

# 《神经网络与深度学习》

## 第二节：线性模型

主讲人：戴金成(副教授，博士生导师)

daijincheng@bupt.edu.cn

神经网络与深度学习课程组



北京郵電大學

Beijing University of Posts and Telecommunications

# 内容导览

---



线性回归：多项式拟合



线性分类：概率生成模型



Logistic回归：神经元模型



多分类问题



线性模型的局限与非线性模型

# 内容导览

---



线性回归：多项式拟合



线性分类：概率生成模型



Logistic回归：神经元模型



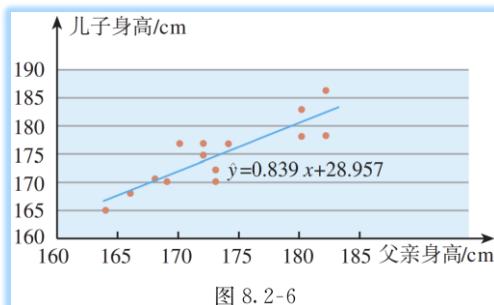
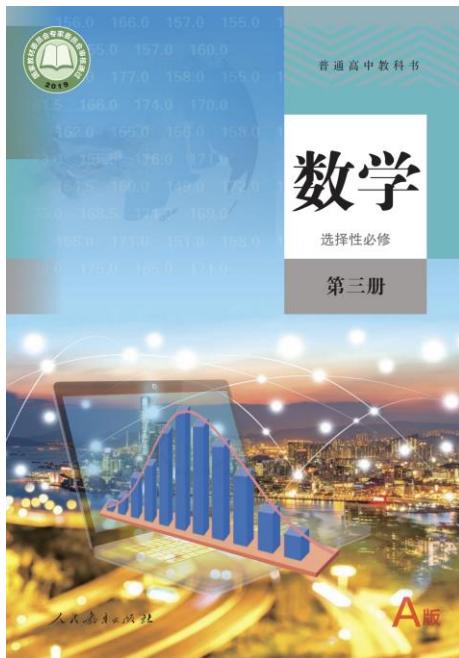
多分类问题



线性模型的局限与非线性模型

# 线性回归

## 高中时期的经典问题



## 8.2 一元线性回归模型及其应用

通过前面的学习我们已经了解到，根据成对样本数据的散点图和样本相关系数，可以推断两个变量是否存在相关关系、是正相关还是负相关，以及线性相关程度的强弱等。进

### 8.2.2 一元线性回归模型参数的最小二乘估计

在一元线性回归模型中，表达式  $Y=bx+a+e$  刻画的是变量  $Y$  与变量  $x$  之间的线性相关关系，其中参数  $a$  和  $b$  未知，需要根据成对样本数据进行估计。由模型的建立过程可知，参数  $a$  和  $b$  刻画了变量  $Y$  与变量  $x$  的线性关系，因此通过成对样本数据估计这两个参数，相当于寻找一条适当的直线，使表示成对样本数据的这些散点在整体上与这条直线最接近。



#### 探究

利用散点图 8.2-1 找出一条直线，使各散点在整体上与此直线尽可能接近。

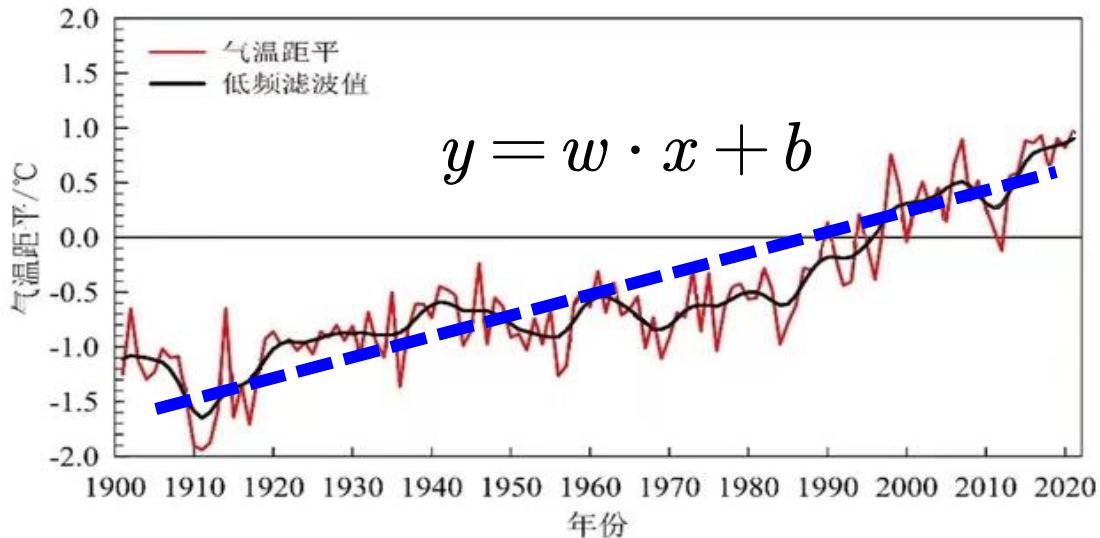
有的同学可能会想，可以采用测量的方法，先画出一条直线，测量出各点与它的距离，然后移动直线，到达一个使距离的和最小的位置。测量出此时的斜率和截距，就可得到一条直线，如图 8.2-2 所示。

# 回归问题的解决

步骤1：定义  
一组函数集合

步骤2：定义函  
数好坏的标准

步骤3：选取  
最优的函数



$$L(w, b) = \frac{1}{2N} \sum_{i=1}^N (\hat{y}_i - (w \cdot x_i + b))^2$$

高中的解法

通过最小二乘法确定最优的参数  $w$  和  $b$

人工智能方法

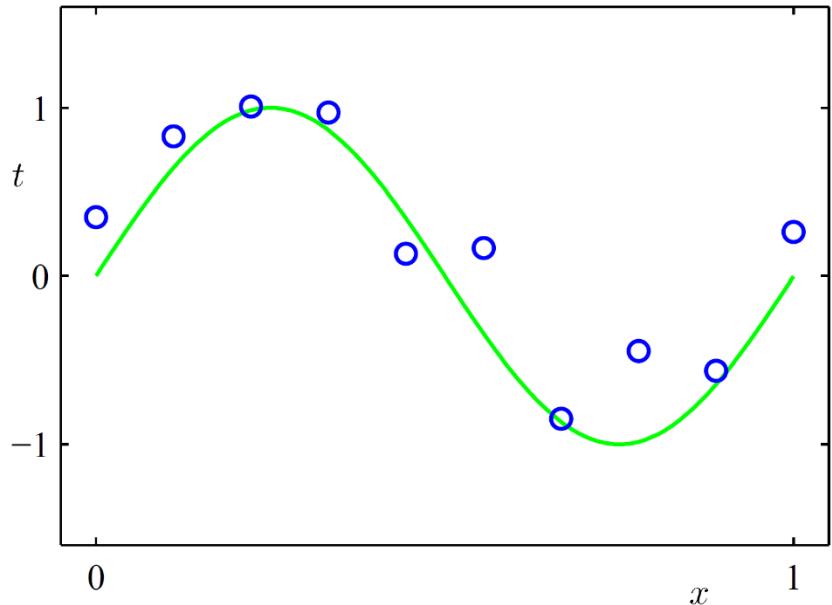
利用梯度下降法求解参数

# 多项式拟合问题

步骤1：定义  
一组函数集合

步骤2：定义函  
数好坏的标准

步骤3：选取  
最优的函数



模型

$$y(x, \omega) = \omega_0 + \omega_1 x + \omega_2 x^2 + \cdots + \omega_M x^M$$

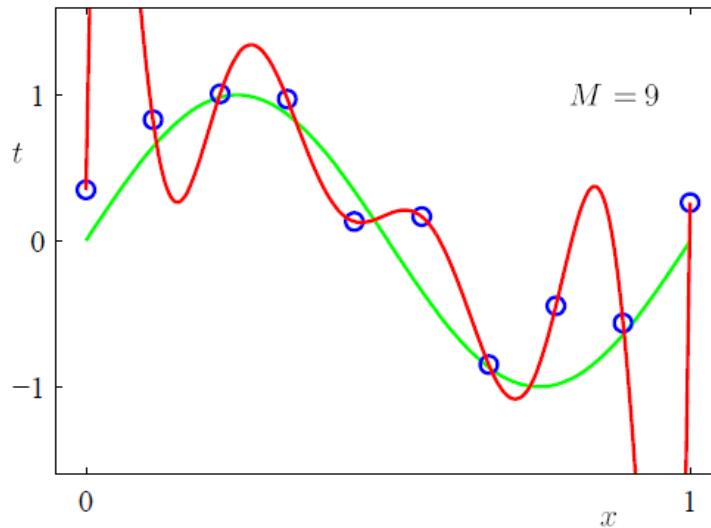
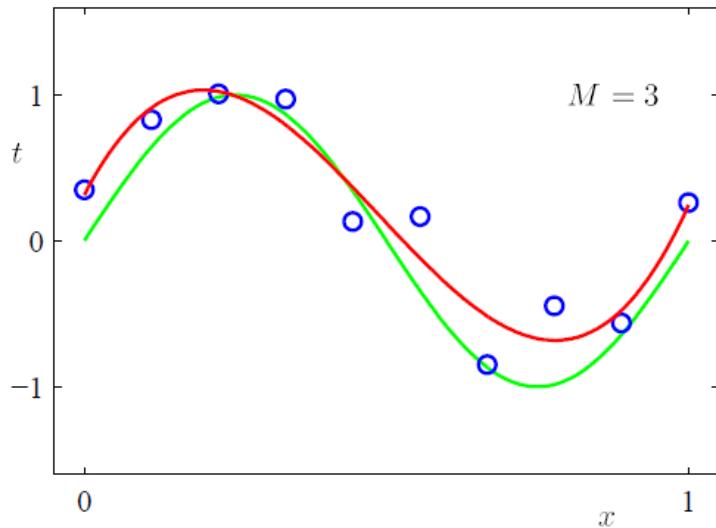
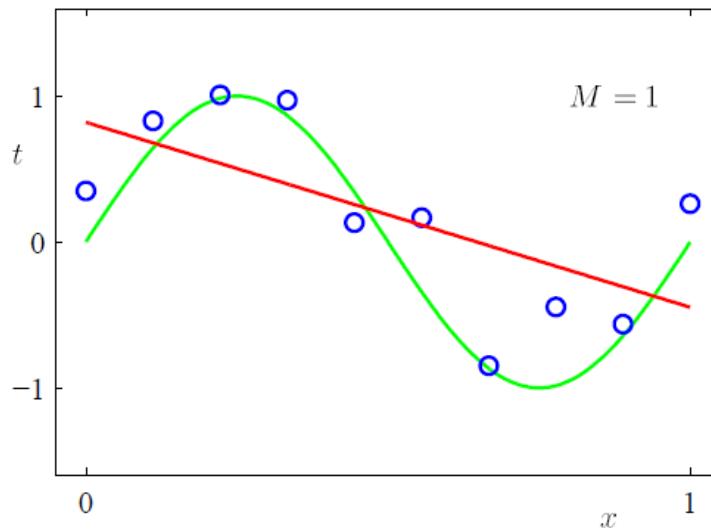
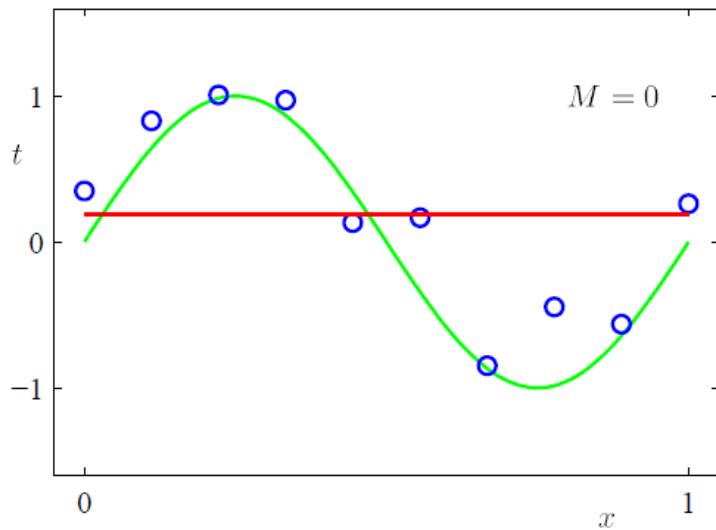
损失函数

$$L(\omega) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \omega) - t_n\}^2$$

机器学习本质：利用数据样本，估计模型中的参数取值

# 多项式拟合问题

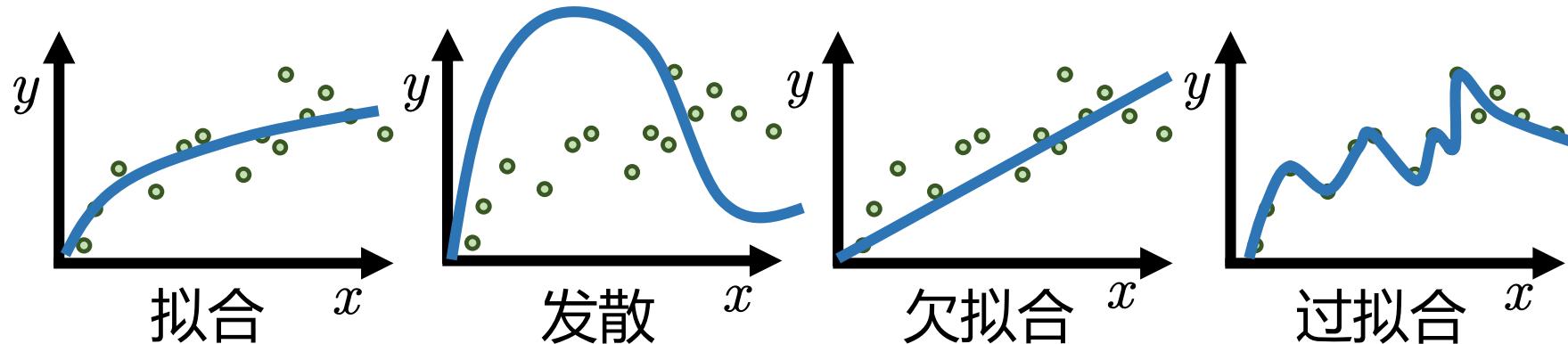
如何在拟合能力和模型复杂度  
之间取得一个好的平衡



# 模型性能的评价

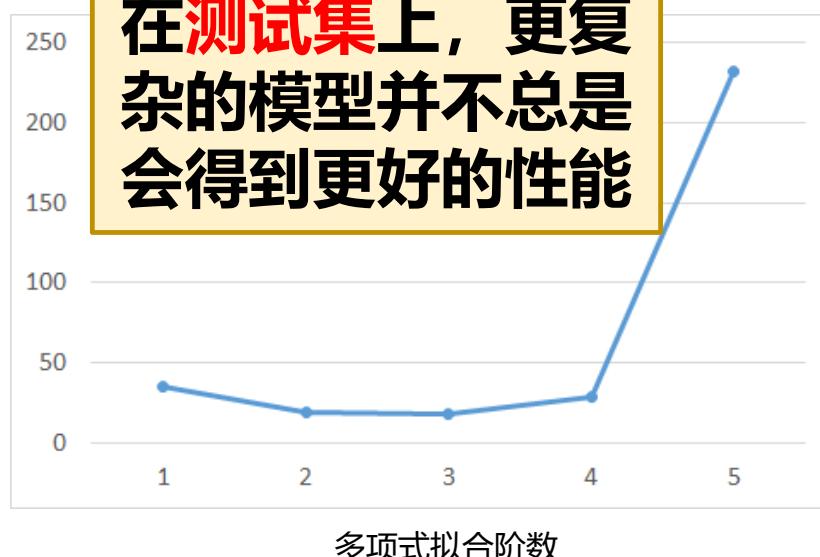
机器学习的本质是找一个函数

利用样本估计模型参数



在测试集上，更复杂的模型并不总是会得到更好的性能

测试集损失



如何在拟合能力和复杂度之间取得一个好的平衡



# 模型选择

$$y = b + w_1 \cdot x$$

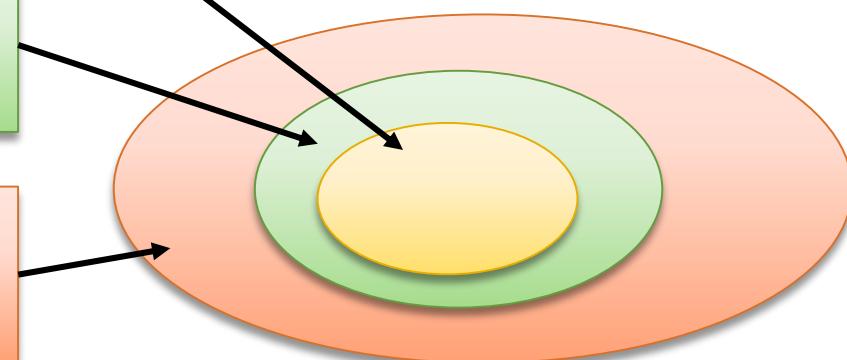
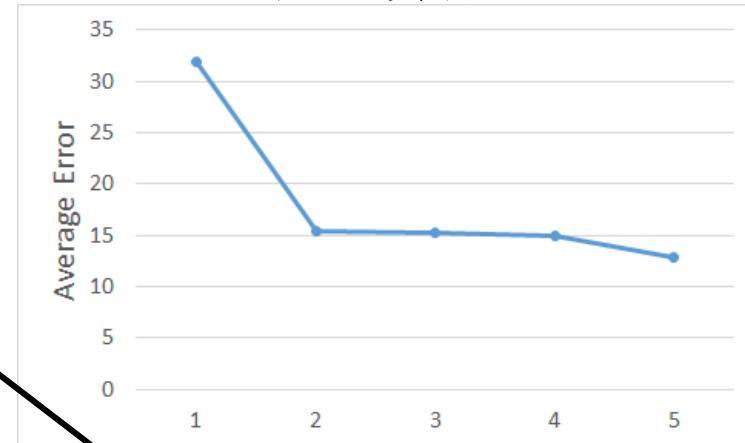
$$y = b + w_1 \cdot x + w_2 \cdot x^2$$

$$y = b + w_1 \cdot x + w_2 \cdot x^2 + w_3 \cdot x^3$$

$$y = b + w_1 \cdot x + w_2 \cdot x^2 + w_3 \cdot x^3 + w_4 \cdot x^4$$

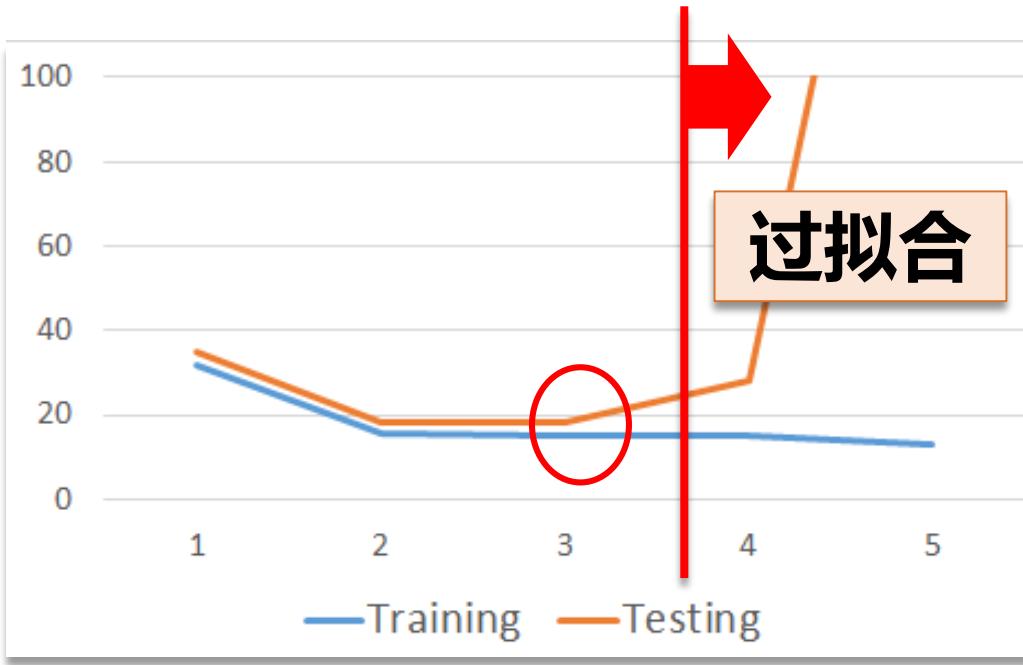
$$y = b + w_1 \cdot x + w_2 \cdot x^2 + w_3 \cdot x^3 + w_4 \cdot x^4 + w_5 \cdot x^5$$

训练数据



在存在最优函数的情况下，更复杂的模型会降低训练数据的误差。

# 模型选择



最高阶数	训练集	测试集
1	31.9	35.0
2	15.4	18.4
3	15.3	18.1
4	14.9	28.2
5	12.8	232.1

更复杂的模型并不总是能够在测试数据上获得更好的性能

如何避免过拟合的情况



选择合适的模型

# 内容导览

---



线性回归：多项式拟合



线性分类：概率生成模型



Logistic回归：神经元模型



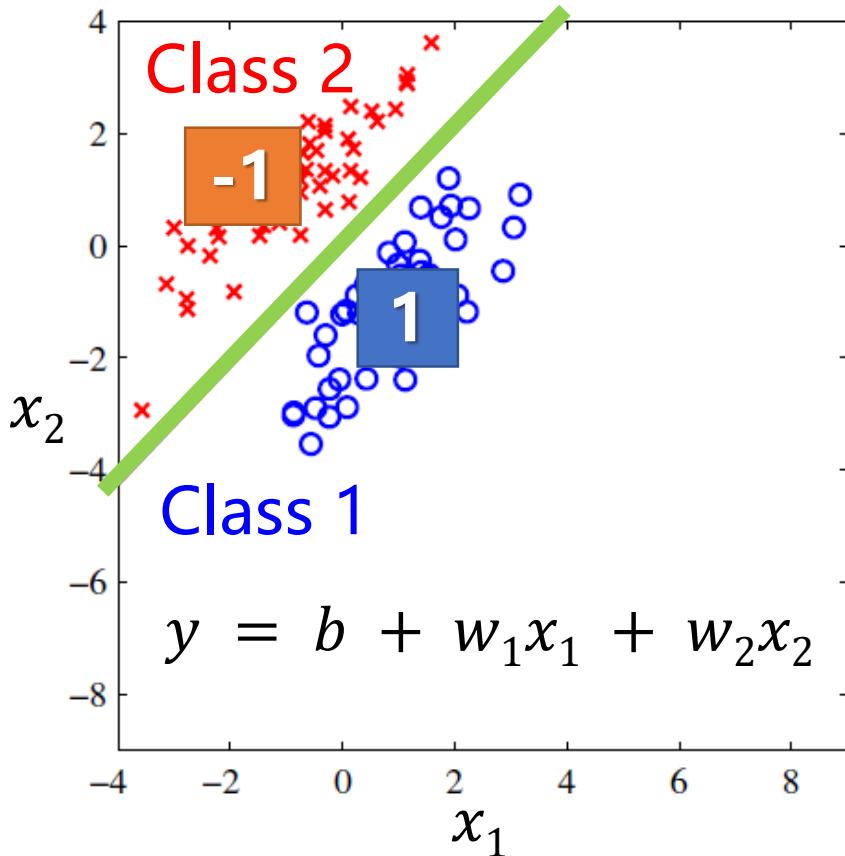
多分类问题



线性模型的局限与非线性模型

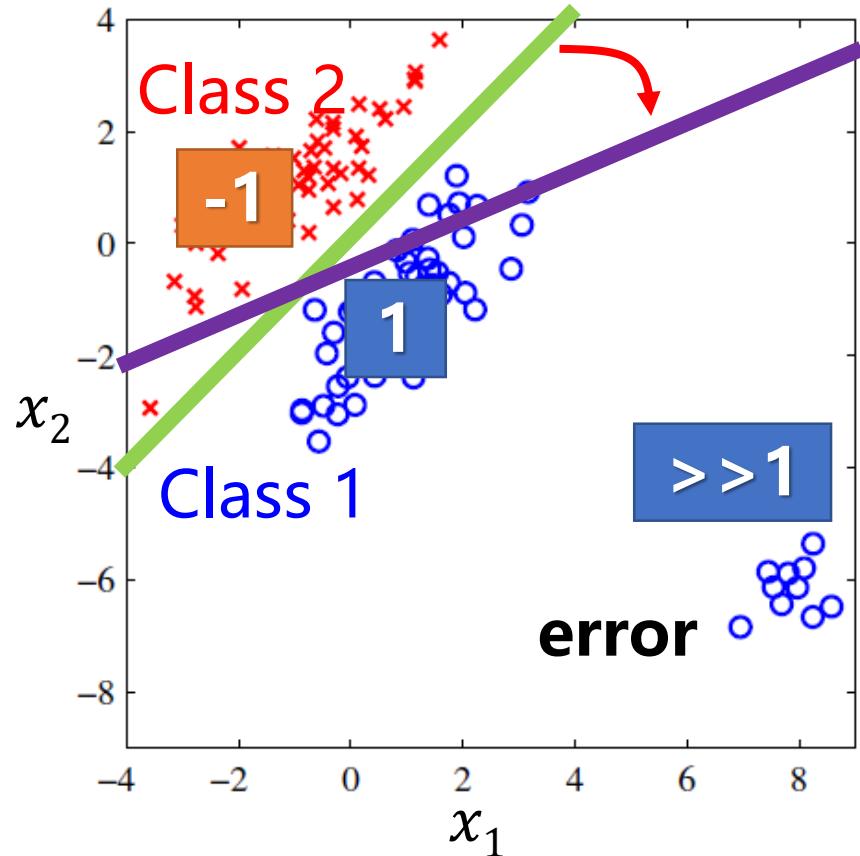
# 能否利用回归思路解决分类问题?

$$b + w_1x_1 + w_2x_2 = 0$$



$$y = b + w_1x_1 + w_2x_2$$

to decrease error

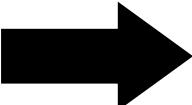


Penalize to the examples that are “too correct” ...

(Bishop, P186)

# 解决分类问题的三步骤

## □ 函数(模型)

$x$  

$$g(x) > 0$$

else

输出 = 类别1

输出 = 类别2

$$f(x)$$

## □ 模型度量(损失函数)

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

$f$ 在训练数据上得到  
错误结果的次数

## □ 寻找最佳函数

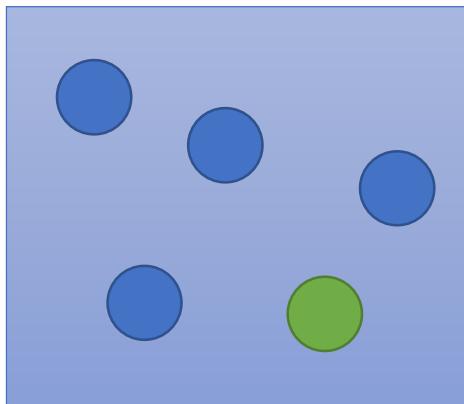
- 示例：感知器，SVM

不是当前使用的  
主流方法

# 概率生成问题

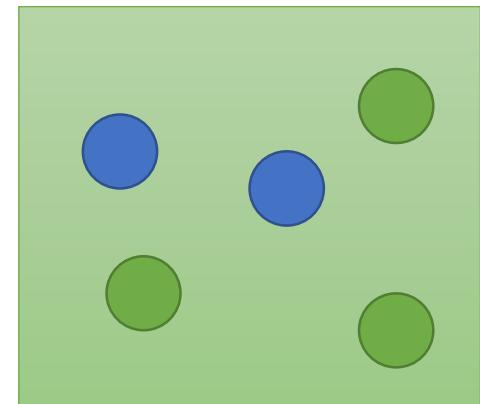
**Box 1**

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



**Box 2**

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



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

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

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

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

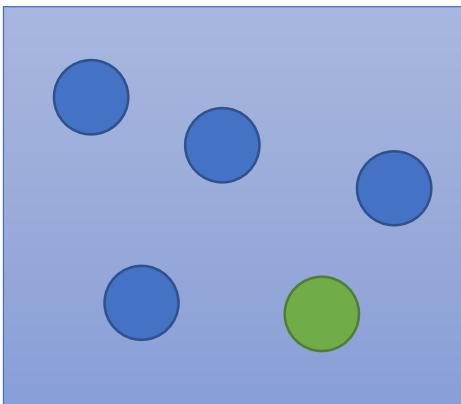
- 是从某一个盒子取出的蓝色球  
**是从哪个盒子取出?**

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

# 分类问题实质：概率生成模型

**Class 1**

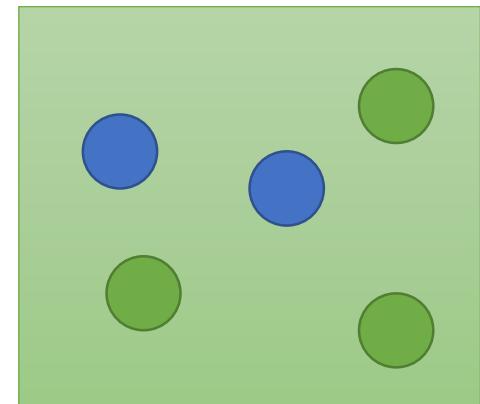
$$P(C_1)$$



$$P(x|C_1)$$

**Class 2**

$$P(C_2)$$



$$P(x|C_2)$$

对于给定样本 $x$ , 其属于某一类别的概率为

贝叶斯公式

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

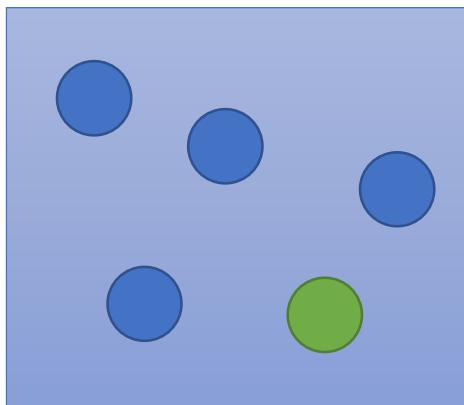
概率生成模型:  $P(x) = P(x|C_1)P(C_1) + P(x|C_2)P(C_2)$

生成模型: 获得数据概率分布

# 概率生成模型：先验概率获取

**Class 1**

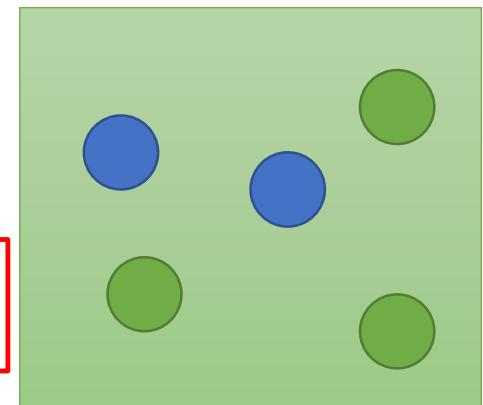
$$P(C_1)$$



$$P(x|C_1)$$

**Class 2**

$$P(C_2)$$



$$P(x|C_2)$$

以Pokémon分类为例，水系(class 1)和常规系(class 2)，  
训练数据包含79只水系Pokémon，61只常规系Pokémon

**先验概率**

$$P(C_1) = \frac{79}{79 + 61} = 0.56$$

$$P(C_2) = \frac{61}{79 + 61} = 0.44$$

# 新样本来自某一类的概率

$$P(x|C_1) = ?$$

$$P($$



$$| \text{水系} ) = ?$$

每个Pokémon都通过其属性表示为一个向量

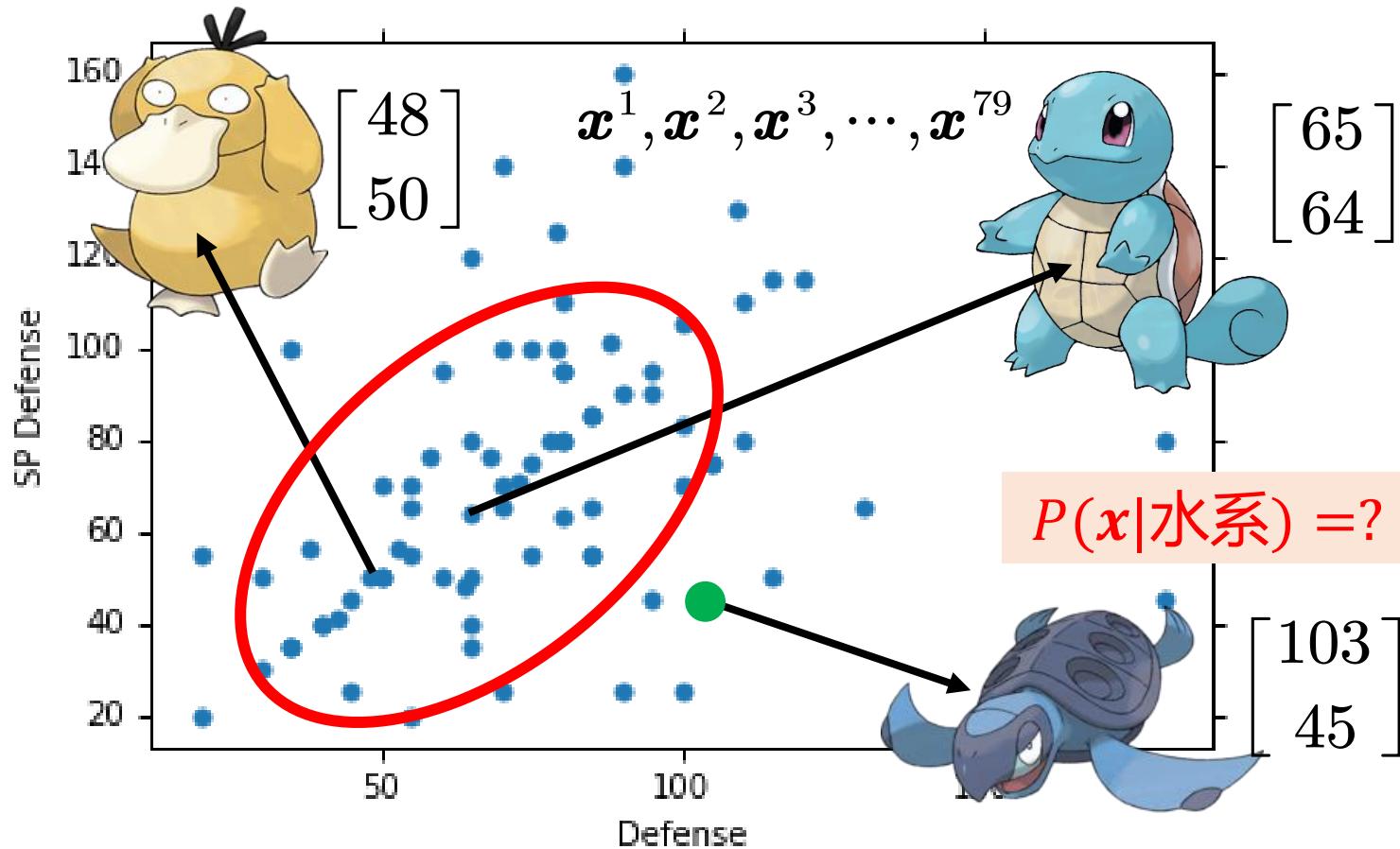


Pokémon属性组成的向量可看作是**特征**



- Pokémon实例源自李宏毅机器学习公开课

# 新样本来自某一类的概率



假设所有点采样自某个高斯分布

# 高斯分布

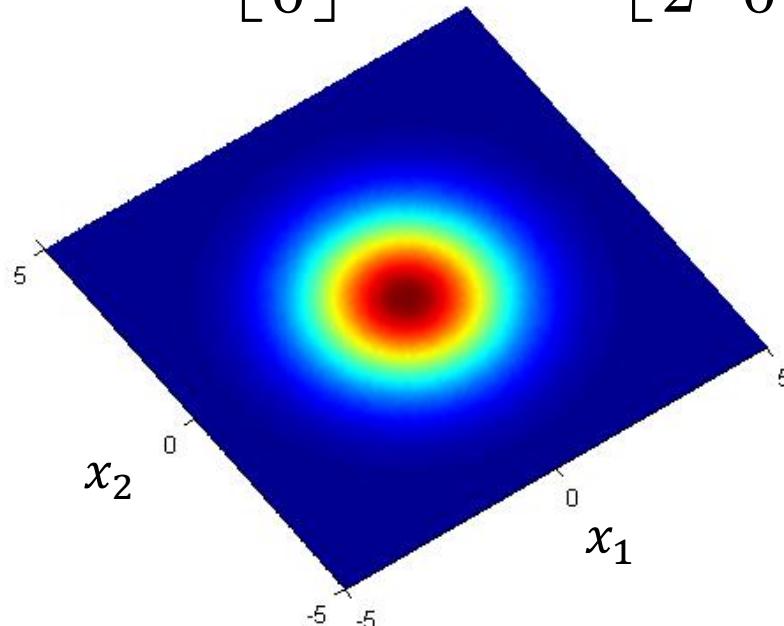
$$f_{\mu, \Sigma}(x) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

□ 输入：向量  $x$       输出：抽样得到  $x$  的概率

□ 均值  $\mu$  和 协方差  $\Sigma$  决定了函数的形状

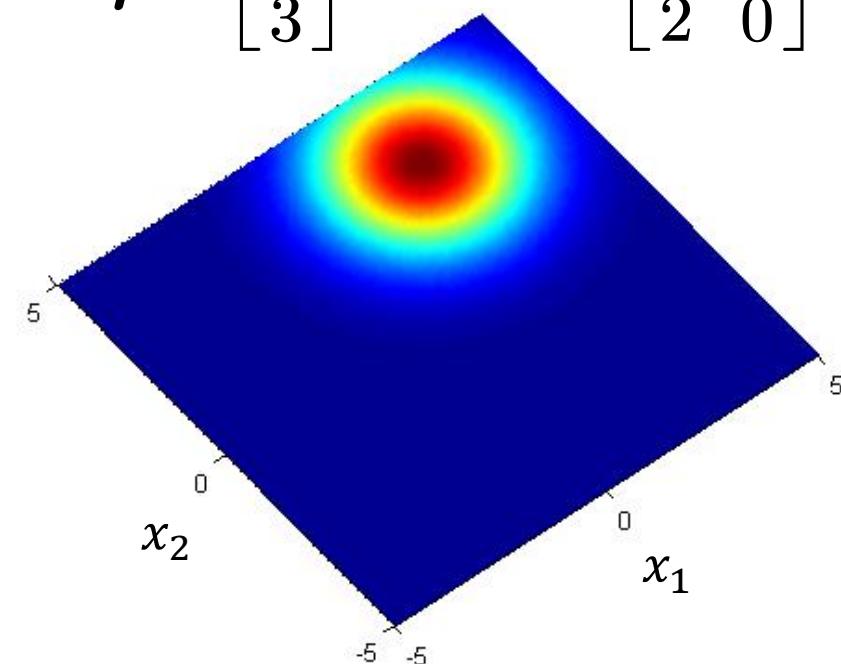
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

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



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

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

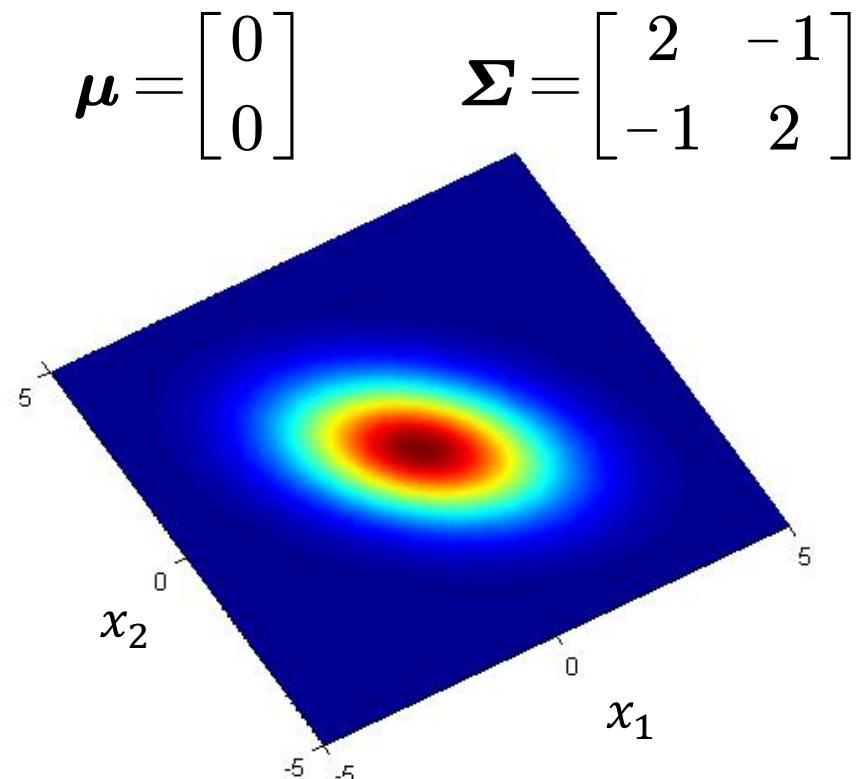
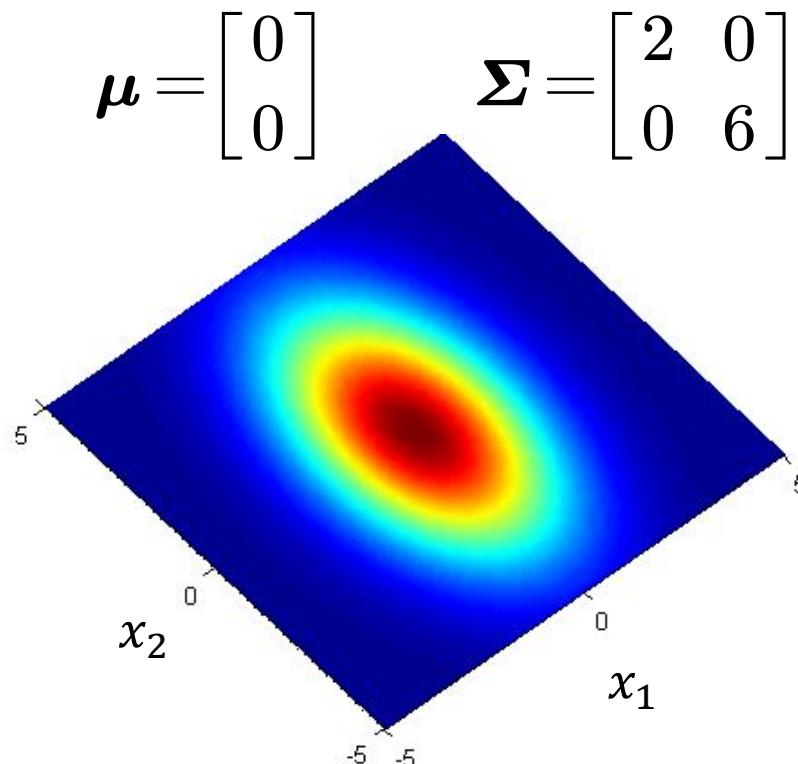


# 高斯分布

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

□ 输入：向量  $x$       输出：抽样得到  $x$  的概率

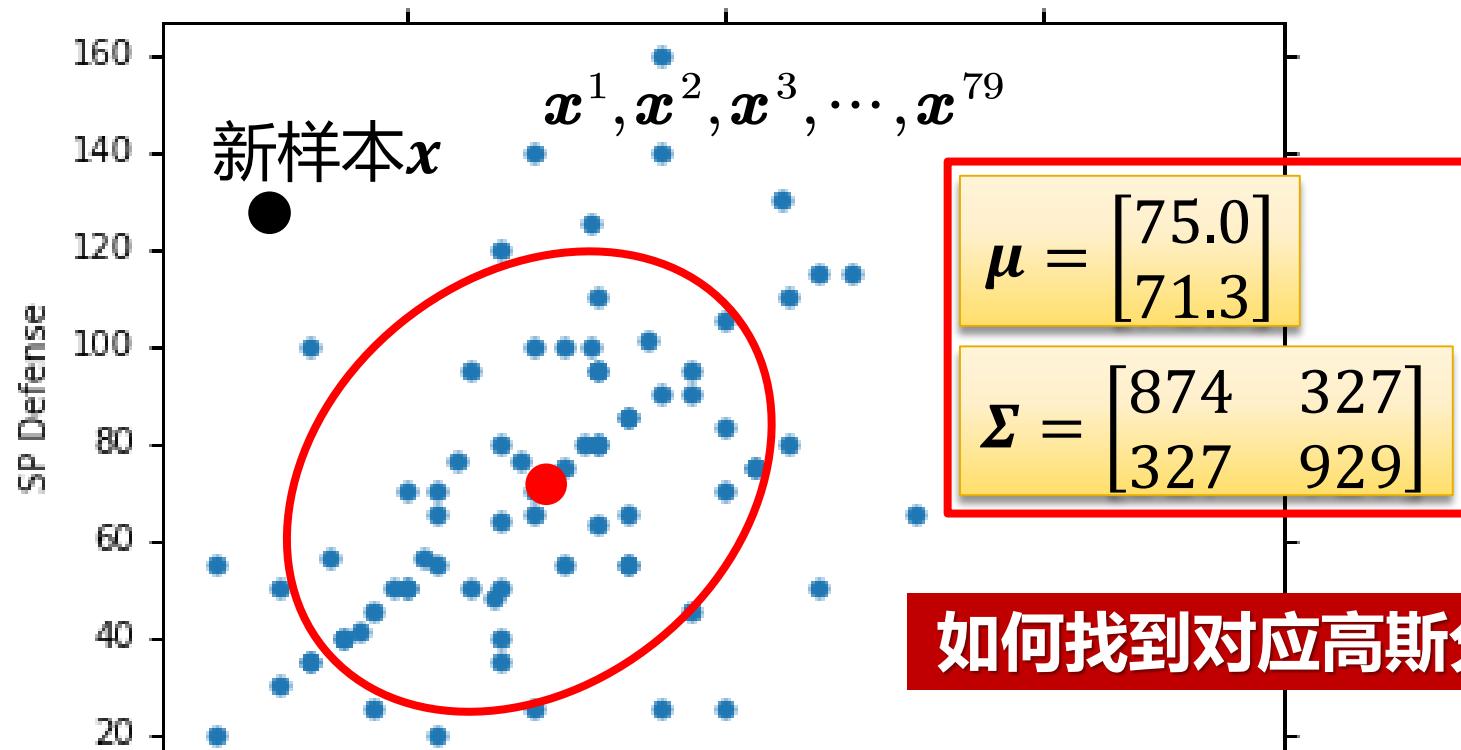
□ 均值  $\mu$  和 协方差  $\Sigma$  决定了函数的形状



# 新样本来自某一类的概率

假设所有点均采样于某一高斯分布

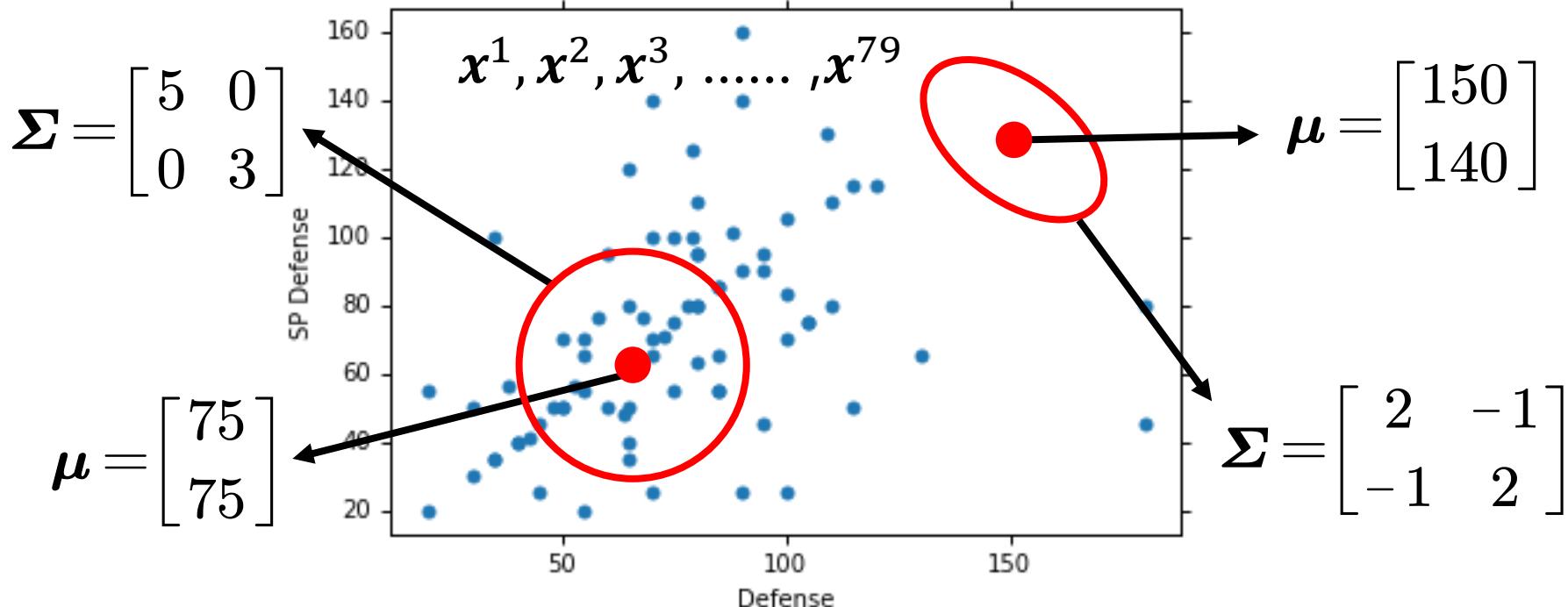
➤ 找出它们背后的高斯分布 → 得到生成新样本的概率



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

# 最大似然估计

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



- 任意均值  $\mu$  和任意协方差矩阵  $\Sigma$  所对应的高斯分布都能生成这些点，**只是生成的概率(似然函数)不同**
- 具有均值  $\mu$  和协方差矩阵  $\Sigma$  的某个高斯分布的似然函数=生成高斯样本  $x^1, x^2, x^3, \dots, x^{79}$  的概率

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

# 最大似然估计

训练样本：水系宝可梦  $x^1, x^2, x^3, \dots, x^{79}$

假设  $x^1, x^2, x^3, \dots, x^{79}$  由高斯分布似然函数  $(\mu, \Sigma)$  采样生成

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

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

$$\mu^*, \Sigma^* = \underset{\mu, \Sigma}{\operatorname{arg\,max}} L(\mu, \Sigma)$$

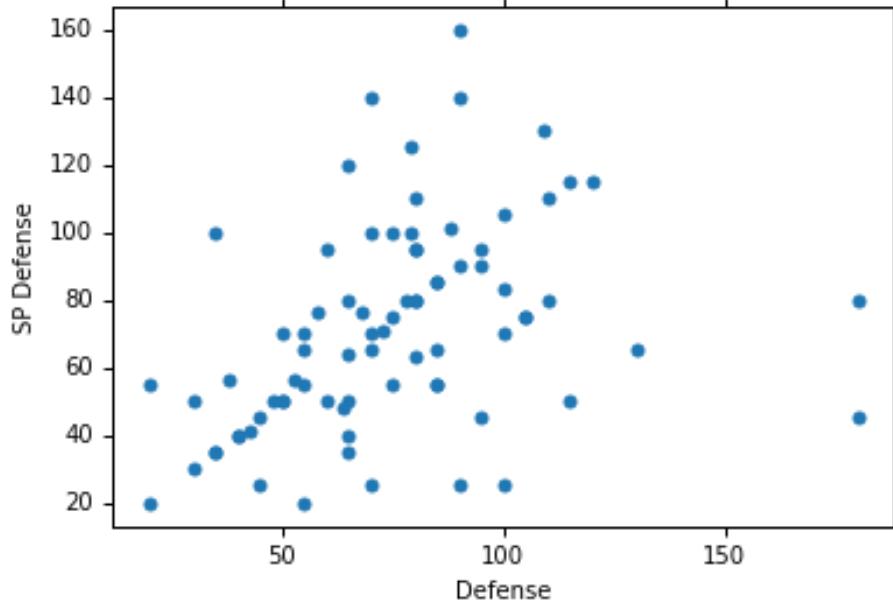
$$\mu^* = \frac{1}{79} \sum_{n=1}^{79} x^n \quad \Sigma^* = \frac{1}{79} \sum_{n=1}^{79} (x^n - \mu^*) (x^n - \mu^*)^T$$

均值

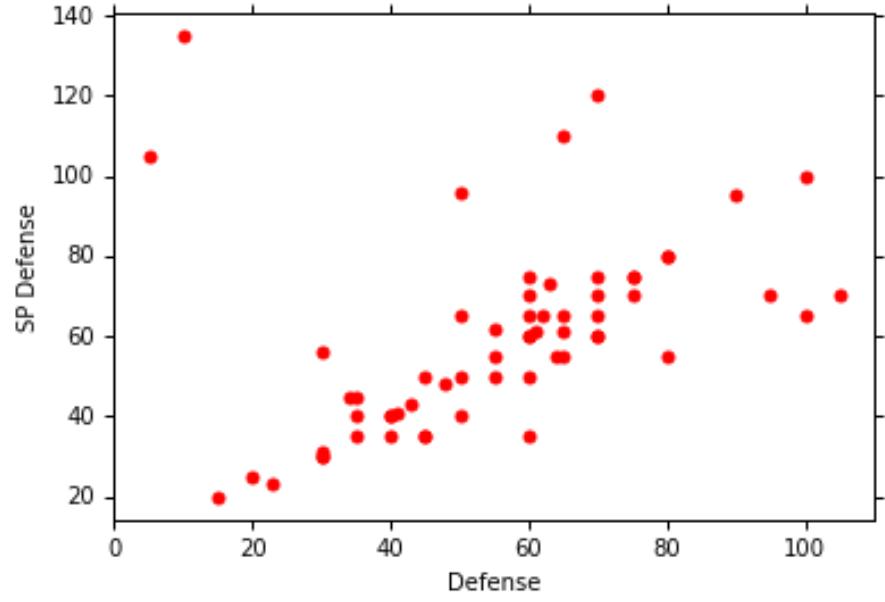
协方差矩阵

# 最大似然估计

Class 1



Class 2



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

$$\boldsymbol{\mu}^2 = \begin{bmatrix} 55.6 \\ 59.8 \end{bmatrix} \quad \boldsymbol{\Sigma}^2 = \begin{bmatrix} 874 & 422 \\ 422 & 685 \end{bmatrix}$$

# 依后验概率进行分类

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

$$\boldsymbol{\mu}^1 = \begin{bmatrix} 75.0 \\ 71.3 \end{bmatrix} \quad \Sigma^1 = \begin{bmatrix} 874 & 327 \\ 327 & 929 \end{bmatrix} \quad P(C_1) = 79/(79 + 61) = 0.56$$

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

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

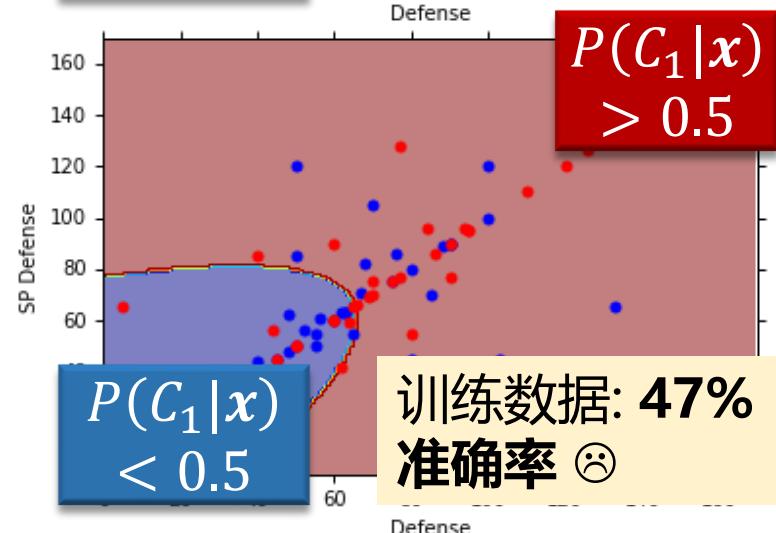
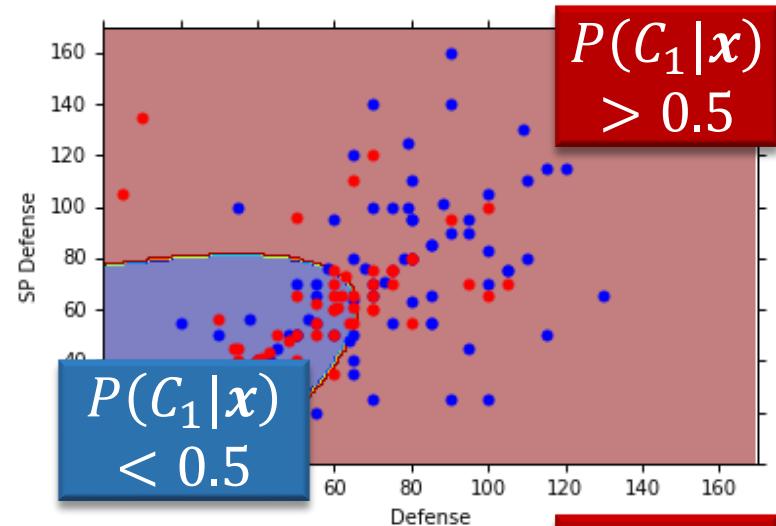
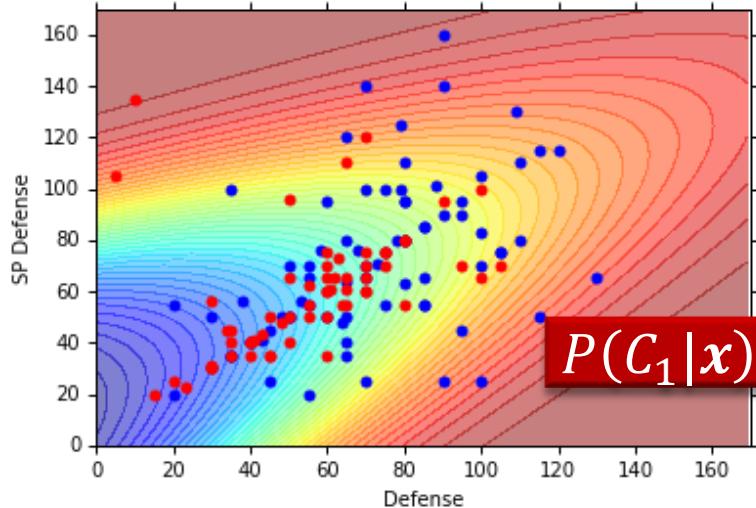
$$f_{\mu^2, \Sigma^2}(\mathbf{x}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma^2|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}^2)^T (\Sigma^2)^{-1} (\mathbf{x} - \boldsymbol{\mu}^2) \right\}$$

$$\boldsymbol{\mu}^2 = \begin{bmatrix} 55.6 \\ 59.8 \end{bmatrix} \quad \Sigma^2 = \begin{bmatrix} 874 & 422 \\ 422 & 685 \end{bmatrix}$$

If  $P(C_1|\mathbf{x}) > 0.5$   
样本  $\mathbf{x}$  属于 Class 1

# 结果分析

口蓝点:  $C_1$  (水系), 红点:  $C_2$  (一般系)

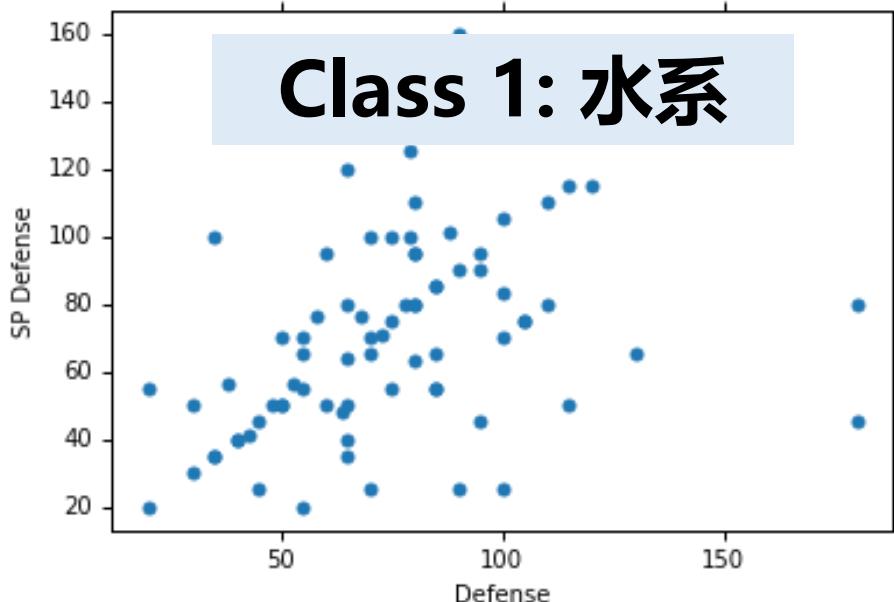


口若考虑总共6个特征

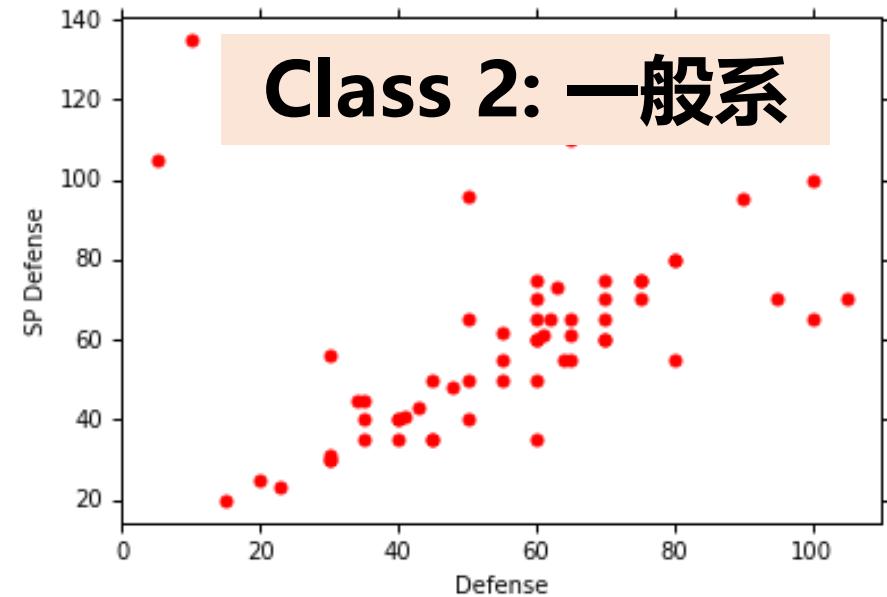
- $\mu^1, \mu^2$ : 6维向量
- $\Sigma^1, \Sigma^2$ :  $6 \times 6$  矩阵

新测试数据:  
64% 准确率 ☺

# 改进模型



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



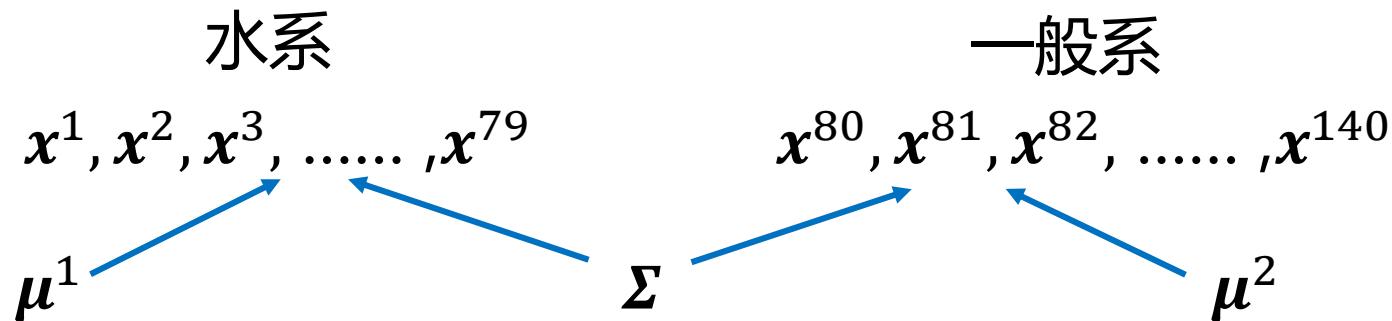
$$\boldsymbol{\mu}^2 = \begin{bmatrix} 55.6 \\ 59.8 \end{bmatrix} \quad \boldsymbol{\Sigma}^2 = \begin{bmatrix} 874 & 422 \\ 422 & 685 \end{bmatrix}$$

## 改进策略

如果两个分布共用同一个 $\Sigma$ , 可以节省参数量

# 改进模型

## 口最大似然估计



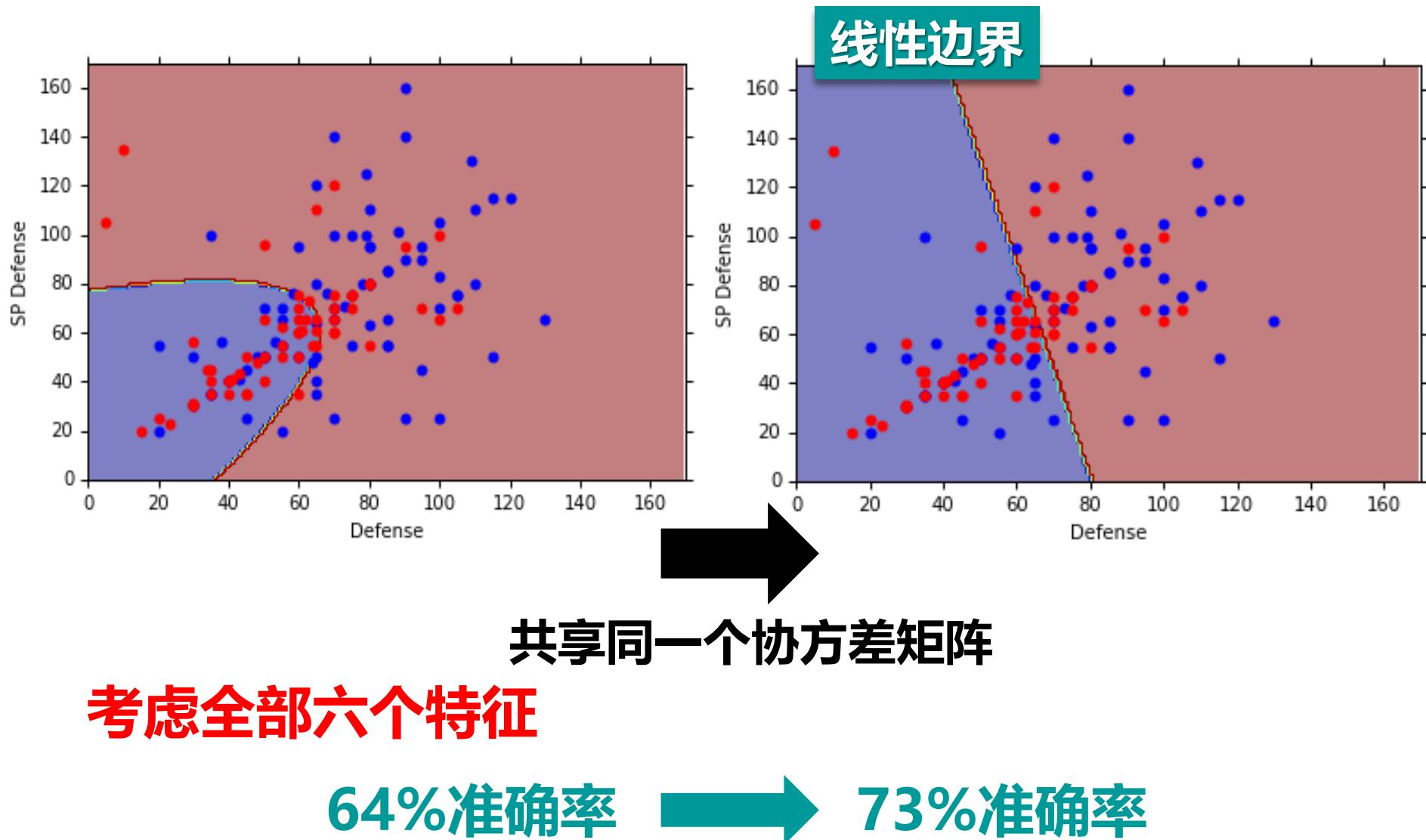
➤ 找到使似然函数  $L(\mu^1, \mu^2, \Sigma)$  最大的  $\mu^1, \mu^2, \Sigma$

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

$$\times f_{\mu^2, \Sigma}(x^{80}) \cdot f_{\mu^2, \Sigma}(x^{81}) \cdots f_{\mu^2, \Sigma}(x^{140})$$

➤  $\mu^1$  和  $\mu^2$  计算方法相同  $\Sigma = \frac{79}{140} \Sigma^1 + \frac{61}{140} \Sigma^2$

# 改进模型



# 内容导览

---



线性回归：多项式拟合



线性分类：概率生成模型



Logistic回归：神经元模型



多分类问题



线性模型的局限与非线性模型

# 从深度学习角度看待分类问题

---

## 口将分类问题看作条件概率估计问题

- 引入非线性函数 $g$ 来预测类别标签的条件概率 $p(y = c|x)$
- 以二分类为例，

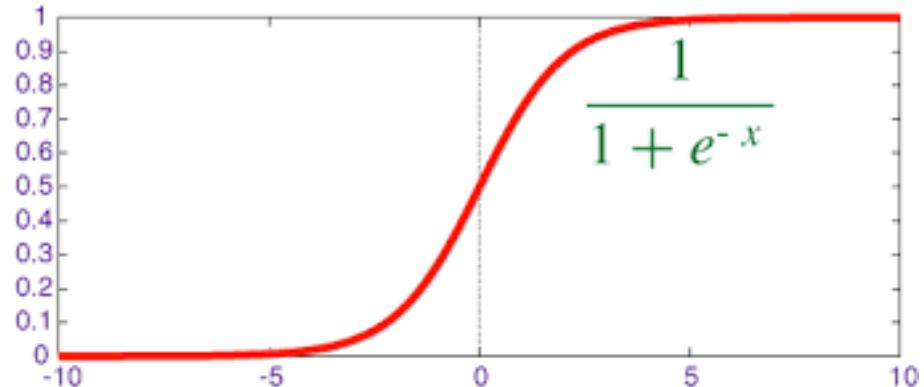
$$p(y = 1|x) = g(f(\mathbf{x}; \mathbf{w}))$$

- 函数 $f$ : 线性函数
- 函数 $g$ : 把线性函数的值域从实数区间“挤压”到了 $(0,1)$ 之间，可以用来表示概率。

# 预备知识：Logistic函数

## □ Logistic函数

$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$



## □ Logistic回归

$$p(y = 1|x) = \sigma(\mathbf{w}^\top \mathbf{x})$$

$$\triangleq \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})}$$

# 预备知识：学习准则

---

□ 模型预测条件概率  $p_\theta(y|x)$

$$p_\theta(y = 1|x) = \sigma(\mathbf{w}^T \mathbf{x})$$

□ 真实条件概率  $p_r(y|x)$

➤ 对于一个样本  $(x, y^*)$ , 其真实条件概率为

$$p_r(y = 1|x) = y^*$$

$$p_r(y = 0|x) = 1 - y^*$$

**如何衡量两个条件分布的差异?**

# 预备知识：信息熵

---

口在信息论中，熵用来衡量一个随机事件的不确定性

➤ 自信息(Self Information)  $I(x) = -\log p(x)$

➤ 熵  $H(X) = \mathbb{E}_X [I(x)]$

$$= \mathbb{E}_X [-\log p(x)]$$

$$= - \sum_{x \in \mathcal{X}} p(x) \log p(x)$$

➤ 熵越高，则随机变量的信息越多；

➤ 熵越低，则随机变量的信息越少

➤ 在对分布  $q(y)$  的符号进行编码时，熵  $H(q)$  也是**理论上最优的平均编码长度**，这种编码方式称为**熵编码**(Entropy Coding)，广泛应用于**数据压缩的理论性能分析**

# 预备知识：交叉熵(Cross Entropy)

口 交叉熵是按照**估计分布**为 $q$ 的最优编码对**真实分布**为 $p$ 的信息进行编码的长度

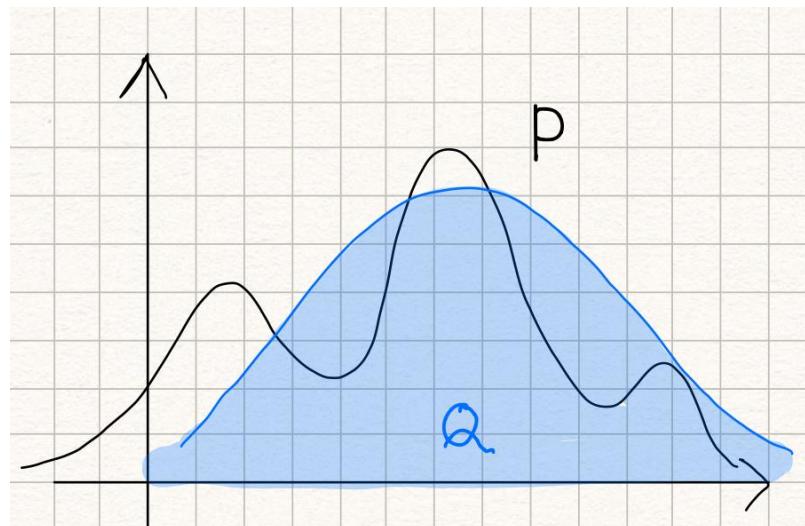
$$\begin{aligned} H(p, q) &= \mathbb{E}_p [-\log q(x)] \\ &= - \sum_x p(x) \log q(x) \end{aligned}$$

- 在给定  $q$  的情况下，如果  $p$  和  $q$  越接近，交叉熵越小
- 如果  $p$  和  $q$  越远，交叉熵就越大

# 预备知识：KL散度(Kullback-Leibler Divergence)

KL散度是用**估计分布** $q$ 来近似**真实分布** $p$ 时所造成的信息损失量

- KL散度的物理含义是：按照估计分布为 $q$ 的最优编码对真实分布为 $p$ 的信息进行编码，其**平均编码长度(即交叉熵)** $H(p, q)$ 和 $p$ 的**最优平均编码长度(即熵)** $H(p)$ 之间的**差异**



$$\begin{aligned} KL(p, q) &= H(p, q) - H(p) \\ &= \sum_x p(x) \log \frac{p(x)}{q(x)} \\ &\geq 0 \end{aligned}$$

# 分类问题的基本步骤

## 口建立模型

后验概率为分类准则

$x$  →

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

If  $P(C_1|x) > 0.5$ , 输出: class 1  
否则, 输出: class 2

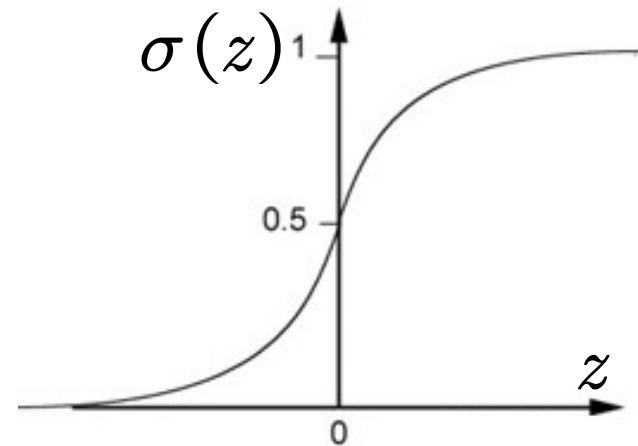
## 口确定评价函数

- 使得似然函数（生成数据的概率）最大的均值  $\mu$  和协方差  $\Sigma$

## 口选择最佳函数

# 依后验概率分类

$$\begin{aligned} P(C_1|\mathbf{x}) &= \frac{P(\mathbf{x}|C_1)P(C_1)}{P(\mathbf{x}|C_1)P(C_1) + P(\mathbf{x}|C_2)P(C_2)} \\ &= \frac{1}{1 + \frac{P(\mathbf{x}|C_2)P(C_2)}{P(\mathbf{x}|C_1)P(C_1)}} \quad \text{记 } z = \ln \frac{P(\mathbf{x}|C_1)P(C_1)}{P(\mathbf{x}|C_2)P(C_2)} \\ &= \frac{1}{1 + \exp(-z)} \\ &= \sigma(z) \quad \text{Sigmoid 函数} \end{aligned}$$



互动

还记得上学期《人工智能数学基础》提到的Sigmoid函数的优点和缺点有哪些吗？

# 依后验概率分类

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

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

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

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

# 依后验概率分类

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

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

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

$$= \frac{1}{2} \ln \frac{|\boldsymbol{\Sigma}^2|}{|\boldsymbol{\Sigma}^1|} - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}^1)^T (\boldsymbol{\Sigma}^1)^{-1} (\mathbf{x} - \boldsymbol{\mu}^1) \\ + \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}^2)^T (\boldsymbol{\Sigma}^2)^{-1} (\mathbf{x} - \boldsymbol{\mu}^2)$$

# 依后验概率分类

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

$$\begin{aligned} &= \mathbf{x}^T (\boldsymbol{\Sigma}^1)^{-1} \mathbf{x} - \mathbf{x}^T (\boldsymbol{\Sigma}^1)^{-1} \boldsymbol{\mu}^1 - (\boldsymbol{\mu}^1)^T (\boldsymbol{\Sigma}^1)^{-1} \mathbf{x} \\ &\quad + (\boldsymbol{\mu}^1)^T (\boldsymbol{\Sigma}^1)^{-1} \boldsymbol{\mu}^1 \end{aligned}$$

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

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

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

综上  $z = \frac{1}{2} \ln \frac{|\boldsymbol{\Sigma}^2|}{|\boldsymbol{\Sigma}^1|} - \frac{1}{2} \mathbf{x}^T (\boldsymbol{\Sigma}^1)^{-1} \mathbf{x} + (\boldsymbol{\mu}^1)^T (\boldsymbol{\Sigma}^1)^{-1} \mathbf{x}$

$$- \frac{1}{2} (\boldsymbol{\mu}^1)^T (\boldsymbol{\Sigma}^1)^{-1} \boldsymbol{\mu}^1 + \frac{1}{2} \mathbf{x}^T (\boldsymbol{\Sigma}^2)^{-1} \mathbf{x} - (\boldsymbol{\mu}^2)^T (\boldsymbol{\Sigma}^2)^{-1} \mathbf{x}$$

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

# 依后验概率分类

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

$$\Sigma^1 = \Sigma^2 = \Sigma$$

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

$$\mathbf{w}^T$$

$$b$$

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

前述概率生成模型做法

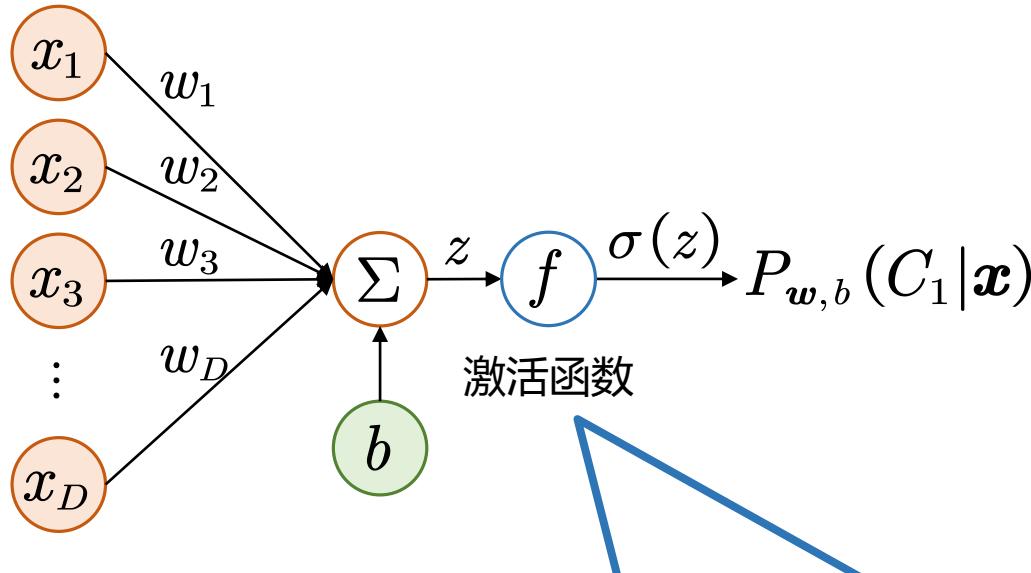
先估计  $N_1, N_2, \boldsymbol{\mu}^1, \boldsymbol{\mu}^2, \Sigma$ , 最大似然估计计算得到  $\mathbf{w}^T, b$

能够直接求解  
 $\mathbf{w}^T$  和  $b$  吗



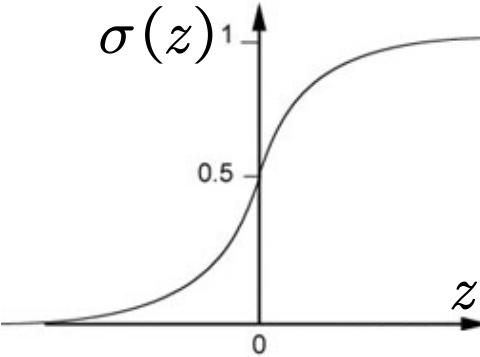
# 神经元解决分类问题：①模型建立

$$P_{\mathbf{w}, b}(C_1 | \mathbf{x}) = \sigma(z) \quad z = \mathbf{w}^T \cdot \mathbf{x} + b = \sum_i w_i x_i + b$$



Sigmoid 函数

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

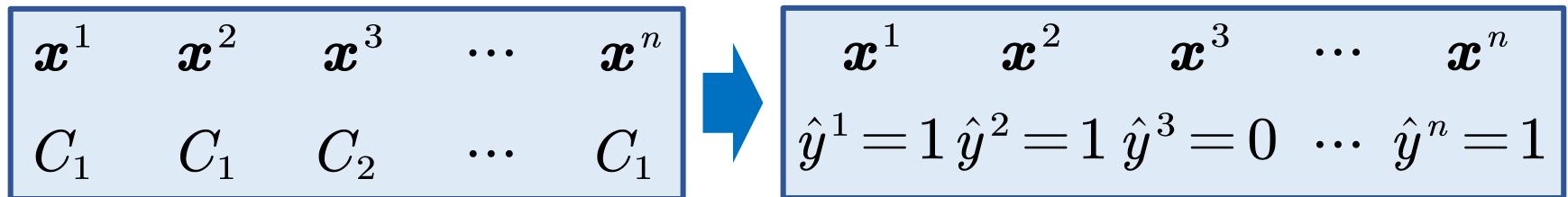


$$\begin{cases} P_{\mathbf{w}, b}(C_1 | \mathbf{x}) \geq 0.5 & \text{class 1} \\ P_{\mathbf{w}, b}(C_1 | \mathbf{x}) < 0.5 & \text{class 2} \end{cases}$$

$$\begin{cases} z \geq 0 & \text{class 1} \\ z < 0 & \text{class 2} \end{cases}$$

# 神经元解决分类问题：②损失函数

口 神经元模型  $f_{w,b}(x) = P_{w,b}(C_1|x)$



► 给定一组  $w$  和  $b$ ，神经元模型估计对应类别的概率是

$$L(w, b) = f_{w,b}(x^1) f_{w,b}(x^2) (1 - f_{w,b}(x^3)) \cdots f_{w,b}(x^n)$$

► 当  $L(w, b)$  取最大值时，能得到最有可能的  $w^*$  和  $b^*$

$$w^*, b^* = \underset{w, b}{\operatorname{argmax}} L(w, b)$$

$$w^*, b^* = \underset{w, b}{\operatorname{argmin}} (-\ln L(w, b))$$

$$-\ln L(w, b) =$$

$$-\ln f_{w,b}(x^1) \quad \rightarrow -(\textcolor{lightgreen}{1} \ln f(x^1) + \textcolor{lightgreen}{0} \ln(1 - f(x^1)))$$

$$-\ln f_{w,b}(x^2) \quad \rightarrow -(\textcolor{lightgreen}{1} \ln f(x^2) + \textcolor{lightgreen}{0} \ln(1 - f(x^2)))$$

$$-\ln(1 - f_{w,b}(x^3)) \quad \rightarrow -(\textcolor{lightorange}{0} \ln f(x^3) + \textcolor{lightorange}{1} \ln(1 - f(x^3)))$$

# 神经元解决分类问题：②损失函数

$$L(\mathbf{w}, b) = f_{\mathbf{w}, b}(\mathbf{x}^1) f_{\mathbf{w}, b}(\mathbf{x}^2) (1 - f_{\mathbf{w}, b}(\mathbf{x}^3)) \cdots f_{\mathbf{w}, b}(\mathbf{x}^n)$$

$$-\ln L(\mathbf{w}, b) = -\ln f_{\mathbf{w}, b}(\mathbf{x}^1) - \ln f_{\mathbf{w}, b}(\mathbf{x}^2) - \ln (1 - f_{\mathbf{w}, b}(\mathbf{x}^3)) \cdots$$

$\hat{y}^n$ : 1 代表 class 1, 0 代表 class 2

$$-\ln L(\mathbf{w}, b) = - \sum_i (\hat{y}^i \ln f_{\mathbf{w}, b}(\mathbf{x}^i) + (1 - \hat{y}^i) \ln (1 - f_{\mathbf{w}, b}(\mathbf{x}^i)))$$

两个伯努利分布的交叉熵

真实分布  $p$ :

$$p(x=1) = \hat{y}^n$$

$$p(x=0) = 1 - \hat{y}^n$$

交叉熵

估计分布  $q$ :

$$q(x=1) = f(\mathbf{x}^n)$$

$$q(x=0) = 1 - f(\mathbf{x}^n)$$

$$H(p, q) = - \sum_{\mathbf{x}} p(\mathbf{x}) \ln q(\mathbf{x})$$

Ground  
Truth  
 $\hat{y}^n = 1$

↔

$\hat{y}^n = 1$        $f(\mathbf{x}^n)$   
 $\hat{y}^n = 0$        $f(\mathbf{x}^n)$  45

# 神经元解决分类问题：③参数优化

$$-\frac{\ln L(\mathbf{w}, b)}{w_i} = -\sum_n \left( \hat{y}^n \frac{\ln f_{\mathbf{w}, b}(\mathbf{x}^n)}{w_i} + (1 - \hat{y}^n) \ln \frac{(1 - f_{\mathbf{w}, b}(\mathbf{x}^n))}{w_i} \right)$$

$$f_{\mathbf{w}, b}(\mathbf{x}) = \sigma(z)$$

$$z = \mathbf{w}^T \cdot \mathbf{x} + b = \sum_i w_i x_i + b$$

$$\frac{\partial \ln f_{\mathbf{w}, b}(\mathbf{x}^n)}{\partial w_i} = \frac{\partial \ln f_{\mathbf{w}, b}(\mathbf{x}^n)}{\partial z} \frac{\partial z}{\partial w_i} \quad \frac{\partial z}{\partial w_i} = x_i^n$$

$$\frac{\partial \ln \sigma(z)}{\partial z} = \frac{1}{\sigma(z)} \sigma(z) (1 - \sigma(z)) = 1 - \sigma(z)$$

$$\frac{\partial \ln (1 - f_{\mathbf{w}, b}(\mathbf{x}^n))}{\partial w_i} = \frac{\partial \ln (1 - f_{\mathbf{w}, b}(\mathbf{x}^n))}{\partial z} \frac{\partial z}{\partial w_i}$$

$$\frac{\partial \ln (1 - \sigma(z))}{\partial z} = -\frac{1}{1 - \sigma(z)} \sigma(z) (1 - \sigma(z)) = -\sigma(z)$$

# 神经元解决分类问题：③参数优化

$$\begin{aligned} -\frac{\partial \ln L(\mathbf{w}, b)}{\partial w_i} &= -\sum_n (\hat{y}^n (1 - f_{\mathbf{w}, b}(\mathbf{x}^n)) - (1 - \hat{y}^n) f_{\mathbf{w}, b}(\mathbf{x}^n)) x_i^n \\ &= -\sum_n (\hat{y}^n - \cancel{\hat{y}^n f_{\mathbf{w}, b}(\mathbf{x}^n)} - f_{\mathbf{w}, b}(\mathbf{x}^n) + \cancel{y^n f_{\mathbf{w}, b}(\mathbf{x}^n)}) x_i^n \\ &= -\sum_n (\hat{y}^n - f_{\mathbf{w}, b}(\mathbf{x}^n)) x_i^n \end{aligned}$$

拓展到向量形式，即有

$$-\frac{\partial \ln L(\mathbf{w}, b)}{\partial \mathbf{w}} = -\sum_n (\hat{y}^n - f_{\mathbf{w}, b}(\mathbf{x}^n)) \mathbf{x}^n$$

参数  $\mathbf{w}$  的梯度下降更新式为

$$\mathbf{w} \leftarrow \mathbf{w} + \eta \sum_n \underline{(\hat{y}^n - f_{\mathbf{w}, b}(\mathbf{x}^n)) \mathbf{x}^n}$$

差异越大，参数更新幅度越大

# Logistic回归与线性回归对比

模型

$$f_{w,b}(x) = \sigma\left(\sum_i w_i x_i + b\right)$$

输出：介于 0~1 之间

$$f_{w,b}(x) = \sum_i w_i x_i + b$$

输出：任意值

# Logistic回归与线性回归对比

模型

$$f_{w,b}(x) = \sigma\left(\sum_i w_i x_i + b\right)$$

输出：介于 0~1 之间

目标

$\hat{y}^n$ : 1、0 代表 class 1、2

$$L(f) = \sum_n l(f(x^n), \hat{y}^n)$$

$$f_{w,b}(x) = \sum_i w_i x_i + b$$

输出：任意值

训练数据:  $(x^n, \hat{y}^n)$

$\hat{y}^n$ : 真实的数值

$$L(f) = \frac{1}{2} \sum_n (f(x^n) - \hat{y}^n)^2$$

交叉熵：

$$l(f(x^n), \hat{y}^n) = -[\hat{y}^n \ln f(x^n) + (1 - \hat{y}^n) \ln(1 - f(x^n))]$$

# Logistic回归与线性回归对比

模型

$$f_{w,b}(x) = \sigma\left(\sum_i w_i x_i + b\right)$$

输出：介于 0~1 之间

目标

$\hat{y}^n$ : 1、0 代表 class 1、2

$$L(f) = \sum_n l(f(x^n), \hat{y}^n)$$

求解

Logistic回归:  $w_i \leftarrow w_i - \eta \sum_n -\underline{\hat{y}^n - f_{w,b}(x^n)} x_i^n$

线性回归:  $w_i \leftarrow w_i - \eta \sum_n -\underline{\hat{y}^n - f_{w,b}(x^n)} x_i^n$

$$f_{w,b}(x) = \sum_i w_i x_i + b$$

输出：任意值

训练数据:  $(x^n, \hat{y}^n)$

$\hat{y}^n$ : 真实的数值

$$L(f) = \frac{1}{2} \sum_n (f(x^n) - \hat{y}^n)^2$$

# 均方误差MSE准则下的Logistic回归

## 口建立模型

$$f_{\mathbf{w}, b}(\mathbf{x}) = \sigma\left(\sum_i w_i x_i + b\right)$$

## 口确定评价函数

➤ 训练数据:  $(x^n, \hat{y}^n)$ ,  $\hat{y}^n$ : **1** for class 1, **0** for class 2

$$L(f) = \frac{1}{2} \sum_n (f_{\mathbf{w}, b}(\mathbf{x}^n) - \hat{y})^2$$

## 口找到最佳的函数：梯度下降

$$\begin{aligned}\frac{\partial (f_{\mathbf{w}, b}(\mathbf{x}^n) - \hat{y})^2}{\partial w_i} &= 2(f_{\mathbf{w}, b}(\mathbf{x}^n) - \hat{y}) \frac{\partial f_{\mathbf{w}, b}(\mathbf{x}^n)}{\partial z} \frac{\partial z}{\partial w_i} \\ &= 2(f_{\mathbf{w}, b}(\mathbf{x}^n) - \hat{y}) f_{\mathbf{w}, b}(\mathbf{x}^n) (1 - f_{\mathbf{w}, b}(\mathbf{x}^n)) x_i^n\end{aligned}$$

# 均方误差MSE准则下的Logistic回归

---

□ 当  $\hat{y}^n = 1$  时

If  $f_{w,b}(x^n) = 1$  (close to target)  $\rightarrow \partial L / \partial w_i = 0$

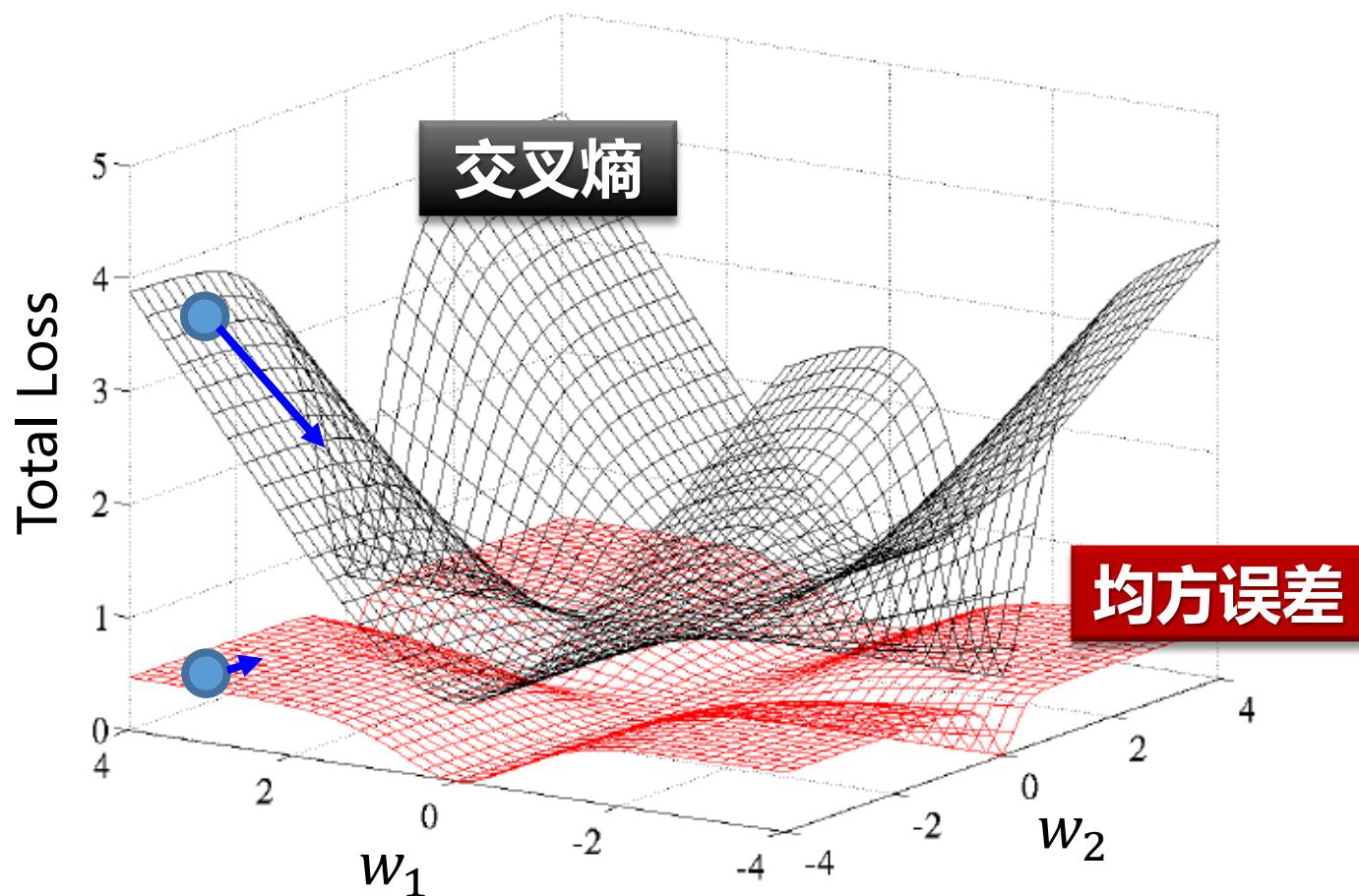
If  $f_{w,b}(x^n) = 0$  (far from target)  $\rightarrow \partial L / \partial w_i = 0$

□ 当  $\hat{y}^n = 0$  时

If  $f_{w,b}(x^n) = 1$  (far from target)  $\rightarrow \partial L / \partial w_i = 0$

If  $f_{w,b}(x^n) = 0$  (close to target)  $\rightarrow \partial L / \partial w_i = 0$

# 交叉熵与均方误差对比



<http://jmlr.org/proceedings/papers/v9/glorot10a/glorot10a.pdf>

# 概率生成模型与神经元模型对比

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

□ 概率生成模型：先找到  $\mu^1, \mu^2, \Sigma^{-1}$ ，再进行计算

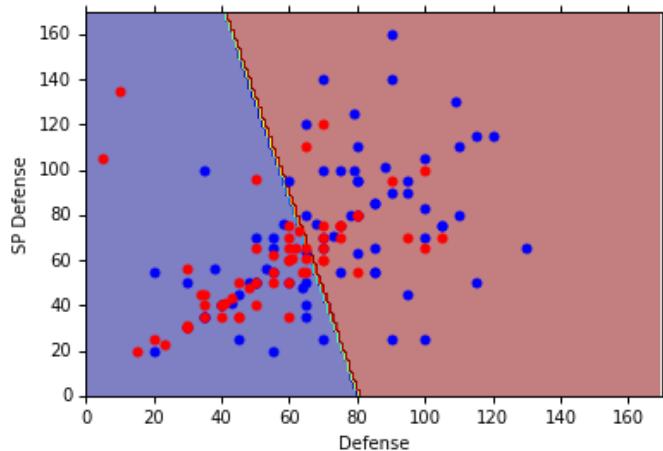
$$\mathbf{w}^T = (\boldsymbol{\mu}^1 - \boldsymbol{\mu}^2)^T (\boldsymbol{\Sigma})^{-1}$$

两种方式能解得同一组  $w^T$  和  $b$  吗

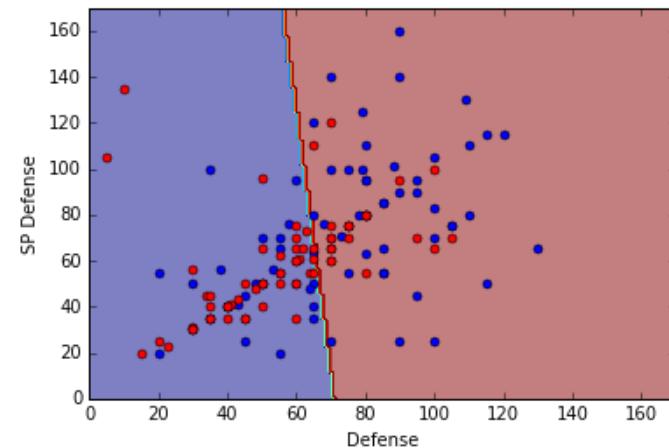
$$b = -\frac{1}{2} (\boldsymbol{\mu}^1)^T (\boldsymbol{\Sigma})^{-1} \boldsymbol{\mu}^1 + \frac{1}{2} (\boldsymbol{\mu}^2)^T (\boldsymbol{\Sigma})^{-1} \boldsymbol{\mu}^2 + \ln \frac{N_1}{N_2}$$

□ 神经元模型：直接找到  $w$  与  $b$

同样的模型，同样的数据也可能优化得到不同函数



概率生成模型 (73%准确度)



神经元模型 (79%准确度)

# 为什么神经元模型结果更好?

## 口朴素贝叶斯分类器

- 假设所有维度都是独立的

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

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

$P(x_k|C_1)$ 服从一维高斯分布

小贴士

对于二进制特征，可以假设它们来自伯努利分布

# 从朴素贝叶斯分类角度来理解

训练数据



Class 1



Class 2



Class 2



Class 2

测试数据



Class 1?  
Class 2?

**朴素贝叶斯的结果是？**

$$P(x|C_i) = P(x_1|C_i)P(x_2|C_i)$$

$$P(C_1) = \frac{1}{13}$$

$$P(x_1 = 1|C_1) = 1$$

$$P(x_2 = 1|C_1) = 1$$

$$P(C_2) = \frac{12}{13}$$

$$P(x_1 = 1|C_2) = \frac{1}{3}$$

$$P(x_2 = 1|C_2) = \frac{1}{3}$$

# 朴素贝叶斯分类

测试数据



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

$\begin{matrix} & 1 \times 1 & & \frac{1}{13} \\ & \swarrow & \downarrow & \downarrow \\ 1 \times 1 & \frac{1}{13} & \frac{1}{3} \times \frac{1}{3} & \frac{12}{13} \end{matrix}$

$$P(C_1) = \frac{1}{13}$$

$$P(x_1 = 1|C_1) = 1$$

$$P(x_2 = 1|C_1) = 1$$

$$P(C_2) = \frac{12}{13}$$

$$P(x_1 = 1|C_2) = \frac{1}{3}$$

$$P(x_2 = 1|C_2) = \frac{1}{3}$$

朴素贝叶斯假设各特征相互独立，即

$P(x|C_2) = \frac{1}{3} \times \frac{1}{3}$ ，而实际上  $P(x|C_2) = 0$

# 概率生成模型与神经元模型对比

---

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

□ 概率生成模型：先找到  $\mu^1, \mu^2, \Sigma^{-1}$ ，再进行计算

$$\mathbf{w}^T = (\boldsymbol{\mu}^1 - \boldsymbol{\mu}^2)^T (\boldsymbol{\Sigma})^{-1}$$

$$b = -\frac{1}{2} (\boldsymbol{\mu}^1)^T (\boldsymbol{\Sigma})^{-1} \boldsymbol{\mu}^1 + \frac{1}{2} (\boldsymbol{\mu}^2)^T (\boldsymbol{\Sigma})^{-1} \boldsymbol{\mu}^2 + \ln \frac{N_1}{N_2}$$

□ 神经元模型：直接找到  $w$  与  $b$

□ 性能对比

➤ 通常**神经元模型效果更好**

➤ **概率生成模型**的优势

■ 在概率分布的假设下只需要更少的数据、面对噪声更鲁棒

■ 先验和类相关概率可以从不同的来源估计

# 内容导览

---



线性回归：多项式拟合



线性分类：概率生成模型



Logistic回归：神经元模型



多分类问题



线性模型的局限与非线性模型

# 多分类问题：以三分类为例

$$C_1: w^1, b_1 \quad z_1 = w^1 \cdot x + b_1$$

$$C_2: w^2, b_2 \quad z_2 = w^2 \cdot x + b_2$$

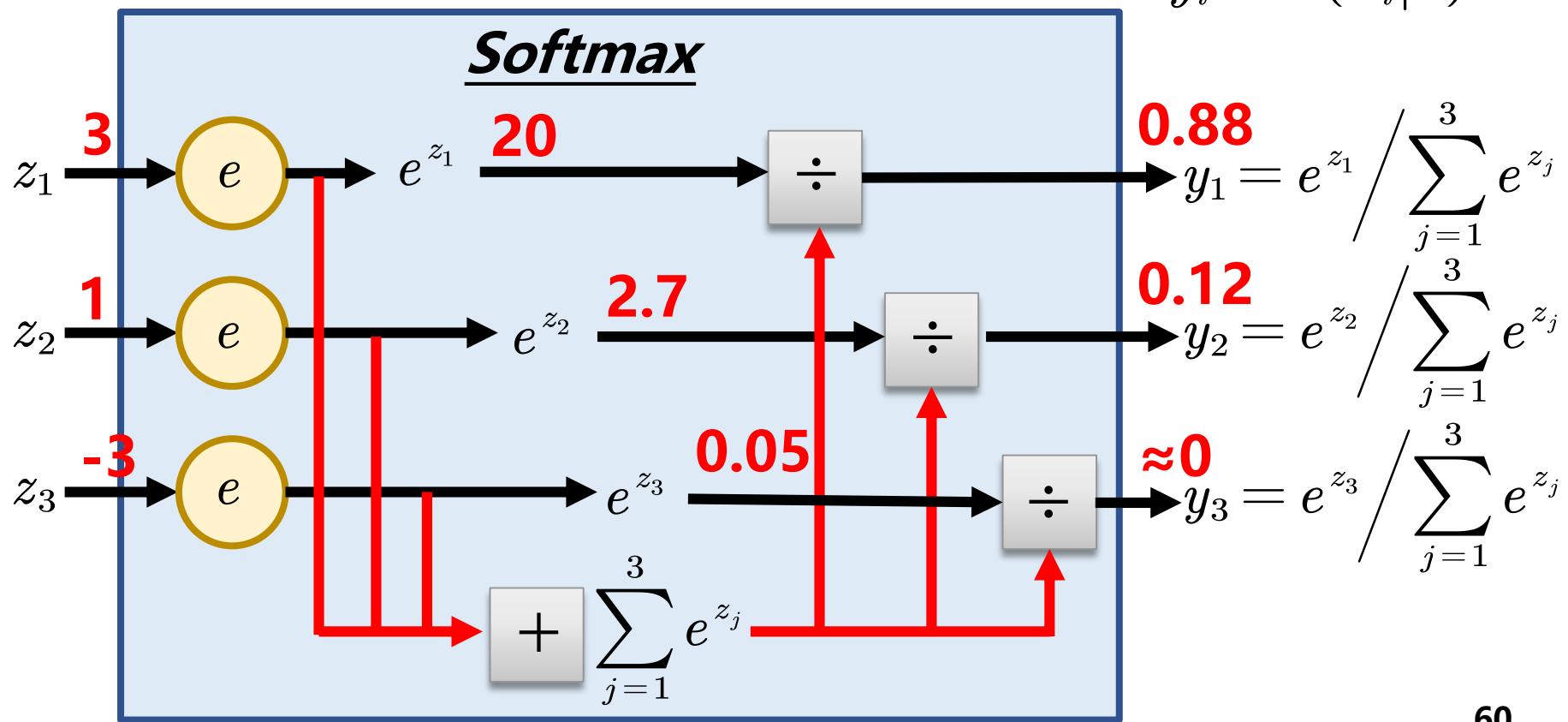
$$C_3: w^3, b_3 \quad z_3 = w^3 \cdot x + b_3$$

概率

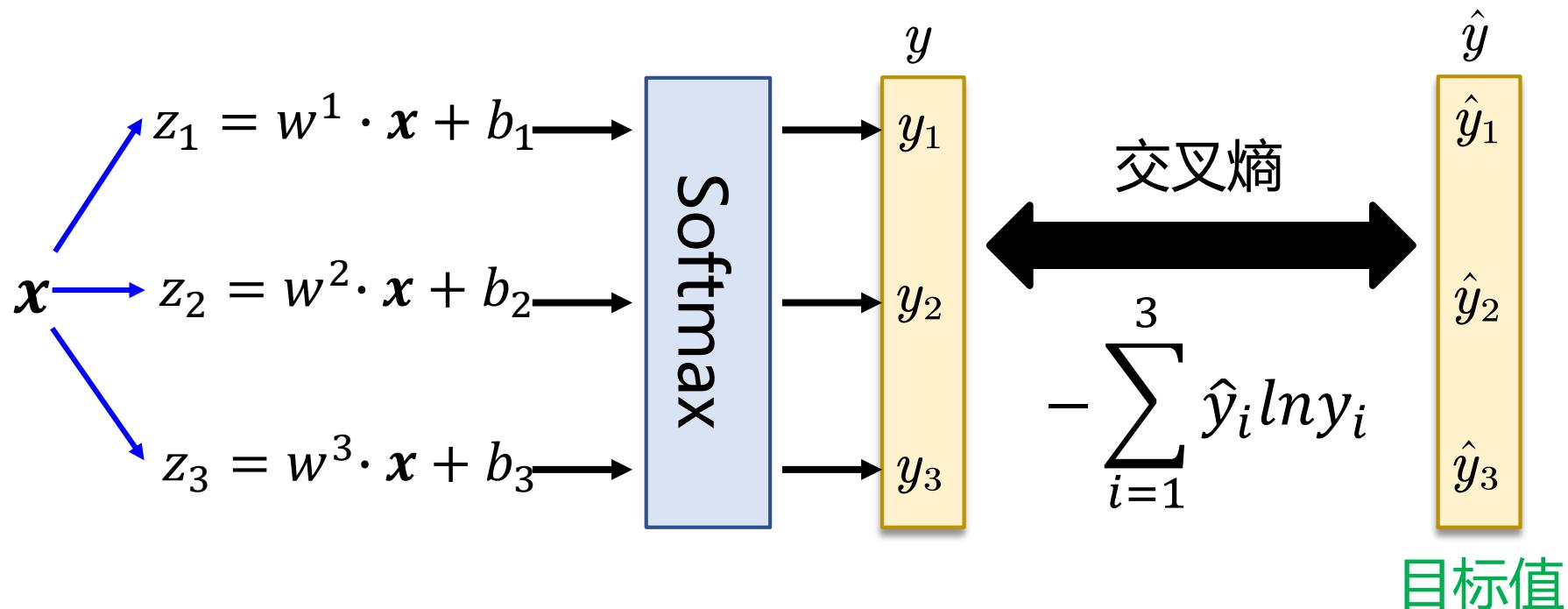
■  $1 > y_i > 0$

■  $\sum_i y_i = 1$

$$y_i = P(C_i | \mathbf{x})$$



# 多分类问题：以三分类为例



If  $x \in \text{class 1}$

$$\hat{y} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

$$-\ln y_1$$

If  $x \in \text{class 2}$

$$\hat{y} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

$$-\ln y_2$$

If  $x \in \text{class 3}$

$$\hat{y} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

$$-\ln y_3$$

# 内容导览

---



线性回归：多项式拟合



线性分类：概率生成模型



Logistic回归：神经元模型

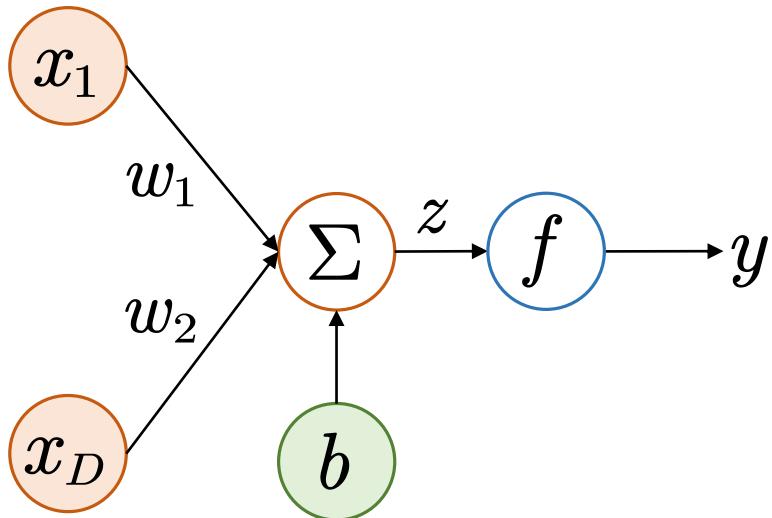


多分类问题



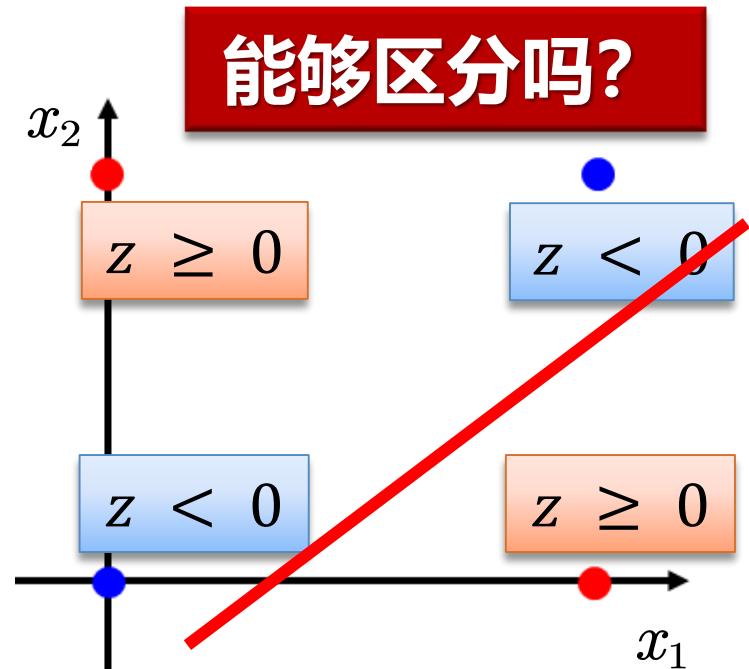
线性模型的局限与非线性模型

# Logistic回归的局限



$$\begin{cases} \text{Class 1} & y \geq 0.5 \quad (z \geq 0) \\ \text{Class 2} & y < 0.5 \quad (z < 0) \end{cases}$$

输入特征		标签
$x_1$	$x_2$	
0	0	Class 2
0	1	Class 1
1	0	Class 1
1	1	Class 2



# Logistic回归的局限

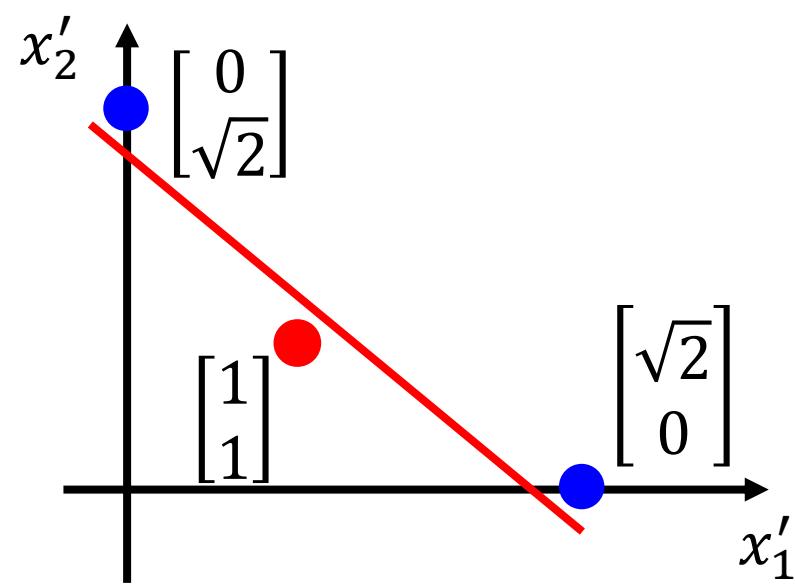
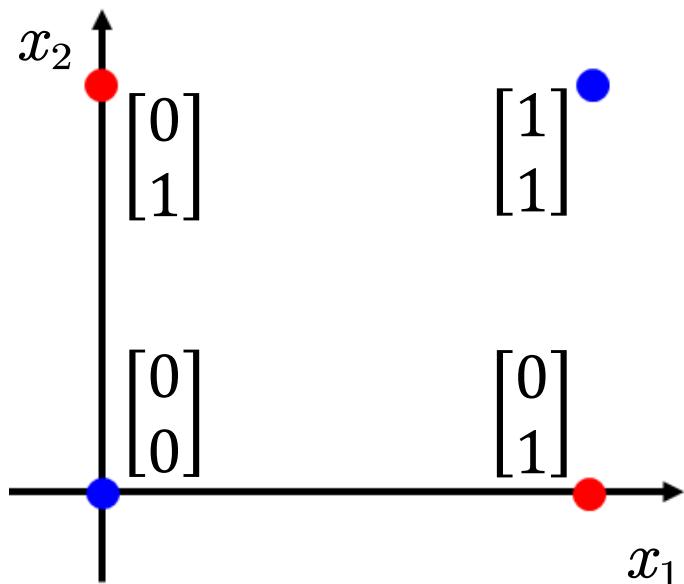
## 口 特征变换

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \longrightarrow \begin{bmatrix} x'_1 \\ x'_2 \end{bmatrix}$$

$x'_1$ : 到  $\begin{bmatrix} 0 \\ 0 \end{bmatrix}$  的距离

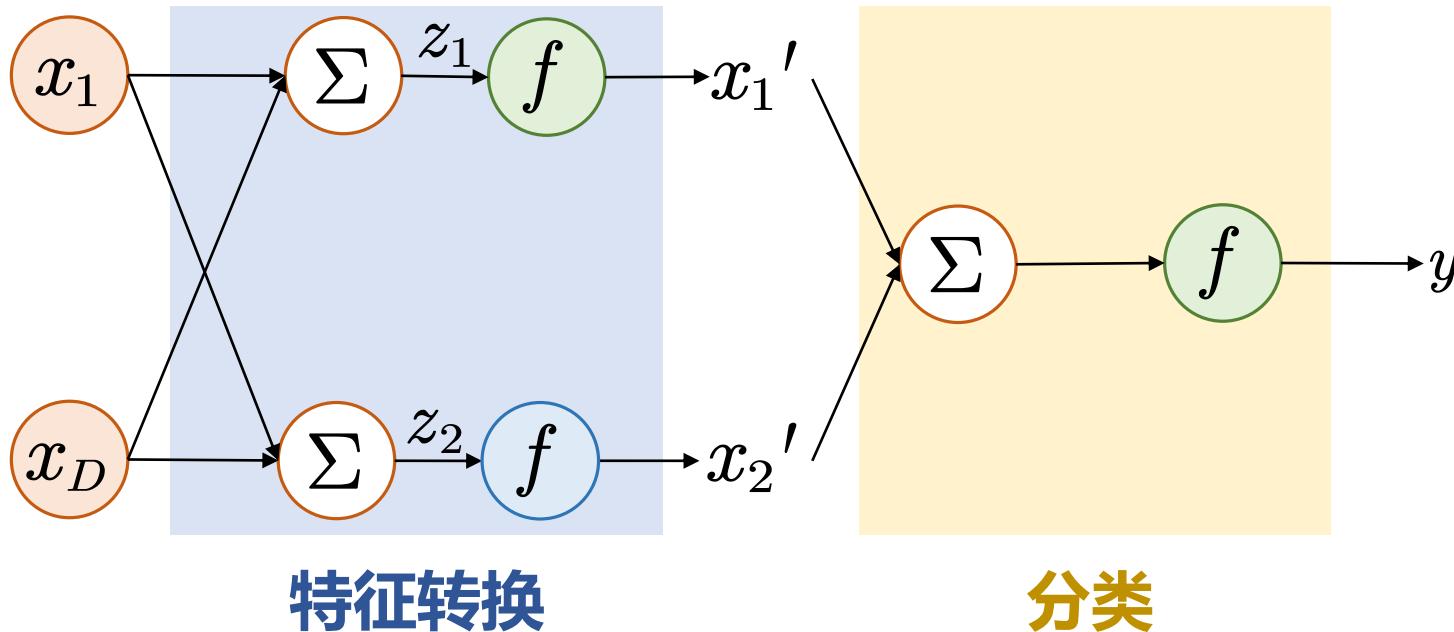
$x'_2$ : 到  $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$  的距离

并不总是有效

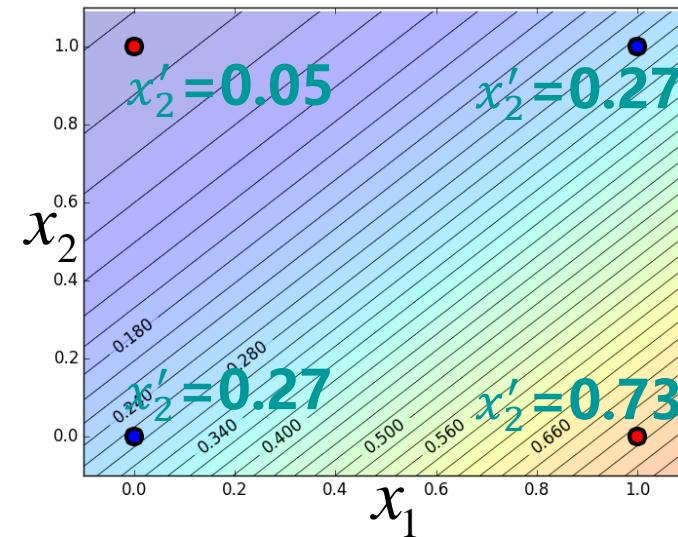
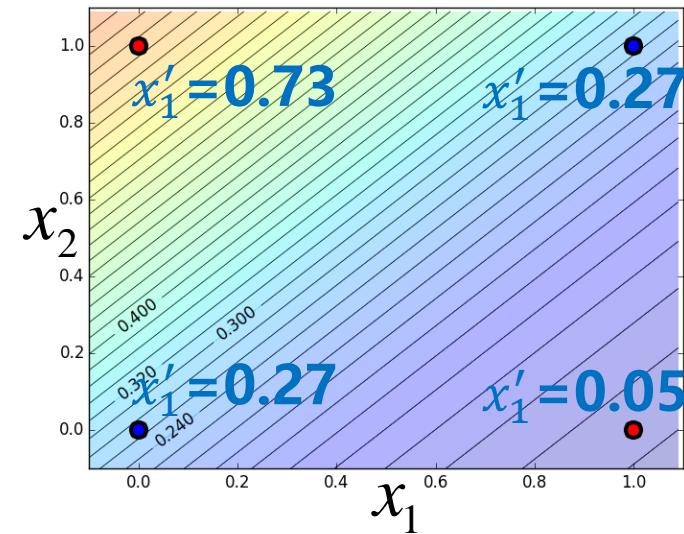
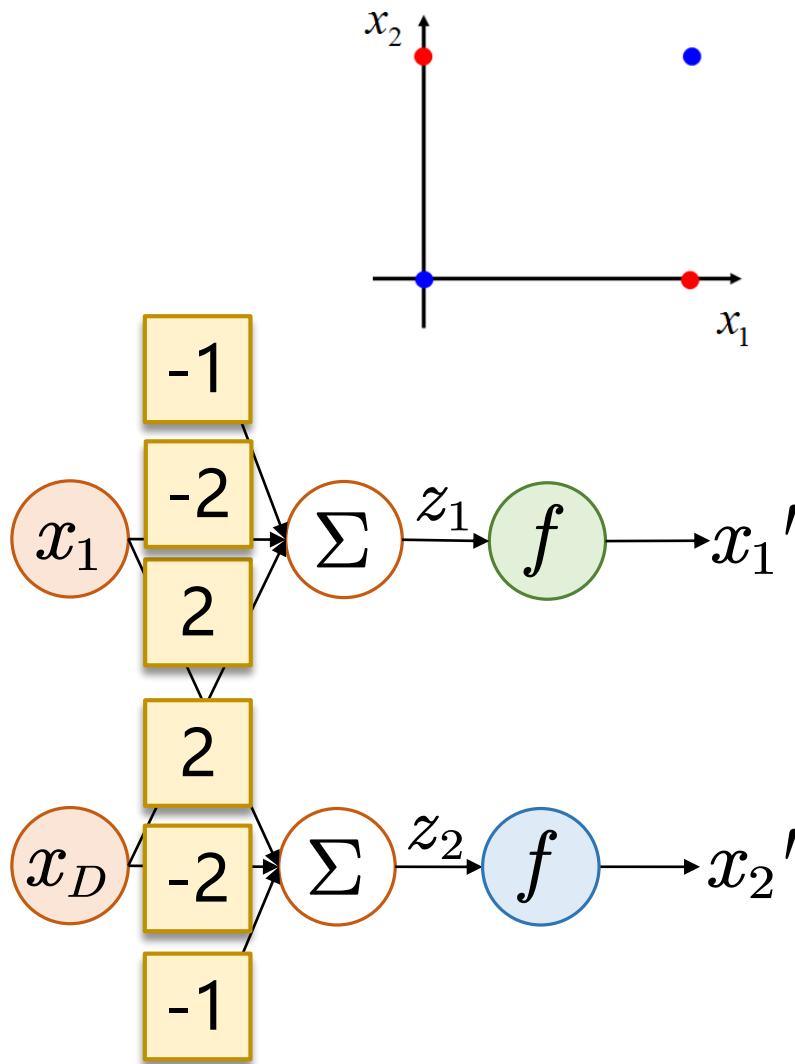


# Logistic回归的局限

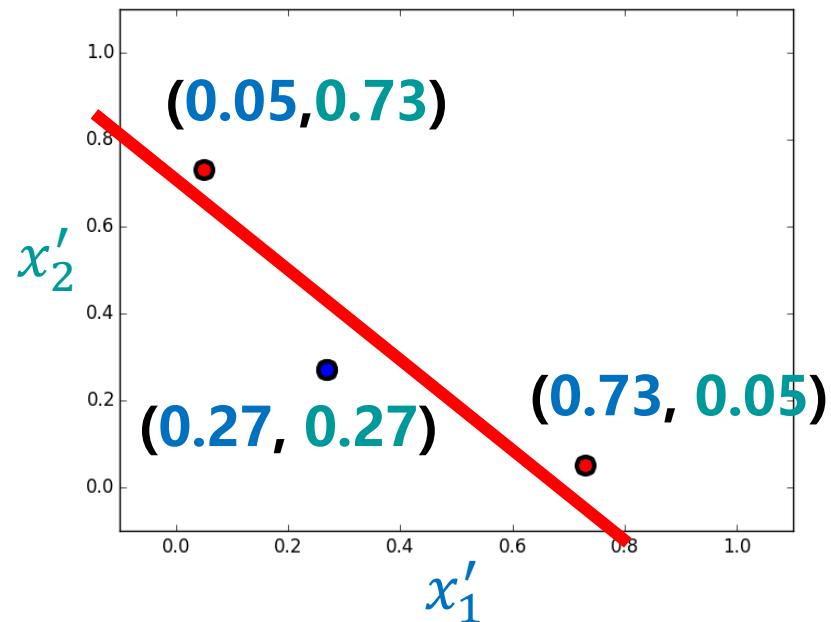
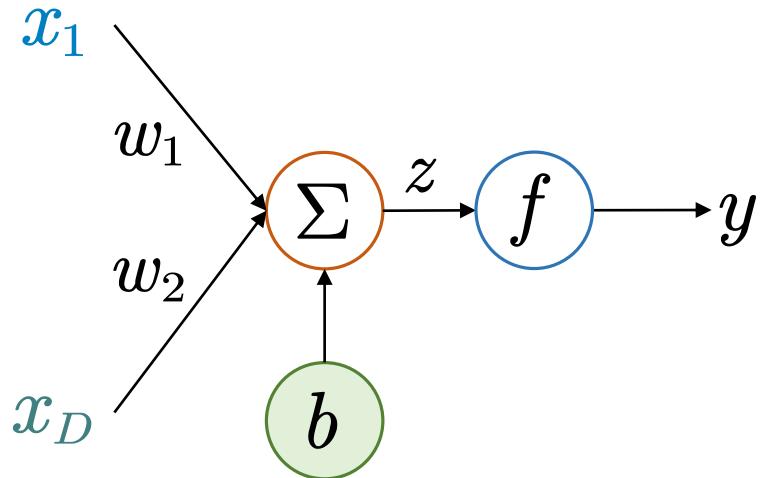
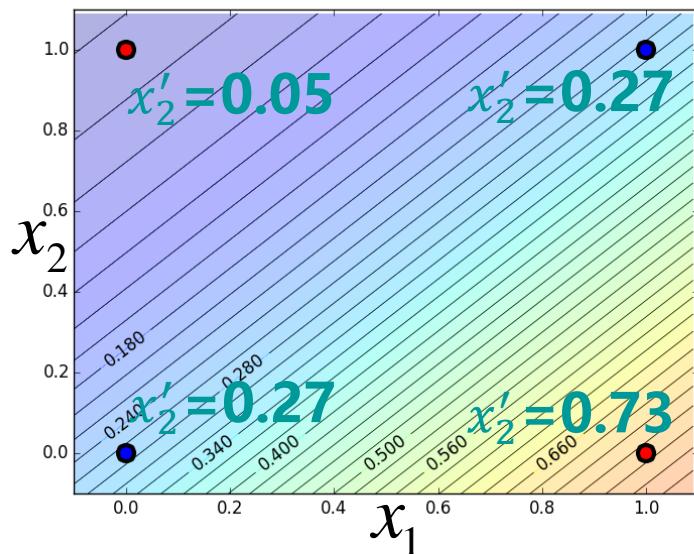
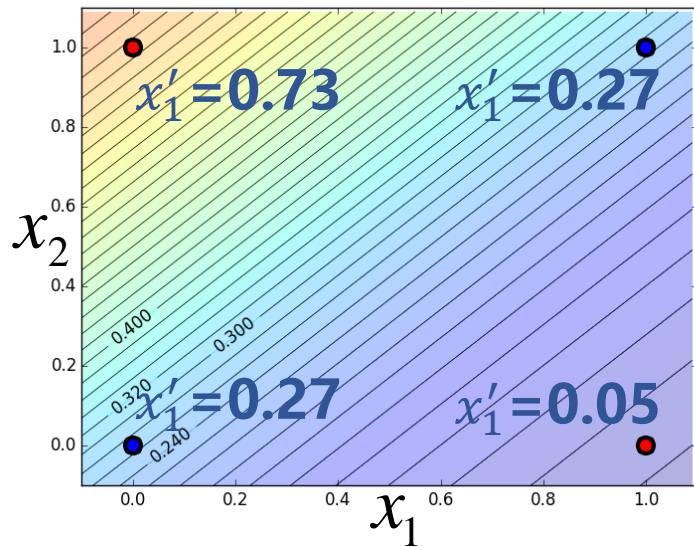
□ 连接多个Logistic回归模型 (图中忽略bias项)



# 连接多个Logistic回归模型

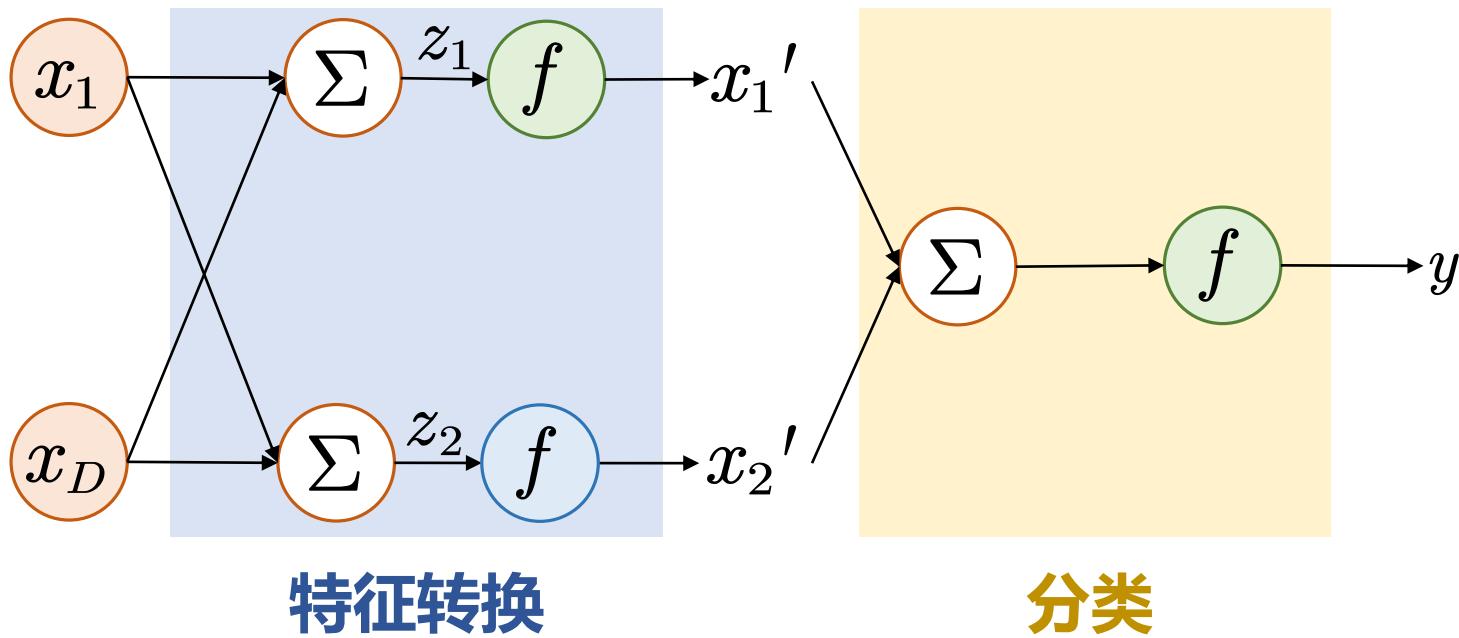
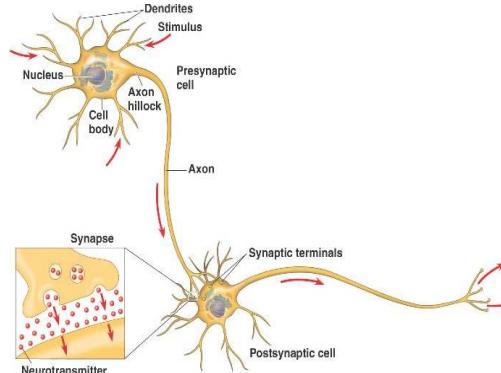


# 连接多个Logistic回归模型



# Logistic回归：深度学习方案

所有Logistic回归模型的参数一起联合学习



神经网络