

# Classification

黃志勝 (Tommy Huang)

義隆電子 人工智慧研發部

國立陽明交通大學 AI學院 合聘助理教授

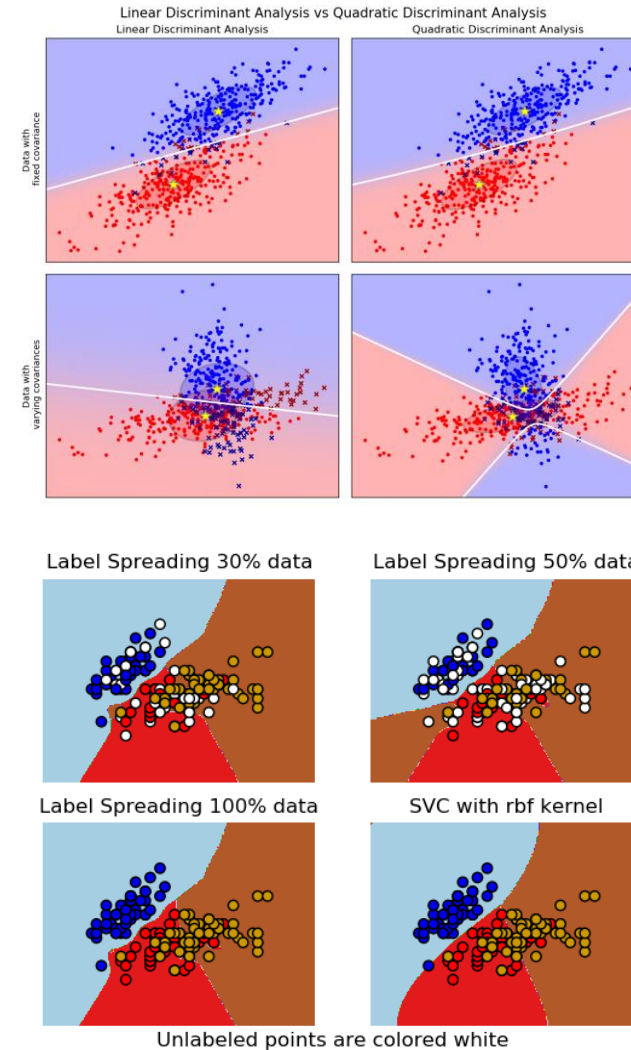
國立台北科技大學 電資學院合聘助理教授



# Classification

Identifying to which category an object belongs to.

- **Logistic Regression**
- **Linear and Quadratic Discriminant Analysis**
- Support Vector Machine
- Nearest neighbors
- Random forest
- Neural Network



# Classification

A Very simple classification problem

“How to classify {male or female} by a measured feature (body fat)?”

**Collected data (body fat(%))**

Female:{22, 25, 30, 33, 35}

Male:{ 10, 15, 20, 25, 30}



# Classification

Female: {22, 25, 30, 33, 35}

Male: { 10, 15, 20, 25, 30}

The simplest way:

Using mean value as decision rule.

$$\frac{\text{Mean value (Female)} + \text{Mean value (Male)}}{2} = \frac{29 + 20}{2} = 24.5$$

Body fat > 24.5 → Female

Body fat < 24.5 → Male



# Classification

Female: {22, 25, 30, 33, 35}

Male: { 10, 15, 20, 25, 30}

The simplest way:

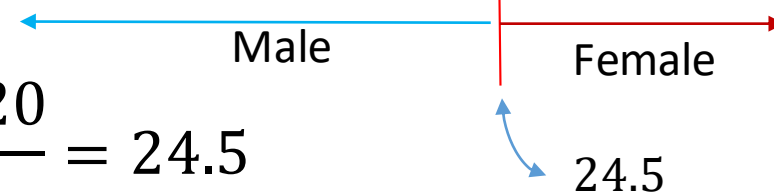
Using mean value as decision rule.

$$\frac{\text{Mean value (Female)} + \text{Mean value (Male)}}{2} = \frac{29 + 20}{2} = 24.5$$

Body fat > 24.5 → Female

Body fat < 24.5 → Male

Male	10	15	20	25	30	
Female			22	25	30	35
分布	-15	15- 20	20- 25	25- 30	30- 35	35-



# Classification

Female with 100 data, Male with 100 data  
(Body fat).

Visualization by histogram.

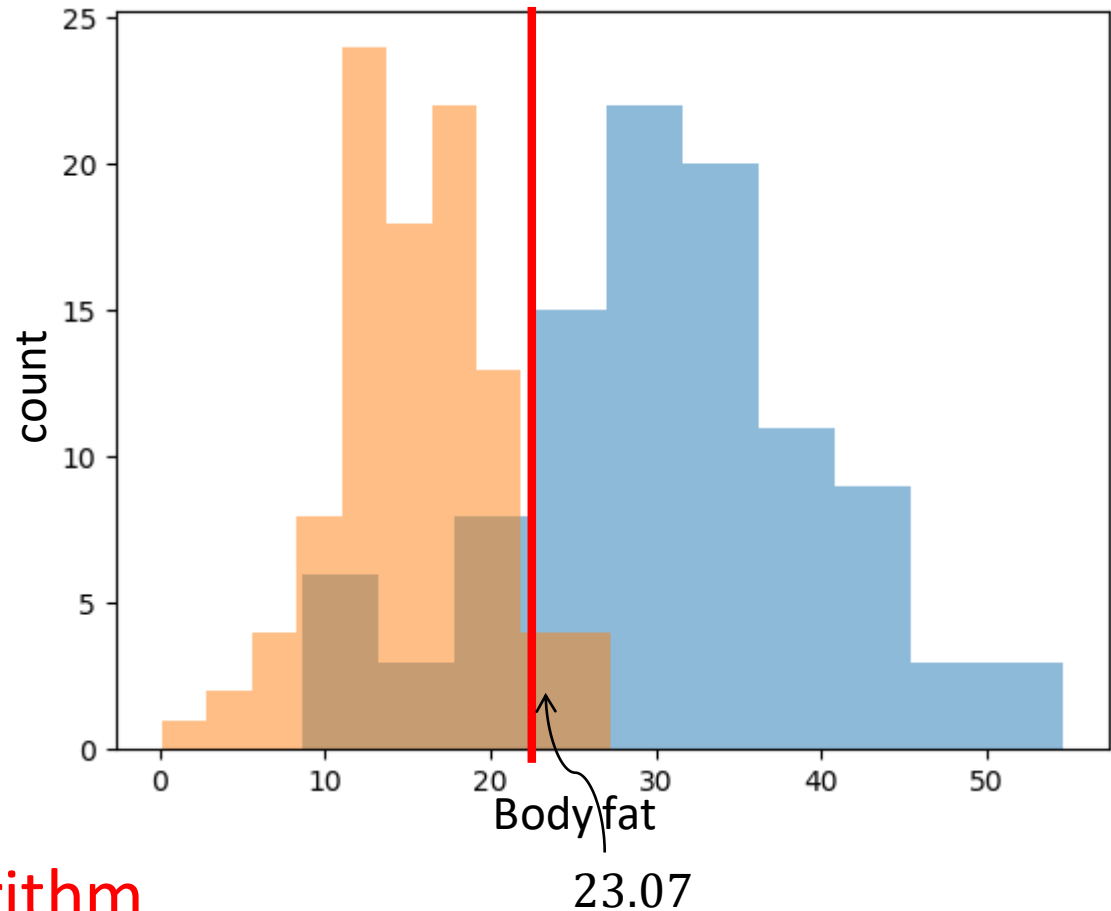
Blue: Male

Red: Female

$$\frac{\text{Mean value (Female)} + \text{Mean value (Male)}}{2}$$

$$= \frac{30.79 + 15.35}{2} = 23.07$$

You Just learn a classification algorithm



# Classification (平均數法)

$\{x_i\}, \forall i, x: \text{baby fat}$

$$\mu_c = \frac{1}{n_c} \sum_{i=1}^{n_c} x_i, c = \{\text{male}, \text{female}\}$$

$$f_{\text{male}}(x) = x - \mu_{\text{male}}$$

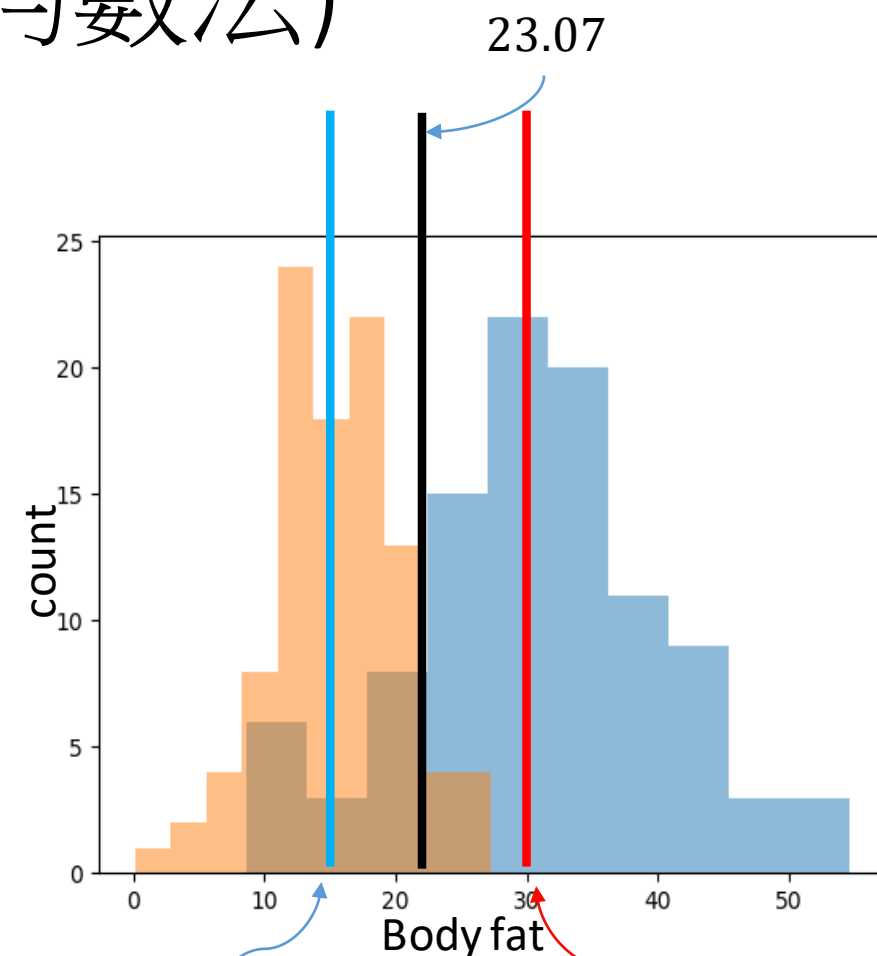
$$f_{\text{female}}(x) = x - \mu_{\text{female}}$$

Decision rule: feature value( $x$ ) is closed to which class, and classify this  $x$  to which class.

Decision rule:

$Decision(x)$

$$= \begin{cases} \text{female} & f_{\text{male}}(x) - f_{\text{female}}(x) \geq 0 \\ \text{male} & f_{\text{male}}(x) - f_{\text{female}}(x) < 0 \end{cases}$$



Mean value (Male)  
15.35

Mean value (Female)  
30.79



# Likelihood function(Single variable)

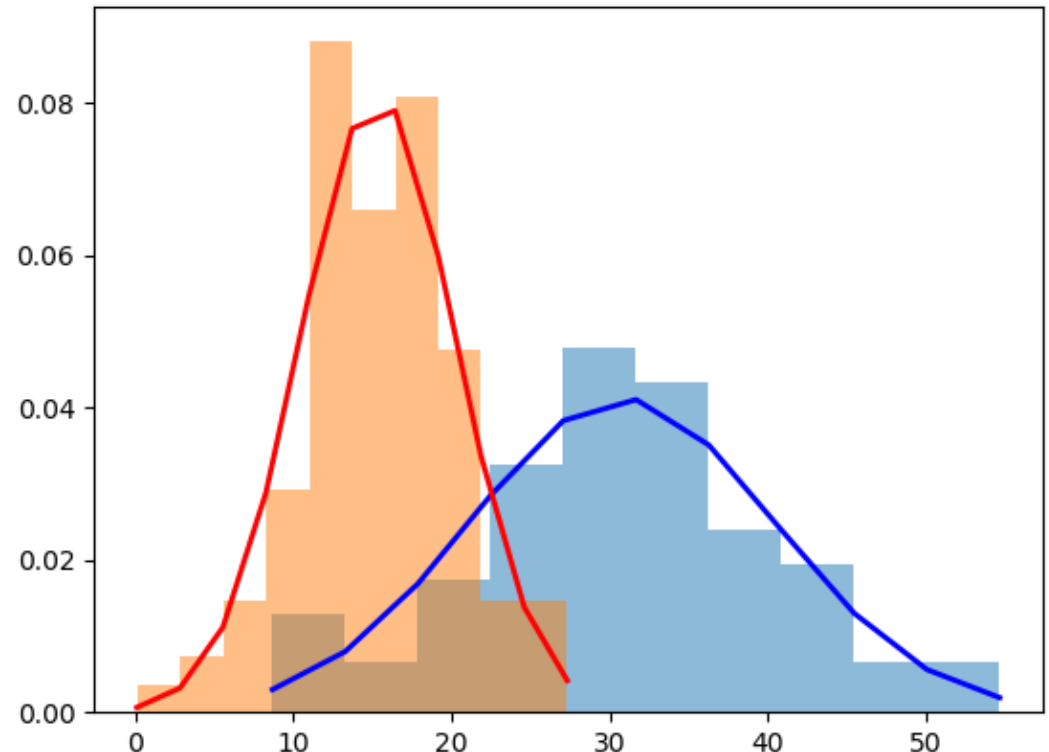
We can assume the histogram (density) is a **Gaussian** (normal)-like distribution.

That means

$$x_{male} \sim N(\mu_{male}, \sigma_{male})$$

$$x_{female} \sim N(\mu_{female}, \sigma_{female})$$

$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$





# Likelihood function(Single variable)

$$x_{male} \sim N(\mu_{male}, \sigma_{male})$$

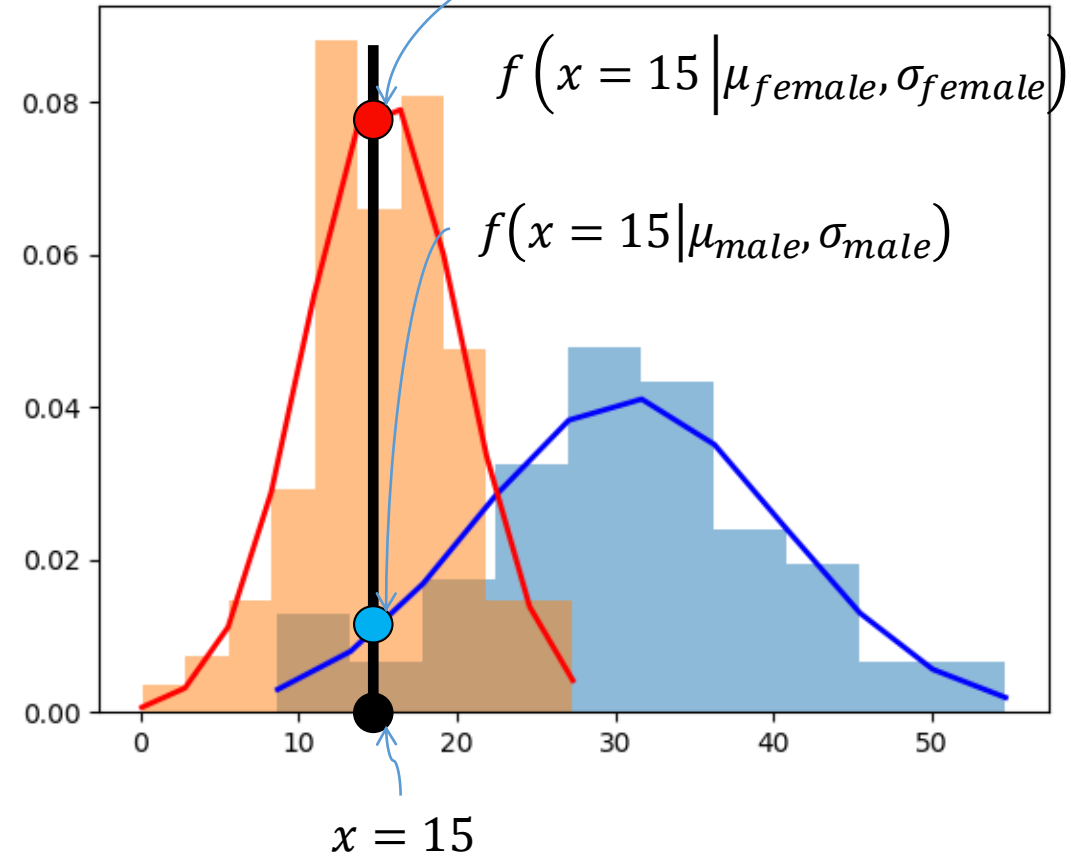
$$x_{female} \sim N(\mu_{female}, \sigma_{female})$$

$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

A unlabeled  $x$ : 15% body fat

$$f(x = 15 | \mu_{female}, \sigma_{female}) > f(x = 15 | \mu_{male}, \sigma_{male})$$

So this unlabeled  $x$  would be classify to Female.



Decision( $x$ )

$$= \begin{cases} \text{female} & f(x | \mu_{female}, \sigma_{female}) \geq f(x | \mu_{male}, \sigma_{male}) \\ \text{male} & f(x | \mu_{female}, \sigma_{female}) < f(x | \mu_{male}, \sigma_{male}) \end{cases}$$



# Classification (Multi-variables)(平均數法)

If we get multi-features (i.e. body fat and height), how to do?

$$\mathbf{x}_i = \begin{bmatrix} x_{bodyfat} \\ x_{height} \end{bmatrix}$$

$$\boldsymbol{\mu}_c = \frac{1}{n_c} \sum_{i=1}^{n_c} \mathbf{x}_i, = \begin{bmatrix} \mu_{bodyfat} \\ \mu_{height} \end{bmatrix}_c, c = \{male, female\}$$

$$f(\mathbf{x}) = \mathbf{x} - \boldsymbol{\mu} = \begin{bmatrix} x_{bodyfat} - \mu_{bodyfat} \\ x_{height} - \mu_{height} \end{bmatrix}$$

向量

純量

純量

We can't make decision with an array.

Decision(x)

$$= \begin{cases} \text{female} & f(x | \mu_{female}, \sigma_{female}) \geq f(x | \mu_{male}, \sigma_{male}) \\ \text{male} & f(x | \mu_{female}, \sigma_{female}) < f(x | \mu_{male}, \sigma_{male}) \end{cases}$$



# Classification (Multi-variables) (平均數法)

$$f(\mathbf{x}) = \mathbf{x} - \boldsymbol{\mu} = \begin{bmatrix} x_{bodyfat} - \mu_{bodyfat} \\ x_{height} - \mu_{height} \end{bmatrix}$$

Quantification ( $\mathbf{x}, \boldsymbol{\mu}$ )

- Euclidean Distance (L2-norm):  $\|\mathbf{x} - \boldsymbol{\mu}\|_{L2} = (\mathbf{x} - \boldsymbol{\mu})^T (\mathbf{x} - \boldsymbol{\mu})$   
 $\qquad\qquad\qquad 1 \times 2 \qquad\qquad\qquad 2 \times 1$
- Mahalanobis Distance



# Likelihood function(Multi-variables)

If we get multi-features (i.e. body fat and height), how to do?

$$\mathbf{x}_i = \begin{bmatrix} x_{bodyfat} \\ x_{height} \end{bmatrix}$$

純量

$$f(\mathbf{x}|\boldsymbol{\mu}, \Sigma)$$

$$= (2\pi)^{-d/2} |\Sigma|^{-0.5} \exp\{-0.5(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})\}$$

$$\text{Mahalanobis Distance} = (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})$$



# Likelihood function(Multi-variables)

$$f(\mathbf{x}|\boldsymbol{\mu}, \Sigma) = (2\pi)^{-d/2} |\Sigma|^{-0.5} \exp\{-0.5(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})\}$$

純量

$|\Sigma|$ : 共變異數的行列式值  $\rightarrow$  純量

$$(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})$$

$$\begin{matrix} 1 \times d & d \times d & d \times 1 \end{matrix}$$

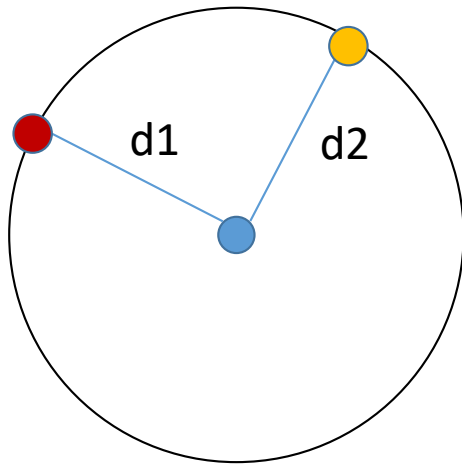
輸出為  $1 \times 1 \rightarrow$  純量



# Distance

Euclidean Distance

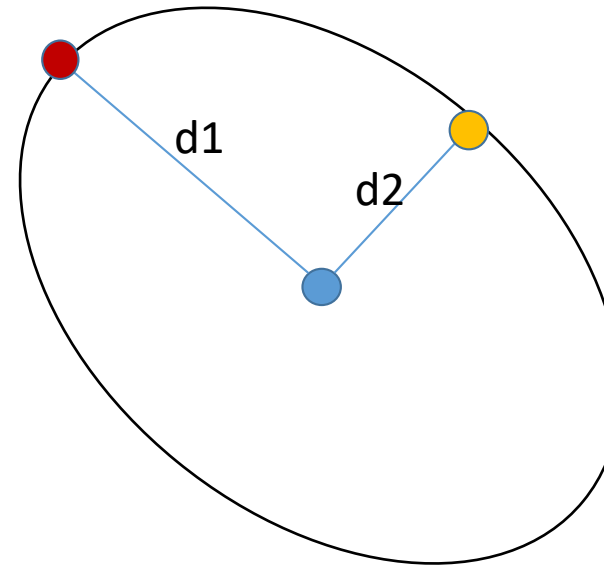
$$(x - \mu)^T (x - \mu)$$



$$d1 = d2$$

Mahalanobis Distance

$$(x - \mu)^T \Sigma^{-1} (x - \mu)$$



$$d1 = d2$$



# Maximum a posterior (MAP)

- In statistics, we hope to make decision by getting a maximum posterior probability for a given  $\mathbf{x}$ .

$p(c|\mathbf{x})$ : posterior probability of data  $\mathbf{x}$  for class  $c$ .

- Likelihood :  $f(\mathbf{x}|\mu,\sigma)$ . Given on a parameter set, the function output for the data  $\mathbf{x}$ .



# Maximum a posterior (MAP)

- $p(c|\mathbf{x})$ : posterior probability of data  $\mathbf{x}$  for class  $c$ .

$$p(c|\mathbf{x}) = \frac{p(c)p(\mathbf{x}|c)}{p(\mathbf{x})}$$

$p(c)$ : prior probability for class  $c$

$p(\mathbf{x}|c)$ : likelihood function for class  $c$

$p(\mathbf{x}) = \sum_{c=1}^L p(c) p(\mathbf{x}|c)$ : normalizing constant





# Maximum a posterior (MAP)

How to make decision?

By checking the posterior probability for all class.

$$\text{Decision}(\mathbf{x}) = \underset{c=\{1,2,\dots,C\}}{\operatorname{argmax}} \{p(c|\mathbf{x})\}$$



$$\underset{c=\{1,2,\dots,C\}}{\operatorname{argmax}} \{p(c)p(\mathbf{x}|c)\}$$



# Prior probability

$p(c)$ , in Chinese = 先驗機率 (在還沒有建模前，得到的先天訊息)



Male : Female = 25: 75



Sample a people.



Female  
75%

**Gender?**



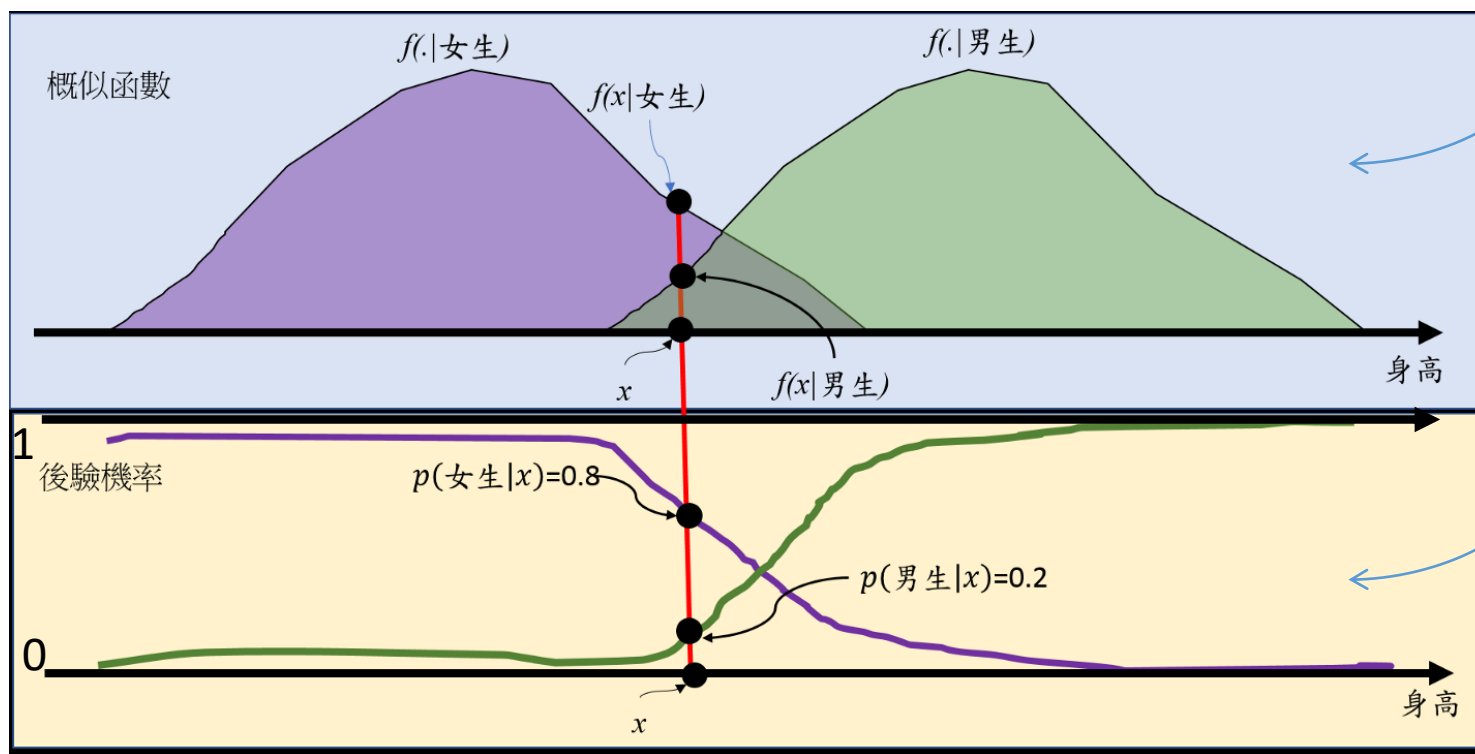
Male  
25%



# Maximum a posterior (MAP)

Posterior probability:  $p(c|\mathbf{x}) = \frac{p(c)p(\mathbf{x}|c)}{p(\mathbf{x})}$

Likelihood :  $f(\mathbf{x}|\mu, \sigma)$



Likelihood function

Likelihood function是積分面積  
總和為1。

Posterior probability

Posterior probability  
是單點在命題的總類別數機率  
合為1。



# Maximum a posterior (MAP)

Decision rule:

$$\begin{aligned} c_{MAP} &= \arg \max_{c=\{1,2,\dots,L\}} \{ p(c|x) \} \\ &= \arg \max_{c=\{1,2,\dots,L\}} \left\{ \frac{p(c)f(x|c)}{p(x)} \right\} = \arg \max_{c=\{1,2,\dots,L\}} \{ p(c)f(x|c) \} \end{aligned}$$



# MAP with Gaussian function

**Gaussian function:**

$$f(x|\mu_c, \Sigma_c) = (2\pi)^{-d/2} |\Sigma_c|^{-0.5} \exp \left\{ -0.5(x - \mu_c)^T \Sigma^{-1} (x - \mu_c) \right\}$$

**MAP:**

$$c_{MAP} = \arg \max_{c=\{1,2,\dots,L\}} \{ p(c|x) \} = \arg \max_{c=\{1,2,\dots,L\}} \{ p(c) f(x|c) \}$$



# MAP with Gaussian function 線性區別分析(linear discriminant classifier, LDC)

Gaussian function + MAP:

$$\begin{aligned} c_{MAP}^{GC} &= \arg \max_{c=\{1,2,\dots,L\}} \{ p(c|x) \} = \arg \max_{c=\{1,2,\dots,L\}} \{ p(c) f(x|c) \} \\ &= \arg \min_{c=\{1,2,\dots,L\}} \{ -2 \ln(p(c)) + \ln(|\Sigma_c|) - 0.5(x - \mu_c)^T \Sigma^{-1}(x - \mu_c) \} \end{aligned}$$



# MAP with Gaussian function

$$c_{MAP}^{GC} = \arg \min_{c=\{1,2,\dots,L\}} \{-2 \ln(p(c)) + \ln(|\Sigma_c|) - 0.5(x - \mu_c)^T \Sigma^{-1} (x - \mu_c)\}$$

The most important term of this formula is measuring the distance between  $x$  and center of distribution

Quadratic discriminant classifier      QDC:  $w_{MAP} = \arg \max_{i=\{1,2,\dots,L\}} \{\ln p(w_i) - 0.5 \ln |\Sigma_i| - 0.5(x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i)\}$

linear discriminant classifier      LDC:  $w_{MAP} = \arg \max_{i=\{1,2,\dots,L\}} \{\ln p(w_i) - 0.5(x - \mu_i)^T \Sigma^{-1} (x - \mu_i)\}$

minimum Euclidean classifier      MEC:  $w_{MAP} = \arg \max_{i=\{1,2,\dots,L\}} \{\ln p(w_i) - 0.5(x - \mu_i)^T I_d (x - \mu_i)\}$

$I_d$ : Identity matrix ( $d \times d$ )

## 範例: 假設男生和女生的先驗機率都是0.5

- 收集的訓練資料中，估計出男生和女生的身高與體重的平均數向量，身高和體重的共變異數矩陣

$$\mathbf{x} \sim N\left(\boldsymbol{\mu}_{\text{男生}} = \begin{bmatrix} 175 \\ 80 \end{bmatrix}, \boldsymbol{\Sigma}_{\text{男生}} = \begin{bmatrix} 10^2 & 30 \\ 30 & 5^2 \end{bmatrix}\right)$$

$$\mathbf{x} \sim N\left(\boldsymbol{\mu}_{\text{女生}} = \begin{bmatrix} 160 \\ 50 \end{bmatrix}, \boldsymbol{\Sigma}_{\text{女生}} = \begin{bmatrix} 8^2 & 50 \\ 50 & 9^2 \end{bmatrix}\right)$$





# 範例: 假設男生和女生的先驗機率都是0.5

- LDA 男生的隸屬函數為：

$$\mathbf{x} \sim N\left(\boldsymbol{\mu}_{\text{男生}} = \begin{bmatrix} 175 \\ 80 \end{bmatrix}, \boldsymbol{\Sigma}_{\text{男生}} = \begin{bmatrix} 10^2 & 30 \\ 30 & 5^2 \end{bmatrix}\right)$$

$$\begin{aligned} & -0.5 \ln(|\boldsymbol{\Sigma}_{C_{\text{男生}}}|) - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_{C_{\text{男生}}})^T \boldsymbol{\Sigma}_{C_{\text{男生}}}^{-1} (\mathbf{x} - \boldsymbol{\mu}_{C_{\text{男生}}}) \\ & = -0.5 \ln(1600) - \frac{1}{2} \left( \mathbf{x} - \begin{bmatrix} 175 \\ 80 \end{bmatrix} \right)^T \begin{bmatrix} 10^2 & 30 \\ 30 & 5^2 \end{bmatrix}^{-1} \left( \mathbf{x} - \begin{bmatrix} 175 \\ 80 \end{bmatrix} \right) \end{aligned}$$

$$\mathbf{x} \sim N\left(\boldsymbol{\mu}_{\text{女生}} = \begin{bmatrix} 160 \\ 50 \end{bmatrix}, \boldsymbol{\Sigma}_{\text{女生}} = \begin{bmatrix} 8^2 & 50 \\ 50 & 9^2 \end{bmatrix}\right)$$

$$\begin{aligned} & -0.5 \ln(|\boldsymbol{\Sigma}_{C_{\text{女生}}}|) - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_{C_{\text{女生}}})^T \boldsymbol{\Sigma}_{C_{\text{女生}}}^{-1} (\mathbf{x} - \boldsymbol{\mu}_{C_{\text{女生}}}) \\ & = -0.5 \ln(2684) - \frac{1}{2} \left( \mathbf{x} - \begin{bmatrix} 160 \\ 50 \end{bmatrix} \right)^T \begin{bmatrix} 8^2 & 50 \\ 50 & 9^2 \end{bmatrix}^{-1} \left( \mathbf{x} - \begin{bmatrix} 160 \\ 50 \end{bmatrix} \right) \end{aligned}$$



範例: 假設男生和女生的先驗機率都是0.5

$$\mathbf{x}^* = \begin{bmatrix} \text{身高} \\ \text{體重} \end{bmatrix} = \begin{bmatrix} 170 \\ 75 \end{bmatrix}$$

男生的隸屬函數值為

$$\begin{aligned} & -0.5 \ln(1600) - \frac{1}{2} \left( \begin{bmatrix} 170 \\ 75 \end{bmatrix} - \begin{bmatrix} 175 \\ 80 \end{bmatrix} \right)^T \begin{bmatrix} 10^2 & 30 \\ 30 & 5^2 \end{bmatrix}^{-1} \left( \begin{bmatrix} 170 \\ 75 \end{bmatrix} - \begin{bmatrix} 175 \\ 80 \end{bmatrix} \right) \\ & = -0.5 \ln(1600) - \frac{1}{2} \frac{1}{1600} \begin{bmatrix} -5 \\ -5 \end{bmatrix}^T \begin{bmatrix} 25 & -30 \\ -30 & 100 \end{bmatrix} \begin{bmatrix} -5 \\ -5 \end{bmatrix} \approx -4.1966915 \end{aligned}$$

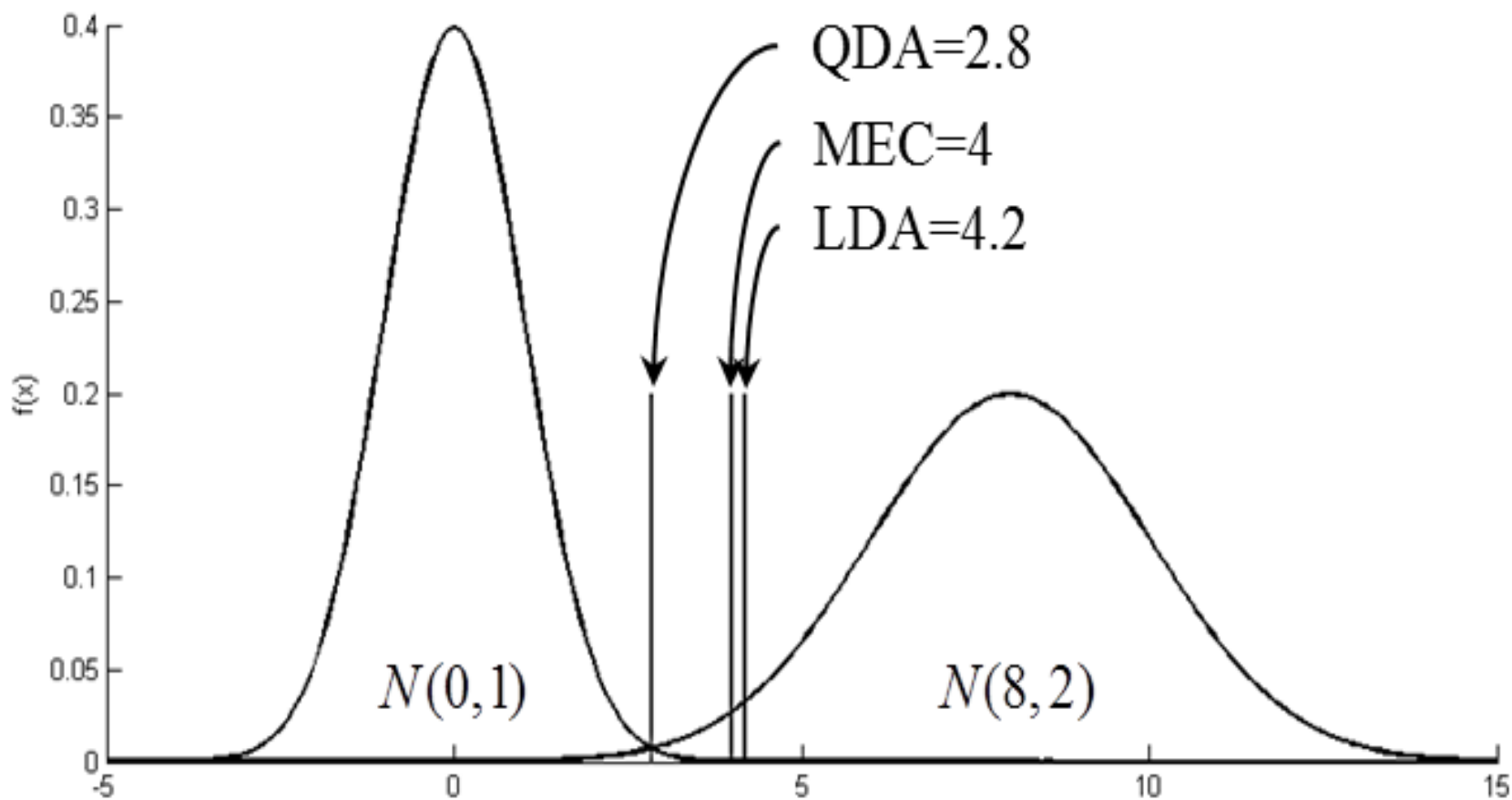
女生的隸屬函數值為

$$\begin{aligned} & -0.5 \ln(2684) - \frac{1}{2} \left( \begin{bmatrix} 170 \\ 75 \end{bmatrix} - \begin{bmatrix} 160 \\ 50 \end{bmatrix} \right)^T \begin{bmatrix} 8^2 & 50 \\ 50 & 9^2 \end{bmatrix}^{-1} \left( \begin{bmatrix} 170 \\ 75 \end{bmatrix} - \begin{bmatrix} 160 \\ 50 \end{bmatrix} \right) \\ & = -0.5 \ln(2684) - \frac{1}{2} \frac{1}{2684} \begin{bmatrix} 10 \\ 25 \end{bmatrix}^T \begin{bmatrix} 81 & -50 \\ -50 & 64 \end{bmatrix} \begin{bmatrix} 10 \\ 25 \end{bmatrix} \approx -8.250811 \end{aligned}$$

$$\begin{aligned} C_{QDA} \left( \mathbf{x}^* = \begin{bmatrix} 170 \\ 75 \end{bmatrix} \right) &= \operatorname{argmax} \{ C_{\text{男生}} = -4.1966915, C_{\text{女生}} = -8.250811 \} \\ &= C_{\text{男生}} \end{aligned}$$



# 範例：一維度的線性區別分析



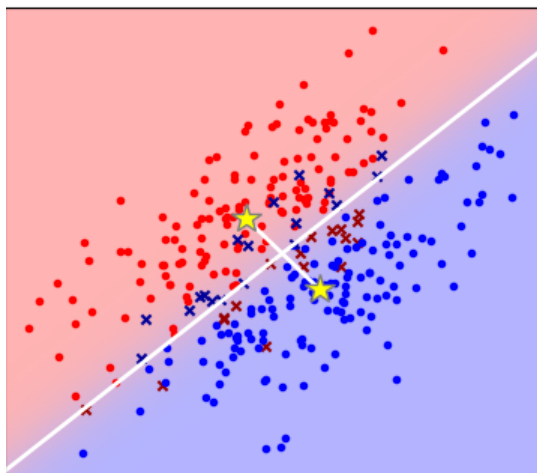
# 範例: 二維度的線性區別分析

類別的共變異數矩陣一致，但平均數不同

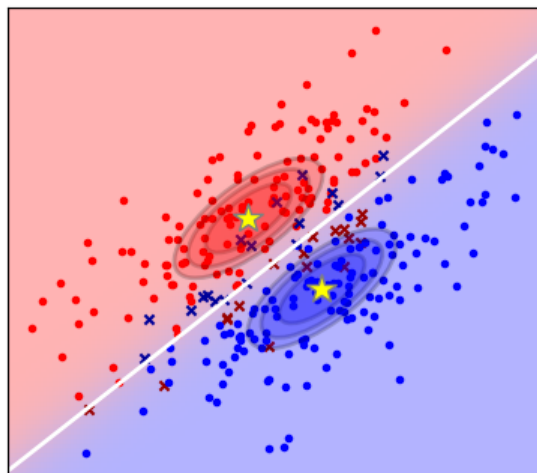
$$\text{類別 1 : } N\left(\boldsymbol{\mu}_{C_1} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \boldsymbol{\Sigma} = \begin{bmatrix} 2 & 1 \\ 1 & 1 \end{bmatrix}\right)$$

$$\text{類別 2 : } N\left(\boldsymbol{\mu}_{C_2} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma} = \begin{bmatrix} 2 & 1 \\ 1 & 1 \end{bmatrix}\right)$$

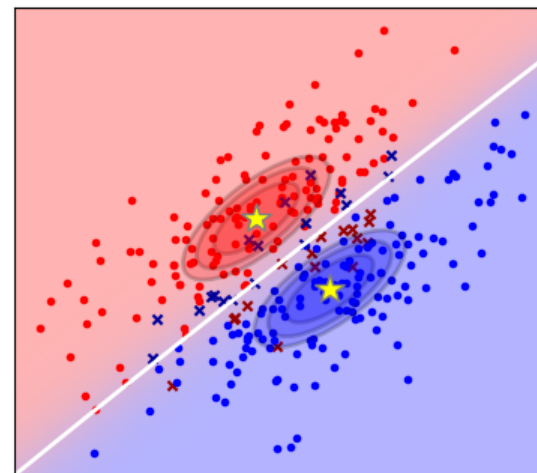
最小歐式距離分類器  
minimum Euclidean classifier



線性區別分析  
linear discriminant analysis



二次式區別分析  
quadratic discriminant analysis



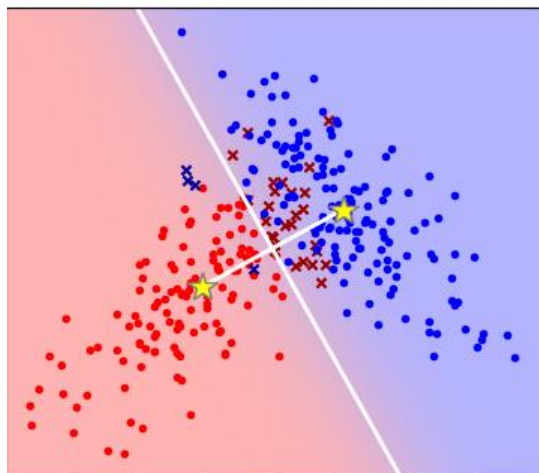
# 範例、二維度的線性區別分析

每個類別的共變異數矩陣和平均數皆不同

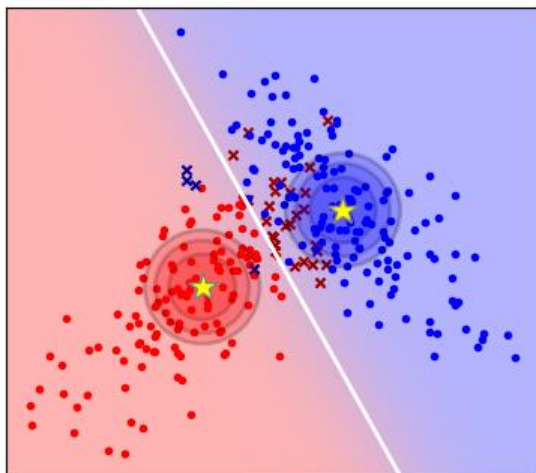
$$\text{類別 1: } N\left(\boldsymbol{\mu}_{C_1} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \boldsymbol{\Sigma}_1 = \begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix}\right)$$

$$\text{類別 2: } N\left(\boldsymbol{\mu}_{C_2} = \begin{bmatrix} 2 \\ 2 \end{bmatrix}, \boldsymbol{\Sigma} = \begin{bmatrix} 1 & -1 \\ -1 & 2 \end{bmatrix}\right)$$

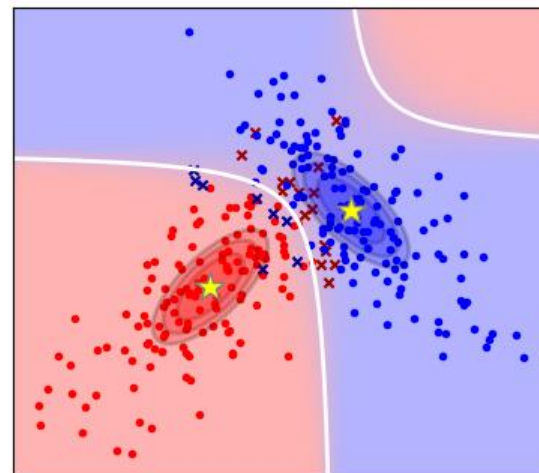
最小歐式距離分類器  
minimum Euclidean classifier



線性區別分析  
linear discriminant analysis



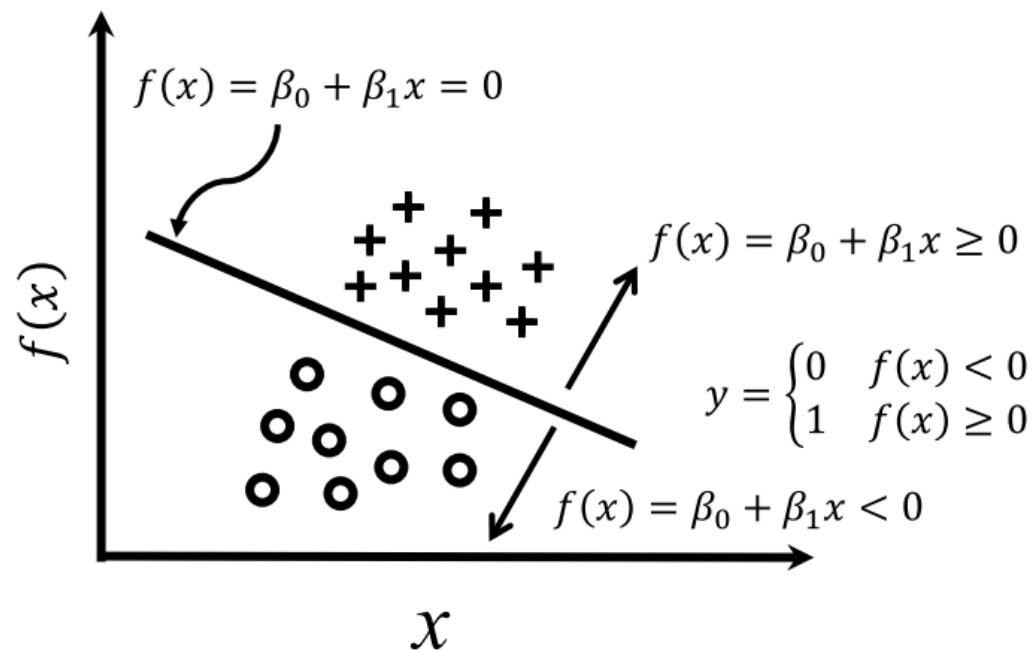
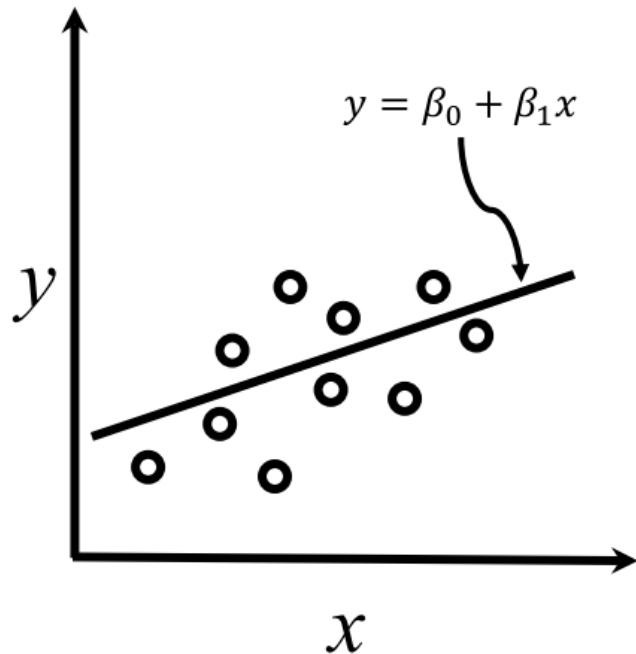
二次式區別分析  
quadratic discriminant analysis



# 羅吉斯迴歸 ( Logistic Regression )

迴歸

羅吉斯迴歸

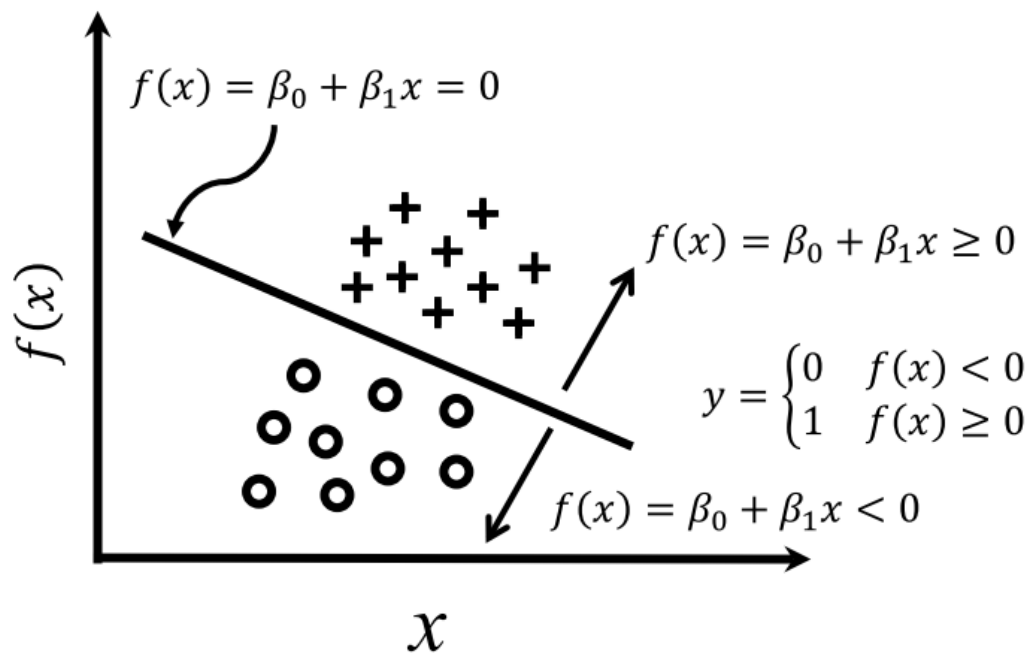


- 線性迴歸跟羅吉斯迴歸公式是一樣的 (但要分清楚，前者在算出數值，後者在做分類)



# 羅吉斯迴歸 (Logistic Regression)

羅吉斯迴歸



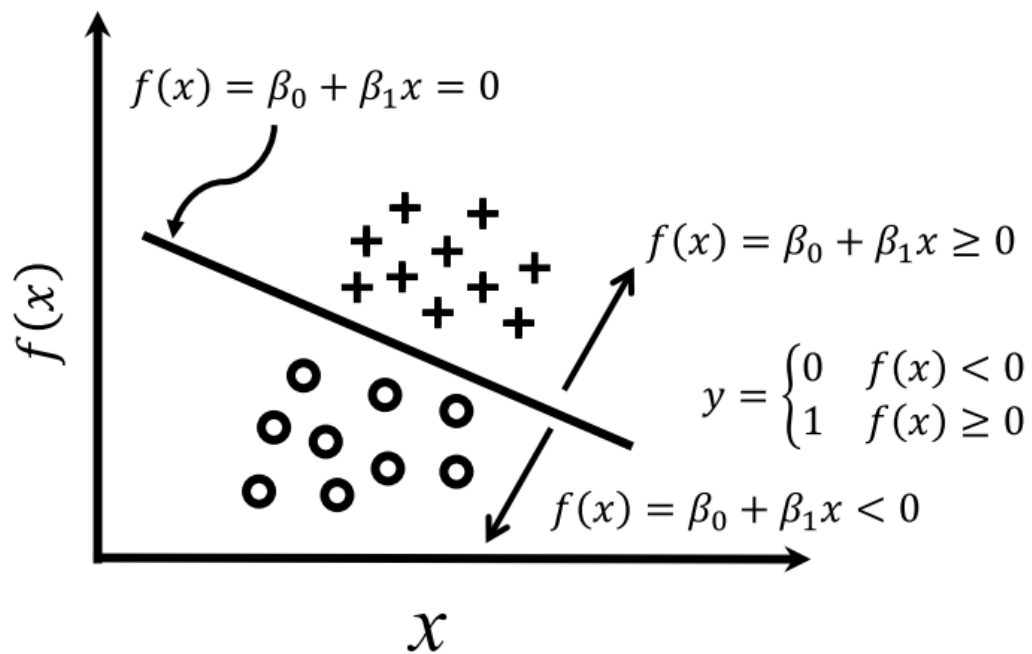
- 羅吉斯迴歸則是希望線性迴歸的輸出可以將兩類的資料能越區隔開越好。
- 最簡單的方式就是任意資料帶入迴歸方程式中判斷輸出值是否大於0，若大於0是一類(類別：1)，小於0則是另一類(類別：0)

$$y = \sigma(f(\mathbf{x})) = \begin{cases} 1 & f(\mathbf{x}) \geq 0 \\ 0 & f(\mathbf{x}) < 0 \end{cases}$$



# 羅吉斯迴歸 ( Logistic Regression )

羅吉斯迴歸



- $\sigma(\cdot)$ 在機器學習上稱為單位階梯函數(unit step function)，大於一個閾值(threshold)是一類，反之為另一類，此例的閾值為0。

$$y = \sigma(f(\mathbf{x})) = \begin{cases} 1 & f(\mathbf{x}) \geq 0 \\ 0 & f(\mathbf{x}) < 0 \end{cases}$$

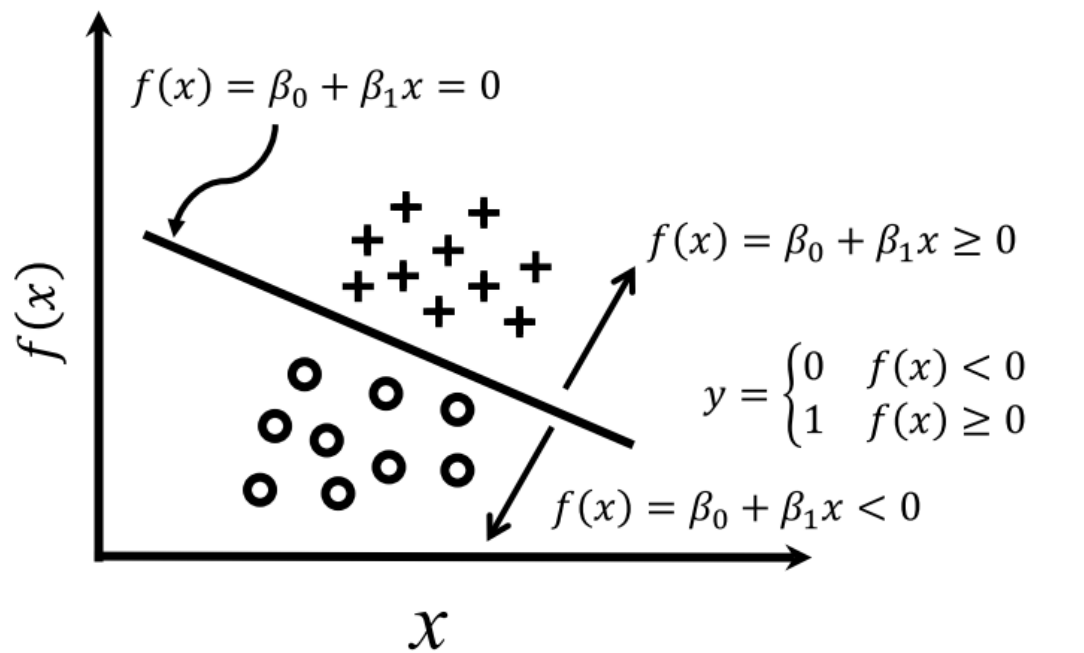




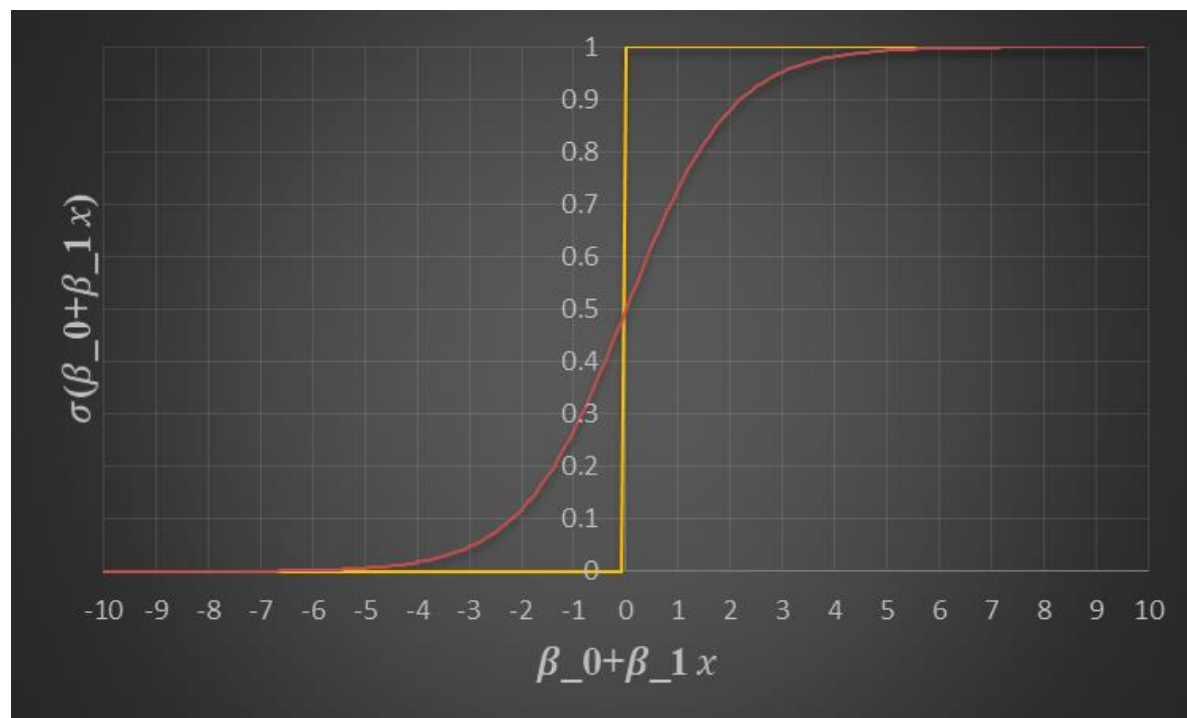
# 羅吉斯迴歸用Sigmoid函數限制值域

羅吉斯迴歸

$$\text{Sigmoid 函數: } s(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{1 + e^x}, x \in [-\infty, \infty], s(x) \in [0, 1]$$



$$\sigma(f(x)) = \frac{1}{1 + e^{-f(x)}} = \frac{e^{f(x)}}{1 + e^{f(x)}}$$



# 羅吉斯迴歸的公式為

$$s(f(\mathbf{x})) = \frac{1}{1 + e^{-f(\mathbf{x})}} = \frac{1}{1 + e^{-\mathbf{x}^T \boldsymbol{\beta}}}$$

或寫成

$$s(f(\mathbf{x})) = \frac{e^{f(\mathbf{x})}}{1 + e^{f(\mathbf{x})}} = \frac{e^{\mathbf{x}^T \boldsymbol{\beta}}}{1 + e^{\mathbf{x}^T \boldsymbol{\beta}}}$$

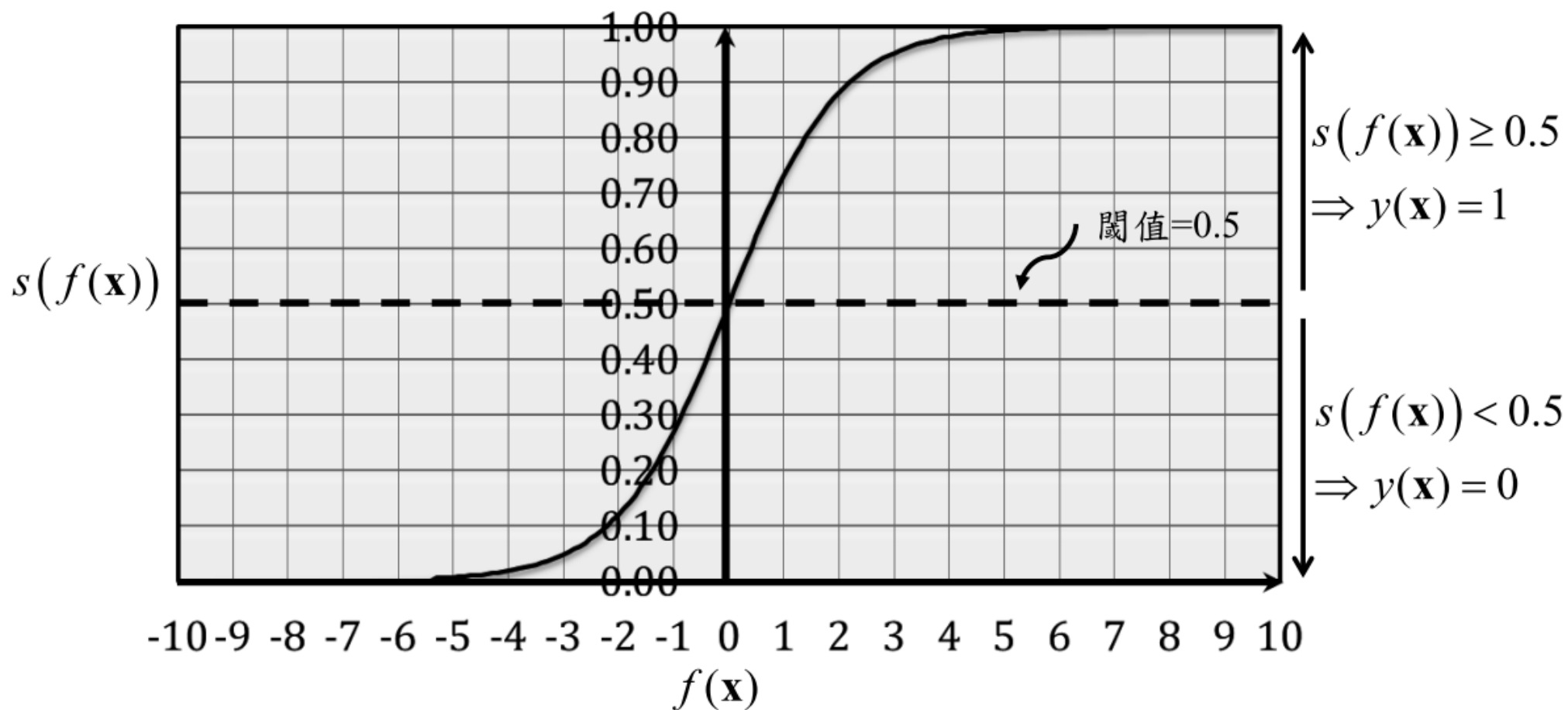
羅吉斯迴歸是二  
分類演算法

$$y = \sigma(s(f(\mathbf{x}))) = \begin{cases} 1 & s(f(\mathbf{x})) \geq 0.5 \\ 0 & s(f(\mathbf{x})) < 0.5 \end{cases}$$



# 羅吉斯迴歸的公式為

$$s(f(\mathbf{x})) = \frac{1}{1 + e^{-f(\mathbf{x})}} = \frac{e^{f(\mathbf{x})}}{1 + e^{f(\mathbf{x})}}$$



# 要怎麼找(求解)羅吉斯回歸參數( $\beta$ )?

$$y = \sigma(s(f(\mathbf{x}))) = \begin{cases} 1 & s(f(\mathbf{x})) \geq 0.5 \\ 0 & s(f(\mathbf{x})) < 0.5 \end{cases}$$

$$s(f(\mathbf{x})) = \frac{e^{f(\mathbf{x})}}{1 + e^{f(\mathbf{x})}} = \frac{e^{\mathbf{x}^T \beta}}{1 + e^{\mathbf{x}^T \beta}}$$

$$f(\mathbf{x}) = \mathbf{x}^T \beta$$

$$\mathbf{x} = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}, \quad \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_d \end{bmatrix}$$



# 要怎麼找(求解)羅吉斯回歸參數( $\beta$ )?

- 回顧一下伯努利機率函數，伯努利試驗結果為成功的機率為 $p$ ，不成功的機率即為 $1-p$

$$f(x) = p^x (1-p)^{1-x} = \begin{cases} p & x=1 \\ 1-p & x=0 \end{cases}$$

- 羅吉斯迴歸的輸出類別為**1**(成功)的機率是

$$p = p(y=1|\mathbf{x})$$

- 輸出類別為 **0** 的機率是

$$p(y=0|\mathbf{x}) = 1 - p(y=1|\mathbf{x}) = 1 - p$$



# 要怎麼找(求解)羅吉斯回歸參數( $\beta$ )?

- 有n組資料，其概似函數為

$$L(\beta) = \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{(1-y_i)}$$

$$p_i = p(y_i = 1 | \mathbf{x}_i), \forall i = 1, \dots, n$$

概似函數最大化→不好做

我們把問題轉成 $-\log$ ，找最小化的問題。

$p_i$ 是什麼

$$\mathcal{L}(\beta) = -\log L(\beta) = -\log \left( \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1-y_i} \right)$$

$$= - \sum_{i=1}^n \log \left( p_i^{y_i} (1 - p_i)^{1-y_i} \right)$$

$$= - \sum_{i=1}^n \left( \log(p_i^{y_i}) + \log((1 - p_i)^{1-y_i}) \right)$$

$$= - \sum_{i=1}^n \left( y_i \log(p_i) + (1 - y_i) \log(1 - p_i) \right)$$

$$= - \sum_{i=1}^n \left( y_i \log(p_i) + \log(1 - p_i) - y_i \log(1 - p_i) \right)$$

$$= - \sum_{i=1}^n \left( y_i \log \left( \frac{p_i}{1 - p_i} \right) + \log(1 - p_i) \right)$$

就是交叉熵的公式



# 要怎麼找(求解)羅吉斯回歸參數( $\beta$ )?

$p_i$ : 羅吉斯回歸的輸出

$$p_i = s(f(\mathbf{x}_i)) = \frac{e^{\mathbf{x}_i^T \boldsymbol{\beta}}}{1 + e^{\mathbf{x}_i^T \boldsymbol{\beta}}}$$

$$\ln(p_i) = \ln\left(\frac{e^{\mathbf{x}_i^T \boldsymbol{\beta}}}{1 + e^{\mathbf{x}_i^T \boldsymbol{\beta}}}\right) = \mathbf{x}_i^T \boldsymbol{\beta} - \ln(1 + e^{\mathbf{x}_i^T \boldsymbol{\beta}})$$

$$\ln(1 - p_i) = \ln\left(1 - \frac{e^{\mathbf{x}_i^T \boldsymbol{\beta}}}{1 + e^{\mathbf{x}_i^T \boldsymbol{\beta}}}\right) = \ln\left(\frac{1}{1 + e^{\mathbf{x}_i^T \boldsymbol{\beta}}}\right) = -\ln(1 + e^{\mathbf{x}_i^T \boldsymbol{\beta}})$$

$$\mathcal{L}(\boldsymbol{\beta}) = -\sum_{i=1}^n \left( y_i \log\left(\frac{p_i}{1 - p_i}\right) + \log(1 - p_i) \right) = -\sum_{i=1}^n \left[ y_i \mathbf{x}_i^T \boldsymbol{\beta} - \ln(1 + e^{\mathbf{x}_i^T \boldsymbol{\beta}}) \right]$$



# 要怎麼找(求解)羅吉斯回歸參數( $\beta$ )?

- 利用偏微分求得此函數的梯度(Gradient)

$$\begin{aligned}
 \frac{\partial \mathcal{L}(\beta)}{\partial \beta} &= \frac{\partial \sum_{i=1}^n \left( \log(1 + e^{\beta^T x_i}) - y_i \beta^T x_i \right)}{\partial \beta} \\
 &= \sum_{i=1}^n \left\{ \frac{\partial \left( \log(1 + e^{\beta^T x_i}) \right)}{\partial \beta} - \frac{\partial (y_i \beta^T x_i)}{\partial \beta} \right\} \\
 &= \sum_{i=1}^n \left\{ \frac{1}{(1 + e^{\beta^T x_i})} \times \frac{\partial (1 + e^{\beta^T x_i})}{\partial \beta} - y_i x_i \right\} \\
 &= \sum_{i=1}^n \left\{ \frac{e^{\beta^T x_i}}{(1 + e^{\beta^T x_i})} x_i - y_i x_i \right\} = \sum_{i=1}^n \{ p_i x_i - y_i x_i \} \\
 &= \sum_{i=1}^n \{ p_i - y_i \} x_i
 \end{aligned}$$





# 要怎麼找(求解)羅吉斯回歸參數( $\beta$ )?

- 梯度下降法

$$\boldsymbol{\beta}^{(t+1)} \leftarrow \boldsymbol{\beta}^{(t)} - \alpha \times \partial \frac{\mathcal{L}(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \boldsymbol{\beta}^{(t)} + \alpha \sum_{i=1}^n (y_i - p_i) \mathbf{x}_i^T$$

$\alpha$  為學習率

- 牛頓法（Newton's Method）求羅吉斯迴歸參數



# Recall

## 線性區別分析: (有closed-form solution)

- Quadratic discriminant classifier  
(距離: 馬氏距離, 不同類別用各自的共變異數矩陣)
- linear discriminant classifier  
(距離: 馬氏距離, 不同類別用共同的共變異數矩陣)
- minimum Euclidean classifier  
(平均數法: 歐式距離)

## 羅吉斯回歸: (沒有closed-form solution)

- 梯度下降法找解
- 牛頓法找解



# Classification

We just learned model-based algorithm.

Model-based: data is assumed as normal distribution. (parameters: mean vector and covariance matrix)

**Can we learn without model?**

ANS:

Yes.

Nearest neighbors, SVM, neural network.

