

# Backpropagation

# Gradient Descent

Network parameters  $\theta = \{w_1, w_2, \dots, b_1, b_2, \dots\}$

Starting Parameters  $\theta^0 \longrightarrow \theta^1 \longrightarrow \theta^2 \longrightarrow \dots$

$$\begin{aligned} \nabla L(\theta) &= \begin{bmatrix} \partial L(\theta) / \partial w_1 \\ \partial L(\theta) / \partial w_2 \\ \vdots \\ \partial L(\theta) / \partial b_1 \\ \partial L(\theta) / \partial b_2 \\ \vdots \end{bmatrix} \end{aligned}$$

Compute  $\nabla L(\theta^0)$        $\theta^1 = \theta^0 - \eta \nabla L(\theta^0)$

Compute  $\nabla L(\theta^1)$        $\theta^2 = \theta^1 - \eta \nabla L(\theta^1)$

Millions of parameters .....

To compute the gradients efficiently,  
we use **backpropagation**.

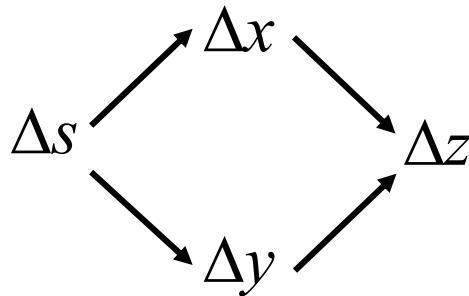
# Chain Rule

**Case 1**       $y = g(x) \quad z = h(y)$

$$\Delta x \rightarrow \Delta y \rightarrow \Delta z \qquad \frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

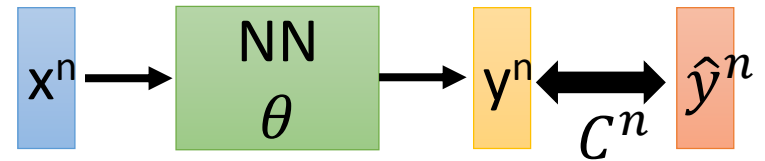
**Case 2**

$$x = g(s) \qquad y = h(s) \qquad z = k(x, y)$$

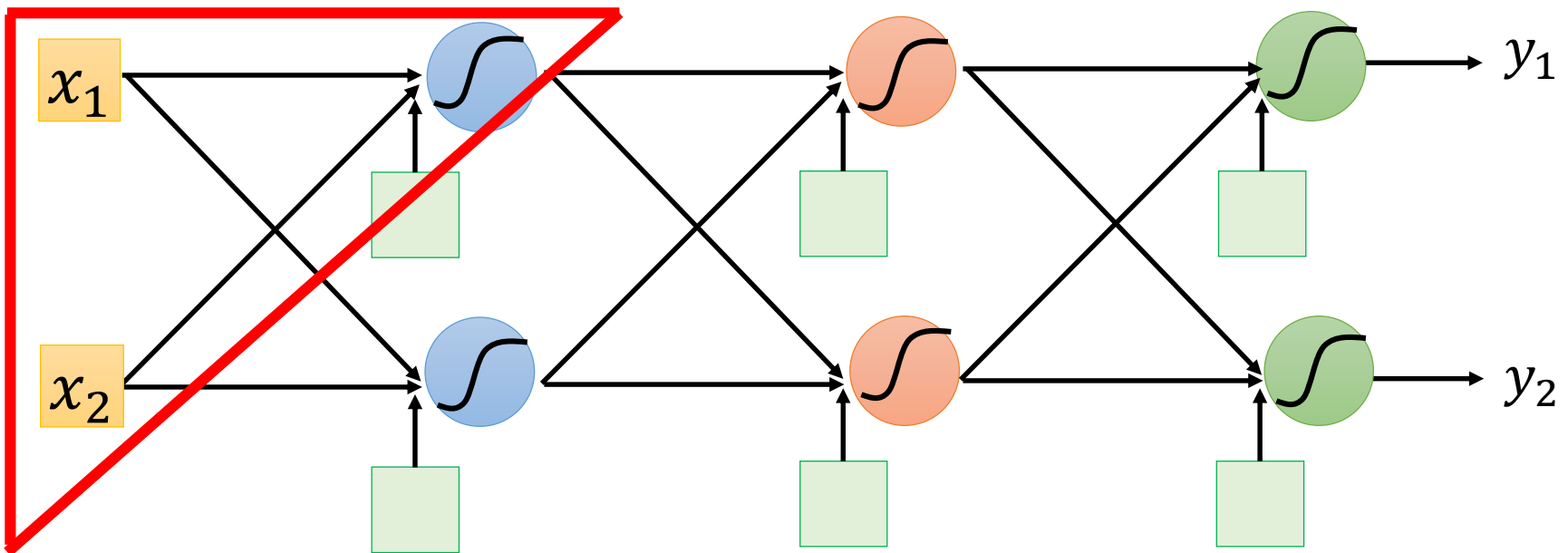


$$\frac{dz}{ds} = \frac{\partial z}{\partial x} \frac{dx}{ds} + \frac{\partial z}{\partial y} \frac{dy}{ds}$$

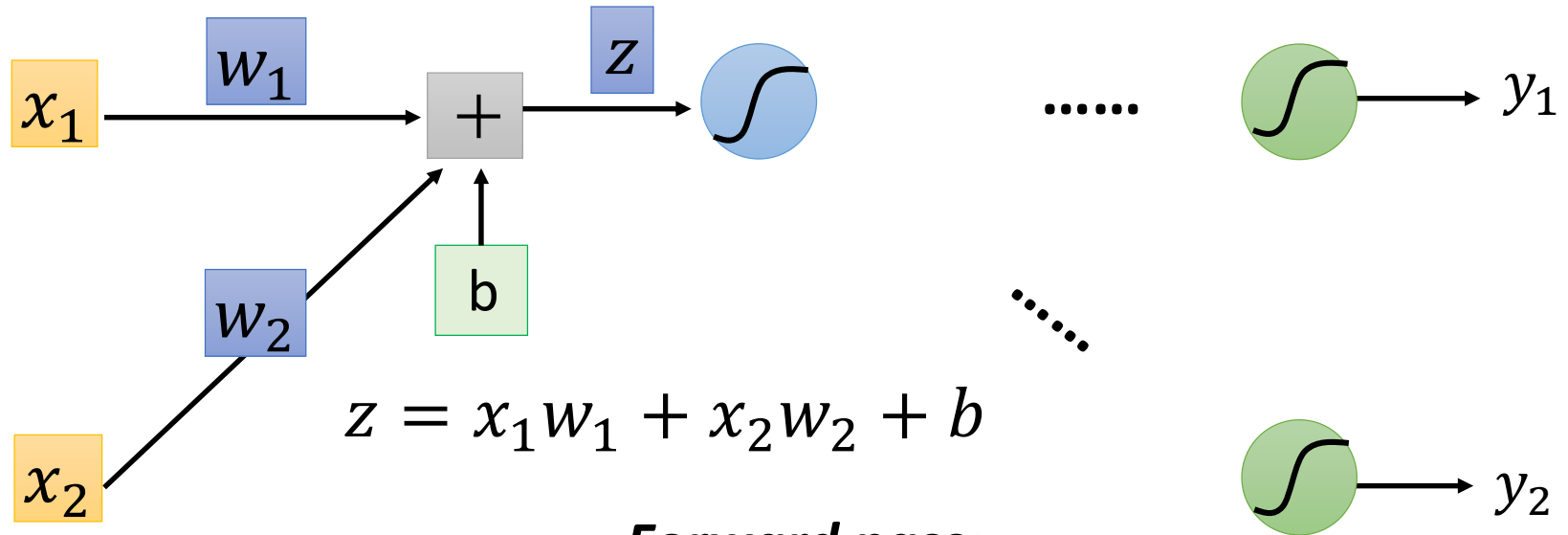
# Backpropagation



$$L(\theta) = \sum_{n=1}^N C^n(\theta) \quad \Rightarrow \quad \frac{\partial L(\theta)}{\partial w} = \sum_{n=1}^N \frac{\partial C^n(\theta)}{\partial w}$$



# Backpropagation



$$z = x_1 w_1 + x_2 w_2 + b$$

**Forward pass:**

Compute  $\partial z / \partial w$  for all parameters

**Backward pass:**

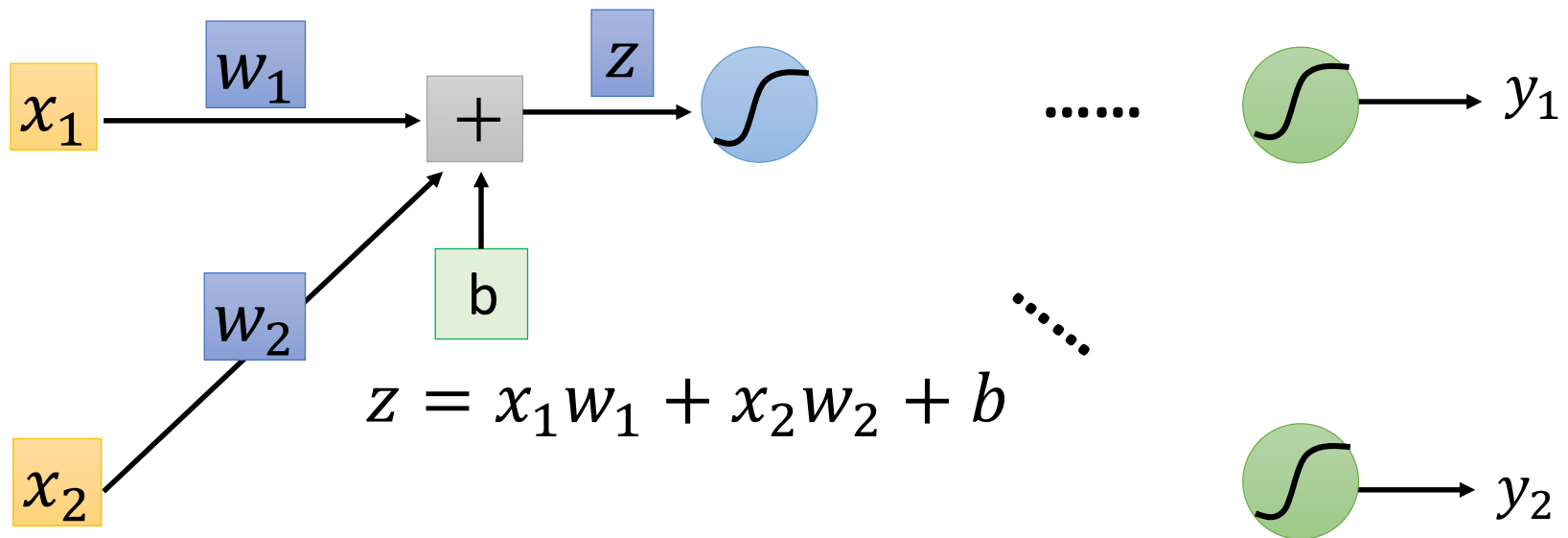
Compute  $\partial C / \partial z$  for all activation function inputs  $z$

$$\frac{\partial C}{\partial w} = ? \quad \frac{\partial z}{\partial w} \frac{\partial C}{\partial z}$$

(Chain rule)

# Backpropagation – Forward pass

Compute  $\partial z / \partial w$  for all parameters



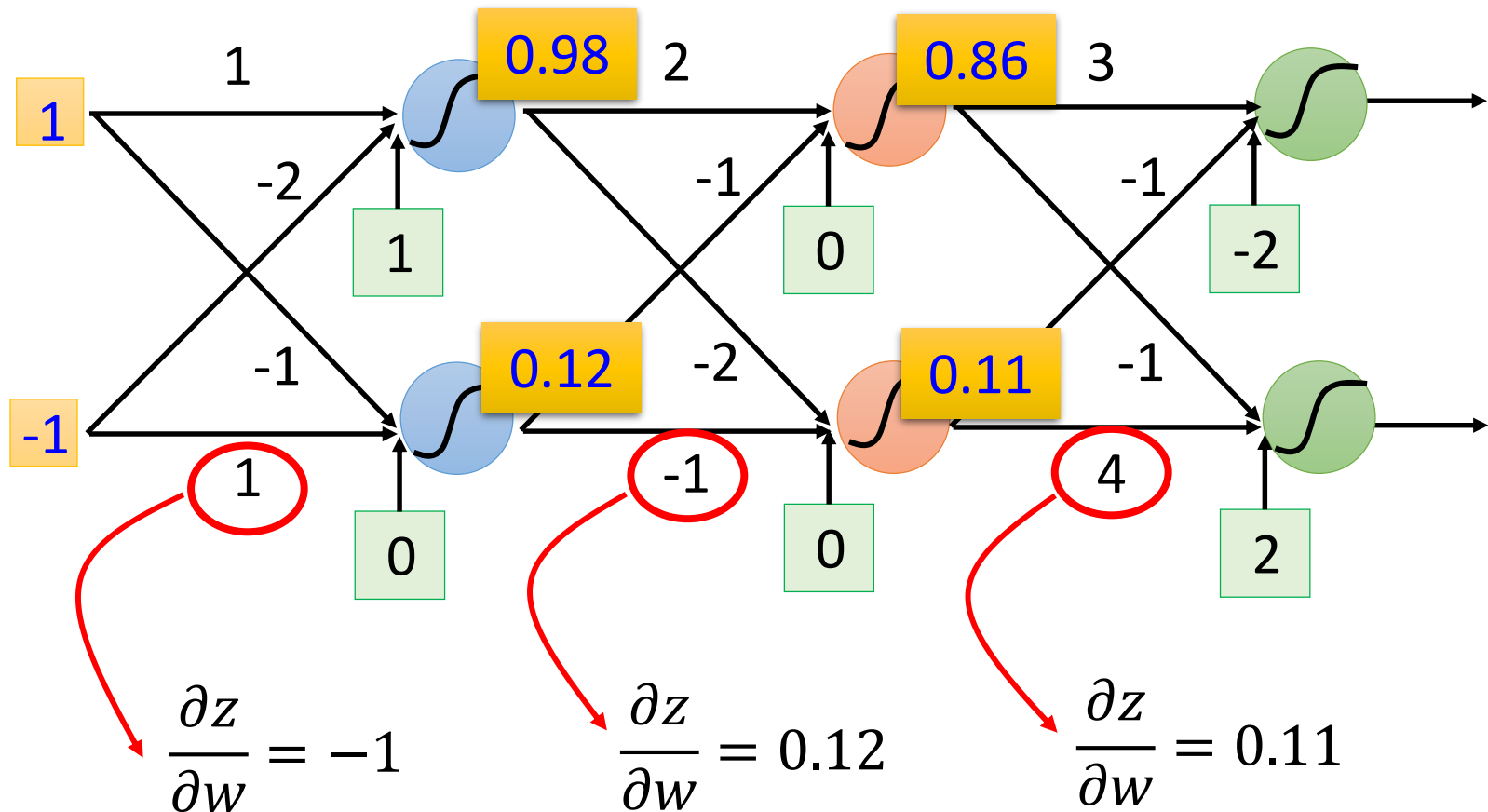
$$\partial z / \partial w_1 = ? \quad x_1$$

$$\partial z / \partial w_2 = ? \quad x_2$$

} The value of the input  
connected by the weight

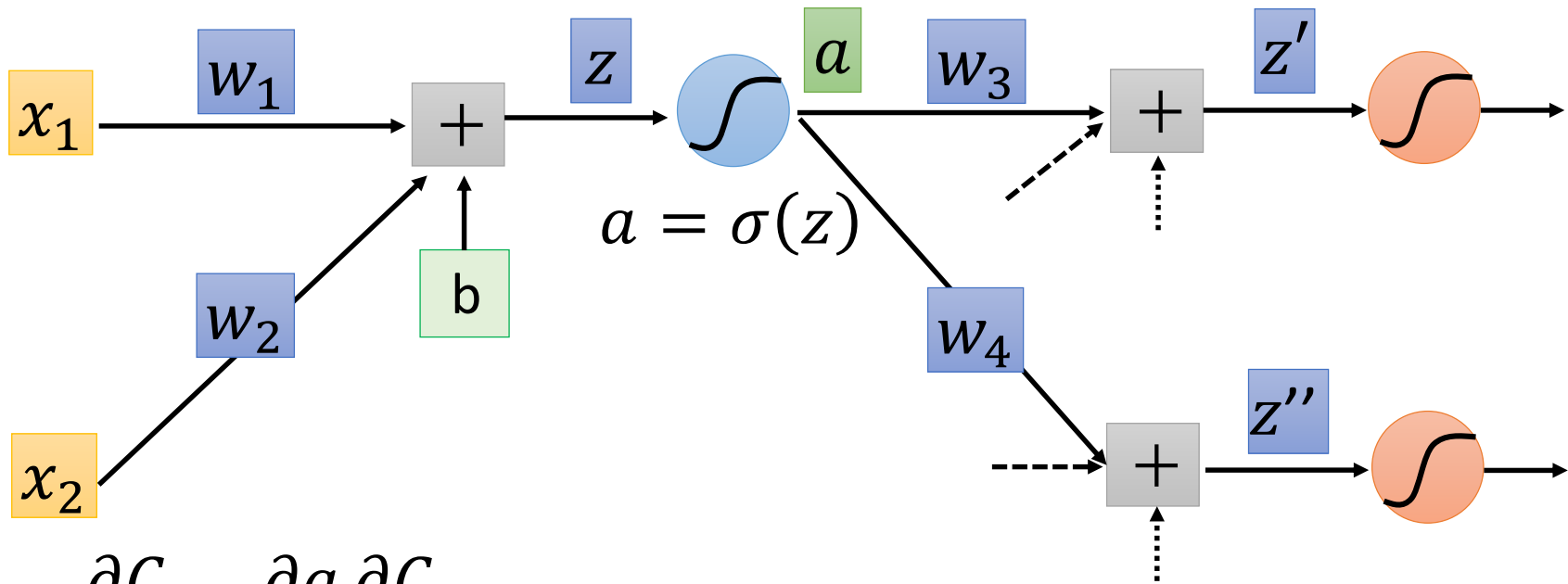
# Backpropagation – Forward pass

Compute  $\partial z / \partial w$  for all parameters



# Backpropagation – Backward pass

Compute  $\partial C / \partial z$  for all activation function inputs  $z$



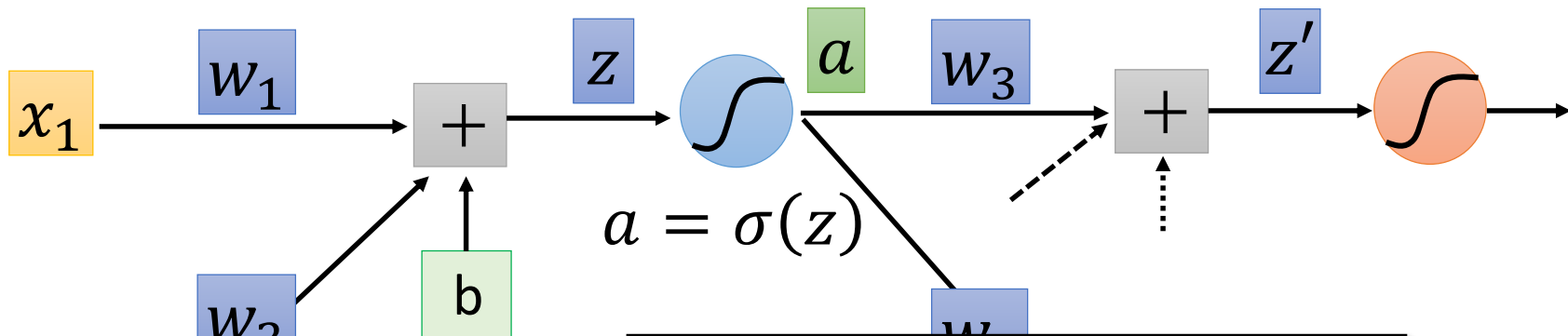
$$\frac{\partial C}{\partial z} = \frac{\partial a}{\partial z} \frac{\partial C}{\partial a}$$

➡  $\sigma'(z)$



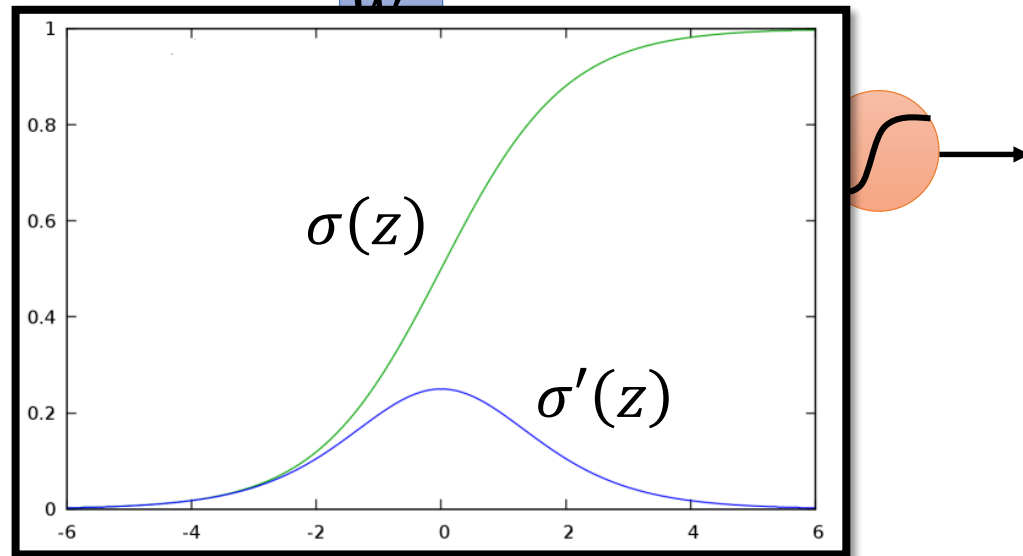
# Backpropagation – Backward pass

Compute  $\partial C / \partial z$  for all activation function inputs  $z$



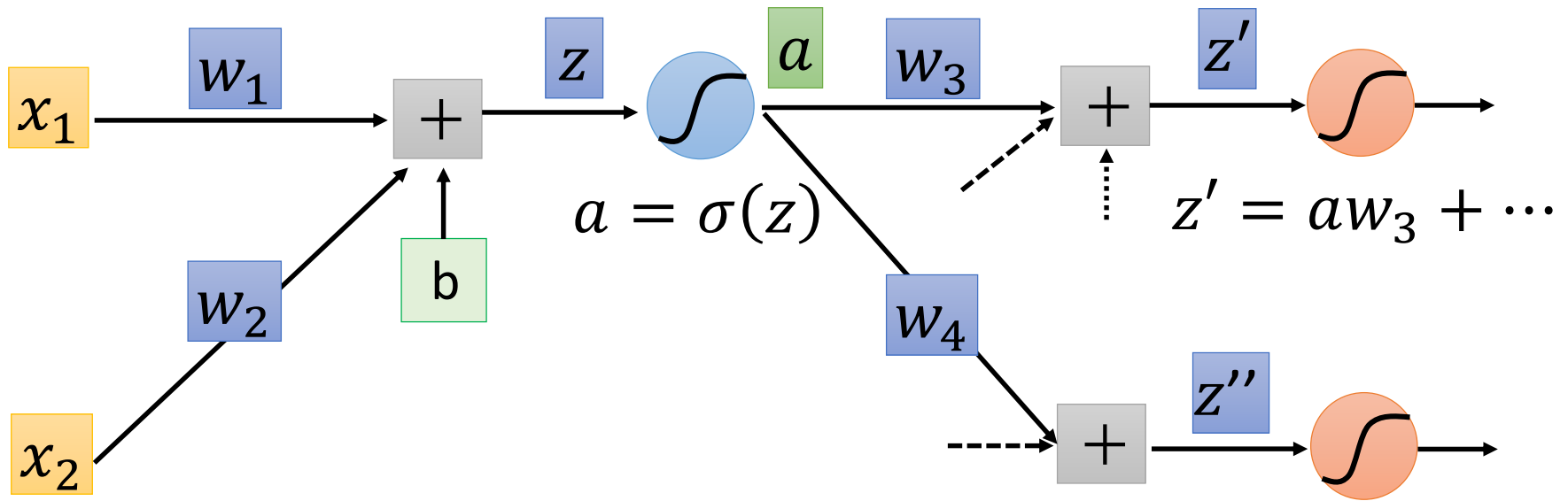
$$\frac{\partial C}{\partial z} = \frac{\partial a}{\partial z} \frac{\partial C}{\partial a}$$

➡  $\sigma'(z)$



# Backpropagation – Backward pass

Compute  $\partial C / \partial z$  for all activation function inputs  $z$

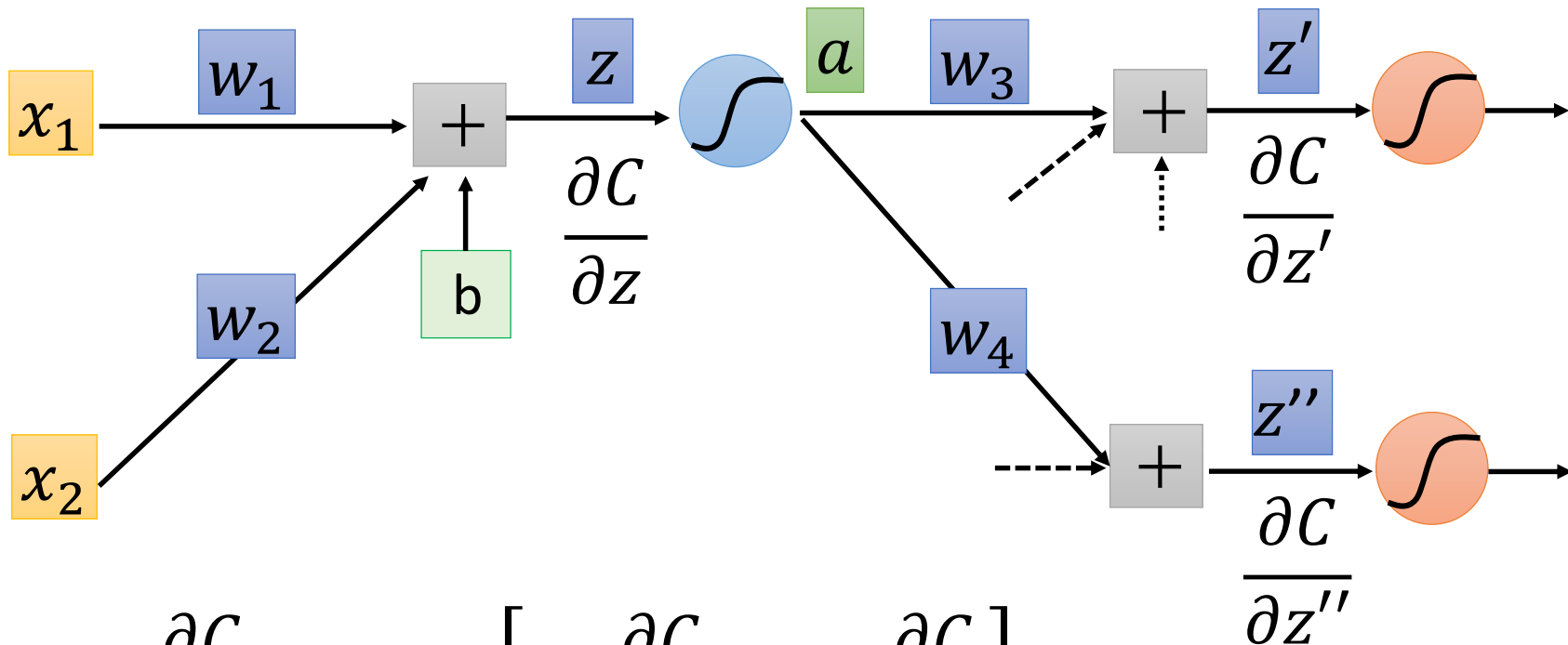


$$\frac{\partial C}{\partial z} = \frac{\partial a}{\partial z} \frac{\partial C}{\partial a}$$
$$\frac{\partial C}{\partial a} = \underbrace{\frac{\partial z'}{\partial a}}_{w_3} \underbrace{\frac{\partial C}{\partial z'}}_{?} + \underbrace{\frac{\partial z''}{\partial a}}_{w_4} \underbrace{\frac{\partial C}{\partial z''}}_{?} \quad (\text{Chain rule})$$

Assumed  
it's known

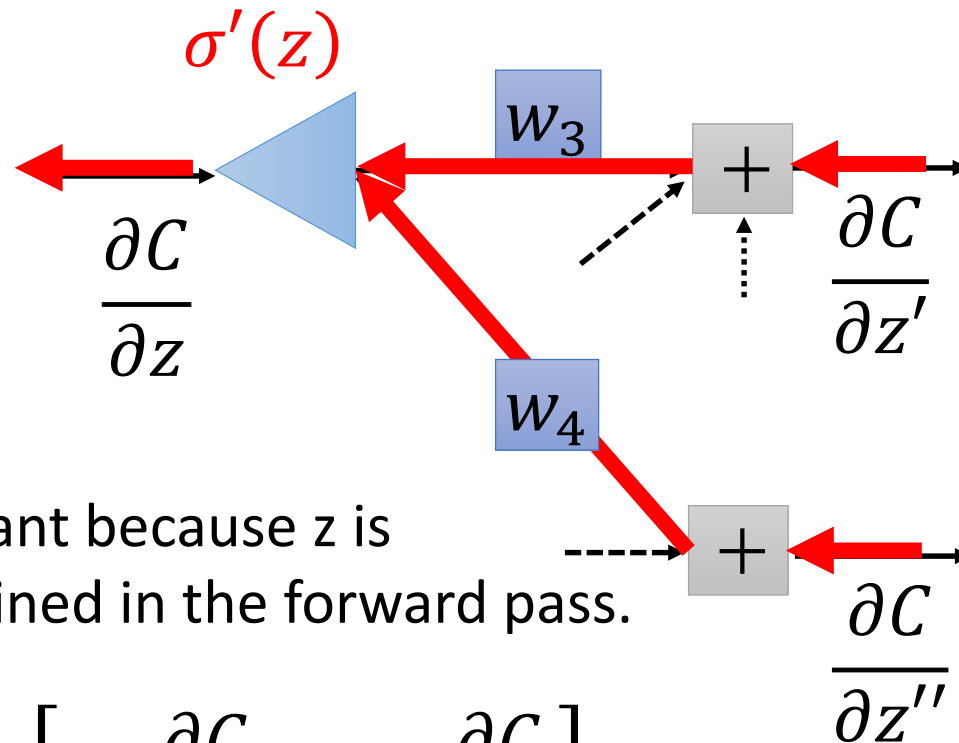
# Backpropagation – Backward pass

Compute  $\partial C / \partial z$  for all activation function inputs  $z$



$$\frac{\partial C}{\partial z} = \sigma'(z) \left[ w_3 \frac{\partial C}{\partial z'} + w_4 \frac{\partial C}{\partial z''} \right]$$

# Backpropagation – Backward pass

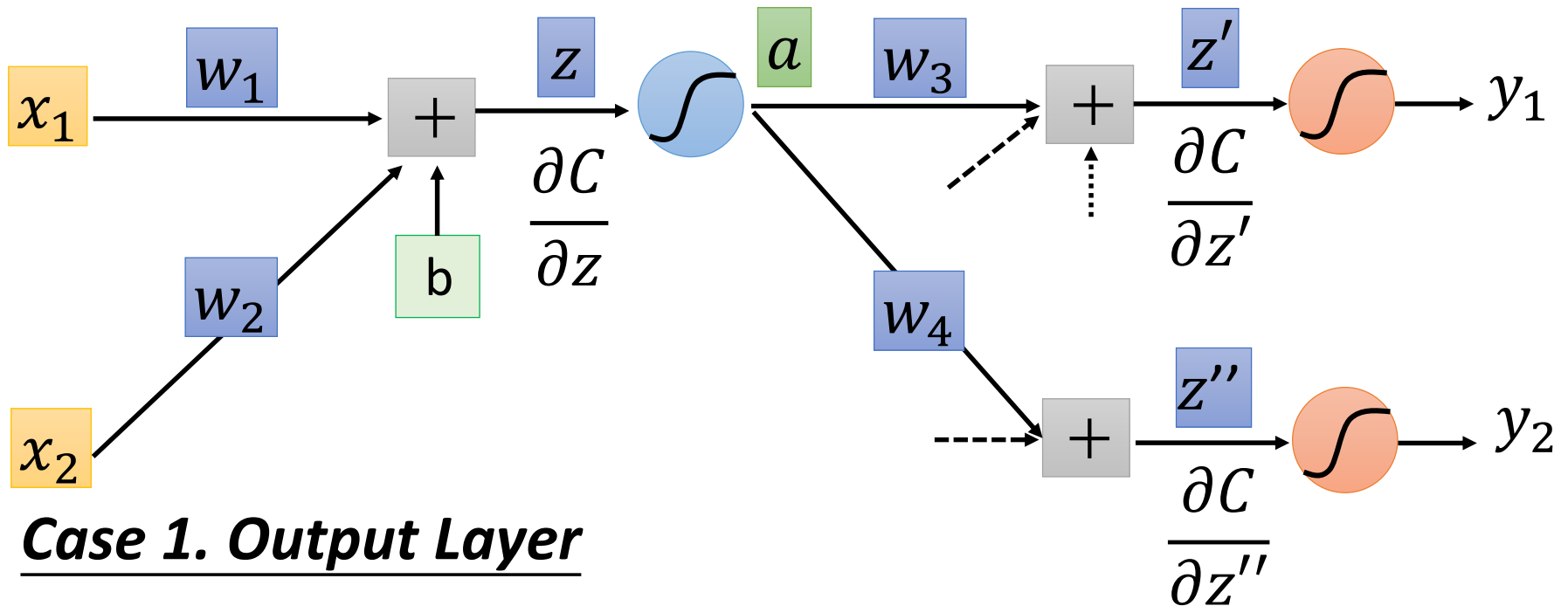


$\sigma'(z)$  is a constant because  $z$  is already determined in the forward pass.

$$\frac{\partial C}{\partial z} = \sigma'(z) \left[ w_3 \frac{\partial C}{\partial z'} + w_4 \frac{\partial C}{\partial z''} \right]$$

# Backpropagation – Backward pass

Compute  $\partial C / \partial z$  for all activation function inputs  $z$



**Case 1. Output Layer**

$$\frac{\partial C}{\partial z'} = \frac{\partial y_1}{\partial z'} \frac{\partial C}{\partial y_1}$$

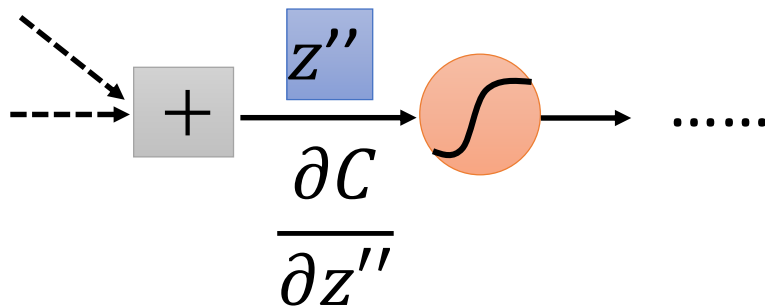
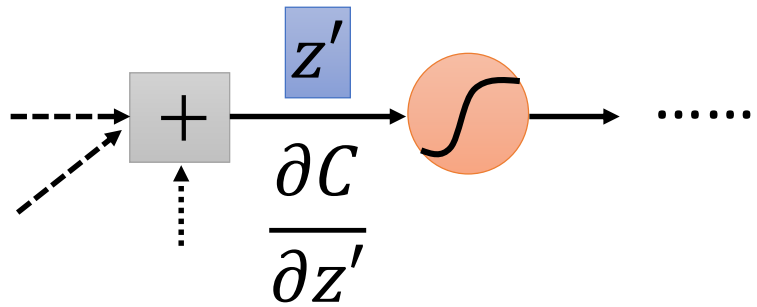
$$\frac{\partial C}{\partial z''} = \frac{\partial y_2}{\partial z''} \frac{\partial C}{\partial y_2}$$

Done!

# Backpropagation – Backward pass

Compute  $\partial C / \partial z$  for all activation function inputs  $z$

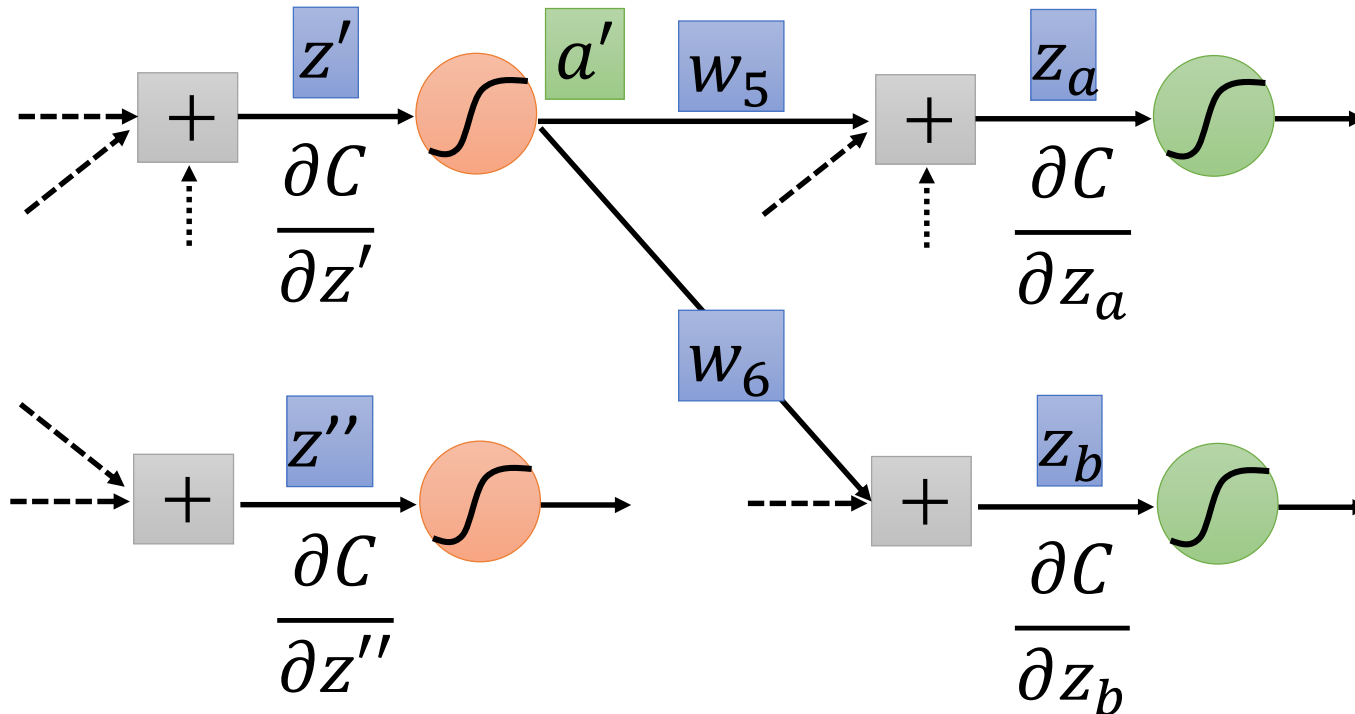
## Case 2. Not Output Layer



# Backpropagation – Backward pass

Compute  $\partial C / \partial z$  for all activation function inputs  $z$

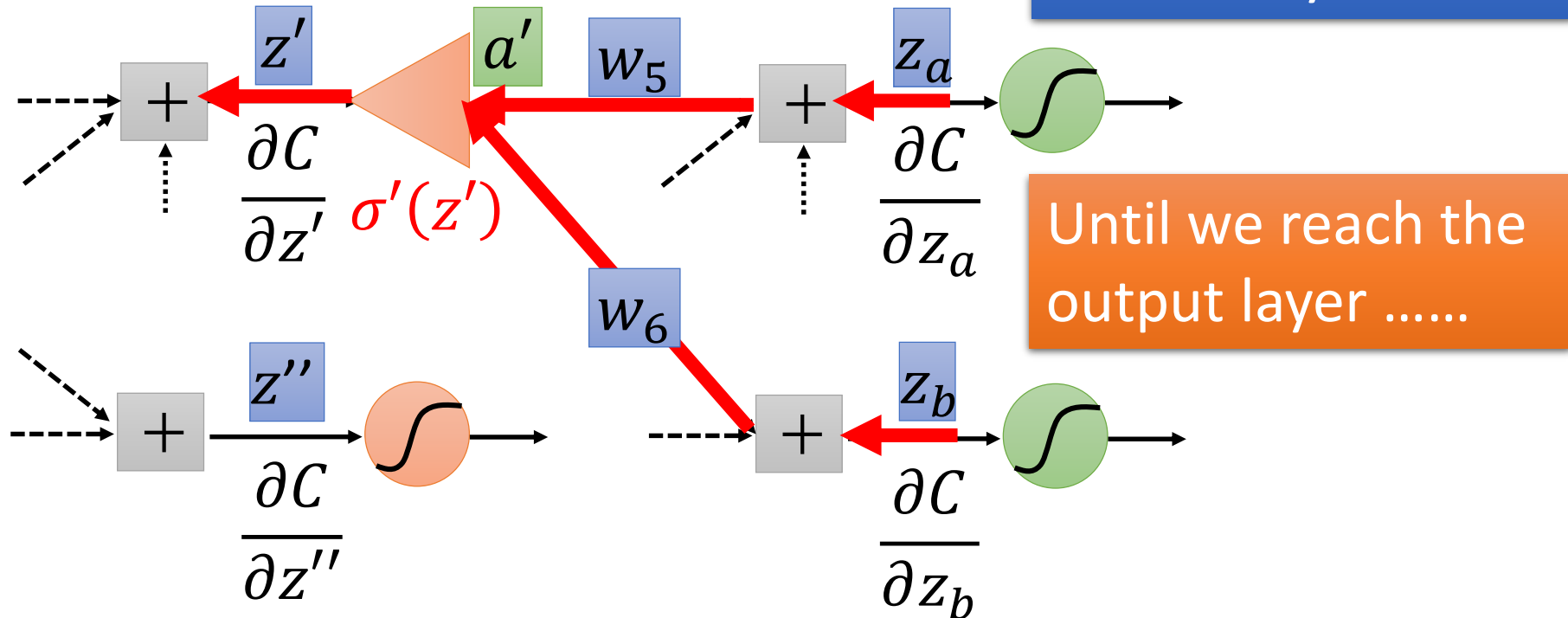
## Case 2. Not Output Layer



# Backpropagation – Backward pass

Compute  $\partial C / \partial z$  for all activation function inputs  $z$

## Case 2. Not Output Layer

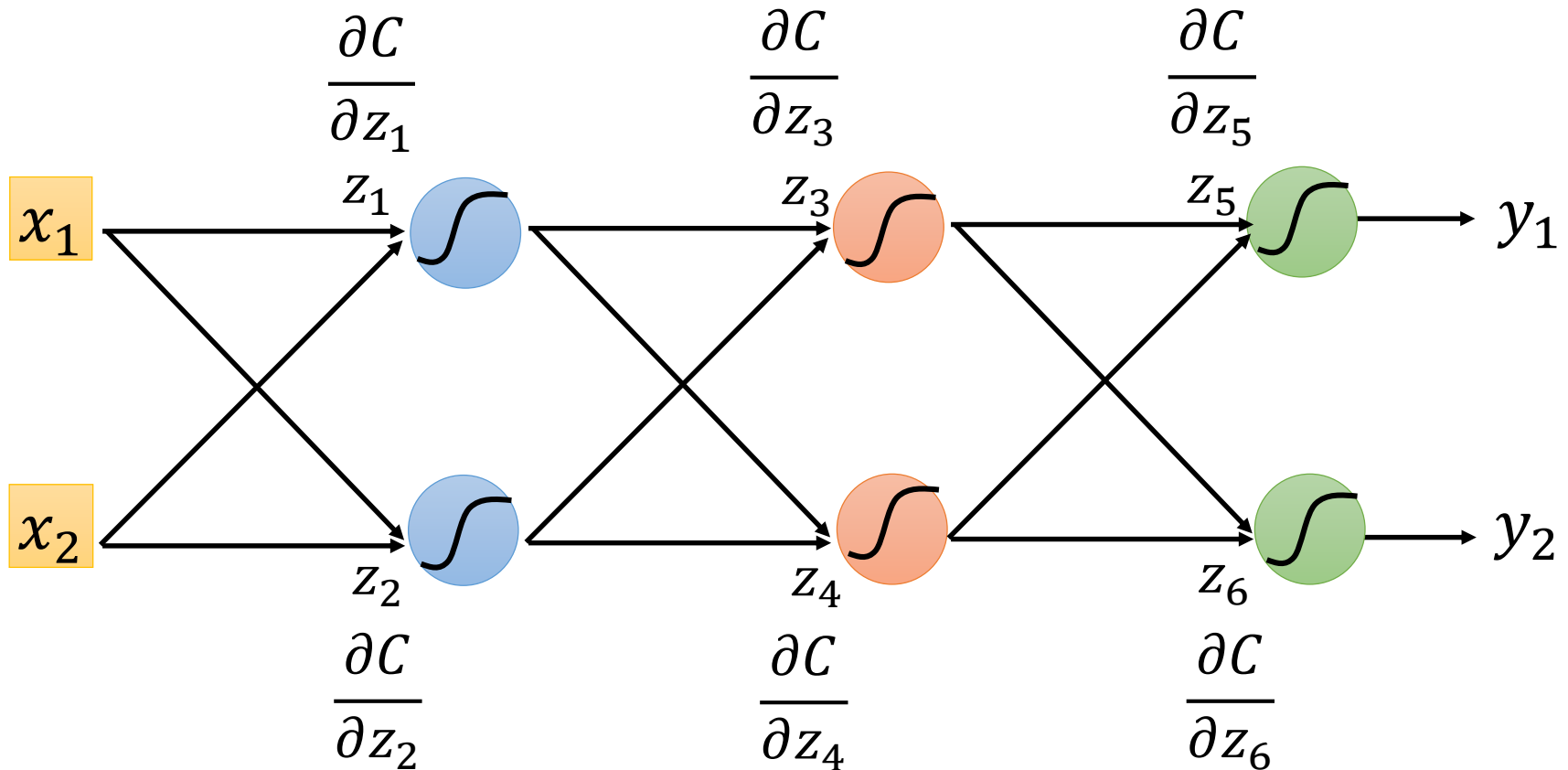




# Backpropagation – Backward Pass

Compute  $\partial C / \partial z$  for all activation function inputs  $z$

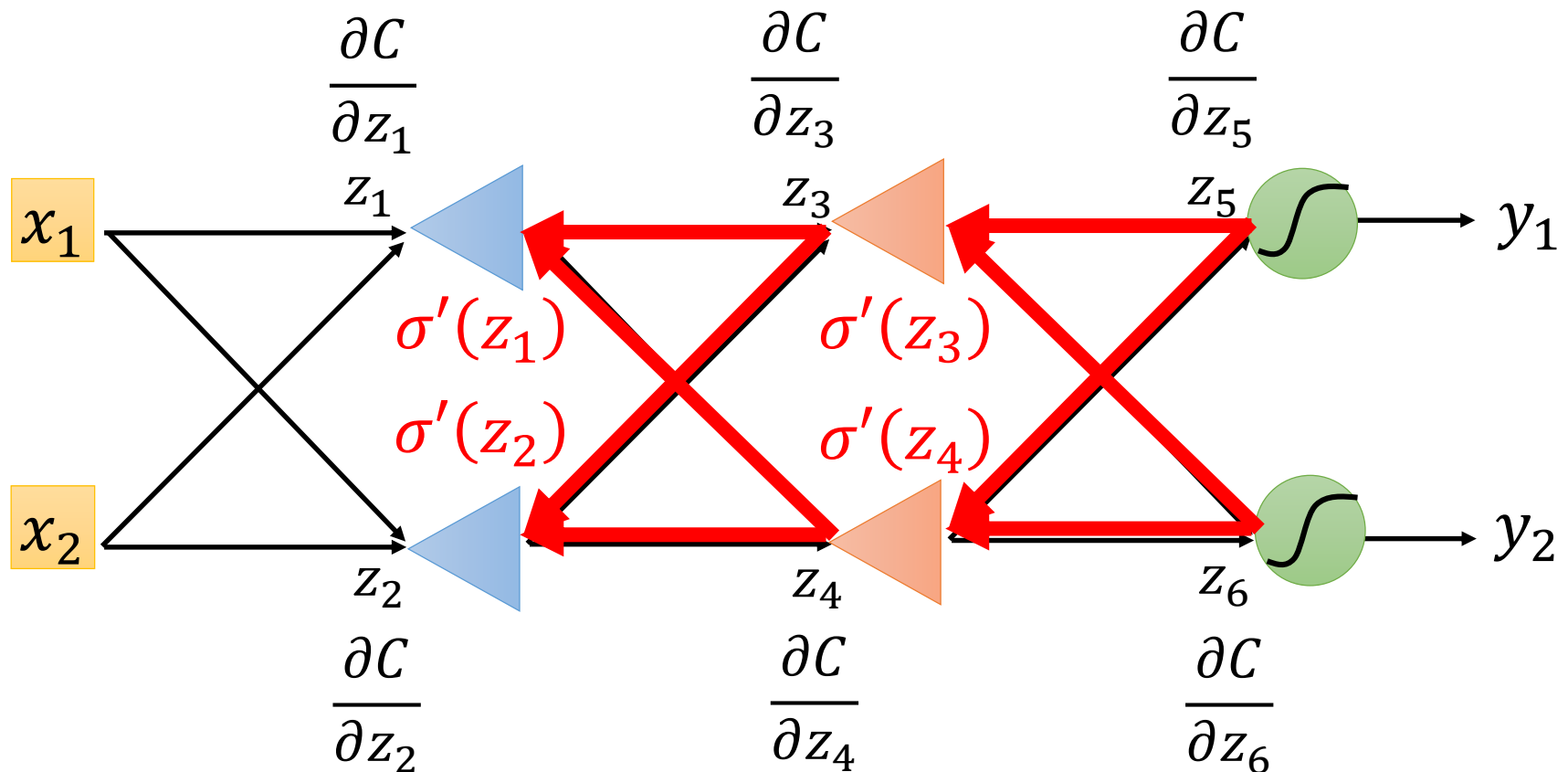
Compute  $\partial C / \partial z$  from the output layer



# Backpropagation – Backward Pass

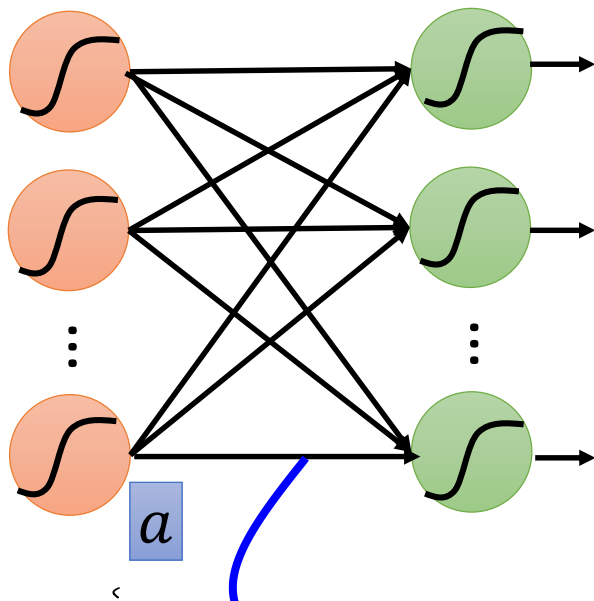
Compute  $\partial C / \partial z$  for all activation function inputs  $z$

Compute  $\partial C / \partial z$  from the output layer



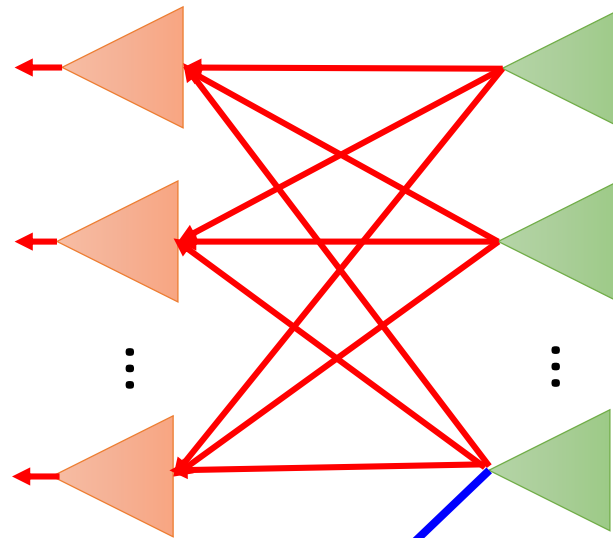
# Backpropagation – Summary

## Forward Pass



$$\frac{\partial z}{\partial w} = a$$

## Backward Pass



$\times$

$$\frac{\partial C}{\partial z}$$

$$= \frac{\partial C}{\partial w}$$

for all  $w$

## 反向传播中

第一级残差  $\delta^L = \frac{\partial C}{\partial a} \odot G'(z^L) = \frac{\partial C}{\partial z^L}$

往后残差每一级乘权重传递  $\delta^l = ((w^{l+1})^T \delta^{l+1}) \odot G'(z^l) = \frac{\partial C}{\partial z^l}$

$$\frac{\partial C}{\partial b_j^l} = \delta_j^l$$

$$\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l$$

$$\begin{matrix} a^{l-1} & z^l \\ 100 \times N & 10 \times N \\ & \begin{matrix} w^l \\ 10 \times 100 \end{matrix} \end{matrix}$$

矩阵形式:

$$a = G(z), z = Wx + b$$

$$\Rightarrow \frac{\partial a}{\partial x} = W^T \frac{\partial a}{\partial z}, \frac{\partial a}{\partial w} = \frac{\partial a}{\partial z} \cdot x^T$$

$$\therefore \frac{\partial C}{\partial w^l} = \delta^l \cdot a^{l-1T} \quad \frac{\partial C}{\partial a^{l-1}} = w^{lT} \cdot \delta^l \quad \frac{\partial C}{\partial b^{l-1}} = \text{sum}(\delta^l, \text{axis}=1)$$