

模式识别与机器学习

**Pattern Recognition
and Machine Learning**

课程内容

■ 模式识别与机器学习概述

■ 模式识别与机器学习的基本方法

- 回归分析、线性判别函数、线性神经网络、核方法和支持向量机、决策树分类、逻辑斯特回归
- 贝叶斯统计决策理论、概率密度函数估计
- 无监督学习和聚类
- 特征选择与提取

逻辑斯特回归

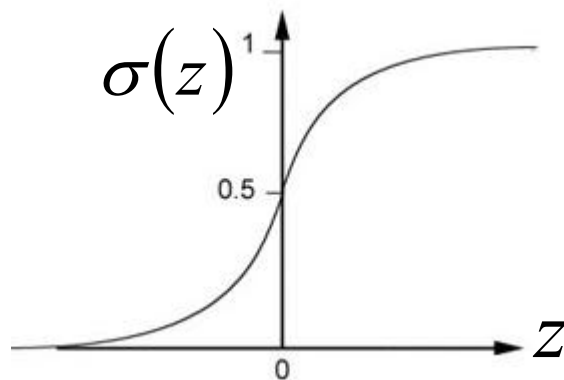
Posterior Probability(后验概率)

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

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

Sigmoid function

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



逻辑斯特回归

Posterior Probability(后验概率)

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

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

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

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

逻辑斯特回归

Posterior Probability(后验概率)

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

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

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

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

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

逻辑斯特回归

Posterior Probability(后验概率)

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

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

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

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

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

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

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

逻辑斯特回归

Posterior Probability(后验概率) $P(C_1|x) = \sigma(z)$

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

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

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

$P(C_1|x) = \sigma(\mathbf{w} \cdot x + b)$ How about directly find \mathbf{w} and b ?

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

逻辑斯特回归

Step 1: Function Set

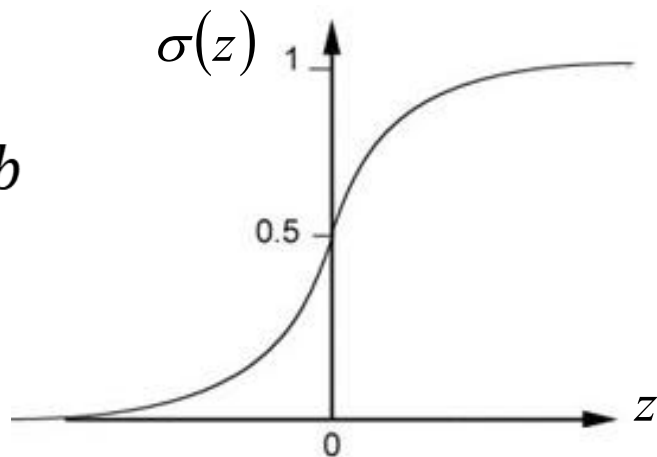
Function set: Including all different w and b

$$\left\{ \begin{array}{ll} z \geq 0 & \text{class 1} \\ z < 0 & \text{class 2} \end{array} \right.$$

$$P_{w,b}(C_1|x) = \sigma(z)$$

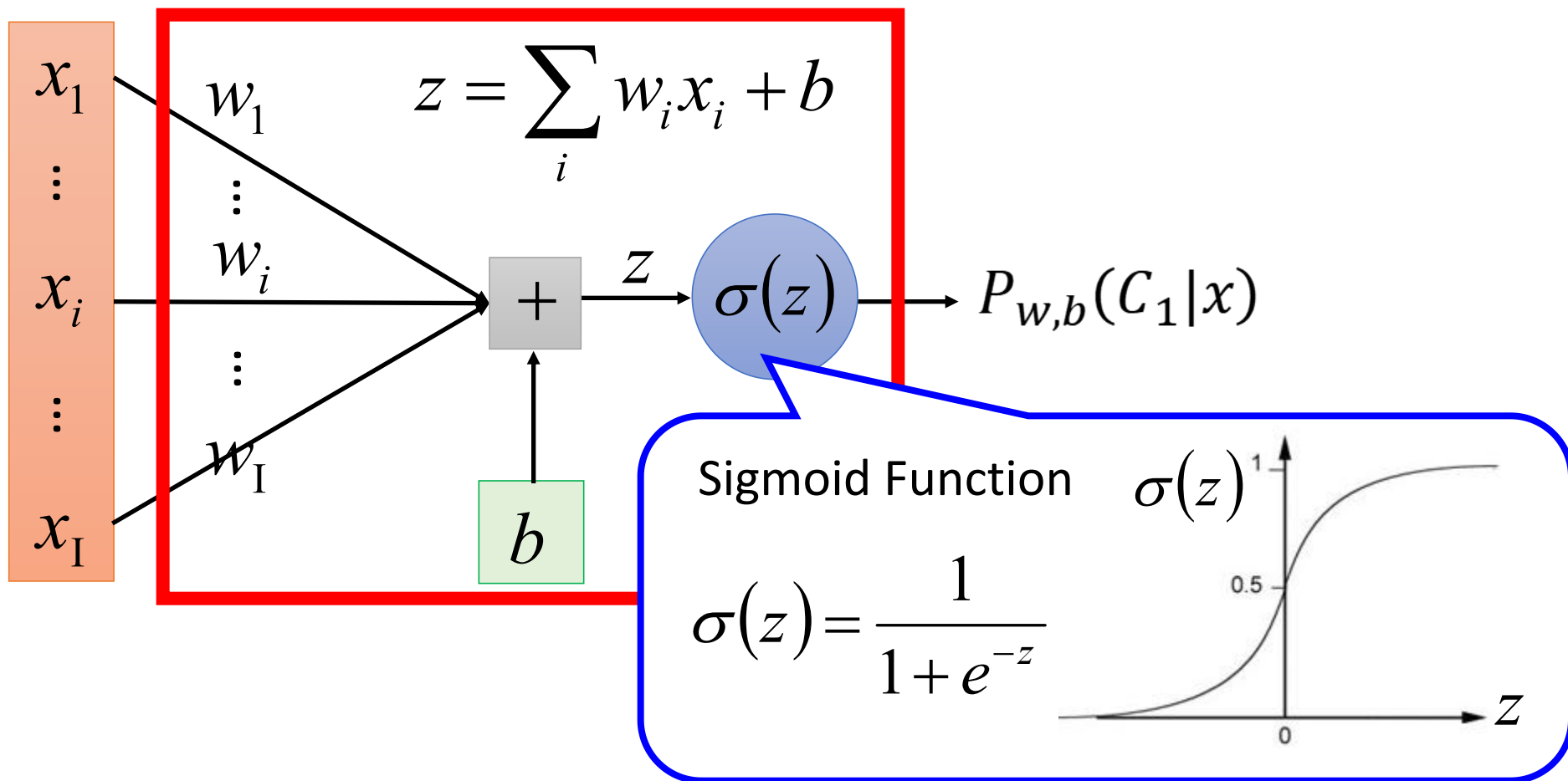
$$z = w \cdot x + b = \sum_i w_i x_i + b$$

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



逻辑斯特回归

Step 1: Function Set



逻辑斯特回归

Step 2: Goodness of a Function

Training Data	x^1	x^2	x^3	x^N
	C_1	C_1	C_2		C_1

Assume the data is generated based on $f_{w,b}(x) = P_{w,b}(C_1|x)$

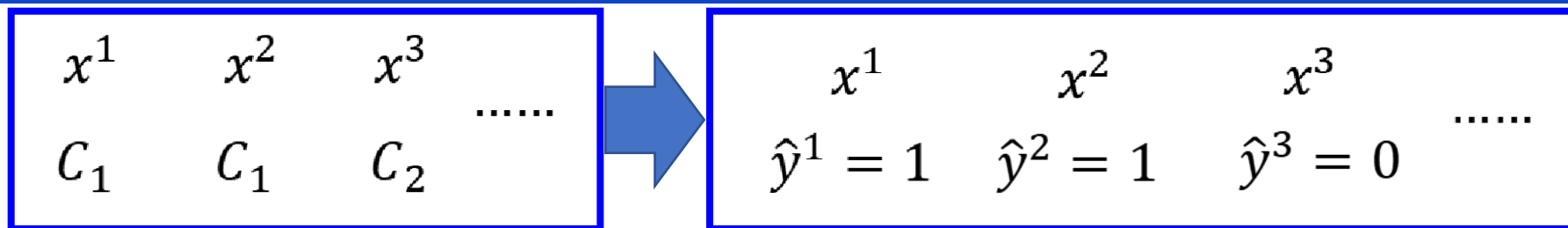
Given a set of w and b , what is its probability of generating the data?

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

The most likely w^* and b^* is the one with the largest $L(w, b)$.

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

逻辑斯特回归



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

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

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

$$= -\ln f_{w,b}(x^1) \rightarrow -[1^1 \ln f(x^1) + \cancel{0 \ln(1 - f(x^1))}]$$

$$-\ln f_{w,b}(x^2) \rightarrow -[1^2 \ln f(x^2) + \cancel{0 \ln(1 - f(x^2))}]$$

$$-\ln(1 - f_{w,b}(x^3)) \rightarrow -[\cancel{0 \ln f(x^3)} + 1 \ln(1 - f(x^3))]$$

⋮

逻辑斯特回归

Step 2: Goodness of a Function

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

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

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

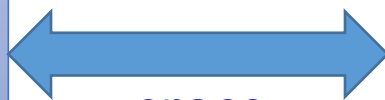
$$= \sum_n - [\hat{y}^n \ln f_{w,b}(x^n) + (1 - \hat{y}^n) \ln(1 - f_{w,b}(x^n))]$$

Cross entropy between two Bernoulli distribution

Distribution p:

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

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



cross
entropy

Distribution q:

$$q(x = 1) = f(x^n)$$

$$q(x = 0) = 1 - f(x^n)$$

$$H(p, q) = - \sum_x p(x) \ln(q(x))$$

逻辑斯特回归

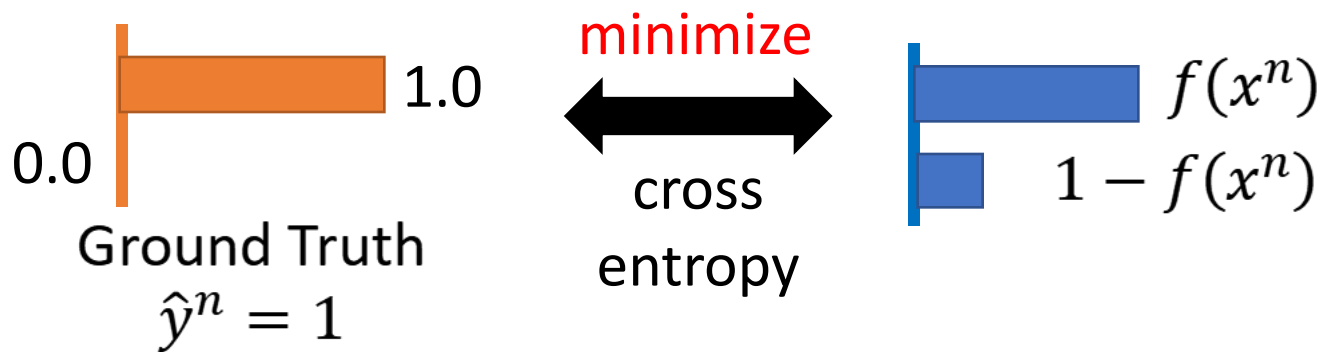
Step 2: Goodness of a Function

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

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

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

$$= \sum_n \underbrace{-[\hat{y}^n \ln f_{w,b}(x^n) + (1 - \hat{y}^n) \ln(1 - f_{w,b}(x^n))]}_{\text{Cross entropy between two Bernoulli distribution}}$$



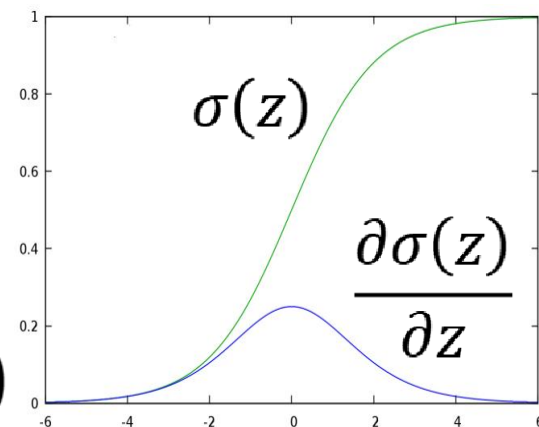
逻辑斯特回归

Step 3: Find the best function

$$\frac{-\ln L(w, b)}{\partial w_i} = \sum_n \frac{(1 - f_{w,b}(x^n))x_i^n}{\partial w_i} - [\hat{y}^n \frac{\ln f_{w,b}(x^n)}{\partial w_i} + (1 - \hat{y}^n) \frac{\ln(1 - f_{w,b}(x^n))}{\partial w_i}]$$

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

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



$$\begin{aligned} f_{w,b}(x) &= \sigma(z) \\ &= 1 / (1 + \exp(-z)) \end{aligned}$$

$$z = w \cdot x + b = \sum_i w_i x_i + b$$

逻辑斯特回归

Step 3: Find the best function

$$\frac{-\ln L(w, b)}{\partial w_i} = \sum_n \frac{(1 - f_{w,b}(x^n))x_i^n}{\partial w_i} - \frac{f_{w,b}(x^n)x_i^n}{\partial w_i}$$
$$= \sum_n - [\hat{y}^n \ln f_{w,b}(x^n) + (1 - \hat{y}^n) \ln(1 - f_{w,b}(x^n))]$$

$$\frac{\partial \ln(1 - f_{w,b}(x))}{\partial w_i} = \frac{\partial \ln(1 - f_{w,b}(x))}{\partial z} \frac{\partial z}{\partial w_i} \quad \frac{\partial z}{\partial w_i} = x_i$$

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

$$f_{w,b}(x) = \sigma(z)$$
$$= 1 / (1 + \exp(-z))$$

$$z = w \cdot x + b = \sum_i w_i x_i + b$$

逻辑斯特回归

Step 3: Find the best function

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

Larger difference, larger update

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

逻辑斯特回归

Logistic Regression

Step 1: $f_{w,b}(x) = \sigma \left(\sum_i w_i x_i + b \right)$

Output: between 0 and 1

Training data: (x^n, \hat{y}^n)

Step 2: \hat{y}^n : 1 for class 1, 0 for class 2

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

Linear Regression

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

Output: any value

Training data: (x^n, \hat{y}^n)

\hat{y}^n : a real number

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

Cross entropy:

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

逻辑斯特回归

Logistic Regression

Step 1: $f_{w,b}(x) = \sigma \left(\sum_i w_i x_i + b \right)$

Output: between 0 and 1

Training data: (x^n, \hat{y}^n)

Step 2: \hat{y}^n : 1 for class 1, 0 for class 2

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

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

Step 3:

Linear regression: $w_i \leftarrow w_i - \eta \sum_n \underbrace{-(\hat{y}^n - f_{w,b}(x^n))x_i^n}$

Linear Regression

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

Output: any value

Training data: (x^n, \hat{y}^n)

\hat{y}^n : a real number

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

逻辑斯特回归

Logistic Regression + Square Error

Step 1: $f_{w,b}(x) = \sigma \left(\sum_i w_i x_i + b \right)$

Step 2: Training data: (x^n, \hat{y}^n) , \hat{y}^n : 1 for class 1, 0 for class 2

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

Step 3:

$$\frac{\partial (f_{w,b}(x) - \hat{y})^2}{\partial w_i} = 2(f_{w,b}(x) - \hat{y}) \frac{\partial f_{w,b}(x)}{\partial z} \frac{\partial z}{\partial w_i}$$
$$= 2(f_{w,b}(x) - \hat{y}) f_{w,b}(x) (1 - f_{w,b}(x)) x_i$$

$\hat{y}^n = 1$ 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) $\Rightarrow \partial L / \partial w_i = 0$

逻辑斯特回归

Logistic Regression + Square Error

Step 1: $f_{w,b}(x) = \sigma \left(\sum_i w_i x_i + b \right)$

Step 2: Training data: (x^n, \hat{y}^n) , \hat{y}^n : 1 for class 1, 0 for class 2

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

Step 3:

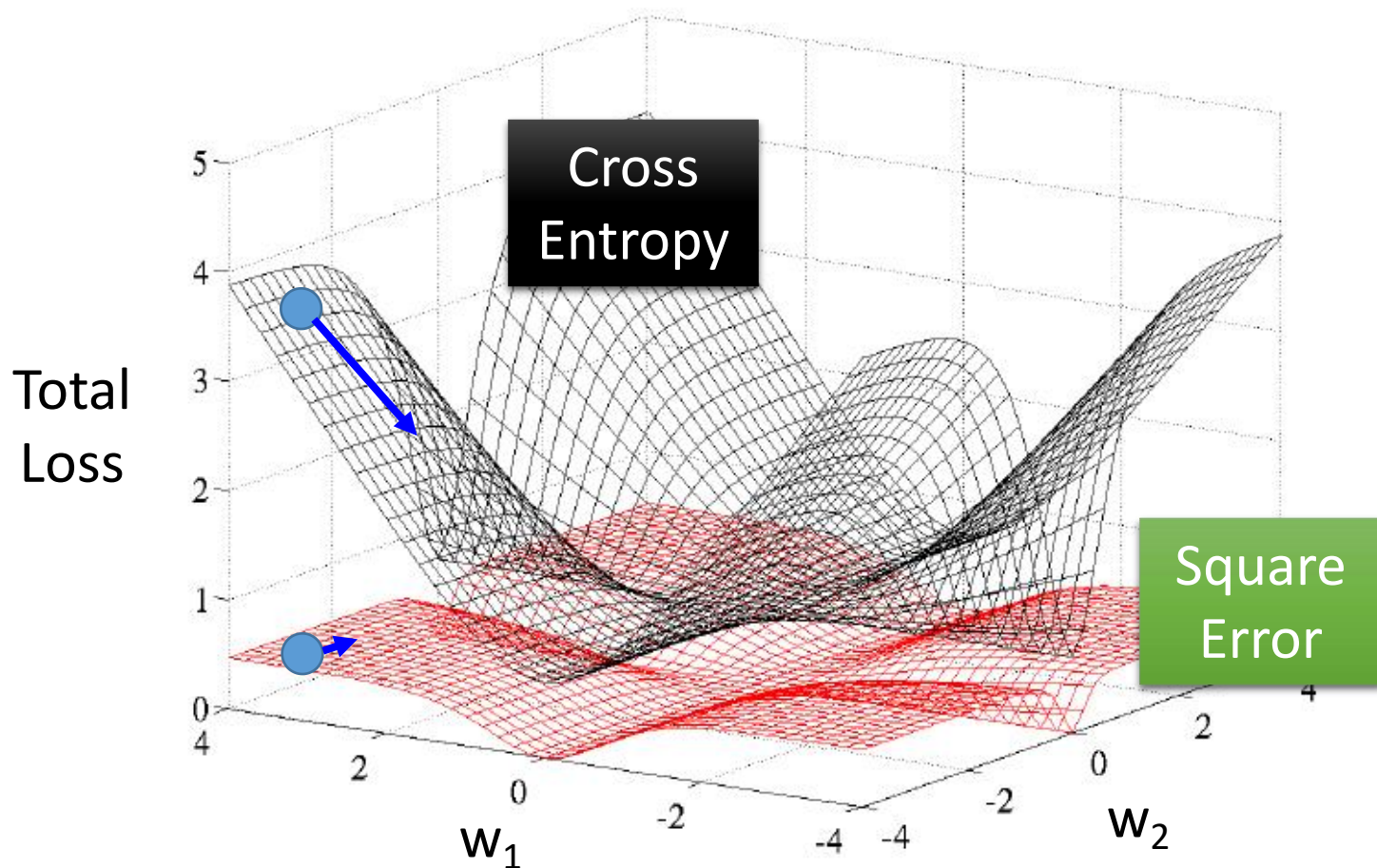
$$\frac{\partial (f_{w,b}(x) - \hat{y})^2}{\partial w_i} = 2(f_{w,b}(x) - \hat{y}) \frac{\partial f_{w,b}(x)}{\partial z} \frac{\partial z}{\partial w_i}$$
$$= 2(f_{w,b}(x) - \hat{y}) f_{w,b}(x) (1 - f_{w,b}(x)) x_i$$

$\hat{y}^n = 1$ 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) $\Rightarrow \partial L / \partial w_i = 0$

逻辑斯特回归

Cross Entropy v.s. Square Error



逻辑斯特回归

Discriminative v.s. Generative

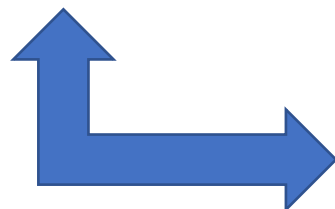
$$P(C_1|x) = \sigma(w \cdot x + b)$$



directly find w and b



Find $\mu^1, \mu^2, \Sigma^{-1}$



Will we obtain the same set of w and b ?

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

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

The same model (function set), but different function may be selected by the same training data.

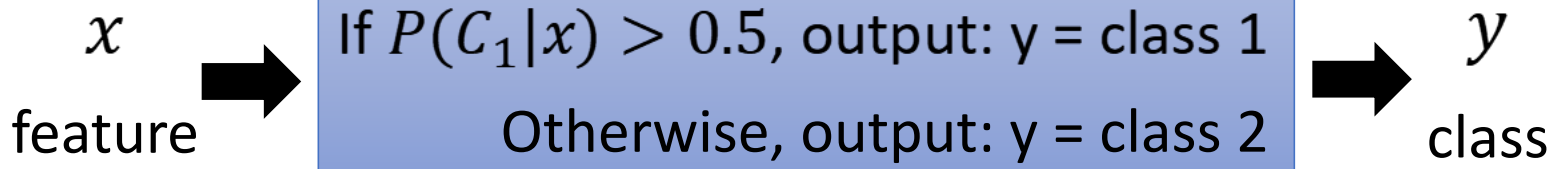
逻辑斯特回归

Three Steps

x^1	x^2	x^3	x^n
\hat{y}^1	\hat{y}^2	\hat{y}^3		\hat{y}^n

- Step 1. Function Set (Model)

$$\hat{y}^n = \text{class 1, class 2}$$



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

w and b are related to N_1, N_2, μ^1, μ^2 ,

- Step 2. Goodness of a function

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

- Step 3. Find the best function: gradient descent

逻辑斯特回归

Multi-class Classification (3 classes as example)

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

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

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

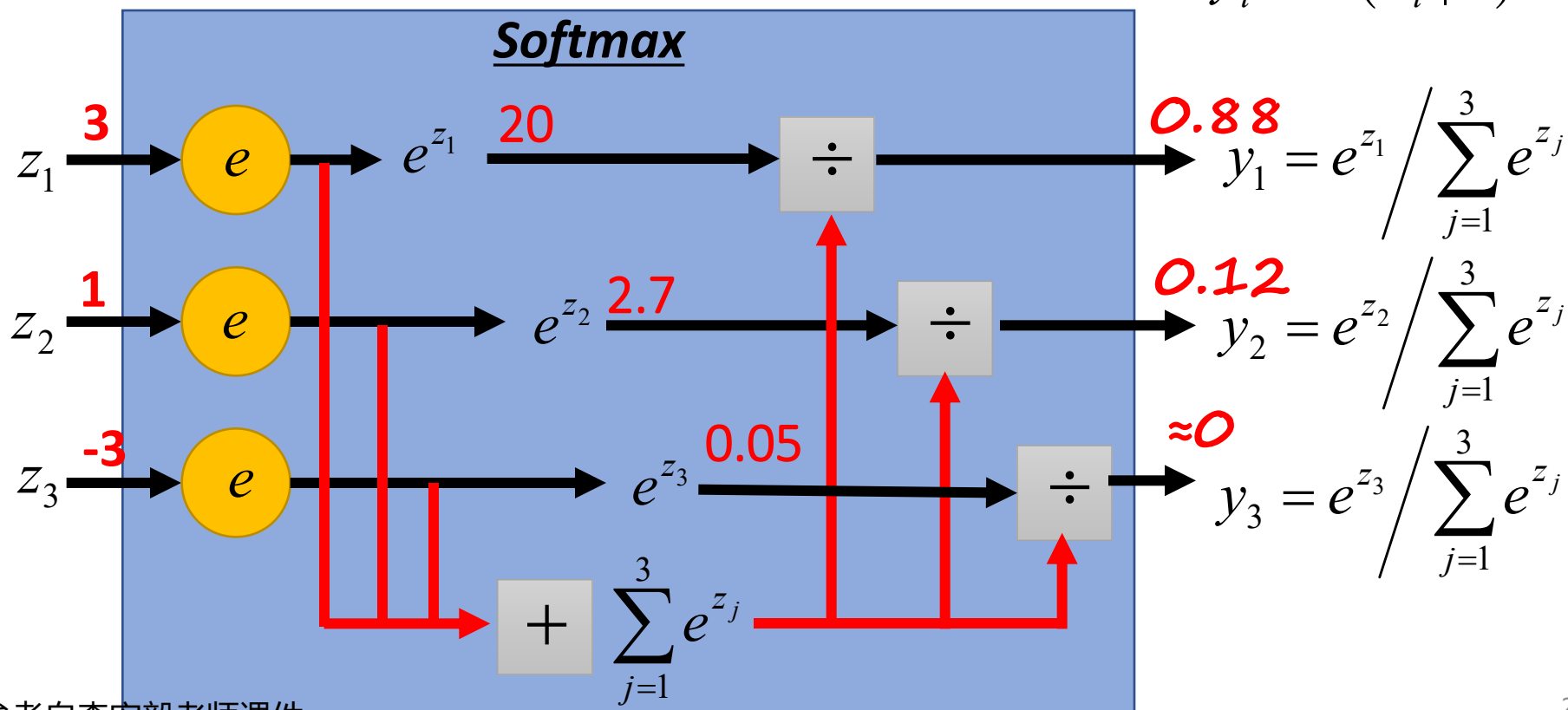
Probability:

$$\blacksquare \quad 1 > y_i > 0$$

$$\blacksquare \quad \sum y_i = 1$$

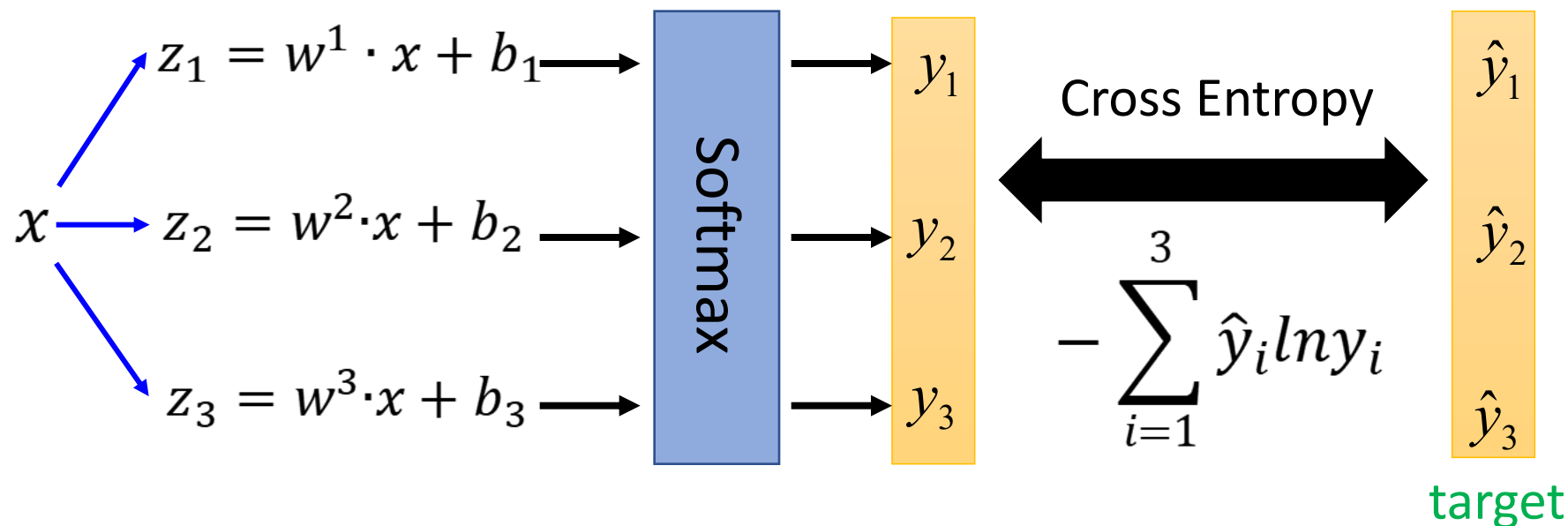
$$y_i = P(C_i | x)$$

Softmax



逻辑斯特回归

Multi-class Classification (3 classes as example)



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

$$\hat{y} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$
$$-\ln y_1$$

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

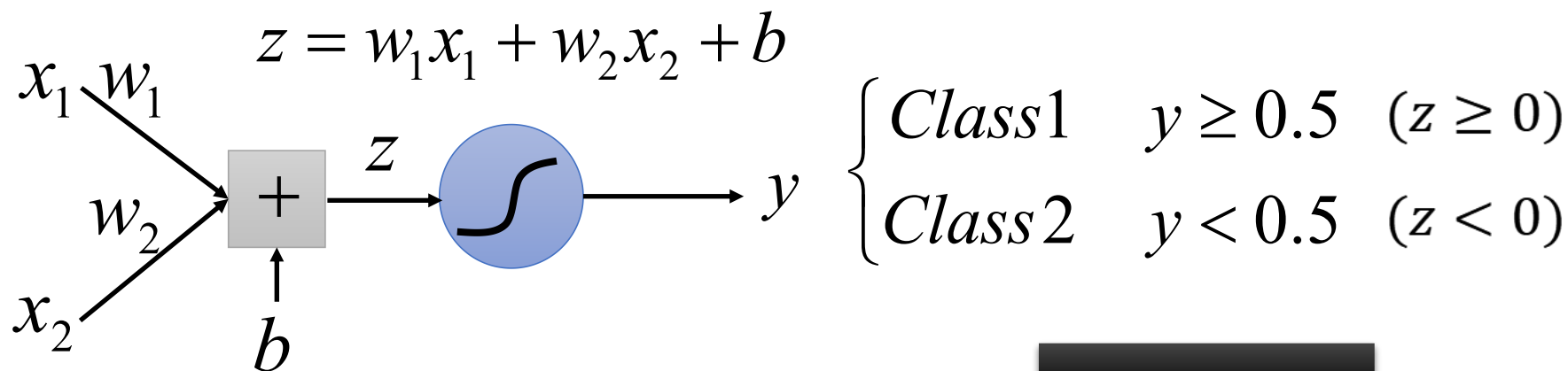
$$\hat{y} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$
$$-\ln y_2$$

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

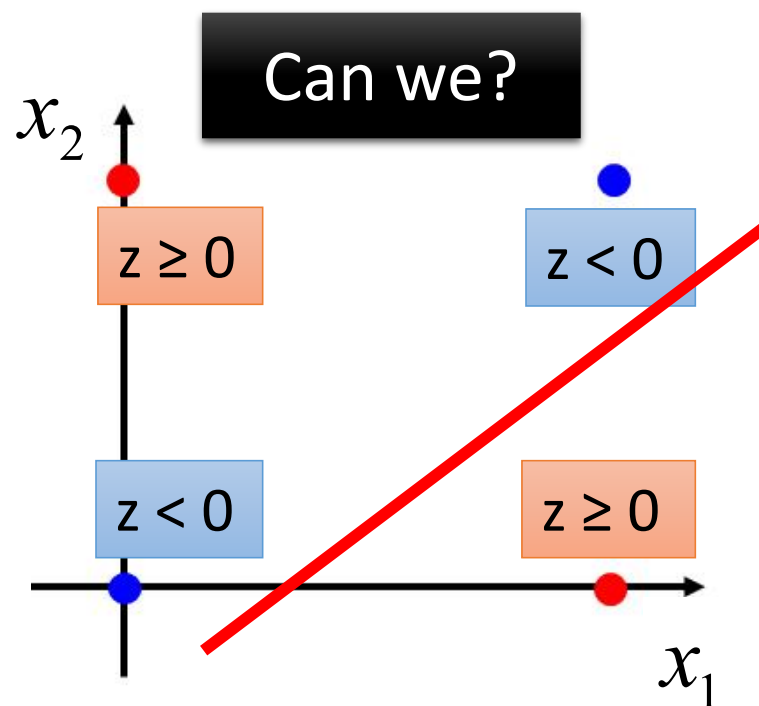
$$\hat{y} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$
$$-\ln y_3$$

逻辑斯特回归

Limitation of Logistic Regression



Input Feature		Label
x_1	x_2	
0	0	Class 2
0	1	Class 1
1	0	Class 1
1	1	Class 2

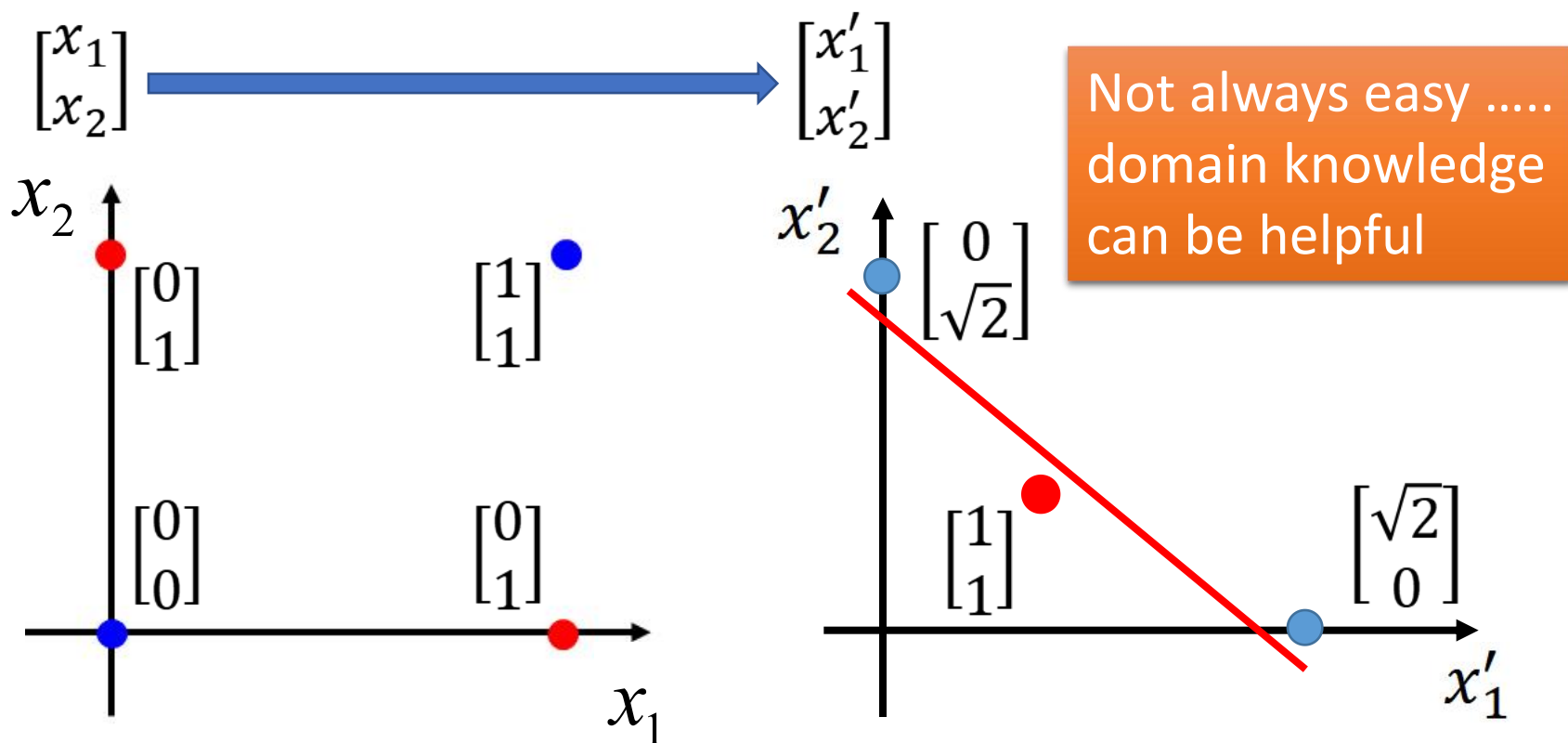


逻辑斯特回归

Limitation of Logistic Regression

- **Feature transformation**

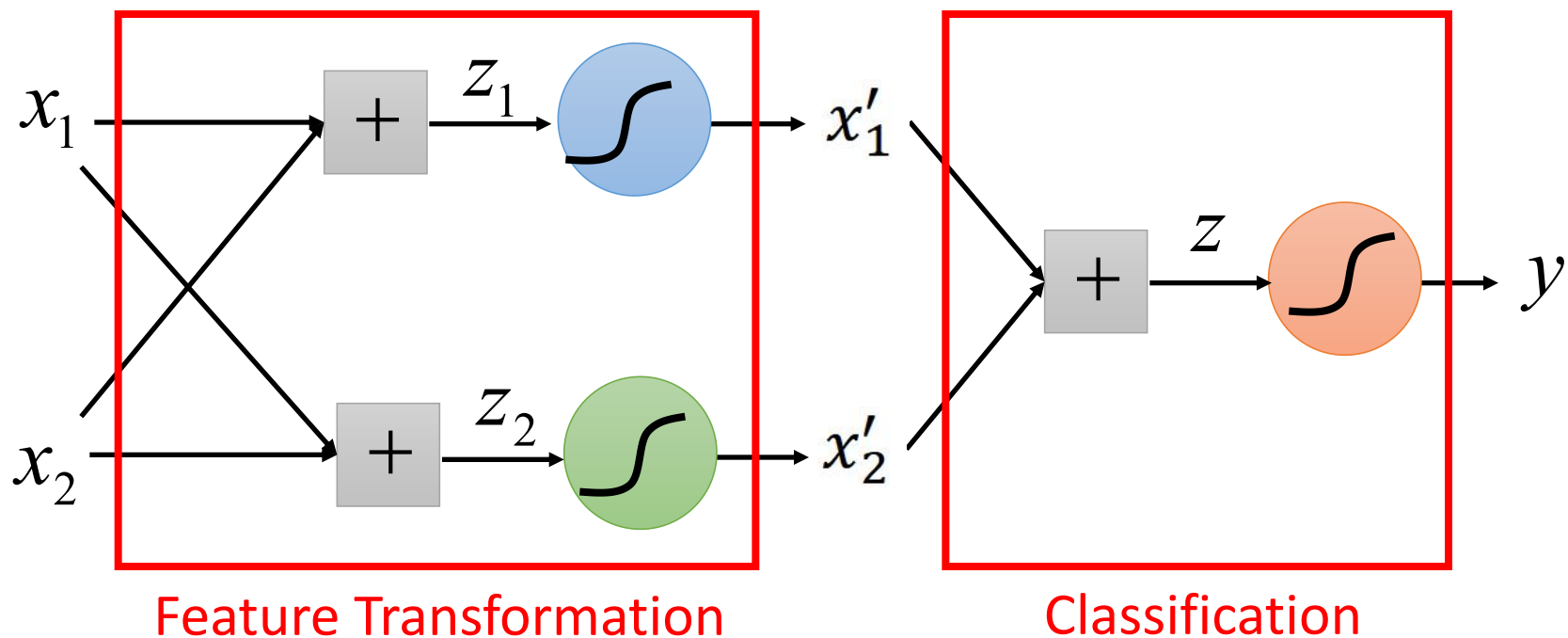
x'_1 : distance to $\begin{bmatrix} 0 \\ 0 \end{bmatrix}$
 x'_2 : distance to $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$



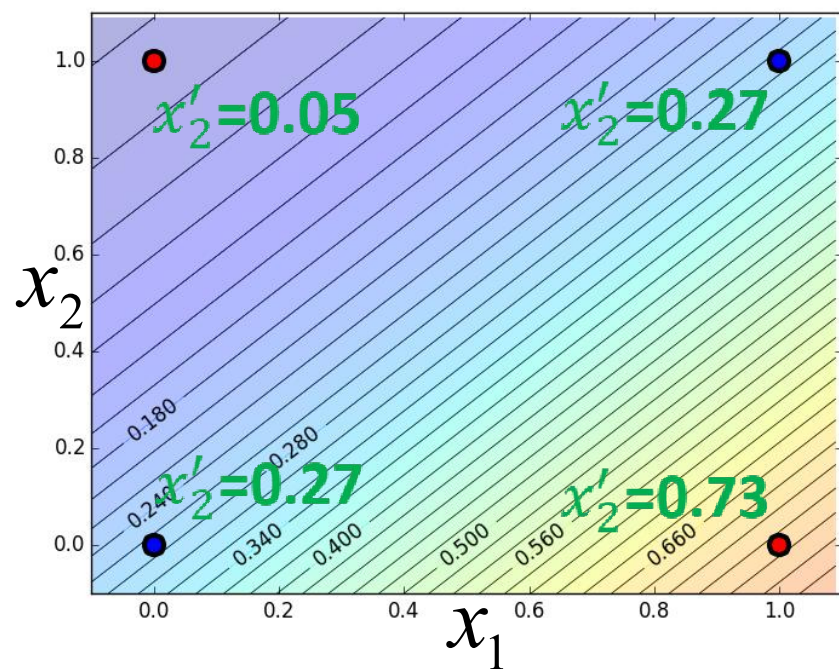
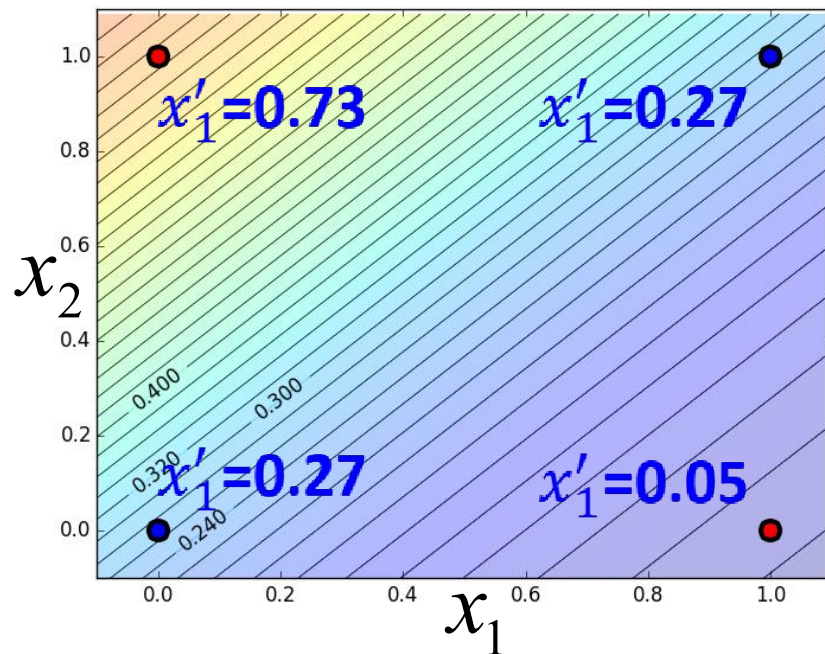
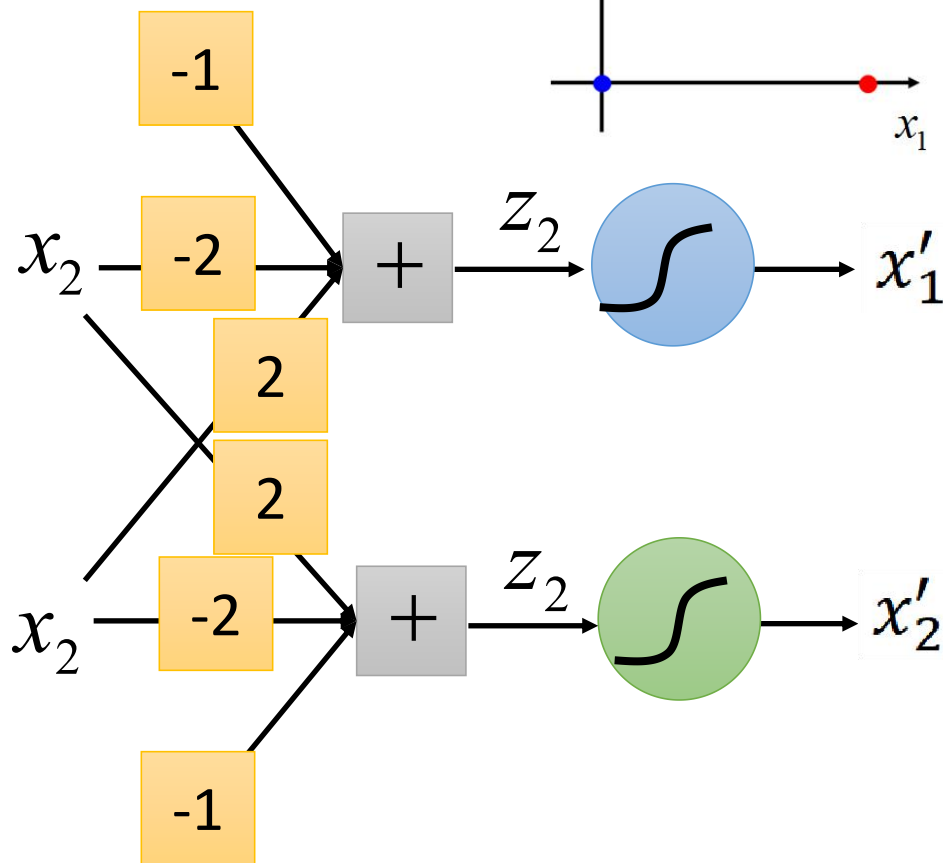
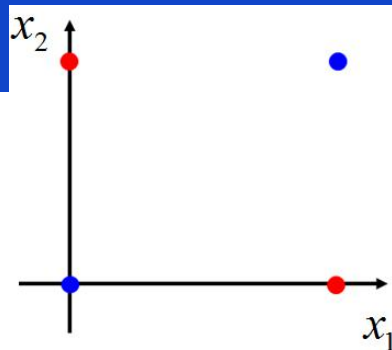
逻辑斯特回归

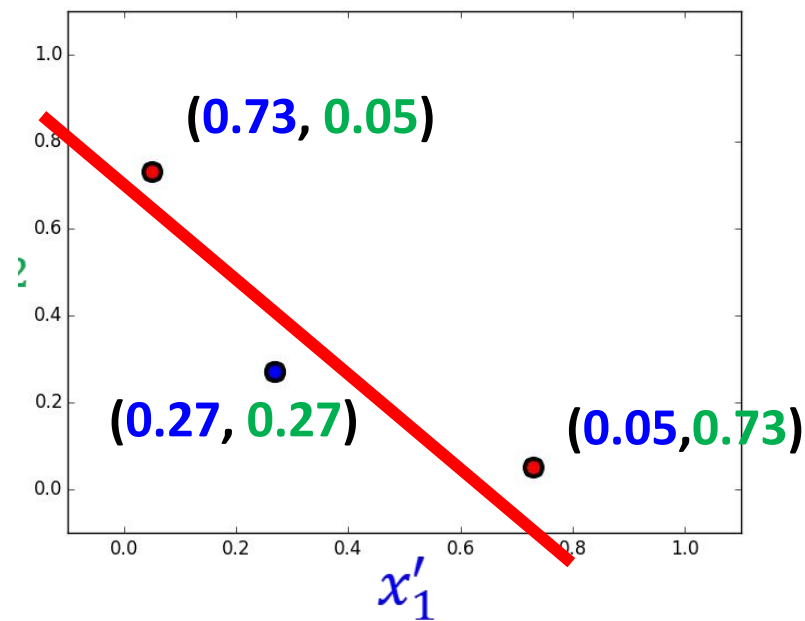
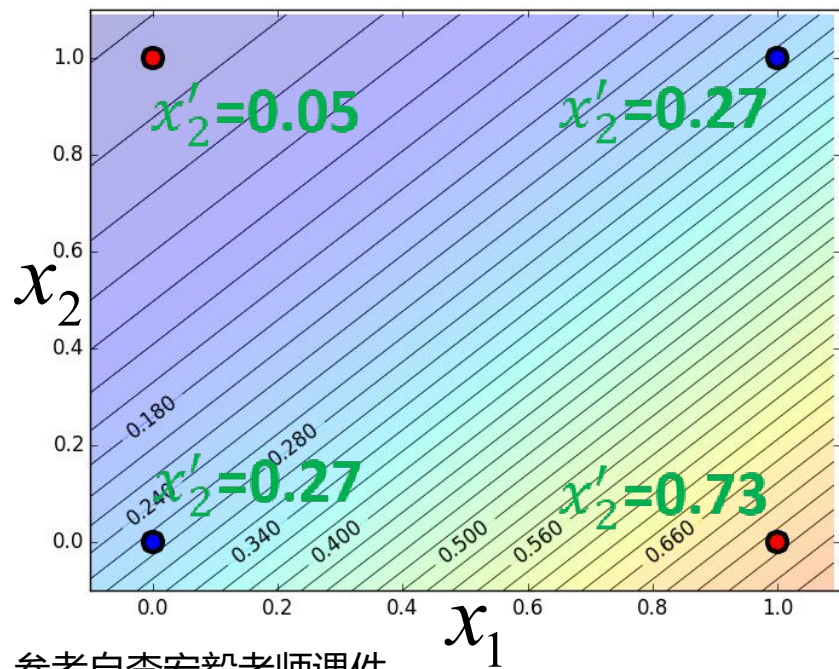
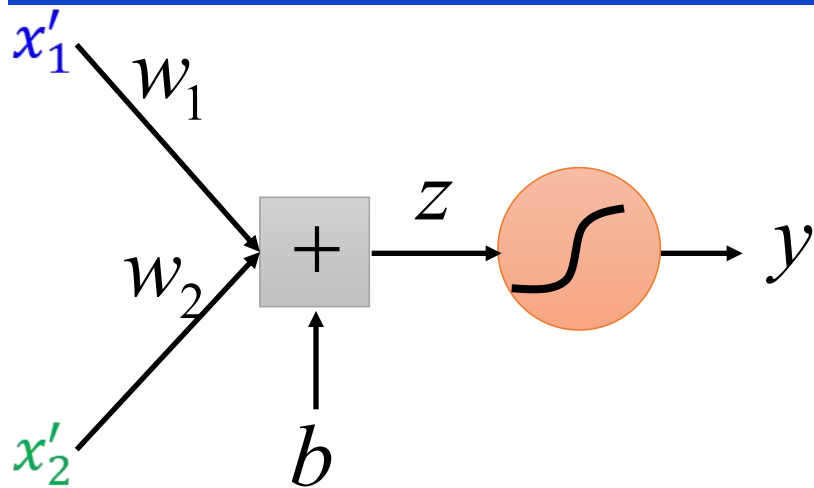
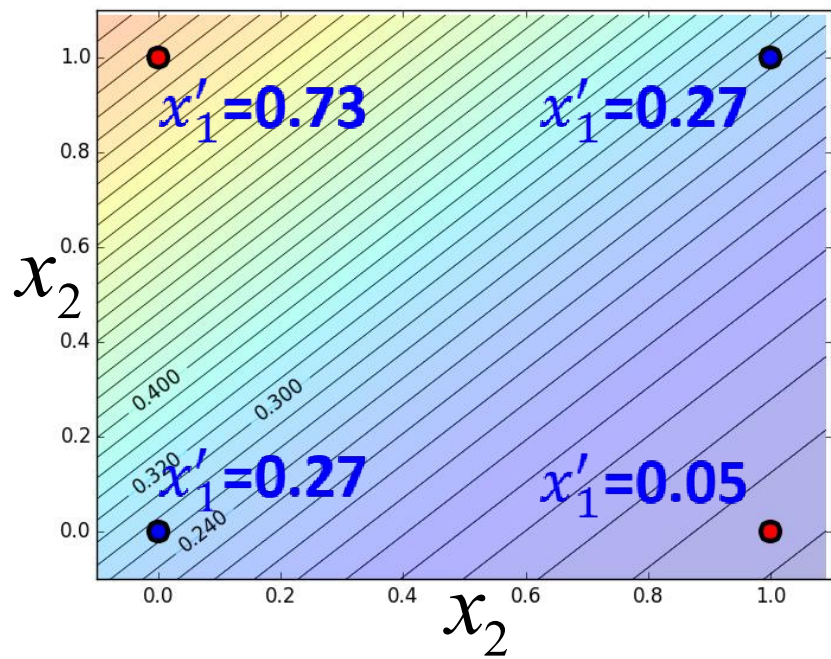
Limitation of Logistic Regression

- Cascading logistic regression models



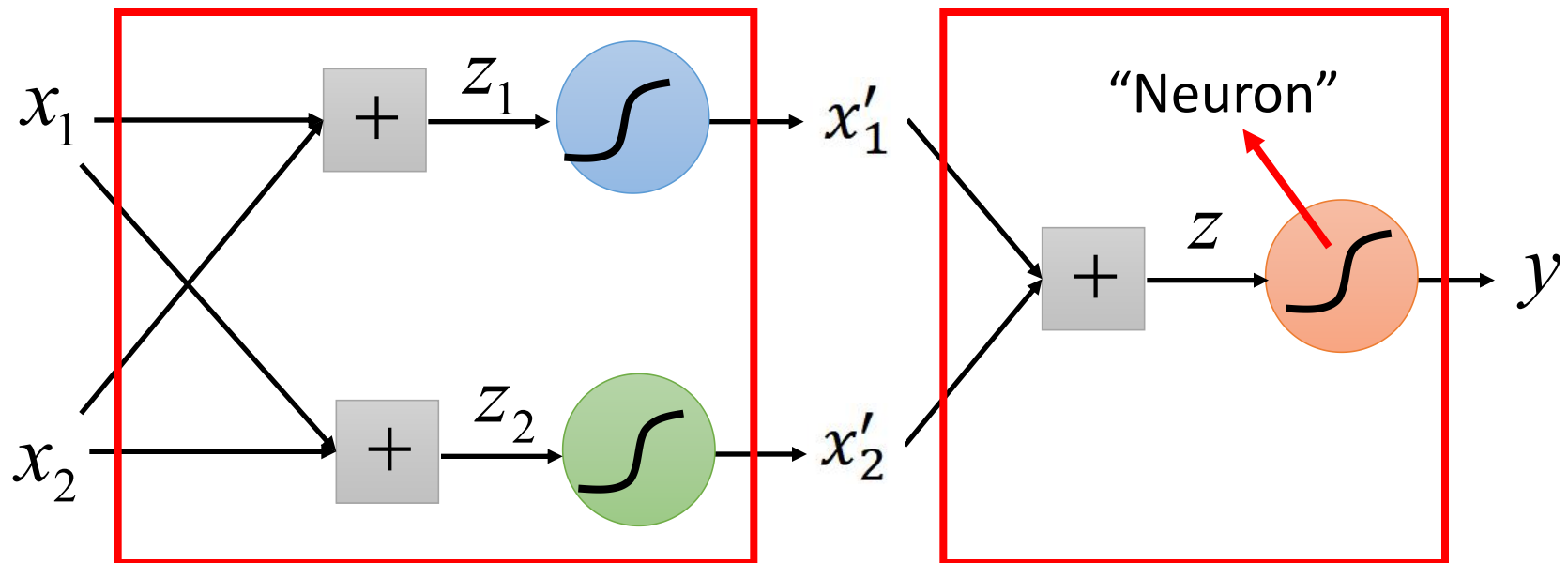
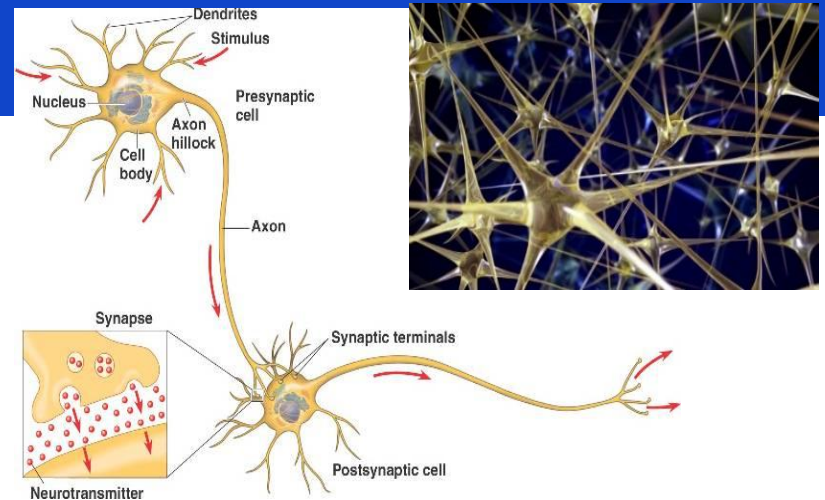
逻辑斯特回归





Deep Learning

All the parameters of the logistic regressions are jointly learned.



Feature Transformation

Classification

Neural Network

Thanks
