

逻辑（对数几率）回归

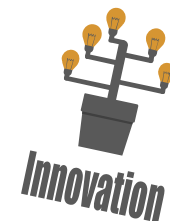
Logistic Regression

郝 奇

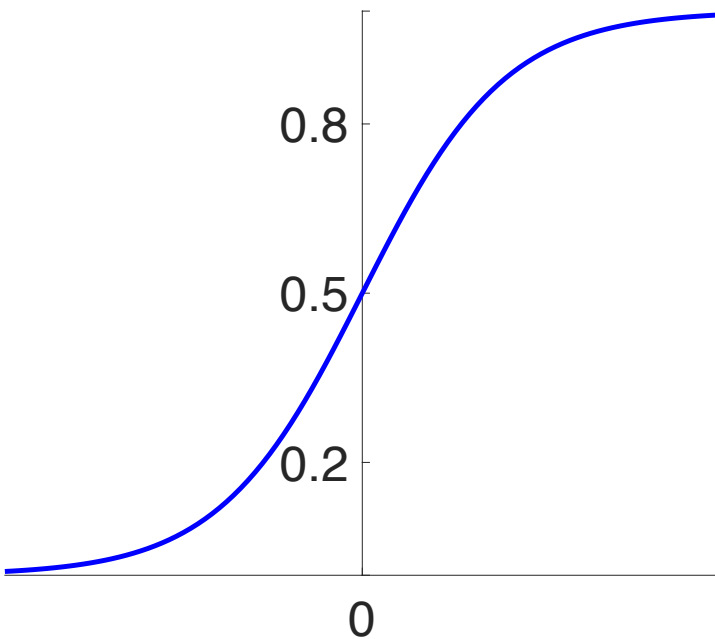
南京大学 天文与空间科学学院

**Application of
Machine Learning
in Astronomy**

机器学习在天文中的应用



Contents

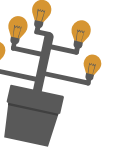


01 Logistic Regression

02 Generative & Discriminative

03 Limitation of LR

Logistic Regression



01



1. Logistic Regression




Logistic regression V.S. Linear Regression

Previous Lecture



1. Logistic Regression

Logistic regression V.S. Linear Regression

	Logistic regression	Linear regression
	?	$f_{w,b}(x) = b + \sum w_i x_i$ <p>Output: any value</p>
	?	Training data: (x^n, \hat{y}^n) \hat{y}^n : a real number $L(f) = \frac{1}{2} \sum_n (\hat{y}^n - f(x^n))^2$
	?	$w_{i+1} \leftarrow w_i - \eta \sum_n -(\hat{y}^n - f_{w,b}(x^n)) x_i^n$

1. Logistic Regression

1.1 Function set

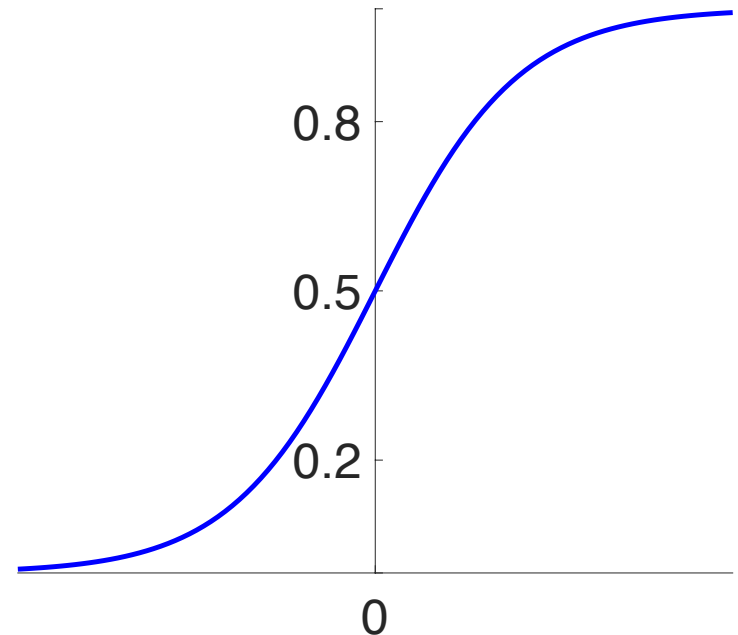
$$P(C_1|x) = \frac{P(x|C_1)P(C_1)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)}$$

If $P(C_1|x) > 0.5$ Output = class 1
else Output = class 2

$$z = wx + b$$

$$P(C_1|x) = \sigma(z) = \sigma(wx + b)$$

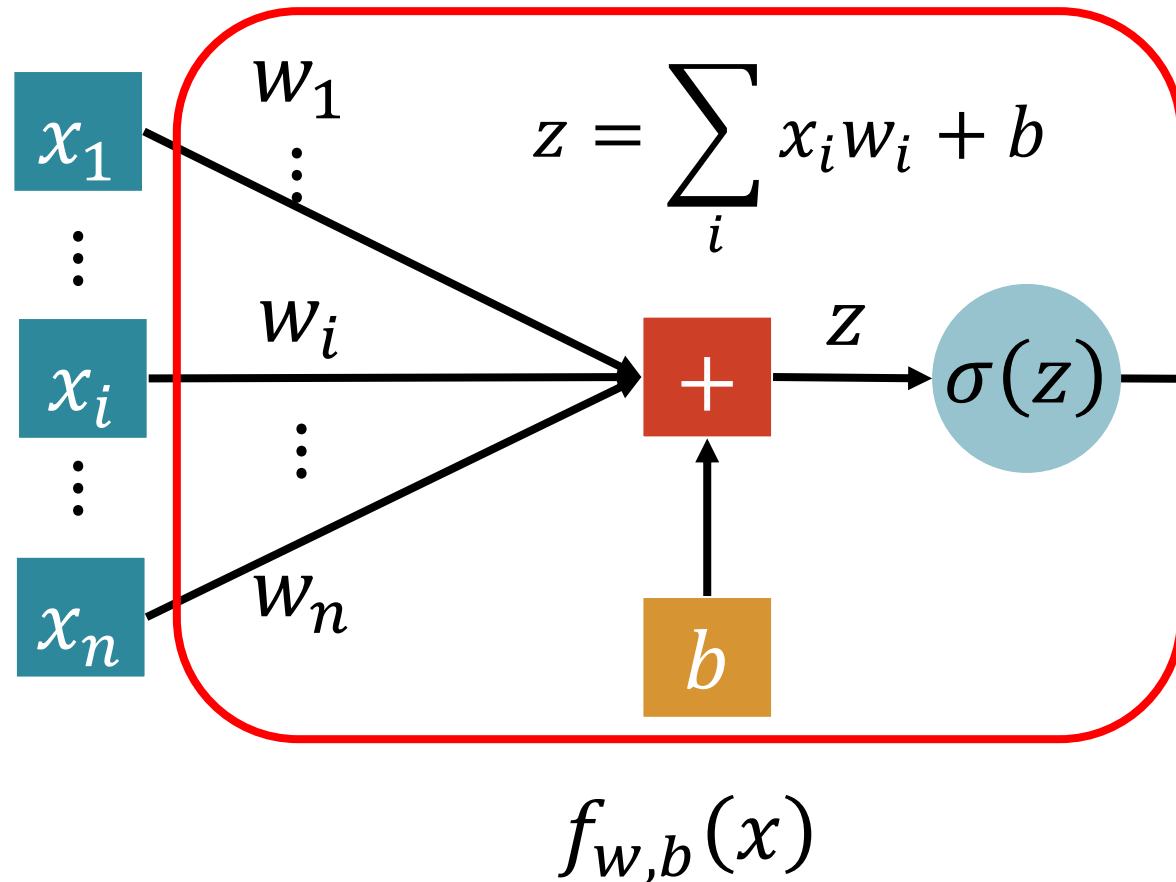
$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$



Function set: $f_{w,b}(x) = P(C_1|x)$ including all different w and b

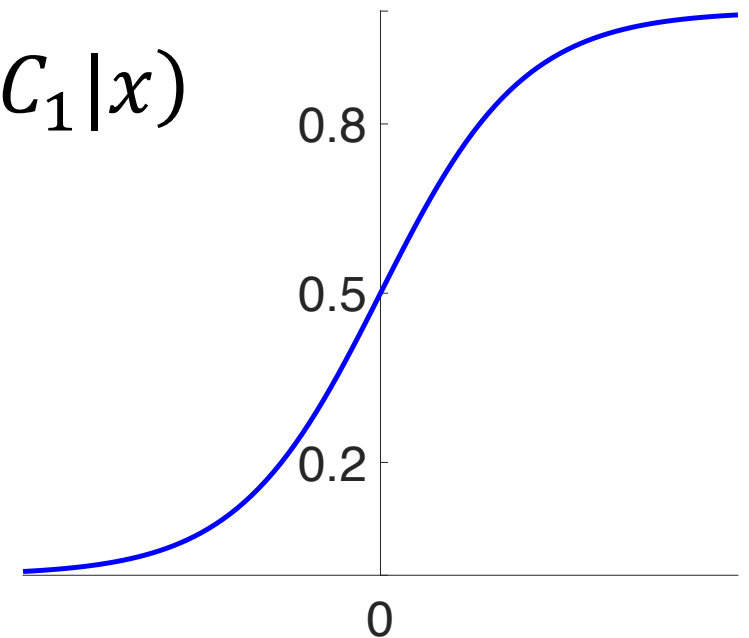
1. Logistic Regression

1.1 Function set



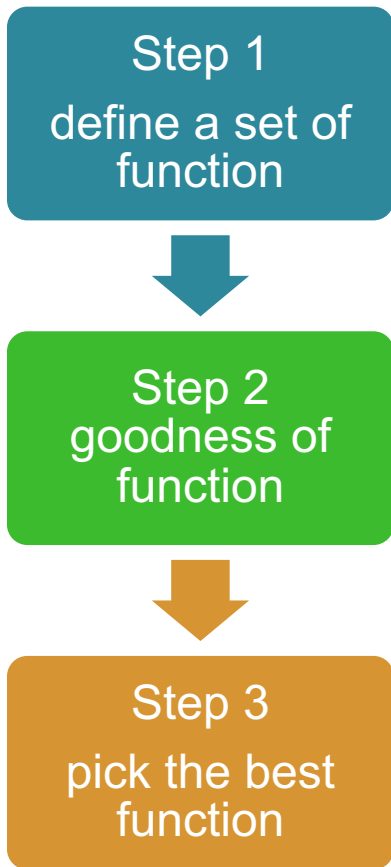
$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

$P_{w,b}(C_1 | x)$



1. Logistic Regression

Logistic regression V.S. Linear Regression



Logistic regression

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

Output: between 0 and 1

?

?

Linear regression

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

Output: any value

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

\hat{y}^n : a real number

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

$$w_{i+1} \leftarrow w_i - \eta \sum_n -(\hat{y}^n - f_{w,b}(x^n)) x_i^n$$

1. Logistic Regression

1.2 Goodness of function

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

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

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 $L(w, b)$:

$$w^*, b^* = \arg \max_{w,b} L(w, b)$$

1. Logistic Regression

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

Training Data	x^1	x^2	x^3	x^n		x^1	x^2	x^3
	C_1	C_1	C_2		C_1		$\hat{y}^1 = 1$	$\hat{y}^2 = 1$	$\hat{y}^3 = 0$	

$$w^*, b^* = \arg \max_{w,b} L(w, b) \iff w^*, b^* = \arg \min_{w,b} -\ln L(w, b)$$

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

$$-\ln f_{w,b}(x^1) \quad \longrightarrow \quad -\left[\hat{y}^1 \ln f_{w,b}(x^1) + (1 - \hat{y}^1) \ln (1 - f_{w,b}(x^1))\right]$$

$$-\ln f_{w,b}(x^2) \quad \longrightarrow \quad -\left[\hat{y}^2 \ln f_{w,b}(x^2) + (1 - \hat{y}^2) \ln (1 - f_{w,b}(x^2))\right]$$

$$-\ln (1 - f_{w,b}(x^3)) \quad \longrightarrow \quad -\left[\hat{y}^3 \ln f_{w,b}(x^3) + (1 - \hat{y}^3) \ln (1 - f_{w,b}(x^3))\right]$$

⋮

1. Logistic Regression

1.2 Goodness of function

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

$$\begin{aligned} -\ln L(w, b) &= \ln f_{w,b}(x^1) + \ln f_{w,b}(x^2) + \ln(1 - f_{w,b}(x^3)) \cdots + \ln f_{w,b}(x^n) \\ &= \sum_n \underbrace{-\left[\hat{y}^n \ln f_{w,b}(x^n) + (1 - \hat{y}^n) \ln(1 - f_{w,b}(x^n)) \right]}_{\text{Cross entropy between two Bernoulli distribution}} \end{aligned}$$

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:

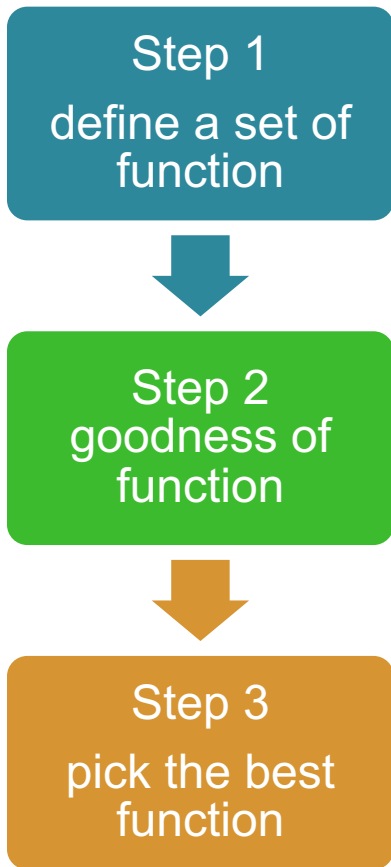
$$q(x = 1) = f(x^n)$$

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

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

1. Logistic Regression

Logistic regression V.S. Linear Regression



Logistic regression

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

Output: between 0 and 1

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

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

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

?

Linear regression

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

Output: any value

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

\hat{y}^n : a real number

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

$$w_{i+1} \leftarrow w_i - \eta \sum_n -(\hat{y}^n - f_{w,b}(x^n)) x_i^n$$

1. Logistic Regression

1.2 Goodness of function

Cross entropy

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

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

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

Why don't we simply use square error as linear regression?

1. Logistic Regression

Step 1
define a set of
function



Step 2
goodness of
function



Step 3
pick the best
function

$$f_{w,b}(x) = \sigma \left(b + \sum w_i x_i \right) \quad z = \sum_i x_i w_i + b \quad \sigma(z) = \frac{1}{1 + \exp(-z)}$$

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

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

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

$\hat{y}^n = 1$

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

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

1. Logistic Regression

Step 1
define a set of
function



Step 2
goodness of
function



Step 3
pick the best
function

$$f_{w,b}(x) = \sigma \left(b + \sum w_i x_i \right) \quad z = \sum_i x_i w_i + b \quad \sigma(z) = \frac{1}{1 + \exp(-z)}$$

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

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

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

$$\hat{y}^n = 0$$

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

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

1. Logistic Regression

1.3 Find the best function

$$z = \sum_i x_i w_i + b \quad \sigma(z) = \frac{1}{1 + \exp(-z)}$$

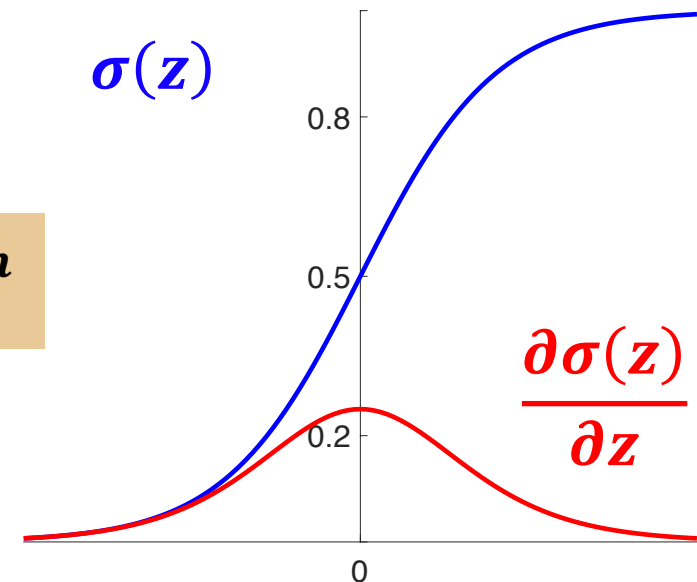
$$\frac{-\ln L(w, b)}{\partial w_i} = \sum_n - \left[\hat{y}^n \frac{\ln f_{w,b}(x^n)}{\partial w_i} + (1 - \hat{y}^n) \frac{\ln(1 - f_{w,b}(x^n))}{\partial w_i} \right]$$

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

$$\frac{\partial z}{\partial w_i} = x_i$$

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

$$(1 - f_{w,b}(x^n)) x_i^n$$



1. Logistic Regression

1.3 Find the best function

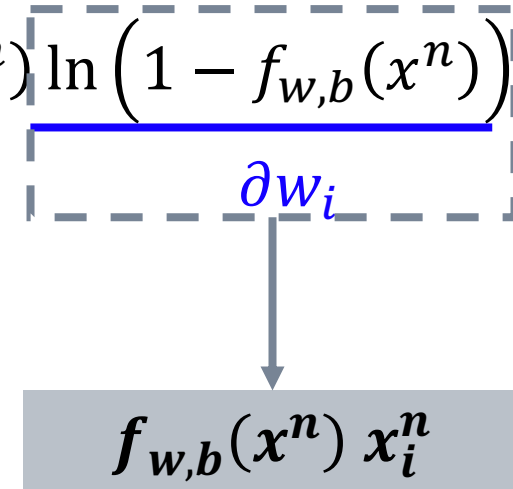
$$z = \sum_i x_i w_i + b \quad \sigma(z) = \frac{1}{1 + \exp(-z)}$$

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

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

$$\frac{\partial z}{\partial w_i} = x_i$$

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



1. Logistic Regression

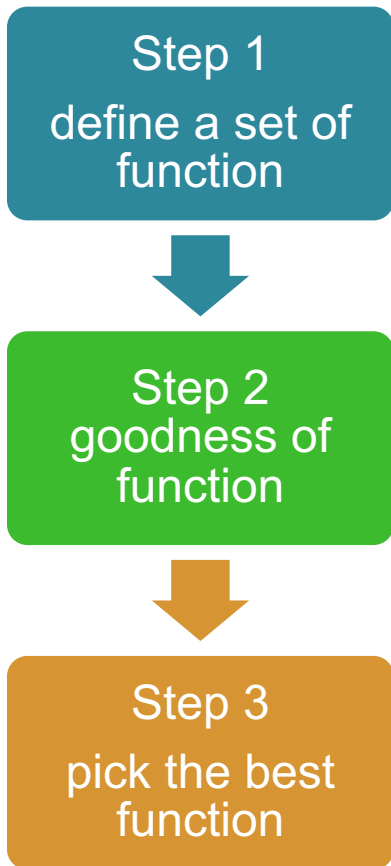
1.3 Find the best function

$$\begin{aligned} \frac{-\ln L(w, b)}{\partial w_i} &= \sum_n - \left[\hat{y}^n \underbrace{\ln f_{w,b}(x^n)}_{\partial w_i} + (1 - \hat{y}^n) \underbrace{\ln (1 - f_{w,b}(x^n))}_{\partial w_i} \right] \\ &= \sum_n - \left[\hat{y}^n (1 - f_{w,b}(x^n)) x_i^n + (1 - \hat{y}^n) f_{w,b}(x^n) x_i^n \right] \\ &= \sum_n - \left[\hat{y}^n - \cancel{\hat{y}^n f_{w,b}(x^n)} - f_{w,b}(x^n) + \cancel{\hat{y}^n f_{w,b}(x^n)} \right] x_i^n \\ &= \sum_n - \left(\hat{y}^n - f_{w,b}(x^n) \right) x_i^n \quad w_{i+1} \leftarrow w_i - \eta \sum_n - \left(\hat{y}^n - f_{w,b}(x^n) \right) x_i^n \end{aligned}$$

**Larger difference,
larger update**

1. Logistic Regression

Logistic regression V.S. Linear Regression



Logistic regression

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

Output: between 0 and 1

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

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

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

$$w_{i+1} \leftarrow w_i - \eta \sum_n -(\hat{y}^n - f_{w,b}(x^n)) x_i^n$$

Linear regression

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

Output: any value

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

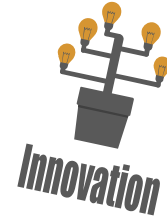
\hat{y}^n : a real number

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

$$w_{i+1} \leftarrow w_i - \eta \sum_n -(\hat{y}^n - f_{w,b}(x^n)) x_i^n$$

Generative & Discriminative

02



2. Generative & Discriminative

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

Find w and b directly

Find μ^1, μ^2, Σ

Shall we obtain the
same set of w and b ?

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

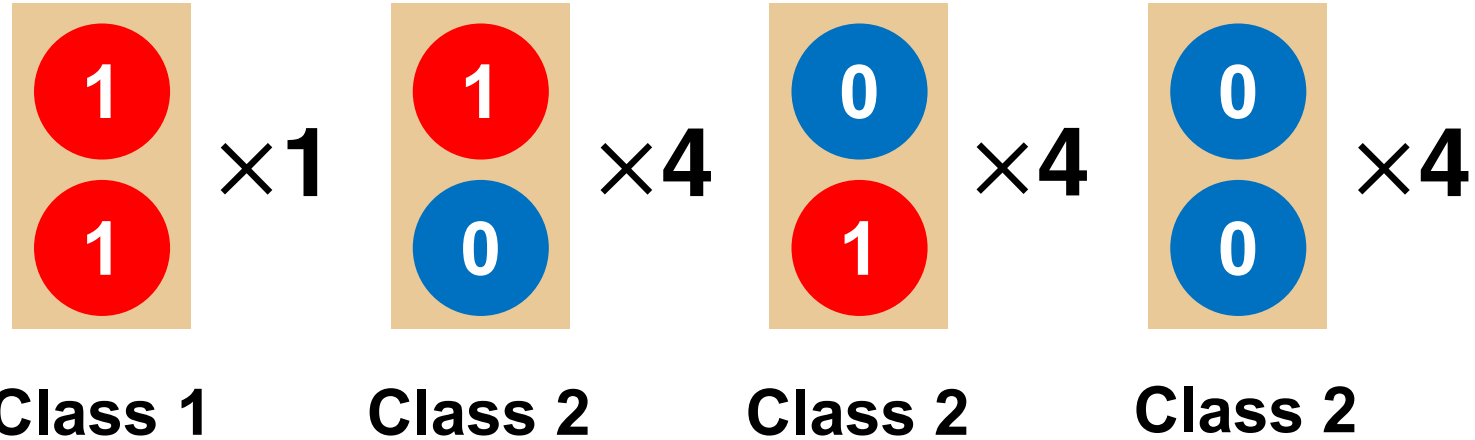
$$b = -\frac{1}{2} (\mu^1)^T (\Sigma^1)^{-1} \mu^1 \\ + \frac{1}{2} (\mu^2)^T (\Sigma^2)^{-1} \mu^2 + \ln \frac{N_1}{N_2}$$

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

2. Generative & Discriminative

2.1 Toy example

Training
Data



Testing
Data

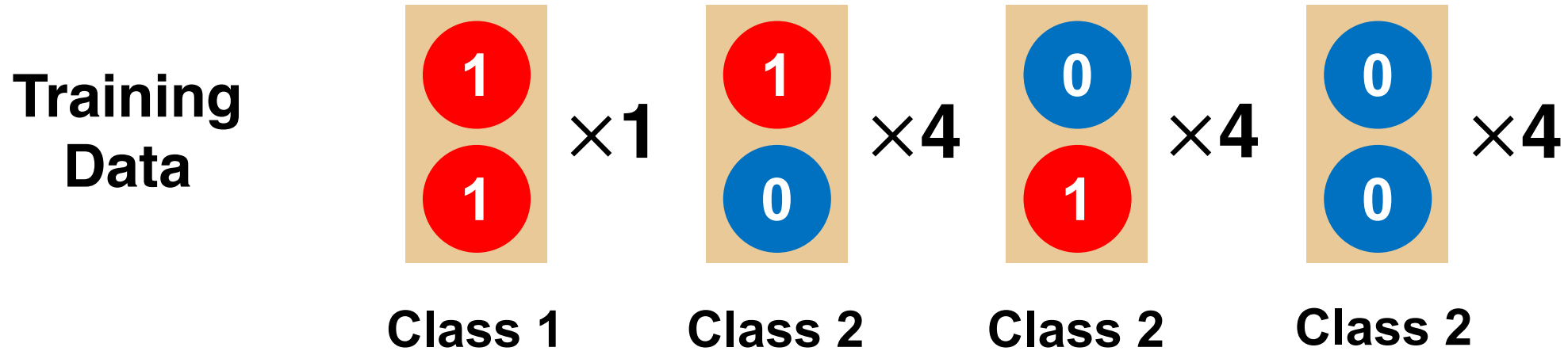


How about Naïve Bayes?

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

2. Generative & Discriminative

2.1 Toy example



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

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

$$P(x_2 = 1|C_2) = 1$$

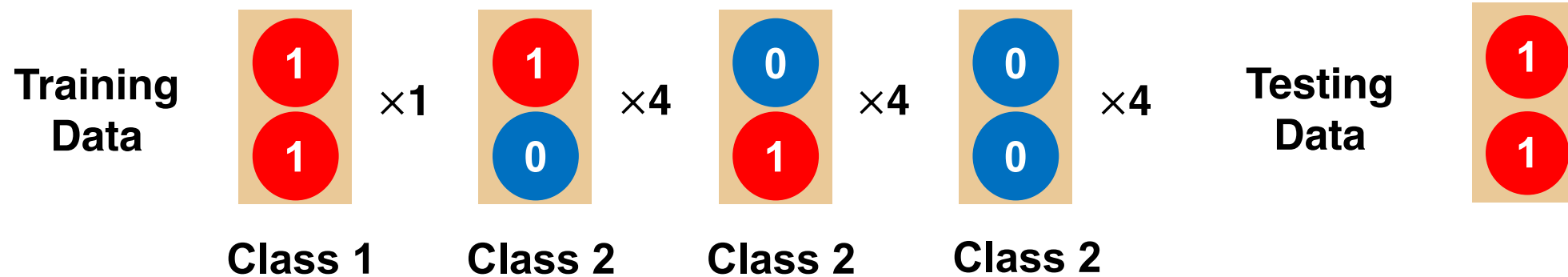
$$P(C_2) = \frac{12}{13}$$

$$P(x_1 = 1|C_2) = \frac{4}{12}$$

$$P(x_2 = 1|C_2) = \frac{4}{12}$$

2. Generative & Discriminative

2.1 Toy example



$$P(C_1) = \frac{1}{13} \quad P(x_1 = 1|C_1) = 1 \quad P(x_2 = 1|C_2) = 1 \quad P(x|C_1) = P(x_1 = 1|C_1)P(x_2 = 1|C_1) = 1 \times 1 = 1$$

$$P(C_2) = \frac{12}{13} \quad P(x_1 = 1|C_2) = \frac{4}{12} \quad P(x_2 = 1|C_2) = \frac{4}{12} \quad P(x|C_2) = P(x_1 = 1|C_2)P(x_2 = 1|C_2) = \frac{4}{12} \times \frac{4}{12} = \frac{1}{9}$$

$$P(C_1|x) = \frac{P(x|C_1)P(C_1)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)} = \frac{1 \times \frac{1}{13}}{1 \times \frac{1}{13} + \frac{1}{9} \times \frac{12}{13}} = \frac{3}{7} < 0.5$$

2. Generative & Discriminative

2.2 Characteristics of Generative models

- **Focus on probability distributions of the data**
 - With the assumption of probability distribution, they can generate samples;
- **More powerful with less of examples**
 - With the assumption of probability distribution, less training data is needed.
- **Can generate new samples**
 - Can be used in semi-supervised learning and unsupervised learning
- **Priors and class-dependent probabilities can be estimated from different sources**

2. Generative & Discriminative

Characteristics of Discriminative models

Previous Lecture

■ Focus on decision boundary

- they do not allow one to generate samples from the joint distribution $\mathbf{P}(x, y)$;

■ More powerful with lot of examples

- Whenever the training data is big ,the accuracy for future data will be good.

■ Not designed to use unlabeled data

- Most discriminative models are inherently supervised and cannot easily be extended to unsupervised learning.

■ Outperform generative models at conditional prediction tasks

- When the number of parameters is limited, for tasks such as classification and regression that do not require the joint distribution.

2. Generative & Discriminative

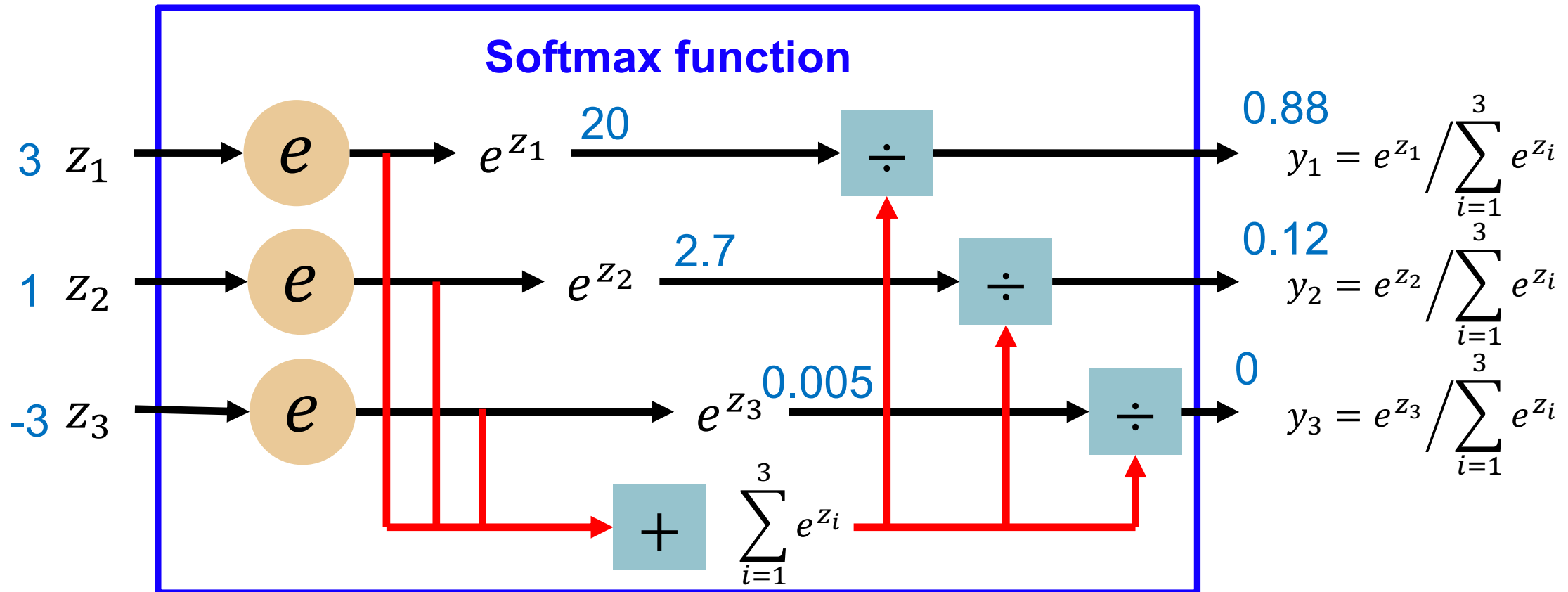
2.3 Multi-class classification

Probability: $0 < y_i < 1; \sum_i y_i = 1.$

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

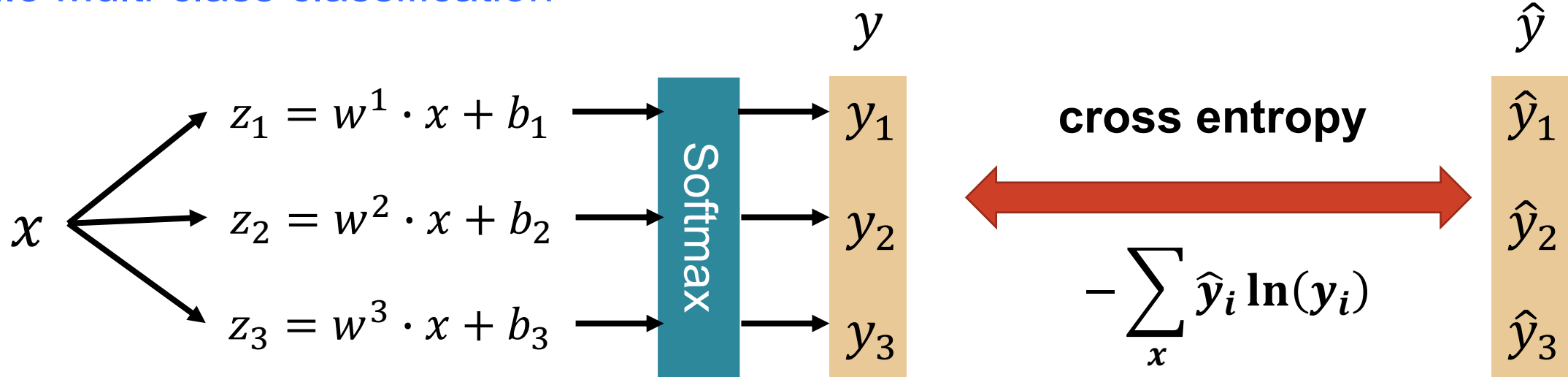
$$C_2: w^2, b_2 \quad z_2 = w^2 \cdot x + b_2$$

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



2. Generative & Discriminative

2.3 Multi-class classification



If $x \in$ class 1

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

If $x \in$ class 2

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

If $x \in$ class 3

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

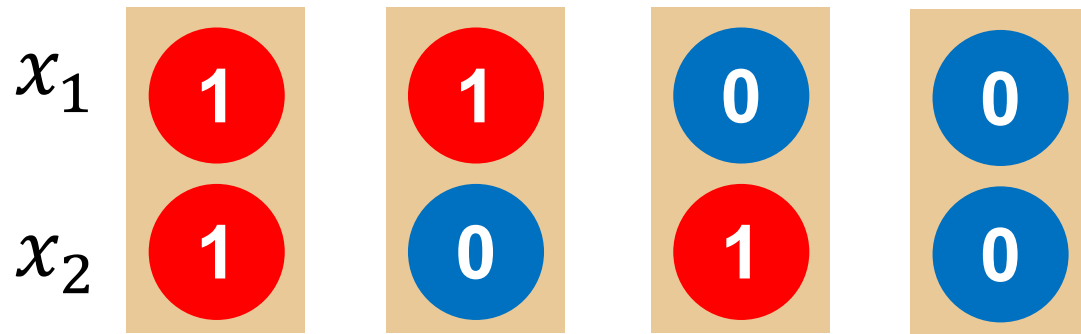
Limitation of Logistic Regression



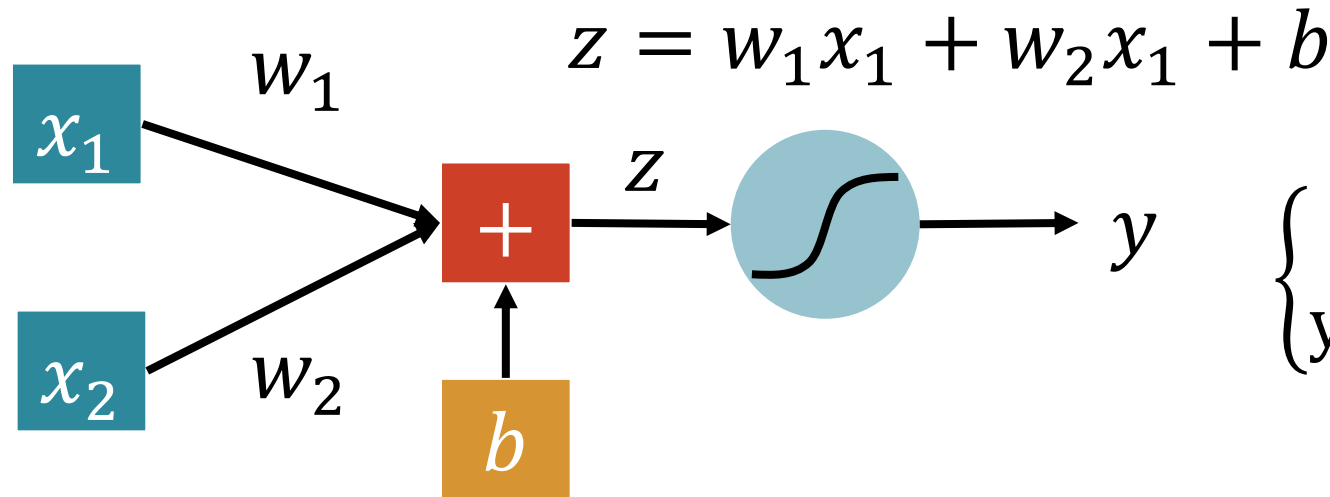
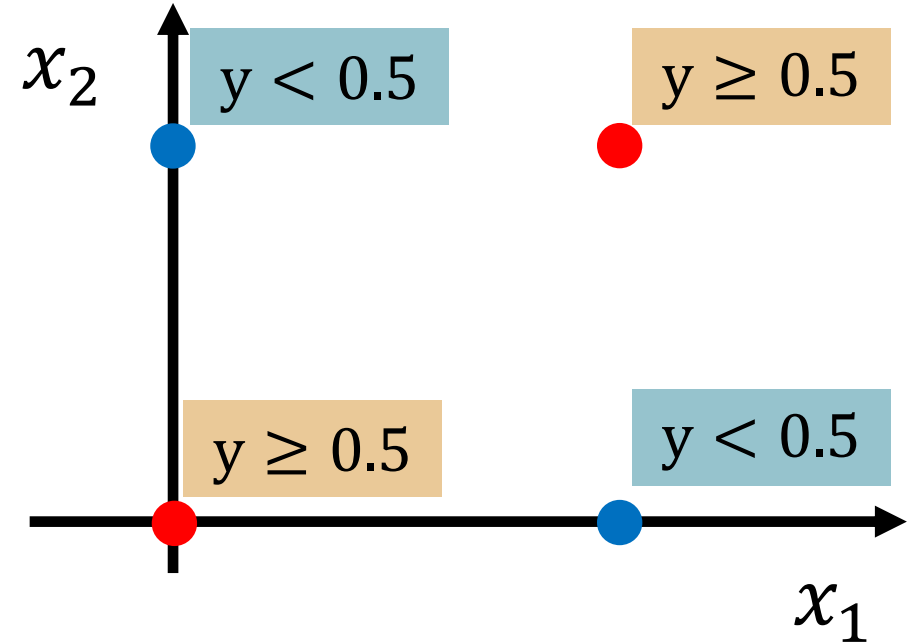
03



3. Limitation of Logistic Regression

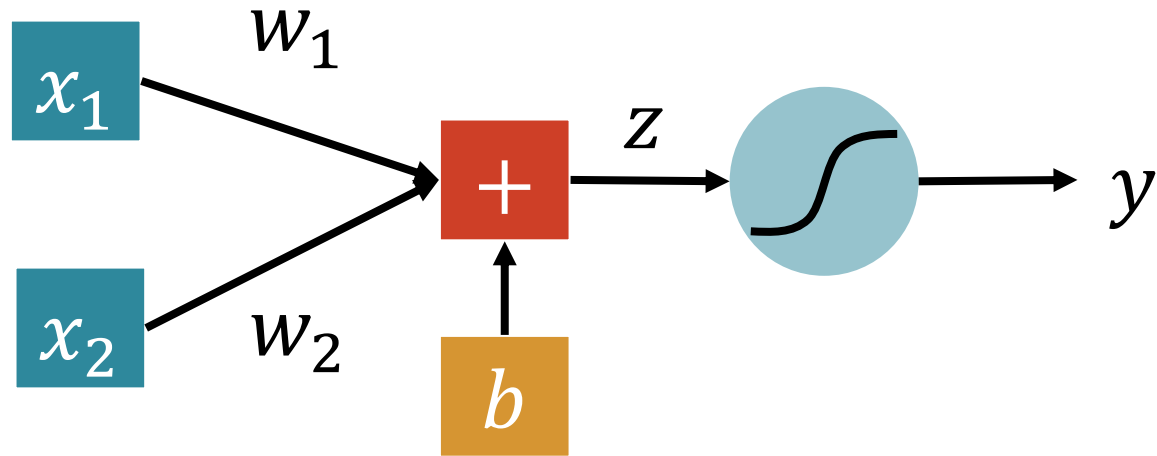


Class 1 Class 2 Class 2 Class 1

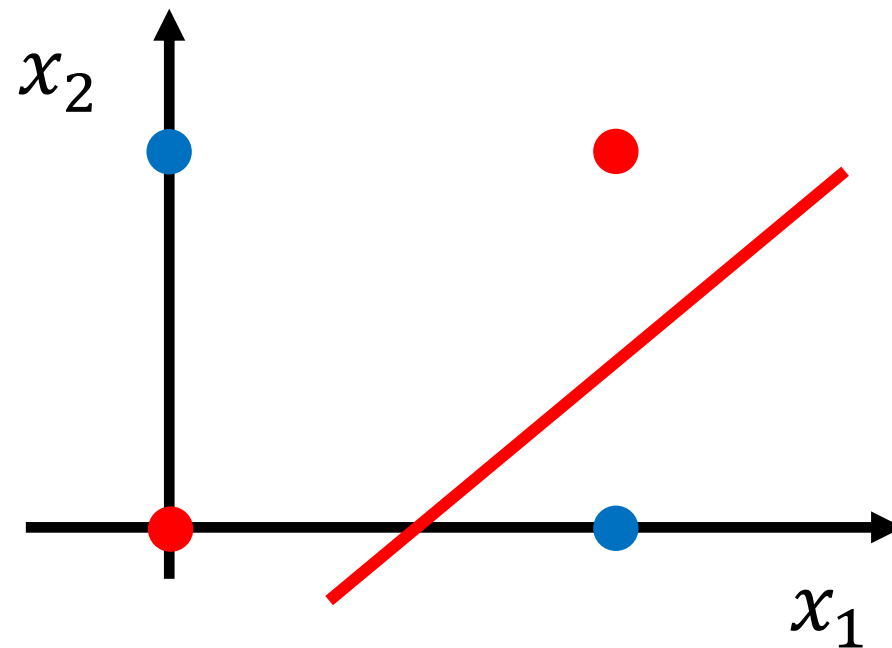
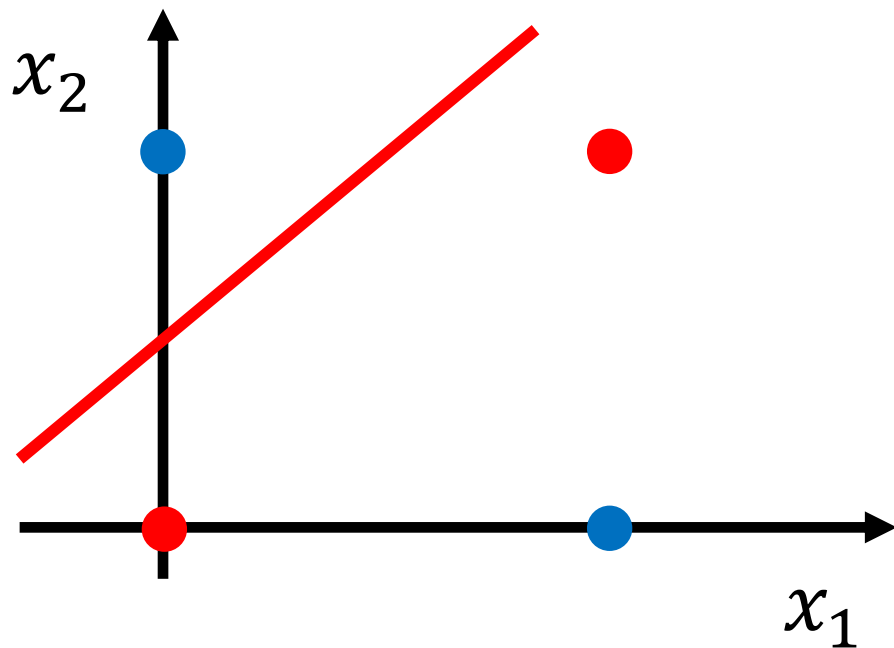


$\begin{cases} y \geq 0.5 & \text{Output} = \text{class 1} \\ y < 0.5 & \text{Output} = \text{class 2} \end{cases}$

3. Limitation of Logistic Regression

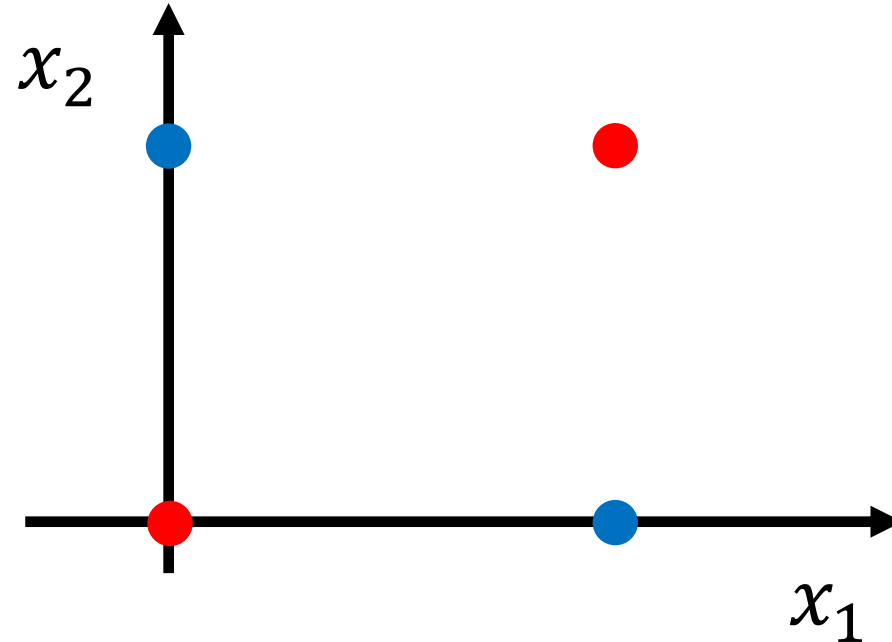


Sorry, we can't...



3. Limitation of Logistic Regression

How to classify this ?



3. Limitation of Logistic Regression

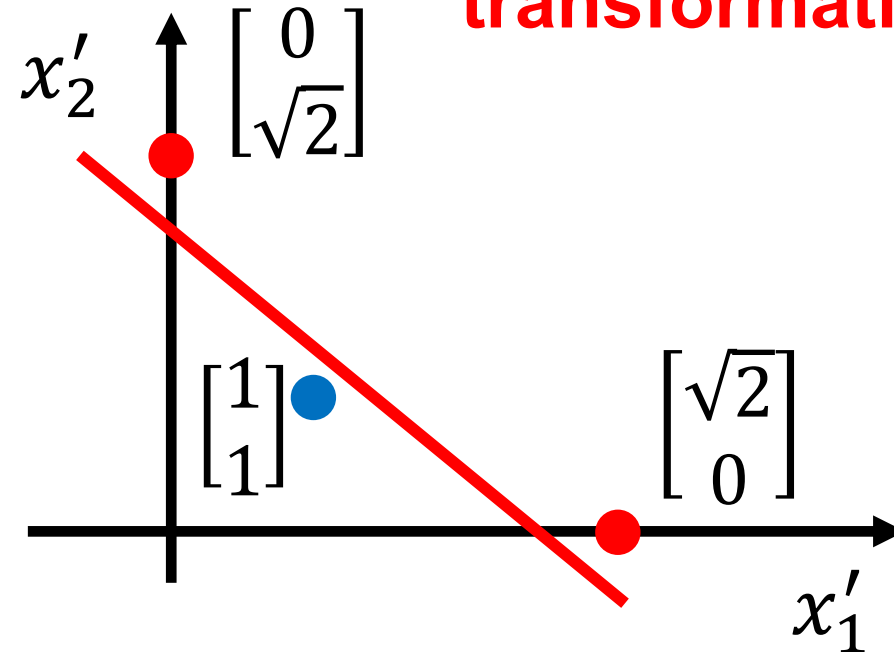
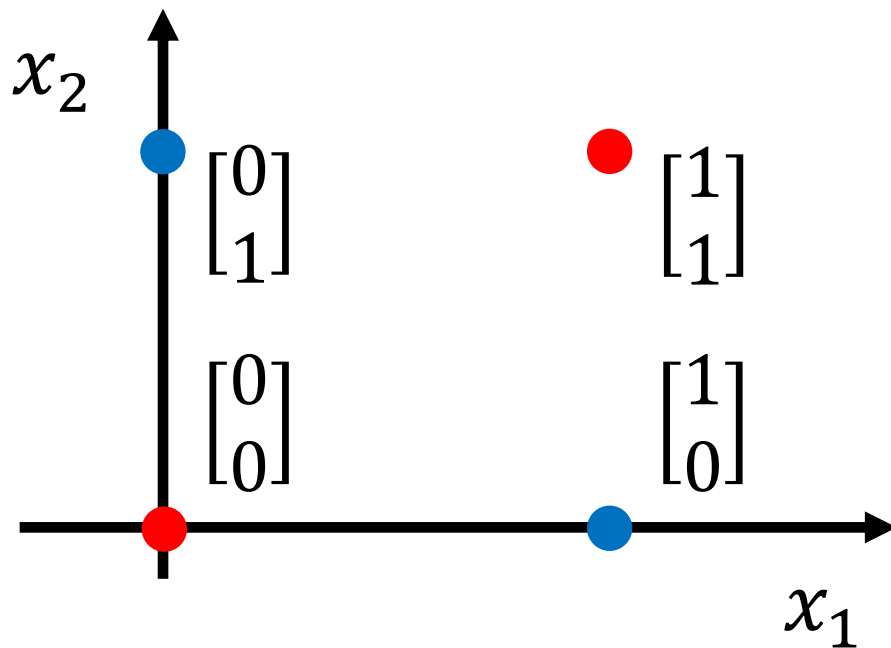
Solution 1: Feature Transformation

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

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

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

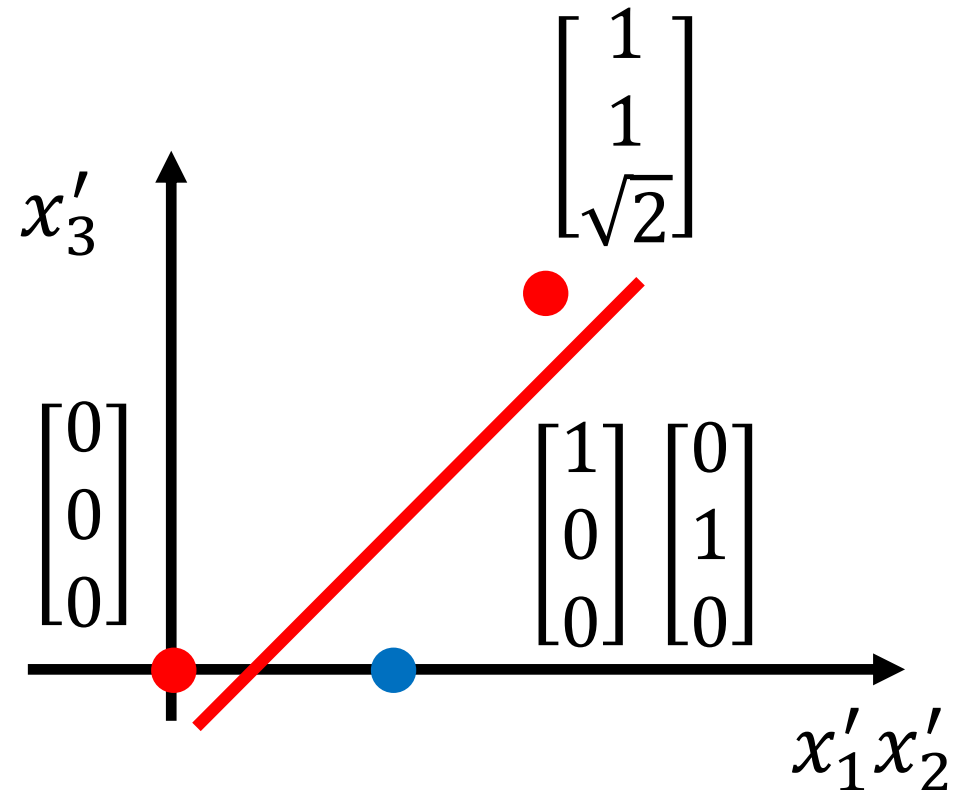
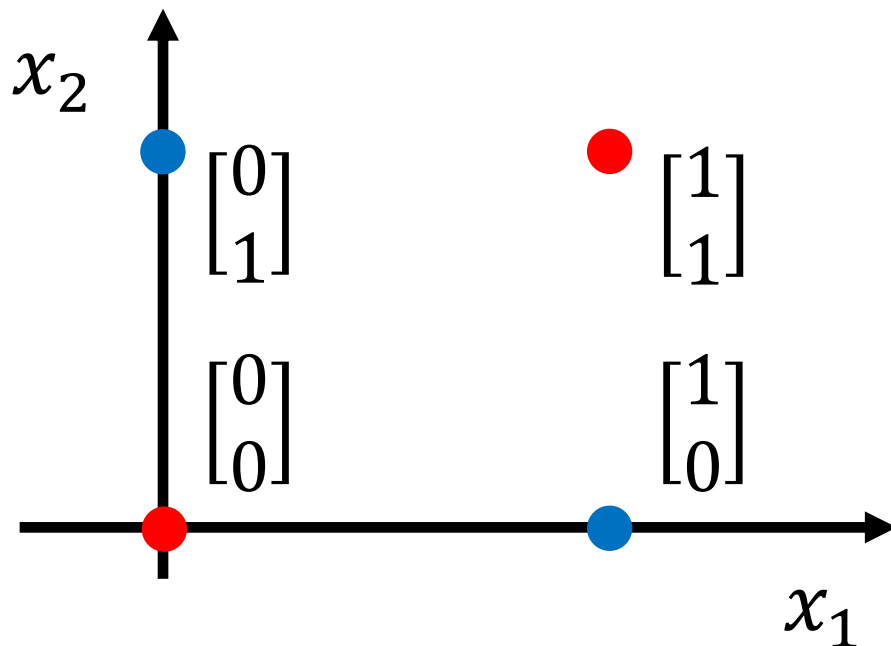
**Not always easy
to find a good
transformation**



3. Limitation of Logistic Regression

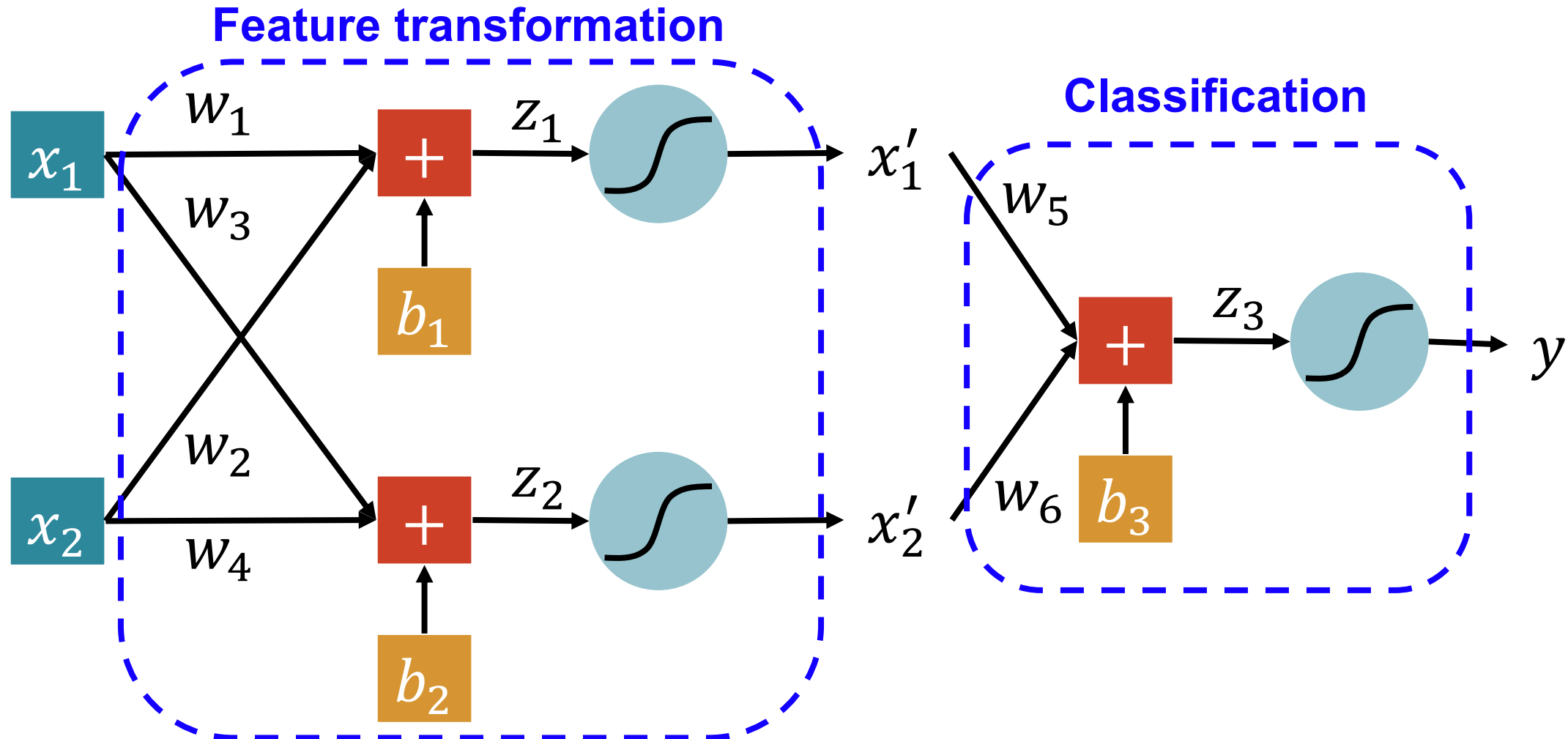
Solution 2: Map data to higher dimension: SVM

$$\Phi: \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \rightarrow \begin{bmatrix} x'_1 \\ x'_2 \\ x'_3 \end{bmatrix} = \begin{bmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2}x_1x_2 \end{bmatrix} \quad \mathbb{R}^2 \rightarrow \mathbb{R}^3$$



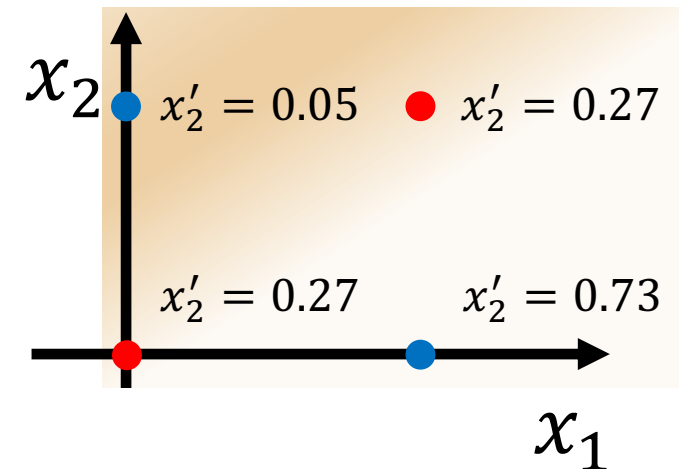
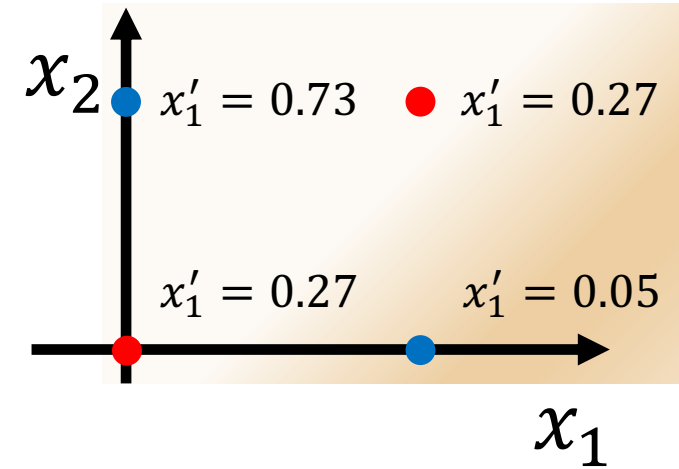
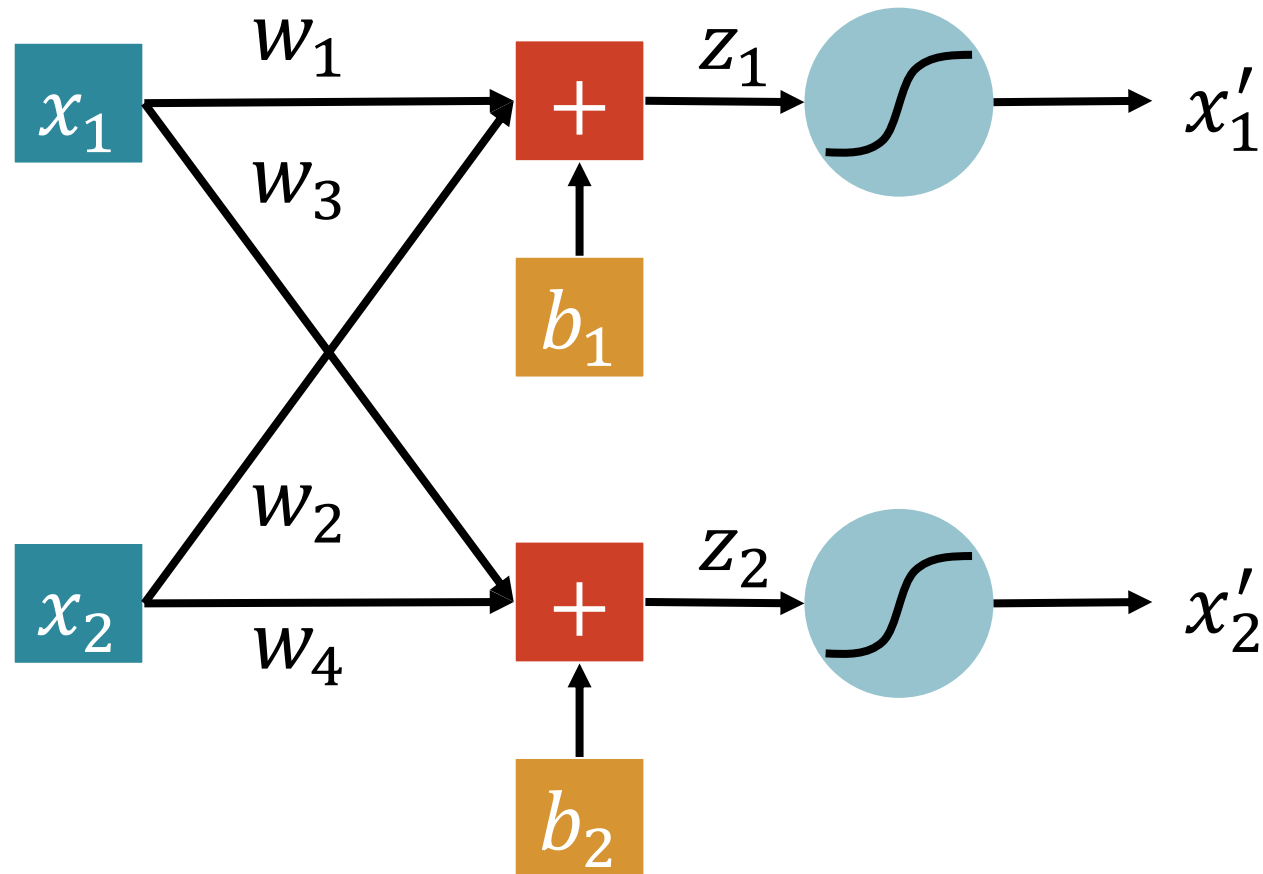
3. Limitation of Logistic Regression

Solution 3: Cascading logistic regression models



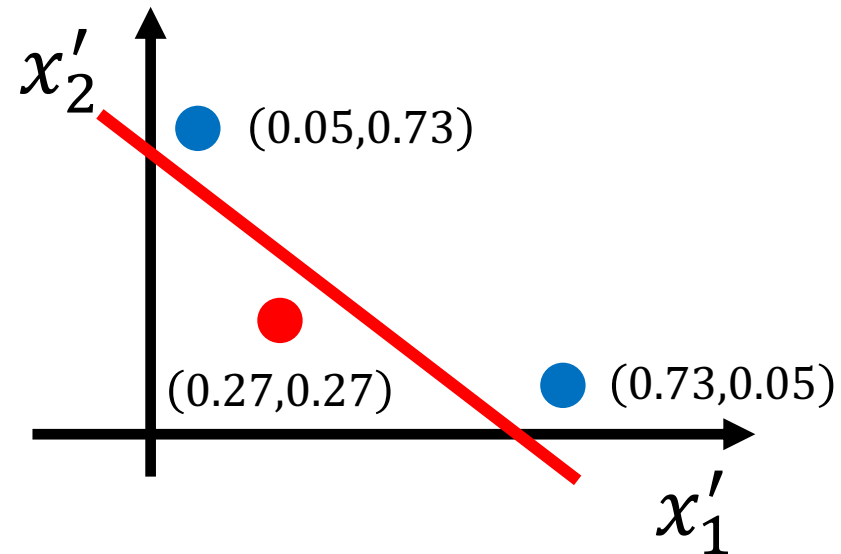
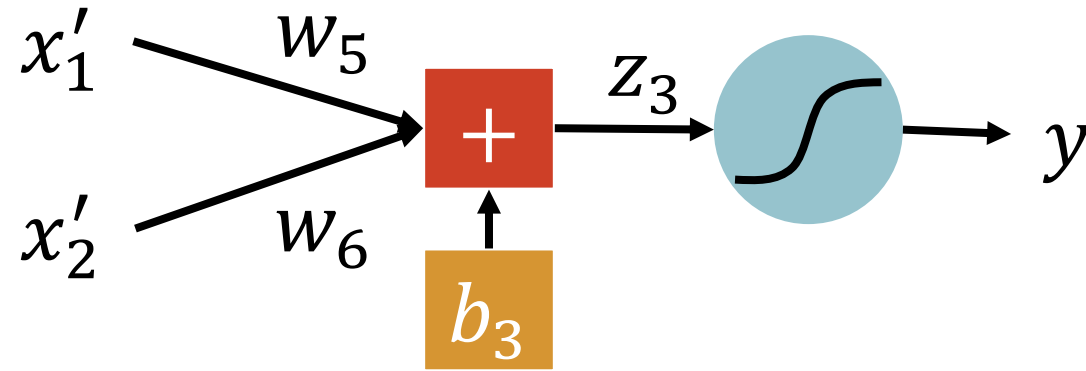
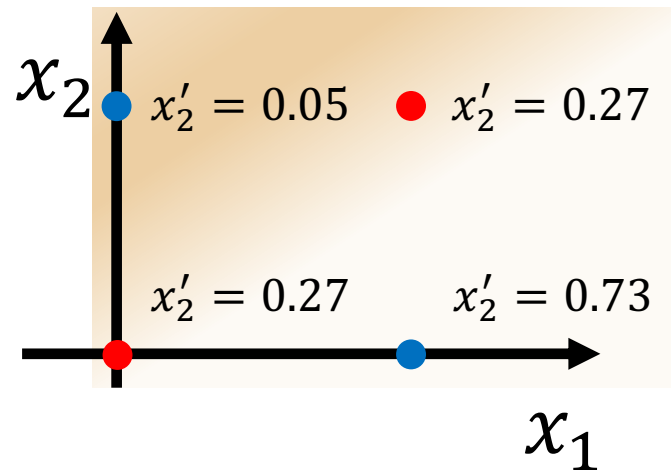
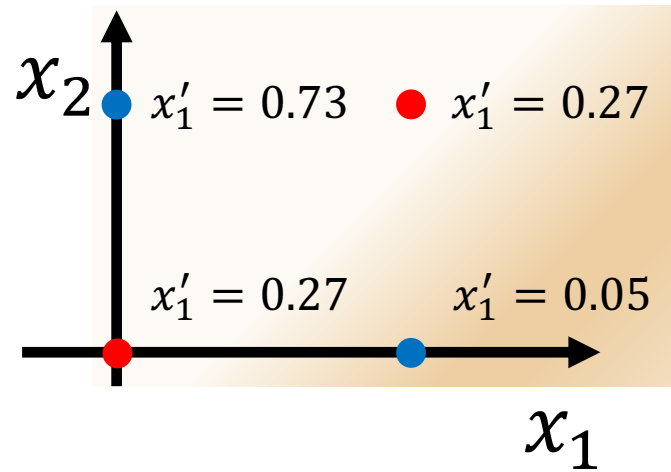
3. Limitation of Logistic Regression

Solution 3: Cascading logistic regression models



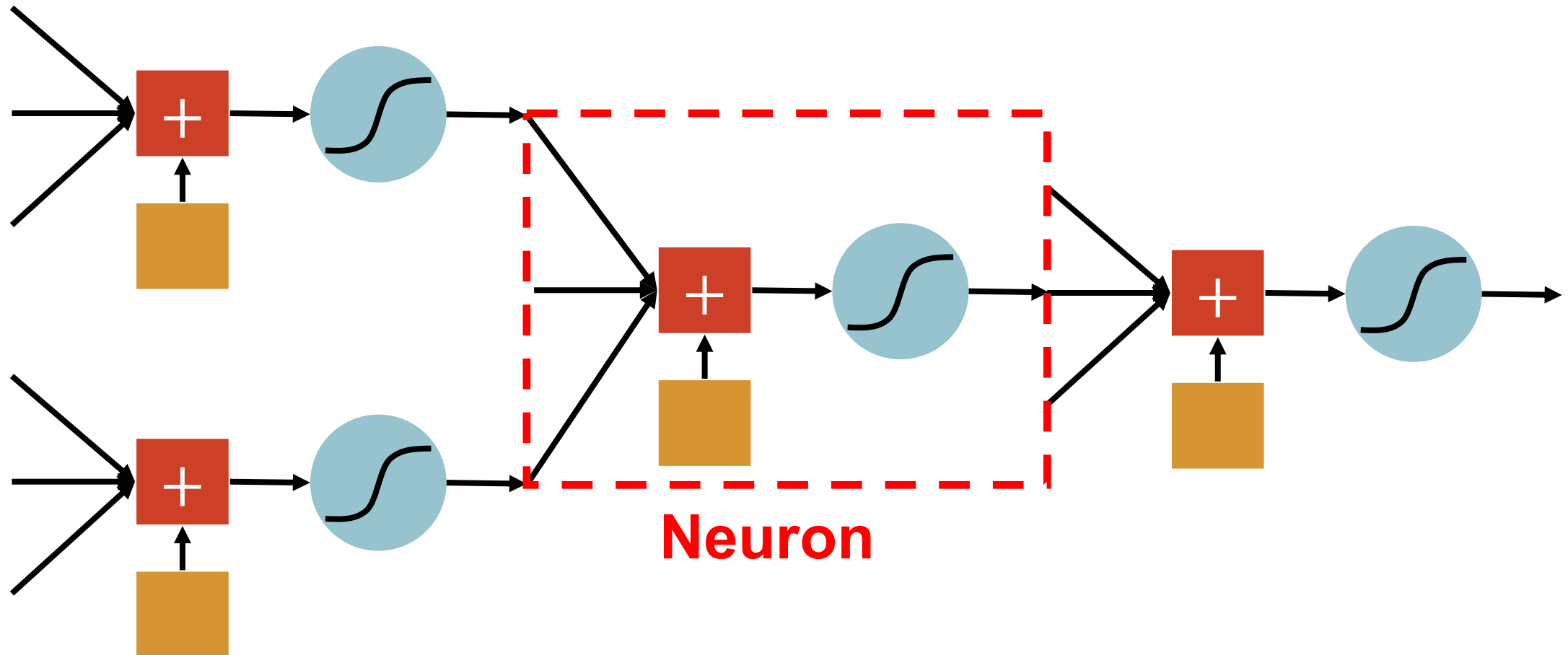
3. Limitation of Logistic Regression

Solution 3: Cascading logistic regression models



3. Limitation of Logistic Regression

Neural Network



Logistic Regression

Hao, Qi

School of Astronomy and Space Science

THANKS

