# CSCE 636 Neural Networks (Deep Learning)

Lecture 3: Gradient Descent and Backpropagation Algorithm

Anxiao (Andrew) Jiang

Based on interesting lecture by Prof. Hung-yi Lee, https://www.youtube.com/watch?v=ibJpTrp5mcE

# Backpropagation

# Gradient Descent

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

Starting Parameters $\theta^0$

$\nabla L(\theta)$

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

$Compute\ \nabla L(\theta^0)$

# Gradient Descent

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

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

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

$Compute\ \nabla L(\theta^0)$

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

learning rate
(such as 0.001)

# Gradient Descent

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

Starting Parameters $\quad \theta^0 \longrightarrow \theta^1 \longrightarrow \theta^2 \longrightarrow \cdots\cdots$

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

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

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

Millions of parameters ......

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

# Chain Rule

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

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

**Case 2**

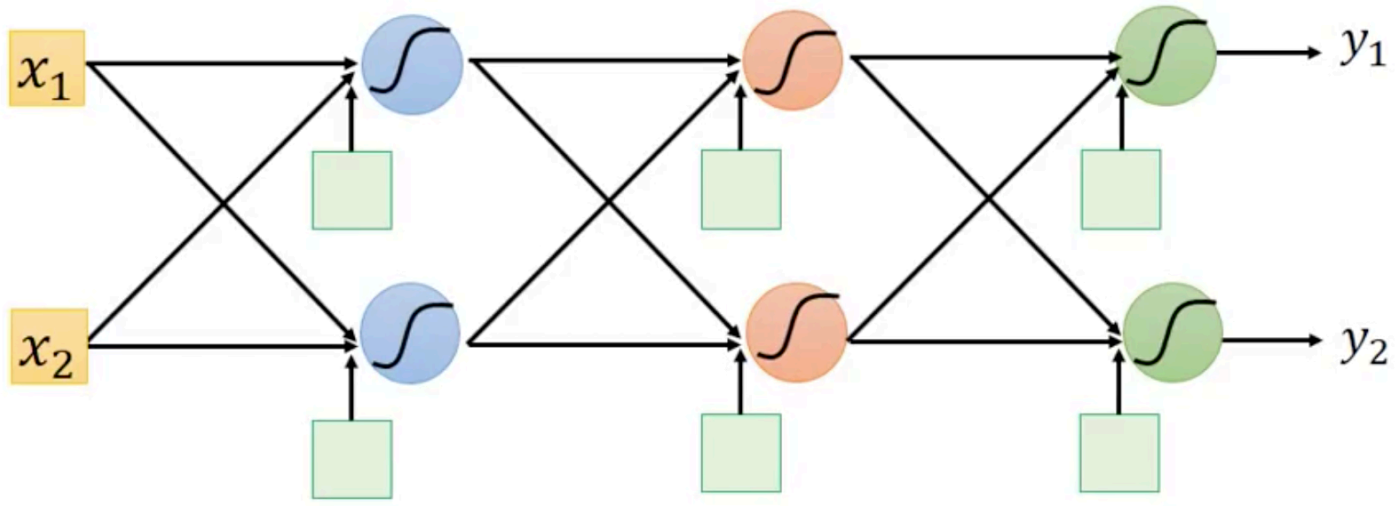$$x = g(s) \qquad y = h(s) \qquad z = k(x, y)$$



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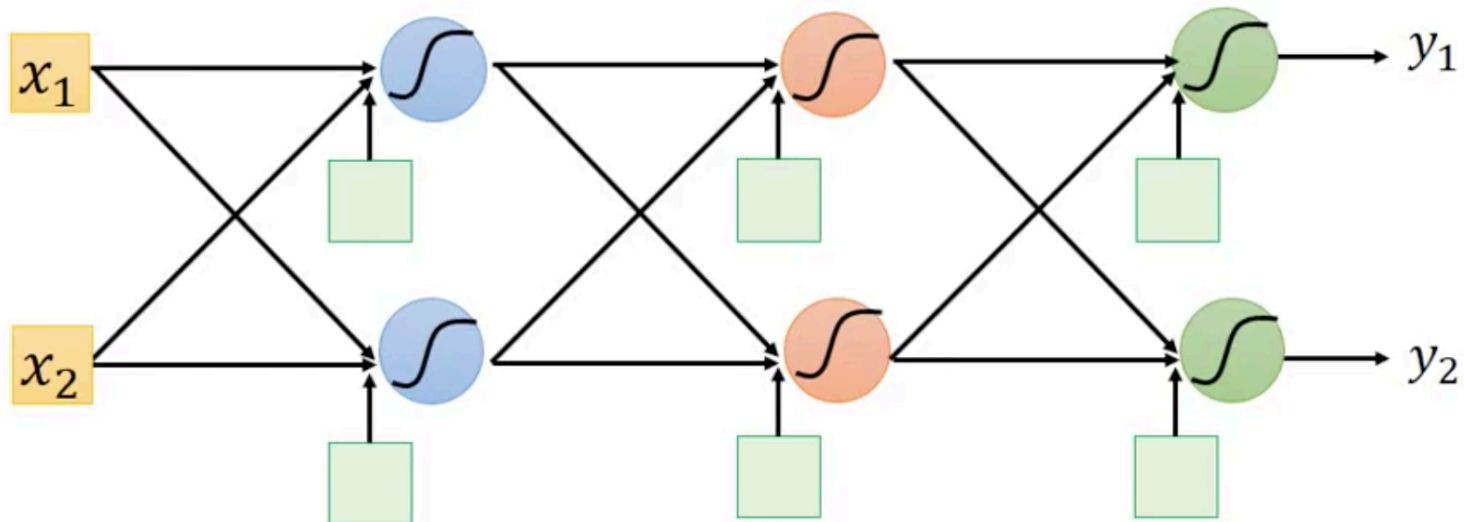
# Backpropagation



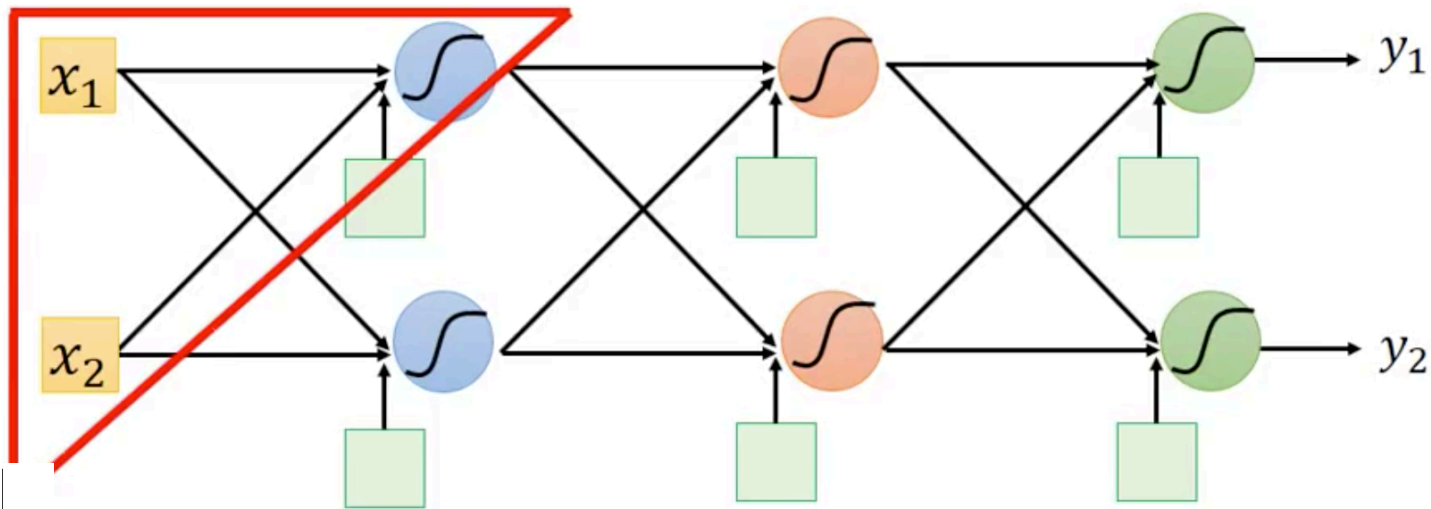$$L(\theta) = \sum_{n=1}^{N} C^n(\theta)$$

# Backpropagation

$$x^n \rightarrow \boxed{\begin{array}{c} \text{NN} \\ \theta \end{array}} \rightarrow y^n \underset{C^n}{\longleftrightarrow} \hat{y}^n$$

$$L(\theta) = \sum_{n=1}^{N} C^n(\theta) \implies \frac{\partial L(\theta)}{\partial w} = \sum_{n=1}^{N} \frac{\partial C^n(\theta)}{\partial w}$$
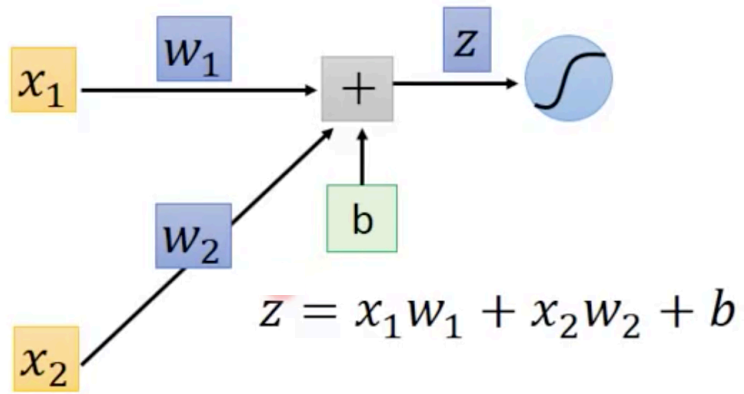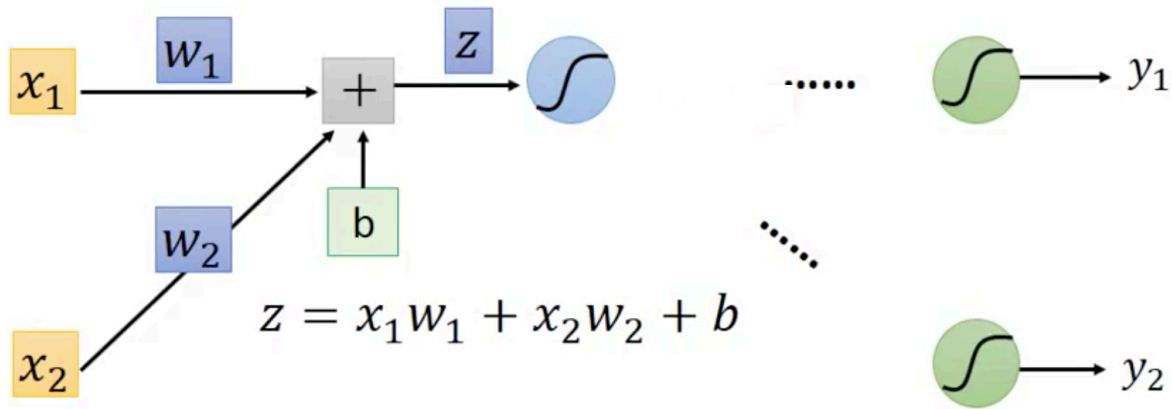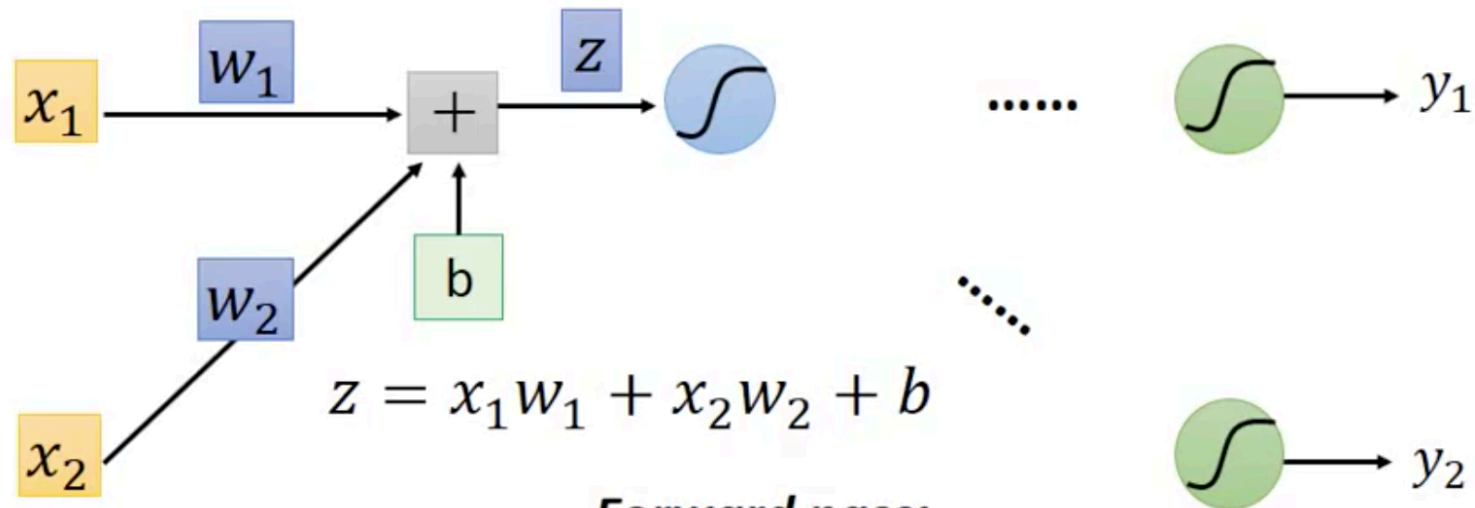
# Backpropagation



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

# Backpropagation

# Backpropagation



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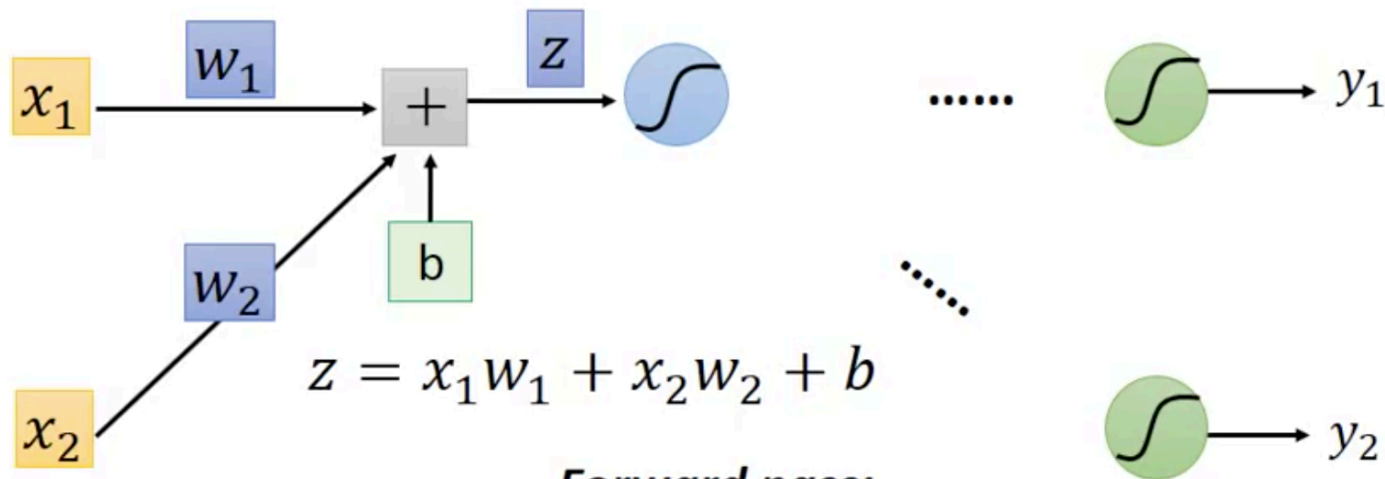
# Backpropagation



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

# Backpropagation



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

**Forward pass:**

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

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

(Chain rule)

# Backpropagation



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

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

(Chain rule)

**Forward pass:**

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

**Backward pass:**

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

# Backpropagation – Forward pass

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



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

$\partial z/\partial w_1 =?$

# Backpropagation – Forward pass

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



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

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

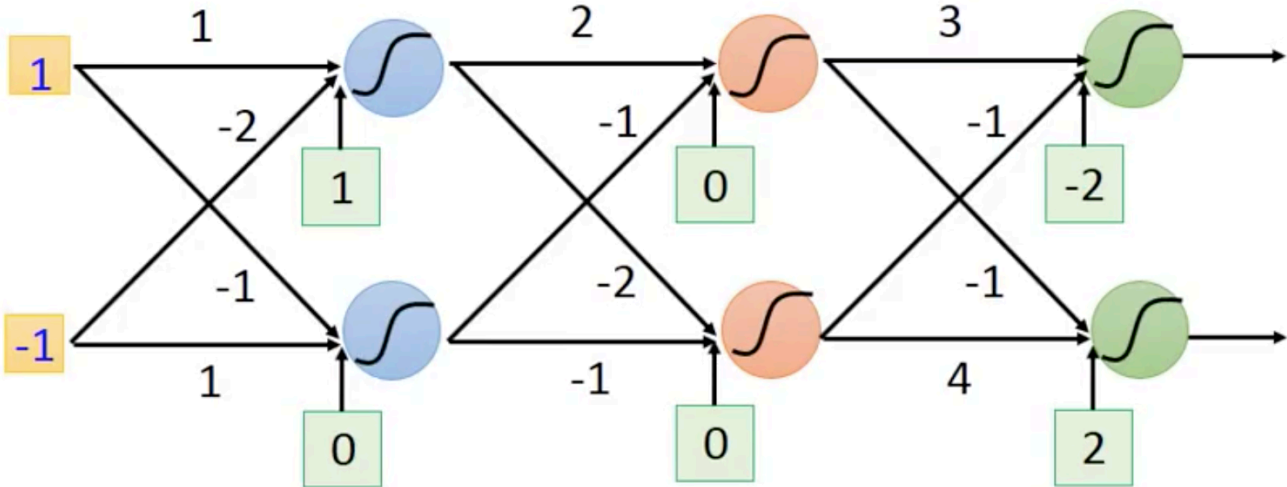$\partial z/\partial w_2 = ? \quad x_2$

# Backpropagation – Forward pass

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



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

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

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

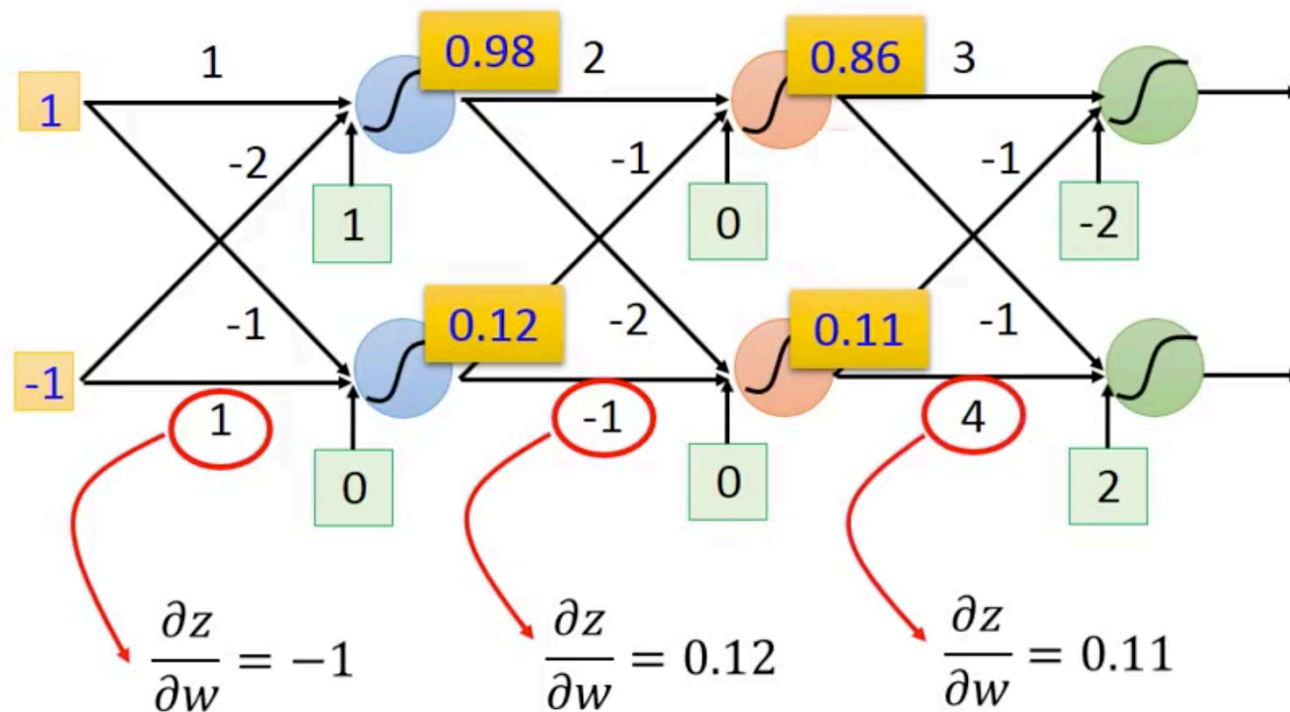The value of the input connected by the weight

# Backpropagation – Forward pass

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

# Backpropagation – Forward pass

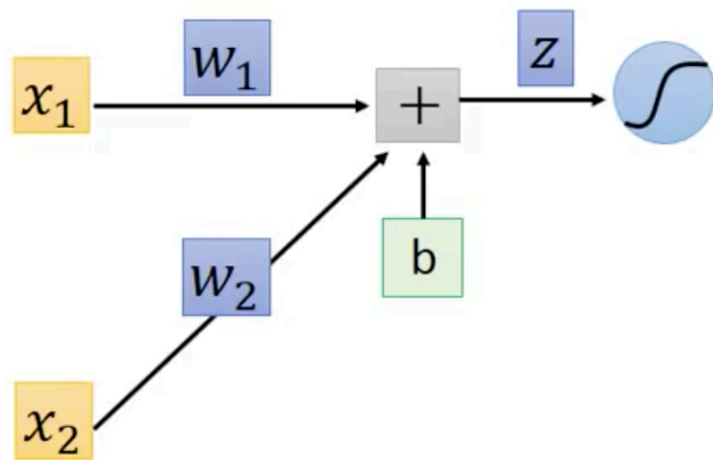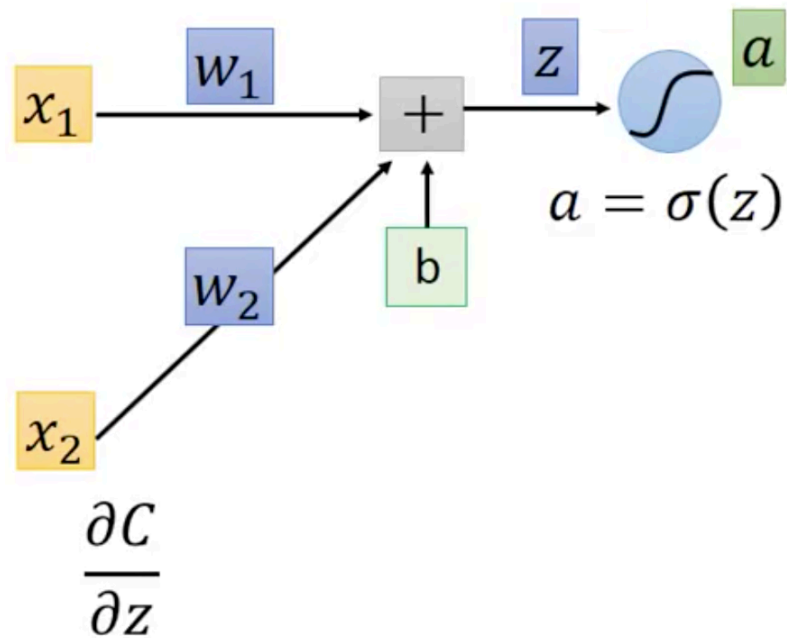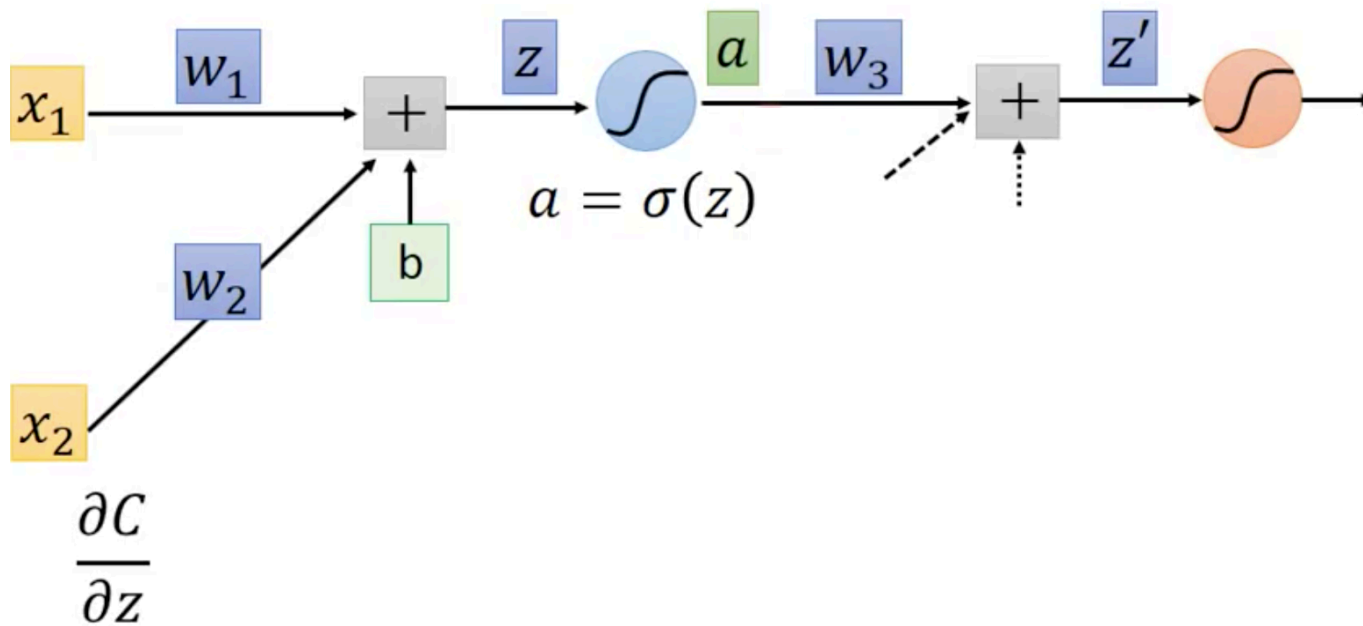Compute $\partial z / \partial w$ for all parameters



$$\frac{\partial z}{\partial w} = -1$$

# Backpropagation – Forward pass

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



$$\frac{\partial z}{\partial w} = -1 \qquad \frac{\partial z}{\partial w} = 0.12$$

# Backpropagation – Forward pass

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



$$\frac{\partial z}{\partial w} = -1 \qquad \frac{\partial z}{\partial w} = 0.12 \qquad \frac{\partial z}{\partial w} = 0.11$$
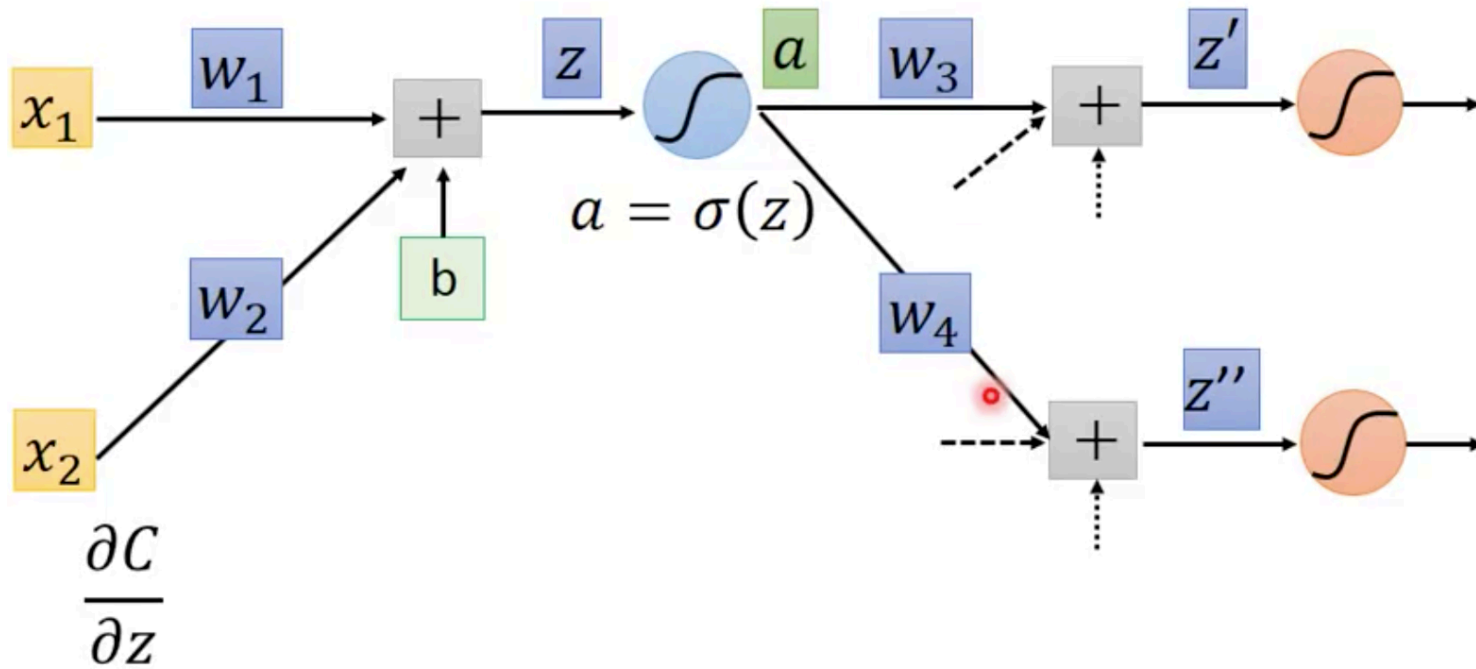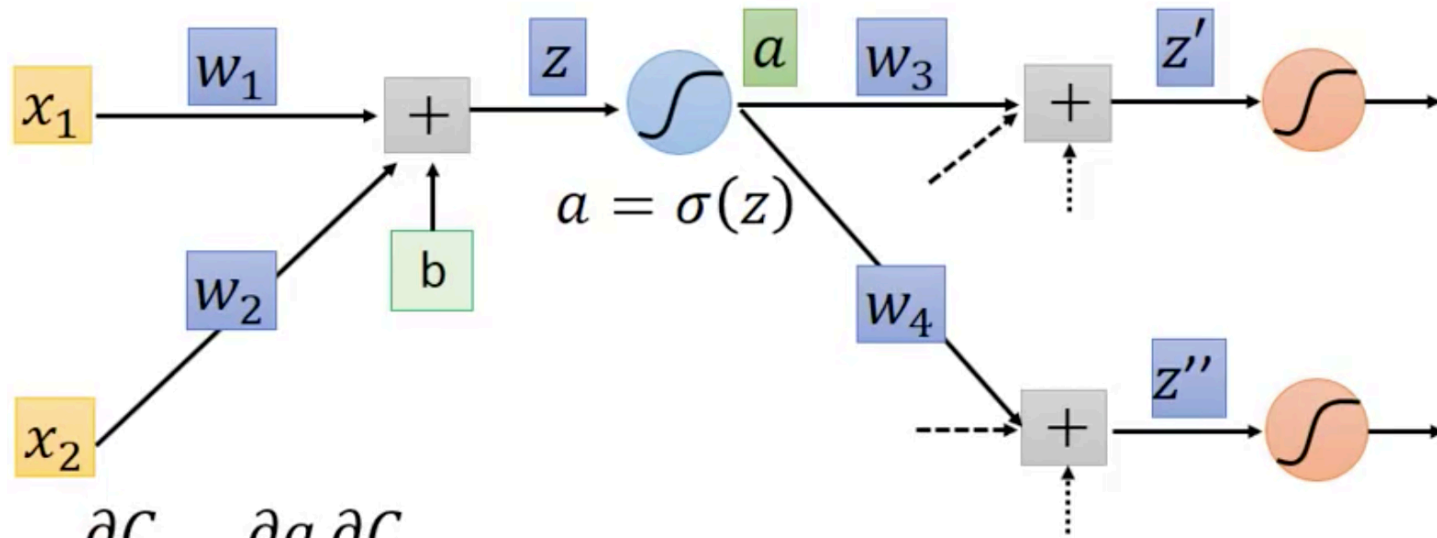
# Backpropagation – Backward pass

Compute $\partial C / \partial z$ for all activation function inputs z
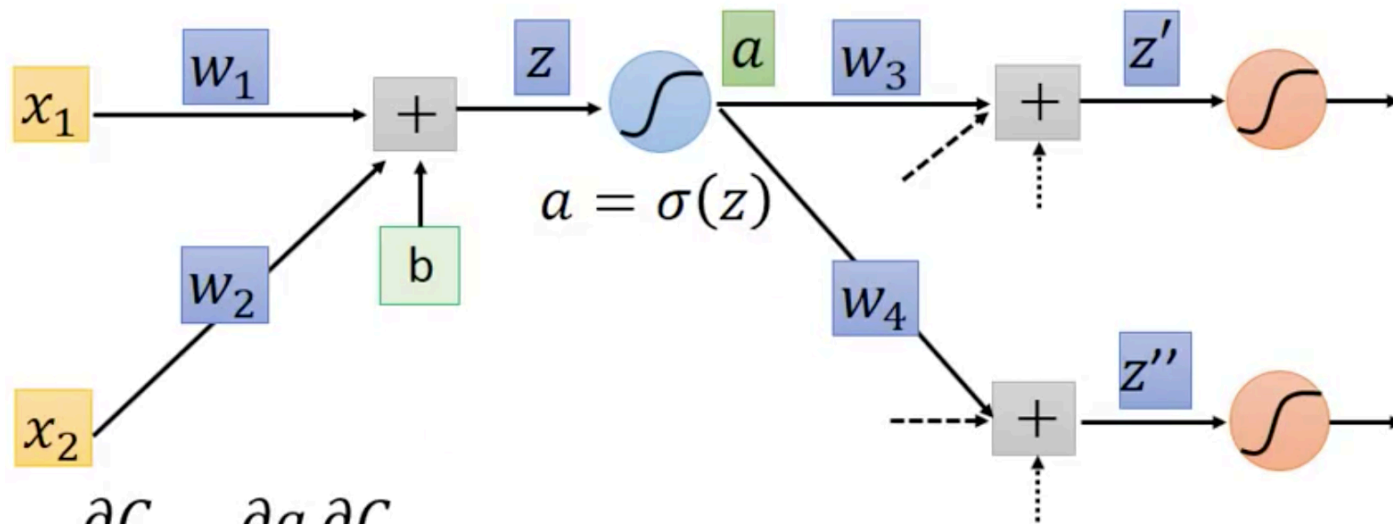
# Backpropagation – Backward pass

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



$$a = \sigma(z)$$

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

# Backpropagation – Backward pass

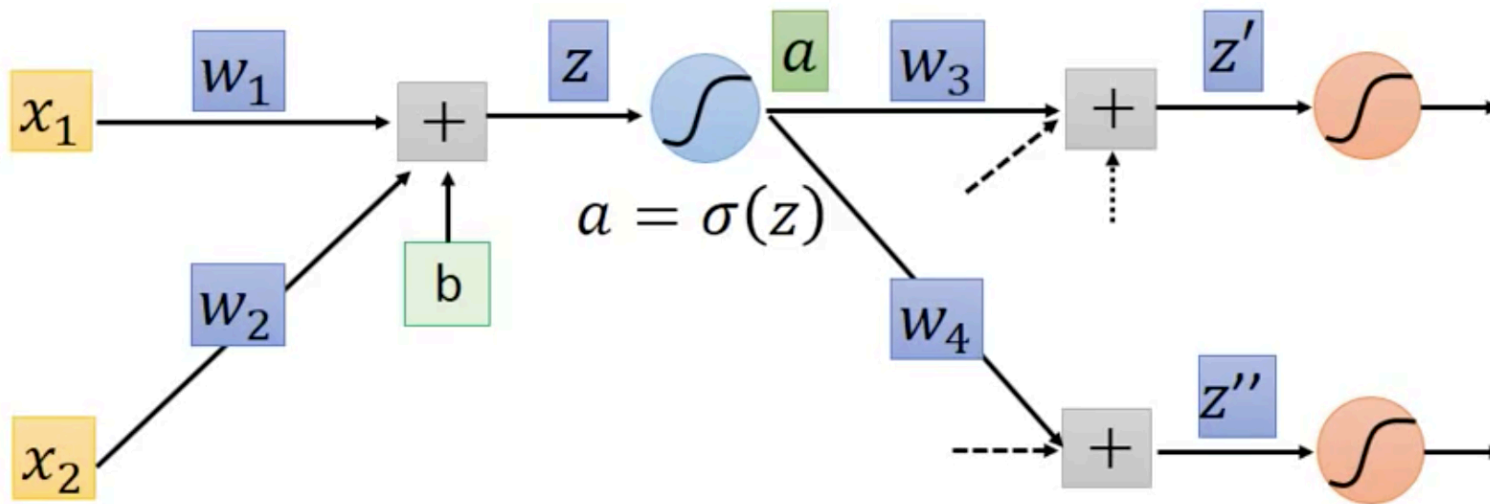Compute $\partial C / \partial z$ for all activation function inputs z



$$a = \sigma(z)$$

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

# Backpropagation – Backward pass

Compute $\partial C / \partial z$ for all activation function inputs z
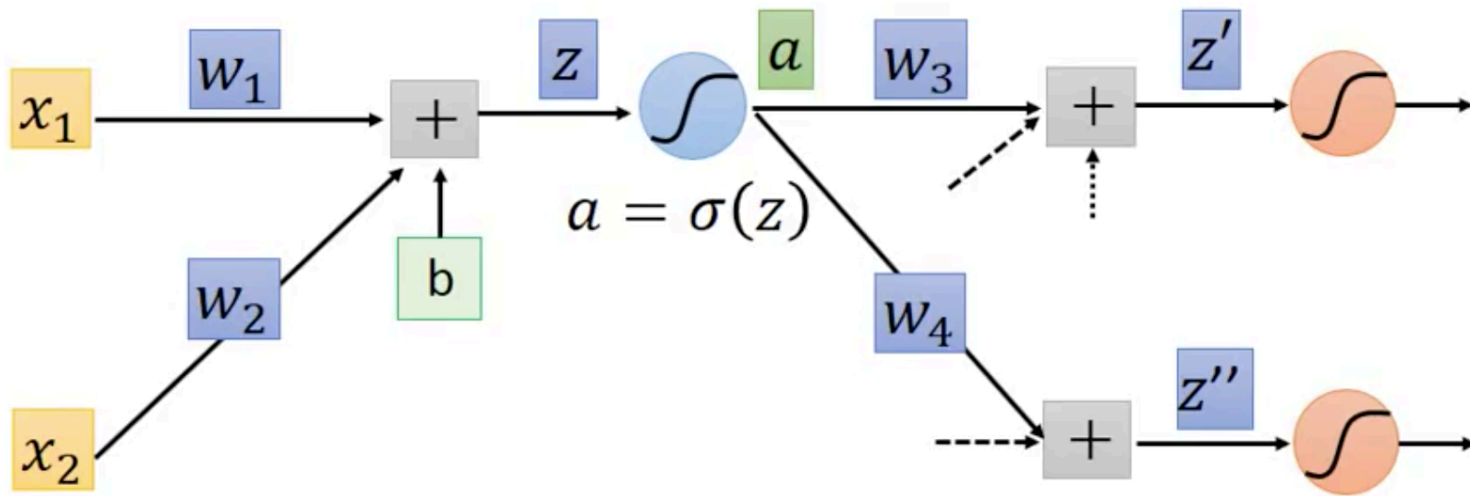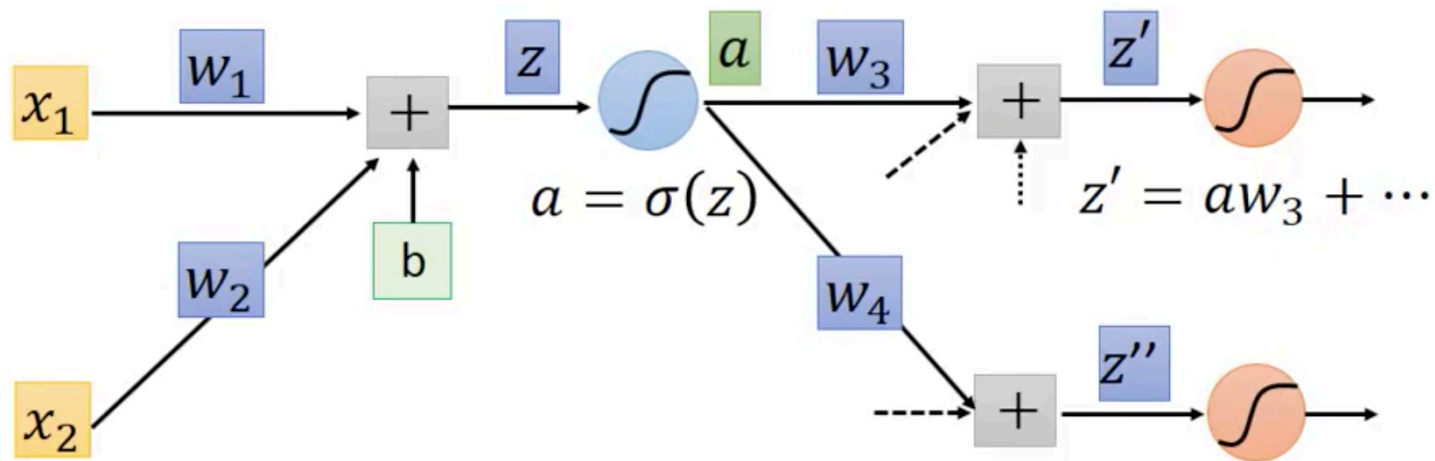
# Backpropagation – Backward pass

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



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

# Backpropagation – Backward pass

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



$$\frac{\partial C}{\partial z} = \frac{\partial a}{\partial z}\frac{\partial C}{\partial a} \qquad \frac{\partial C}{\partial a} = \frac{\partial z'}{\partial a}\frac{\partial C}{\partial z'} + \frac{\partial z''}{\partial a}\frac{\partial C}{\partial z''} \text{ (Chain rule)}$$

# Backpropagation – Backward pass

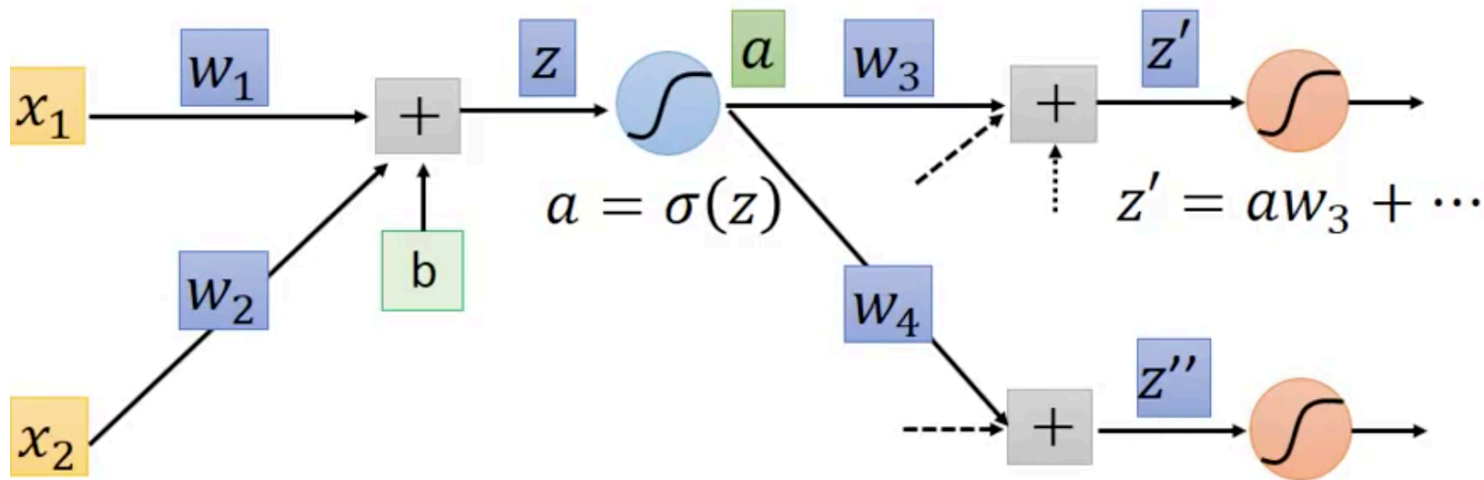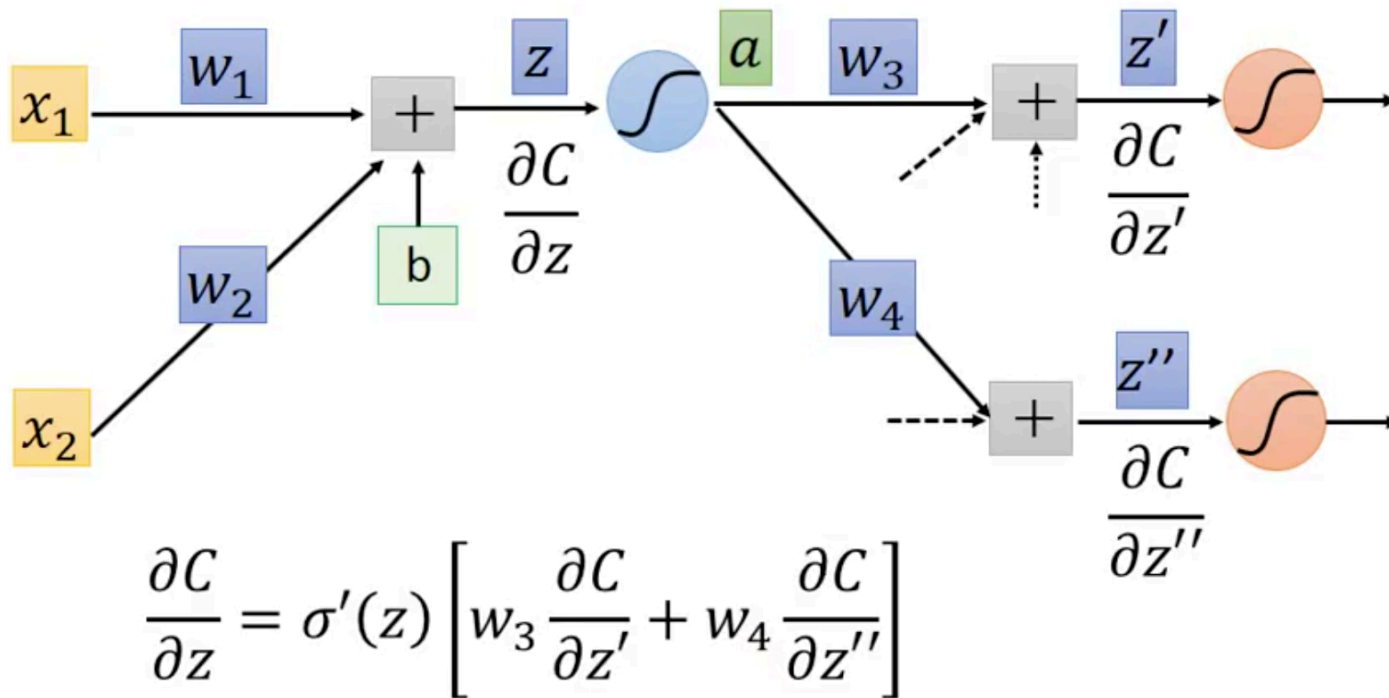Compute $\partial C / \partial z$ for all activation function inputs z



$$\frac{\partial C}{\partial z} = \frac{\partial a}{\partial z}\frac{\partial C}{\partial a} \qquad \frac{\partial C}{\partial a} = \frac{\partial z'}{\partial a}\frac{\partial C}{\partial z'} + \frac{\partial z''}{\partial a}\frac{\partial C}{\partial z''} \quad \text{(Chain rule)}$$
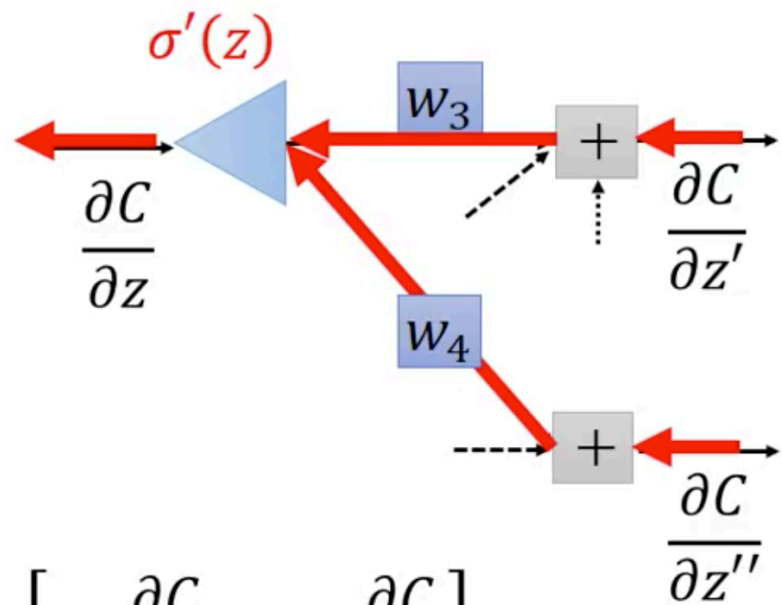
# Backpropagation – Backward pass

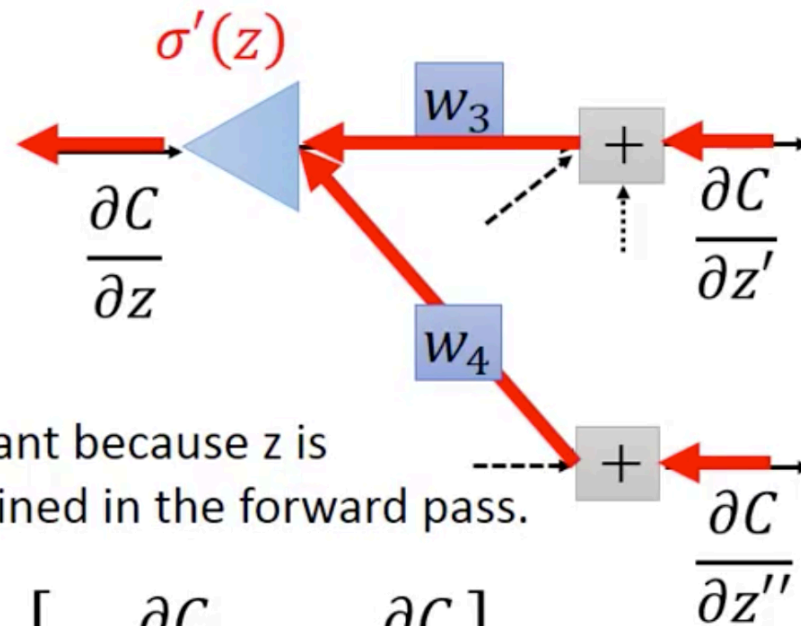Compute $\partial C / \partial z$ for all activation function inputs z



$a = \sigma(z)$

$z' = aw_3 + \cdots$

$$\frac{\partial C}{\partial z} = \frac{\partial a}{\partial z}\frac{\partial C}{\partial a} \qquad \frac{\partial C}{\partial a} = \frac{\partial z'}{\partial a}\frac{\partial C}{\partial z'} + \frac{\partial z''}{\partial a}\frac{\partial C}{\partial z''} \quad \text{(Chain rule)}$$

$w_3 \qquad w_4$

# Backpropagation – Backward pass

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



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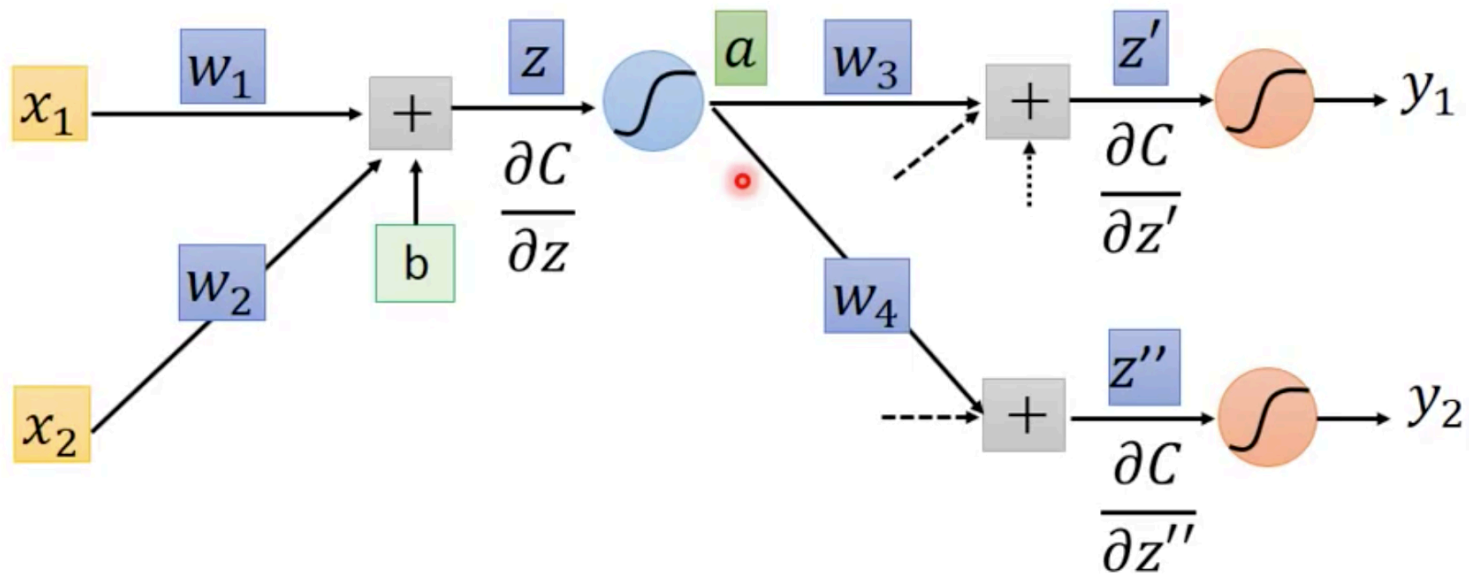
# Backpropagation – Backward pass



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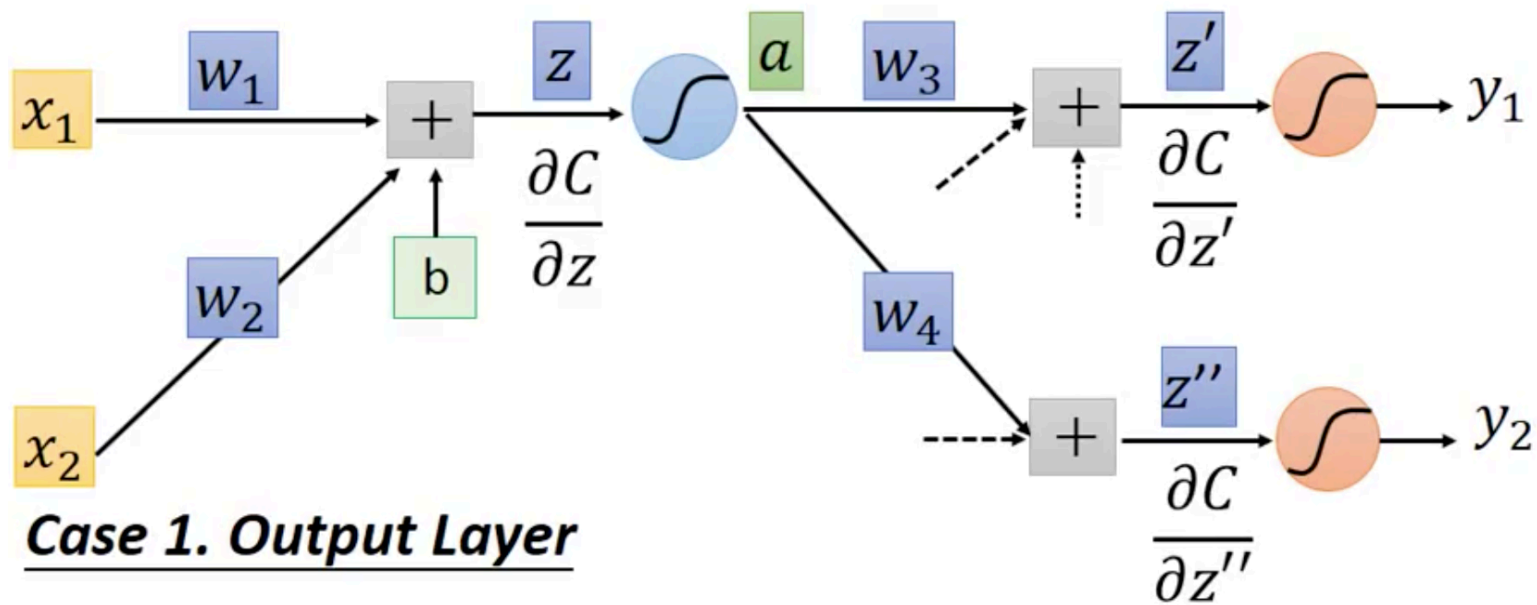
# Backpropagation – Backward pass



$\sigma'(z)$

$w_3$

$+$

$\dfrac{\partial C}{\partial z}$

$\dfrac{\partial C}{\partial z'}$

$w_4$

$+$

$\dfrac{\partial C}{\partial z''}$

$\sigma'(z)$ 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

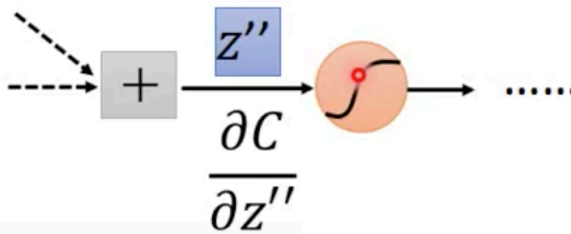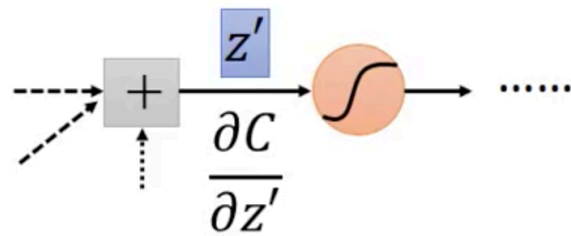# Backpropagation – Backward pass

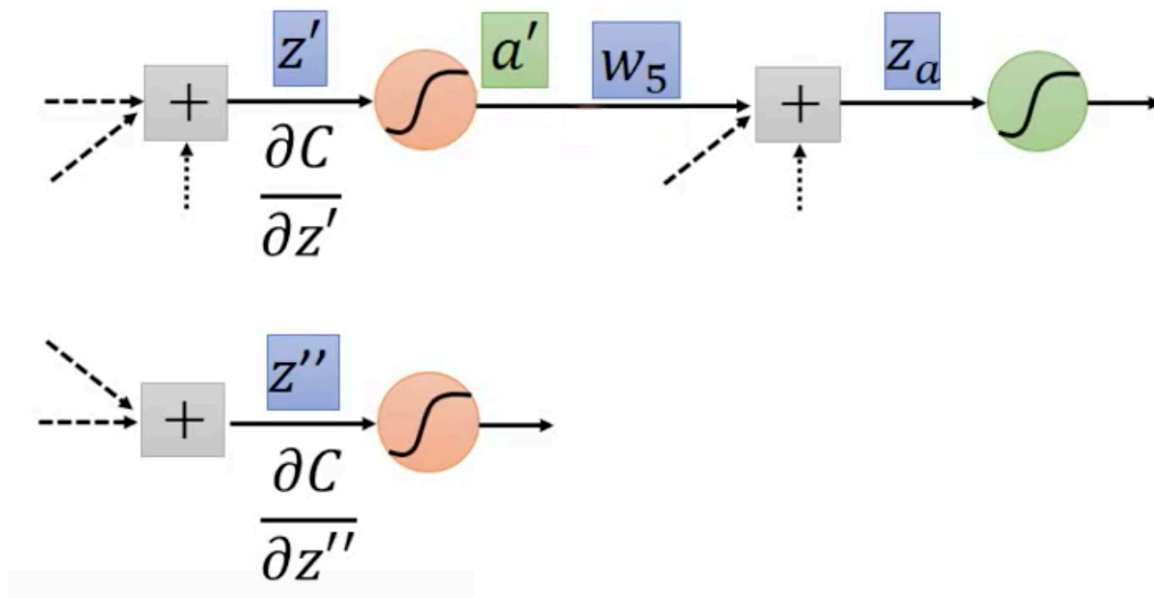Compute $\partial C / \partial z$ for all activation function inputs z



**Case 1. Output Layer**

# Backpropagation – Backward pass

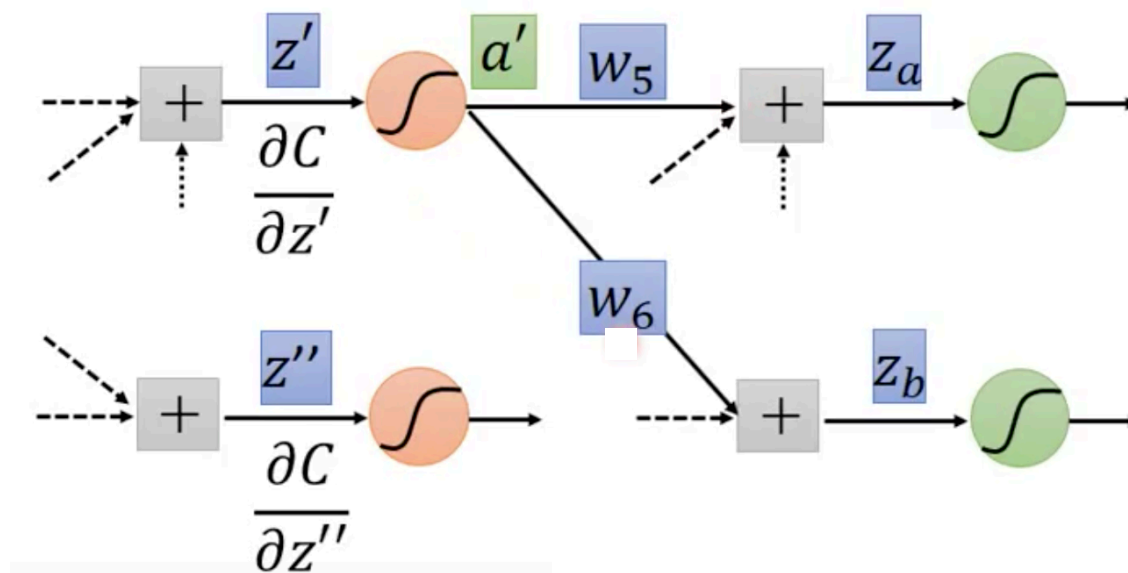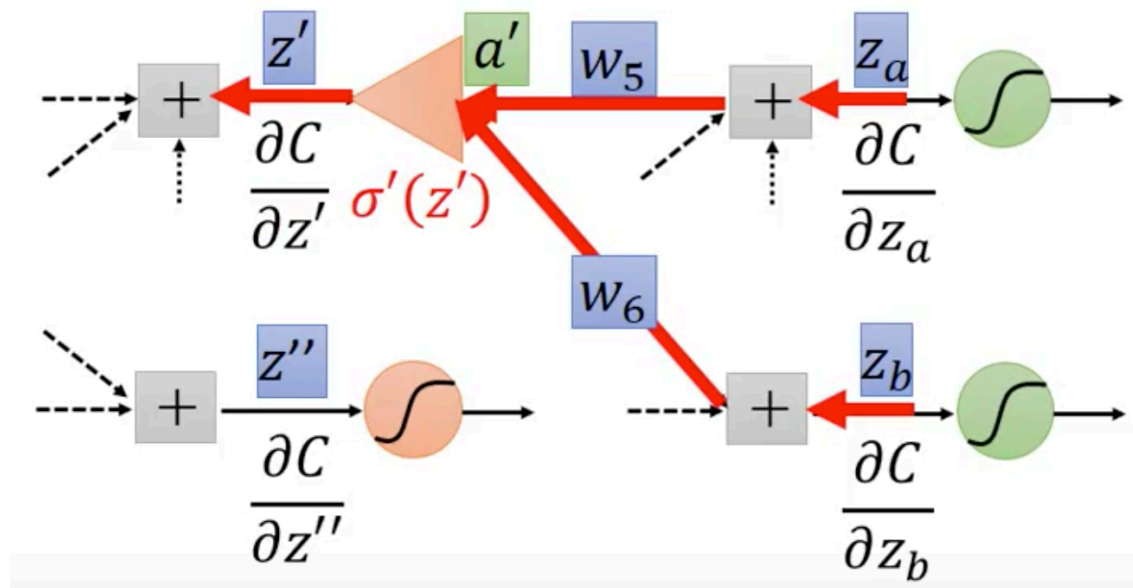Compute $\partial C/\partial z$ for all activation function inputs z



**Case 1. Output Layer**

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

# Backpropagation – Backward pass

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



**Case 1. Output Layer**

$$\frac{\partial C}{\partial z'} = \frac{\partial y_1}{\partial z'} \frac{\partial C}{\partial y_1} \qquad \frac{\partial C}{\partial z''} = \frac{\partial y_2}{\partial z''} \frac{\partial C}{\partial y_2}$$

# Backpropagation – Backward pass

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

## Case 2. Not Output Layer
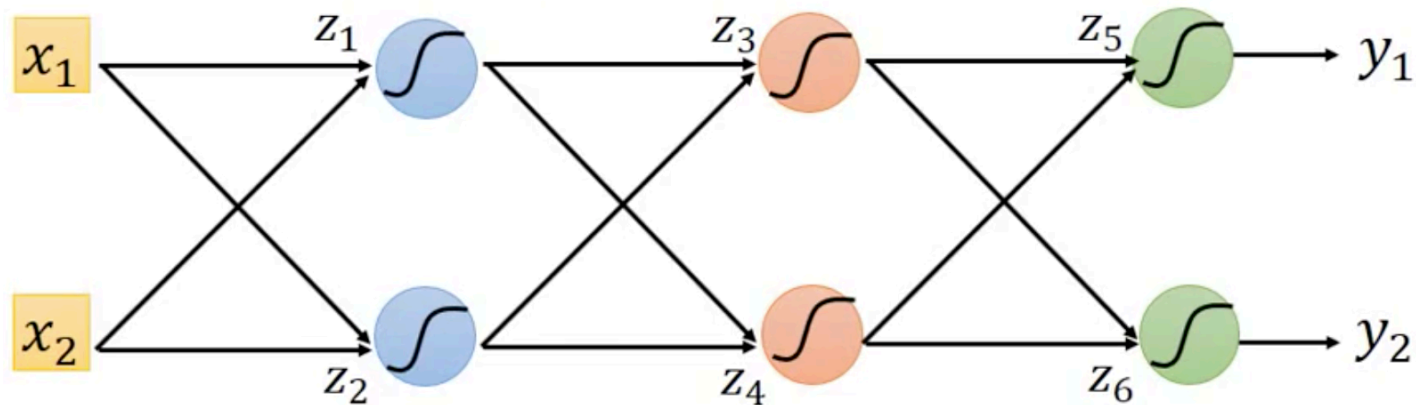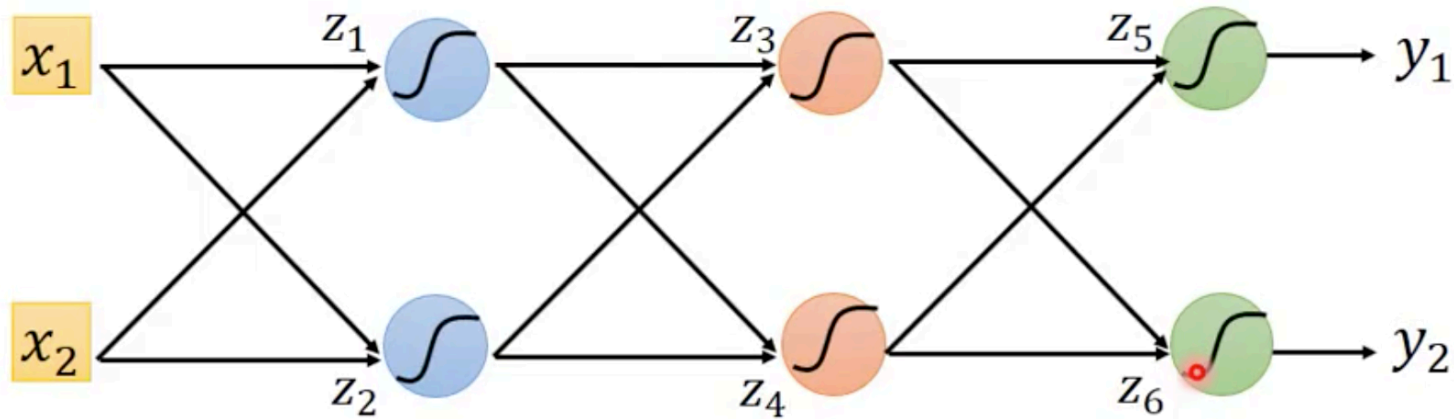
# Backpropagation – Backward pass

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

## Case 2. Not Output Layer
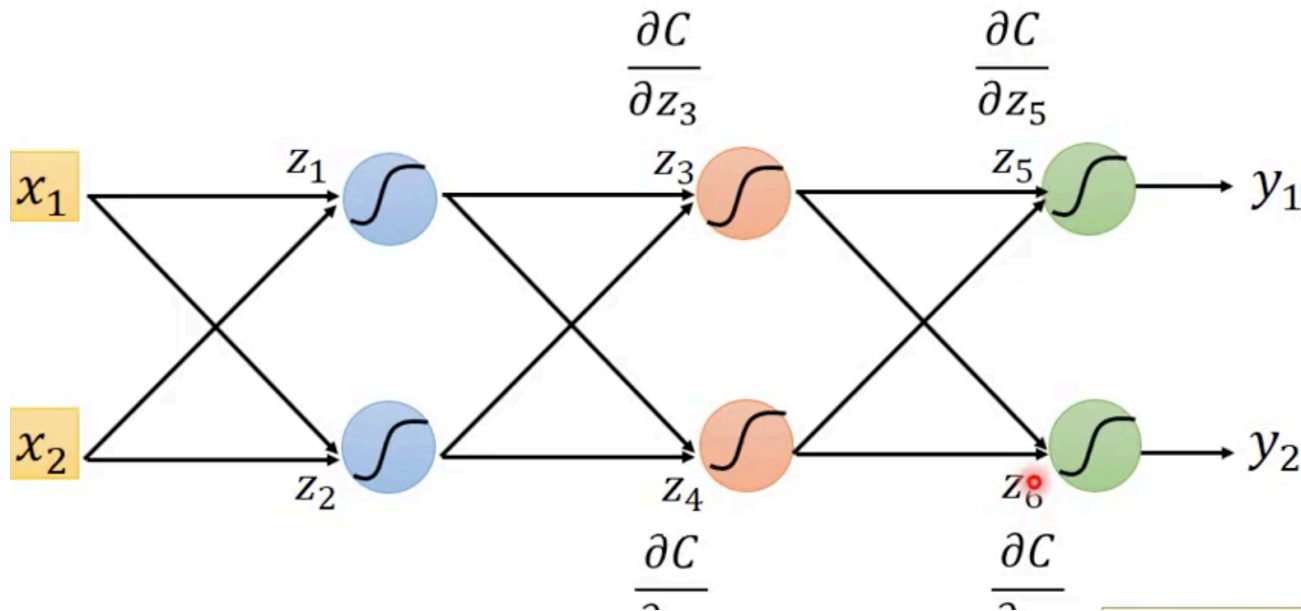
# Backpropagation – Backward Pass

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

# 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

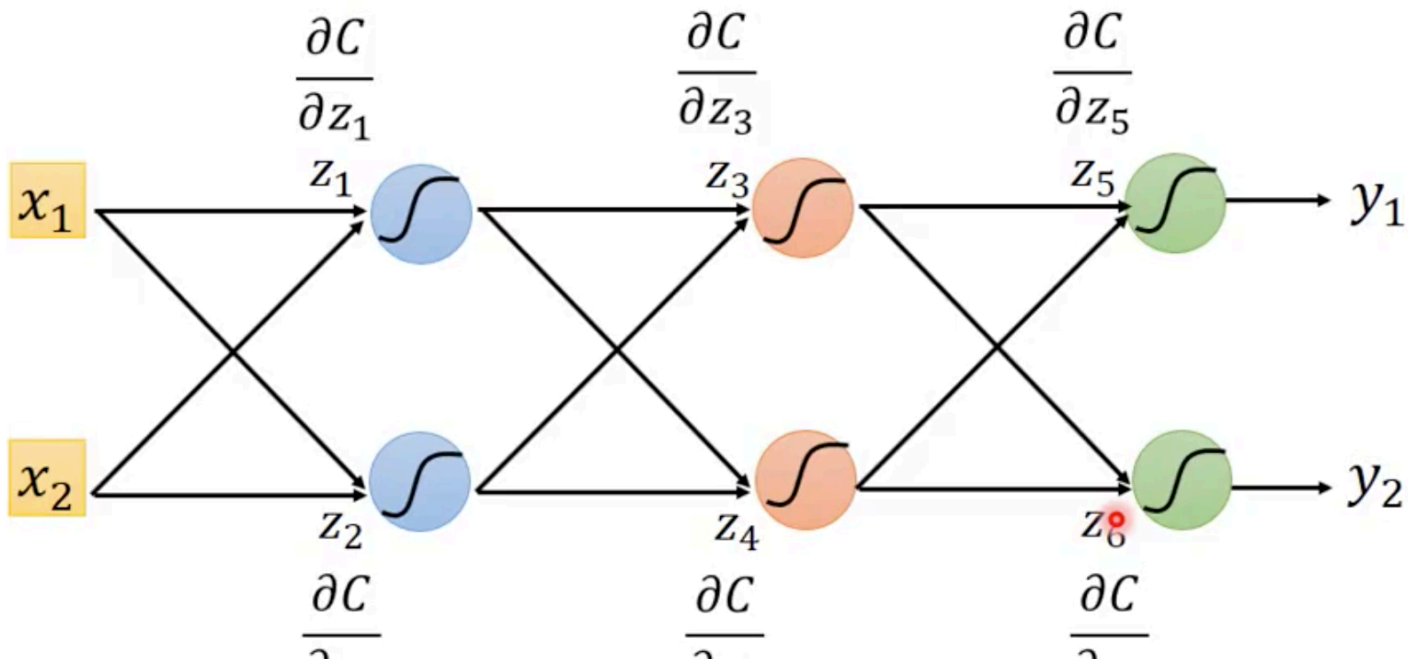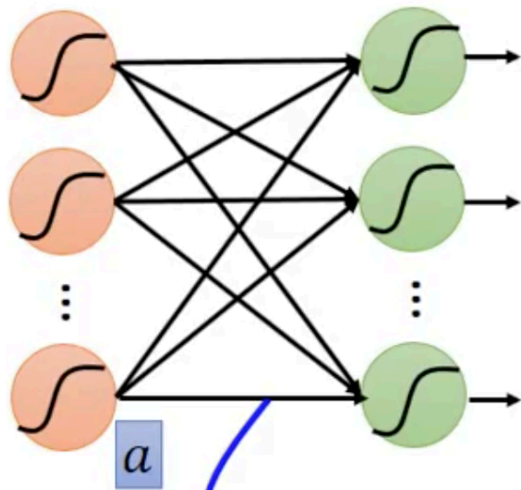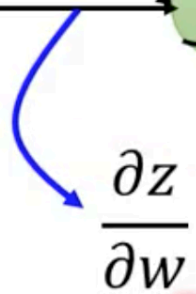Compute $\partial C / \partial z$ for all activation function inputs z

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

# Backpropagation – 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

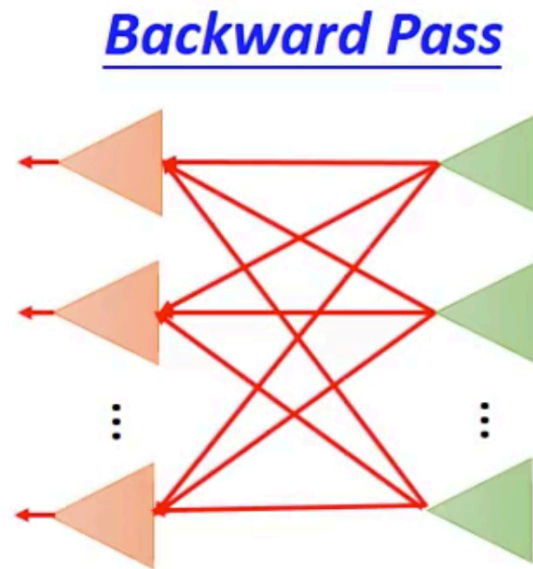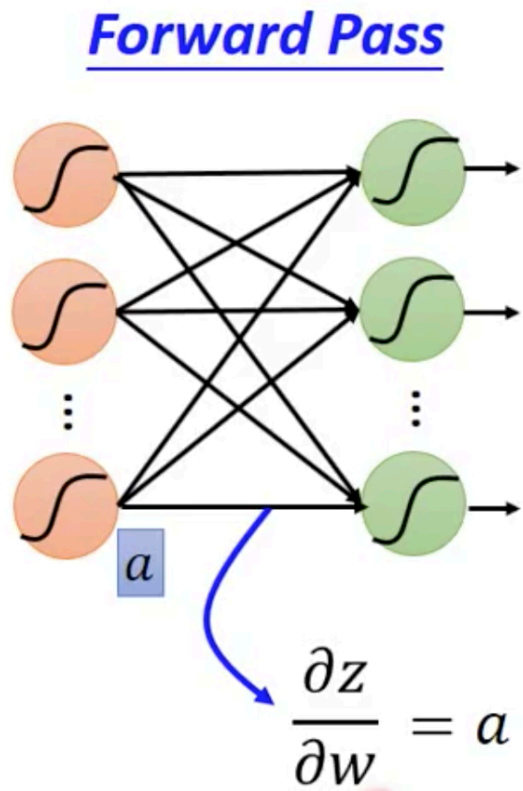Compute $\partial C / \partial z$ from the output layer

# Backpropagation – Summary
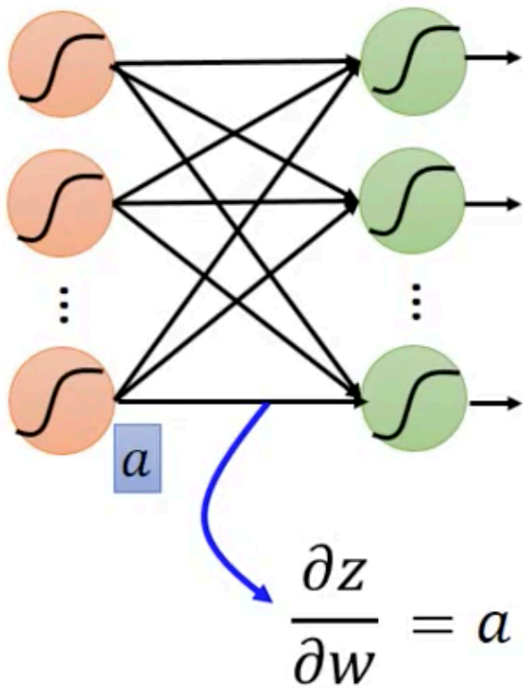
**_Forward Pass_**                    **_Backward Pass_**

# Backpropagation – Summary

**Forward Pass**

**Backward Pass**



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

# Backpropagation – Summary



**Forward Pass**                 **Backward Pass**

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

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