x)' = \sqrt{\frac{k}{\pi}} e^{-kx^2}$$

### △ Step 3 代 $\mu = 0, \sigma^2$

$$\int_{-\infty}^{\infty} x^2 \sqrt{\frac{k}{\pi}} e^{-kx^2} dx = \sigma^2$$

$$\Rightarrow \sqrt{\frac{k}{\pi}} \left[ x \frac{e^{-2kx^2}}{-2k} \right]_{-\infty}^{\infty} - \int_{-\infty}^{\infty} \frac{e^{-2kx^2}}{-2k} dx = \sigma^2$$

$$\Rightarrow \sqrt{\frac{k}{\pi}} \left( 0 + \frac{1}{2k} \int_{-\infty}^{\infty} e^{-2kx^2} dx \right) = \sigma^2$$

$$\Rightarrow \frac{1}{2k} = \sigma^2 \Rightarrow k = \frac{1}{2\sigma^2}$$

## MLE on Gaussian 必考

$D$ : 一组 data  $x_1, x_2, \dots, x_n$

$$L(\theta; \mu, \sigma^2 | D) = P(D | \theta) = \prod_{i=1}^n P(x_i | \theta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

$$\text{求 MLE: } \arg\max_{\theta} \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} \equiv \arg\max_{\theta} \sum_{i=1}^n \ln \left( \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} \right)$$

$$\begin{aligned} \text{令 } L &= \sum_{i=1}^n \ln \left( \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} \right) = -\frac{1}{2} \sum_{i=1}^n \ln(2\pi\sigma^2) + \sum_{i=1}^n \frac{-(x_i - \mu)^2}{2\sigma^2} \\ &= -\frac{n}{2} \ln(2\pi\sigma^2) + \sum_{i=1}^n \frac{-(x_i - \mu)^2}{2\sigma^2} \end{aligned}$$

$$\begin{aligned} \textcircled{1} \frac{dL}{d\mu} &= \frac{d}{d\mu} \sum_{i=1}^n \frac{-(x_i - \mu)^2}{2\sigma^2} \\ &\Rightarrow \frac{d}{d\mu} \sum_{i=1}^n (x_i^2 - 2\mu x_i + \mu^2) = 0 \\ &\Rightarrow 2 \sum_{i=1}^n x_i = 2 \sum_{i=1}^n \mu = 2n\mu \\ &\Rightarrow \mu = \frac{\sum_{i=1}^n x_i}{n} \end{aligned}$$

$$\begin{aligned} \textcircled{2} \frac{dL}{d\sigma} &= \frac{d}{d\sigma} \left( -\frac{n}{2} \ln(2\pi\sigma^2) + \sum_{i=1}^n \frac{-(x_i - \mu)^2}{2\sigma^2} \right) \\ \text{令 } s &= \sigma^2 \\ &= \left( -\frac{n}{2} \frac{2\pi}{2\pi s} \right) + \sum_{i=1}^n \frac{-(x_i - \mu)^2}{2s^2} = 0 \\ &= \frac{n}{2s} = \frac{\sum_{i=1}^n (x_i - \mu)^2}{2s^2} \\ &\Rightarrow s = \frac{\sum_{i=1}^n (x_i - \mu)^2}{n} \end{aligned}$$

## Conjugate Prior for Gaussian

Gaussian 也有 conjugate 性质: Prior: Gaussian  $\rightarrow$  Posterior: Gaussian

令 prior: Gaussian  $\sim N(\mu | \mu_0, \sigma_0^2)$

likelihood: Gaussian  $\sim N(D, \mu, \sigma^2)$ ,  $\sigma^2$  is fix

则 posterior: Gaussian  $\sim N(\mu | \circ \circ \circ)$ ?

解:

$$\textcircled{1} P(\theta | D) = \frac{P(D | \theta) P(\theta)}{P(D)} = \frac{P(D | \mu) P(\mu)}{P(D)}$$

$$\begin{aligned} \textcircled{2} P(D | \mu) P(\mu) &= \prod_{i=1}^n P(x_i | \mu) P(\mu | \mu_0, \sigma_0^2) \\ &= \left( \frac{1}{\sqrt{2\pi}\sigma} \right)^n e^{-\frac{1}{2\sigma^2} \sum (x_i - \mu)^2} \cdot \frac{1}{\sqrt{2\pi}\sigma_0} e^{-\frac{1}{2\sigma_0^2} (\mu - \mu_0)^2} \\ &= \left( \frac{1}{\sqrt{2\pi}\sigma} \right)^{n+1} e^{-\frac{\sum (x_i - \mu)^2}{2\sigma^2} - \frac{(\mu - \mu_0)^2}{2\sigma_0^2}} \rightarrow \text{目标: 符合 Gaussian form} \\ &\quad A e^{-\frac{(x - \mu)^2}{B}} = e^{-\frac{(x - \mu)^2}{B} + c} \end{aligned}$$

$$③ \frac{-\sum (x_i - \mu)^2}{2\sigma^2} + \frac{-(\mu - \mu_0)^2}{2\sigma_0^2} = \frac{-1}{2} \left[ \frac{\sum x_i^2}{\sigma^2} - \frac{2\mu \sum x_i}{\sigma^2} + \frac{\sum x_i^2}{\sigma^2} + \frac{\mu^2}{\sigma_0^2} - \frac{2\mu\mu_0}{\sigma_0^2} + \frac{\mu_0^2}{\sigma_0^2} \right]$$

假設：  
 $k' = \left( \frac{n}{\sigma^2} + \frac{1}{\sigma_0^2} \right)$   
 $\mu_n = \frac{\left( \frac{\sum x_i}{\sigma^2} + \frac{\mu_0}{\sigma_0^2} \right)}{k'}$   
 $k = \frac{k'}{2}$

$$= \frac{-1}{2} \left[ \mu^2 \left( \frac{n}{\sigma^2} + \frac{1}{\sigma_0^2} \right) - 2\mu \left( \frac{\sum x_i}{\sigma^2} + \frac{\mu_0}{\sigma_0^2} \right) + \left( \frac{\sum x_i^2}{\sigma^2} + \frac{\mu_0^2}{\sigma_0^2} \right) \right]$$

$$= \frac{-k'}{2} \left[ \mu^2 - 2\mu\mu_n + \mu_n^2 - \mu_n^2 + \left( \frac{\sum x_i^2}{\sigma^2} + \frac{\mu_0^2}{\sigma_0^2} \right) \right]$$

$$= -k \left[ (\mu - \mu_n)^2 \right] + C, \quad C = -k \left( \mu_n^2 - \frac{\sum x_i^2}{\sigma^2} + \frac{\mu_0^2}{\sigma_0^2} \right)$$

$$\therefore e^{\frac{-\sum (x_i - \mu)^2}{2\sigma^2} + \frac{-(\mu - \mu_0)^2}{2\sigma_0^2}} = e^{-k(\mu - \mu_n)^2 + C} = A e^{-k(\mu - \mu_n)^2}, \quad A = e^C$$

$$④ \text{ Marginal} = P(D) = \int_{-\infty}^{\infty} P(D|\mu') P(\mu') d\mu'$$

$$= \int_{-\infty}^{\infty} A e^{-k(\mu' - \mu_n)^2} d\mu' = \int_{-\infty}^{\infty} A e^{-k\mu'^2} d\mu'$$

$$= A \sqrt{\frac{\pi}{k}}$$

$$③④ \Rightarrow \frac{P(D|\mu) P(\mu)}{P(D)} = \frac{A e^{-k(\mu - \mu_n)^2}}{A \sqrt{\frac{\pi}{k}}} = \frac{e^{-\frac{k'}{2}(\mu - \mu_n)^2}}{\sqrt{\pi k'}}$$

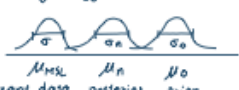
假設  $\sigma_n^2 = \frac{1}{k'} = \frac{1}{\left( \frac{n}{\sigma^2} + \frac{1}{\sigma_0^2} \right)}$

$$= \frac{e^{-\frac{1}{2\sigma_n^2}(\mu - \mu_n)^2}}{\sqrt{2\pi \sigma_n^2}} = N \sim (\mu_n, \sigma_n^2)$$

已知  $\mu_n = \sigma_n^2 \left( \frac{\sum x_i}{\sigma^2} + \frac{\mu_0}{\sigma_0^2} \right)$

$$= \frac{\sigma_n^2 n}{\sigma^2} \mu_{MLE} + \frac{\sigma_n^2}{\sigma_0^2} (\mu_0)$$

$\Rightarrow$  介於  $\mu_{MLE}$  和  $\mu_0$  間

$$\therefore \frac{\sigma_n^2 n}{\sigma^2} + \frac{\sigma_n^2}{\sigma_0^2} = 1$$


$\mu_{MLE}$  current data     $\mu_n$  posterior     $\mu_0$  prior

● 分析

$$P(\theta|D) = N \sim (\mu_n, \sigma_n^2) = \frac{e^{-\frac{1}{2\sigma_n^2}(\mu - \mu_n)^2}}{\sqrt{2\pi \sigma_n^2}}$$

其中  $\sigma_n^2 = \frac{1}{\left( \frac{n}{\sigma^2} + \frac{1}{\sigma_0^2} \right)}$ ,  $\mu_n = \sigma_n^2 \left( \frac{n\mu_{MLE}}{\sigma^2} + \frac{\mu_0}{\sigma_0^2} \right)$

$\therefore$  當  $n \rightarrow 0$  :  $\sigma_n^2 = \frac{1}{0 + \frac{1}{\sigma_0^2}} = \sigma_0^2$ ,  $\mu_n = \sigma_0^2 \left( 0 + \frac{\mu_0}{\sigma_0^2} \right) = \mu_0$  資料量=0 時  $\rightarrow N \sim (\mu_0, \sigma_0^2)$

當  $n \rightarrow \infty$  :  $\sigma_n^2 = 0$ ,  $\mu_n = \left( \frac{1}{\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2}} \right) \left( \frac{n\mu_{MLE}}{\sigma^2} + \frac{\mu_0}{\sigma_0^2} \right) = \mu_{MLE} = \frac{\sum x_i}{n}$  資料量和 sample space 一樣大  $\Rightarrow$  不需要先驗了