

# 计算机视觉

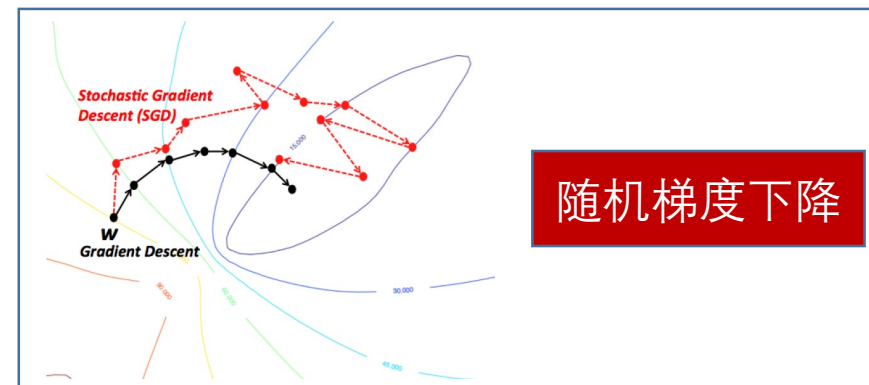
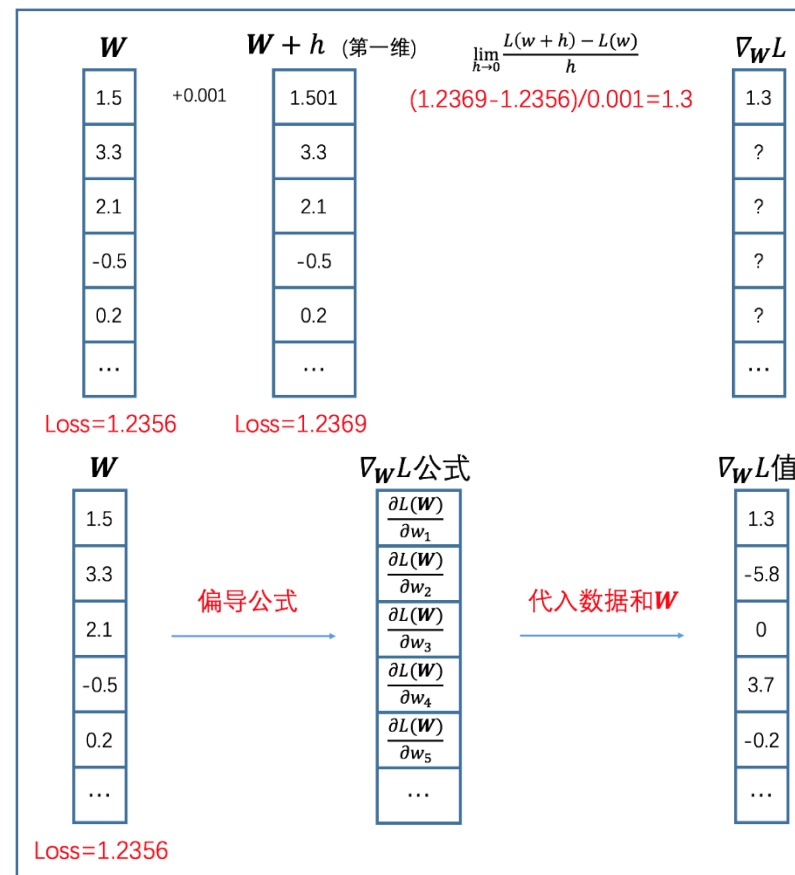
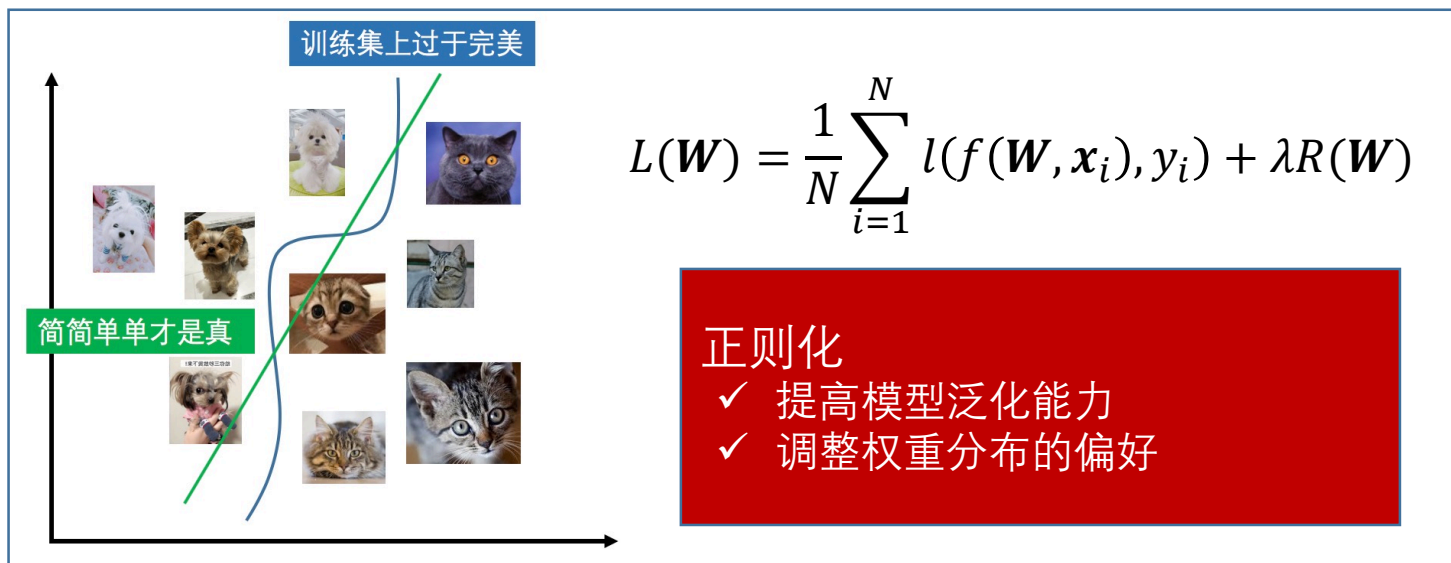
# Computer Vision

Lecture 4: 神经网络和反向传播

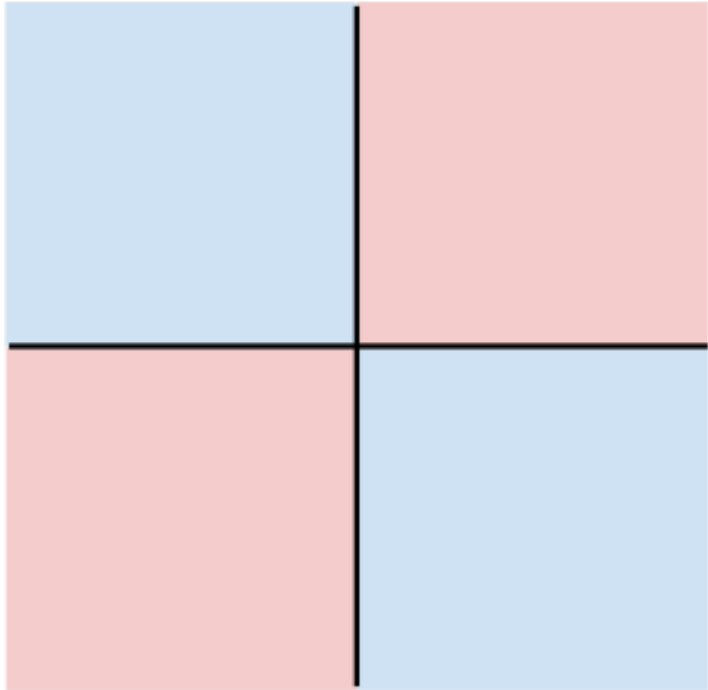


# L03: 损失函数和优化

- multiclass SVM loss (Hinge loss) :  $L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$   
✓ SVM分类器
- cross-entropy loss :  $L_i = -\log(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}})$   
✓ Softmax分类器 (多类别逻辑回归)

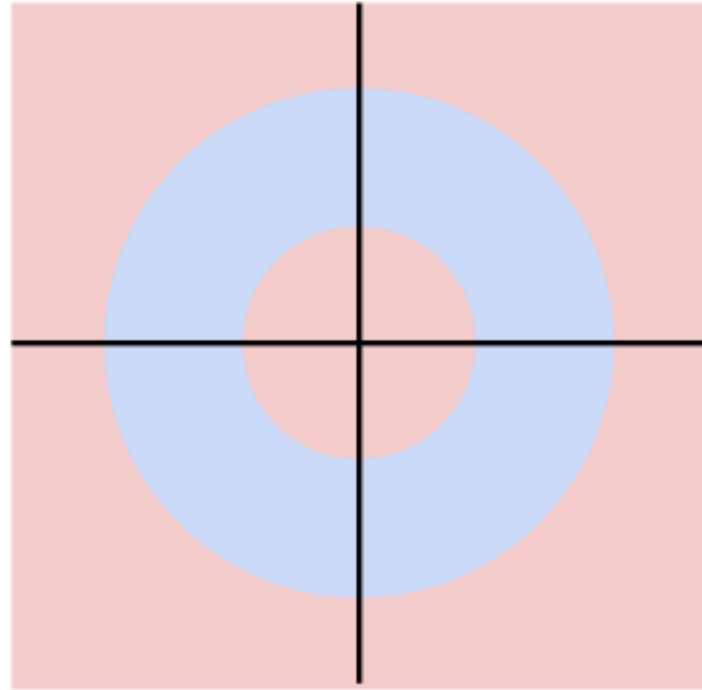


# 线性分类器无法处理的情况 (L02)



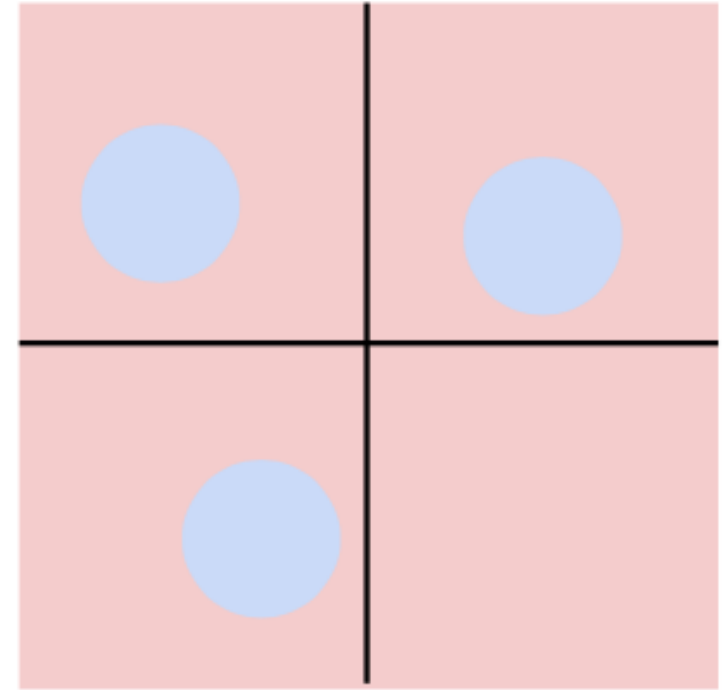
蓝色:  $x_1 \times x_2 < 0$

红色:  $x_1 \times x_2 > 0$



蓝色:  $1 \leq L2 \text{ norm} < 2$

红色: 其他点



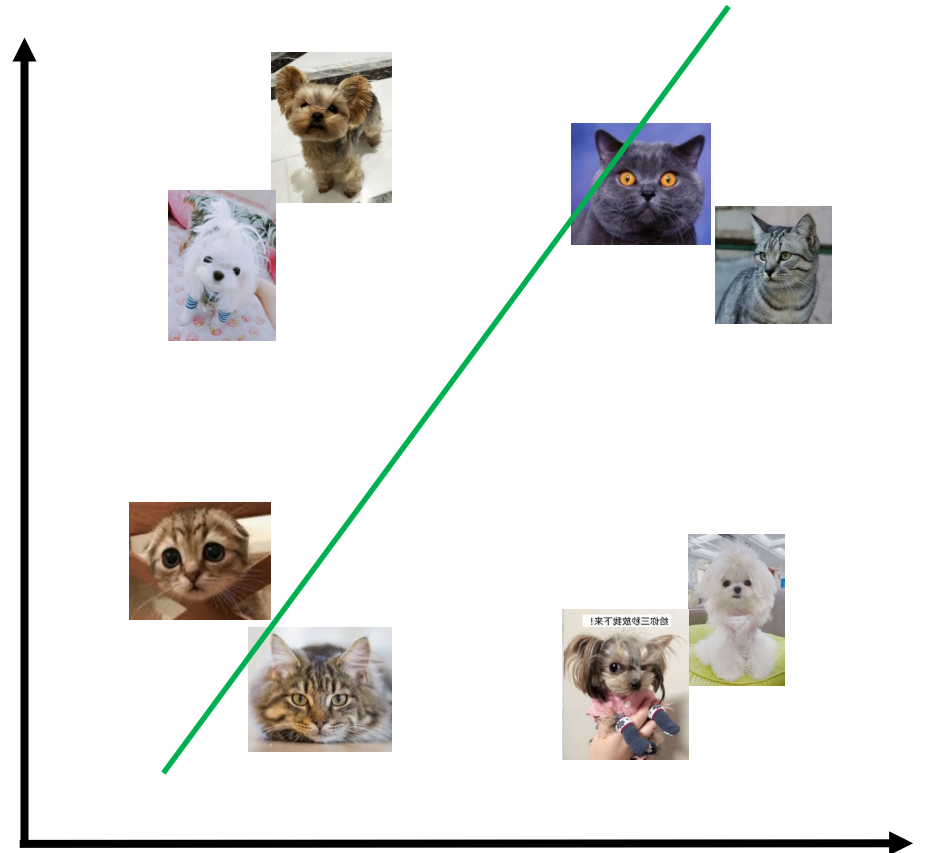
蓝色: 三个离散区域

红色: 其他点

# 线性分类器无法处理的情况 (L02)

$$\begin{array}{|c|c|c|c|} \hline w_1 & w_2 & w_3 & w_4 \\ \hline \end{array}
 + 
 \begin{array}{|c|} \hline x_1 \\ \hline x_2 \\ \hline x_3 \\ \hline x_4 \\ \hline \end{array}
 + 
 \begin{array}{|c|} \hline b \\ \hline \end{array}
 = 
 \begin{array}{|c|} \hline \text{score}(f(W,x)) \\ \hline \end{array}$$

$$f(W, x) = w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + b$$

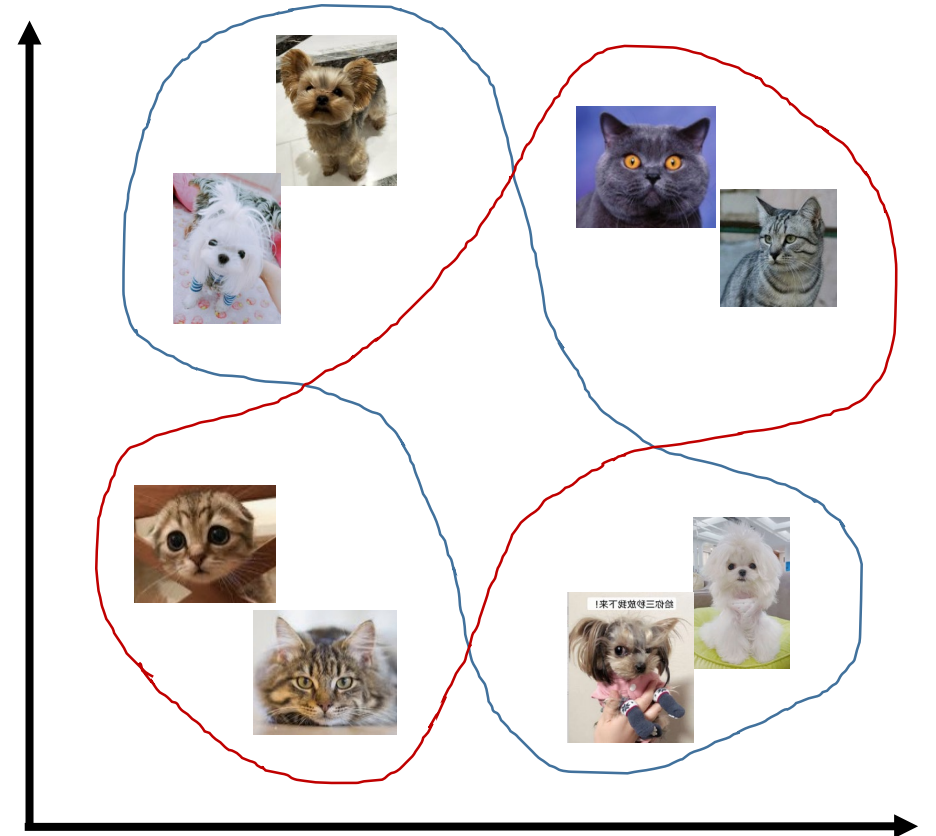


# 线性分类器无法处理的情况 (L02)

$$\begin{array}{|c|c|c|c|} \hline w_1 & w_2 & w_3 & w_4 \\ \hline \end{array}
 + 
 \begin{array}{|c|} \hline x_1 \\ \hline x_2 \\ \hline x_3 \\ \hline x_4 \\ \hline \end{array}
 + 
 \begin{array}{|c|} \hline b \\ \hline \end{array}
 = 
 \begin{array}{|c|} \hline \text{score } (f(W,x)) \\ \hline \end{array}$$

$$f(W, x) = w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + b$$

可能的方案：增加高阶多项式项

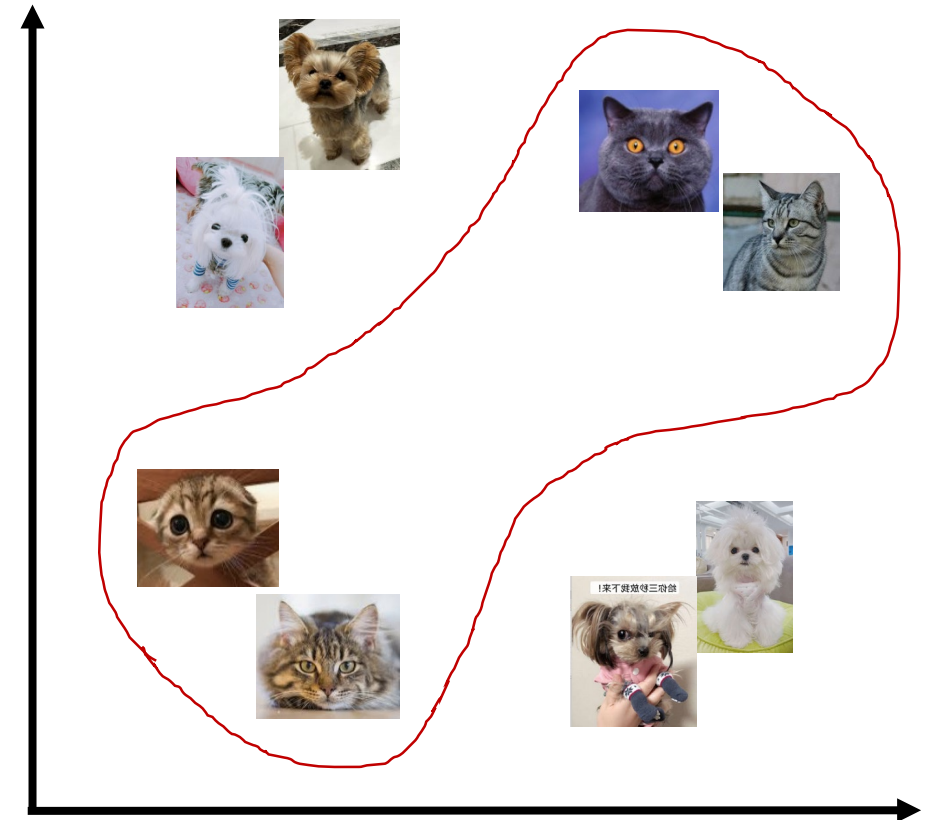


# 非线性函数

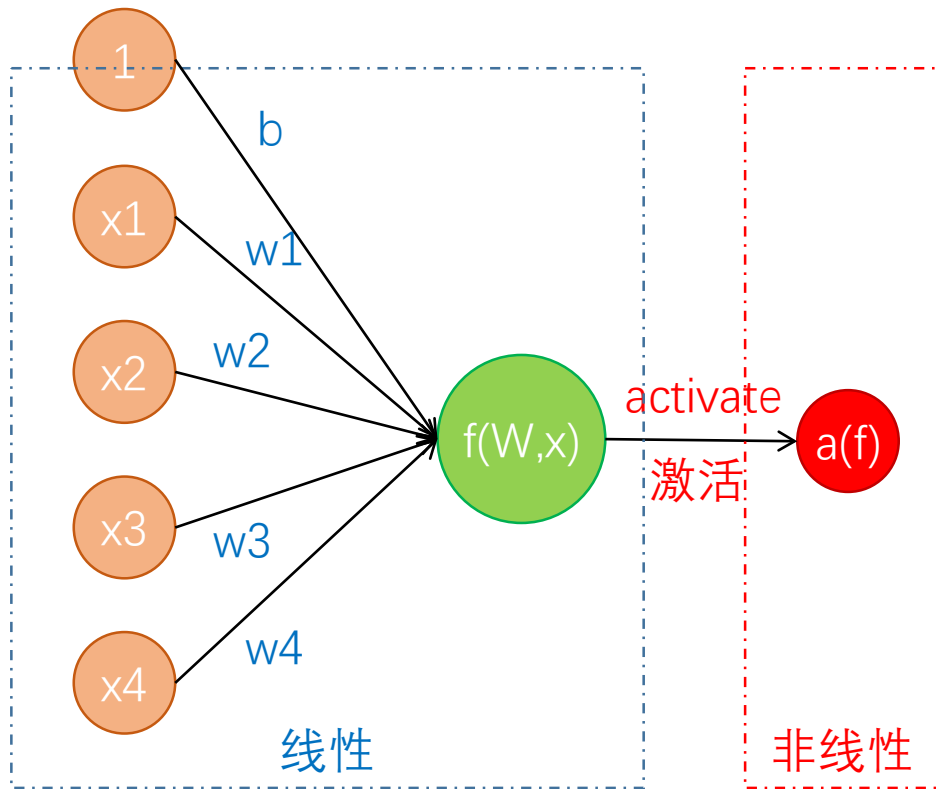
$$\begin{aligned}
 f(W, x) = & w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + b \\
 & + w_{11}x_1^2 + w_{12}x_1x_2 + w_{13}x_1x_3 + w_{14}x_1x_4 + \dots \\
 & + w_{111}x_1^3 + w_{112}x_1^2x_2 + \dots \\
 & + w_{1111}x_1^4 + \dots
 \end{aligned}$$

Q: 输入变量有n个, 最多有多少参数?

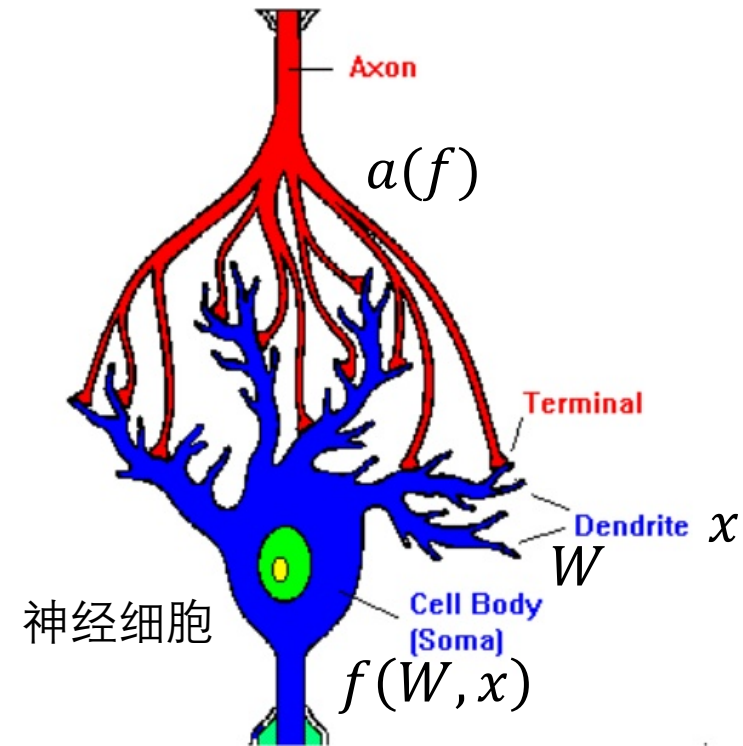
A:  $O(n^1) + O(n^2) + O(n^3) + \dots O(n^n)!!!$



# 神经网络 (Neural Network)



$a$ 称为激活函数  
(activation function)



例如Sigmoid:  $\sigma(x) = \frac{1}{1+e^{-x}}$

# 激活函数 (Activation functions)

## Sigmoid

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

较为通用

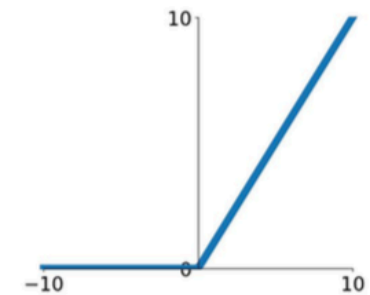
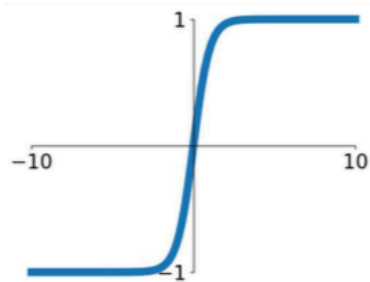
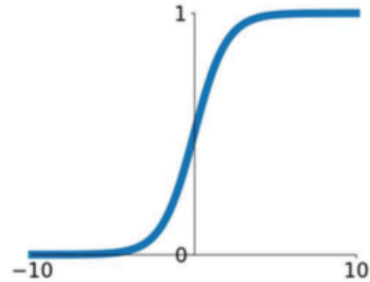
## tanh

$$\tanh(x)$$

## ReLU

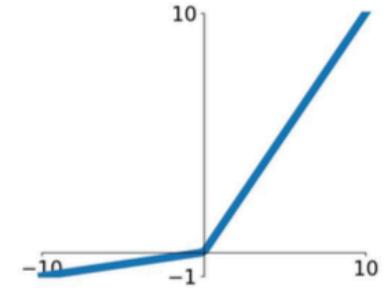
$$\max(0, x)$$

CV中最常用



## Leaky ReLU

$$\max(0.1x, x)$$

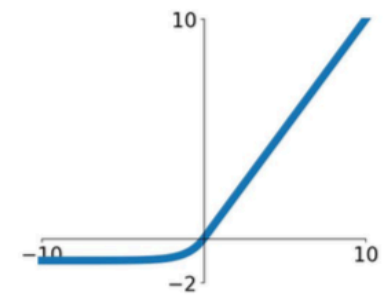


## Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

## ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$

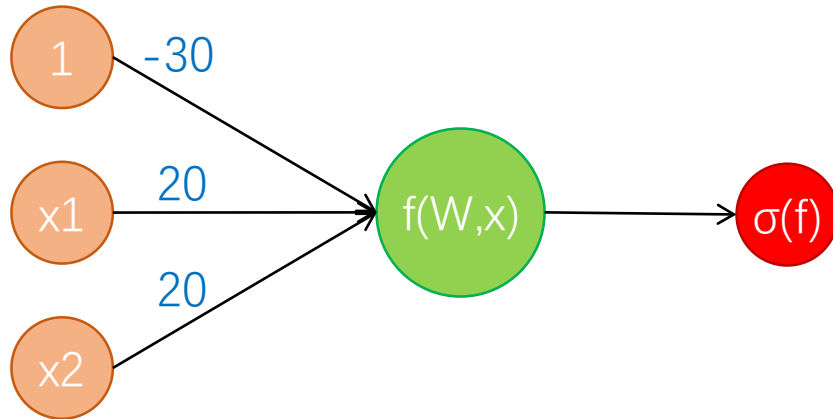
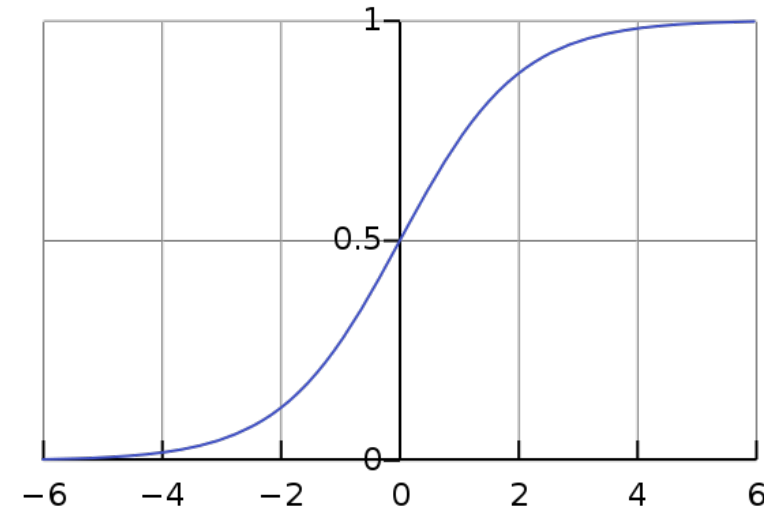
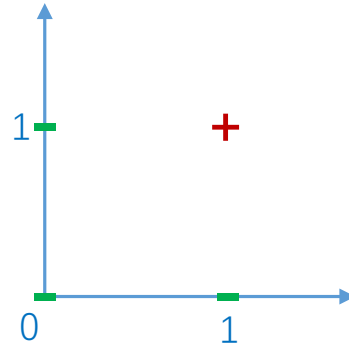




# 非线性例子

$$x_1, x_2 \in \{0,1\}$$

$$y = x_1 \text{ AND } x_2$$



$$f(W, x) = 20x_1 + 20x_2 - 30$$

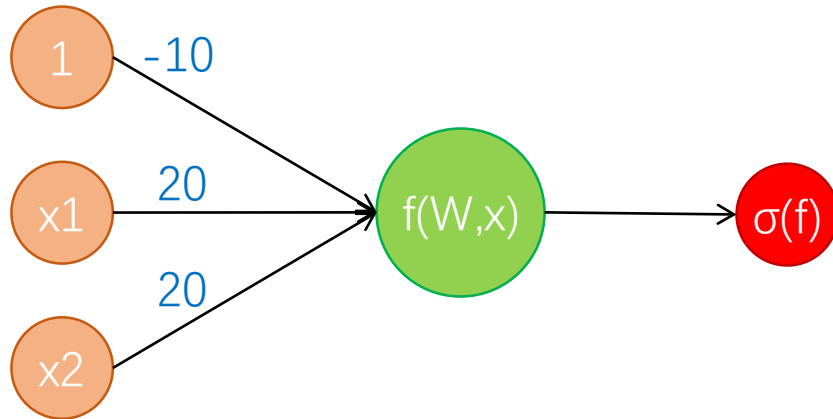
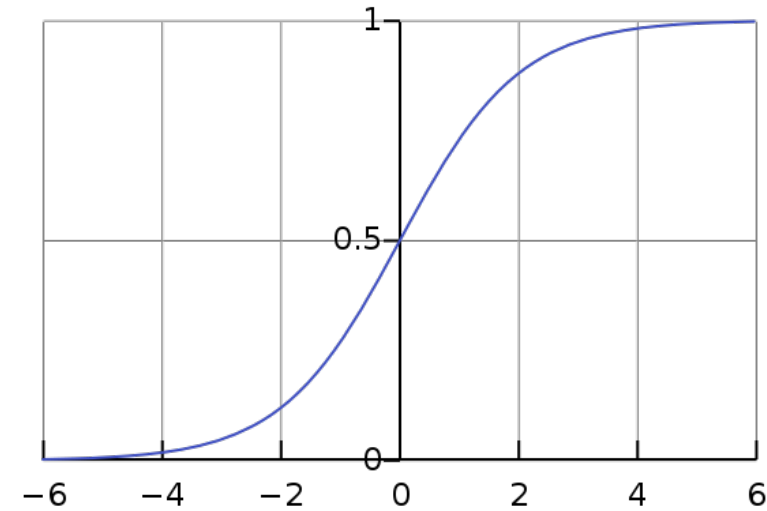
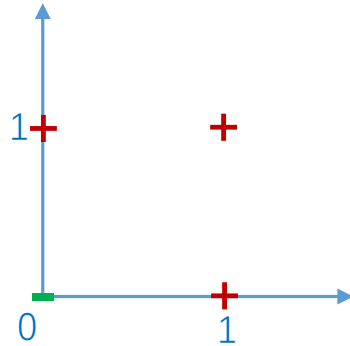
$$\sigma(f) = \frac{1}{1 + e^{-f}}$$

x1	x2	$\sigma(f(W,x))$
0	0	$\sigma(-30) \approx 0$
0	1	$\sigma(-10) \approx 0$
1	0	$\sigma(-10) \approx 0$
1	1	$\sigma(10) \approx 1$

# 非线性例子

$$x_1, x_2 \in \{0,1\}$$

$$y = x_1 \text{ OR } x_2$$

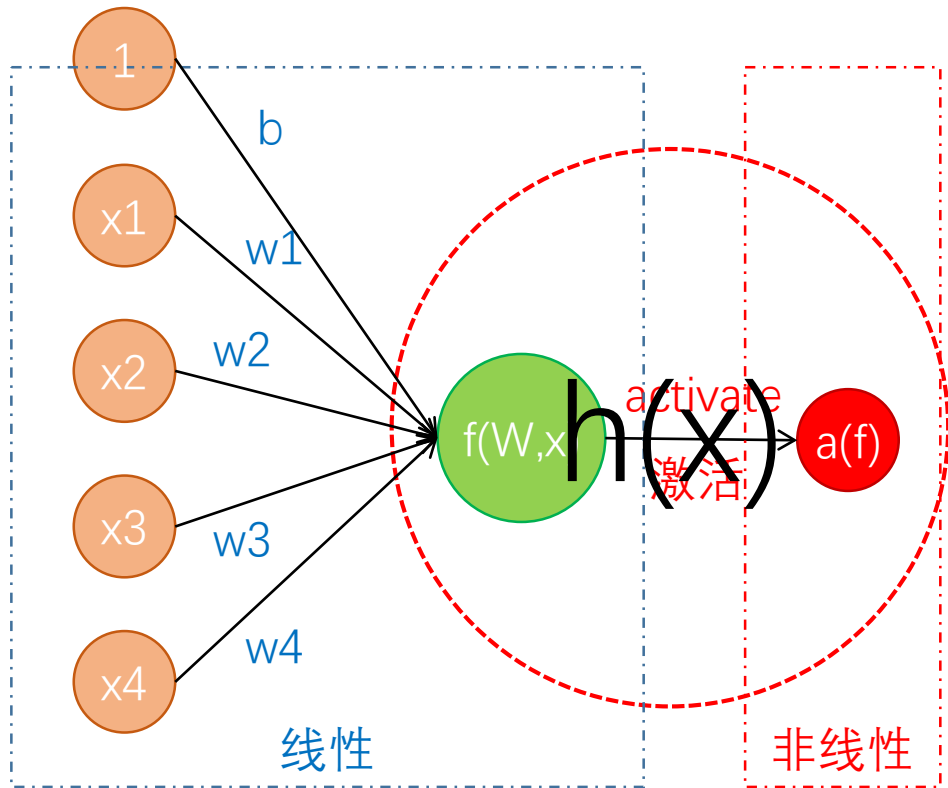


$$f(W, x) = 20x_1 + 20x_2 - 10$$

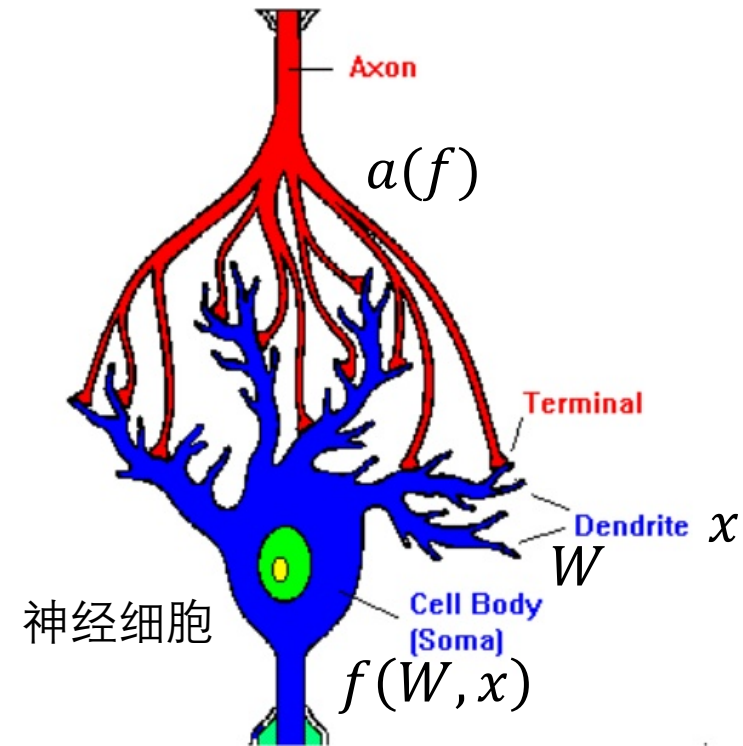
$$\sigma(f) = \frac{1}{1 + e^{-f}}$$

x1	x2	$\sigma(f(W,x))$
0	0	$\sigma(-10) \approx 0$
0	1	$\sigma(10) \approx 1$
1	0	$\sigma(10) \approx 1$
1	1	$\sigma(30) \approx 1$

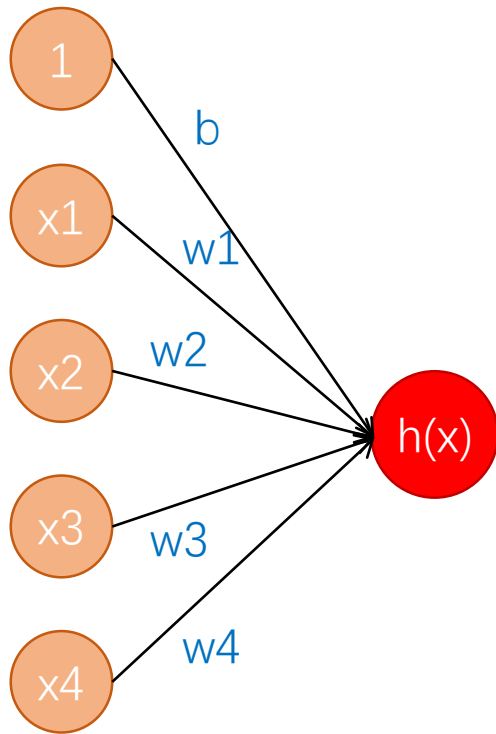
# 神经网络



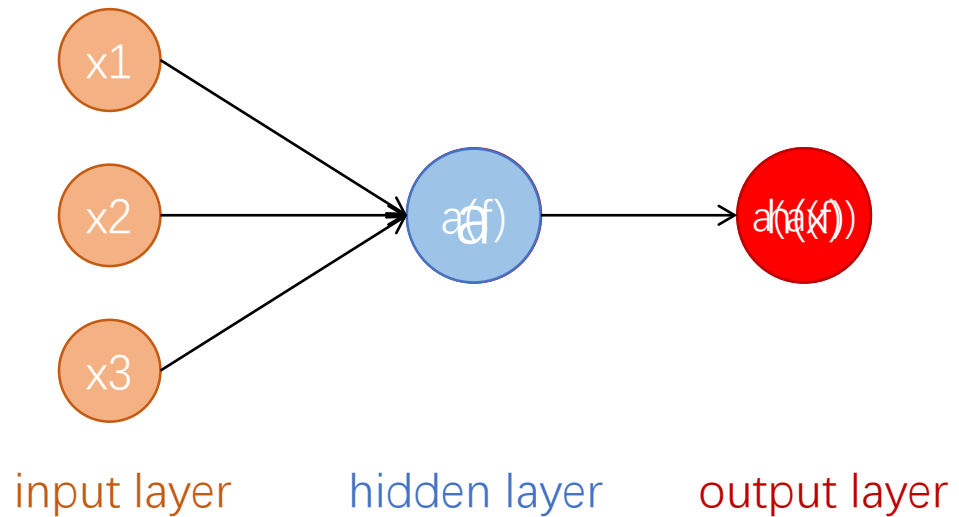
a称为激活函数  
(activation function)



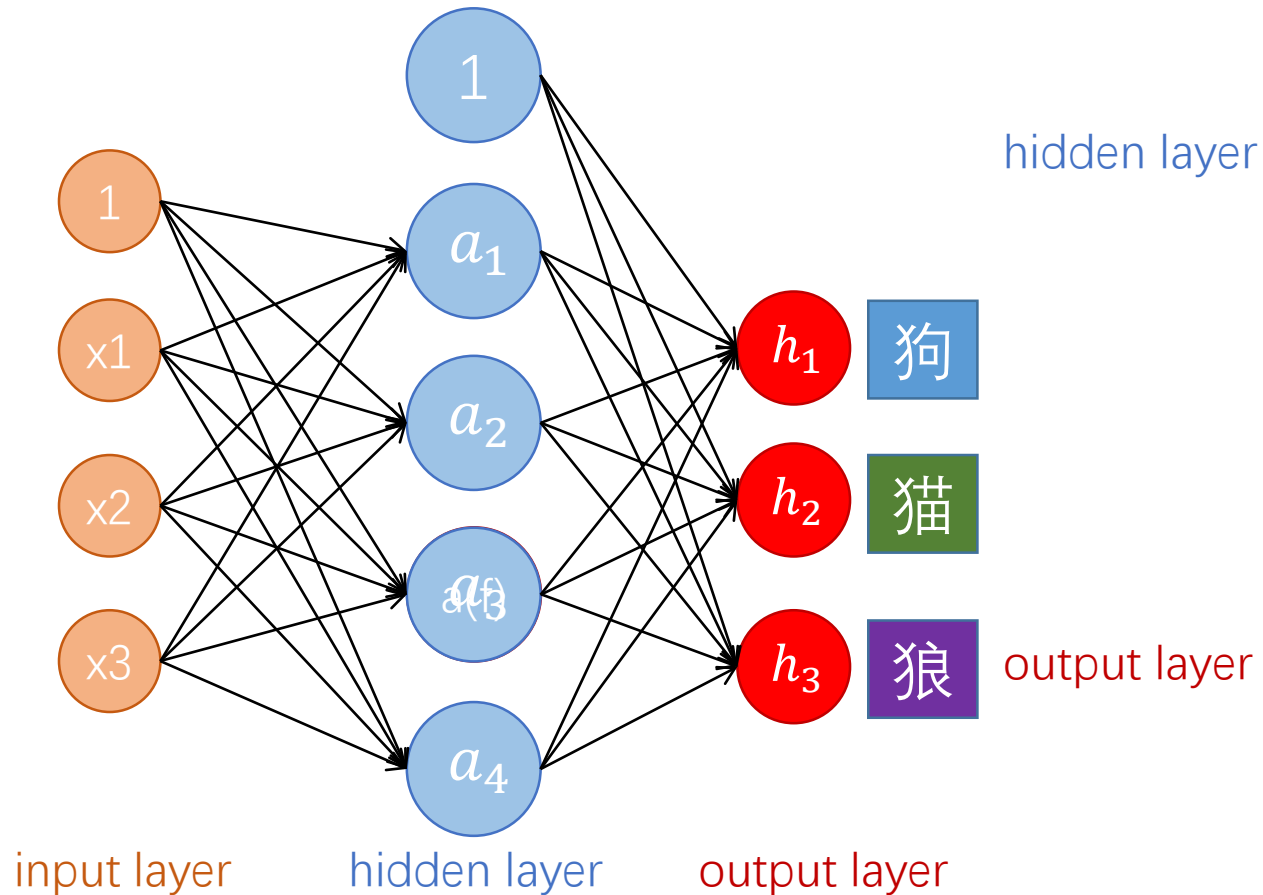
# 神经网络



# 神经网络: hidden layer



# 神经网络



2-layer neural network

假设使用Sigmoid激活

$$a_1 = \sigma(w_{11}x_1 + w_{12}x_2 + w_{13}x_3 + b_1)$$

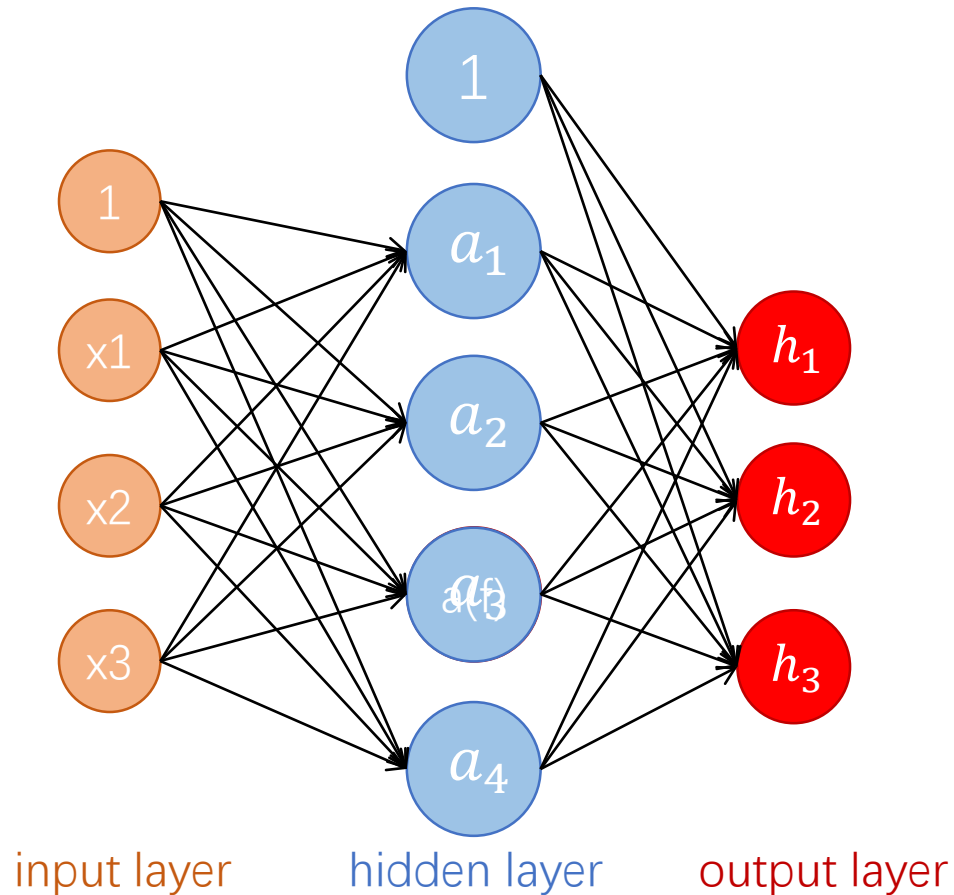
$$a_2 = \sigma(w_{21}x_1 + w_{22}x_2 + w_{23}x_3 + b_2)$$

$$a_3 = \sigma(w_{31}x_1 + w_{32}x_2 + w_{33}x_3 + b_3)$$

$$a_4 = \sigma(w_{41}x_1 + w_{42}x_2 + w_{43}x_3 + b_4)$$

$$h_1 = \sigma(w_{11}a_1 + w_{12}a_2 + w_{13}a_3 + w_{14}a_4 + b_1)$$

# 神经网络



2-layer neural network

假设使用Sigmoid激活

$$a_1 = \sigma(w_{11}^1 x_1 + w_{12}^1 x_2 + w_{13}^1 x_3 + b_1^1)$$

$$a_2 = \sigma(w_{21}^1 x_1 + w_{22}^1 x_2 + w_{23}^1 x_3 + b_2^1)$$

$$a_3 = \sigma(w_{31}^1 x_1 + w_{32}^1 x_2 + w_{33}^1 x_3 + b_3^1)$$

$$a_4 = \sigma(w_{41}^1 x_1 + w_{42}^1 x_2 + w_{43}^1 x_3 + b_4^1)$$

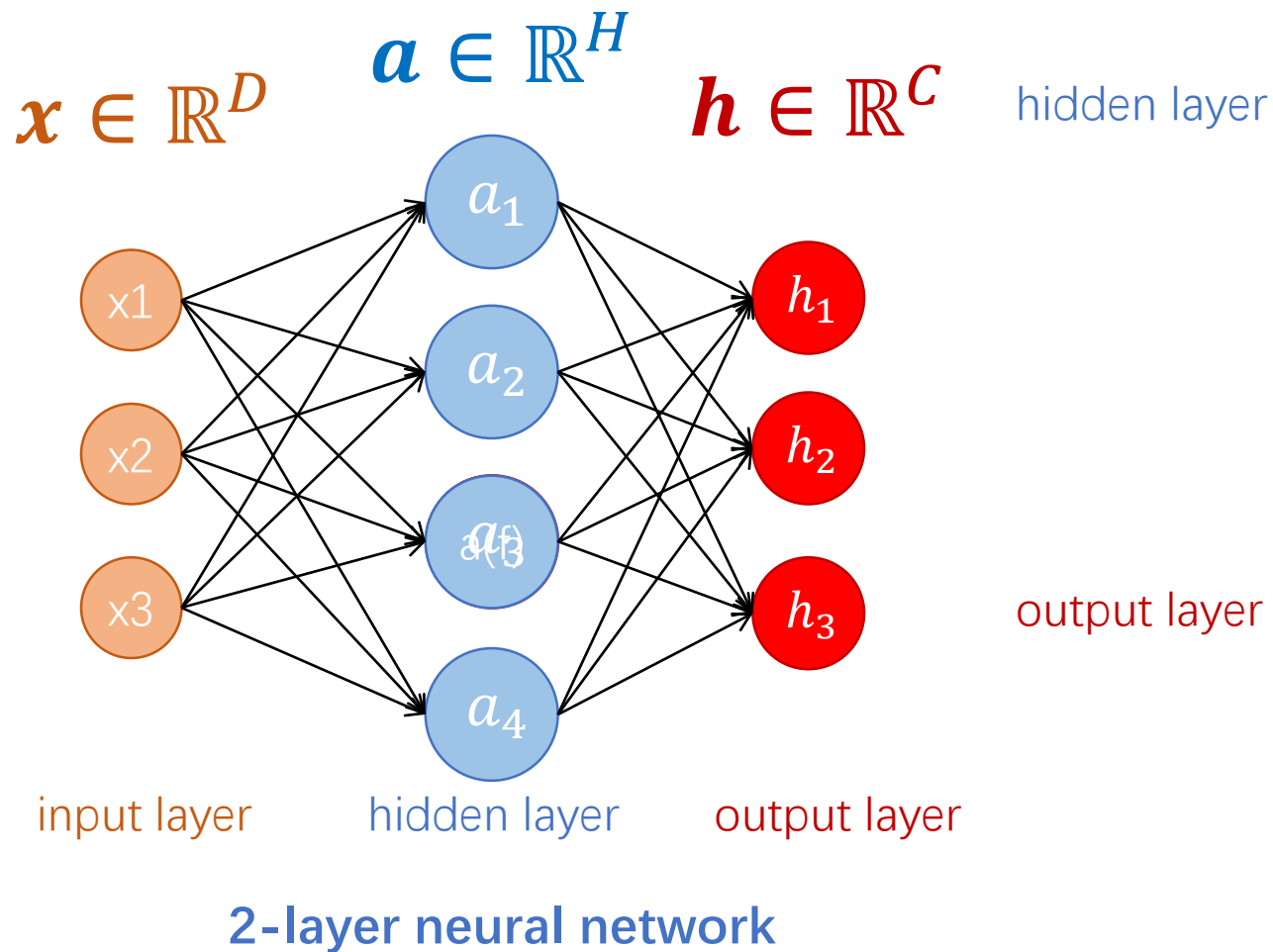
hidden layer

$$h_1 = \sigma(w_{11}^2 a_1 + w_{12}^2 a_2 + w_{13}^2 a_3 + w_{14}^2 a_4 + b_1^2)$$

$$\text{output layer } h_2 = \sigma(w_{21}^2 a_1 + w_{22}^2 a_2 + w_{23}^2 a_3 + w_{24}^2 a_4 + b_2^2)$$

$$h_3 = \sigma(w_{31}^2 a_1 + w_{32}^2 a_2 + w_{33}^2 a_3 + w_{34}^2 a_4 + b_3^2)$$

# 神经网络



假设使用Sigmoid激活

$$a_1 = \sigma(w_{11}^1 x_1 + w_{12}^1 x_2 + w_{13}^1 x_3 + b_1^1)$$

$$a_2 = \sigma(w_{21}^1 x_1 + w_{22}^1 x_2 + w_{23}^1 x_3 + b_2^1)$$

$$a_3 = \sigma(w_{31}^1 x_1 + w_{32}^1 x_2 + w_{33}^1 x_3 + b_3^1)$$

$$a_4 = \sigma(w_{41}^1 x_1 + w_{42}^1 x_2 + w_{43}^1 x_3 + b_4^1)$$

$$W^1 \in \mathbb{R}^{H \times D}$$

$$h_1 = \sigma(w_{11}^2 a_1 + w_{12}^2 a_2 + w_{13}^2 a_3 + w_{14}^2 a_4 + b_1^2)$$

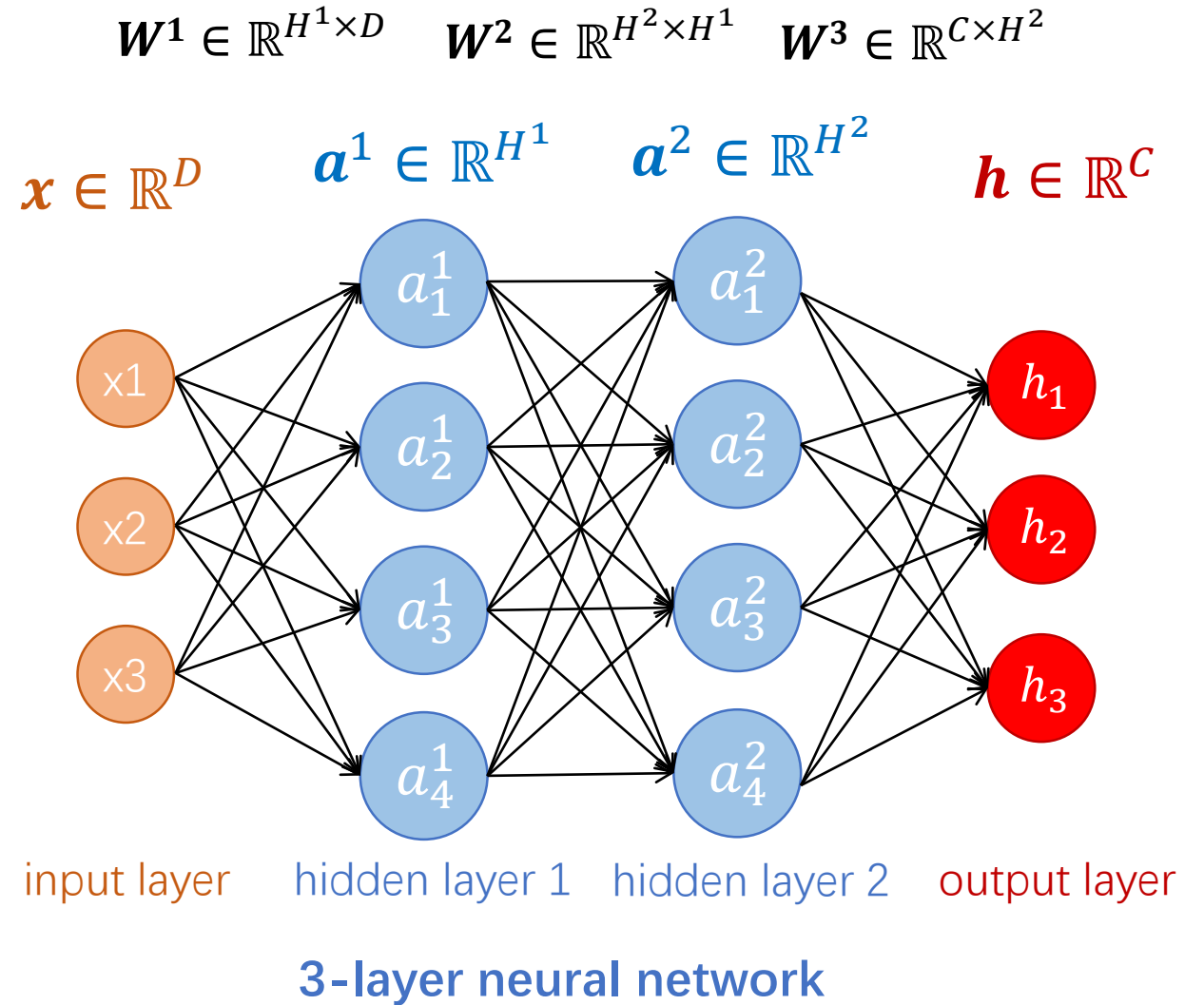
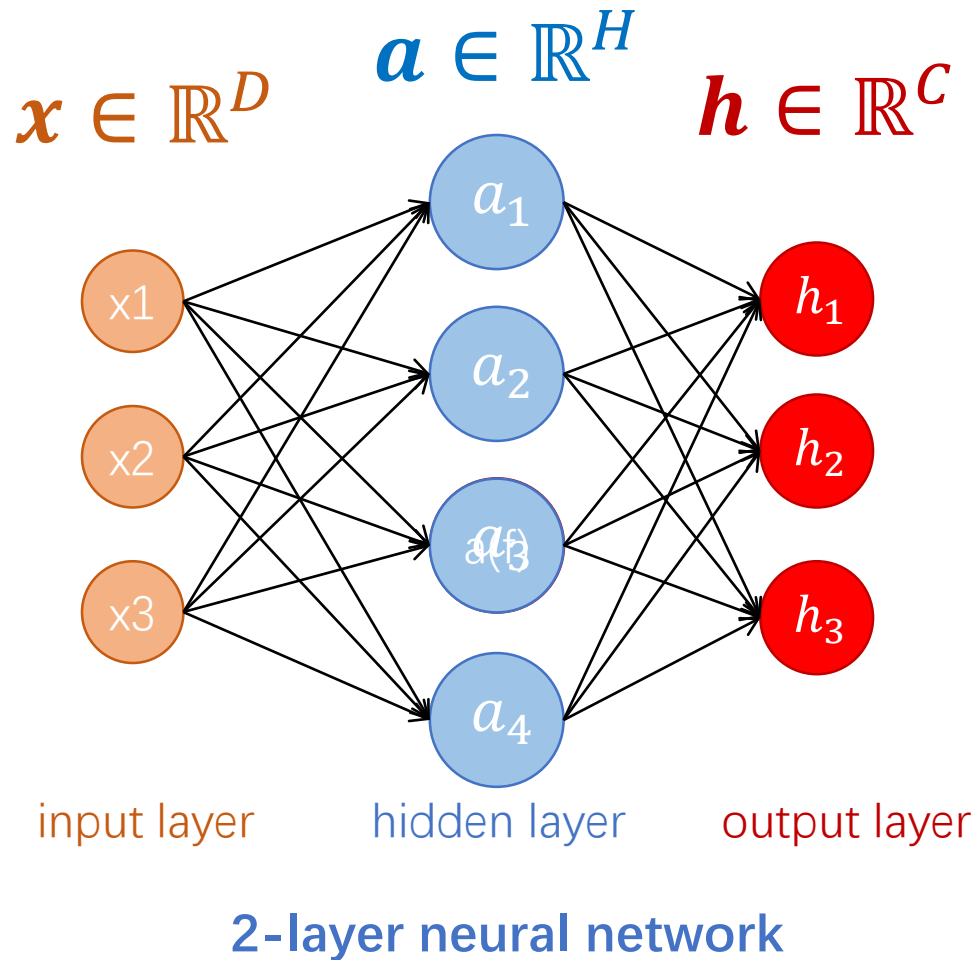
output layer    $h_2 = \sigma(w_{21}^2 a_1 + w_{22}^2 a_2 + w_{23}^2 a_3 + w_{24}^2 a_4 + b_2^2)$

$$h_3 = \sigma(w_{31}^2 a_1 + w_{32}^2 a_2 + w_{33}^2 a_3 + w_{34}^2 a_4 + b_3^2)$$

$$W^2 \in \mathbb{R}^{C \times H}$$

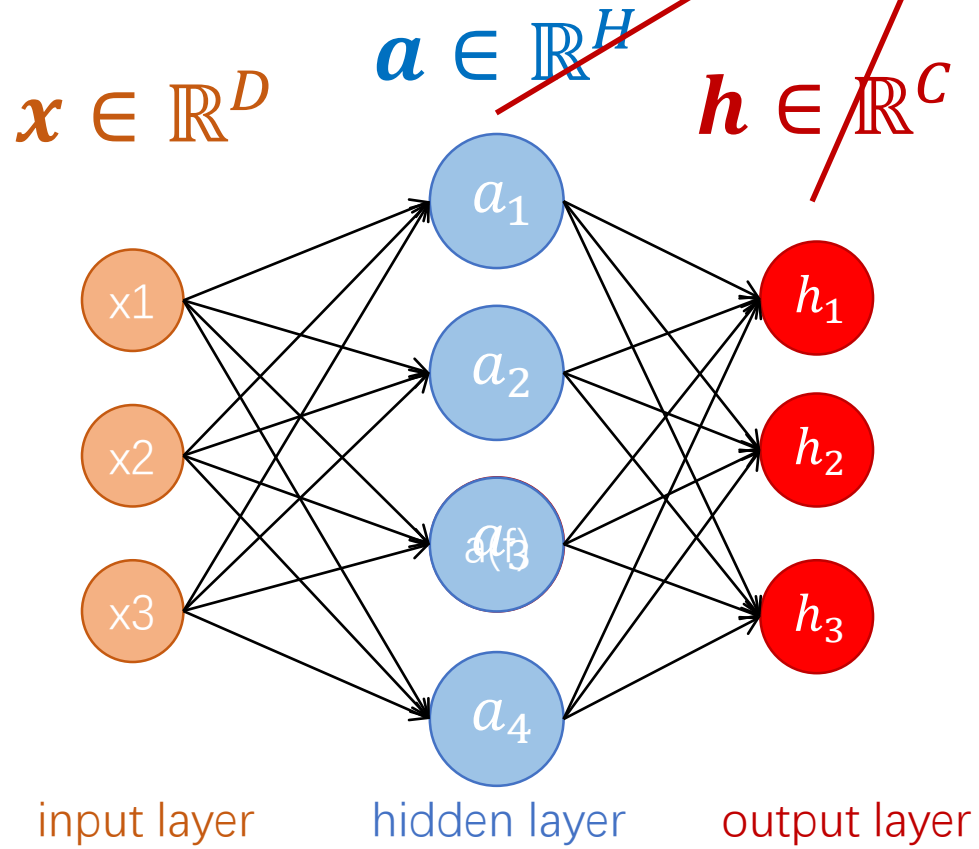


# 神经网络

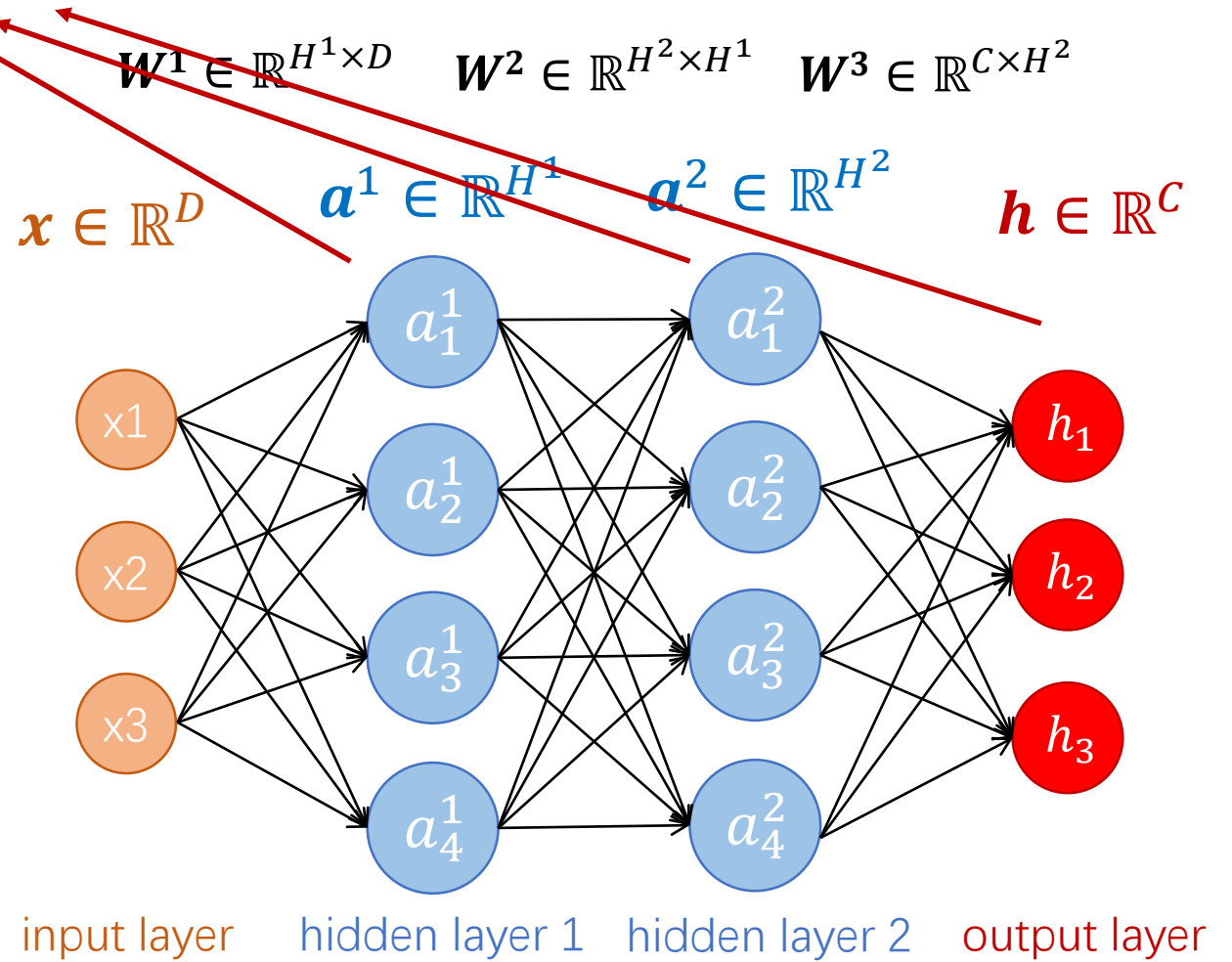


# 神经网络

全连接层



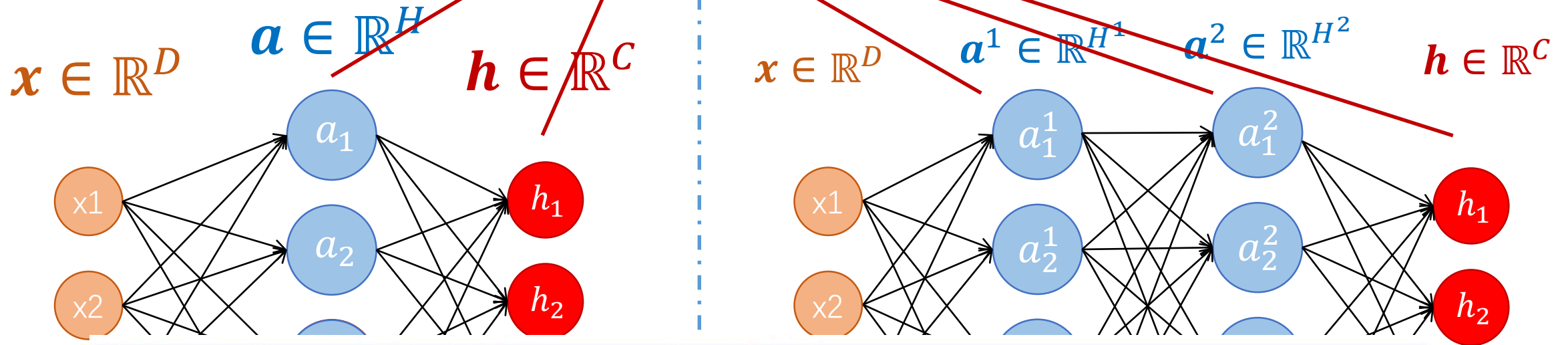
2-layer neural network



3-layer neural network

# 神经网络

全连接层



# forward-pass of a 3-layer neural network:

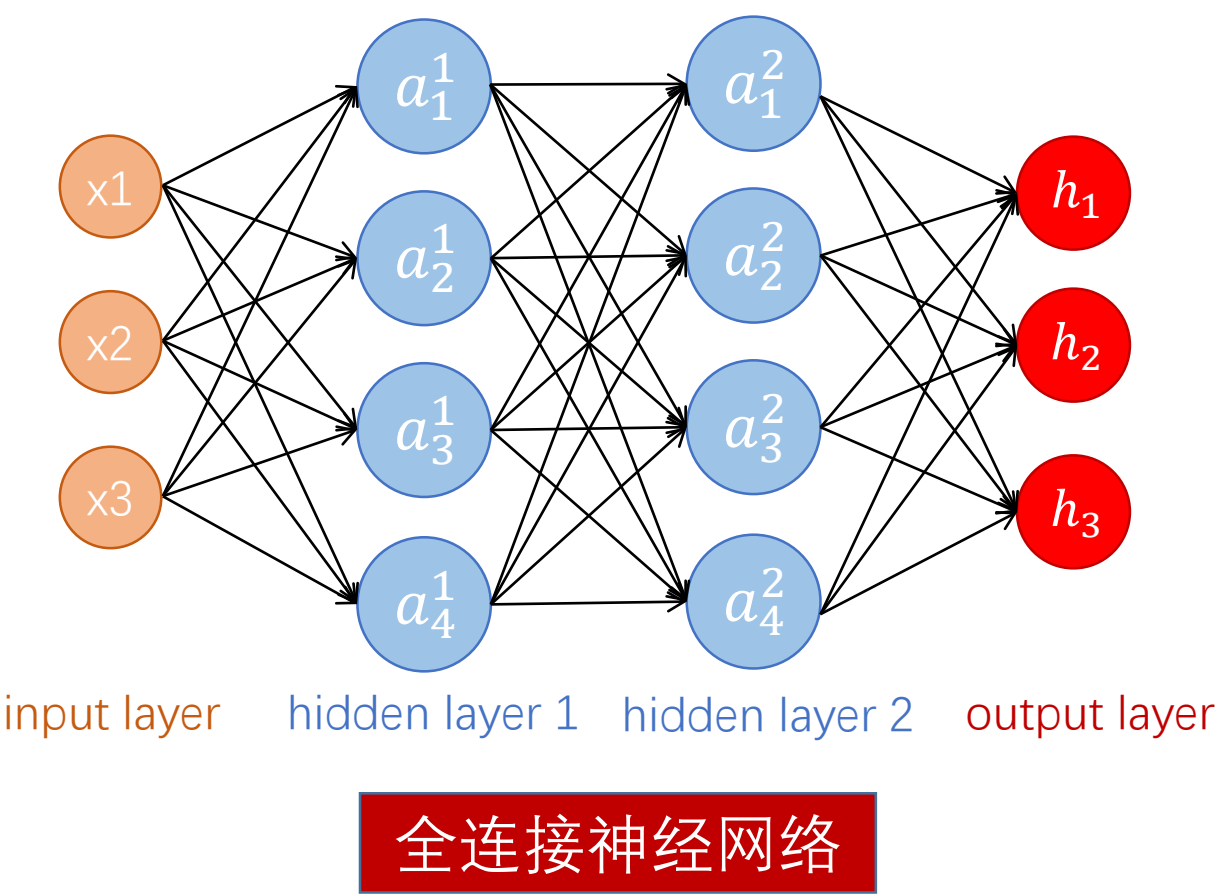
```
f = lambda x: 1.0/(1.0 + np.exp(-x)) # activation function (use sigmoid)
x = np.random.randn(3, 1) # random input vector of three numbers (3x1)
h1 = f(np.dot(W1, x) + b1) # calculate first hidden layer activations (4x1)
h2 = f(np.dot(W2, h1) + b2) # calculate second hidden layer activations (4x1)
out = np.dot(W3, h2) + b3 # output neuron (1x1)
```

2-layer neural network

3-layer neural network



# 随机连接的神经网络

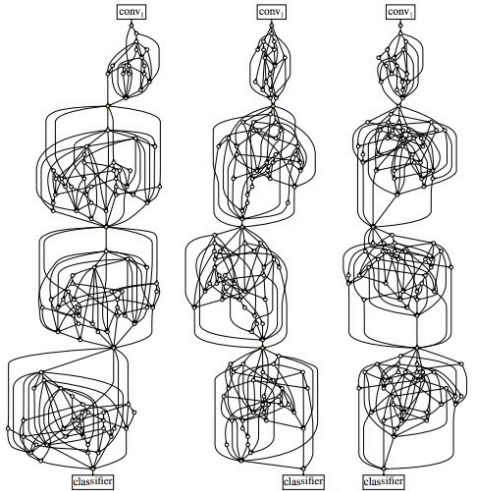


## Exploring Randomly Wired Neural Networks for Image Recognition

Saining Xie    Alexander Kirillov    Ross Girshick    Kaiming He  
Facebook AI Research (FAIR)

### Abstract

Neural networks for image recognition have evolved through extensive manual design from simple chain-like models to structures with multiple wiring paths. The success of ResNets [11] and DenseNets [16] is due in large part to their innovative wiring plans. Now, neural architecture search (NAS) studies are exploring the joint optimization of wiring and operation types, however, the space of possible wirings is constrained and still driven by manual design despite being searched. In this paper, we explore a more diverse set of connectivity patterns through the lens of randomly wired neural networks. To do this, we first define the concept of a stochastic network generator that encapsulates the entire network generation process. Encapsulation provides a unified view of NAS and randomly wired networks. Then, we use three classical random graph models to generate randomly wired graphs for networks. The results are surprising: several variants of these random generators yield network instances that have competitive accuracy on the ImageNet benchmark. These results suggest



This image is CC0 Public Domain

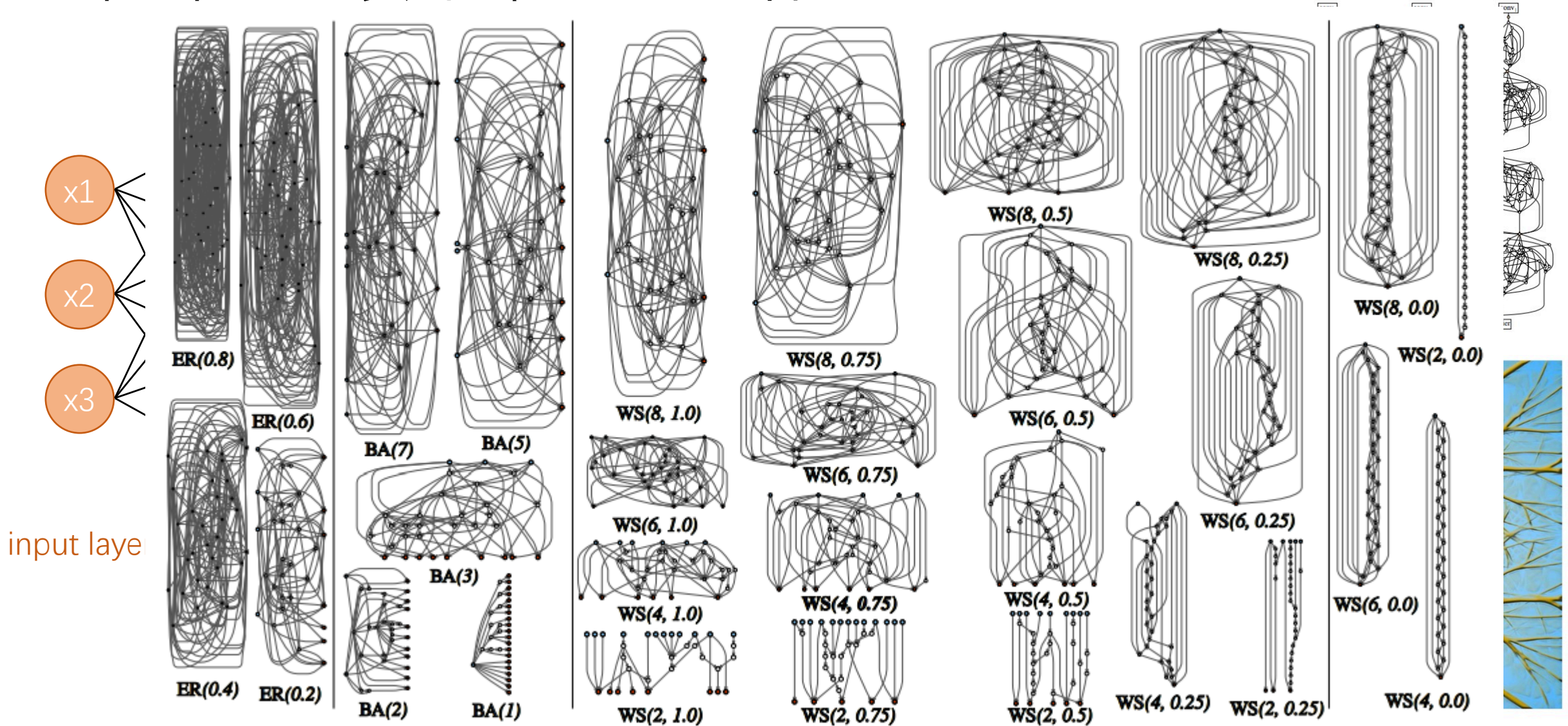


## Exploring Randomly Wired Neural Networks for Image Recognition

Saining Xie Alexander Kirillov Ross Girshick Kaiming He

Facebook AI Research (FAIR)

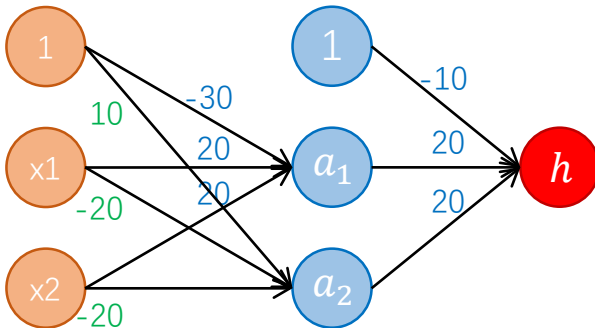
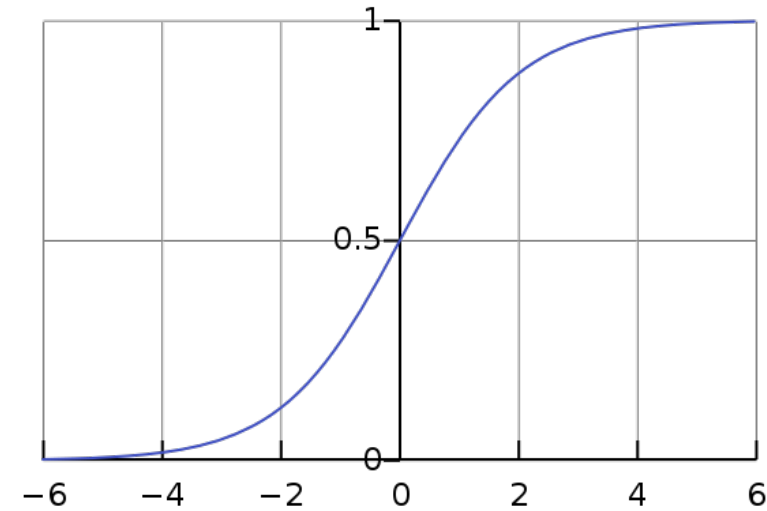
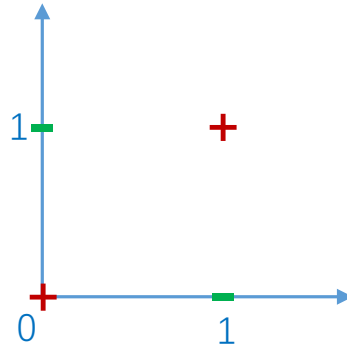
# 随机连接的神经网络



# 两层神经网络例子

$$x_1, x_2 \in \{0,1\}$$

$$y = x_1 \text{ XNOR } x_2$$

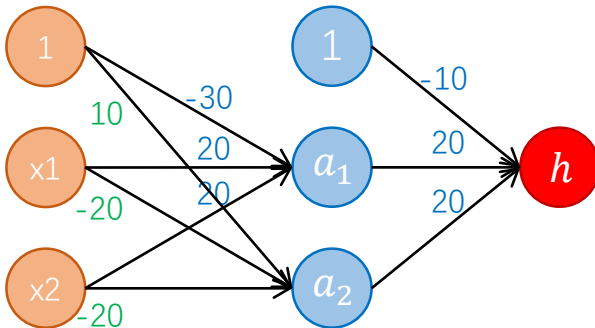
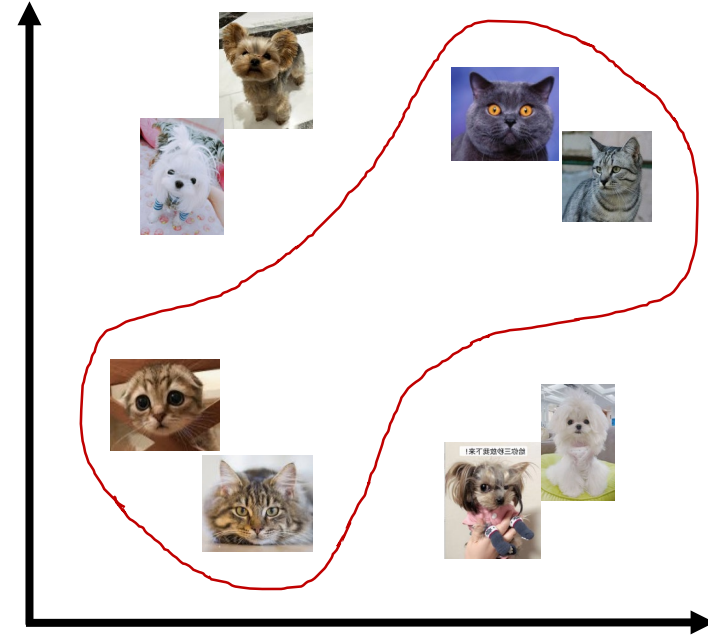
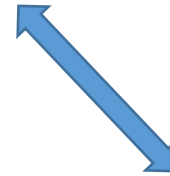
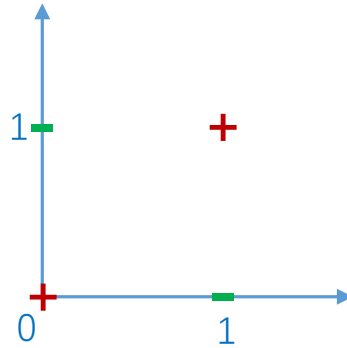


$x_1$	$x_2$	$a_1$	$a_2$	$h$
0	0	$\sigma(-30) \approx 0$	$\sigma(10) \approx 1$	$\sigma(10) \approx 1$
0	1	$\sigma(-10) \approx 0$	$\sigma(-10) \approx 0$	$\sigma(-10) \approx 0$
1	0	$\sigma(-10) \approx 0$	$\sigma(-10) \approx 0$	$\sigma(-10) \approx 0$
1	1	$\sigma(10) \approx 1$	$\sigma(-30) \approx 0$	$\sigma(10) \approx 1$

# 两层神经网络例子

$$x_1, x_2 \in \{0,1\}$$

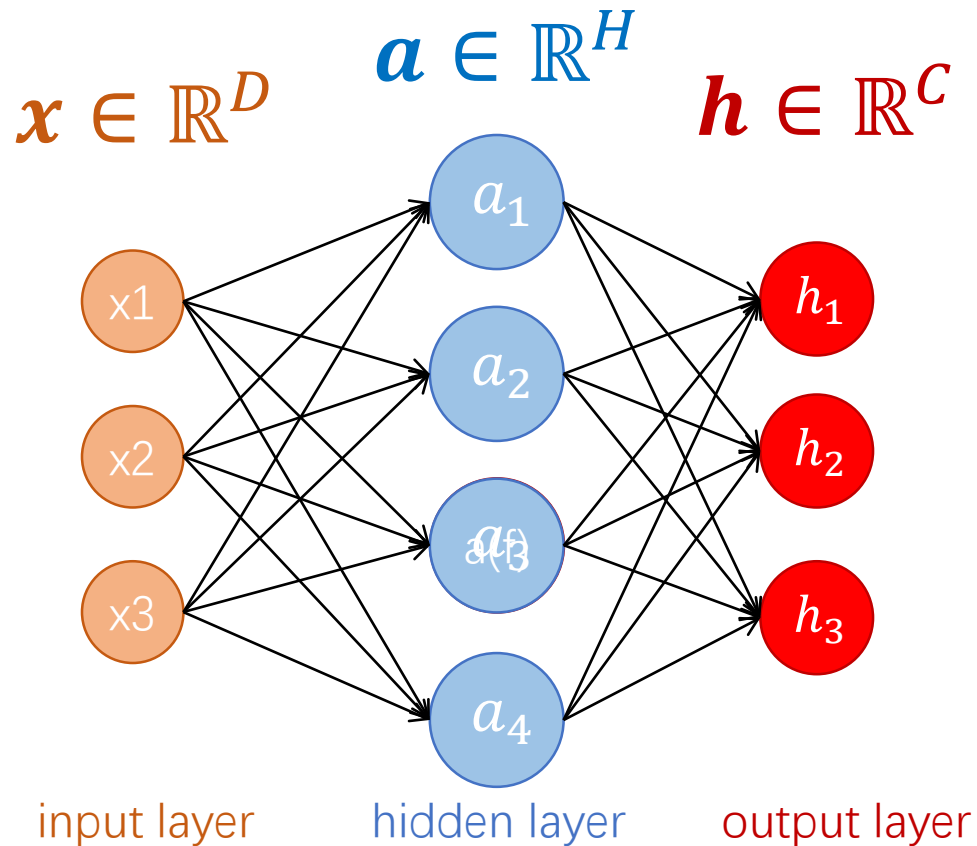
$$y = x_1 \text{ XNOR } x_2$$



x1	x2	$a_1$	$a_2$	$h$
0	0	$\sigma(-30) \approx 0$	$\sigma(10) \approx 1$	$\sigma(10) \approx 1$
0	1	$\sigma(-10) \approx 0$	$\sigma(-10) \approx 0$	$\sigma(-10) \approx 0$
1	0	$\sigma(-10) \approx 0$	$\sigma(-10) \approx 0$	$\sigma(-10) \approx 0$
1	1	$\sigma(10) \approx 1$	$\sigma(-30) \approx 0$	$\sigma(10) \approx 1$

# 神经网络计算

训练集:  $(x_i, y_i)_{i=1}^N$



2-layer neural network

ReLU  
 $\max(0, x)$

$$a = \max(0, f(x; W^1)); \quad s = f(a, W^2);$$

Softmax

$$h_k = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

Loss func

$$L = \frac{1}{N} \sum_{i=1}^N l(h_i, y_i) + \lambda R(W^1) + \lambda R(W^2)$$

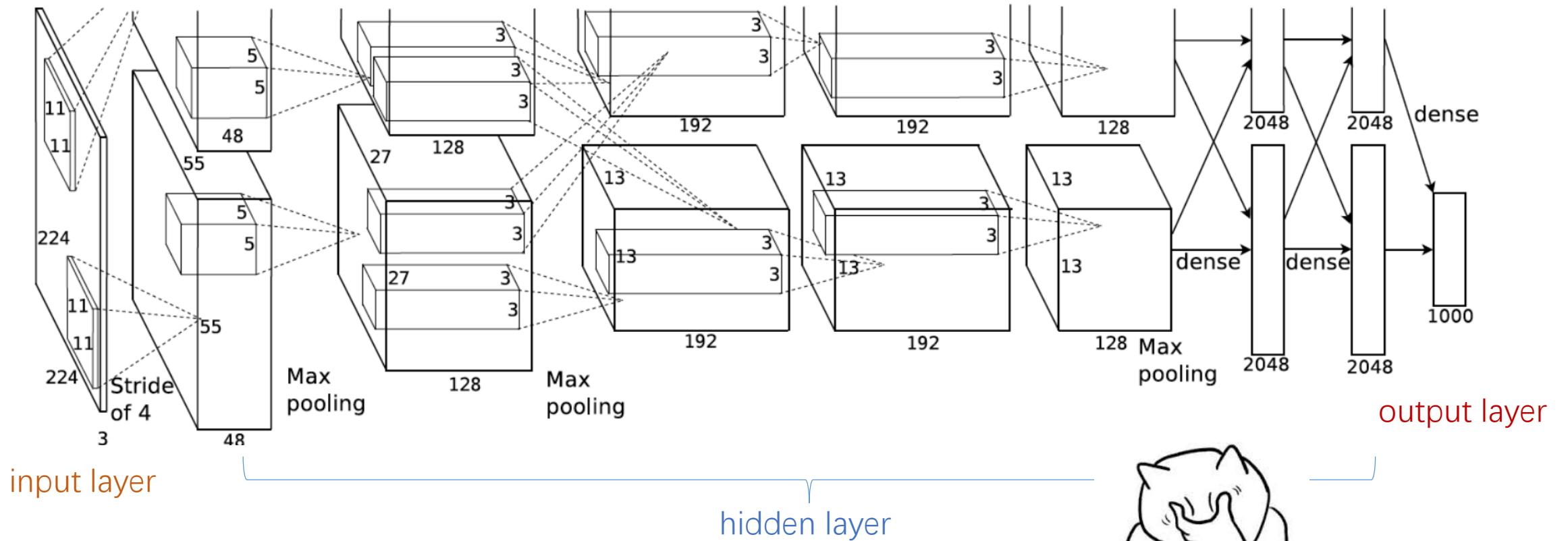
优化

计算  $\frac{\partial L}{\partial W^1}, \frac{\partial L}{\partial W^2} \quad \frac{\partial L}{\partial W_{ij}^k}, k = 1, 2$

$$\nabla_W L = \nabla_W \left[ \frac{1}{N} \sum_{i=1}^N l(h_i, y_i) + \lambda R(W^1) + \lambda R(W^2) \right]$$



# AlexNet

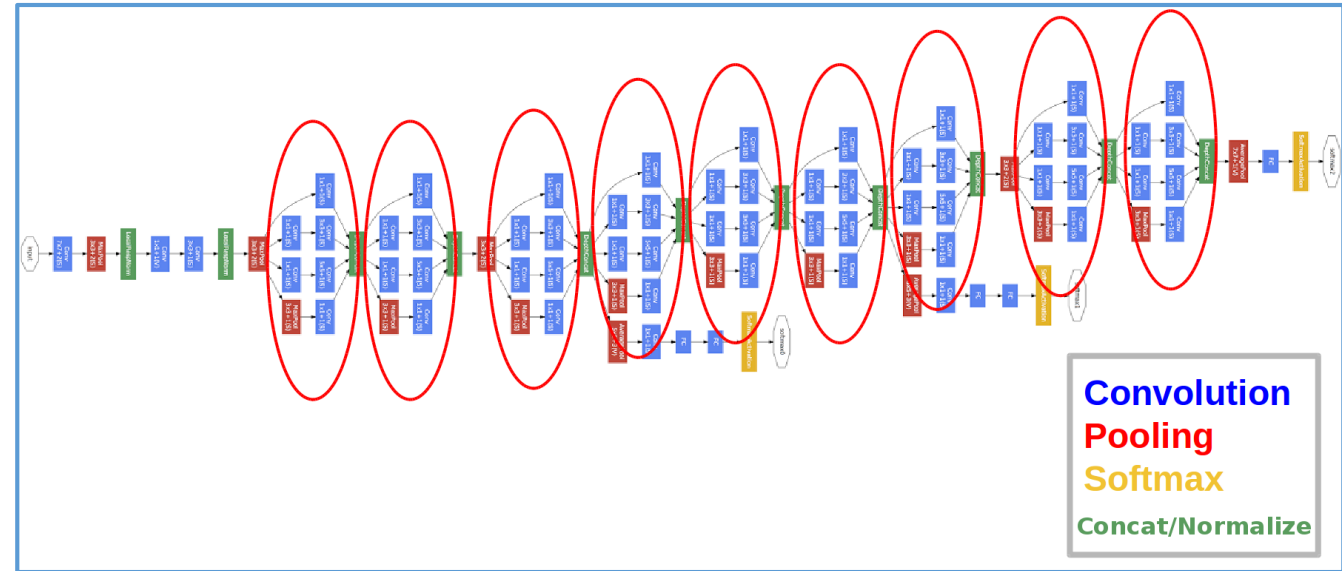
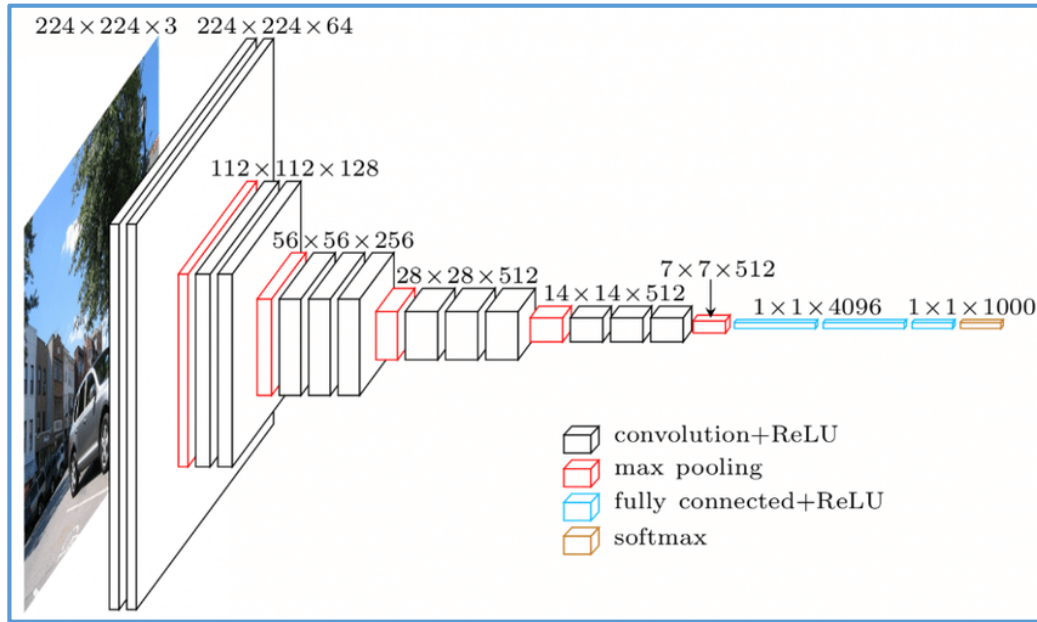


$$\nabla_W L$$



心好累

# VGG, GoogleNet



$\nabla_W L$



# 神经网络

训练集:  $(x_i, y_i)_{i=1}^N$

前向传递 (forward pass)



ReLU  
 $\max(0, x)$

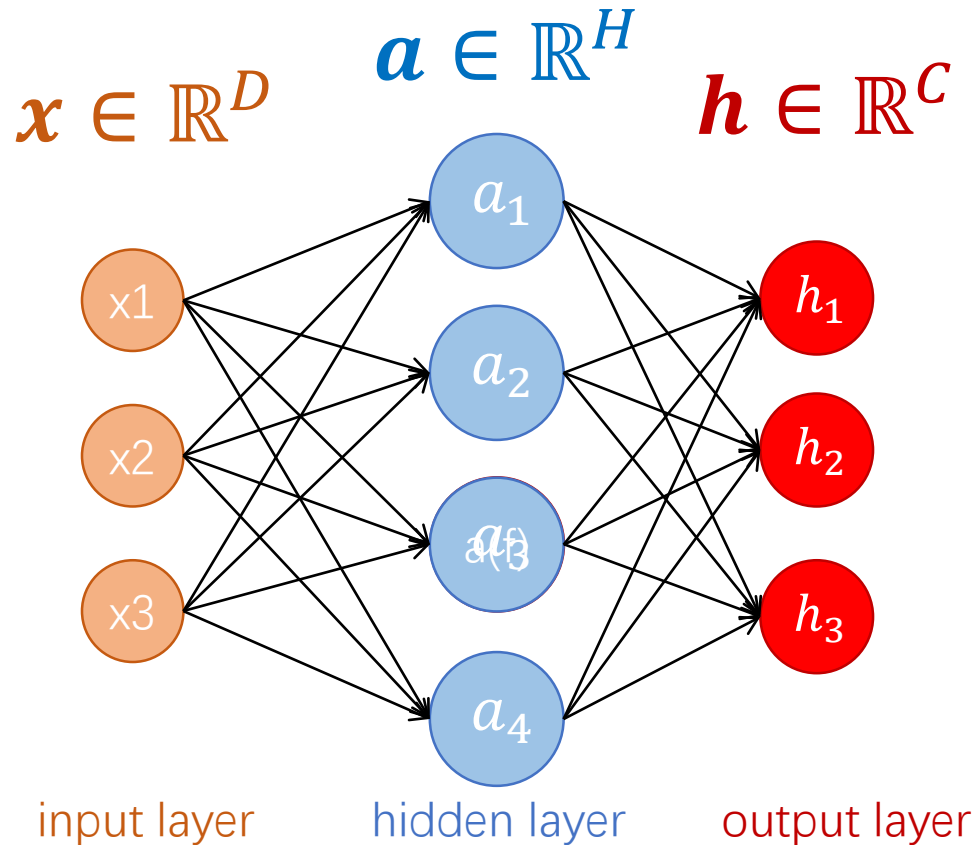
$$\mathbf{a} = \max(0, f(\mathbf{x}; \mathbf{W}^1)); \quad \mathbf{s} = f(\mathbf{a}, \mathbf{W}^2);$$

Softmax

$$h_k = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

Loss func

$$L = \frac{1}{N} \sum_{i=1}^N l(\mathbf{h}_i, y_i) + \lambda R(\mathbf{W}^1) + \lambda R(\mathbf{W}^2)$$



2-layer neural network

反向传播

(Backpropagation)

反向传递 (backward pass)

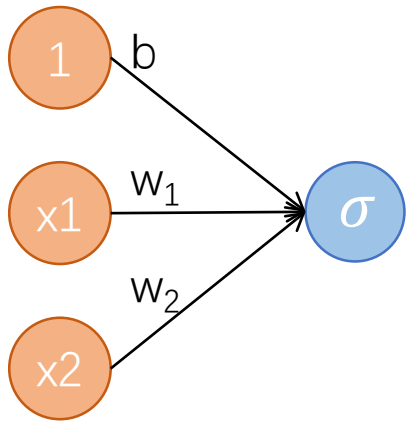


$$\frac{\partial L}{\partial \mathbf{W}^2} = \frac{\partial L}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{W}^2}, \quad \frac{\partial L}{\partial \mathbf{W}^1} = \frac{\partial L}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{a}} \frac{\partial \mathbf{a}}{\partial \mathbf{W}^1}$$

Chain rule

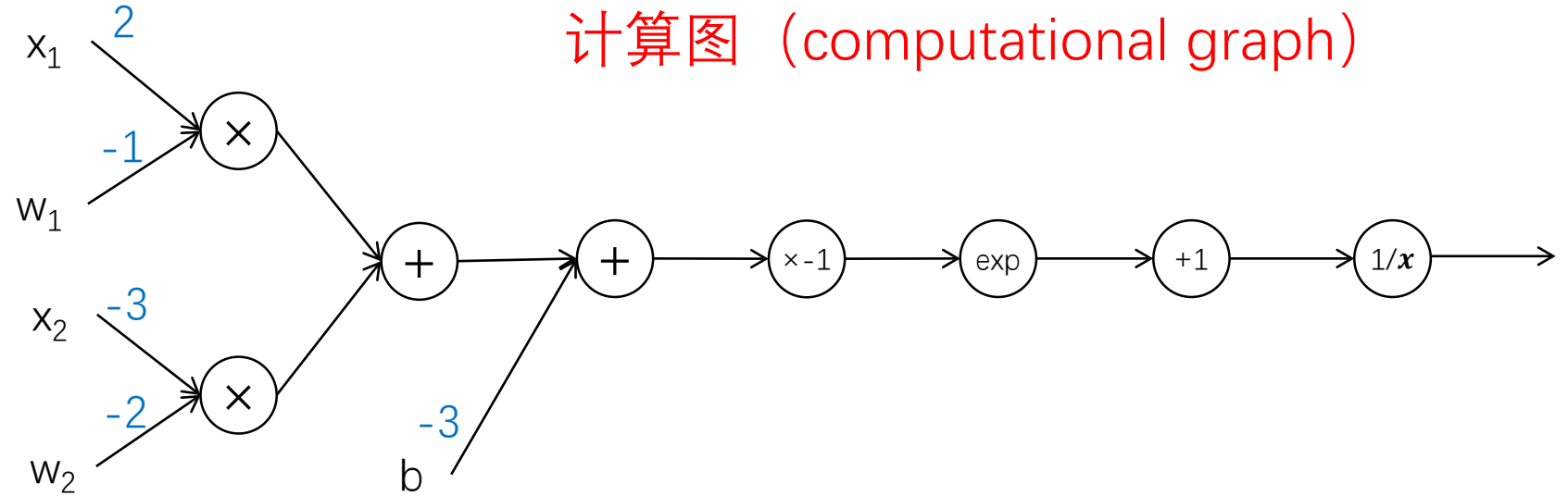
# 反向传播 (Backpropagation)

$$\sigma(W, x) = \frac{1}{1 + e^{-(w_1x_1 + w_2x_2 + b)}}$$



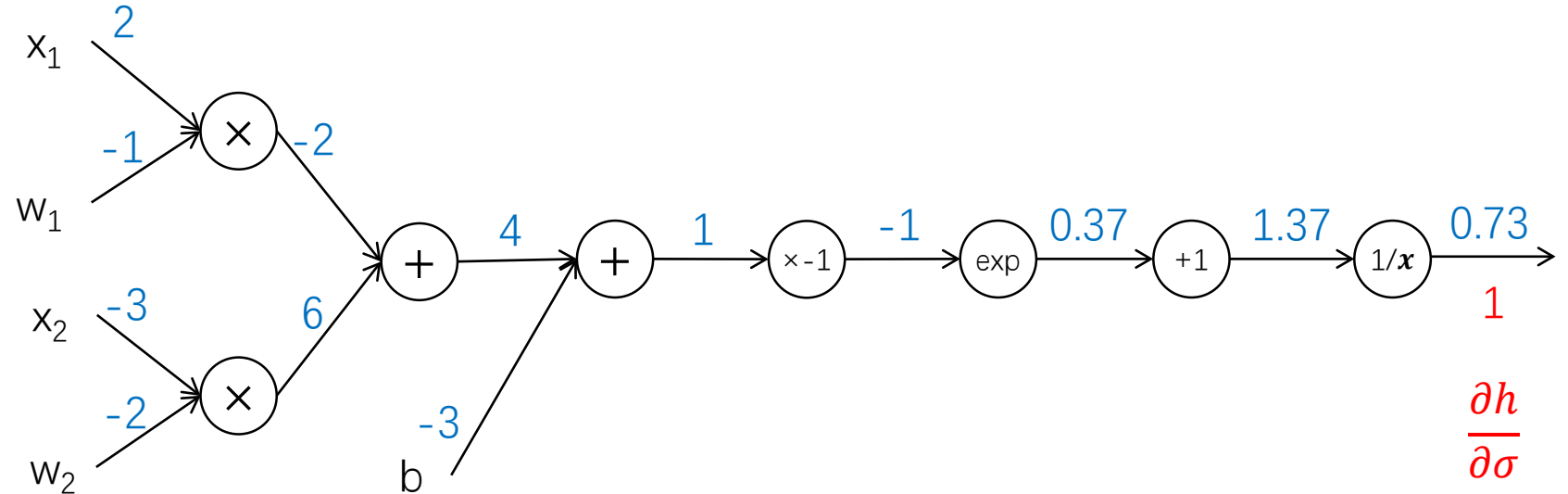
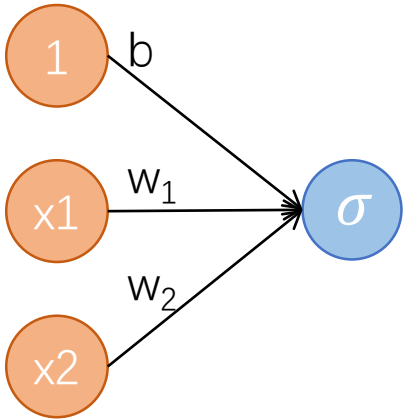
$$\frac{\partial \sigma}{\partial w_1}, \frac{\partial \sigma}{\partial w_2}, \frac{\partial \sigma}{\partial b}, \frac{\partial \sigma}{\partial x_1}, \frac{\partial \sigma}{\partial x_2}$$

计算图 (computational graph)



# 反向传播 (Backpropagation)

$$\sigma(W, x) = \frac{1}{1 + e^{-(w_1x_1 + w_2x_2 + b)}}$$

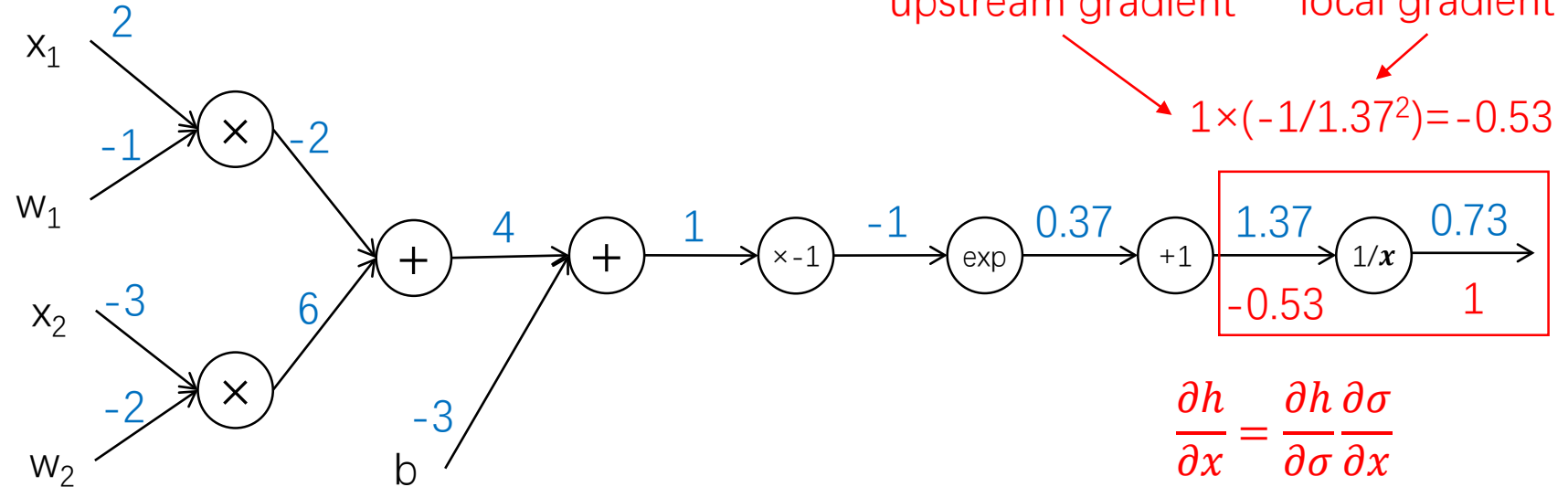
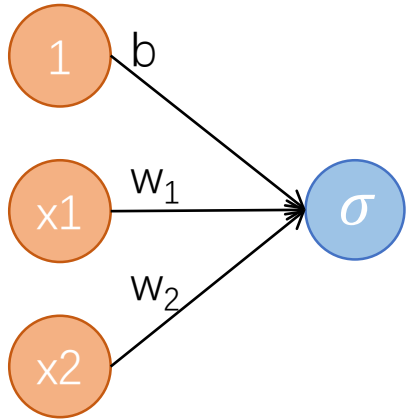


$$\frac{\partial \sigma}{\partial w_1}, \frac{\partial \sigma}{\partial w_2}, \frac{\partial \sigma}{\partial b}, \frac{\partial \sigma}{\partial x_1}, \frac{\partial \sigma}{\partial x_2}$$

$$\begin{aligned} f(x) = e^x &\longrightarrow \frac{df}{dx} = e^x & f(x) = \frac{1}{x} &\longrightarrow \frac{df}{dx} = -1/x^2 \\ f_a(x) = ax &\longrightarrow \frac{df}{dx} = a & f_c(x) = c + x &\longrightarrow \frac{df}{dx} = 1 \end{aligned}$$

# 反向传播 (Backpropagation)

$$\sigma(W, x) = \frac{1}{1 + e^{-(w_1 x_1 + w_2 x_2 + b)}}$$



$$\frac{\partial \sigma}{\partial w_1}, \frac{\partial \sigma}{\partial w_2}, \frac{\partial \sigma}{\partial b}, \frac{\partial \sigma}{\partial x_1}, \frac{\partial \sigma}{\partial x_2}$$

$$f(x) = e^x \rightarrow \frac{df}{dx} = e^x$$

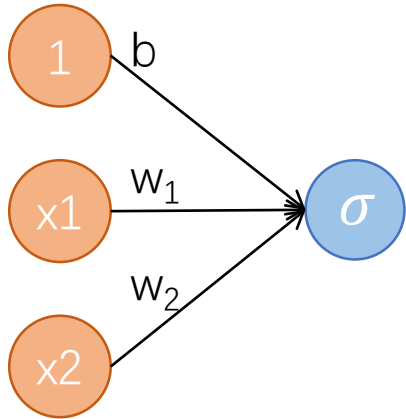
$$f_a(x) = ax \rightarrow \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \rightarrow \frac{df}{dx} = -1/x^2$$

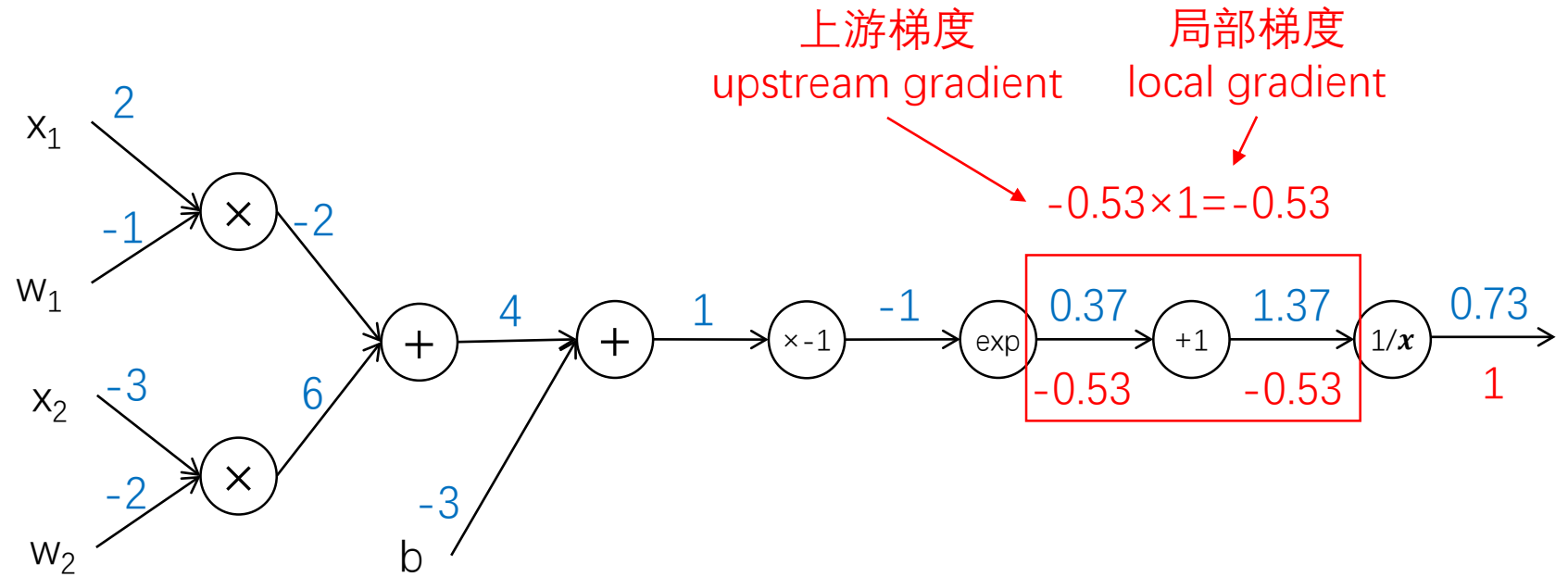
$$f_c(x) = c + x \rightarrow \frac{df}{dx} = 1$$

# 反向传播 (Backpropagation)

$$\sigma(W, x) = \frac{1}{1 + e^{-(w_1x_1 + w_2x_2 + b)}}$$



$$\frac{\partial \sigma}{\partial w_1}, \frac{\partial \sigma}{\partial w_2}, \frac{\partial \sigma}{\partial b}, \frac{\partial \sigma}{\partial x_1}, \frac{\partial \sigma}{\partial x_2}$$



$$f(x) = e^x \rightarrow \frac{df}{dx} = e^x$$

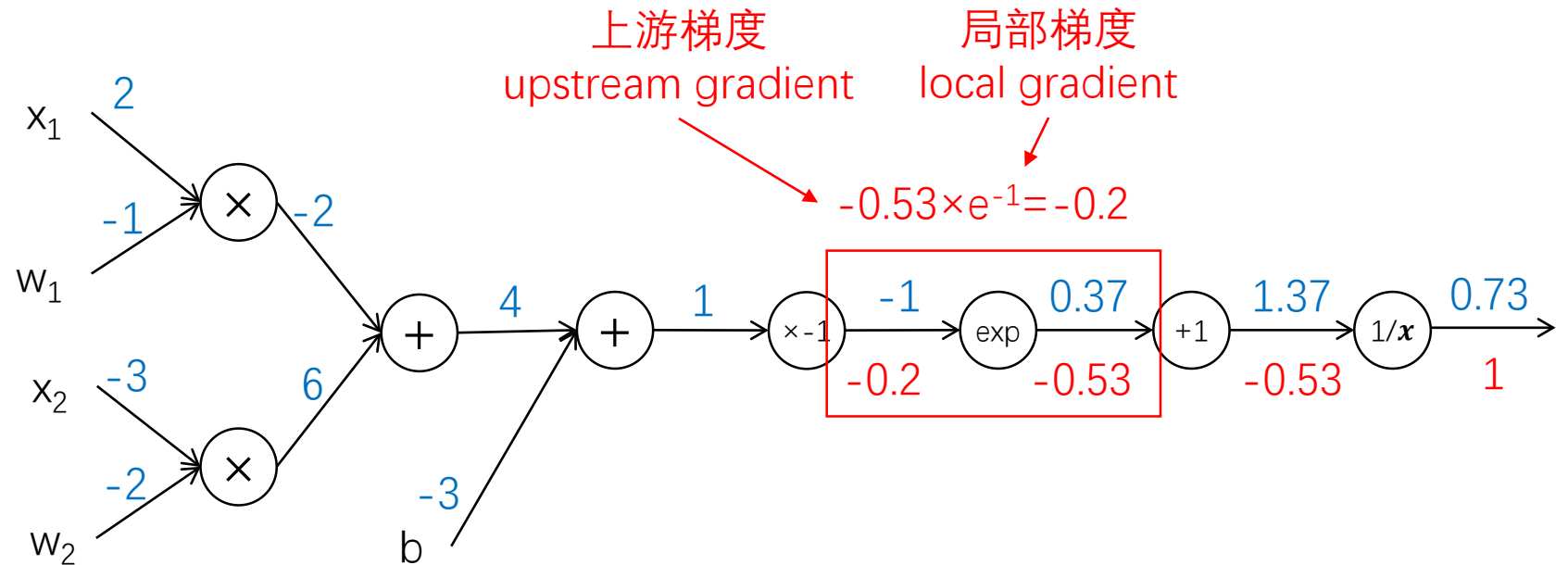
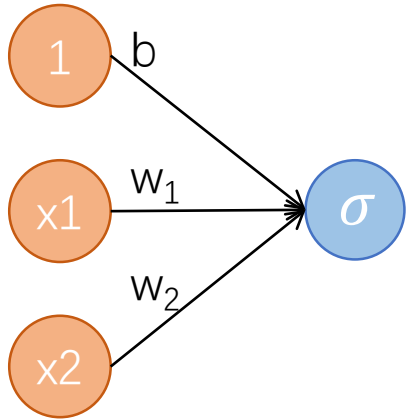
$$f_a(x) = ax \rightarrow \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \rightarrow \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \rightarrow \frac{df}{dx} = 1$$

# 反向传播 (Backpropagation)

$$\sigma(W, x) = \frac{1}{1 + e^{-(w_1x_1 + w_2x_2 + b)}}$$



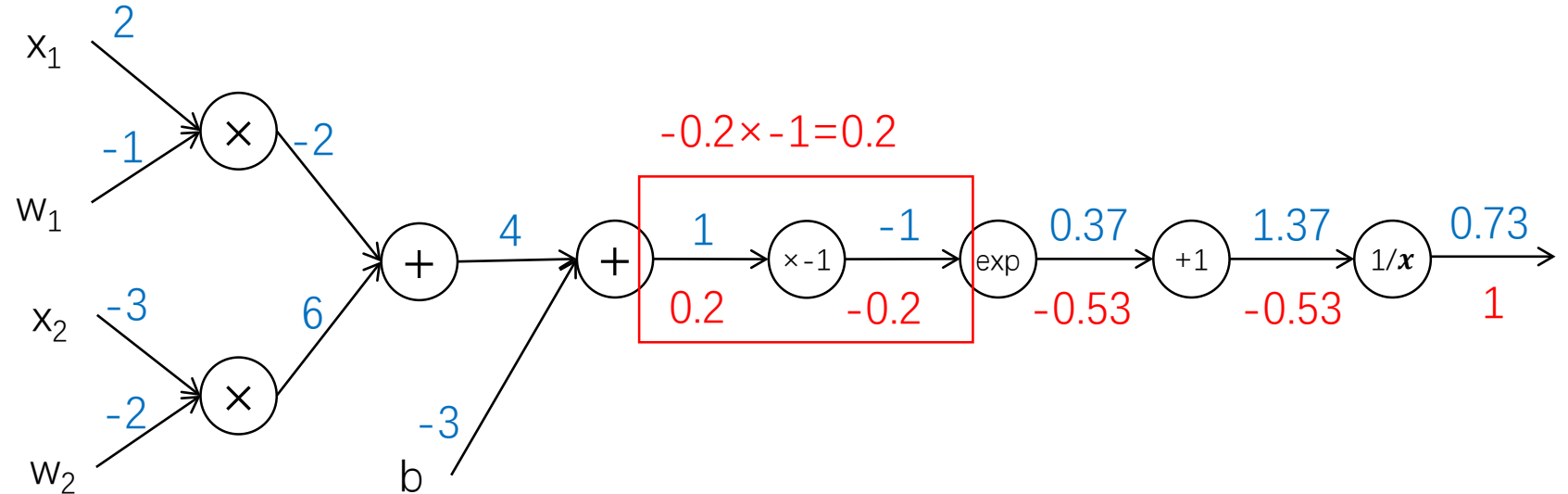
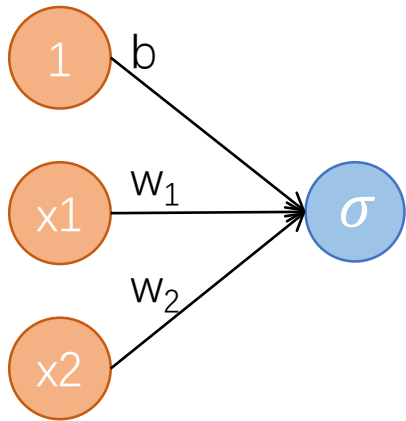
$$\frac{\partial \sigma}{\partial w_1}, \frac{\partial \sigma}{\partial w_2}, \frac{\partial \sigma}{\partial b}, \frac{\partial \sigma}{\partial x_1}, \frac{\partial \sigma}{\partial x_2}$$

$f(x) = e^x \rightarrow \frac{df}{dx} = e^x$	$f(x) = \frac{1}{x} \rightarrow \frac{df}{dx} = -1/x^2$
$f_a(x) = ax \rightarrow \frac{df}{dx} = a$	$f_c(x) = c + x \rightarrow \frac{df}{dx} = 1$



# 反向传播 (Backpropagation)

$$\sigma(W, x) = \frac{1}{1 + e^{-(w_1x_1 + w_2x_2 + b)}}$$



$$\frac{\partial \sigma}{\partial w_1}, \frac{\partial \sigma}{\partial w_2}, \frac{\partial \sigma}{\partial b}, \frac{\partial \sigma}{\partial x_1}, \frac{\partial \sigma}{\partial x_2}$$

$$f(x) = e^x \rightarrow \frac{df}{dx} = e^x$$

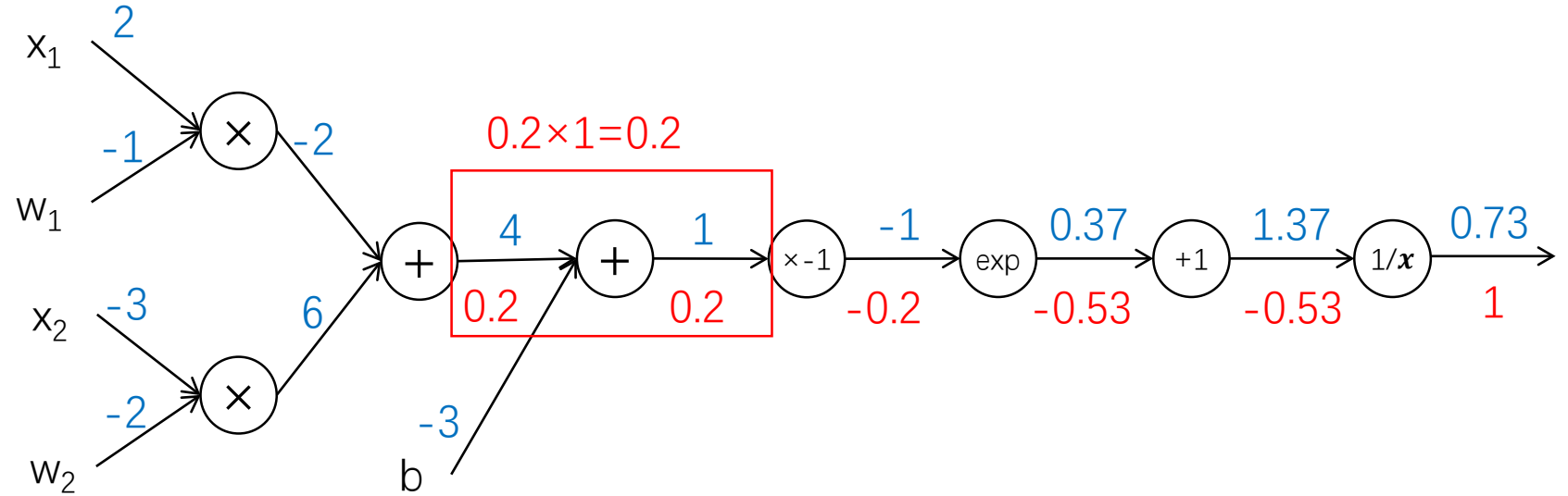
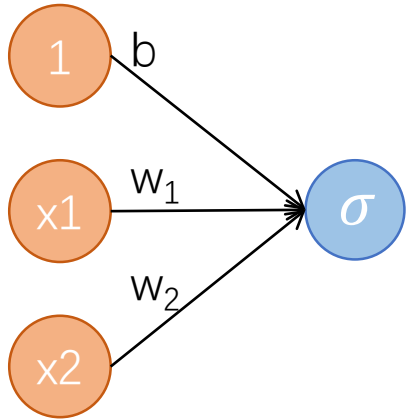
$$f_a(x) = ax \rightarrow \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \rightarrow \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \rightarrow \frac{df}{dx} = 1$$

# 反向传播 (Backpropagation)

$$\sigma(W, x) = \frac{1}{1 + e^{-(w_1x_1 + w_2x_2 + b)}}$$



$$\frac{\partial \sigma}{\partial w_1}, \frac{\partial \sigma}{\partial w_2}, \frac{\partial \sigma}{\partial b}, \frac{\partial \sigma}{\partial x_1}, \frac{\partial \sigma}{\partial x_2}$$

$$f(x) = e^x \rightarrow \frac{df}{dx} = e^x$$

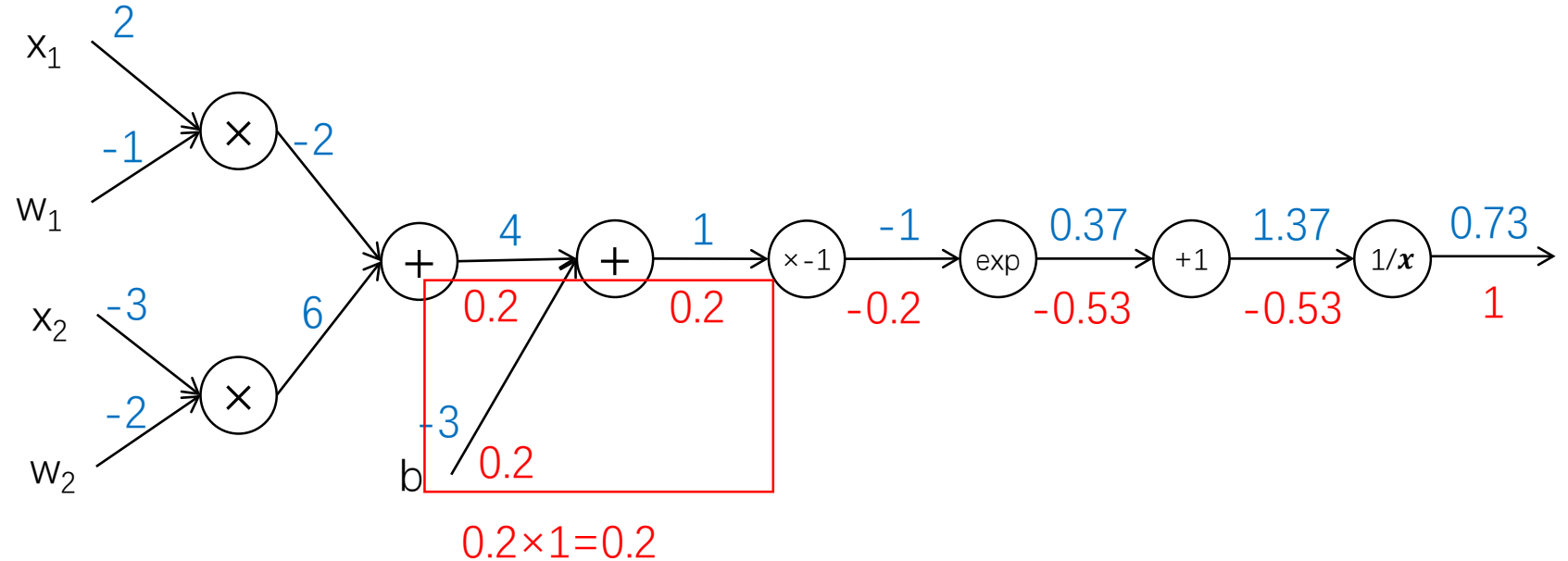
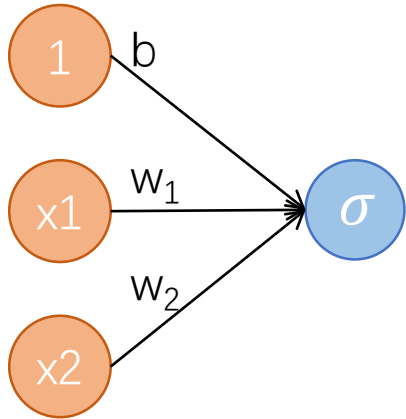
$$f_a(x) = ax \rightarrow \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \rightarrow \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \rightarrow \frac{df}{dx} = 1$$

# 反向传播 (Backpropagation)

$$\sigma(W, x) = \frac{1}{1 + e^{-(w_1x_1 + w_2x_2 + b)}}$$



$$\frac{\partial \sigma}{\partial w_1}, \frac{\partial \sigma}{\partial w_2}, \frac{\partial \sigma}{\partial b}, \frac{\partial \sigma}{\partial x_1}, \frac{\partial \sigma}{\partial x_2}$$

$$f(x) = e^x \rightarrow \frac{df}{dx} = e^x$$

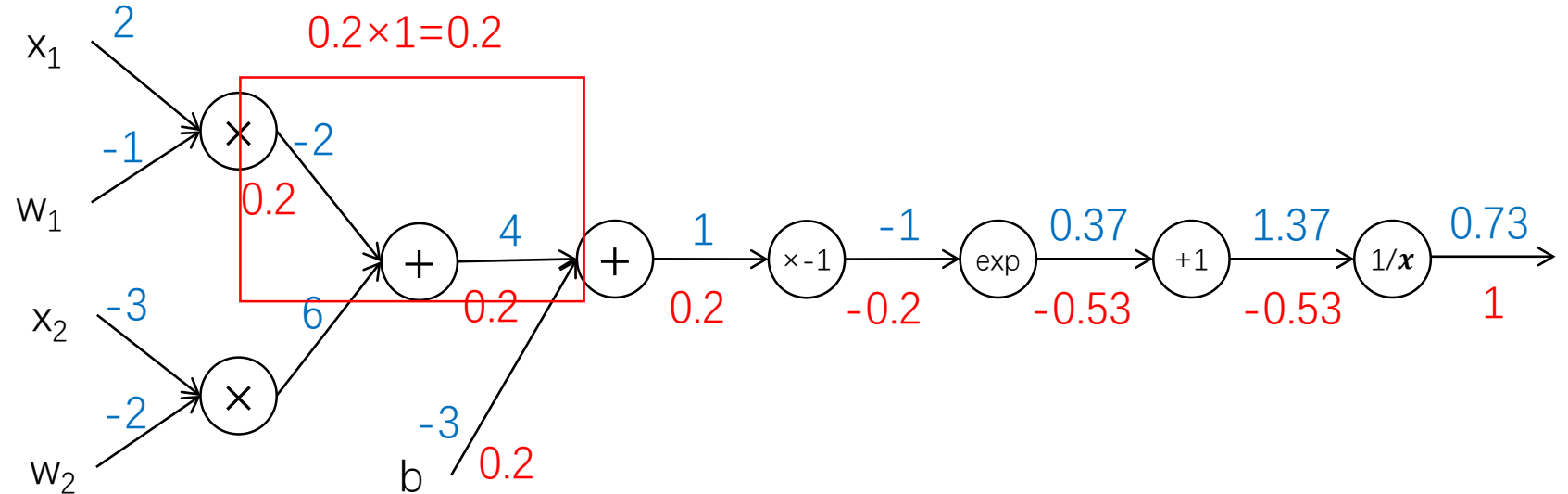
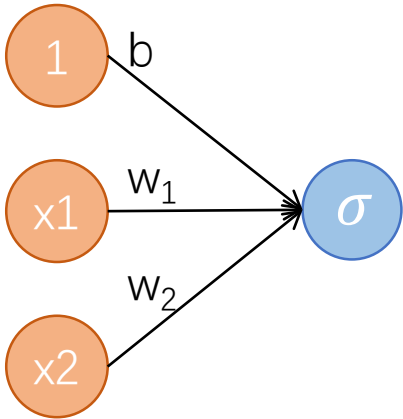
$$f_a(x) = ax \rightarrow \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \rightarrow \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \rightarrow \frac{df}{dx} = 1$$

# 反向传播 (Backpropagation)

$$\sigma(W, x) = \frac{1}{1 + e^{-(w_1x_1 + w_2x_2 + b)}}$$



$$\frac{\partial \sigma}{\partial w_1}, \frac{\partial \sigma}{\partial w_2}, \frac{\partial \sigma}{\partial b}, \frac{\partial \sigma}{\partial x_1}, \frac{\partial \sigma}{\partial x_2}$$

$$f(x) = e^x \rightarrow \frac{df}{dx} = e^x$$

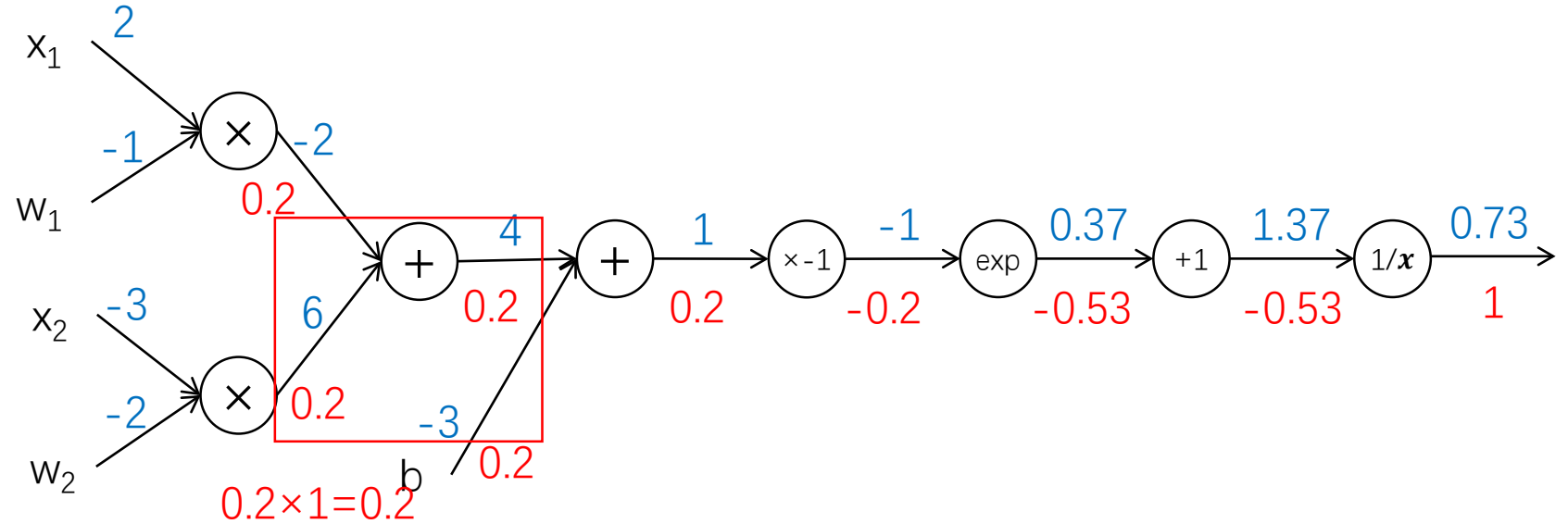
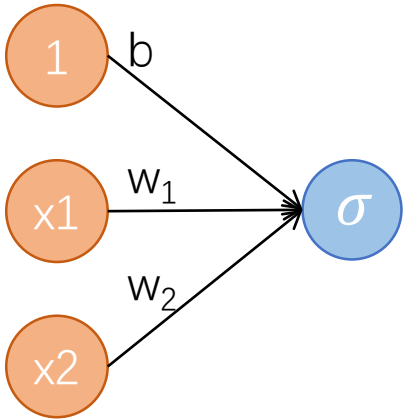
$$f_a(x) = ax \rightarrow \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \rightarrow \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \rightarrow \frac{df}{dx} = 1$$

# 反向传播 (Backpropagation)

$$\sigma(W, x) = \frac{1}{1 + e^{-(w_1x_1 + w_2x_2 + b)}}$$



$$\frac{\partial \sigma}{\partial w_1}, \frac{\partial \sigma}{\partial w_2}, \frac{\partial \sigma}{\partial b}, \frac{\partial \sigma}{\partial x_1}, \frac{\partial \sigma}{\partial x_2}$$

$$f(x) = e^x \rightarrow \frac{df}{dx} = e^x$$

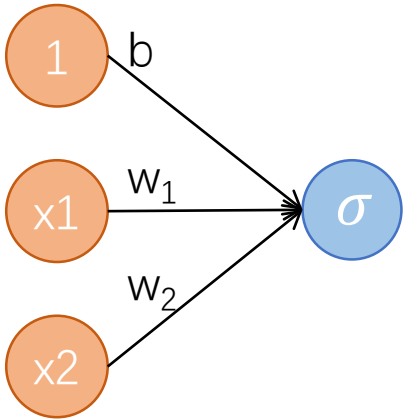
$$f_a(x) = ax \rightarrow \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \rightarrow \frac{df}{dx} = -1/x^2$$

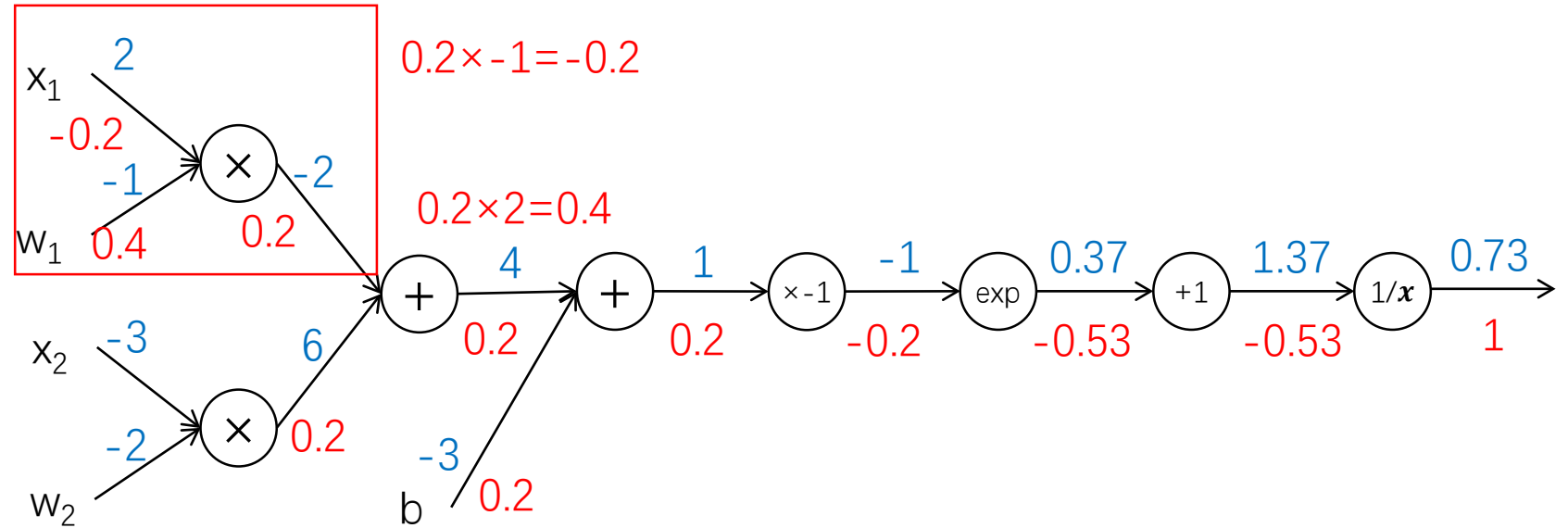
$$f_c(x) = c + x \rightarrow \frac{df}{dx} = 1$$

# 反向传播 (Backpropagation)

$$\sigma(W, x) = \frac{1}{1 + e^{-(w_1 x_1 + w_2 x_2 + b)}}$$



$$\frac{\partial \sigma}{\partial w_1}, \frac{\partial \sigma}{\partial w_2}, \frac{\partial \sigma}{\partial b}, \frac{\partial \sigma}{\partial x_1}, \frac{\partial \sigma}{\partial x_2}$$



$$f(x) = e^x \quad \longrightarrow \quad \frac{df}{dx} = e^x$$

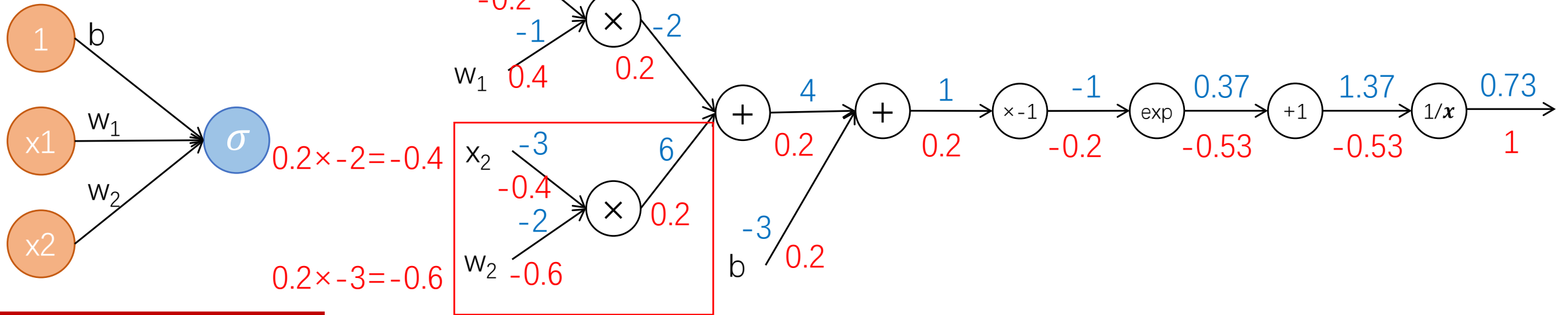
$$f_a(x) = ax \longrightarrow \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \quad \longrightarrow \quad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \longrightarrow \frac{df}{dx} = 1$$

# 反向传播 (Backpropagation)

$$\sigma(W, x) = \frac{1}{1 + e^{-(w_1x_1 + w_2x_2 + b)}}$$



$$\frac{\partial \sigma}{\partial w_1}, \frac{\partial \sigma}{\partial w_2}, \frac{\partial \sigma}{\partial b}, \frac{\partial \sigma}{\partial x_1}, \frac{\partial \sigma}{\partial x_2}$$

$$f(x) = e^x \rightarrow \frac{df}{dx} = e^x$$

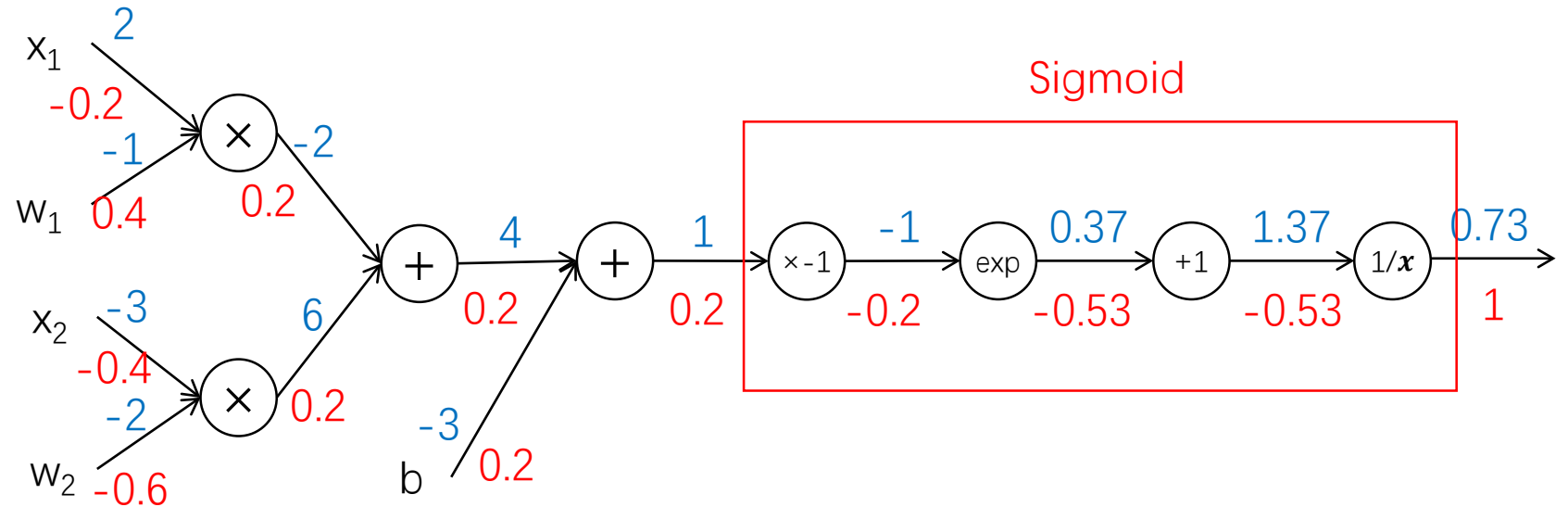
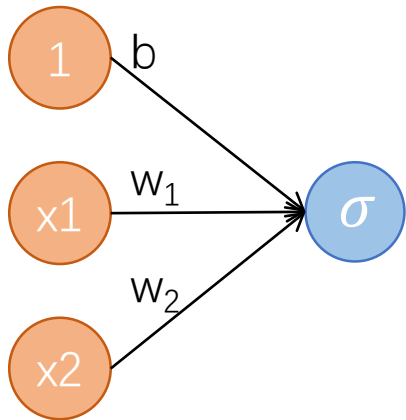
$$f_a(x) = ax \rightarrow \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \rightarrow \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \rightarrow \frac{df}{dx} = 1$$

# 反向传播 (Backpropagation)

$$\sigma(W, x) = \frac{1}{1 + e^{-(w_1x_1 + w_2x_2 + b)}}$$



$$\frac{\partial \sigma}{\partial w_1}, \frac{\partial \sigma}{\partial w_2}, \frac{\partial \sigma}{\partial b}, \frac{\partial \sigma}{\partial x_1}, \frac{\partial \sigma}{\partial x_2}$$

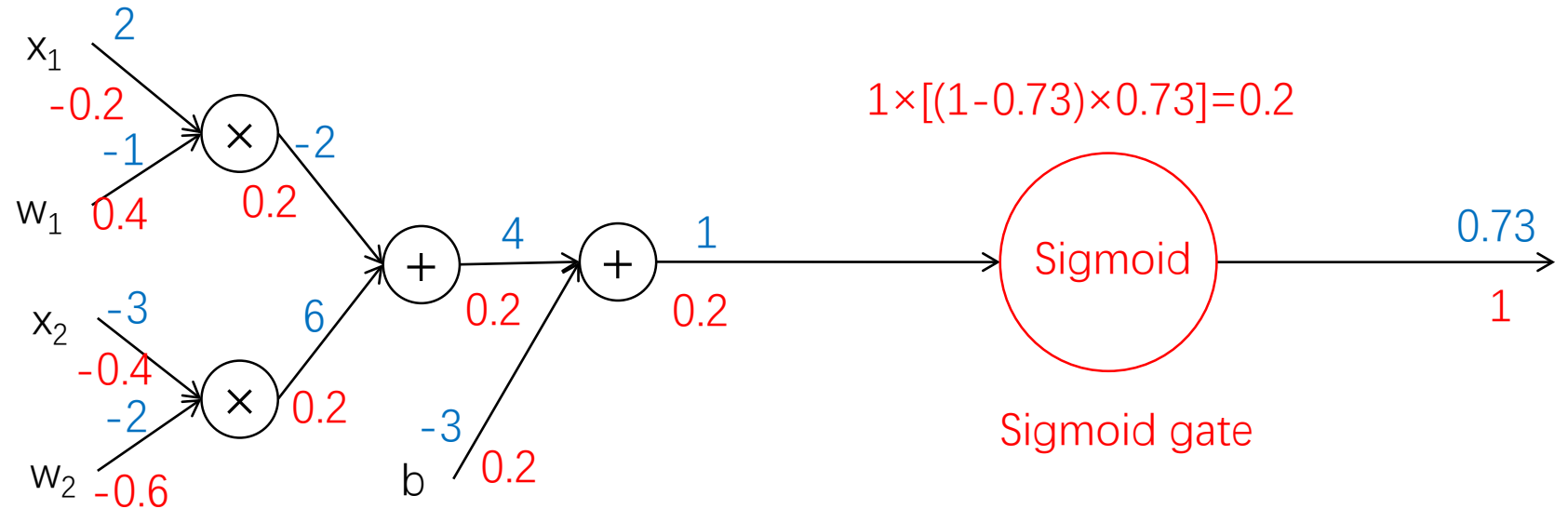
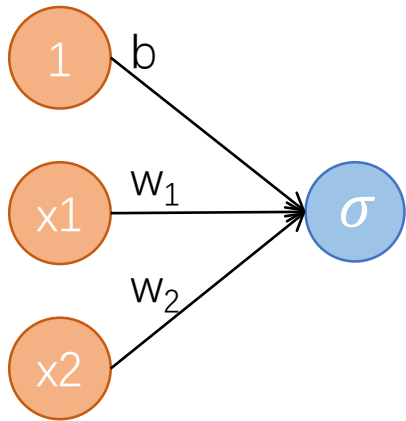
Sigmoid本地梯度

$$\frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2} = \left( \frac{1 + e^{-x} - 1}{1 + e^{-x}} \right) \left( \frac{1}{1 + e^{-x}} \right) = (1 - \sigma(x)) \sigma(x)$$



# 反向传播 (Backpropagation)

$$\sigma(W, x) = \frac{1}{1 + e^{-(w_1x_1 + w_2x_2 + b)}}$$



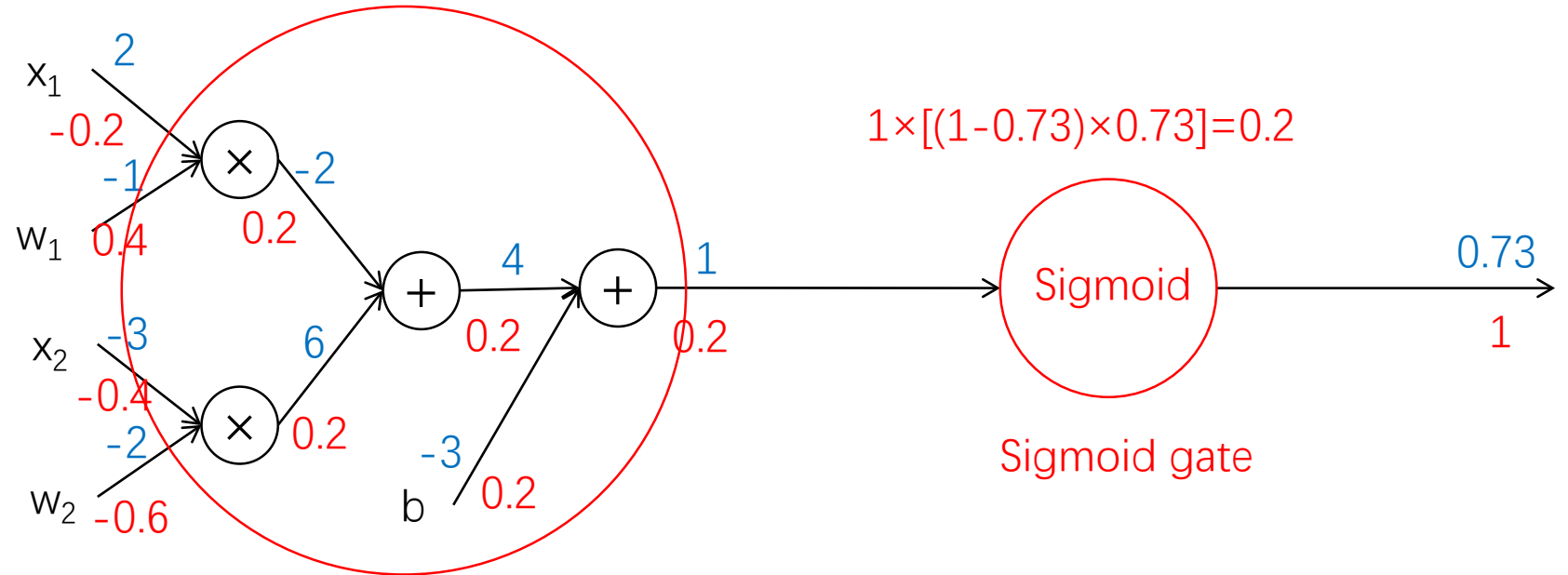
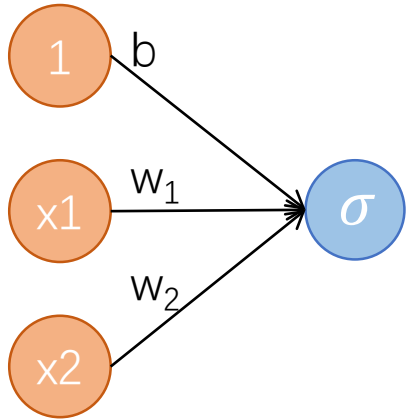
$$\frac{\partial \sigma}{\partial w_1}, \frac{\partial \sigma}{\partial w_2}, \frac{\partial \sigma}{\partial b}, \frac{\partial \sigma}{\partial x_1}, \frac{\partial \sigma}{\partial x_2}$$

Sigmoid本地梯度

$$\frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2} = \left( \frac{1 + e^{-x} - 1}{1 + e^{-x}} \right) \left( \frac{1}{1 + e^{-x}} \right) = (1 - \sigma(x)) \sigma(x)$$

# 反向传播 (Backpropagation)

$$\sigma(W, x) = \frac{1}{1 + e^{-(w_1x_1 + w_2x_2 + b)}}$$



$$\frac{\partial \sigma}{\partial w_1}, \frac{\partial \sigma}{\partial w_2}, \frac{\partial \sigma}{\partial b}, \frac{\partial \sigma}{\partial x_1}, \frac{\partial \sigma}{\partial x_2}$$

$$y = w_1x_1 + w_2x_2 + b$$

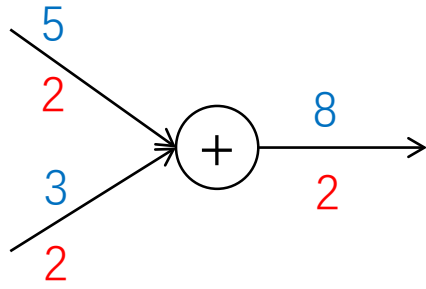
Sigmoid本地梯度

$$\frac{\partial \sigma}{\partial w_i} = \frac{\partial \sigma}{\partial y} \frac{\partial y}{\partial w_i} \quad \frac{\partial \sigma}{\partial x_i} = \frac{\partial \sigma}{\partial y} \frac{\partial y}{\partial x_i}$$

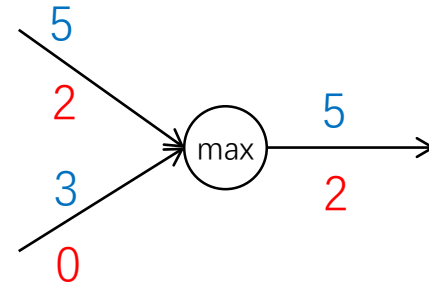
$$\frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2} = \left( \frac{1 + e^{-x} - 1}{1 + e^{-x}} \right) \left( \frac{1}{1 + e^{-x}} \right) = (1 - \sigma(x)) \sigma(x)$$

# 梯度流的常见模式

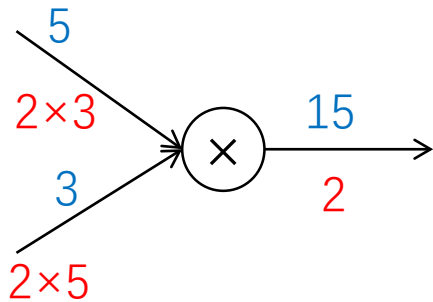
add gate: 拷贝上游梯度



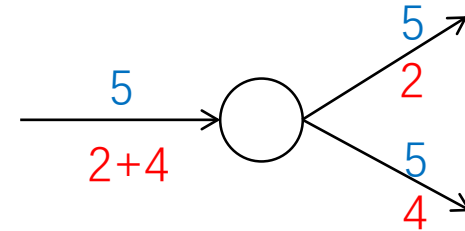
max gate: 上游梯度路由给较大变量



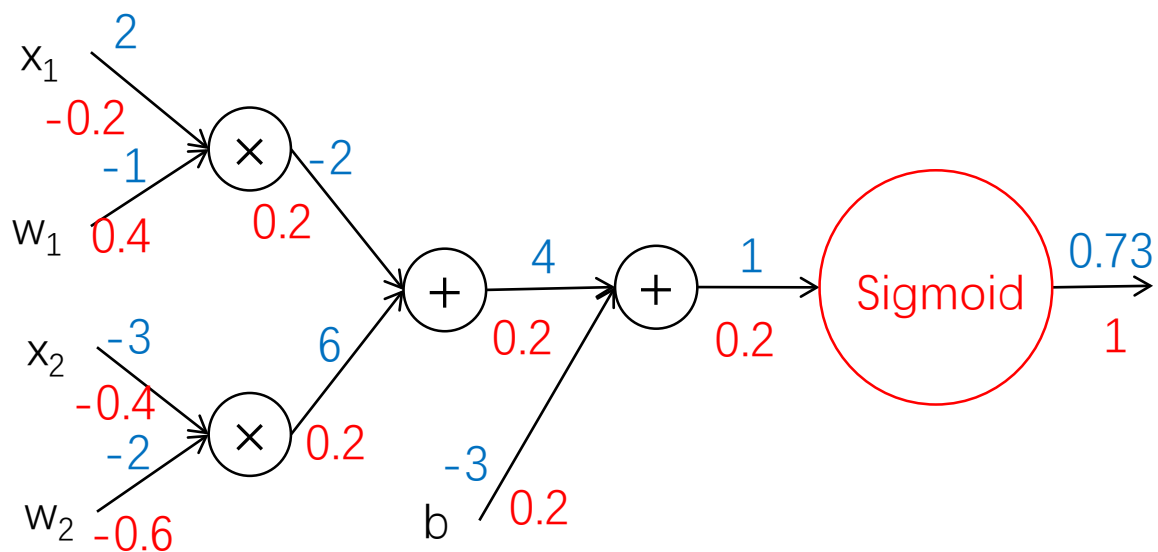
mul gate: 上游梯度  $\times$  互换变量



copy gate: 上游梯度相加



# 代码实现



Python 3.6.8 (v3.6.8:3c6b436a57, Dec 24 2018) [AMD64] flat.py

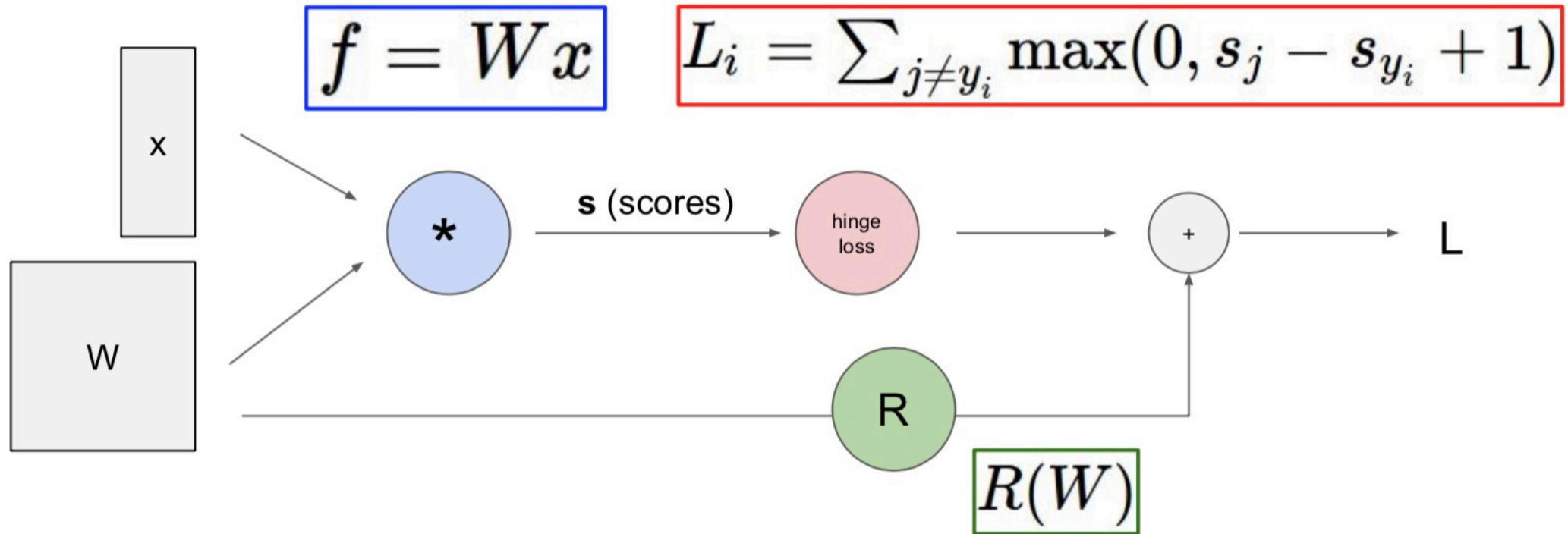
```
def f(w1,x1,w2,x2,b)
    s1 = w1 * x1
    s2 = w2 * x2
    s3 = s1 + s2
    s4 = s3 + b
    L = sigmoid(s4)
```

前向传递, 计算输出

```
grad_L = 1.0
grad_s4 = grad_L * (1-L) * L
grad_b = grad_s4
grad_s3 = grad_s4
grad_s1 = grad_s3
grad_s2 = grad_s3
grad_w1 = grad_s1 * x1
grad_x1 = grad_s1 * w1
grad_w2 = grad_s2 * x2
grad_x2 = grad_s2 * w2
```

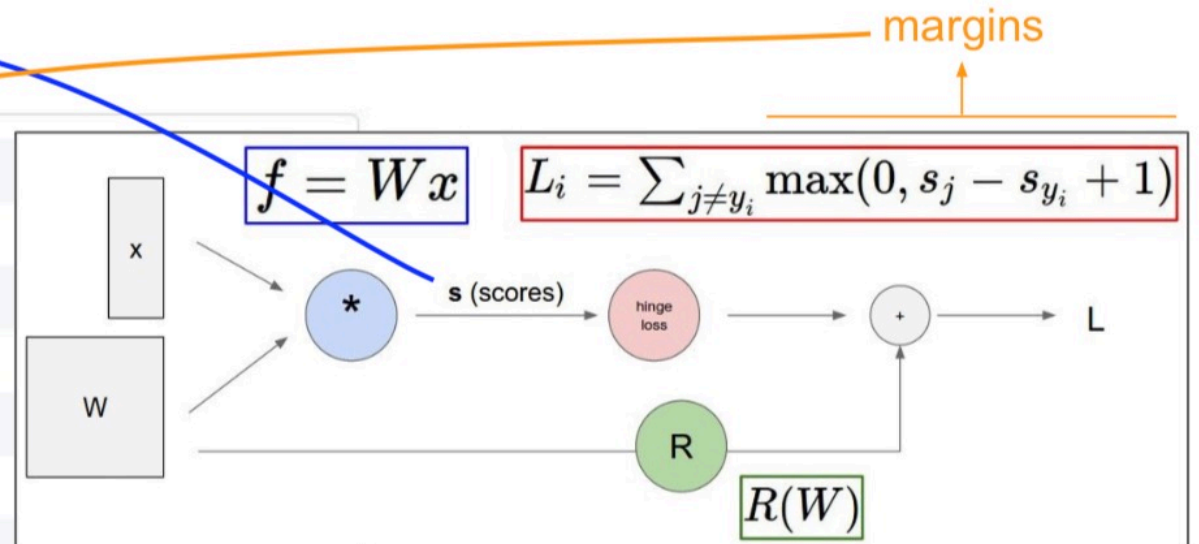
反向传递, 计算梯度

# SVM分类器计算图+反向传播

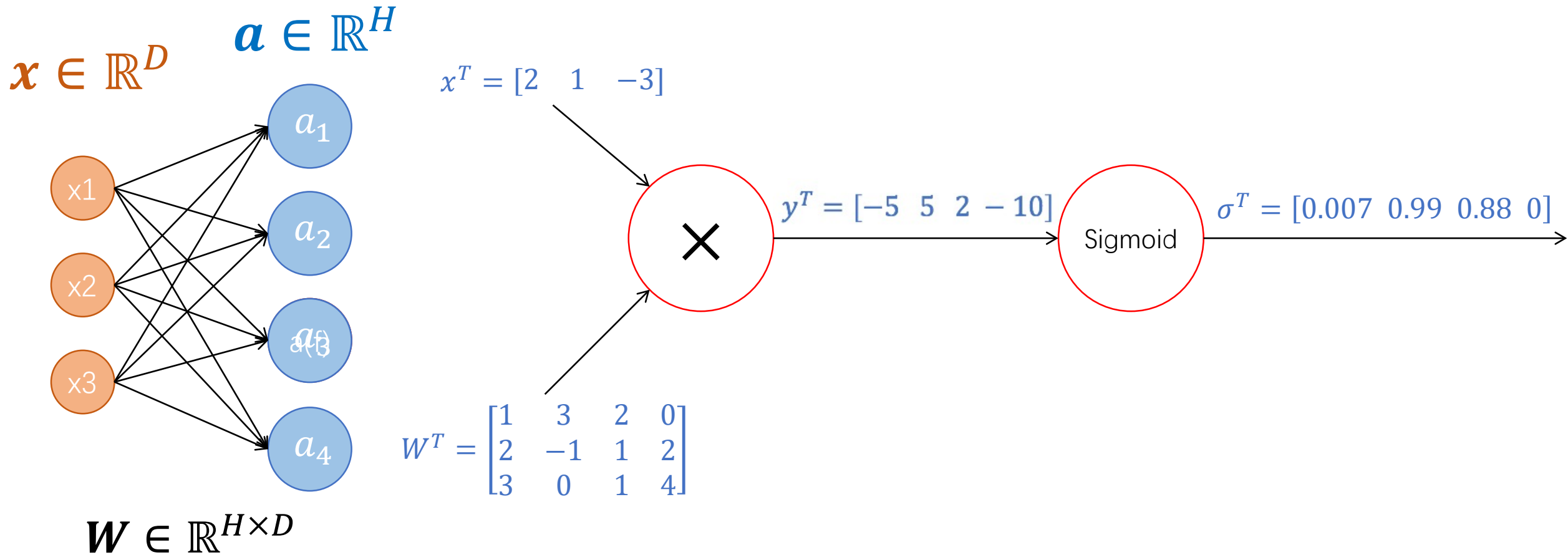


# SVM分类器计算图+反向传播实现

```
# receive W (weights), X (data)
# forward pass (we have 8 lines)
scores = #...
margins = #...
data_loss = #...
reg_loss = #...
loss = data_loss + reg_loss
# backward pass (we have 5 lines)
dmargins = # ... (optionally, we go direct to dscores)
dscores = #...
dW = #...
```

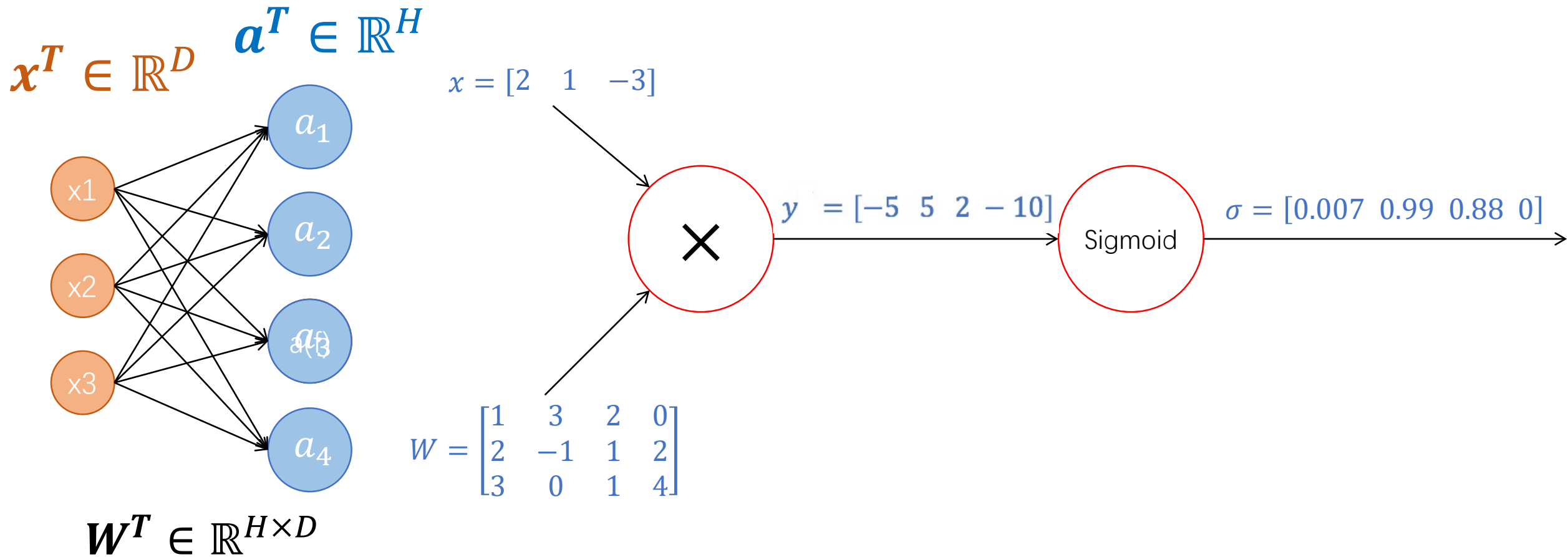


# 反向传播的矩阵运算





# 反向传播的矩阵运算

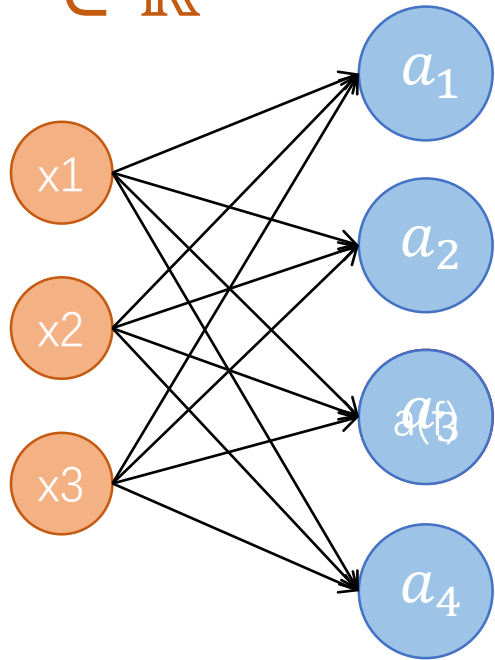


# 反向传播的矩阵运算

$$\mathbf{x}^T \in \mathbb{R}^D \quad \mathbf{a}^T \in \mathbb{R}^H$$

雅可比矩阵  
(Jacobian matrix)

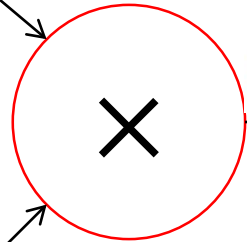
$$\mathbf{J} = \begin{bmatrix} \frac{\partial \mathbf{f}}{\partial x_1} & \cdots & \frac{\partial \mathbf{f}}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}.$$



$$\mathbf{W}^T \in \mathbb{R}^{H \times D}$$

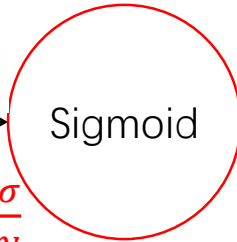
$$\mathbf{x} = [2 \quad 1 \quad -3]$$

$$\mathbf{W} = \begin{bmatrix} 1 & 3 & 2 & 0 \\ 2 & -1 & 1 & 2 \\ 3 & 0 & 1 & 4 \end{bmatrix}$$



$$\mathbf{y} = [-5 \quad 5 \quad 2 \quad -10]$$

计算  $\frac{\partial \sigma}{\partial y}$



$$\sigma = [0.007 \quad 0.99 \quad 0.88 \quad 0]$$

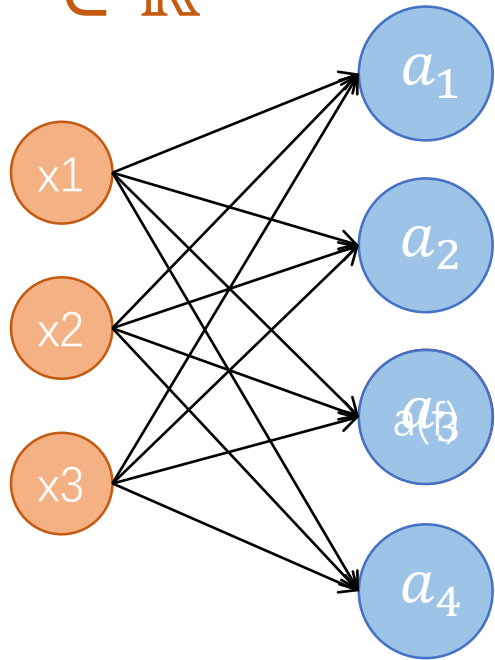
$$\frac{\partial h}{\partial \sigma} = [4 \quad 1 \quad 2 \quad -1]$$

# 反向传播的矩阵运算

$$\mathbf{x}^T \in \mathbb{R}^D \quad \mathbf{a}^T \in \mathbb{R}^H$$

雅可比矩阵  
(Jacobian matrix)

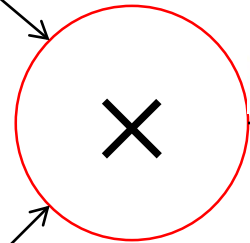
$$\mathbf{J} = \begin{bmatrix} \frac{\partial \mathbf{f}}{\partial x_1} & \cdots & \frac{\partial \mathbf{f}}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}.$$



$$\mathbf{W}^T \in \mathbb{R}^{H \times D}$$

$$\mathbf{x} = [2 \quad 1 \quad -3]$$

$$\mathbf{W} = \begin{bmatrix} 1 & 3 & 2 & 0 \\ 2 & -1 & 1 & 2 \\ 3 & 0 & 1 & 4 \end{bmatrix}$$



$$\mathbf{y} = [-5 \quad 5 \quad 2 \quad -10]$$

Sigmoid

$$\sigma = [0.007 \quad 0.99 \quad 0.88 \quad 0]$$

计算  $\frac{\partial \sigma}{\partial \mathbf{y}}$

$$\frac{\partial h}{\partial \sigma} = [4 \quad 1 \quad 2 \quad -1]$$

$$\left(\frac{\partial \sigma}{\partial \mathbf{y}}\right)_{n,m} = \left(\frac{\partial \sigma_n}{\partial y_m}\right) \longrightarrow \frac{\partial \sigma}{\partial \mathbf{y}} \in \mathbb{R}^{H \times H}$$

假设  $H=4096$ ，需要至少  $16777216 \approx 16\text{M}$  内存

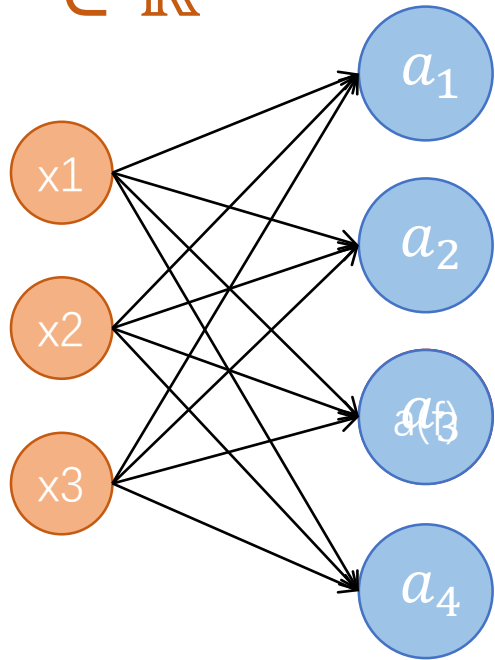
实际上我们不需要显示实现 Jacobian Matrix

# 反向传播的矩阵运算

$$\mathbf{x}^T \in \mathbb{R}^D \quad \mathbf{a}^T \in \mathbb{R}^H$$

雅可比矩阵  
(Jacobian matrix)

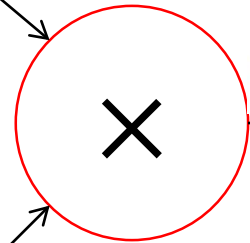
$$\mathbf{J} = \begin{bmatrix} \frac{\partial \mathbf{f}}{\partial x_1} & \cdots & \frac{\partial \mathbf{f}}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}.$$



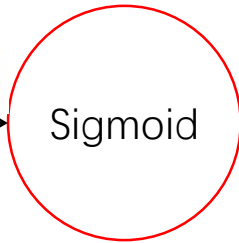
$$\mathbf{W}^T \in \mathbb{R}^{H \times D}$$

$$\mathbf{x} = [2 \quad 1 \quad -3]$$

$$\mathbf{W} = \begin{bmatrix} 1 & 3 & 2 & 0 \\ 2 & -1 & 1 & 2 \\ 3 & 0 & 1 & 4 \end{bmatrix}$$



$$\mathbf{y} = [-5 \quad 5 \quad 2 \quad -10]$$



$$\boldsymbol{\sigma} = [0.007 \quad 0.99 \quad 0.88 \quad 0]$$

$$\frac{\partial \mathbf{h}}{\partial \boldsymbol{\sigma}} = [4 \quad 1 \quad 2 \quad -1]$$

$$\left(\frac{\partial \sigma}{\partial y}\right)_{1,1} = (1 - 0.007) \times 0.007 = 0.007$$

$$\left(\frac{\partial \sigma}{\partial y}\right)_{1,2} = 0$$

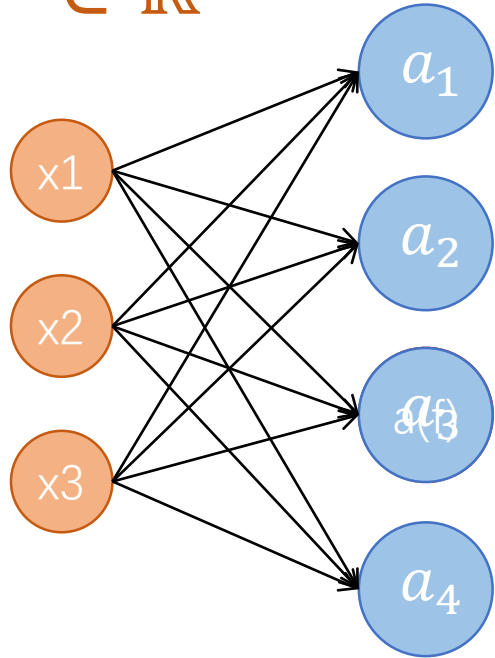
$$\left(\frac{\partial \sigma}{\partial y}\right)_{1,3} = 0$$

$$\left(\frac{\partial \sigma}{\partial y}\right)_{1,4} = 0$$

$$\frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2} = \left(\frac{1 + e^{-x} - 1}{1 + e^{-x}}\right) \left(\frac{1}{1 + e^{-x}}\right) = (1 - \sigma(x))\sigma(x)$$

# 反向传播的矩阵运算

$$\mathbf{x}^T \in \mathbb{R}^D \quad \mathbf{a}^T \in \mathbb{R}^H$$



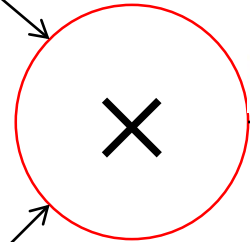
$$\mathbf{W}^T \in \mathbb{R}^{H \times D}$$

$$\mathbf{x} = [2 \quad 1 \quad -3]$$

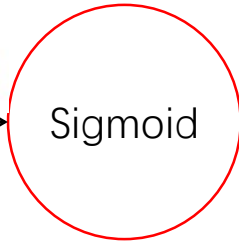
$$\mathbf{W} = \begin{bmatrix} 1 & 3 & 2 & 0 \\ 2 & -1 & 1 & 2 \\ 3 & 0 & 1 & 4 \end{bmatrix}$$

雅可比矩阵  
(Jacobian matrix)

$$\mathbf{J} = \begin{bmatrix} \frac{\partial \mathbf{f}}{\partial x_1} & \cdots & \frac{\partial \mathbf{f}}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}.$$



$$\mathbf{y} = [-5 \quad 5 \quad 2 \quad -10]$$



$$\sigma = [0.007 \quad 0.99 \quad 0.88 \quad 0]$$

$$\frac{\partial h}{\partial \sigma} = [4 \quad 1 \quad 2 \quad -1]$$

$$\left(\frac{\partial \sigma}{\partial y}\right)_{1,1} = (1 - 0.007) \times 0.007 = 0.007$$

$$\left(\frac{\partial \sigma}{\partial y}\right)_{1,2} = 0 \quad \left(\frac{\partial \sigma}{\partial y}\right)_{1,4} = 0 \quad y_i \text{ 只影响 } \sigma_i$$

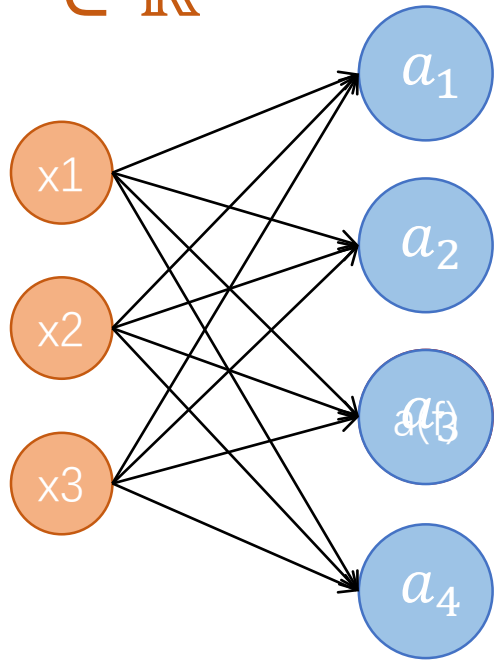
$$\frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2} = \left(\frac{1 + e^{-x} - 1}{1 + e^{-x}}\right) \left(\frac{1}{1 + e^{-x}}\right) = (1 - \sigma(x))\sigma(x)$$

# 反向传播的矩阵运算

$$\mathbf{x}^T \in \mathbb{R}^D \quad \mathbf{a}^T \in \mathbb{R}^H$$

雅可比矩阵  
(Jacobian matrix)

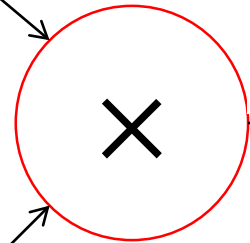
$$\mathbf{J} = \begin{bmatrix} \frac{\partial \mathbf{f}}{\partial x_1} & \cdots & \frac{\partial \mathbf{f}}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}.$$



$$\mathbf{W}^T \in \mathbb{R}^{H \times D}$$

$$\mathbf{x} = [2 \quad 1 \quad -3]$$

$$\mathbf{W} = \begin{bmatrix} 1 & 3 & 2 & 0 \\ 2 & -1 & 1 & 2 \\ 3 & 0 & 1 & 4 \end{bmatrix}$$



$$\mathbf{y} = [-5 \quad 5 \quad 2 \quad -10]$$

Sigmoid

$$\sigma = [0.007 \quad 0.99 \quad 0.88 \quad 0]$$

$$\frac{\partial h}{\partial \sigma} = [4 \quad 1 \quad 2 \quad -1]$$

$$\frac{\partial \sigma}{\partial \mathbf{y}} = \begin{bmatrix} 0.007 & 0 & 0 & 0 \\ 0 & 0.01 & 0 & 0 \\ 0 & 0 & 0.11 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

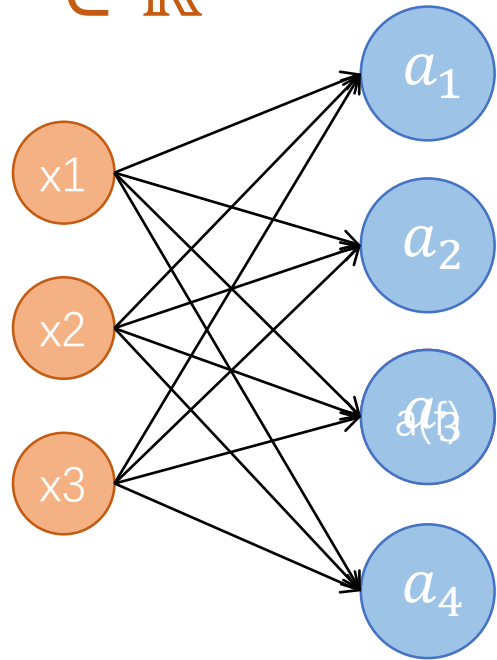
$$\frac{\partial h}{\partial y_i} = \frac{\partial h}{\partial \sigma} \frac{\partial \sigma}{\partial y_i} = \sum_j \frac{\partial h}{\partial \sigma_j} \frac{\partial \sigma_j}{\partial y_i}$$

# 反向传播的矩阵运算

$$\mathbf{x}^T \in \mathbb{R}^D \quad \mathbf{a}^T \in \mathbb{R}^H$$

雅可比矩阵  
(Jacobian matrix)

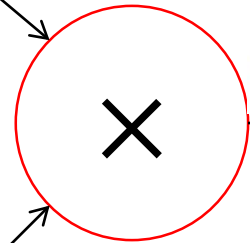
$$\mathbf{J} = \begin{bmatrix} \frac{\partial \mathbf{f}}{\partial x_1} & \cdots & \frac{\partial \mathbf{f}}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}.$$



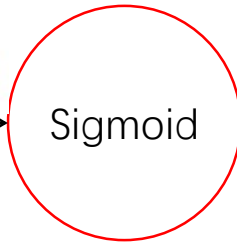
$$\mathbf{W}^T \in \mathbb{R}^{H \times D}$$

$$\mathbf{x} = [2 \quad 1 \quad -3]$$

$$\mathbf{W} = \begin{bmatrix} 1 & 3 & 2 & 0 \\ 2 & -1 & 1 & 2 \\ 3 & 0 & 1 & 4 \end{bmatrix}$$



$$\mathbf{y} = [-5 \quad 5 \quad 2 \quad -10]$$



$$\boldsymbol{\sigma} = [0.007 \quad 0.99 \quad 0.88 \quad 0]$$

$$\frac{\partial h}{\partial \boldsymbol{\sigma}} = [4 \quad 1 \quad 2 \quad -1]$$

$$\frac{\partial \boldsymbol{\sigma}}{\partial \mathbf{y}} = \begin{bmatrix} 0.007 & 0 & 0 & 0 \\ 0 & 0.01 & 0 & 0 \\ 0 & 0 & 0.11 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

$$\frac{\partial h}{\partial y_i} = \frac{\partial h}{\partial \boldsymbol{\sigma}} \frac{\partial \boldsymbol{\sigma}}{\partial y_i} = \sum_j \frac{\partial h}{\partial \sigma_j} \frac{\partial \sigma_j}{\partial y_i}$$

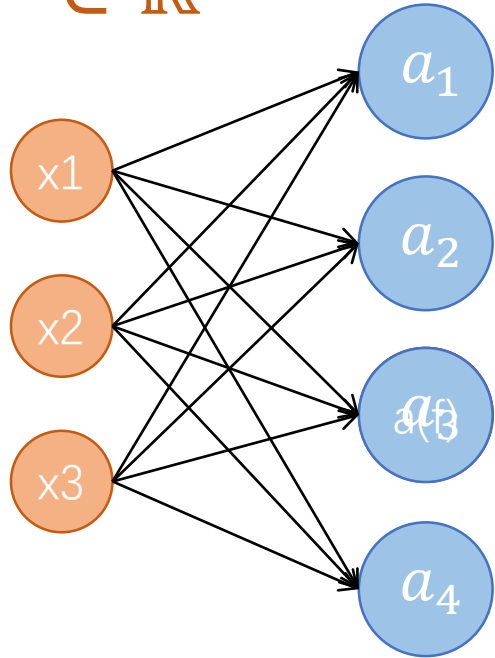
$$\frac{\partial h}{\partial \mathbf{y}} = \frac{\partial h}{\partial \boldsymbol{\sigma}} \frac{\partial \boldsymbol{\sigma}}{\partial \mathbf{y}} = [4 \quad 1 \quad 2 \quad -1] \begin{bmatrix} 0.007 & 0 & 0 & 0 \\ 0 & 0.01 & 0 & 0 \\ 0 & 0 & 0.11 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} = [0.028 \quad 0.01 \quad 0.22 \quad 0]$$

# 反向传播的矩阵运算

$$\mathbf{x}^T \in \mathbb{R}^D \quad \mathbf{a}^T \in \mathbb{R}^H$$

雅可比矩阵  
(Jacobian matrix)

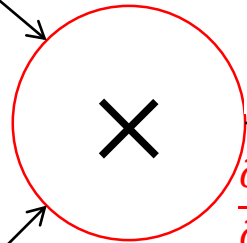
$$\mathbf{J} = \begin{bmatrix} \frac{\partial \mathbf{f}}{\partial x_1} & \cdots & \frac{\partial \mathbf{f}}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}.$$



$$\mathbf{W}^T \in \mathbb{R}^{H \times D}$$

$$\mathbf{x} = [2 \quad 1 \quad -3]$$

$$\mathbf{W} = \begin{bmatrix} 1 & 3 & 2 & 0 \\ 2 & -1 & 1 & 2 \\ 3 & 0 & 1 & 4 \end{bmatrix}$$



$$\mathbf{y} = [-5 \quad 5 \quad 2 \quad -10]$$

$$\frac{\partial h}{\partial \mathbf{y}} = [0.028 \quad 0.01 \quad 0.22 \quad 0]$$

$$\frac{\partial \sigma}{\partial \mathbf{y}} = \begin{bmatrix} 0.007 & 0 & 0 & 0 \\ 0 & 0.01 & 0 & 0 \\ 0 & 0 & 0.11 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

$$\frac{\partial h}{\partial \mathbf{y}} = \frac{\partial h}{\partial \sigma} \frac{\partial \sigma}{\partial \mathbf{y}} = [4 \quad 1 \quad 2 \quad -1] \begin{bmatrix} 0.007 & 0 & 0 & 0 \\ 0 & 0.01 & 0 & 0 \\ 0 & 0 & 0.11 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} = [0.028 \quad 0.01 \quad 0.22 \quad 0]$$

Sigmoid

$$\sigma = [0.007 \quad 0.99 \quad 0.88 \quad 0]$$

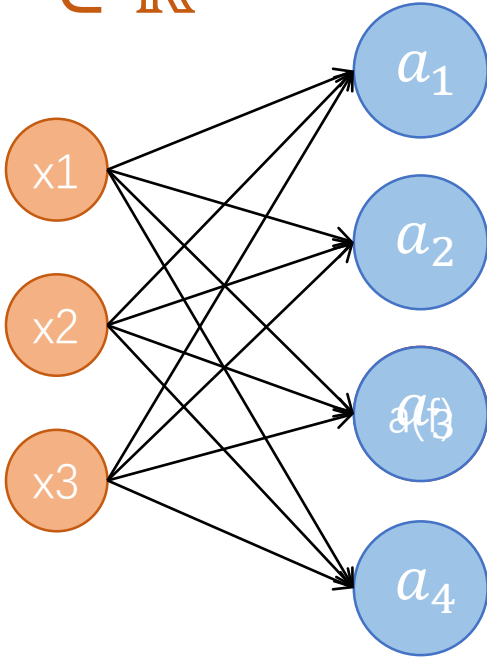
$$\frac{\partial h}{\partial \sigma} = [4 \quad 1 \quad 2 \quad -1]$$

$$\frac{\partial h}{\partial y_i} = \frac{\partial h}{\partial \sigma} \frac{\partial \sigma}{\partial y_i} = \sum_j \frac{\partial h}{\partial \sigma_j} \frac{\partial \sigma_j}{\partial y_i}$$



# 反向传播的矩阵运算

$x^T \in \mathbb{R}^D$      $a^T \in \mathbb{R}^H$



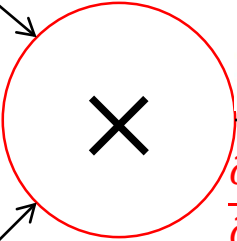
$W^T \in \mathbb{R}^{H \times D}$

$x = [2 \quad 1 \quad -3]$

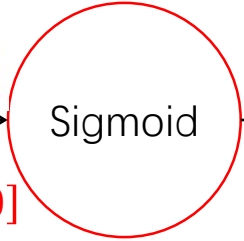
$W = \begin{bmatrix} 1 & 3 & 2 & 0 \\ 2 & -1 & 1 & 2 \\ 3 & 0 & 1 & 4 \end{bmatrix}$

雅可比矩阵  
(Jacobian matrix)

$J = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \cdots & \frac{\partial f}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$



$y = [-5 \quad 5 \quad 2 \quad -10]$



$\sigma = [0.007 \quad 0.99 \quad 0.88 \quad 0]$

$\frac{\partial h}{\partial y} = [0.028 \quad 0.01 \quad 0.22 \quad 0]$

$\frac{\partial h}{\partial \sigma} = [4 \quad 1 \quad 2 \quad -1]$

$\frac{\partial \sigma}{\partial y} = \begin{bmatrix} 0.007 & 0 & 0 & 0 \\ 0 & 0.01 & 0 & 0 \\ 0 & 0 & 0.11 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$

$\frac{\partial h}{\partial y_i} = \frac{\partial h}{\partial \sigma} \frac{\partial \sigma}{\partial y_i} = \sum_j \frac{\partial h}{\partial \sigma_j} \frac{\partial \sigma_j}{\partial y_i}$

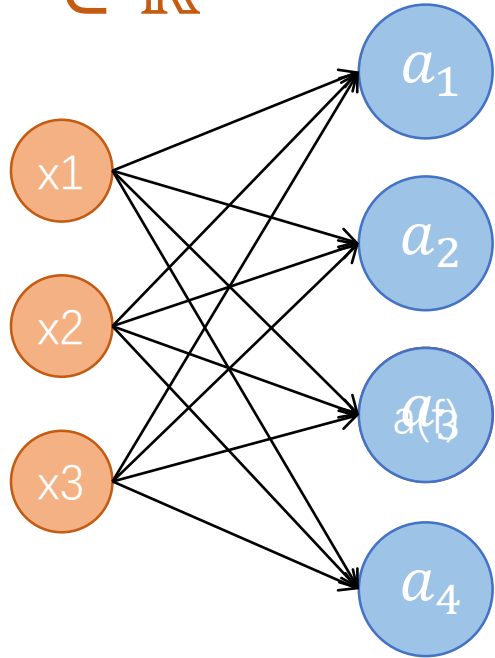
$\frac{\partial h}{\partial y} = \frac{\partial h}{\partial \sigma} \frac{\partial \sigma}{\partial y} = [4 \quad 1 \quad 2 \quad -1]$

只需要进行pairwise的梯度计算

$\begin{bmatrix} 0.007 & 0 & 0 & 0 \\ 0 & 0.01 & 0 & 0 \\ 0 & 0 & 0.11 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0.028 & 0.01 & 0.22 & 0 \\ 0.007 & 0 & 0 & 0 \\ 0 & 0.01 & 0 & 0 \\ 0 & 0 & 0.11 & 0 \end{bmatrix}$

# 反向传播的矩阵运算

$$\mathbf{x}^T \in \mathbb{R}^D \quad \mathbf{a}^T \in \mathbb{R}^H$$



$$\mathbf{W}^T \in \mathbb{R}^{H \times D}$$

$$\mathbf{x} = [2 \quad 1 \quad -3]$$

$$\mathbf{W} = \begin{bmatrix} 1 & 3 & 2 & 0 \\ 2 & -1 & 1 & 2 \\ 3 & 0 & 1 & 4 \end{bmatrix}$$

雅可比矩阵  
(Jacobian matrix)

$$\mathbf{J} = \begin{bmatrix} \frac{\partial \mathbf{f}}{\partial x_1} & \cdots & \frac{\partial \mathbf{f}}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

$$\mathbf{y}^T = \begin{bmatrix} x_1 w_{11} + x_2 w_{21} + x_3 w_{31} \\ x_1 w_{12} + x_2 w_{22} + x_3 w_{32} \\ x_1 w_{13} + x_2 w_{23} + x_3 w_{33} \\ x_1 w_{14} + x_2 w_{24} + x_3 w_{34} \end{bmatrix}$$

$$\mathbf{y} = [-5 \quad 5 \quad 2 \quad -10]$$

$$\frac{\partial \mathbf{h}}{\partial \mathbf{y}} = [0.028 \quad 0.01 \quad 0.22 \quad 0]$$

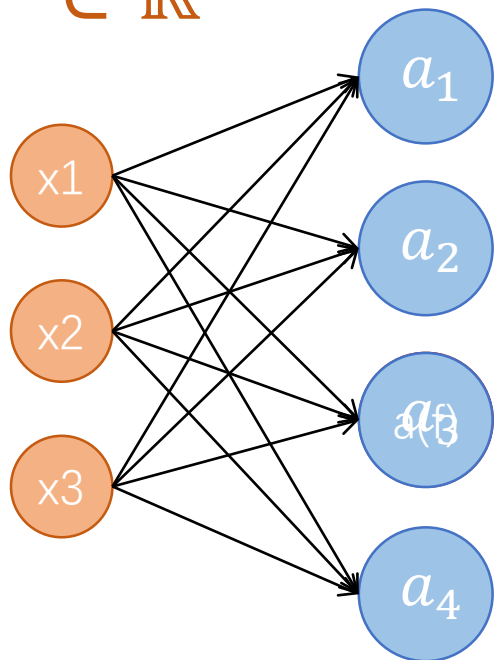
Sigmoid

$$\sigma = [0.007 \quad 0.99 \quad 0.88 \quad 0]$$

$$\frac{\partial \mathbf{h}}{\partial \sigma} = [4 \quad 1 \quad 2 \quad -1]$$

# 反向传播的矩阵运算

$$\mathbf{x}^T \in \mathbb{R}^D \quad \mathbf{a}^T \in \mathbb{R}^H$$



$$\mathbf{W}^T \in \mathbb{R}^{H \times D}$$

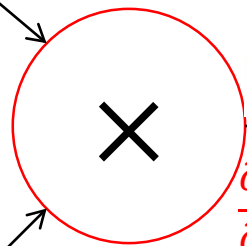
$$\mathbf{x} = [2 \quad 1 \quad -3]$$

雅可比矩阵  
(Jacobian matrix)

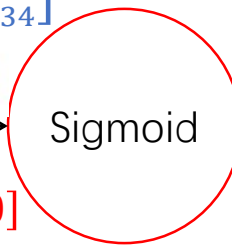
$$\mathbf{J} = \begin{bmatrix} \frac{\partial \mathbf{f}}{\partial x_1} & \cdots & \frac{\partial \mathbf{f}}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

$$\mathbf{y}^T = \begin{bmatrix} x_1 w_{11} + x_2 w_{21} + x_3 w_{31} \\ x_1 w_{12} + x_2 w_{22} + x_3 w_{32} \\ x_1 w_{13} + x_2 w_{23} + x_3 w_{33} \\ x_1 w_{14} + x_2 w_{24} + x_3 w_{34} \end{bmatrix}$$

$$\mathbf{y} = [-5 \quad 5 \quad 2 \quad -10]$$



$$\frac{\partial h}{\partial \mathbf{y}} = [0.028 \quad 0.01 \quad 0.22 \quad 0]$$



$$\sigma = [0.007 \quad 0.99 \quad 0.88 \quad 0]$$

$$\frac{\partial h}{\partial \sigma} = [4 \quad 1 \quad 2 \quad -1]$$

$$\mathbf{W} = \begin{bmatrix} 1 & 3 & 2 & 0 \\ 2 & -1 & 1 & 2 \\ 3 & 0 & 1 & 4 \end{bmatrix}$$

$$\left(\frac{\partial y}{\partial x}\right)_{n,m} = \left(\frac{\partial y_n}{\partial x_m}\right)$$

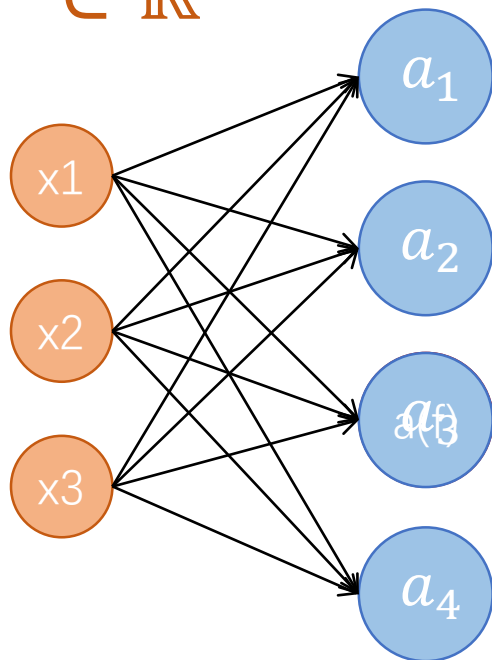


$$\frac{\partial y}{\partial x} \in \mathbb{R}^{H \times D}$$

和  $\mathbf{W}^T$  形状一样!

# 反向传播的矩阵运算

$$\mathbf{x}^T \in \mathbb{R}^D \quad \mathbf{a}^T \in \mathbb{R}^H$$



$$\mathbf{W}^T \in \mathbb{R}^{H \times D}$$

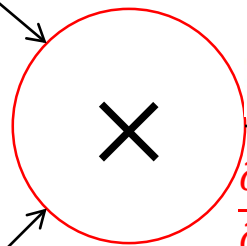
$$\mathbf{x} = [2 \quad 1 \quad -3]$$

雅可比矩阵  
(Jacobian matrix)

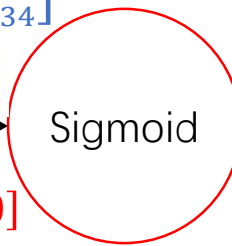
$$\mathbf{J} = \begin{bmatrix} \frac{\partial \mathbf{f}}{\partial x_1} & \cdots & \frac{\partial \mathbf{f}}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

$$\mathbf{y}^T = \begin{bmatrix} x_1 w_{11} + x_2 w_{21} + x_3 w_{31} \\ x_1 w_{12} + x_2 w_{22} + x_3 w_{32} \\ x_1 w_{13} + x_2 w_{23} + x_3 w_{33} \\ x_1 w_{14} + x_2 w_{24} + x_3 w_{34} \end{bmatrix}$$

$$\mathbf{y} = [-5 \quad 5 \quad 2 \quad -10]$$



$$\frac{\partial \mathbf{h}}{\partial \mathbf{y}} = [0.028 \quad 0.01 \quad 0.22 \quad 0]$$



$$\sigma = [0.007 \quad 0.99 \quad 0.88 \quad 0]$$

$$\frac{\partial \mathbf{h}}{\partial \sigma} = [4 \quad 1 \quad 2 \quad -1]$$

$$\left( \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)_{n,m} = \left( \frac{\partial y_n}{\partial x_m} \right)$$



$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \in \mathbb{R}^{H \times D}$$

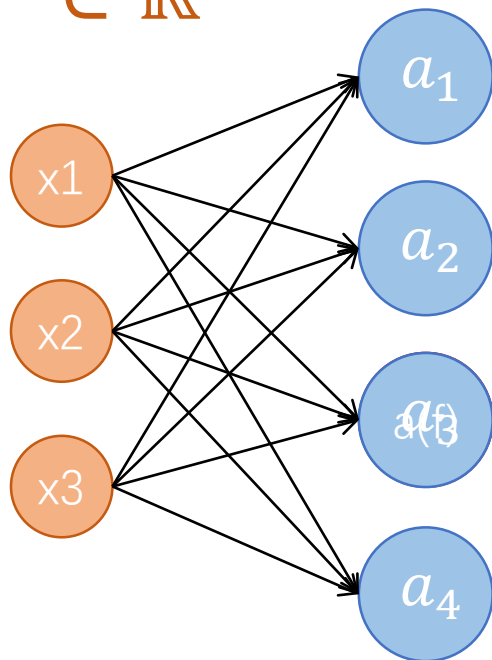
和  $\mathbf{W}^T$  形状一样!

$$\mathbf{W} = \begin{bmatrix} 1 & 3 & 2 & 0 \\ 2 & -1 & 1 & 2 \\ 3 & 0 & 1 & 4 \end{bmatrix}$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} w_{11} & w_{21} & w_{31} \\ w_{12} & w_{22} & w_{32} \\ w_{13} & w_{23} & w_{33} \\ w_{14} & w_{24} & w_{34} \end{bmatrix} = \mathbf{W}^T !$$

# 反向传播的矩阵运算

$$\mathbf{x}^T \in \mathbb{R}^D \quad \mathbf{a}^T \in \mathbb{R}^H$$



$$\mathbf{W}^T \in \mathbb{R}^{H \times D}$$

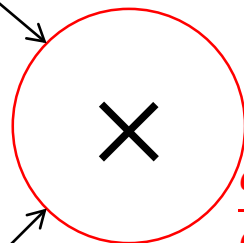
$$\mathbf{x} = [2 \quad 1 \quad -3]$$

雅可比矩阵  
(Jacobian matrix)

$$\mathbf{J} = \begin{bmatrix} \frac{\partial \mathbf{f}}{\partial x_1} & \cdots & \frac{\partial \mathbf{f}}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

$$\mathbf{y}^T = \begin{bmatrix} x_1 w_{11} + x_2 w_{21} + x_3 w_{31} \\ x_1 w_{12} + x_2 w_{22} + x_3 w_{32} \\ x_1 w_{13} + x_2 w_{23} + x_3 w_{33} \\ x_1 w_{14} + x_2 w_{24} + x_3 w_{34} \end{bmatrix}$$

$$\mathbf{y} = [-5 \quad 5 \quad 2 \quad -10]$$



$$\frac{\partial \mathbf{h}}{\partial \mathbf{y}} = [0.028 \quad 0.01 \quad 0.22 \quad 0]$$

Sigmoid

$$\sigma = [0.007 \quad 0.99 \quad 0.88 \quad 0]$$

$$\frac{\partial \mathbf{h}}{\partial \sigma} = [4 \quad 1 \quad 2 \quad -1]$$

$$\mathbf{W} = \begin{bmatrix} 1 & 3 & 2 & 0 \\ 2 & -1 & 1 & 2 \\ 3 & 0 & 1 & 4 \end{bmatrix}$$

$$\left( \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)_{n,m} = \left( \frac{\partial y_n}{\partial x_m} \right)$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \in \mathbb{R}^{H \times D} \quad \text{和 } \mathbf{W}^T \text{ 形状一样!}$$

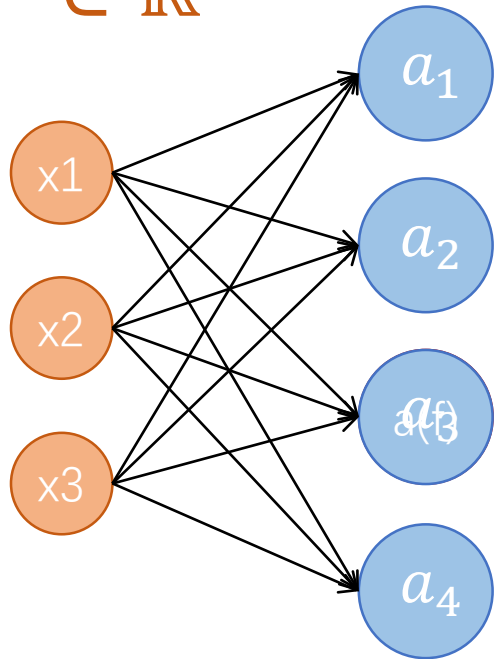
$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} w_{11} & w_{21} & w_{31} \\ w_{12} & w_{22} & w_{32} \\ w_{13} & w_{23} & w_{33} \\ w_{14} & w_{24} & w_{34} \end{bmatrix} = \mathbf{W}^T !$$

$$\frac{\partial \mathbf{h}}{\partial x_i} = \frac{\partial \mathbf{h}}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial x_i} = \sum_j \frac{\partial \mathbf{h}}{\partial y_j} \frac{\partial y_j}{\partial x_i}$$

$$\frac{\partial \mathbf{h}}{\partial \mathbf{x}} = \frac{\partial \mathbf{h}}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \frac{\partial \mathbf{h}}{\partial \mathbf{y}} \mathbf{W}^T$$

# 反向传播的矩阵运算

$$\mathbf{x}^T \in \mathbb{R}^D \quad \mathbf{a}^T \in \mathbb{R}^H$$



$$\mathbf{W}^T \in \mathbb{R}^{H \times D}$$

$$\mathbf{x} = [2 \quad 1 \quad -3]$$

$$\frac{\partial h}{\partial \mathbf{x}} = [0.50 \quad 0.27 \quad 0.30]$$

和 $\mathbf{x}$ 的形状一样

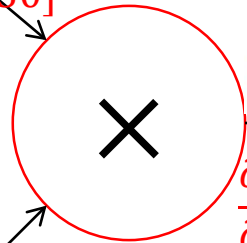
$$\mathbf{W} = \begin{bmatrix} 1 & 3 & 2 & 0 \\ 2 & -1 & 1 & 2 \\ 3 & 0 & 1 & 4 \end{bmatrix}$$

雅可比矩阵  
(Jacobian matrix)

$$\mathbf{J} = \begin{bmatrix} \frac{\partial \mathbf{f}}{\partial x_1} & \cdots & \frac{\partial \mathbf{f}}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

$$\mathbf{y}^T = \begin{bmatrix} x_1 w_{11} + x_2 w_{21} + x_3 w_{31} \\ x_1 w_{12} + x_2 w_{22} + x_3 w_{32} \\ x_1 w_{13} + x_2 w_{23} + x_3 w_{33} \\ x_1 w_{14} + x_2 w_{24} + x_3 w_{34} \end{bmatrix}$$

$$\mathbf{y} = [-5 \quad 5 \quad 2 \quad -10]$$



$$\frac{\partial h}{\partial \mathbf{y}} = [0.028 \quad 0.01 \quad 0.22 \quad 0]$$

$$\left( \frac{\partial y}{\partial x} \right)_{n,m} = \left( \frac{\partial y_n}{\partial x_m} \right)$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} w_{11} & w_{21} & w_{31} \\ w_{12} & w_{22} & w_{32} \\ w_{13} & w_{23} & w_{33} \\ w_{14} & w_{24} & w_{34} \end{bmatrix} = \mathbf{W}^T !$$

Sigmoid

$$\sigma = [0.007 \quad 0.99 \quad 0.88 \quad 0]$$

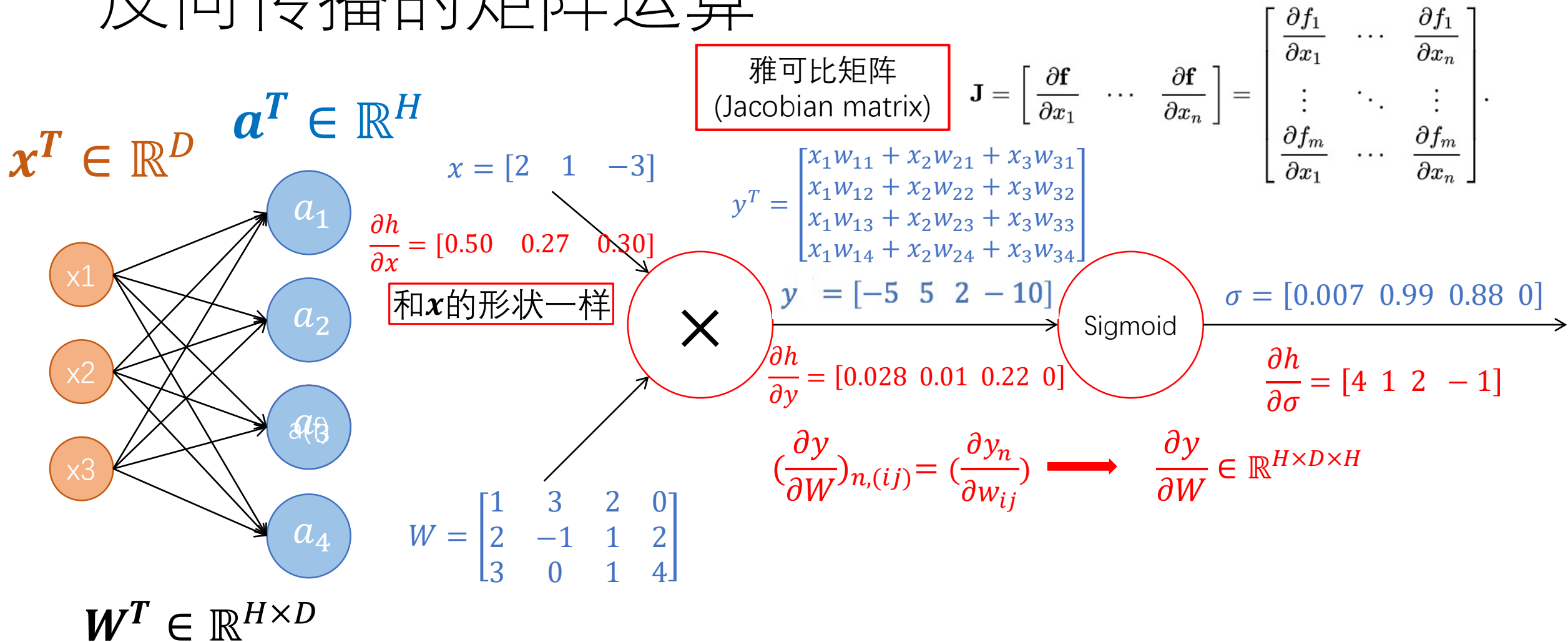
$$\frac{\partial h}{\partial \sigma} = [4 \quad 1 \quad 2 \quad -1]$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \in \mathbb{R}^{H \times D} \quad \text{和} \mathbf{W}^T \text{形状一样!}$$

$$\frac{\partial h}{\partial x_i} = \frac{\partial h}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial x_i} = \sum_j \frac{\partial h}{\partial y_j} \frac{\partial y_j}{\partial x_i}$$

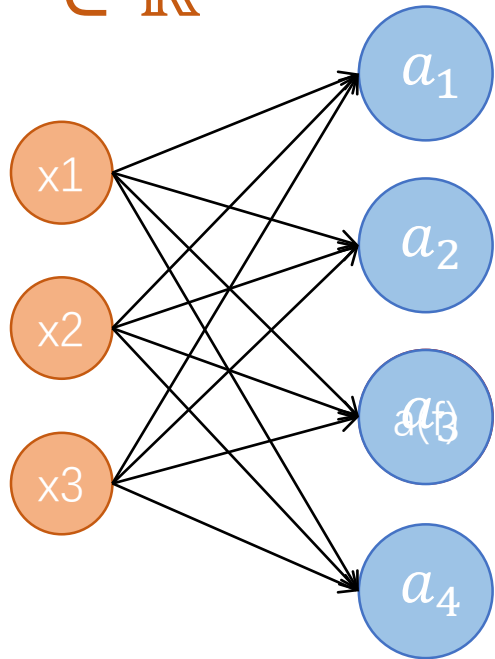
$$\frac{\partial h}{\partial \mathbf{x}} = \frac{\partial h}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \frac{\partial h}{\partial \mathbf{y}} \mathbf{W}^T$$

# 反向传播的矩阵运算



# 反向传播的矩阵运算

$$\mathbf{x}^T \in \mathbb{R}^D \quad \mathbf{a}^T \in \mathbb{R}^H$$



$$\mathbf{W}^T \in \mathbb{R}^{H \times D}$$

$$\mathbf{x} = [2 \quad 1 \quad -3]$$

$$\frac{\partial h}{\partial \mathbf{x}} = [0.50 \quad 0.27 \quad 0.30]$$

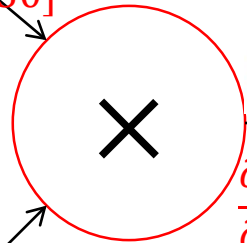
和 $\mathbf{x}$ 的形状一样

雅可比矩阵  
(Jacobian matrix)

$$\mathbf{J} = \begin{bmatrix} \frac{\partial \mathbf{f}}{\partial x_1} & \cdots & \frac{\partial \mathbf{f}}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

$$\mathbf{y}^T = \begin{bmatrix} x_1 w_{11} + x_2 w_{21} + x_3 w_{31} \\ x_1 w_{12} + x_2 w_{22} + x_3 w_{32} \\ x_1 w_{13} + x_2 w_{23} + x_3 w_{33} \\ x_1 w_{14} + x_2 w_{24} + x_3 w_{34} \end{bmatrix}$$

$$\mathbf{y} = [-5 \quad 5 \quad 2 \quad -10]$$



$$\frac{\partial h}{\partial \mathbf{y}} = [0.028 \quad 0.01 \quad 0.22 \quad 0]$$

Sigmoid

$$\sigma = [0.007 \quad 0.99 \quad 0.88 \quad 0]$$

$$\frac{\partial h}{\partial \sigma} = [4 \quad 1 \quad 2 \quad -1]$$

$$\mathbf{W} = \begin{bmatrix} 1 & 3 & 2 & 0 \\ 2 & -1 & 1 & 2 \\ 3 & 0 & 1 & 4 \end{bmatrix}$$

$$\left( \frac{\partial \mathbf{y}}{\partial \mathbf{W}} \right)_{n,(ij)} = \left( \frac{\partial y_n}{\partial w_{ij}} \right) \longrightarrow \frac{\partial \mathbf{y}}{\partial \mathbf{W}} \in \mathbb{R}^{H \times D \times H}$$

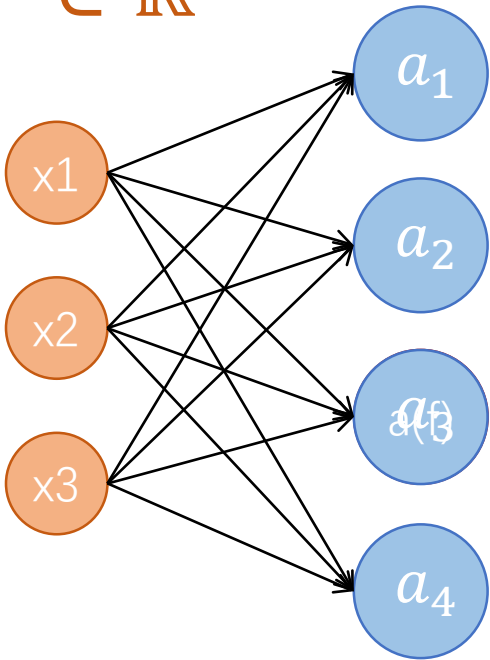
假设 $H=4096$ ,  $D=3072$ , 需要至少 $51539607552 \approx 5\text{G}$ 内存!!

$w_{ij}$ 只影响 $y_j$ !!!



# 反向传播的矩阵运算

$$\mathbf{x}^T \in \mathbb{R}^D \quad \mathbf{a}^T \in \mathbb{R}^H$$



$$\mathbf{W}^T \in \mathbb{R}^{H \times D}$$

$$\mathbf{x} = [2 \quad 1 \quad -3]$$

$$\frac{\partial h}{\partial \mathbf{x}} = [0.50 \quad 0.27 \quad 0.30]$$

和 $\mathbf{x}$ 的形状一样

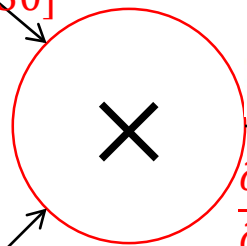
$$\mathbf{W} = \begin{bmatrix} 1 & 3 & 2 & 0 \\ 2 & -1 & 1 & 2 \\ 3 & 0 & 1 & 4 \end{bmatrix}$$

雅可比矩阵  
(Jacobian matrix)

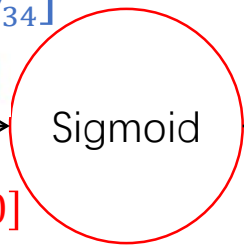
$$\mathbf{J} = \begin{bmatrix} \frac{\partial \mathbf{f}}{\partial x_1} & \cdots & \frac{\partial \mathbf{f}}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

$$\mathbf{y}^T = \begin{bmatrix} x_1 w_{11} + x_2 w_{21} + x_3 w_{31} \\ x_1 w_{12} + x_2 w_{22} + x_3 w_{32} \\ x_1 w_{13} + x_2 w_{23} + x_3 w_{33} \\ x_1 w_{14} + x_2 w_{24} + x_3 w_{34} \end{bmatrix}$$

$$\mathbf{y} = [-5 \quad 5 \quad 2 \quad -10]$$



$$\frac{\partial h}{\partial \mathbf{y}} = [0.028 \quad 0.01 \quad 0.22 \quad 0]$$



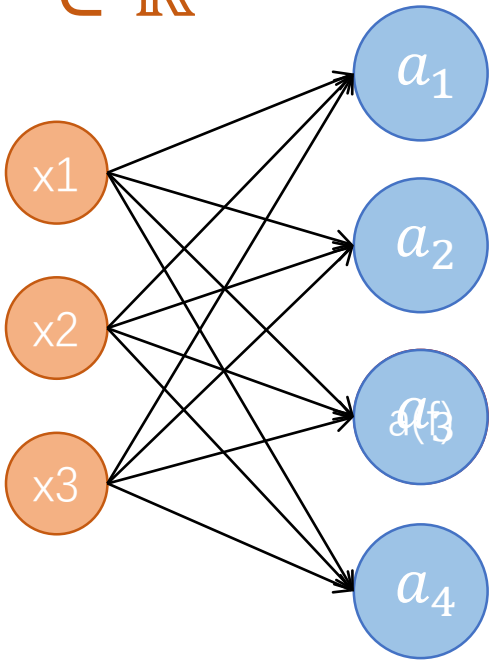
$$\sigma = [0.007 \quad 0.99 \quad 0.88 \quad 0]$$

$$\frac{\partial h}{\partial \sigma} = [4 \quad 1 \quad 2 \quad -1]$$

$$\frac{\partial h}{\partial w_{ij}} = \frac{\partial h}{\partial y_j} \frac{\partial y_j}{\partial w_{ij}} = \frac{\partial h}{\partial y_j} x_i = x_i \frac{\partial h}{\partial y_j}$$

# 反向传播的矩阵运算

$$\mathbf{x}^T \in \mathbb{R}^D \quad \mathbf{a}^T \in \mathbb{R}^H$$



$$\mathbf{W}^T \in \mathbb{R}^{H \times D}$$

$$\mathbf{x} = [2 \quad 1 \quad -3]$$

$$\frac{\partial h}{\partial \mathbf{x}} = [0.50 \quad 0.27 \quad 0.30]$$

和 $\mathbf{x}$ 的形状一样

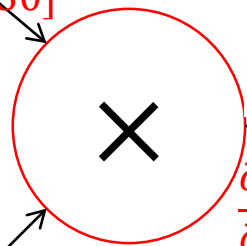
$$\mathbf{W} = \begin{bmatrix} 1 & 3 & 2 & 0 \\ 2 & -1 & 1 & 2 \\ 3 & 0 & 1 & 4 \end{bmatrix}$$

雅可比矩阵  
(Jacobian matrix)

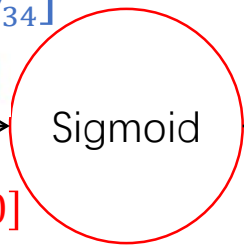
$$\mathbf{J} = \begin{bmatrix} \frac{\partial \mathbf{f}}{\partial x_1} & \cdots & \frac{\partial \mathbf{f}}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

$$\mathbf{y}^T = \begin{bmatrix} x_1 w_{11} + x_2 w_{21} + x_3 w_{31} \\ x_1 w_{12} + x_2 w_{22} + x_3 w_{32} \\ x_1 w_{13} + x_2 w_{23} + x_3 w_{33} \\ x_1 w_{14} + x_2 w_{24} + x_3 w_{34} \end{bmatrix}$$

$$\mathbf{y} = [-5 \quad 5 \quad 2 \quad -10]$$



$$\frac{\partial h}{\partial \mathbf{y}} = [0.028 \quad 0.01 \quad 0.22 \quad 0]$$



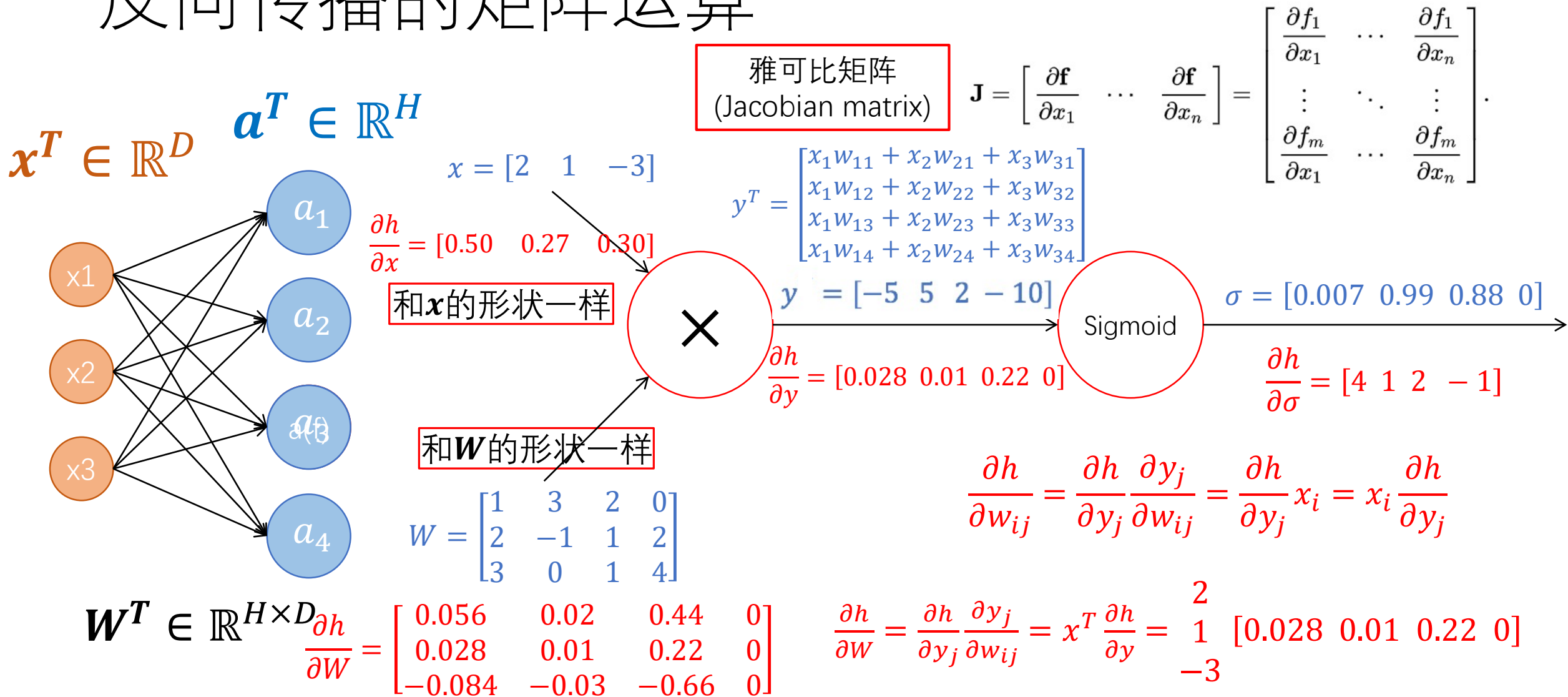
$$\sigma = [0.007 \quad 0.99 \quad 0.88 \quad 0]$$

$$\frac{\partial h}{\partial \sigma} = [4 \quad 1 \quad 2 \quad -1]$$

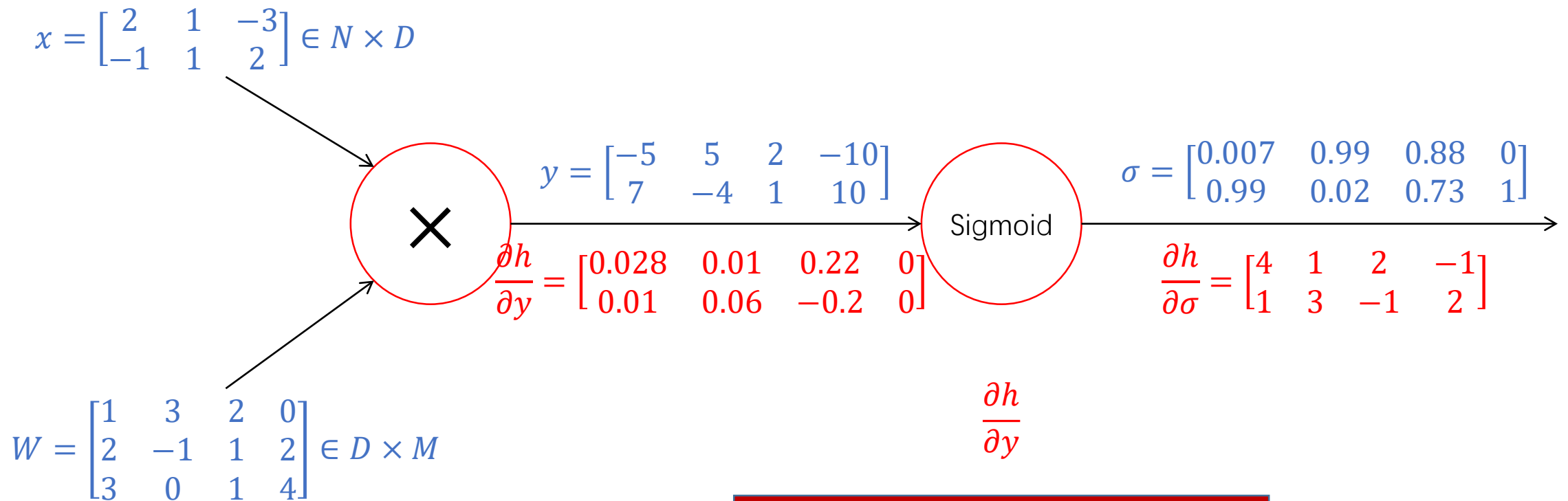
$$\frac{\partial h}{\partial w_{ij}} = \frac{\partial h}{\partial y_j} \frac{\partial y_j}{\partial w_{ij}} = \frac{\partial h}{\partial y_j} x_i = x_i \frac{\partial h}{\partial y_j}$$

$$\frac{\partial h}{\partial \mathbf{W}} = \frac{\partial h}{\partial y_j} \frac{\partial y_j}{\partial w_{ij}} = \mathbf{x}^T \frac{\partial h}{\partial \mathbf{y}} = \begin{bmatrix} 2 \\ 1 \\ -3 \end{bmatrix} [0.028 \quad 0.01 \quad 0.22 \quad 0]$$

# 反向传播的矩阵运算

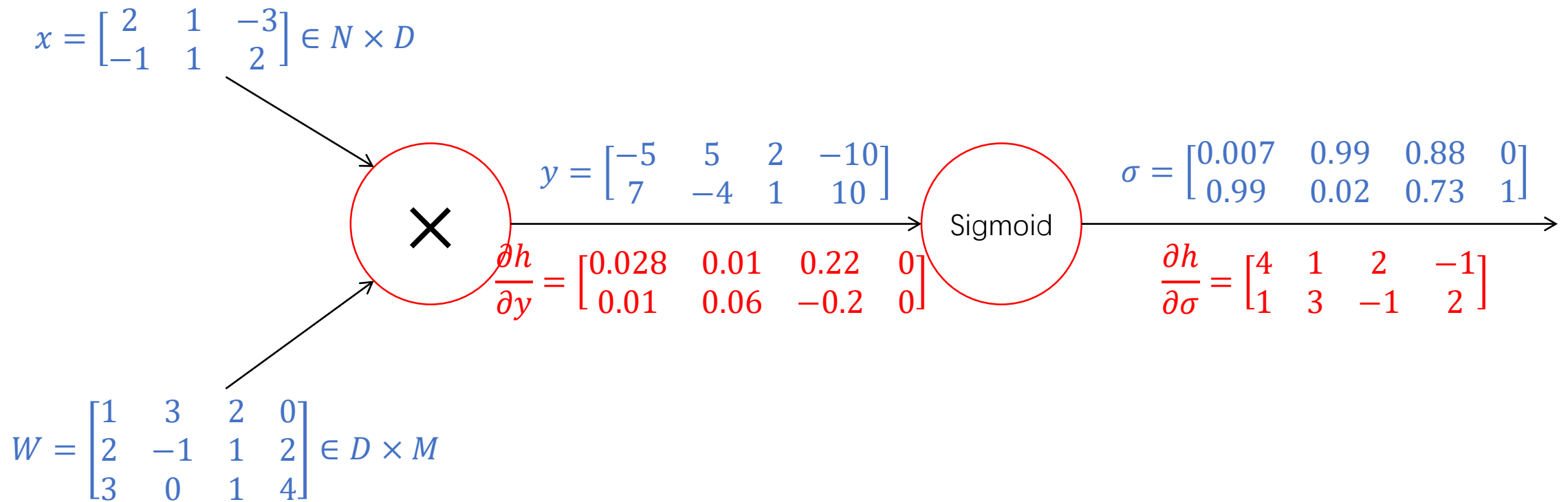


# 反向传播的矩阵运算：minibatch



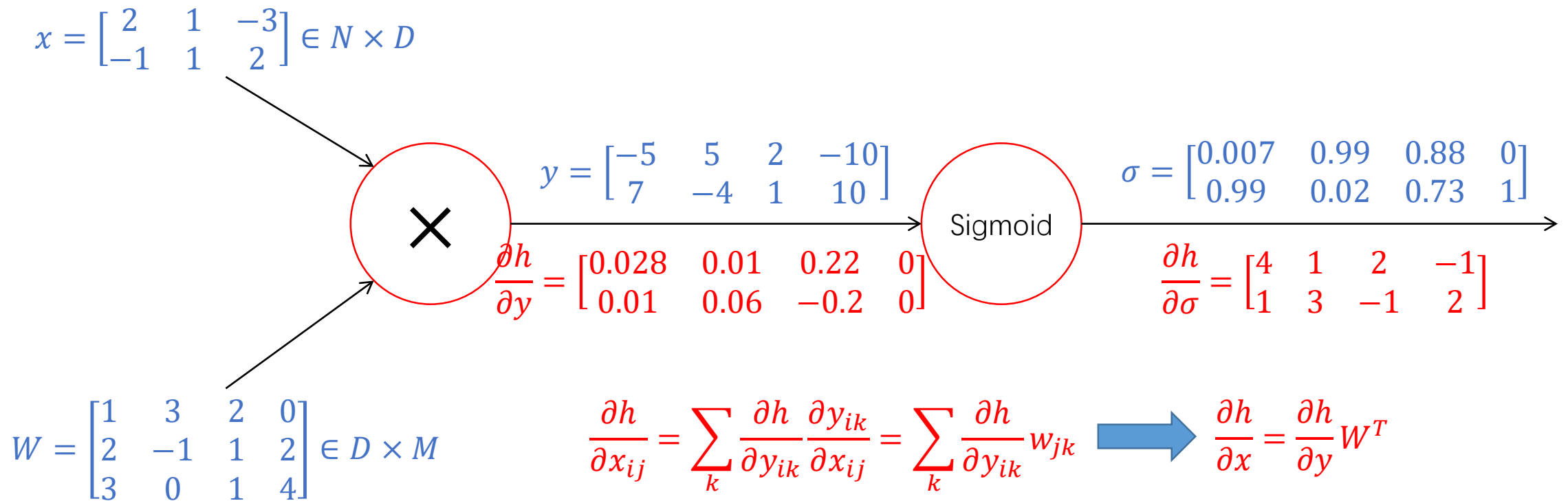
$y_{ij}$  只影响  $\sigma_{ij}$ , 只需 pairwise 计算

# 反向传播的矩阵运算：minibatch



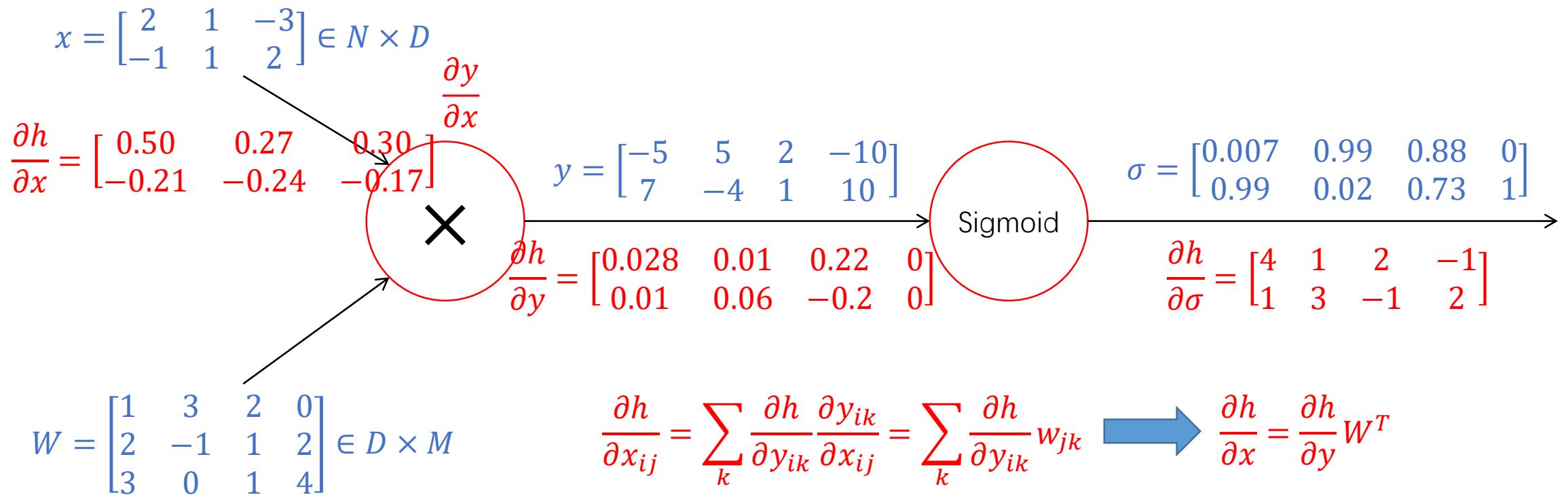
$x_{ij}$  只影响  $y$  的第  $i$  行

# 反向传播的矩阵运算：minibatch



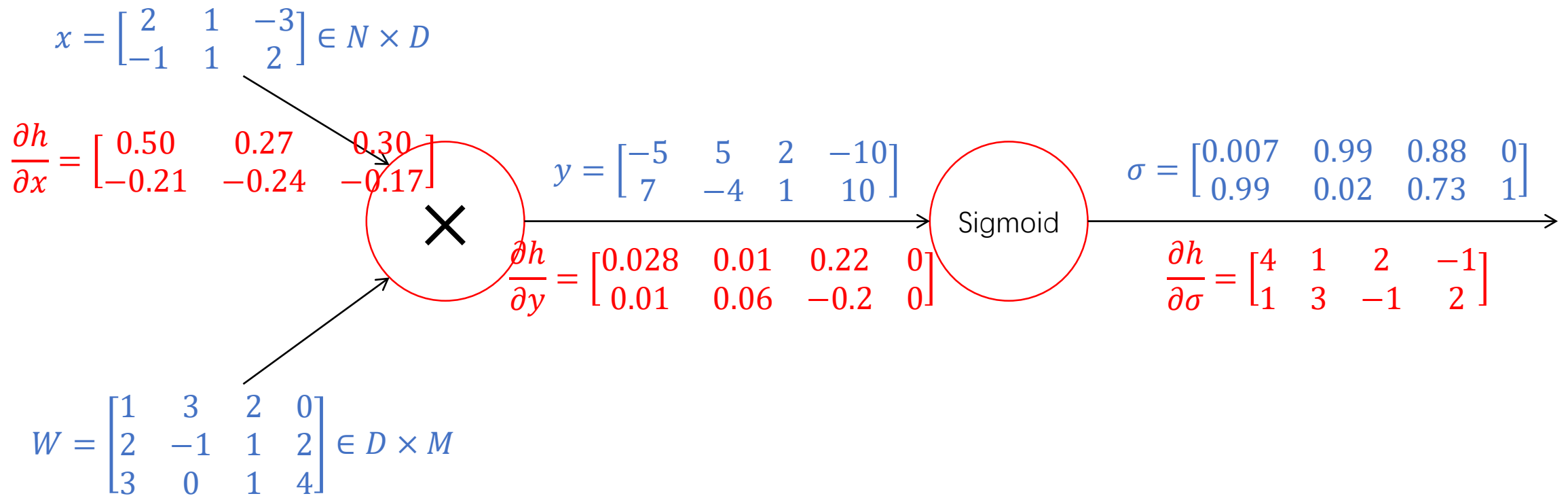
$x_{ij}$  只影响  $y$  的第  $i$  行

# 反向传播的矩阵运算：minibatch



$x_{ij}$  只影响  $y$  的第  $i$  行

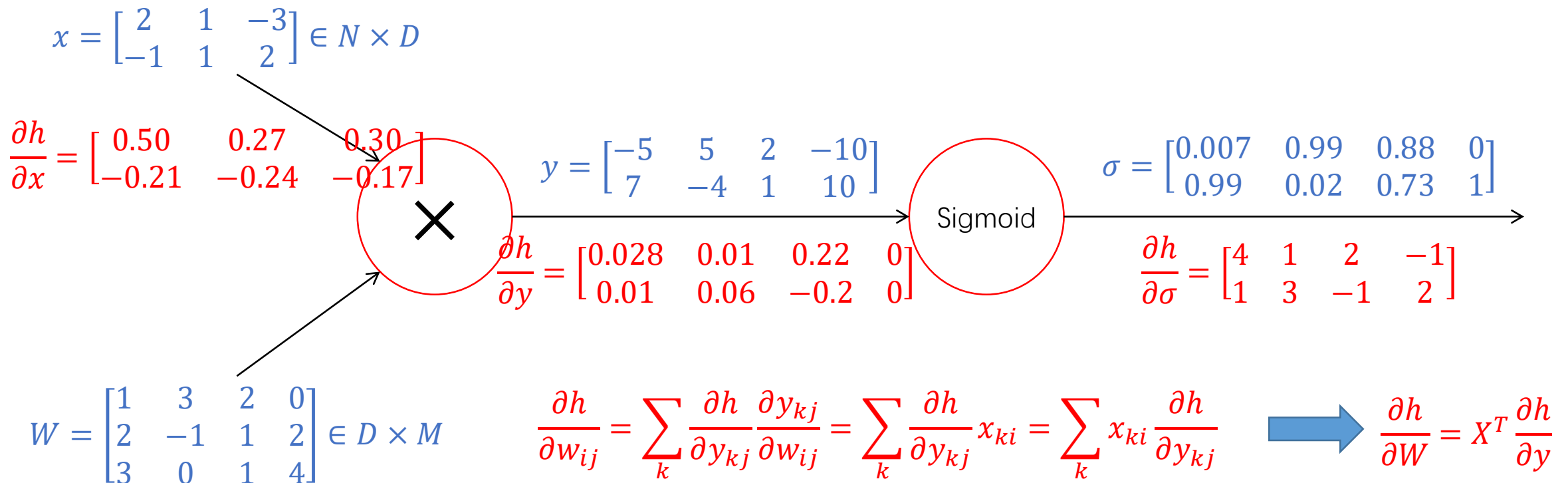
# 反向传播的矩阵运算：minibatch



$w_{ij}$  只影响  $y$  的第  $j$  列

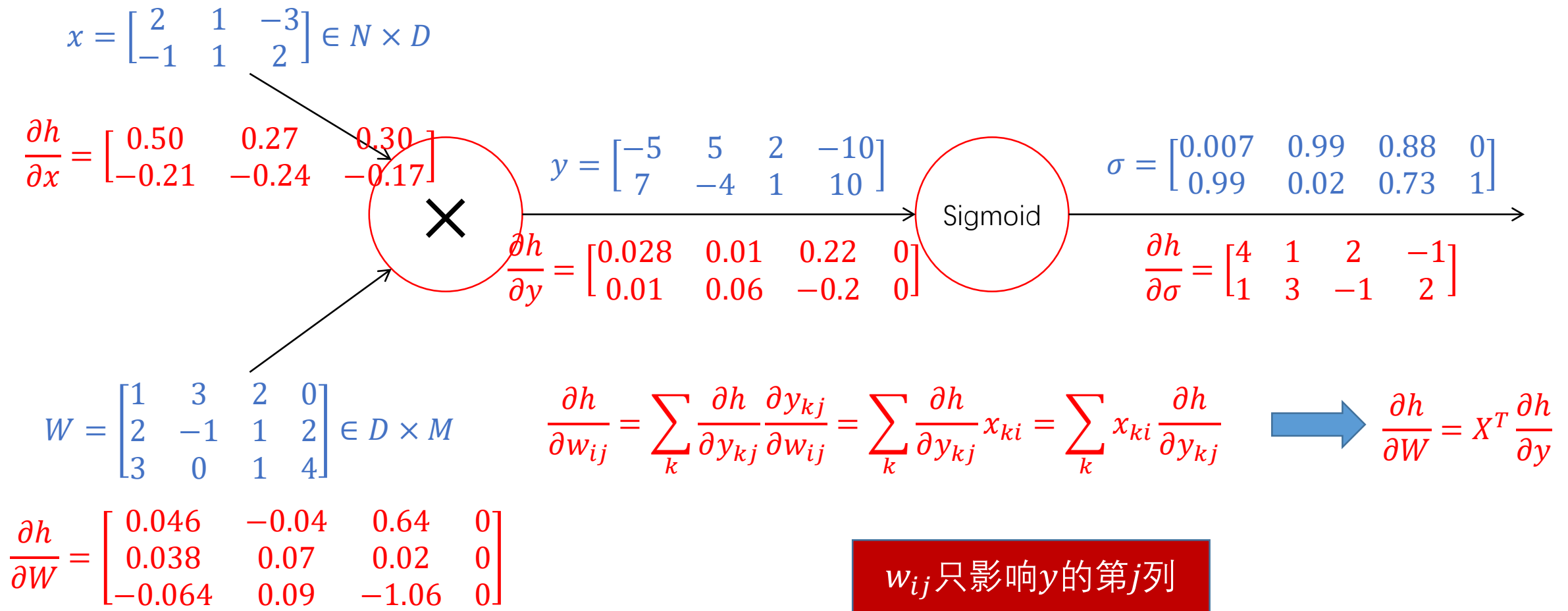


# 反向传播的矩阵运算：minibatch

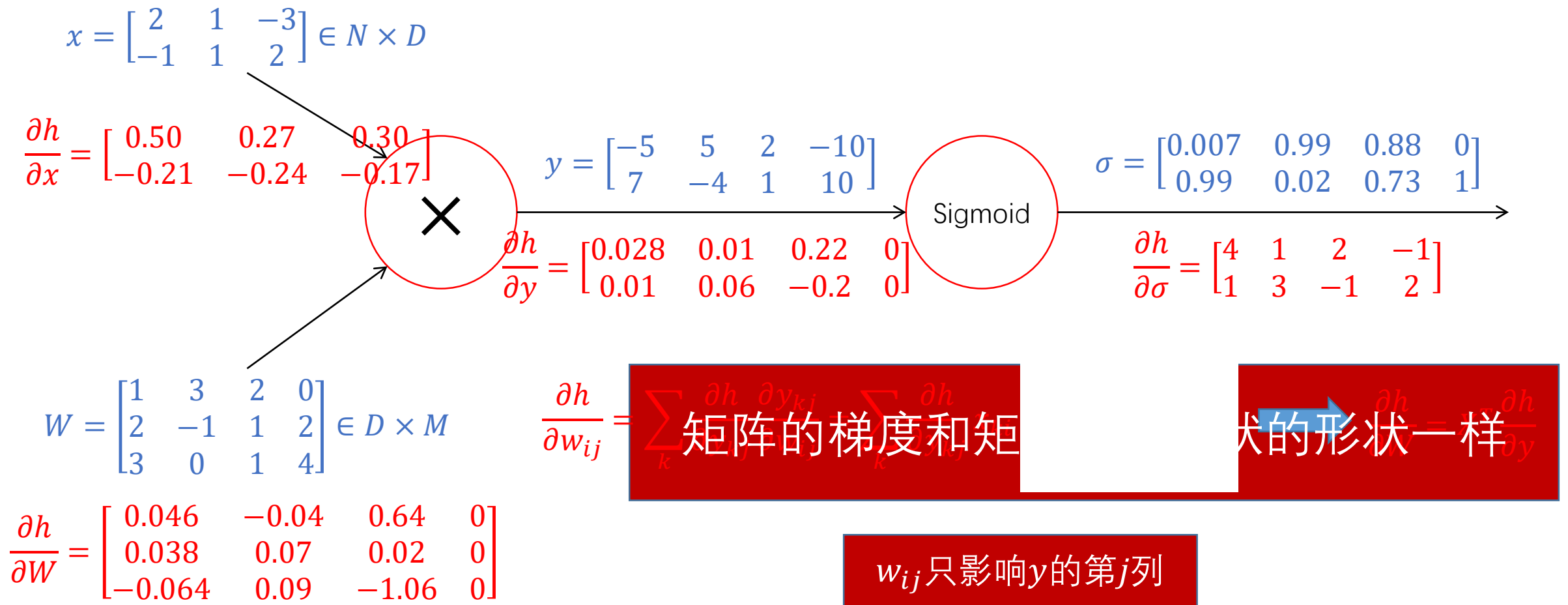


$w_{ij}$  只影响  $y$  的第  $j$  列

# 反向传播的矩阵运算：minibatch

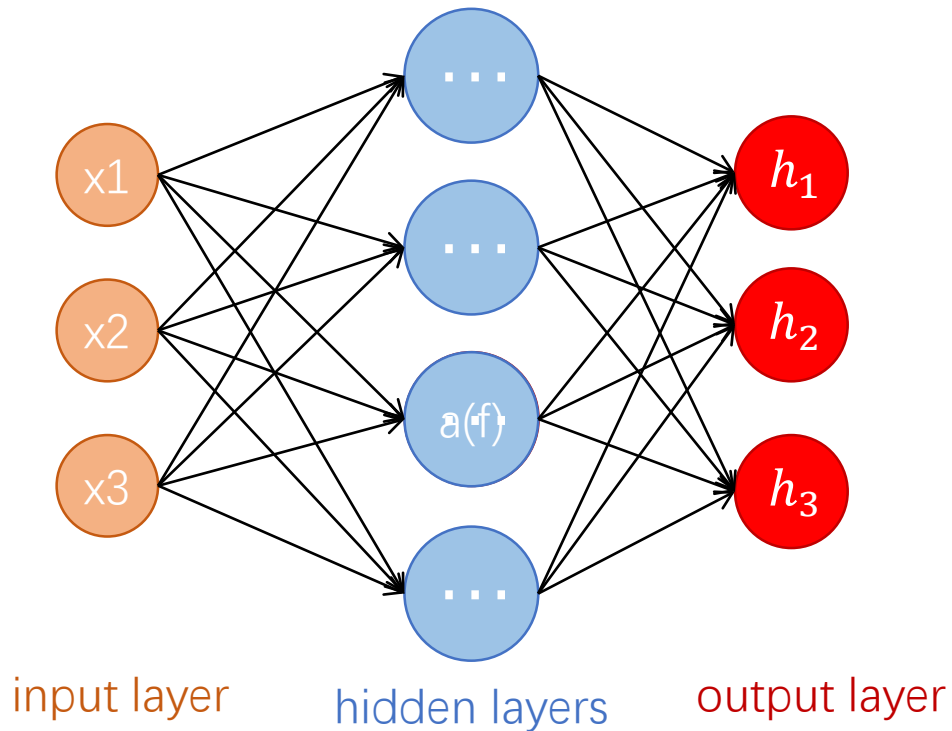


# 反向传播的矩阵运算：minibatch



# 多层神经网络反向传播

Chain rule



$$\frac{\partial L}{\partial \mathbf{W}^l} = \frac{\partial L}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{W}^l}$$

$$\frac{\partial L}{\partial \mathbf{W}^{l-1}} = \frac{\partial L}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{a}^l} \frac{\partial \mathbf{a}^l}{\partial \mathbf{W}^{l-1}}$$

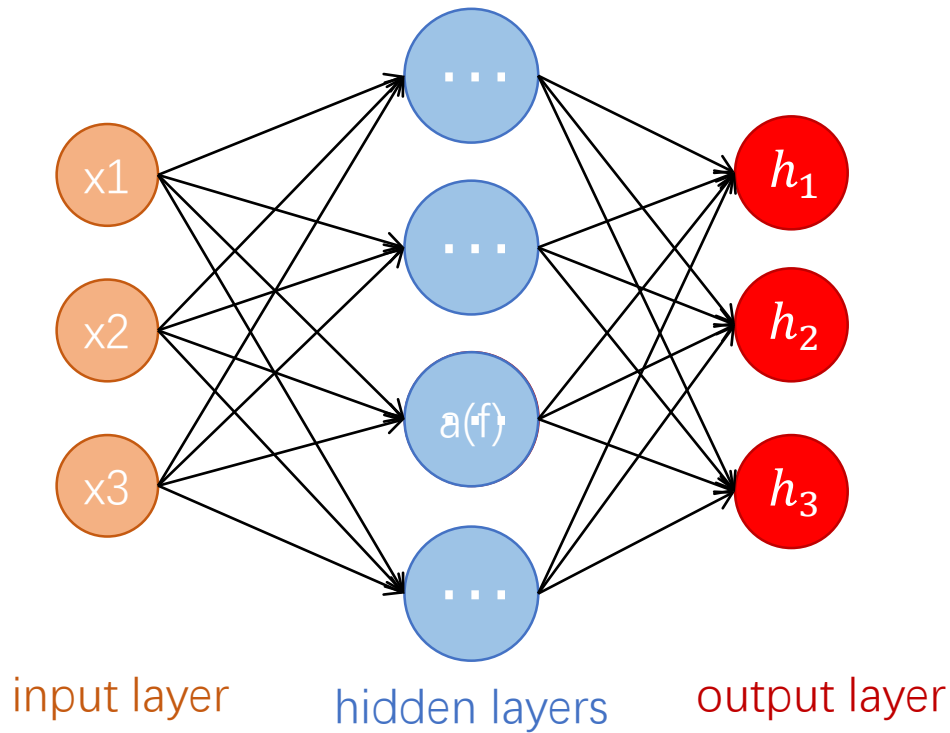
$$\frac{\partial L}{\partial \mathbf{W}^{l-2}} = \frac{\partial L}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{a}^l} \frac{\partial \mathbf{a}^l}{\partial \mathbf{a}^{l-1}} \frac{\partial \mathbf{a}^{l-1}}{\partial \mathbf{W}^{l-2}}$$

⋮

$$\frac{\partial L}{\partial \mathbf{W}^1} = \frac{\partial L}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{a}^l} \frac{\partial \mathbf{a}^l}{\partial \mathbf{a}^{l-1}} \frac{\partial \mathbf{a}^{l-1}}{\partial \mathbf{a}^{l-2}} \cdots \cdots \frac{\partial \mathbf{a}^2}{\partial \mathbf{W}^1}$$

# 多层神经网络反向传播

Chain rule



$$\frac{\partial L}{\partial \mathbf{W}^l} = \frac{\partial L}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{W}^l}$$

$$\frac{\partial L}{\partial \mathbf{W}^{l-1}} = \frac{\partial L}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{a}^l} \frac{\partial \mathbf{a}^l}{\partial \mathbf{W}^{l-1}}$$

$$\frac{\partial L}{\partial \mathbf{W}^{l-2}} = \frac{\partial L}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{a}^l} \frac{\partial \mathbf{a}^l}{\partial \mathbf{a}^{l-1}} \frac{\partial \mathbf{a}^{l-1}}{\partial \mathbf{W}^{l-2}}$$

⋮

不要忘记每一层的偏置项权重 $b$ , 以及正则项!

$$\frac{\partial L}{\partial \mathbf{W}^1} = \frac{\partial L}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{a}^l} \frac{\partial \mathbf{a}^l}{\partial \mathbf{a}^{l-1}} \frac{\partial \mathbf{a}^{l-1}}{\partial \mathbf{a}^{l-2}} \cdots \cdots \frac{\partial \mathbf{a}^2}{\partial \mathbf{W}^1}$$

# Assignment1: 两层神经网络实现

✓loss()函数： 返回loss和梯度

✓train()函数： 梯度下降优化， 调整超参数

✓Predict()函数： 测试集分类

```

# Compute the forward pass
scores = None
#####
# TODO: Perform the forward pass, computing the class scores for the input. #
# Store the result in the scores variable, which should be an array of #
# shape (N, C). #
#####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

pass

# *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

# If the targets are not given then jump out, we're done
if y is None:
    return scores

# Compute the loss
loss = None
#####
# TODO: Finish the forward pass, and compute the loss. This should include #
# both the data loss and L2 regularization for W1 and W2. Store the result #
# in the variable loss, which should be a scalar. Use the Softmax #
# classifier loss. #
#####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

pass

# *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

# Backward pass: compute gradients
grads = {}
#####
# TODO: Compute the backward pass, computing the derivatives of the weights #
# and biases. Store the results in the grads dictionary. For example, #
# grads['W1'] should store the gradient on W1, and be a matrix of same size #
#####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

pass

# *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

return loss, grads

```

# 小结

- 神经网络
  - ✓ 激活函数
  - ✓ 神经网络的层次
  - ✓ 前馈计算
- 优化
  - ✓ 计算图
  - ✓ 反向传播
  - ✓ 链式规则

# L05

- Convolutional Neural Networks

