

# Classification: Probabilistic Generative Model

# Classification



例子

- Credit Scoring
  - Input: income, savings, profession, age, past financial history .....
  - Output: accept or refuse

- Medical Diagnosis
  - Input: current symptoms, age, gender, past medical history .....
  - Output: which kind of diseases

- Handwritten character recognition

Input:  output: 金

- Face recognition
  - Input: image of a face, output: person

# Example Application



$$f(\text{Pikachu}) = \text{Electric} \quad f(\text{Squirtle}) = \text{Water} \quad f(\text{Bulbasaur}) = \text{Grass}$$

pokemon games (*NOT* pokemon cards or Pokemon Go)



# Example Application

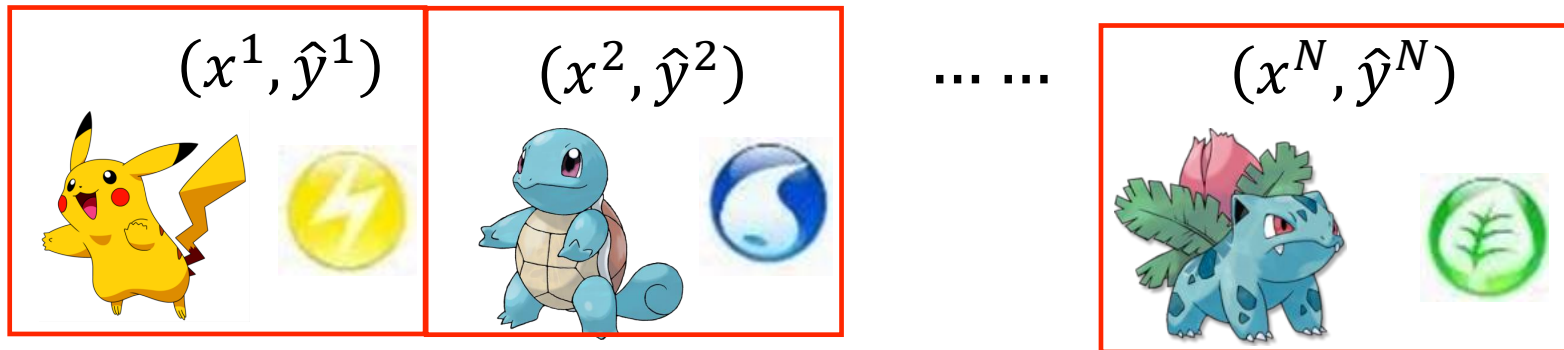
想要將寶可夢丟進 Function 中，就必須將寶可夢「數值化」=> 用 Vector 表示

- **Total:** sum of all stats that come after this, a general guide to how strong a pokemon is **320**
- **HP:** hit points, or health, defines how much damage a pokemon can withstand before fainting **35**
- **Attack:** the base modifier for normal attacks (eg. Scratch, Punch) **55**
- **Defense:** the base damage resistance against normal attacks **40**
- **SP Atk:** special attack, the base modifier for special attacks (e.g. fire blast, bubble beam) **50**
- **SP Def:** the base damage resistance against special attacks **50**
- **Speed:** determines which pokemon attacks first each round **90**

Can we predict the “type” of pokemon based on the information?

# How to do Classification

- Training data for Classification

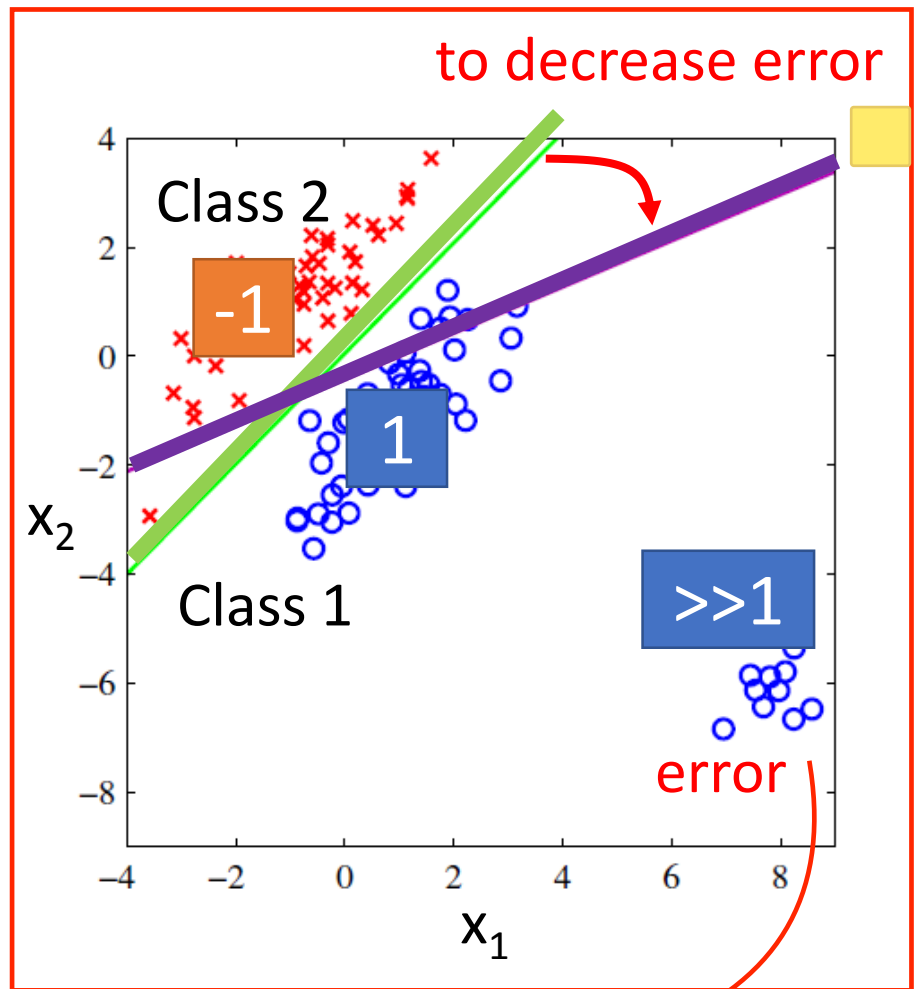
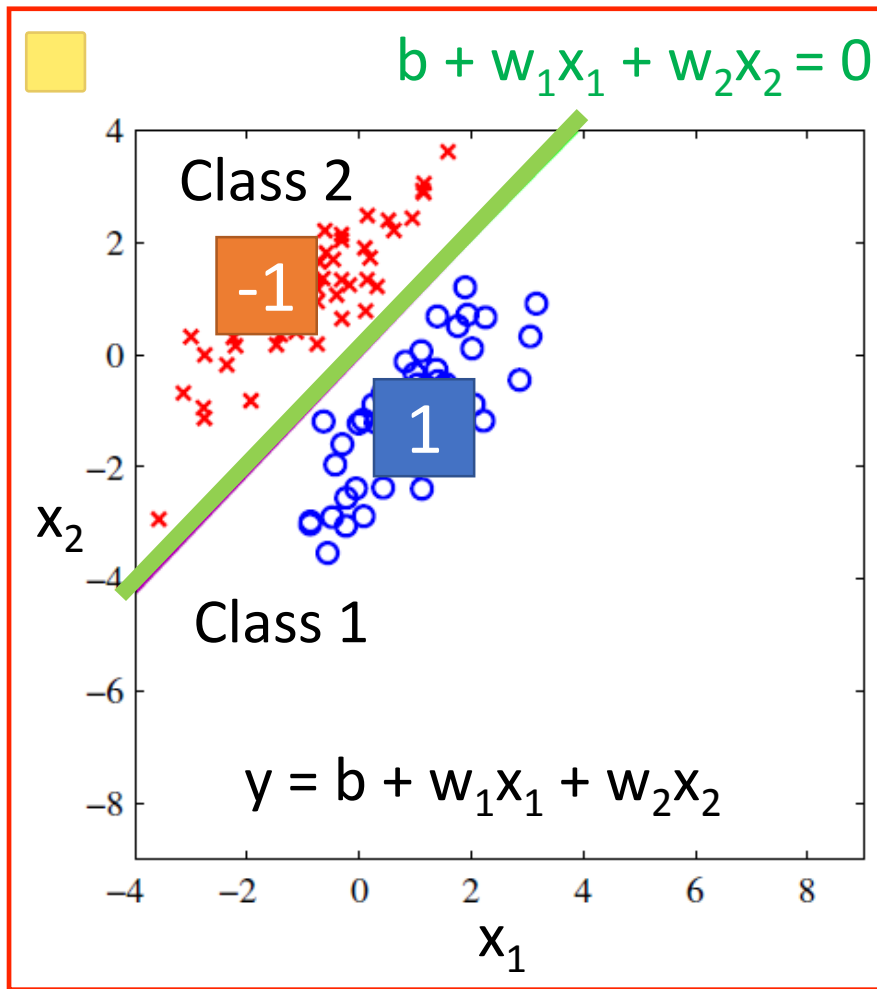


## Classification as Regression? ■

Binary classification as example

Training: Class 1 means the target is 1; Class 2 means the target is -1

Testing: closer to 1  $\rightarrow$  class 1; closer to -1  $\rightarrow$  class 2



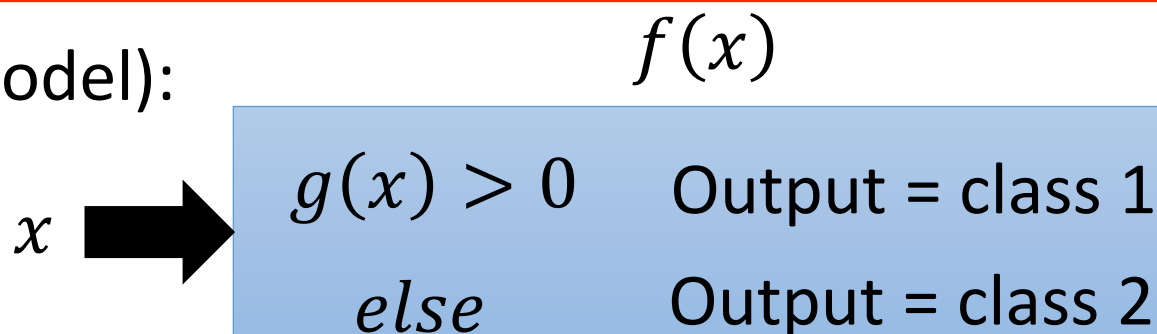
Penalize to the examples that are “too correct” ...

(Bishop, P186)

Multiple class: Class 1 means the target is 1; Class 2 means the target is 2; Class 3 means the target is 3 ..... problematic

# Ideal Alternatives

- Function (Model):



- Loss function:

$$L(f) = \sum_n \delta(f(x^n) \neq \hat{y}^n)$$

The number of times  $f$  get incorrect results on training data.

- Find the best function:

- Example: Perceptron, SVM

Not Today

若  $A, B$  為樣本空間  $\Omega$  中二事件, 且  $P(B) > 0$ 。則在給定  $B$  發生之下,  $A$  之條件機

率, 以  $P(A|B)$  表之, 定義為

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \circ \text{-----(1)}$$

在上述條件機率的定義中,  $B$  成為新的樣本空間:  $P(B|B) = 1$ 。也就是原先的樣本空間  $\Omega$

修正為  $B$ 。所有事件發生之機率, 都要先將其針對與  $B$  的關係做修正。例如, 若  $A$  與  $B$  為互斥事件, 且  $P(B) > 0$ , 則因  $P(A \cap B) = 0$ , 故  $P(A|B) = 0$ ; 若  $P(A)$  亦為正, 則此時亦有

$$P(B|A) = 0 \circ$$

條件機率也可用來求非條件下的機率。由(1)式得

$$P(A \cap B) = P(A|B)P(B) \circ \text{-----(2)}$$

故若知道  $P(A|B)$  及  $P(B)$ 。則可得到  $P(A \cap B)$ 。當然亦有

$$P(A \cap B) = P(B|A)P(A), \text{-----(3)}$$

只要  $P(A) > 0$ 。結合(2)式與(3)式, 得

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \circ \text{-----(4)}$$



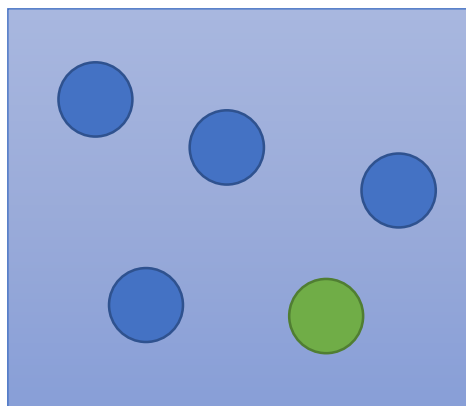
例如,  $A_1$  ,  $A_2$  為樣本空間  $\Omega$  中之一分割, 在給定一事件  $B$ , 且  $P(B) > 0$ , 則

$$P(A_1|B) = \frac{P(B|A_1)P(A_1)}{P(B|A_1)P(A_1) + P(B|A_2)P(A_2)} .$$

# Two Boxes

Box 1

$$P(B_1) = 2/3$$

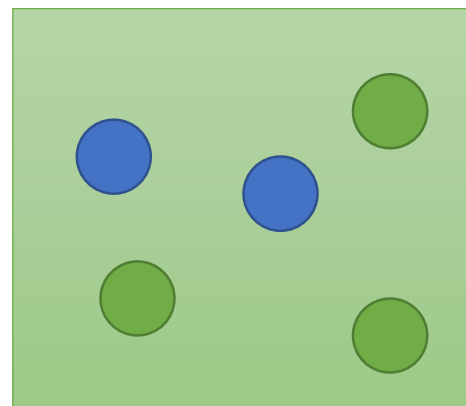


$$P(\text{Blue} | B_1) = 4/5$$

$$P(\text{Green} | B_1) = 1/5$$

Box 2

$$P(B_2) = 1/3$$



$$P(\text{Blue} | B_2) = 2/5$$

$$P(\text{Green} | B_2) = 3/5$$

 from one of the boxes

Where does it come from?

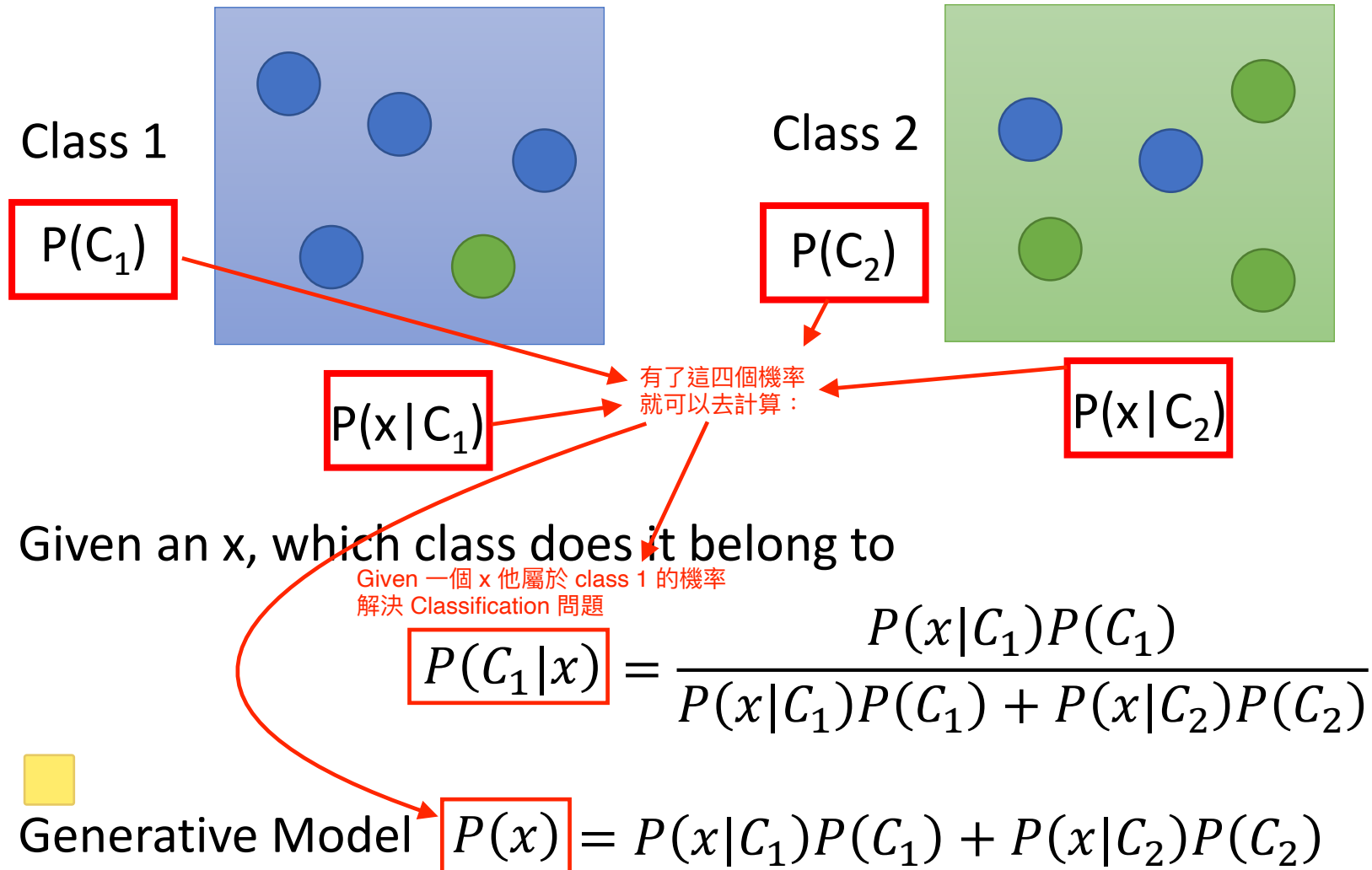
由條件機率公式推得而來

將 Box 換成 Class

$$P(B_1 | \text{Blue}) = \frac{P(\text{Blue} | B_1)P(B_1)}{P(\text{Blue} | B_1)P(B_1) + P(\text{Blue} | B_2)P(B_2)}$$

# Two Classes

Estimating the Probabilities  
From training data

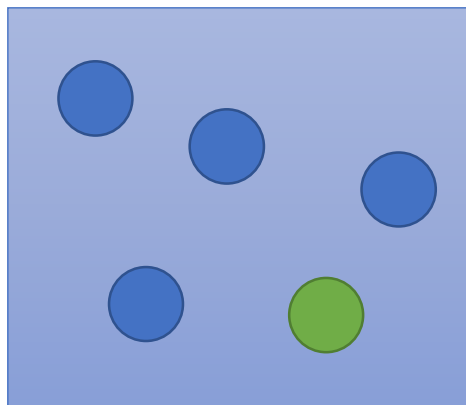


Class 的「機率」稱為 Prior (容易計算)

Prior

Class 1

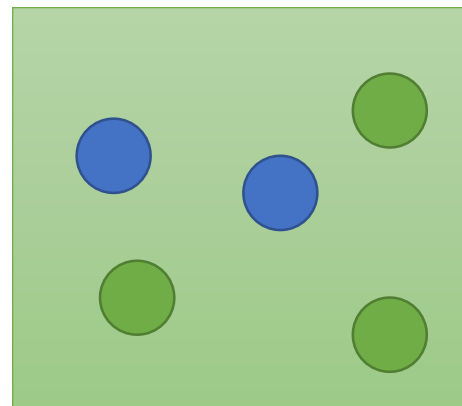
$P(C_1)$



Water

Class 2

$P(C_2)$



Normal

Water and Normal type with ID < 400 for training,  
rest for testing

Training: 79 Water, 61 Normal

例子：Training Data 有 79 隻水系寶可夢與 61 隻一般系寶可夢

$$P(C_1) = 79 / (79 + 61) = 0.56$$

$$P(C_2) = 61 / (79 + 61) = 0.44$$

# Probability from Class

Given 一個 Class 要 Sample 出一隻寶可夢的「機率」稱為 Likelihood (難計算)

$$P(x | C_1) = ? \quad P(\text{ } | \text{Water}) = ?$$



Each Pokémon is represented as  
a vector by its attribute.

➡ feature

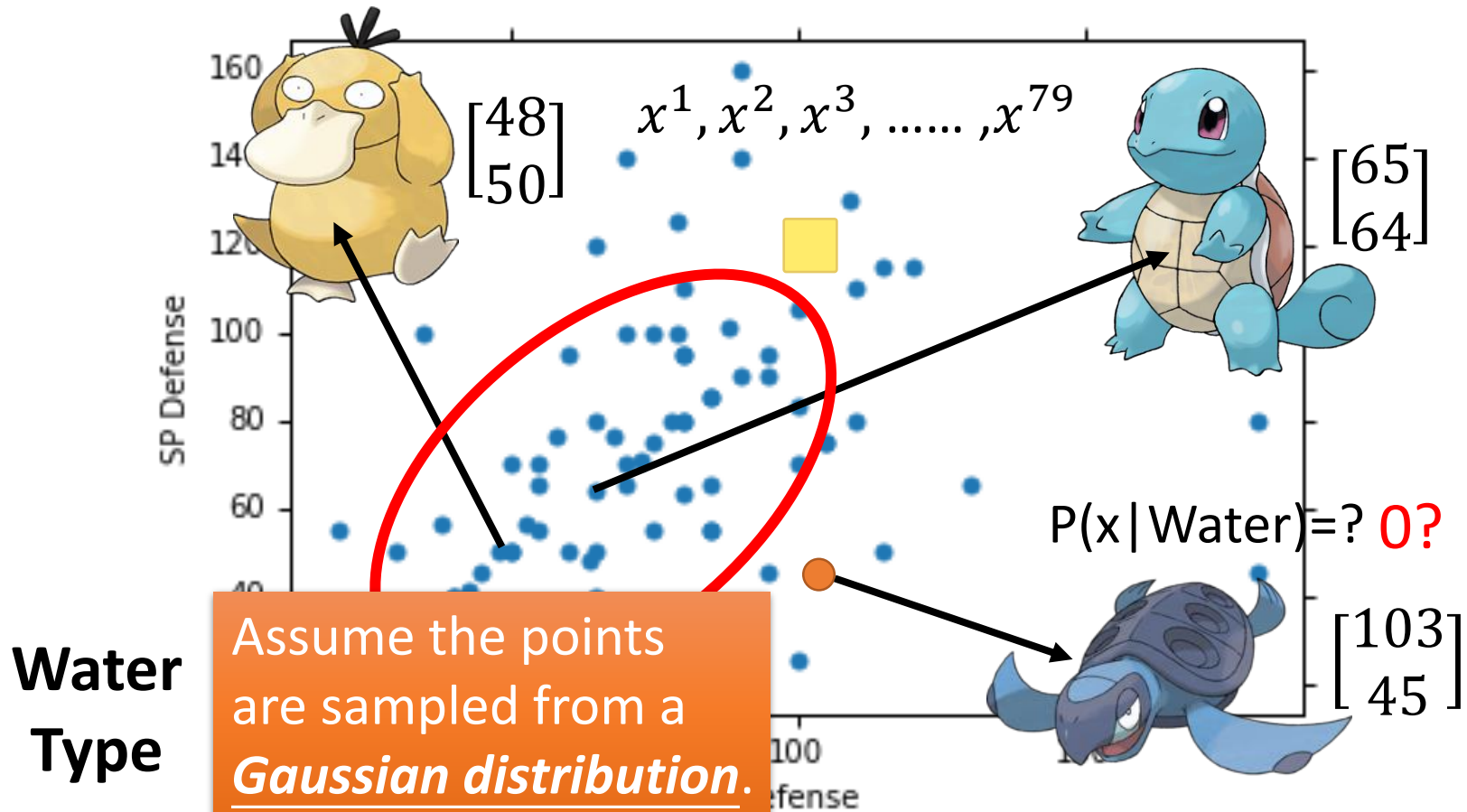
**Water  
Type**



79 in total

# Probability from Class - Feature

- Considering **Defense** and **SP Defense**



# Gaussian Distribution

<https://blog.slinuxer.com/tag/pca>

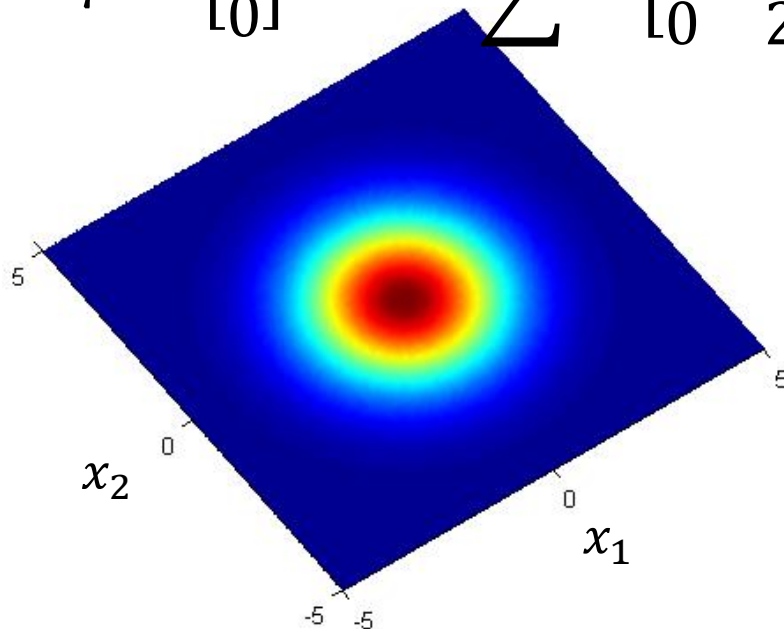
$$f_{\mu, \Sigma}(x) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

Input: vector  $x$ , output: probability of sampling  $x$

The shape of the function determines by **mean  $\mu$**  and **covariance matrix  $\Sigma$**

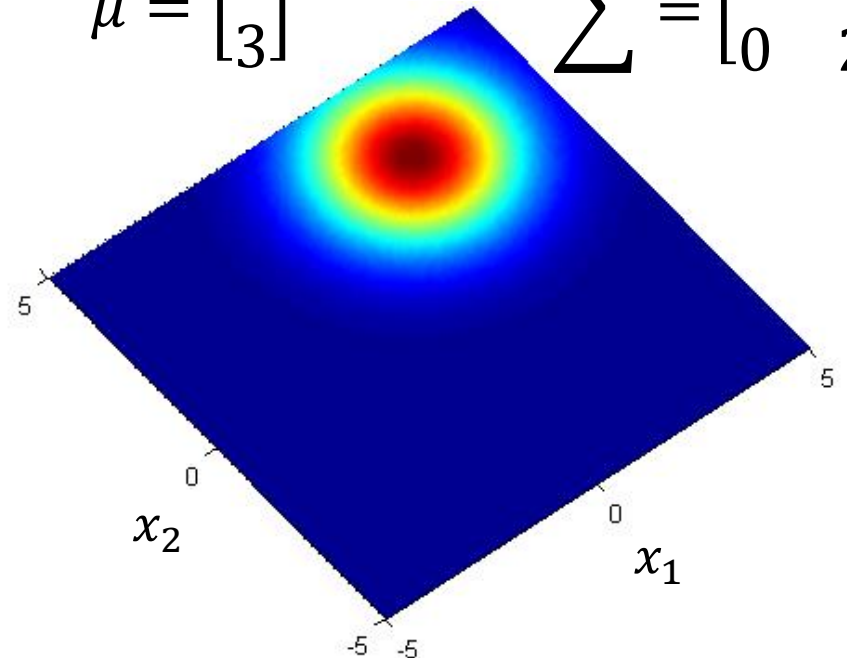
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$



$$\mu = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

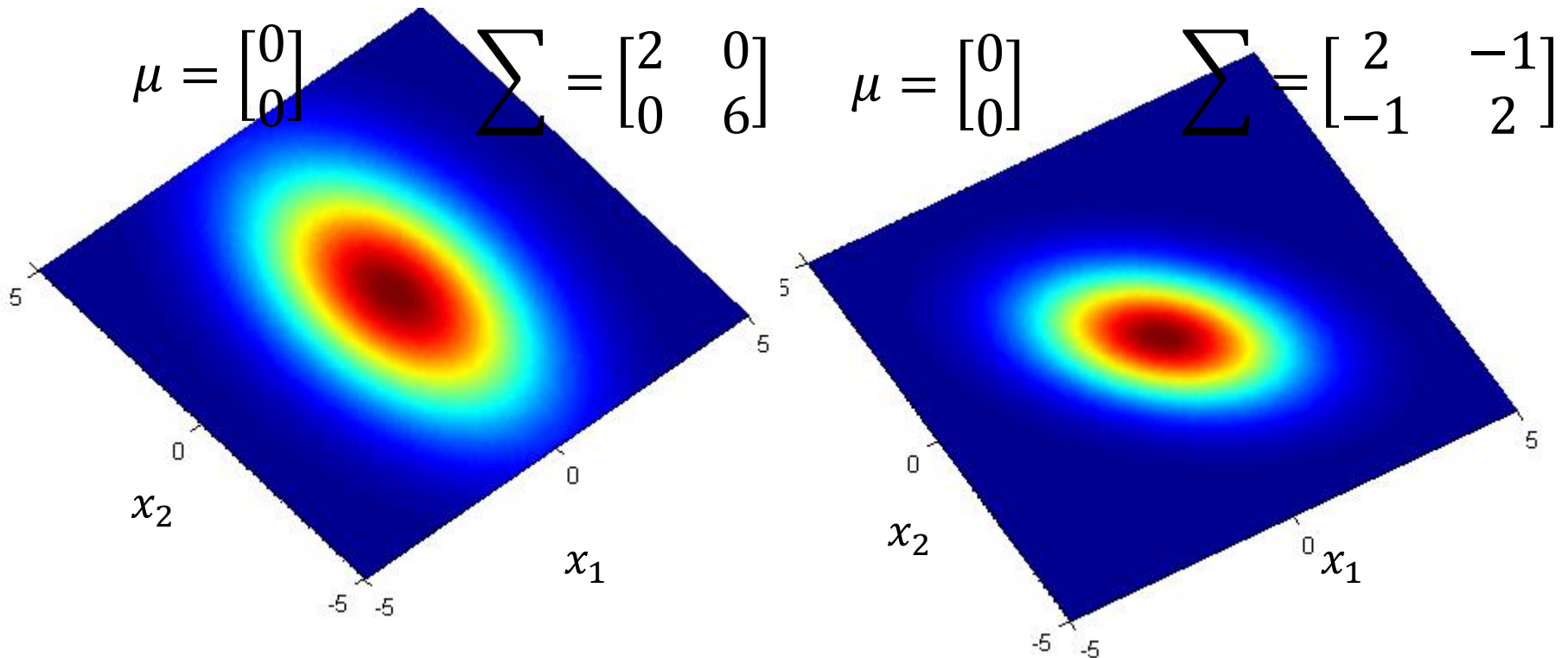


# Gaussian Distribution

$$f_{\mu,\Sigma}(x) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

Input: vector  $x$ , output: probability of sampling  $x$

The shape of the function determines by **mean  $\mu$**  and **covariance matrix  $\Sigma$**

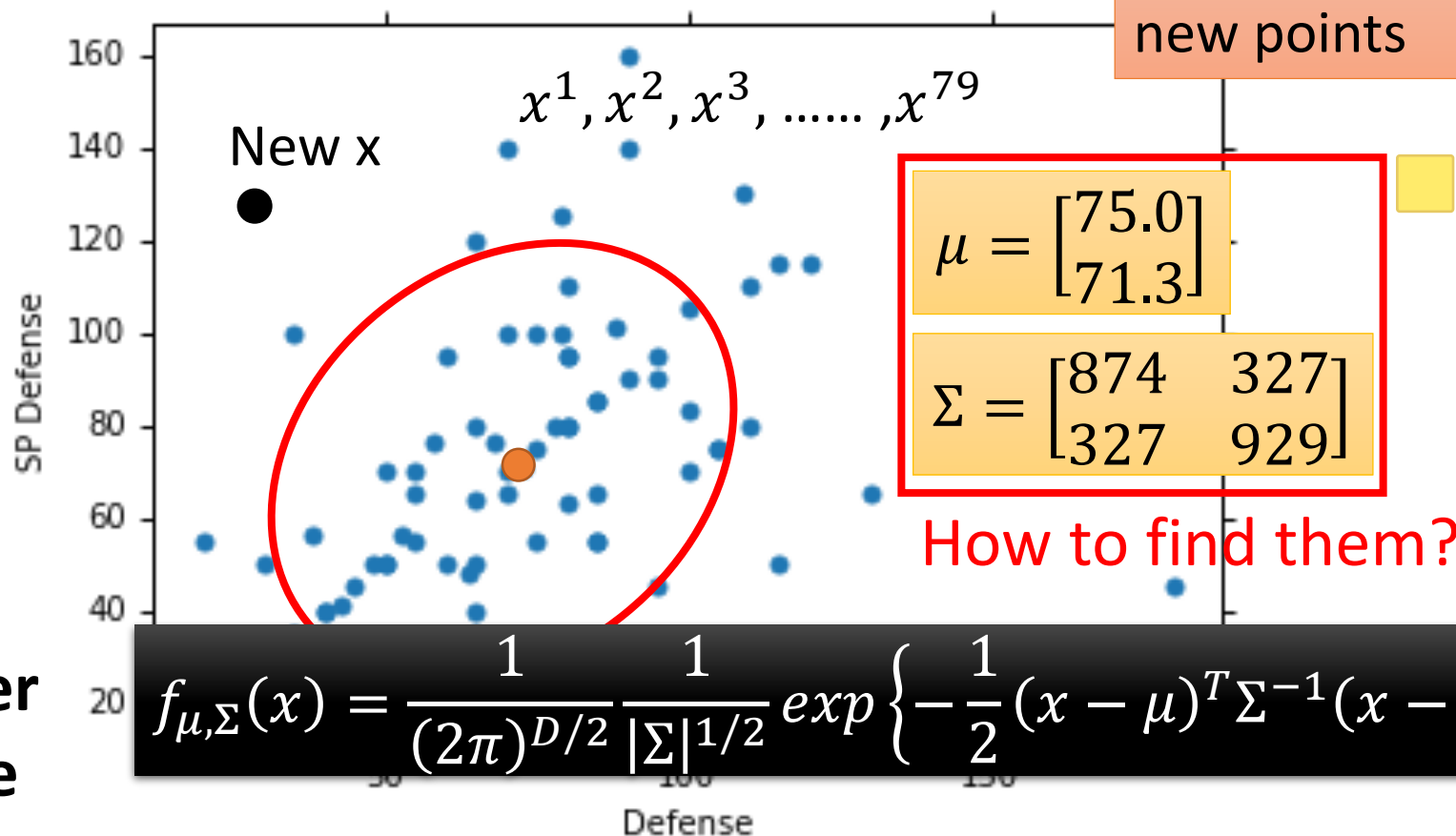




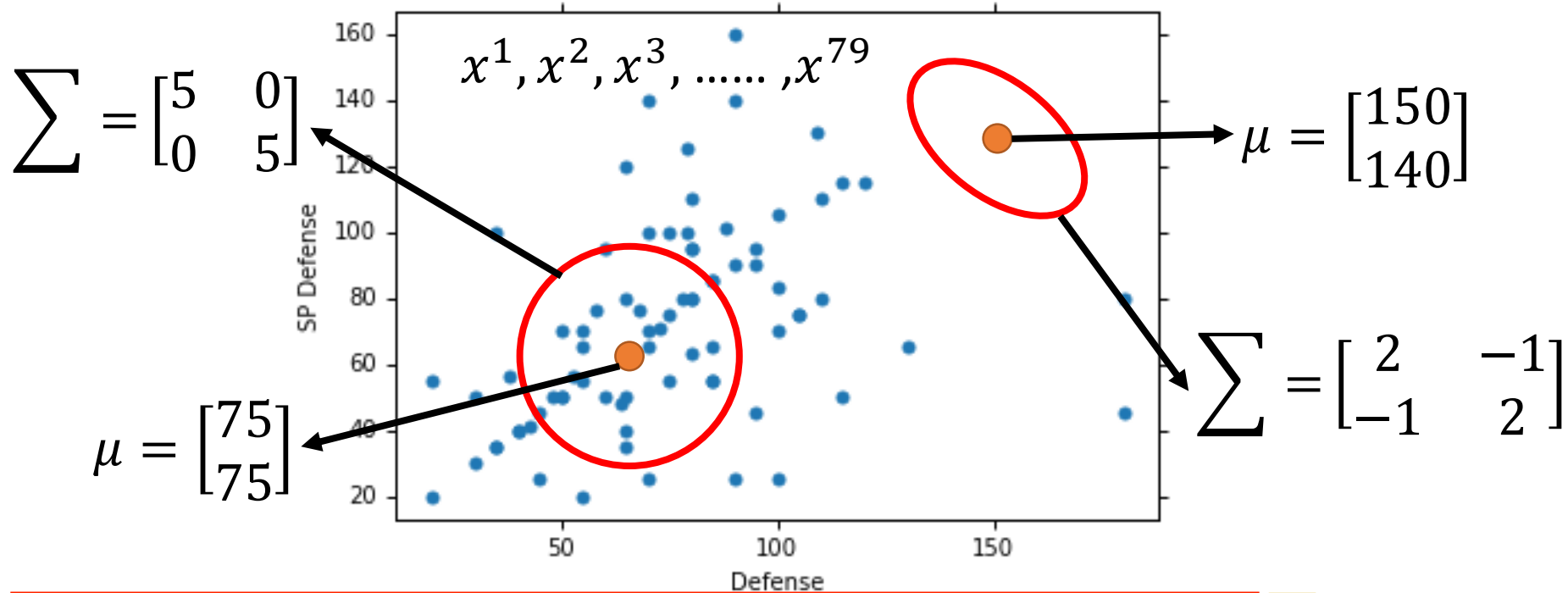
# Probability from Class

Assume the points are sampled from a Gaussian distribution

Find the Gaussian distribution behind them → Probability for new points



**Maximum Likelihood**  $f_{\mu, \Sigma}(x) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$



The Gaussian with any mean  $\mu$  and covariance matrix  $\Sigma$  can generate these points. ➡

Different Likelihood

Likelihood of a Gaussian with mean  $\mu$  and covariance matrix  $\Sigma$   
= the probability of the Gaussian samples  $x^1, x^2, x^3, \dots, x^{79}$

$$L(\mu, \Sigma) = f_{\mu, \Sigma}(x^1) f_{\mu, \Sigma}(x^2) f_{\mu, \Sigma}(x^3) \dots f_{\mu, \Sigma}(x^{79})$$

# Maximum Likelihood

We have the “Water” type Pokémons:  $x^1, x^2, x^3, \dots, x^{79}$

We assume  $x^1, x^2, x^3, \dots, x^{79}$  generate from the Gaussian  $(\mu^*, \Sigma^*)$  with the **maximum likelihood**

$$L(\mu, \Sigma) = f_{\mu, \Sigma}(x^1) f_{\mu, \Sigma}(x^2) f_{\mu, \Sigma}(x^3) \dots f_{\mu, \Sigma}(x^{79})$$

$$f_{\mu, \Sigma}(x) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

$$\mu^*, \Sigma^* = \arg \max_{\mu, \Sigma} L(\mu, \Sigma)$$



窮舉所有的 Mean 與 Covariance Matrix (Sigma) 使 Likelihood 最大

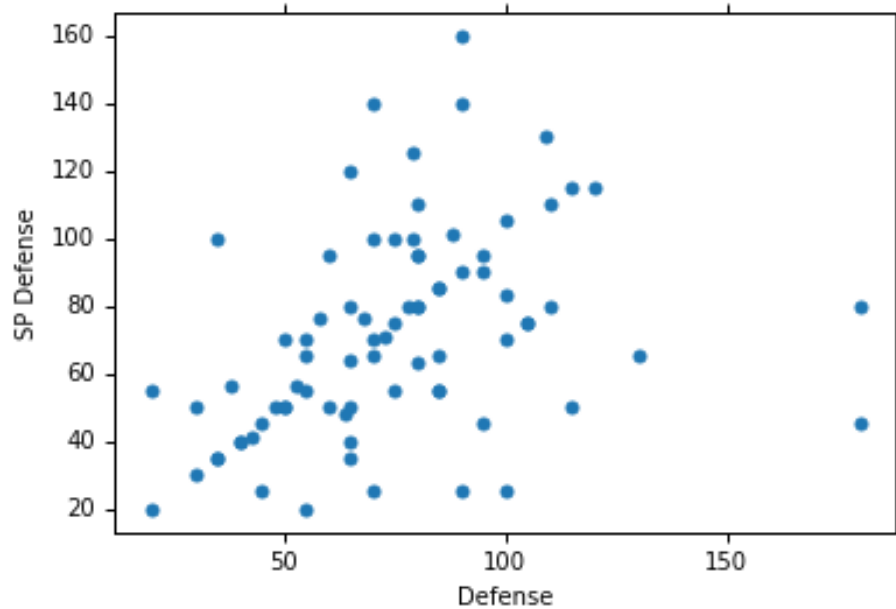
$$\mu^* = \frac{1}{79} \sum_{n=1}^{79} x^n$$

average

$$\Sigma^* = \frac{1}{79} \sum_{n=1}^{79} (x^n - \mu^*) (x^n - \mu^*)^T$$

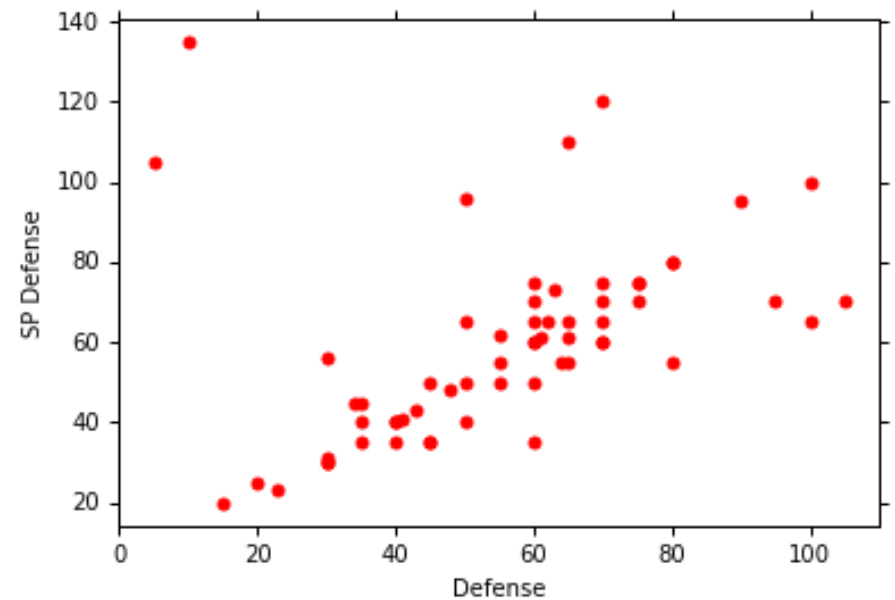
# Maximum Likelihood

Class 1: Water



$$\mu^1 = \begin{bmatrix} 75.0 \\ 71.3 \end{bmatrix} \quad \Sigma^1 = \begin{bmatrix} 874 & 327 \\ 327 & 929 \end{bmatrix}$$

Class 2: Normal



$$\mu^2 = \begin{bmatrix} 55.6 \\ 59.8 \end{bmatrix} \quad \Sigma^2 = \begin{bmatrix} 847 & 422 \\ 422 & 685 \end{bmatrix}$$

# Now we can do classification 😊

$$f_{\mu^1, \Sigma^1}(x) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma^1|^{1/2}} \exp\left\{-\frac{1}{2}(x - \mu^1)^T (\Sigma^1)^{-1} (x - \mu^1)\right\}$$

$P(C_1) = 79 / (79 + 61) = 0.56$

$$\mu^1 = \begin{bmatrix} 75.0 \\ 71.3 \end{bmatrix} \quad \Sigma^1 = \begin{bmatrix} 874 & 327 \\ 327 & 929 \end{bmatrix}$$

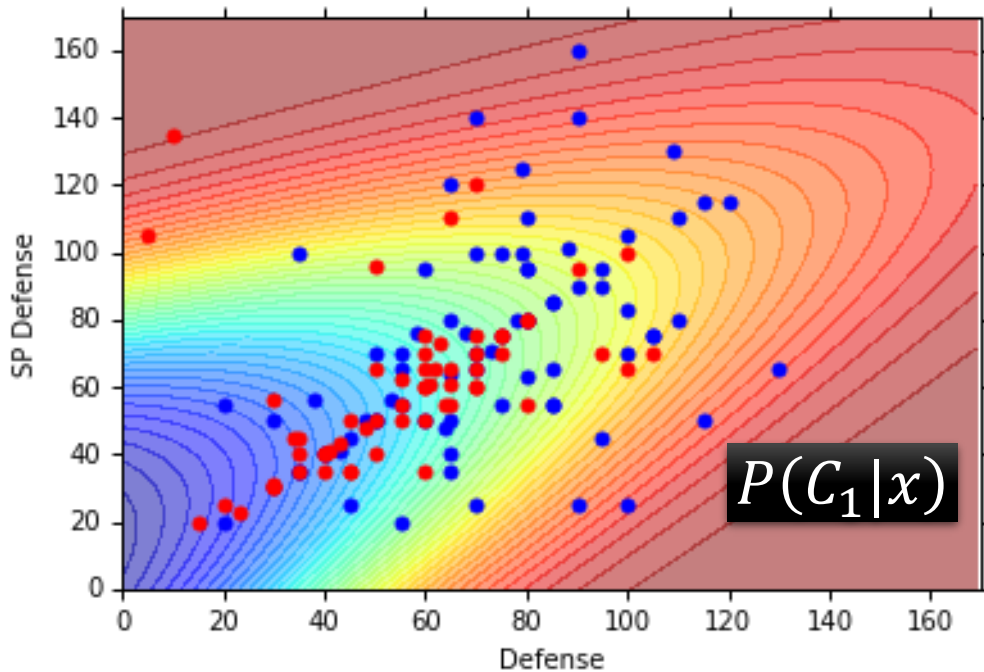
$$\boxed{P(C_1|x)}^{\text{posterior}} = \frac{\boxed{P(x|C_1)}^{\text{likelihood}} \boxed{P(C_1)}^{\text{prior}}}{P(x|C_1)P(C_1) + \boxed{P(x|C_2)}^{\text{likelihood}} \boxed{P(C_2)}^{\text{prior}}}$$

$$f_{\mu^2, \Sigma^2}(x) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma^2|^{1/2}} \exp\left\{-\frac{1}{2}(x - \mu^2)^T (\Sigma^2)^{-1} (x - \mu^2)\right\}$$

$P(C_2) = 61 / (79 + 61) = 0.44$

$$\mu^2 = \begin{bmatrix} 55.6 \\ 59.8 \end{bmatrix} \quad \Sigma^2 = \begin{bmatrix} 847 & 422 \\ 422 & 685 \end{bmatrix}$$

If  $P(C_1|x) > 0.5$  ➡ x belongs to class 1 (Water)



Blue points:  $C_1$  (Water), Red points:  $C_2$  (Normal)

How's the results?

Model 的表现不好

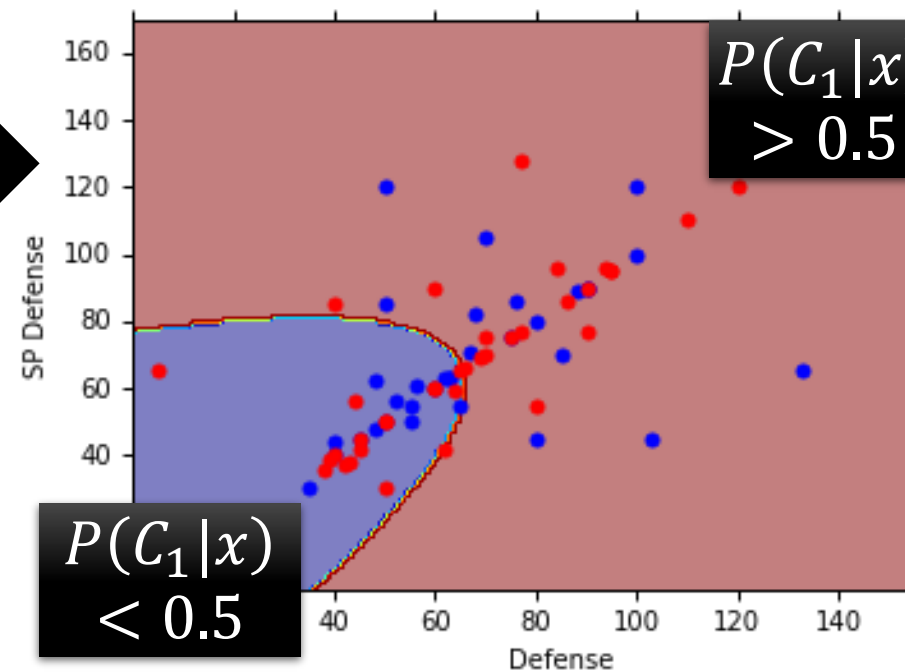
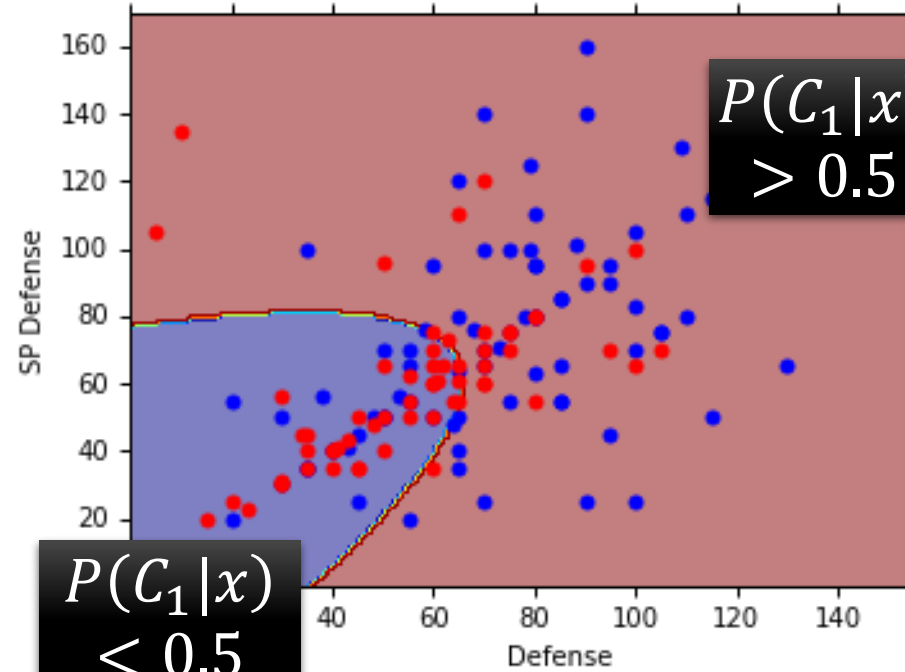
Testing data: 47% accuracy

All: total, hp, att, sp att,  
de, sp de, speed (7 features)

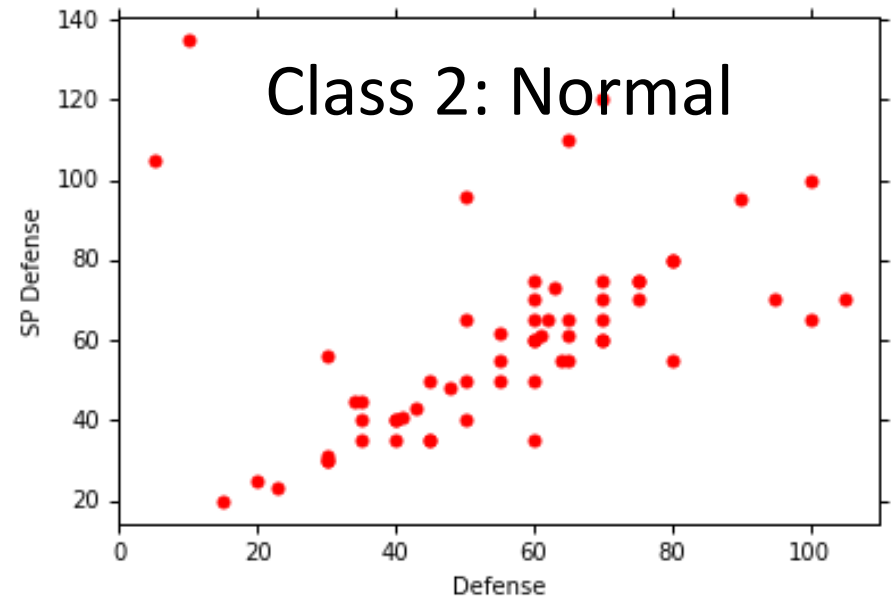
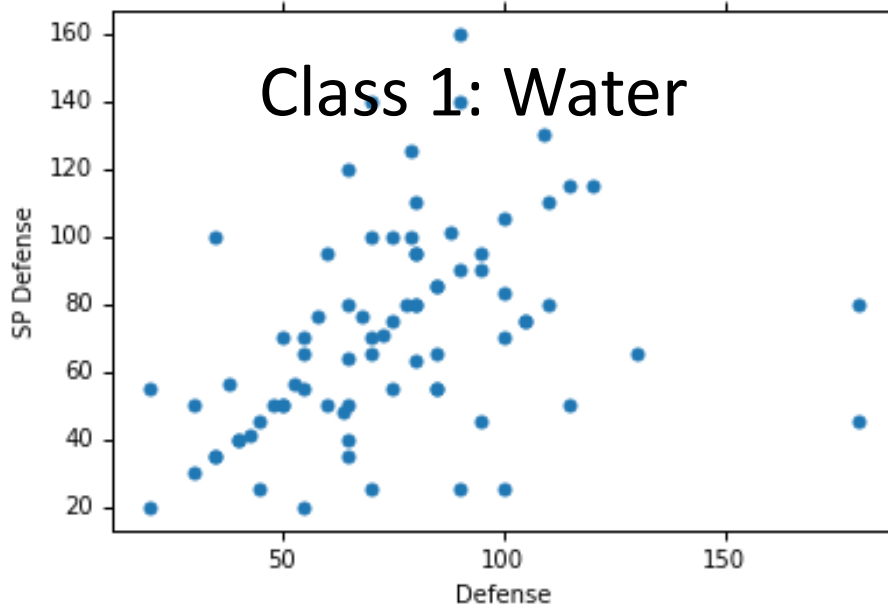
$\mu^1, \mu^2$ : 7-dim vector

$\Sigma^1, \Sigma^2$ : 7 x 7 matrices

54% accuracy ... ☹️



# Modifying Model



$$\mu^1 = \begin{bmatrix} 75.0 \\ 71.3 \end{bmatrix} \quad \Sigma^1 = \begin{bmatrix} 874 & 327 \\ 327 & 929 \end{bmatrix}$$

$$\mu^2 = \begin{bmatrix} 55.6 \\ 59.8 \end{bmatrix} \quad \Sigma^2 = \begin{bmatrix} 847 & 422 \\ 422 & 685 \end{bmatrix}$$

將兩個 Class 的 Gaussian Distribution 的 Covariance Matrix 設為相同的 => 減少 Model 的參數量 => 避免 Overfitting

The same  $\Sigma$

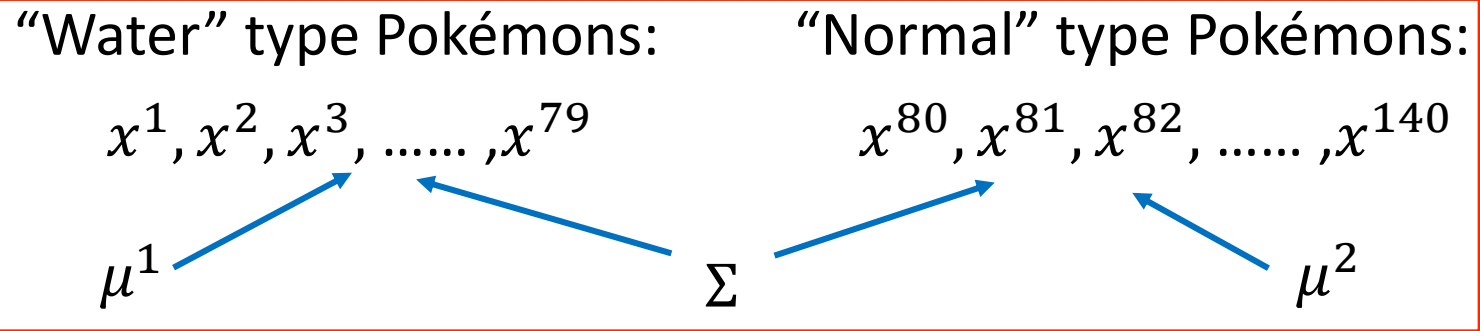
Less parameters

# Modifying Model

Ref: Bishop,  
chapter 4.2.2

- Maximum likelihood

兩個 Gaussian Distribution 使用各自的 Mean 但是使用共同的 Covariance Matrix



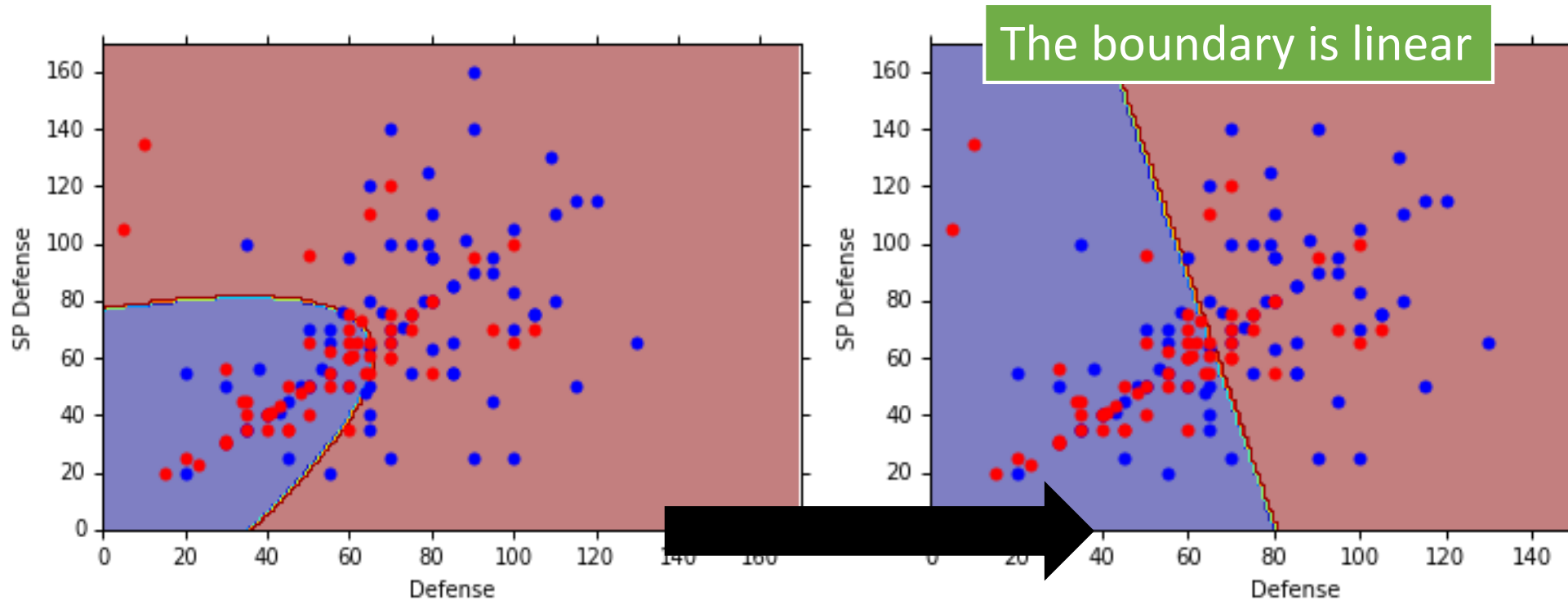
Find  $\mu^1, \mu^2, \Sigma$  maximizing the likelihood  $L(\mu^1, \mu^2, \Sigma)$

$$L(\mu^1, \mu^2, \Sigma) = f_{\mu^1, \Sigma}(x^1) f_{\mu^1, \Sigma}(x^2) \cdots f_{\mu^1, \Sigma}(x^{79}) \\ \times f_{\mu^2, \Sigma}(x^{80}) f_{\mu^2, \Sigma}(x^{81}) \cdots f_{\mu^2, \Sigma}(x^{140})$$

$$\mu^1 \text{ and } \mu^2 \text{ is the same} \quad \Sigma = \frac{79}{140} \Sigma^1 + \frac{61}{140} \Sigma^2$$



# Modifying Model



The same covariance matrix

All: total, hp, att, sp att, de, sp de, speed

54% accuracy



73% accuracy

# Three Steps

- Function Set (Model): Prior 與 Likelihood 就「參數」

$x$  

$$P(C_1|x) = \frac{P(x|C_1)P(C_1)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)}$$

If  $P(C_1|x) > 0.5$ , output: class 1  
Otherwise, output: class 2

- Goodness of a function:
  - The mean  $\mu$  and covariance  $\Sigma$  that maximizing the likelihood (the probability of generating data)
- Find the best function: easy

# Probability Distribution

- You can always use the distribution you like 😊



$$P(x|C_1) = P(x_1|C_1) P(x_2|C_1) \cdots P(x_k|C_1) \cdots$$

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_k \\ \vdots \\ x_K \end{bmatrix}$$

1-D Gaussian

For binary features, you may assume they are from Bernoulli distributions.

If you assume all the dimensions are independent, then you are using *Naive Bayes Classifier*.

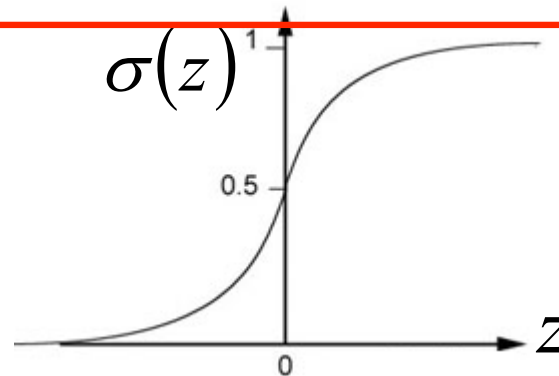
# Posterior Probability

關鍵：Posterior Probability 可以轉為 Sigmoid Function (z)

$$\begin{aligned} P(C_1|x) &= \frac{P(x|C_1)P(C_1)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)} \\ &= \frac{1}{1 + \frac{P(x|C_2)P(C_2)}{P(x|C_1)P(C_1)}} = \frac{1}{1 + \exp(-z)} = \sigma(z) \end{aligned}$$

Sigmoid function

$$z = \ln \frac{P(x|C_1)P(C_1)}{P(x|C_2)P(C_2)}$$



Warning of Math

# Posterior Probability

$$P(C_1|x) = \sigma(z) \quad \text{sigmoid} \quad z = \ln \frac{P(x|C_1)P(C_1)}{P(x|C_2)P(C_2)}$$

$$z = \ln \frac{P(x|C_1)}{P(x|C_2)} + \ln \frac{P(C_1)}{P(C_2)} \rightarrow \frac{\frac{N_1}{N_1 + N_2}}{\frac{N_2}{N_1 + N_2}} = \frac{N_1}{N_2}$$

$$P(x|C_1) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma^1|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu^1)^T (\Sigma^1)^{-1} (x - \mu^1) \right\}$$

$$P(x|C_2) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma^2|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu^2)^T (\Sigma^2)^{-1} (x - \mu^2) \right\}$$

$$z = \ln \frac{P(x|C_1)}{P(x|C_2)} + \ln \frac{P(C_1)}{P(C_2)} = \frac{N_1}{N_2}$$

$$P(x|C_1) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma^1|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu^1)^T (\Sigma^1)^{-1} (x - \mu^1) \right\}$$

$$P(x|C_2) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma^2|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu^2)^T (\Sigma^2)^{-1} (x - \mu^2) \right\}$$

$$\ln \frac{\cancel{\frac{1}{(2\pi)^{D/2}}} \frac{1}{|\Sigma^1|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu^1)^T (\Sigma^1)^{-1} (x - \mu^1) \right\}}{\cancel{\frac{1}{(2\pi)^{D/2}}} \frac{1}{|\Sigma^2|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu^2)^T (\Sigma^2)^{-1} (x - \mu^2) \right\}}$$

$$= \ln \frac{|\Sigma^2|^{1/2}}{|\Sigma^1|^{1/2}} \exp \left\{ -\frac{1}{2} [(x - \mu^1)^T (\Sigma^1)^{-1} (x - \mu^1) - (x - \mu^2)^T (\Sigma^2)^{-1} (x - \mu^2)] \right\}$$

$$= \ln \frac{|\Sigma^2|^{1/2}}{|\Sigma^1|^{1/2}} - \frac{1}{2} [(x - \mu^1)^T (\Sigma^1)^{-1} (x - \mu^1) - (x - \mu^2)^T (\Sigma^2)^{-1} (x - \mu^2)]$$

$$z = \ln \frac{P(x|C_1)}{P(x|C_2)} + \ln \frac{P(C_1)}{P(C_2)} = \frac{N_1}{N_2}$$

$$= \ln \frac{|\Sigma^2|^{1/2}}{|\Sigma^1|^{1/2}} - \frac{1}{2} [ \underbrace{(x - \mu^1)^T (\Sigma^1)^{-1} (x - \mu^1)}_{\text{red}} - \underbrace{(x - \mu^2)^T (\Sigma^2)^{-1} (x - \mu^2)}_{\text{red}} ]$$

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

$$= x^T (\Sigma^1)^{-1} x - \underbrace{x^T (\Sigma^1)^{-1} \mu^1 - (\mu^1)^T (\Sigma^1)^{-1} x}_{\text{blue}} + (\mu^1)^T (\Sigma^1)^{-1} \mu^1$$

$$= x^T (\Sigma^1)^{-1} x - \underbrace{2(\mu^1)^T (\Sigma^1)^{-1} x}_{\text{blue}} + (\mu^1)^T (\Sigma^1)^{-1} \mu^1$$

$$(x - \mu^2)^T (\Sigma^2)^{-1} (x - \mu^2)$$

$$= x^T (\Sigma^2)^{-1} x - 2(\mu^2)^T (\Sigma^2)^{-1} x + (\mu^2)^T (\Sigma^2)^{-1} \mu^2$$

$$z = \ln \frac{|\Sigma^2|^{1/2}}{|\Sigma^1|^{1/2}} - \frac{1}{2} x^T (\Sigma^1)^{-1} x + (\mu^1)^T (\Sigma^1)^{-1} x - \frac{1}{2} (\mu^1)^T (\Sigma^1)^{-1} \mu^1 \\ + \frac{1}{2} x^T (\Sigma^2)^{-1} x - (\mu^2)^T (\Sigma^2)^{-1} x + \frac{1}{2} (\mu^2)^T (\Sigma^2)^{-1} \mu^2 + \ln \frac{N_1}{N_2}$$



End of Warning

$$P(C_1|x) = \sigma(z)$$

$$z = \cancel{\ln \frac{|\Sigma^2|^{1/2}}{|\Sigma^1|^{1/2}}} - \cancel{\frac{1}{2} x^T (\Sigma^1)^{-1} x} + (\mu^1)^T (\Sigma^1)^{-1} x - \frac{1}{2} (\mu^1)^T (\Sigma^1)^{-1} \mu^1 \\ + \cancel{\frac{1}{2} x^T (\Sigma^2)^{-1} x} - (\mu^2)^T (\Sigma^2)^{-1} x + \frac{1}{2} (\mu^2)^T (\Sigma^2)^{-1} \mu^2 + \ln \frac{N_1}{N_2}$$

$$\Sigma_1 = \Sigma_2 = \Sigma$$

$$z = \underbrace{(\mu^1 - \mu^2)^T \Sigma^{-1} x}_{\mathbf{w}^T} - \underbrace{\frac{1}{2} (\mu^1)^T \Sigma^{-1} \mu^1 + \frac{1}{2} (\mu^2)^T \Sigma^{-1} \mu^2}_{b} + \ln \frac{N_1}{N_2}$$

$$P(C_1|x) = \sigma(\mathbf{w} \cdot x + b)$$

How about directly find  $\mathbf{w}$  and  $b$ ?

In generative model, we estimate  $N_1, N_2, \mu^1, \mu^2, \Sigma$

Then we have  $\mathbf{w}$  and  $b$