# Backpropagation

## Hung-yi Lee
李宏毅

# Gradient Descent

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

Starting Parameters $\qquad \theta^0 \longrightarrow \theta^1 \longrightarrow \theta^2 \longrightarrow$ ......

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

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

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

Millions of parameters ......

To compute the gradients efficiently, we use ***backpropagation***.
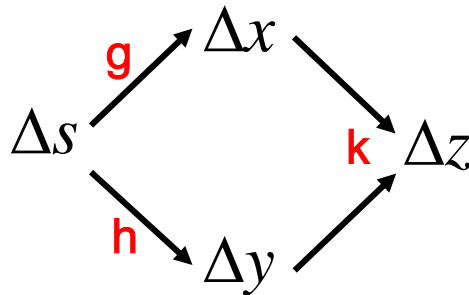
# Chain Rule <span style="color:red">Some preliminaries</span>

**_Case 1_**    $y = g(x)$    $z = h(y)$

$$\Delta x \xrightarrow{\text{ g }} \Delta y \xrightarrow{\text{ h }} \Delta z \qquad \frac{dz}{dx} = \frac{dz}{dy}\frac{dy}{dx}$$
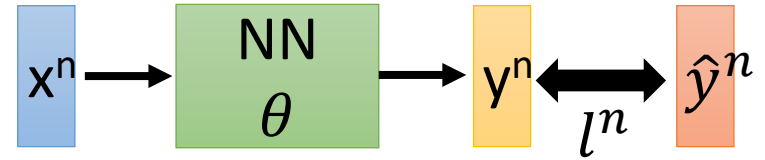
**_Case 2_**

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



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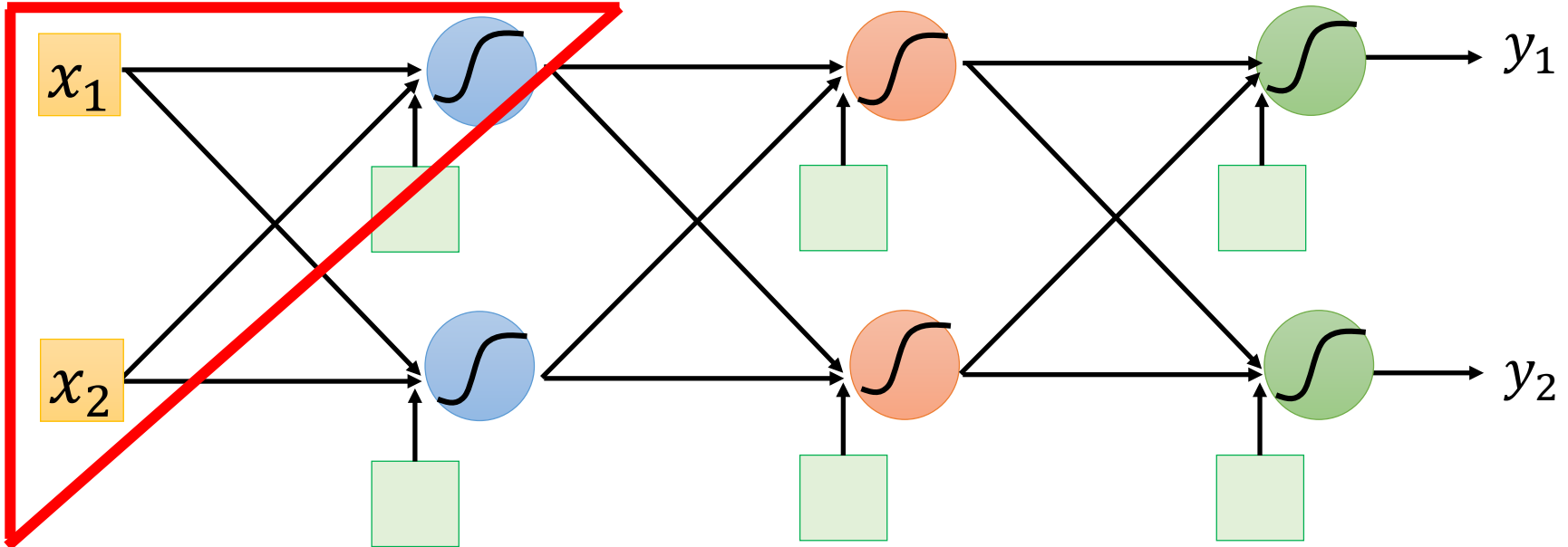
# Backpropagation

$$x^n \longrightarrow \boxed{\begin{array}{c} NN \\ \theta \end{array}} \longrightarrow y^n \longleftrightarrow \underset{l^n}{\hat{y}^n}$$

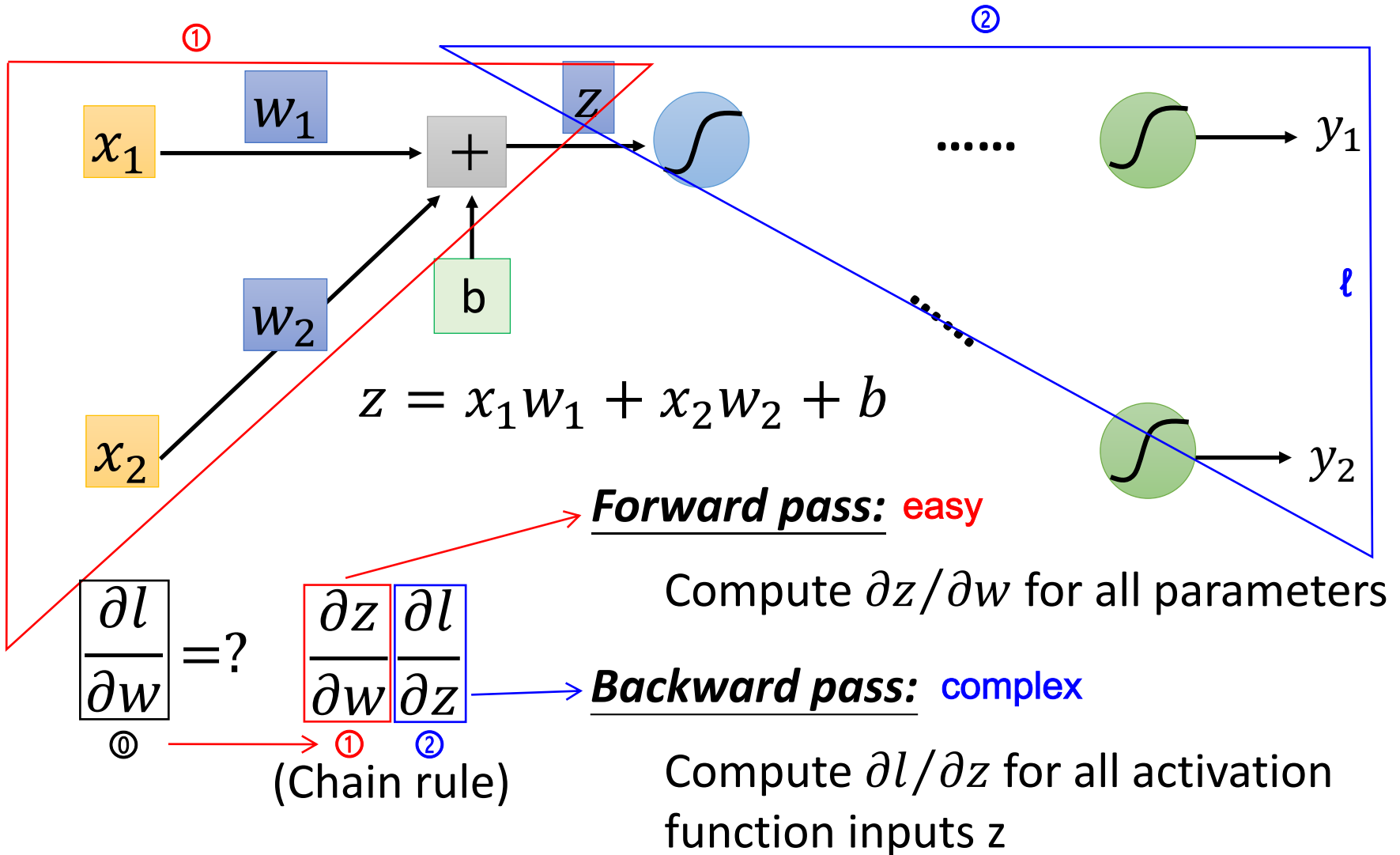<span style="color:blue">Ignore the superscript n here. (One data at a time.)</span>

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

<span style="color:red">We can just understand this part.</span>
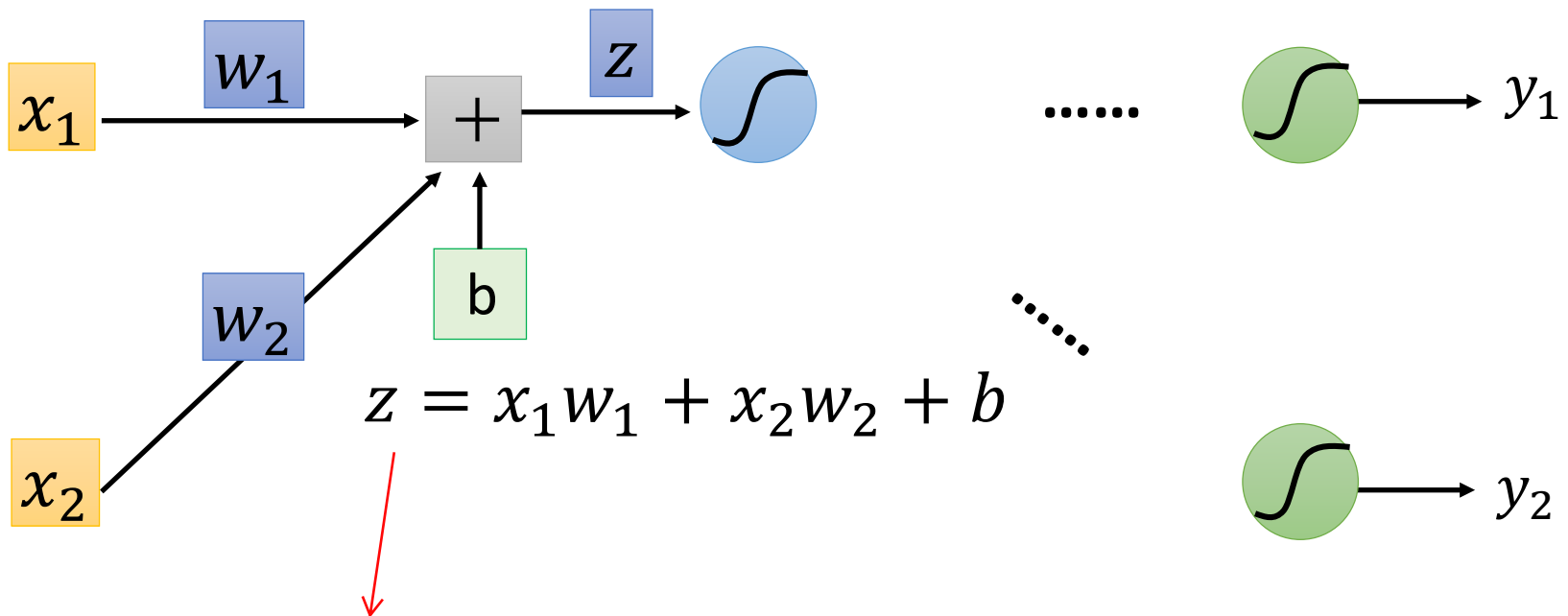
<span style="color:red">Look at this part first.</span>

$x_1$   $x_2$   $y_1$   $y_2$



4/17

# Backpropagation

①

②

$w_1$

$z$

$x_1$ → $+$ → ∫ → ...... → ∫ → $y_1$

$w_2$

b

∫ → $y_2$

$\ell$

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

**Forward pass:** easy

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

$$\left|\frac{\partial l}{\partial w}\right| = ? \quad \frac{\partial z}{\partial w} \frac{\partial l}{\partial z}$$

⓪ → ① ②

(Chain rule)

**Backward pass:** complex

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

# Backpropagation − <u>Forward pass</u>

Compute $\boxed{\partial z/\partial w}$ ① for all parameters

$$z = x_1 w_1 + x_2 w_2 + b$$
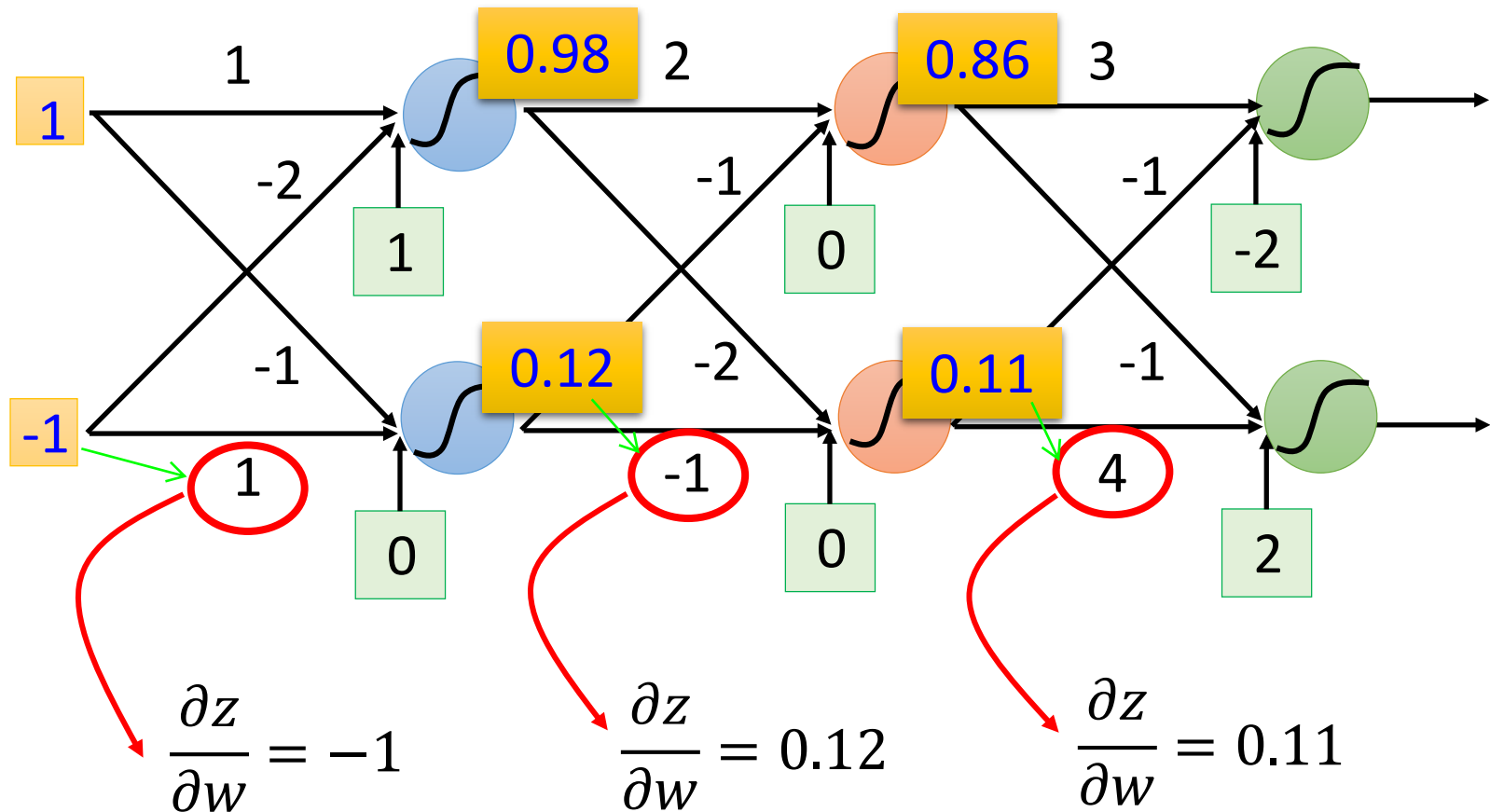
$\partial z/\partial w_1 =?\ x_1$

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

The value of the input connected by the weight

# Backpropagation − Forward pass
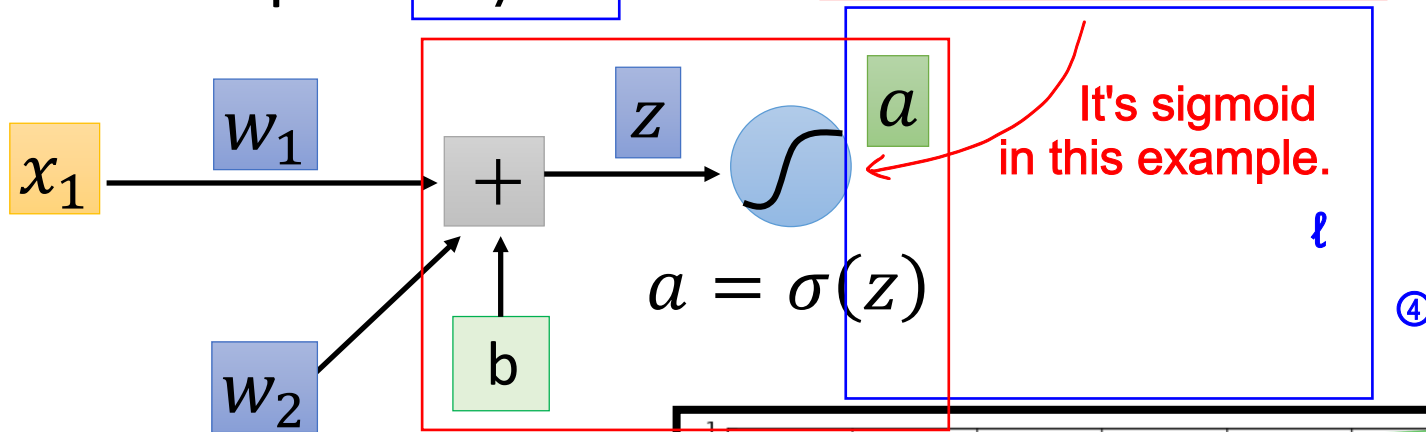
Compute $\partial z / \partial w$ ① for all parameters    It's just the input.



$$\frac{\partial z}{\partial w} = -1 \qquad \frac{\partial z}{\partial w} = 0.12 \qquad \frac{\partial z}{\partial w} = 0.11$$

# Backpropagation − <u>Backward pass</u>

Compute $\boxed{\partial l / \partial z}$ ② for all <u>activation function</u> inputs z ⇒ ③

$$a = \sigma(z)$$

It's sigmoid in this example.

$x_1$   $w_1$   +   $z$   ∫   $a$   $\ell$

b

$x_2$   $w_2$

$$\underbrace{\frac{\partial l}{\partial z}}_{②} = \underbrace{\frac{\partial a}{\partial z}}_{③} \underbrace{\frac{\partial l}{\partial a}}_{④}$$

complex

↘ easy

➡ $\sigma'(z)$

$\sigma(z)$

$\sigma'(z)$

# Backpropagation − Backward pass

②

Compute $\boxed{\partial l / \partial z}$ from the output layer ⇒ ④

$w_1$  $z$  $a$  $w_3$  $z'$

$x_1$

$a = \sigma(z)$

b

$z' = aw_3 + \cdots$

$w_2$

In z', the only term that is related to a is w₃.

$x_2$

We assume that Δa will only influence Δz' and Δz".

$w_4$

$z''$

④

$$\frac{\partial l}{\partial z} = \frac{\partial a}{\partial z} \frac{\partial l}{\partial a} \qquad \boxed{\frac{\partial l}{\partial a}} = \frac{\partial z'}{\partial a} \frac{\partial l}{\partial z'} + \frac{\partial z''}{\partial a} \frac{\partial l}{\partial z''} \text{ (Chain rule)}$$
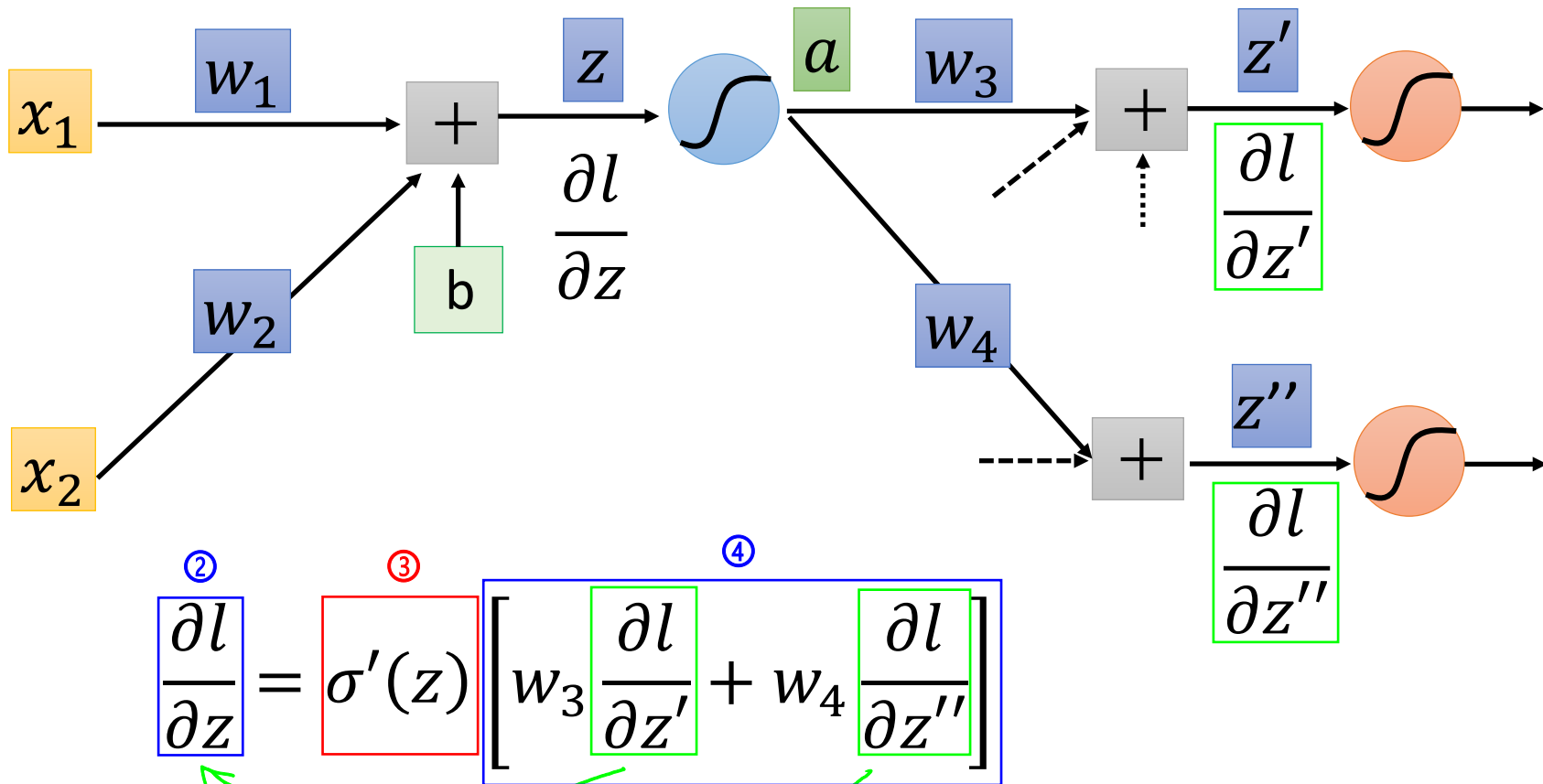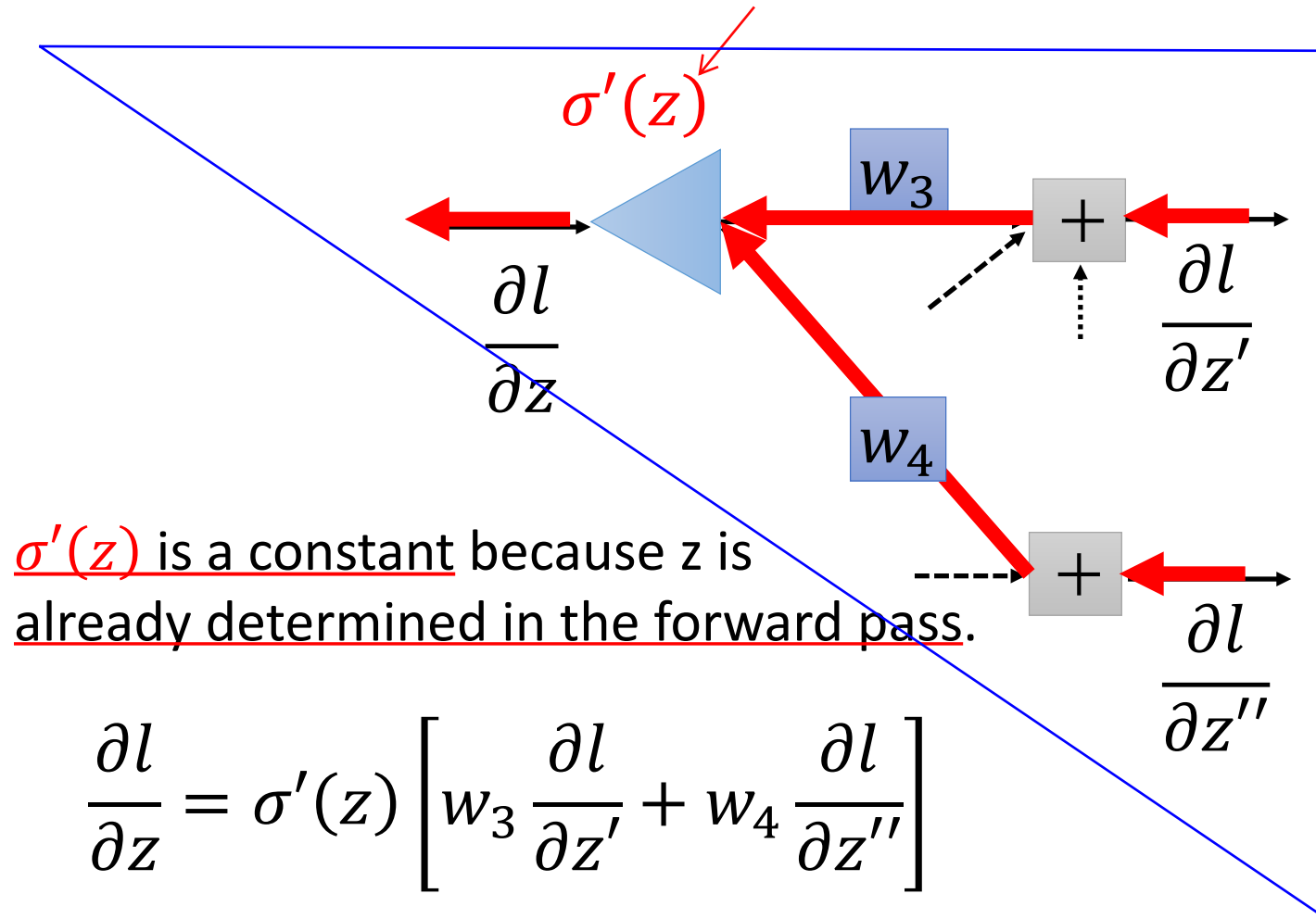
$w_3$  ?  $w_4$  ?

Assumed it's known

Forward pass

Backward pass

9/17

# Backpropagation − Backward pass

②

Compute $\boxed{\partial l / \partial z}$ from the output layer $\Rightarrow$ ④



②     ③     ④

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

If we know the latter, we can calculate the front.

# Backpropagation − Backward pass

Amplifier (It's multiplication, not addition.)

$\sigma'(z)$

$w_3$

$w_4$

$\dfrac{\partial l}{\partial z}$

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

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

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

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

# Backpropagation − Backward pass

Compute $\partial l / \partial z$ from the output layer $\Rightarrow$ ④

$x_1$

$w_1$

$w_2$

$x_2$

$+$

b

$z$

$\dfrac{\partial l}{\partial z}$

$\int$

$a$

$w_3$

$w_4$

$+$

$z'$

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

$\int$

$y_1$

$+$

$z''$

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

$\int$

$y_2$

If we are at the last layer.

## *Case 1. Output Layer*

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

easy · known

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

easy · known

Done!

# Backpropagation − Backward pass

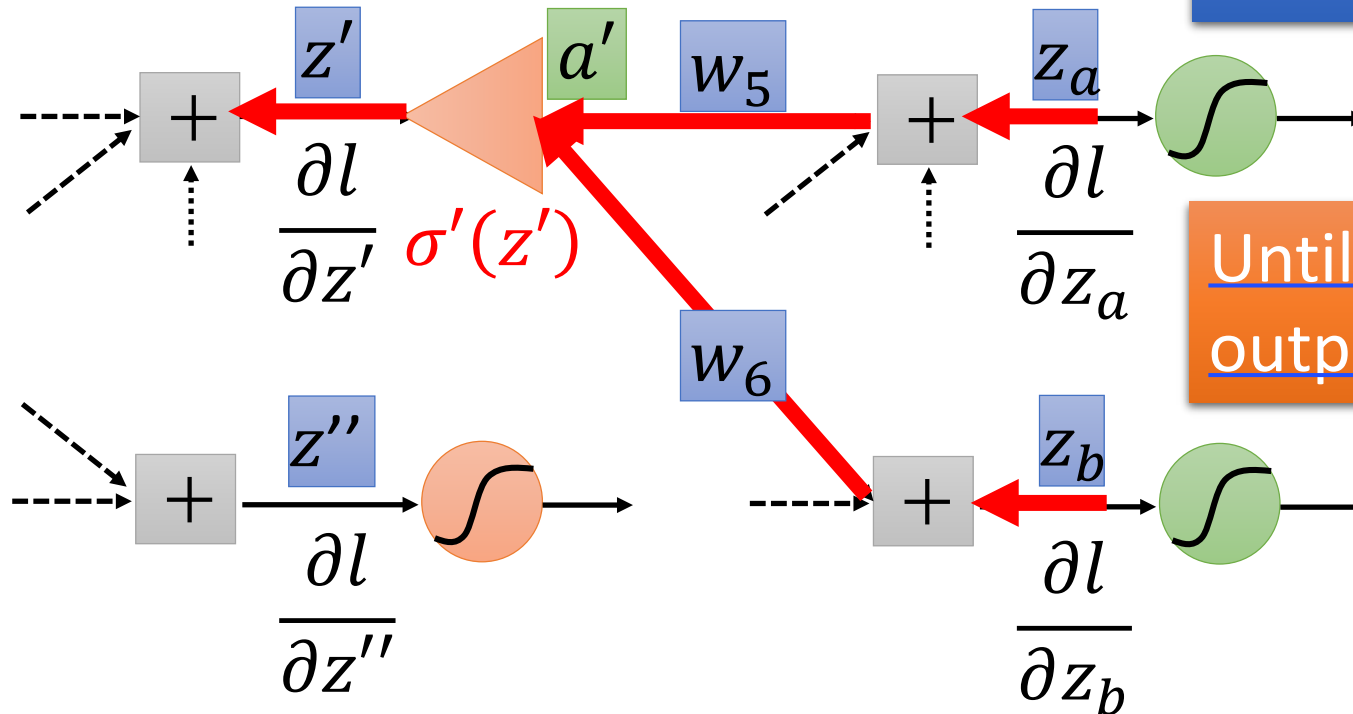Compute $\partial l/\partial z$ from the output layer ⇒ ④

## *Case 2. Not Output Layer*  If we are NOT at the last layer.

# Backpropagation − Backward pass

Compute $\partial l / \partial z$ from the output layer $\Rightarrow$ ④

**_Case 2. Not Output Layer_**
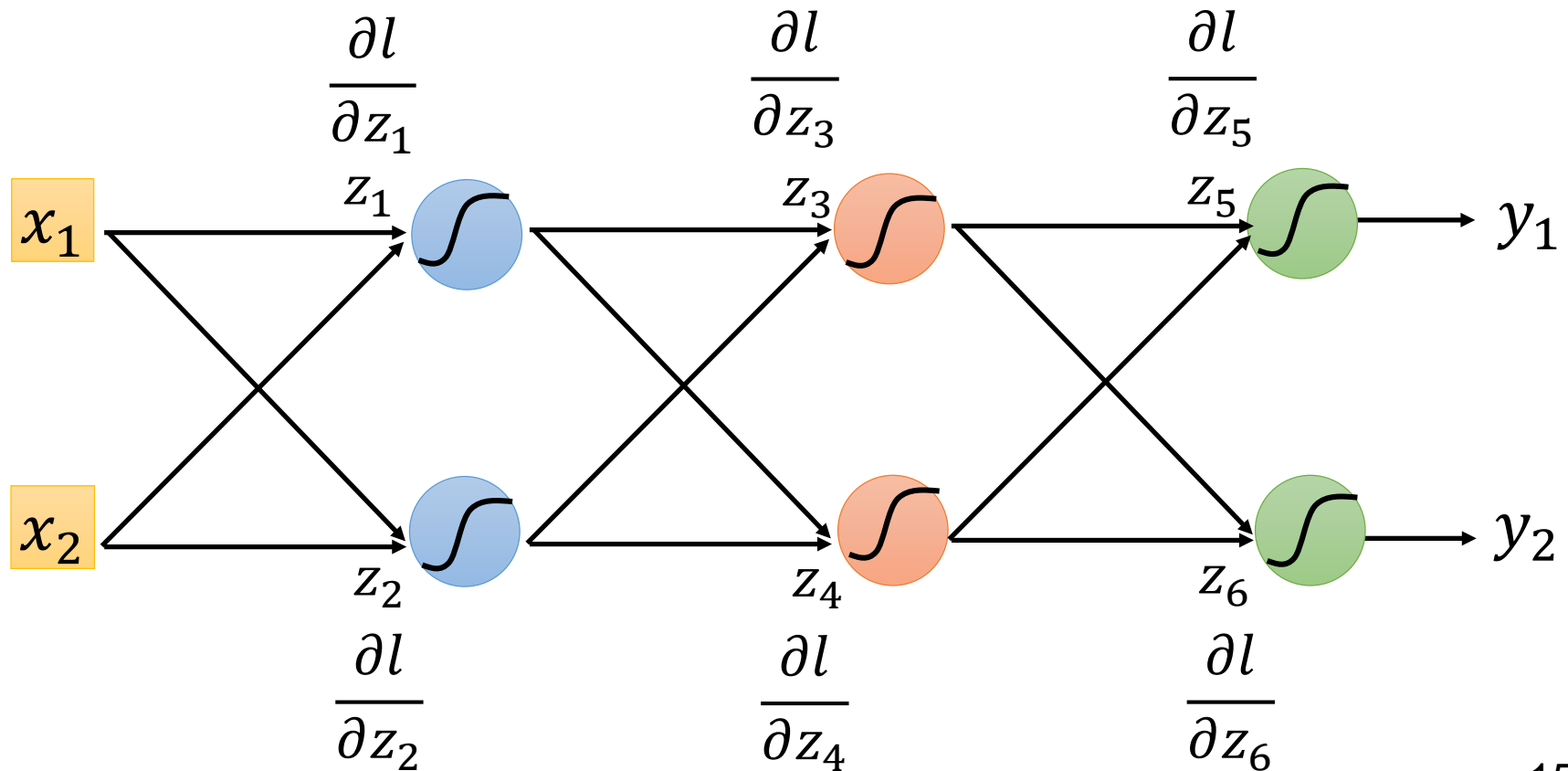


Compute $\partial l / \partial z$ recursively

Until we reach the output layer ……

# Backpropagation − Backward Pass
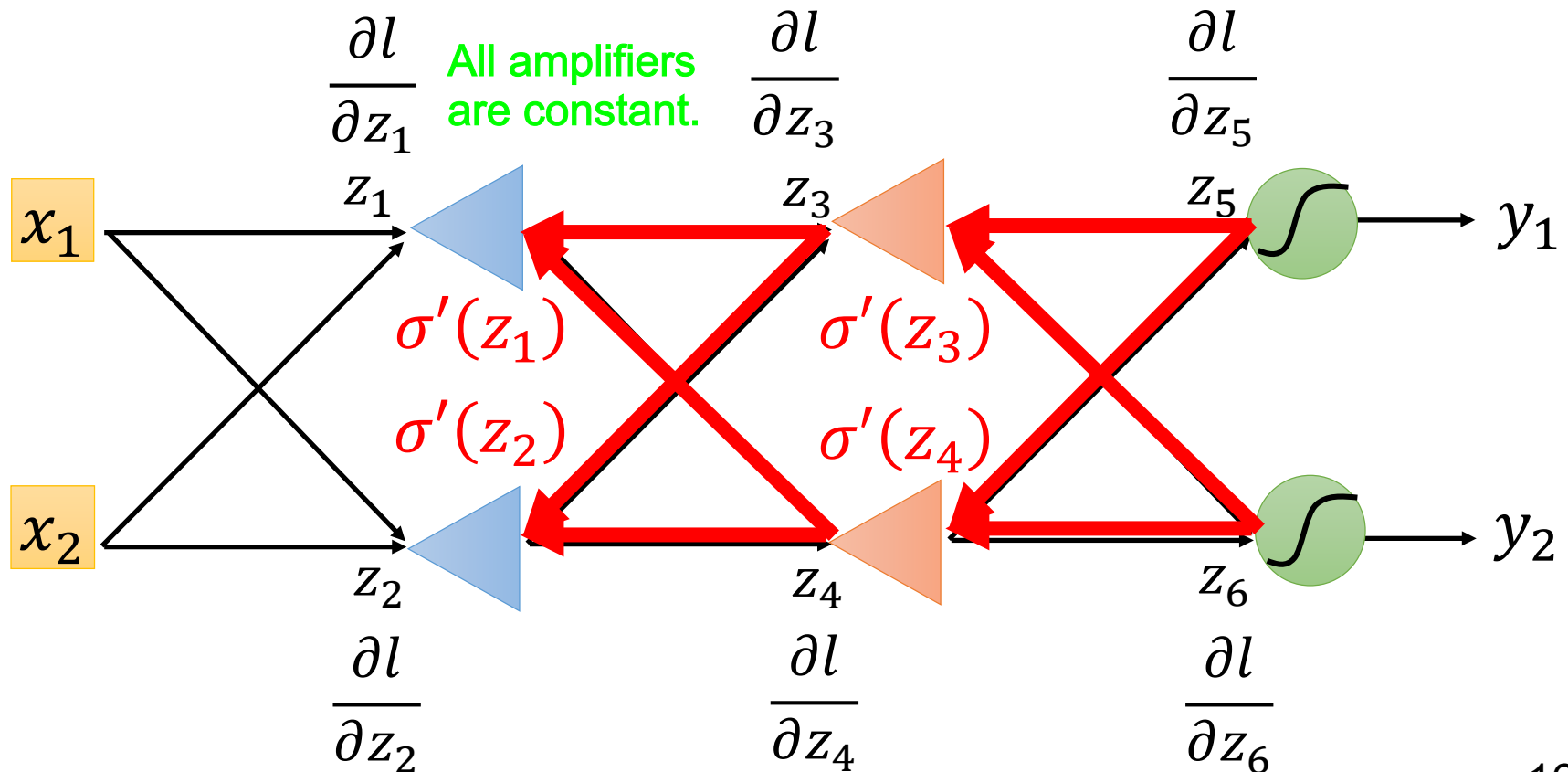
Compute $\partial l / \partial z$ for all activation function inputs z ⇒ ③

Compute $\partial l / \partial z$ from the output layer ⇒ ④ Backward pass

# Backpropagation − Backward Pass
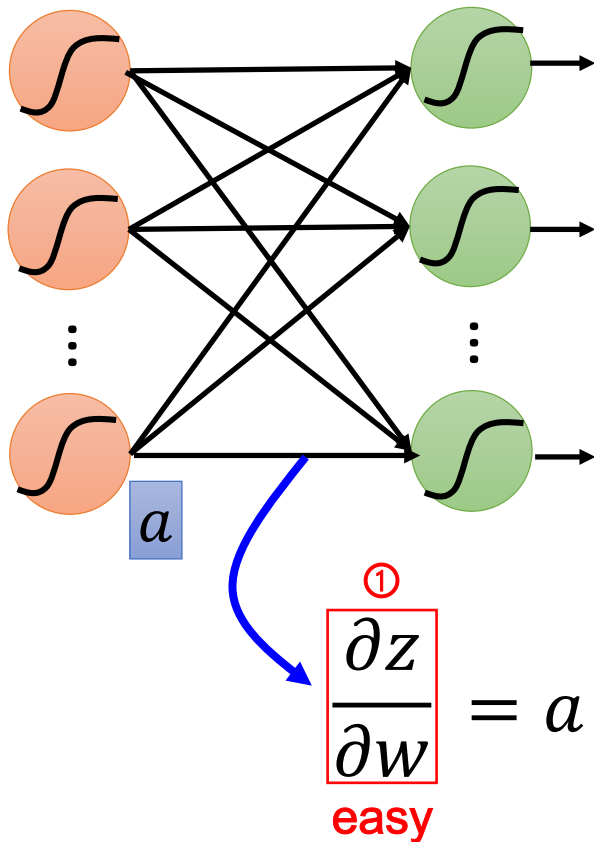
Compute $\partial l / \partial z$ for all activation function inputs z ⇒ ③

Compute $\partial l / \partial z$ from the output layer ⇒ ④

# Backpropagation – Summary

**_Forward Pass_**

**_Backward Pass_**



$$\textcircled{1} \quad \frac{\partial z}{\partial w} = a$$

easy

$$\textbf{X}$$

$$\textcircled{2} \quad \frac{\partial l}{\partial z}$$

complex

$$= \quad \textcircled{0} \quad \frac{\partial l}{\partial w}$$

for all w