



## Classification

$X \rightarrow \text{function} \rightarrow \text{class n}$

- Ex.: ① Credit Scoring  
 ② Medical Diagnosis  
 ③ Handwritten character recognition  
 ④ face recognition

Example: class. 電子商務類別.

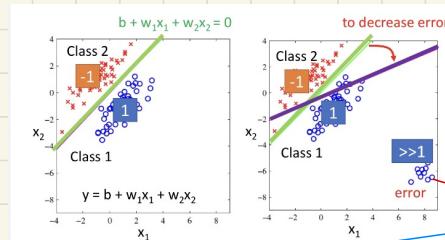
How to do classification?

training data for classification

considering Binary classification as example  
 二分类

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



若用 classification, 则分界线是绿色的;

若用 regression, 则是紫色的, 因为 regression 会惩罚这些点

所以需要改变!!! 不能用 regression

Ideal Alternatives (理想方法)

• function (model):

$$x \Rightarrow \begin{cases} f(x) > 0 & \text{output = class 1} \\ \text{else} & \text{output = class 2} \end{cases}$$

• Loss function

$$L(f) = \frac{1}{n} \sum_{i=1}^n \delta(f(x_i) \neq y_i)$$

无法区分

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

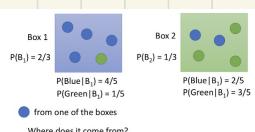
• find the best function

example: perceptron sum

感知机

怎么寻找 the best function

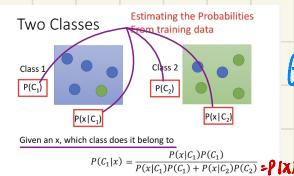
Two Boxes



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

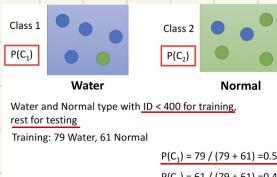
全概率率公式

## Two Classes



Generative model  $P(x) = P(x|C_1)P(C_1) + P(x|C_2)P(C_2)$

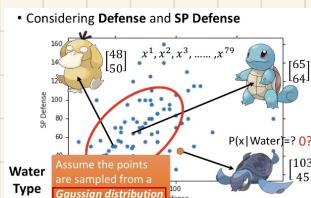
Then we need to compute  
Prior



$P(x|C_i) ??$

each Pokemon is represented as a vector by its attribute.  $\Rightarrow$  feature

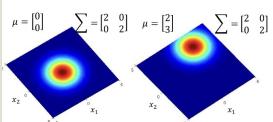
↓ probability from class - feature



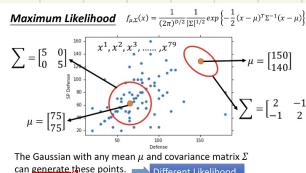
Gaussian Distribution

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

Input: vector  $x$ , output: probability of sampling  $x$ . The shape of the function determines by mean  $\mu$  and covariance matrix  $\Sigma$ .



根据 79 个点求得  $\mu, \Sigma$  则  $f_{\mu, \Sigma}(x)$  可写出来  
怎么计算  $\mu, \Sigma$  = Maximum Likelihood



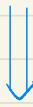
Likelihood of a Gaussian with mean  $\mu$  and covariance matrix  $\Sigma$

= the probability of the Gaussian samples  $x^1, x^2, x^3, \dots, x^{79}$

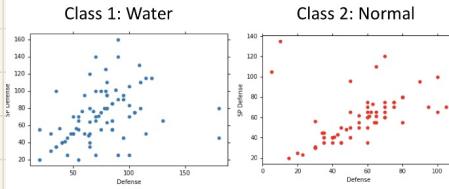
$$L(\mu, \Sigma) = f_{\mu, \Sigma}(x^1) f_{\mu, \Sigma}(x^2) \dots f_{\mu, \Sigma}(x^{79})$$

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

$$\begin{aligned}
 L(\mu, \Sigma) &= f_{\mu, \Sigma}(x^1) f_{\mu, \Sigma}(x^2) \cdots f_{\mu, \Sigma}(x^n) \\
 f_{\mu, \Sigma}(x) &= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\} \\
 \mu^*, \Sigma^* &= \arg \max_{\mu, \Sigma} L(\mu, \Sigma) \\
 \mu^* &= \frac{1}{n} \bar{x} \\
 \Sigma^* &= \frac{1}{n} \sum (x^i - \mu^*) (x^i - \mu^*)^T
 \end{aligned}$$



计算结果



$$\mu^1 = \begin{bmatrix} 75.0 \\ 71.3 \end{bmatrix}, \Sigma^1 = \begin{bmatrix} 874 & 327 \\ 327 & 929 \end{bmatrix}, \mu^2 = \begin{bmatrix} 55.6 \\ 59.8 \end{bmatrix}, \Sigma^2 = \begin{bmatrix} 847 & 422 \\ 422 & 685 \end{bmatrix}$$



Now we can do classification

$$\begin{aligned}
 f_{\mu^1, \Sigma^1}(x) &= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma^1|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu^1)^T (\Sigma^1)^{-1} (x - \mu^1) \right\} & P(C_1) \\
 \mu^1 &= \begin{bmatrix} 75.0 \\ 71.3 \end{bmatrix}, \Sigma^1 = \begin{bmatrix} 874 & 327 \\ 327 & 929 \end{bmatrix} & = 79 / (79 + 61) = 0.56 \\
 P(C_1|x) &= \frac{P(x|C_1)P(C_1)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)} & \\
 f_{\mu^2, \Sigma^2}(x) &= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma^2|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu^2)^T (\Sigma^2)^{-1} (x - \mu^2) \right\} & P(C_2) \\
 \mu^2 &= \begin{bmatrix} 55.6 \\ 59.8 \end{bmatrix}, \Sigma^2 = \begin{bmatrix} 847 & 422 \\ 422 & 685 \end{bmatrix} & = 61 / (79 + 61) = 0.44
 \end{aligned}$$

If  $P(C_1|x) > 0.5 \Rightarrow x \text{ belongs to class 1 (water)}$

How's the results?

Testing data: 47% accuracy

All: total, hp, att, sp att,

de, sp de, speed (7 features)

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

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

实际上每个宝可梦都有7个特征，可能在7维空间里

以正常识别，但结果仍然是54%的正确率

54% accuracy ...



Modifying Model

将不同类别的协方差设为相同的工。原因：由于一个协方差矩阵的参数是 $D$ 个，如果设置不同，就会造成参数过多。

• Maximum likelihood

"Water" type Pokémons:

$$x^1, x^2, x^3, \dots, x^{79}$$

$$\mu^1$$

"Normal" type Pokémons:

$$x^{80}, x^{81}, x^{82}, \dots, x^{140}$$

$$\mu^2$$

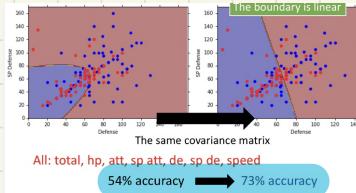
$$\Sigma$$

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

$$\mu^1, \mu^2 \text{ 与之前算法一致} \quad \Sigma = \frac{29}{140} \bar{\Sigma} + \frac{61}{140} \bar{\Sigma}^2$$



Result:



Conclusion:

① Three step

- Function Set (Model)

$$x \longrightarrow$$

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

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

question? Why choose this probability distribution

You can always use the distribution you like

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

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

1-D Gaussian

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

To be continued

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

朴素贝叶斯

# posterior probability

后验概率

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

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

Sigmoid function

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



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

$z = \ln$

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

$\frac{\frac{N_1}{N_1+N_2}}{\frac{N_2}{N_1+N_2}} = \frac{N_1}{N_2}$

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

$$\begin{aligned} P(x|C_1) &= \frac{1}{(2\pi)^{D/2} |\Sigma_1|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu^1)^T (\Sigma_1^{-1}) (x - \mu^1) \right\} \\ P(x|C_2) &= \frac{1}{(2\pi)^{D/2} |\Sigma_2|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu^2)^T (\Sigma_2^{-1}) (x - \mu^2) \right\} \\ \ln \frac{1}{(2\pi)^{D/2} |\Sigma_1|^{1/2}} &\exp \left\{ -\frac{1}{2} (x - \mu^1)^T (\Sigma_1^{-1}) (x - \mu^1) \right\} \\ \ln \frac{1}{(2\pi)^{D/2} |\Sigma_2|^{1/2}} &\exp \left\{ -\frac{1}{2} (x - \mu^2)^T (\Sigma_2^{-1}) (x - \mu^2) \right\} \\ &= \ln \frac{|\Sigma_1|^{1/2}}{|\Sigma_2|^{1/2}} \exp \left\{ -\frac{1}{2} [(x - \mu^1)^T (\Sigma_1^{-1}) (x - \mu^1) \right. \\ &\quad \left. - (x - \mu^2)^T (\Sigma_2^{-1}) (x - \mu^2)] \right\} \\ &= \ln \frac{|\Sigma_1|^{1/2}}{|\Sigma_2|^{1/2}} - \frac{1}{2} [(x - \mu^1)^T (\Sigma_1^{-1}) (x - \mu^1) - (x - \mu^2)^T (\Sigma_2^{-1}) (x - \mu^2)] \end{aligned}$$

$(x - \mu^1)^T (\Sigma_1^{-1}) (x - \mu^1)$

$= x^T (\Sigma_1^{-1}) x - x^T (\Sigma_1^{-1}) \mu^1 - (\mu^1)^T (\Sigma_1^{-1}) x + (\mu^1)^T (\Sigma_1^{-1}) \mu^1$ 
 $= x^T (\Sigma_1^{-1}) x - 2(\mu^1)^T (\Sigma_1^{-1}) x + (\mu^1)^T (\Sigma_1^{-1}) \mu^1$

$(x - \mu^2)^T (\Sigma_2^{-1}) (x - \mu^2)$

$= x^T (\Sigma_2^{-1}) x - 2(\mu^2)^T (\Sigma_2^{-1}) x + (\mu^2)^T (\Sigma_2^{-1}) \mu^2$

$z = \ln \frac{|\Sigma_1|^{1/2}}{|\Sigma_2|^{1/2}} - \frac{1}{2} x^T (\Sigma_1^{-1}) x + (\mu^1)^T (\Sigma_1^{-1}) x - \frac{1}{2} (\mu^1)^T (\Sigma_1^{-1}) \mu^1$

$- \frac{1}{2} x^T (\Sigma_2^{-1}) x - (\mu^2)^T (\Sigma_2^{-1}) x + \frac{1}{2} (\mu^2)^T (\Sigma_2^{-1}) \mu^2 + \ln \frac{N_1}{N_2}$

$\bar{z}_1 = \bar{z}_2 = \bar{z}$

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

$P(C_1|x) = \delta(w \cdot x + b)$ 

为时 boundary 是 linear

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

How about directly find  $w$  and  $b$ ?

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

Then we have  $w$  and  $b$

生成模型

模型表示为给定输入  $x$

产生输出  $y$  生成关系 基于统计学与 Bayes 理论