前言: Gaussin 的优点.
① symmetric
② unimodal : singal peak , mean = mode
③ localization : 有極值 at $\mu$ , Distance $(x_i . \mu)$ ↑ $P(x_i)$ ↓

## ♀ Central limit theorem 中央 極限定理

- $X \sim D(\mu, \sigma^2)$ 取 n 个 的 sample mean $\bar{X}$
  $\Rightarrow y \sim N(\mu_y, \sigma_y)$ ? $\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$ $\Rightarrow$ ⊔ if n→∞
  (D 為任一分佈 such as ⊔, ⋒, ⊑ )
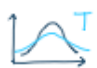
- that is : sample 1 的 mean = $\bar{X}_1$   sample 1 = $\{x_{11} ... x_{1n}\}$
            sample 2 的 mean = $\bar{X}_2$   sample 2 = $\{x_{21} ... x_{2n}\}$
            ...
            $\bar{X}_1, \bar{X}_2, ... \bar{X}_k \sim$ Normal Distribution

- 補充 MGF (moment generating function)
  把机率描述成 1st moment + 2nd moment + ...
  所以 (likelihood, prior) 可寫成 1 function $\Rightarrow$ 作 conjugate

## ♀ Students T distribution

- ⌢ᵀ 一個 Var 更大 (more tolerant) 的 distribution

- 因 Z-test 假設 $N \sim (\mu, \sigma^2)$ 需要知道 $\mu, \sigma^2$
  但取小樣本時, 只用 Ex 10 人建分佈很不準
  所以使用 $T = \int N(x|\mu, \sigma^2) \cdot T(\sigma^2|a,b) \, dab$

- 应用 : t-SNE    ? $\int N(x|\mu, \gamma^{-1}) \Gamma (\gamma|a, b) d\gamma$
                        Gamma

## ∮ multivariate

- univariate : $N = \frac{1}{\sqrt{2\pi}\sigma} e^{\frac{-1}{2\sigma^2}(x-\mu)^2}$

- multivariate : $P(\mu, \Sigma) = \frac{1}{\sqrt{2\pi}^k |\Sigma|^{\frac{1}{2}}} e^{\frac{-1}{2}(x-\mu)^T \boxed{\Sigma^{-1}} (x-\mu)}$  ※正定  ╮ precision matrix

$$\Sigma = \text{covariance matrix} = \begin{array}{c} x_1 \\ x_2 \\ \vdots \\ x_k \end{array} \begin{bmatrix} \frac{(x_1-\mu)^2}{n} & & \cdots & x_k \\ & \frac{(x_2-\mu)^2}{n} & \ddots & \\ & & \ddots & \\ & & & \frac{(x_k-\mu)^2}{n} \end{bmatrix}$$
otherwise $\frac{(x_m-\mu_m)(x_n-\mu_n)}{N}$

對角線 variance

$\Rightarrow$ If $\Sigma$ = diagonal matrix = $\begin{bmatrix} \diagdown & 0 \\ 0 & \diagdown \end{bmatrix}$ . it means anti-corelated

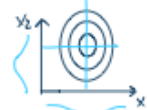## ∮ shape of data

- Isotropic
  $\Sigma = I$



- Orthogonal
  $\Sigma$ = Diagonal



  - $P(\mu, \Sigma) = \frac{1}{\sqrt{2\pi}^k |\Sigma|^{\frac{1}{2}}} e^{\frac{-1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$
    Euclidean distance
    欧式距離
  - 2 變数独立

- General
  $\Sigma = \begin{bmatrix} \diagdown & x \\ x & \diagdown \end{bmatrix}$
  non-zero



  - $P(\mu, \Sigma) = \frac{1}{\sqrt{2\pi}^k |\Sigma|^{\frac{1}{2}}} e^{\frac{-1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$
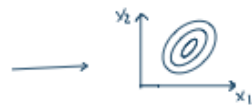    Mahalanobis distance
    馬式距離
  - 2 變数不独立

## ∮ Affine property : linear Trans (放大小/rotate) + 平移 = Affine Trans

- $X \sim N(\mu, \overset{c}{\sigma^2}) \xrightarrow[\text{affine transformation}]{f(x)} Ax+b \sim N(A\mu+b, ACA^T)$

$E(x) = \mu = \int x p(x) dx$
$\Rightarrow E(Ax+b) = \int (Ax+b) p(x) dx$
$\quad = A\int x p(x) dx + b\int P(x) dx$
$\quad = A E(x) + b \cdot 1$
$\quad = A\mu + b$

- Ex :



  Standard Normal
  $X \sim N(0, I)$ $\longrightarrow$ $Ax+\mu \sim N(\mu, c)$

$Cov(x) = \Sigma = E\{(x-\mu)(x-\mu)^T\}$
$Cov(Ax+b) = E\{[(Ax+b)-(A\mu+b)][(Ax+b)-(A\mu+b)]^T\}$
$\quad = E\{[A(x-\mu)][A(x-\mu)]^T\}$
$\quad = E\{A(x-\mu)(x-\mu)^T A^T\}$
$\quad = A E\{(x-\mu)(x-\mu)^T\} A^T$
$\quad = A\Sigma A^T$

- Any Gaussin can be derived from Isotropic by



  scale    rotate    relocate

# ♀ Univariate to multivariate Gaussion

可以線性變換要滿足 2 條件

$T(x+Y) = T(x) + T(Y) \sim T(ax) = aT(x)$

● ① $N(x+Y) = N(x) + N(Y)$



$X_1 \sim N(\mu_1, \sigma_1^2)$  $\qquad$ $X_2 \sim N(\mu_2, \sigma_2^2)$  $\qquad$ $Y = X_1 + X_2$

令 $Y = X_1 + X_2 \Rightarrow X = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix}$, $E(x) = \mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$  $\Sigma = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}$

取 $A = [1 \ 1]$, $b = 0$

則 by Affine : $AX + b = [1 \ 1] \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} + 0 = X_1 + X_2 \sim N(A\mu+b, A\Sigma A^T)$

∵ $A\mu + b = [1 \ 1] \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} + 0 = \mu_1 + \mu_2$

$\quad A\Sigma A^T = [1 \ 1] \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \sigma_1^2 + \sigma_2^2$

∴ $Y = X_1 + X_2 \sim N(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$

$\Rightarrow Y = \Sigma X_i \sim N(\Sigma \mu_i, \Sigma \sigma_i^2)$

● ② $N(Ax) = AN(x)$

令 $Y = b_1 X_1 + b_2 X_2 + \dots$

取 $A = [b_1 \ b_2 \dots]$

則 $AX + b = Y \sim N(A\mu+b, A\Sigma A^T)$

∵ $Y \sim N(b_1\mu_1 + b_2\mu_2 \dots, b_1\sigma_1^2 + b_2\sigma_2^2 + \dots)$

$\Rightarrow Y \sim N(B\Sigma\mu_i, B\Sigma\sigma_i^2)$

# ♀ Marginal Gaussin (muti to univariate)

令 $X = \begin{bmatrix} X_a \\ X_b \end{bmatrix} \begin{matrix} \rightarrow X_a = \begin{bmatrix} X_1 \\ \vdots \\ X_k \end{bmatrix} \\ \rightarrow X_b = \begin{bmatrix} X_{k+1} \\ \vdots \\ X_n \end{bmatrix} \end{matrix}$  $\mu = \begin{bmatrix} \mu_a \\ \mu_b \end{bmatrix}$  $\Sigma = \begin{bmatrix} \sigma_a^2 & \sigma_{ab}^2 \\ \sigma_{ba}^2 & \sigma_b^2 \end{bmatrix}$

取 $A = [I \ 0]$, $b = 0$

則 $AX + b = [I \ 0] \begin{bmatrix} X_a \\ X_b \end{bmatrix} = X_a \sim N(A\mu+b, A\Sigma A^T)$
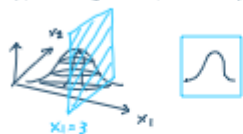
$\quad A\mu + b = \mu_a$

$\quad A\Sigma A^T = [I \ 0] \begin{bmatrix} \sigma_a^2 & x \\ x & x \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \sigma_a^2$
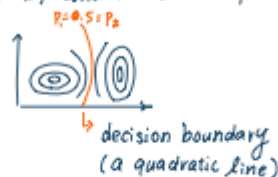
∴ $X_a \sim N(\mu_a, \sigma_a^2)$

$\quad X_b \sim N(\mu_b, \sigma_b^2)$

# $\varphi$ Conditional Gaussin

multivariate gaussin 中, 一部份變取為定值
其餘變取形成之分佈仍為Gaussin

註: If conditional independence



decision boundary
(a quadratic line)

令 $X = \begin{bmatrix} X_a \\ X_b \end{bmatrix}$, $P(\mu, \Sigma) = \frac{1}{\sqrt{2\pi}^k |\Sigma|^{0.5}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$

1. $-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu) = \frac{-1}{2} \begin{bmatrix} X_a-\mu_a & X_b-\mu_b \end{bmatrix} \begin{bmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{bmatrix} \begin{bmatrix} X_a-\mu_a \\ X_b-\mu_b \end{bmatrix}$

$= \frac{-1}{2} \begin{bmatrix} (x_a-\mu_a)\Lambda_{aa}+(x_b-\mu_b)\Lambda_{ba}, & (x_a-\mu_a)\Lambda_{ab}+(x_b-\mu_b)\Lambda_{bb} \end{bmatrix} \begin{bmatrix} X_a-\mu_a \\ X_b-\mu_b \end{bmatrix}$

$= \frac{-1}{2}(x_a-\mu_a)\Lambda_{aa}(x_a-\mu_a) + \frac{-1}{2}(x_b-\mu_b)^T\Lambda_{ba}(x_a-\mu_a) + \frac{-1}{2}(x_a-\mu_a)^T\Lambda_{ab}(x_b-\mu_b) + const$

$= \frac{-1}{2} X_a^T\Lambda_{aa}X_a + X_a^T\Lambda_{aa}\mu_a - X_a^T\Lambda_{ba}X_b + X_a^T\Lambda_{ba}\mu_b + const$

2. $\frac{-1}{2}(x-\mu)^T \Sigma_{xa|xb}^{-1}(x-\mu)$

$= \frac{-1}{2} X^T \Sigma_{xa|xb}^{-1} X + X^T \Sigma_{xa|xb}^{-1} \mu + const$

3. 比較 1., 2. 得 $\begin{cases} \Sigma_{xa|xb}^{-1} = \Lambda_{aa} \Rightarrow \Sigma_{xa|xb} = \Lambda_{aa}^{-1} \\ \\ X^T\Sigma_{xa|xb}^{-1}\mu = X_a^T\Lambda_{aa}\mu_a - X_a^T\Lambda_{ab}(X_b-\mu_b) \end{cases}$
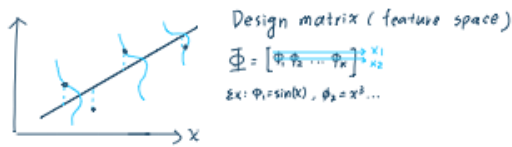
$\therefore \Sigma_{xa|xb}^{-1}\mu = \Lambda_{aa}\mu_a - \Lambda_{ab}(X_b-\mu_b)$

$\Rightarrow \mu = \underset{\Lambda_{aa}^{-1}}{\underline{\Sigma_{xa|xb}}}(\Lambda_{aa}\mu_a - \Lambda_{ab}(X_b-\mu_b))$

$= \mu_a - \Lambda_{aa}^{-1}\Lambda_{ab}(X_b-\mu_b)$

4. $\begin{bmatrix} A & B \\ C & D \end{bmatrix}^{-1} = \begin{bmatrix} M & -MBD^{-1} \\ \vee & \times \end{bmatrix}$    $M = (A-BD^{-1}C)^{-1}$

$\begin{bmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{bmatrix}^{-1} = \begin{bmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{bmatrix} \Rightarrow \begin{cases} \Lambda_{aa} = M \Rightarrow (\Lambda_{aa})^{-1} = M^{-1} \\ \\ \Lambda_{ab} = -MBD^{-1} \end{cases}$

5. 結合 3., 4. 得 $\mu = \mu_a - (-BD^{-1})(X_b-\mu_b) = \mu_a + \Sigma_{ab}\Sigma_{bb}(X_b-\mu_b)$

## ◊ Probability view of linear regression

● LSE 前情提要



Design matrix (feature space)
$$\Phi = [\phi_1 \phi_2 \cdots \phi_k] \begin{matrix} x_1 \\ x_2 \end{matrix}$$
Ex: $\phi_1 = \sin(x)$, $\phi_2 = x^3 \cdots$

Line : $w^\top \phi(x)$, $D = \{(x_1, y_1) \cdots (x_d, y_d)\}$

$$y = w^\top \phi(x) + \varepsilon \sim N \sim (y \mid \underline{w^\top \Phi(x)}, \sigma^2) \quad {}^{0 \to w^\top \Phi(x) = \text{line}}$$
or
$$y = Xw = \begin{bmatrix} 1 & x_1 \cdots y_1^n \\ & \vdots \\ 1 & x_d \cdots x_d^n \end{bmatrix} \begin{bmatrix} w_0 \\ \vdots \\ w_n \end{bmatrix}$$

● 重點提要

| frequentist | bayesian |
|---|---|
| LSE | ⟺ MLE |
| rLSE | ⟺ MAP |
| | fully - gaussian ↓ |
| | Predictive distribution |

## ◊ LSE ⟺ MLE

○ 目標 : 找出最大可能 fit 目前 data 的 $w$

上課版: Likelihood = $P(D \mid w)$
$$= \prod_{\forall x} N(w^\top \Phi(x_d), \sigma^2)$$
$$= (\frac{1}{\sqrt{2\pi}\sigma})^n \prod_D e^{\frac{-1}{2\sigma^2}[y_d - w^\top \Phi(x_d)]^2} \quad \propto \prod_D e^{\frac{-1}{2\sigma^2}[y_d - w^\top \Phi(x_d)]^2}$$
$$\xrightarrow{\log} \propto \sum_D [\underbrace{y_k}_{b} - \underbrace{w^\top \Phi(x_d)}_{Ax}]^2$$

講義版: Likelihood = $P(D_y \mid D_x, W)$
$$= \prod_i \frac{1}{\sqrt{2\pi}\sigma_i} e^{\frac{-1}{2\sigma_i^2}(y_i - X_i W)^2} \propto \prod_i e^{\frac{-1}{2\sigma_i^2}(y_i - X_i W)^2}$$

$$\xrightarrow{\log} \sum_i \lg \frac{1}{\sqrt{2\pi}\sigma_i^2} + \sum_i \underwave{\frac{-1}{2\sigma_i^2}(y_i - X_i W)^2} \quad \text{对 } W \text{ 微分後剩此項}$$

求 MLE ⟺ 求 $\max \left( \sum_i (X_i W - y_i)^2 \right)$ ⟺ 求 $\max (\|A\vec{x} - \vec{b}\|^2)$

## ◊ rLSE ⟺ MAP
{ 但比 rLSE 更好
∵ 利用 conjugate gaussian 性質
此法可持續更新 $w$
直至收斂

● 目標 : $\min E(w) = \sum_{n=1}^{N} [y(x_n, w) - t_n]^2 + \lambda \|w\|^2$
找出 $w$ 使 posterior 最大化，$\lambda \|w\|^2$ 相當 prior

- 前言
  posterior $P(w|D) = \dfrac{P(D|w)\,P(w)}{P(D)}$

  prior $P(w) \sim N(0, b^{-1}I)$, $b^{-1}$ 為 precision matrix

  $b$ 其實是 prior 之 cov-matrix 的倒数. 故

  $b^{-1} = \begin{bmatrix} \sigma^2_{w_0} & 0 & \cdots & 0 \\ 0 & \sigma^2_{w_1} & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma^2_{w_m} \end{bmatrix}$, $X = [1, x \ldots x^m]$, $y = \sum\limits_{k=0}^{n} w_i x^i = Xw$

- 求 max (prior · likelihood), 注意 prior 為 multivar, likelihood 為 univar

  已知 $\overset{\text{likelihood}}{P(D|w)} \propto e^{\sum\limits_i \frac{-1}{2\sigma_k^2}(y_i - x_i w)^2}$  　注意 posterior 為 multivar gaussian 分佈

  故 $P(w|D) \propto \overset{\text{likelihood}}{P(D|w)} \overset{\text{prior}}{P(w)}$ : MAP 和 MLE 即差在有 prior

  $\propto e^{\sum\limits_i \frac{-a}{2}(y_i - x_i w)^2} \cdot e^{\frac{-b}{2}(w-\vec{0})^T I (w-\vec{0})}$  　cov-matrix 倒数

  $= e^{\sum\limits_i \frac{-a}{2}(y_i - x_i w)^2 + \frac{-1}{2} w^T b I w}$

  取 $\ln$ 得 $P(w|D) \propto \sum\limits_i \frac{-a}{2}(y_i - x_i w)^2 + \frac{-1}{2} w^T b I w$

  matrix form 為 $\frac{-a}{2}\left( \|Xw - y\|^2 + \frac{b}{a} w^T w \right)$

  故可發現 $\frac{b}{a}$ 相当 rLSE 之 $\lambda$

  remind that $\begin{cases} x \text{ 解} = (A^T A + \lambda I)^{-1} A^T b \\ \lambda \text{大} \Rightarrow b\text{大}(var\text{小}) \quad \lambda\text{小} \Rightarrow b\text{小}(b\text{大}) \therefore b^{-1} \propto var \end{cases}$

- 求 $a\|Xw - y\|^2 + b w^T w$ quadratic form

  是為了証明 posterior 為 multivar gaussian distribution

  $a\|Xw - y\|^2 + b w^T w$

  $= a(Xw - y)^T(Xw - y) + b w^T w$

  $= a(w^T X^T X w - 2 w^T X^T y + y^T y) + b w^T w$

  $= w^T(a X^T X + b I)w - 2a w^T X^T y + a y^T y$

  对照 quadratic form $(x - \mu)^T \Lambda (x - \mu)$
  $\qquad\qquad = x^T \Lambda x - 2 x^T \Lambda \mu + \mu^T \mu$

  得出 $\Lambda = a X^T X + b I$
  $\qquad \mu = a \Lambda^{-1} X^T y$ ∵ $a w^T X^T y = w^T \Lambda \mu$

  故 poterior $\sim N(\mu, \Lambda^{-1})$

  註① $e^{w^T(a X^T X + b I)w - 2a w^T X^T y + a y^T y}$
  $\qquad = e^{(w^T \Lambda w - 2 w^T \Lambda \mu + \mu^T \mu) - \mu^T \mu + a y^T y}$  　当係数
  $\qquad = A\, e^{(w-\mu)^T \Lambda (w - \mu)}$

  註② $\mu = a \Lambda^{-1} X^T y$
  $\qquad = a(a X^T X - b I)^{-1} X^T y \qquad = (X^T X + \lambda I) X^T y$
  $\qquad = (\frac{a}{a} X^T X - \frac{a}{b} I) X^T y \qquad$ 故 rLSE 即 frequentist 版的 MAP

L09-2 : online learning

Iter 1 : 令 prior $\sim N(0, b^{-1}I)$，$b^{-1}$ 可隨意假設

則 posterior $\sim N(\mu, \Lambda^{-1})$

其中. $\begin{cases} \Lambda = ax^Tx + bI \\ \mu = a\Lambda^{-1}x^Ty \end{cases}$

Iter 2 : Prior 更新為 $(m, s^{-1})$

則 posterior $\sim N(\mu', \Lambda'^{-1})$

其中 $\begin{cases} \Lambda' = ax^Tx + bI \\ \mu' = \Lambda'^{-1}(ax^Ty + s \cdot m) \end{cases}$

推導过程 :

$P(w|D) \propto P(D|w) P(w)$

$\propto e^{\frac{-a}{2}\sum(x_iw \cdot y_i) + \frac{-1}{2}(w-m)^Ts(w-m)}$

取 ln 得 : $\frac{-a}{2}\sum(x_iw \cdot y_i) + \frac{-1}{2}(w-m)^Ts(w-m)$

矩陣 form : $\frac{-a}{2}\left[ \|Xw - y\|^2 + \frac{1}{a}(w-m)^Ts(w-m) \right]$

忽略係數 : $a(w^Tx^Txw - 2w^Tx^Ty + y^Ty) + (w^Tsw - 2w^Tsm + m^Tsm)$

$= w^T(ax^Tx + s)w - 2w^T(ax^Ty + sm) + ay^Ty + m^Tsm$

对照 quadratic form of $(w-\mu')^T\Lambda'(w-\mu')$

$= w^T\Lambda'w - 2w^T\Lambda'\mu' + \mu'^T\Lambda'\mu'$

故得 $\begin{cases} \Lambda' = ax^Tx + s \\ \mu' = \Lambda'^{-1}(ax^Ty + sm) \end{cases}$

註: $(ax^Ty - sm) = \Lambda'\mu'$

# ∮ Fully Baysian (Prediction distribution)

- review ⎡ MLE
         ⎣ MAP

- 



MLE: point estimation. $x \to y$

MAP: $x \to$ select $p(x)$ 最大者的 $y$

$$P(y|D) = \int \overset{likelihood}{P(y|w,P)} \overset{prior}{P(w|P_x)} d\omega$$

$$= \int N(y|Xw, \sigma^2) \overset{MAP}{N(w|\mu, \Lambda^{-1})} d\omega$$

① By Affine, $P(y|D)$ is gaussin $\Leftarrow$ $P(y|w,P)$ is gaussin

② $W$ 變換代表不同 possible line

③ marginalize $w$ 得 $P(y|D)$

$$\propto \int e^{-\frac{a}{2}(Xw-y)^2} e^{-\frac{1}{2}(w-\mu)^T \Lambda (w-\mu)} d\omega$$

$$= \int e^{-\frac{a}{2}(w^T X^T X w - 2w^T X^T y + y^T y) + -\frac{1}{2}(w^T \Lambda w - 2w^T \Lambda \mu + \mu^T \Lambda \mu)} d\omega$$

$$= \int e^{-\frac{1}{2}[w^T(aX^TX + \Lambda)w - 2w^T(aX^Ty + \Lambda\mu) + ay^Ty + \mu^T\Lambda\mu]} d\omega$$

> 對照 quadratic form of $(x-\mu)^T \Lambda (x-\mu)$
> $$= x^T \Lambda x - 2x^T \Lambda \mu + \mu^T \Lambda \mu$$
> 可令 $\begin{cases} C = aX^TX + \Lambda \\ \mu' = C^{-1}(aX^Ty + \Lambda\mu) \end{cases}$

$$= \int e^{-\frac{1}{2}[(w-\mu')^T C(w-\mu') - \mu'^T C\mu' + ay^Ty + \mu^T\Lambda\mu]} d\omega$$

$$= \int \underset{\text{multivar gaussian 積分後=1}}{\underbrace{e^{-\frac{1}{2}[(w-\mu')^T C(w-\mu')]}}} \cdot \underset{\text{当常數項提出}}{\underbrace{e^{\frac{1}{2}(\mu'^T C\mu' - ay^Ty - \mu^T\Lambda\mu)}}} d\omega$$

$$= e^{\frac{1}{2}(\cdot\mu'^T C\mu' + ay^Ty + \mu^T\Lambda\mu)}$$

$$= e^{-\frac{1}{2}(y-\mu'')^T C'(y-\mu'')}$$

其中 $\mu'' = X\mu$ , $(C')^{-1} = \frac{1}{a} + X\Lambda^{-1}X^T$

故分佈為 $N \sim (X\mu, \frac{1}{a} + X\Lambda^{-1}X^T)$

指數項展開：

$$-[c^{-1}(axy+\Lambda\mu)]^{\top}c'[c^{-1}(axy+\Lambda\mu)] + ay^{\top}y + \mu^{\top}\Lambda\mu$$

$$= -(axy+\Lambda\mu)^{\top}\underbrace{(c^{-1})^{\top}c'c^{-1}}_{=c^{-1}}(axy+\Lambda\mu) + ay^{\top}y + \mu^{\top}\Lambda\mu$$

$$= -a^2y^{\top}xc^{-1}x^{\top}y - 2ay^{\top}xc^{-1}\Lambda\mu - \mu^{\top}\Lambda c^{-1}\Lambda\mu + ay^{\top}y + \mu^{\top}\Lambda\mu$$

$$= y^{\top}\underline{(a-a^2xc^{-1}x^{\top})}y - 2\underline{ay^{\top}xc^{-1}\Lambda}\mu - \mu^{\top}\Lambda c^{-1}\Lambda\mu + \mu^{\top}\Lambda\mu$$

對照 quadratic form of $(y-\mu'')c'(y-\mu'')$

可得 
$$\begin{cases} c' = a - a^2xc^{-1}x^{\top} \\ \mu'' = c'^{-1}(axc^{-1}\Lambda\mu) = ac'^{-1}x(ax^{\top}x+\Lambda)^{-1}\Lambda\mu = \color{orange}{x\mu} \\ \because c'\mu'' = axc^{-1}\Lambda\mu \end{cases}$$

利用 sherman-Morrison formula

若 $C = ax^{\top}x + \Lambda$

則 $C^{-1} = \Lambda^{-1} - \dfrac{\Lambda^{-1}ax^{\top}x\Lambda^{-1}}{1+ax\Lambda^{-1}x^{\top}}$   check: $CC^{-1}=1$

故 $c' = a - a^2xc^{-1}x^{\top}$

$$= a - a^2x\left(\Lambda^{-1} - \dfrac{\Lambda^{-1}ax^{\top}x\Lambda^{-1}}{1+ax\Lambda^{-1}x^{\top}}\right)x^{\top}$$

令為 $\alpha$

$$= a - a^2x\left(\Lambda^{-1} - \dfrac{\Lambda^{-1}ax^{\top}\alpha(x^{\top})^{-1}}{1+a\alpha}\right)x^{\top}$$

$$= a - a^2(\underset{\alpha}{x\Lambda^{-1}x^{\top}} - \dfrac{x\Lambda^{-1}ax^{\top}\alpha}{1+a\alpha})$$

$$= a - a^2(\alpha - \dfrac{a\alpha^2}{1+a\alpha})$$

$$= a - a^2\left(\dfrac{\alpha}{1+a\alpha}\right)$$

$$= \dfrac{a}{1+a\alpha}$$

故 $(c')^{-1} = \dfrac{1+a\alpha}{a}$ $\color{orange}{= \dfrac{1}{a} + x\Lambda^{-1}x^{\top}}$

故 $\mu'' = a(c')^{-1}xc^{-1}\Lambda\mu$

$$= a\left[(c')^{-1}xc^{-1}\Lambda\mu\right]^{\top} \quad \text{取T}$$

$$= a\left[\mu^{\top}\Lambda c^{-1}x^{\top}(c')^{-1}\right]$$

$$= a\mu^{\top}\Lambda c^{-1}x^{\top}(c')^{-1}$$

$$= a\mu^{\top}\Lambda\left(\Lambda^{-1} - \dfrac{\Lambda^{-1}ax^{\top}x\Lambda^{-1}}{1+ax\Lambda^{-1}x^{\top}\alpha}\right)x^{\top}(c')^{-1}$$

$$= a\mu^{\top}\left(x^{\top} - \dfrac{ax^{\top}\alpha}{1+a\alpha}\right)\left(\dfrac{1+a\alpha}{a}\right)$$

$$= a\mu^{\top}\left(\dfrac{x^{\top}}{1+a\alpha}\right)\left(\dfrac{1+a\alpha}{a}\right) = \mu^{\top}x^{\top}$$

$$= x\mu \quad \text{取T}$$

## ♀ Decision Theory (classification)

- review : In either MLE or MAP,
  we estimate the params of $N(\mu, \sigma^2)$ according to observation

## ♀ Estimator

- This is a statistic (統計量) to approximate the property of a distribution
  也就是從抽樣樣本逼似母體樣本

1. 抽樣樣本 random var

 ① $\mu_{MLE} = \hat{\mu} = \frac{1}{n} \sum_i x_i$ if $D = x_1, \dots x_n$ is i.i.d.

 ② $\sigma_{MLE} = \hat{\sigma}$

 $\begin{cases} \sigma_{MLE}^2 = \hat{\sigma}^2 = \frac{1}{n} \sum (x_i - \mu)^2 \\ \sigma_{unbiased}^2 = \frac{1}{n-1} \sum (x_i - \mu)^2 \end{cases}$ , $\mu$ 為母體平均
 
 自由度

 Bias 來自 $E(\hat{\theta}) - \theta \neq 0$ , 用 $\sigma_{unbiased}$ , $E(\hat{\theta}) - \theta$ 才 $= 0$

2. 逼似母體

 ① $E(\hat{\mu}) = $ 樣本 mean 再取 mean $\Rightarrow \hat{\mu}$ 呈 gaussin ∵ 中央極限 th

  $= E(\frac{1}{n} \sum x_i) = \frac{1}{n} E(\sum x_i) = \frac{1}{n} n\mu = \mu$

 ② $Var(\hat{\mu}) = $ 樣本 mean 再取 var

  $= Var(\frac{1}{n} \sum x_i) = Var(\frac{x_1}{n} + \dots \frac{x_n}{n}) = \frac{1}{n^2} Var(x_1^2 + \dots x_n^2) = \frac{1}{n^2} n Var(x)$

  $= \frac{Var(x)}{n} = \frac{\sigma^2}{n}$

  ∴ 取樣次數越多, $\hat{\mu}$ 分佈越陡 越能精估母體 $\mu$

  $P(\hat{\mu})$

  $E(\hat{\mu}) = \mu$

L11-2 : Bias

## ↯ Bias

- Estimator 和母体差 $E(\hat{\theta})-\theta$

  $E(\hat{\mu})$ 前面推过, $= E(\frac{1}{n}\Sigma x_i) = \frac{1}{n}\Sigma E(x_i) = \frac{1}{n}n\mu = \mu$
  $\therefore$ bias $= 0$

  $Var(\hat{\mu})$ 前面推过 $= Var(\frac{x_1+\cdots x_n}{n}) = \frac{1}{n^2}Var(x_1+\cdots x_n)$
  $$= \frac{1}{n^2}(n\sigma^2) = \frac{\sigma^2}{n}$$

  $E(\sigma^2_{MLE}) = E[\frac{1}{n}\Sigma(x_i-\hat{\mu})^2] = \frac{1}{n}E[\Sigma x_i^2 - 2\hat{\mu}\overset{n\hat{\mu}}{\underline{\Sigma x_i}} + \hat{\mu}^2]$
  $$= \frac{1}{n}E[\Sigma x_i^2 - n\hat{\mu}^2]$$
  $$= \frac{1}{n}E(\Sigma x_i^2) - E(\hat{\mu}^2)$$
  $$= \underline{E(x_i^2)} - \underline{E(\hat{\mu}^2)}$$
  $$\quad\quad Var(x_i)+E^2(x_i) \quad Var(\hat{\mu})+E^2(\hat{\mu})$$
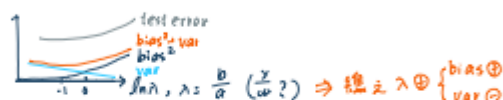  $$= (\sigma^2 + \mu^2) - (\frac{\sigma^2}{n} + \mu^2)$$
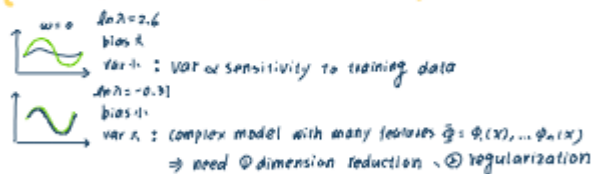  $$= \sigma^2 - \frac{\sigma^2}{n} \Rightarrow 总之 E(\sigma^2_{MLE}) = \sigma^2 - \frac{\sigma^2}{n} \neq \sigma^2$$

  故改为 $E(\sigma^2_{modified}) = E[\frac{1}{n-1}\Sigma(x_i-\hat{\mu})^2] = \frac{1}{n-1}E(\Sigma x_i^2 - n\hat{\mu}^2)$
  $$= \frac{1}{n-1}n(\sigma^2 + \hat{\mu}^2) - \frac{1}{n-1}n(\frac{\sigma^2}{n} - \hat{\mu}^2)$$
  $$= \frac{1}{n-1}(n-1)\sigma^2 = \sigma^2 \Rightarrow 总之 E(\sigma^2_{mod}) = \sigma^2$$
  $$\therefore bias = 0$$

## ↯ $MsE(\hat{\theta}) = bias^2(\hat{\theta}) + Var(\hat{\theta})$

- $E[(\hat{\theta}-\theta)^2] = E[(\hat{\theta}-\mu-(\theta-\mu))^2]$
  $$= E[(\hat{\theta}-\mu)^2] - 2E[(\hat{\theta}-\mu)(\theta-\mu)] + E[(\theta-\mu)^2]$$
  $$\quad\quad\quad\quad = E[(\hat{\theta}-\mu)](\theta-\mu)$$
  $$\quad\quad\quad\quad = [E(\hat{\theta})-\mu](\theta-\mu)$$
  $$= E[(\hat{\theta}-\mu)^2] + \underset{E(\hat{\theta})}{(\theta-\mu)^2}$$
  $$= Var(\hat{\theta}) + bias^2(\hat{\theta})$$
  $$= Var(\hat{\theta}) \text{ if } bias(\hat{\theta})=0$$

  

  test error
  $bias^2+var$
  $bias^2$
  $var$
  $\ell n\lambda$ , $\lambda \approx \frac{b}{a}$ $(\frac{\chi}{\omega}?)$ $\Rightarrow$ 总之 $\lambda$ 田 $\begin{cases} bias\oplus \\ var\ominus \end{cases}$

## ↯ Bias-Variance tradeoff



$w \approx 0$ $\ell n\lambda = 2.6$
bias $\nearrow$
var $\searrow$ : var $\propto$ sensitivity to training data

$\ell n\lambda = -0.31$
bias $\searrow$
var $\nearrow$ : complex model with many features $\hat{\phi} = \phi_1(x),\cdots \phi_n(x)$
$\Rightarrow$ need ① dimension reduction 、② regularization
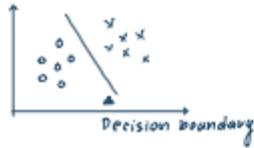
结论 : ① low variance : underfitting (simple)
- Regression
- Naive Bayes
- Linear model

② low bias : overfitting (complex)
- non-linear
- non parametric (no assumption)
- KNN

**L12-1：檢驗名詞**

## 前言



learn from data / Decision boundary

● Confusion matrix

$\theta_{label}$ , $\hat{\theta}_{model}$

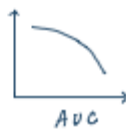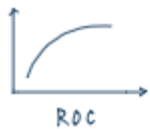|  | $\hat{\theta}=Y$ | $\hat{\theta}=N$ |  |
|---|---|---|---|
| $\theta=Y$ | TP | FN | Yes |
| $\theta=N$ | FP | TN | No |
|  | P | N |  |

$accuracy = \dfrac{TP+TN}{total}$

$1 \cdot accuracy = \dfrac{FP+FN}{total} = error\ rate$

$Mcc : \dfrac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP+FP)(TP+FN)(TN+FP)(TN+FN)}}$

$specificity = \dfrac{TN}{TN+FP}$ ( No)

$sensitivity = \dfrac{TP}{TP+FN}$ (power/recall) ( Yes)

$\Rightarrow F_1 - score = 2\dfrac{Precision \cdot Recall}{Precision + Recall}$

False positive rate FPR $\dfrac{TN}{TN+FP}$

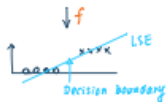Positive Prediction Rate ( Precision ) $= \dfrac{TP}{TP+FP}$

● ROC, AUC



ROC / AUC

**L12-2 : Regressian to classification**

## Regression to classification

● one-coding

indicator function $f\begin{cases}\hat{\theta}=1\\ \hat{\theta}=0\end{cases}$

↓f

LSE

Decision boundary

缺点 ① affected by outliers

有 bias !!

● 2 dimension

② 不能分成

● one-k-coding : multi class

$f\begin{cases}\theta=0\\ \hat{\theta},x=0\end{cases}$

課本 :

ideal

總之, we found that linear model is not that good

# ⚡ Loss Function Alternative

- Fisher Linear Discriminant (FLD)



# ⚡ Perceptron

perceptron ⟶ Logistic regression ⟶ nested regression
(used in neural network)



- Perception criterium

$$J = \sum max [\underbrace{w^T\phi}_{wrong} \overset{t \in |v-|}{(-t)}, \underbrace{0}_{right}]$$

∴ J 代表 wrong prediction



$w^T\phi$

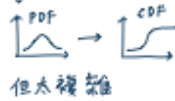註: NN  

# Logistic Regression

- 利用 sigmoid 作 activation function

  想法 : gaussion CDF

  

  但太複雜

  總之找到 logistic function $f(x) = \frac{1}{1+e^{-kx}}$

  

  $K=1$    $K=10$

- Probability point of view

  o Given $D = \{(x_i, y_i) \mid y = \{0,1\}\}$ s.t. $y_i \sim \text{Bernoulli}(f(\overset{x\,w}{w^T \phi}))$

  o 藉找 MLE , 找 $w$ :

  MLE function $= \prod \text{Bernoulli}(y \mid \overset{\theta}{w_i \phi})$

1. $\underset{w}{\arg\max}\ P(D \mid \overset{\theta}{w})$

   $= \underset{w}{\arg\max}\ \prod_i \left[ \left(\frac{1}{1+e^{-x_i w}}\right)^{y_i} \left(\frac{e^{-x_i w}}{1+e^{-x_i w}}\right)^{1-y_i} \right]$

2. $\Rightarrow J = \sum_{i=1}^{n} \left[ y_i \log\left(\frac{1}{1+e^{-x_i w}}\right) + (1-y_i) \log\left(\frac{e^{-x_i w}}{1+e^{-x_i w}}\right) \right]$

   $\Rightarrow \frac{\delta J}{\delta w_j} = \sum \left[ y_i \frac{\delta}{\delta w_j} \log\left(\frac{1}{1+e^{-x_i w}}\right) + (1-y_i)\frac{\delta}{\delta w_j} \log\frac{e^{-x_i w}}{1+e^{-x_i w}} \right]$

   ① $\frac{\delta}{\delta w_j} \log\left(\frac{1}{1+e^{-x_i w}}\right) = \frac{-\delta}{\delta w_j} \log(1+e^{-x_i w}) = \frac{(x_{ij})e^{-x_i w}}{1+e^{-x_i w}}$

   ② $\frac{\delta}{\delta w_j}(1-y_i)\log\left(\frac{e^{-x_i w}}{1+e^{-x_i w}}\right) = (1-y_i)\frac{\delta}{\delta w_j}\left[ \log e^{-x_i w} - \log(1+e^{-x_i w}) \right]$

   $= (1-y_i)\left( \frac{-x_{ij}e^{-x_i w}}{e^{-x_i w}} + \frac{x_{ij}e^{-x_i w}}{1+e^{-x_i w}} \right)$

   $= (1-y_i)\frac{-x_{ij}}{1+e^{-x_i w}}$

3. $\frac{\delta J}{\delta w_j} = \sum_{i=1}^{n} \left( \frac{y_i x_{ij} e^{-x_i w} - x_{ij} + y_i x_{ij}}{1+e^{-x_i w}} \right)$

   $= \sum_{i=1}^{n} \left[ x_{ij}\left( y_i - \frac{1}{1+e^{-x_i w}} \right) \right]$

4. 令 $\frac{\delta J}{\delta w_j} = 0$，以 Newton's method $x_{n+1} = x_n - H^{-1}f(x_n) \nabla f(x_n)$

① $\nabla f = \begin{bmatrix} \frac{\delta J}{\delta w_1} \\ \vdots \\ \frac{\delta J}{\delta w_k} \end{bmatrix} = 0$，$\Phi = \begin{bmatrix} \Phi_1 \cdots \Phi_d \\ x_{11} & x_{1d} \\ x_{21} & x_{2d} \\ \vdots & \vdots \\ x_{n1} & x_{nd} \end{bmatrix}$

$= \Phi^T \left( \frac{1}{1 + e^{-w^T\Phi}} - y \right)$

$\underset{\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}}{\underbrace{\quad}} \quad \underset{\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}}{\underbrace{\quad}}$

② $H = \begin{bmatrix} \frac{\delta^2 J}{\delta w_1 \delta w_1} & \frac{\delta^2 J}{\delta w_1 w_2} \\ \frac{\delta^2 J}{\delta w_1 \delta w_2} & \frac{\delta^2 J}{\delta w_2 \delta w_2} \\ & & \ddots \end{bmatrix}$

For entry $\frac{\delta}{\delta w_k} \frac{\delta}{\delta w_j} J$ :

$= \frac{\delta}{\delta w_k} \left[ \sum_{i=1}^{n} x_{ij} \left( y_i - \frac{1}{1 + e^{-x_i w}} \right) \right]$

$= \frac{-\delta}{\delta w_k} \sum_{i=1}^{n} \frac{x_{ij}}{1 + e^{-x_i w}}$

$= \frac{-\delta}{\delta w_k} \left( \frac{x_{ij}}{1 + e^{-x_{i1}w_1}} + \frac{x_{ij}}{1 + e^{-x_{i2}w_2}} + \cdots + \frac{x_{ij}}{1 + e^{-x_{ik}w_k}} + \cdots \right)$

$= \frac{x_{ij} x_{ik} e^{-x_{ik}w_k}}{\underline{(1 + e^{-x_{ik}w_k})^2}}$

可寫成 $H = \Phi^T D \Phi$

③ $\sum_{i=1}^{n} x_{ij} x_{ik} (1 + e^{-x_{ik}w_k})^{-2} e^{-x_{ik}w_k}$ 寫成 $H = \Phi^T D \Phi$ 的 matrix form

$\Phi = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1d} \\ x_{21} & x_{22} & \cdots & x_{2d} \\ \vdots & & & \vdots \\ x_{d1} & \cdots & \cdots & x_{dd} \end{bmatrix}$ $\quad D = \begin{bmatrix} \frac{e^{-x_{i1}w_1}}{(1 + e^{-x_{i1}w_1})^2} & & 0 \\ & \ddots & \\ 0 & & \frac{e^{-x_{id}w_d}}{(1 + e^{-x_{id}w_d})^2} \end{bmatrix}$

$x_{ij}$ : jth column $\begin{bmatrix} & \\ \end{bmatrix}^{i=1,2,\cdots,d}$ $\quad x_{ik}$ : kth column $\begin{bmatrix} x_{1k} \\ x_{2k} \\ \vdots \\ x_{dk} \end{bmatrix}$

最後取 inverse 代入 $x_{n+1} = x_n - H^{-1}f(x_n) \nabla f(x_n)$