# Backpropagation

## Chun-Hsiang Chan

Undergraduate Program in Intelligent Computing and Big Data, Chung Yuan Christian University
Master Program in Intelligent Computing and Big Data, Chung Yuan Christian University
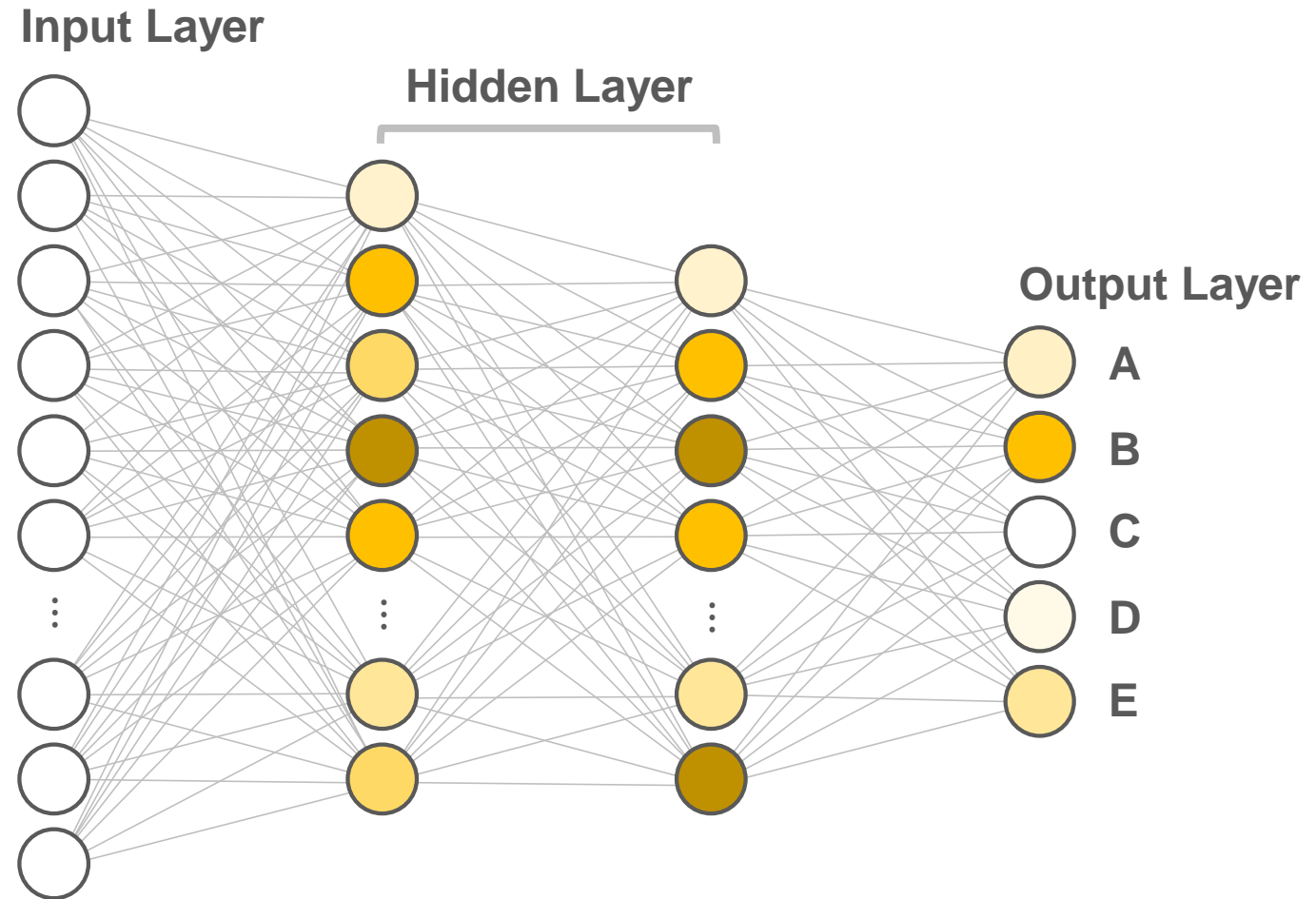
**Slide credit:** Prof. Hung-Yi Lee & Prof. Yun-Nung Chen

# Outline

- Forward and Back Propagation
- Chain Rule
- $\partial C(\theta)/\partial w_{ij}^l$
- $\partial C(\theta)/\partial z_i^l$
- $\partial C(\theta)/\partial z_i^l = \delta_i^l$
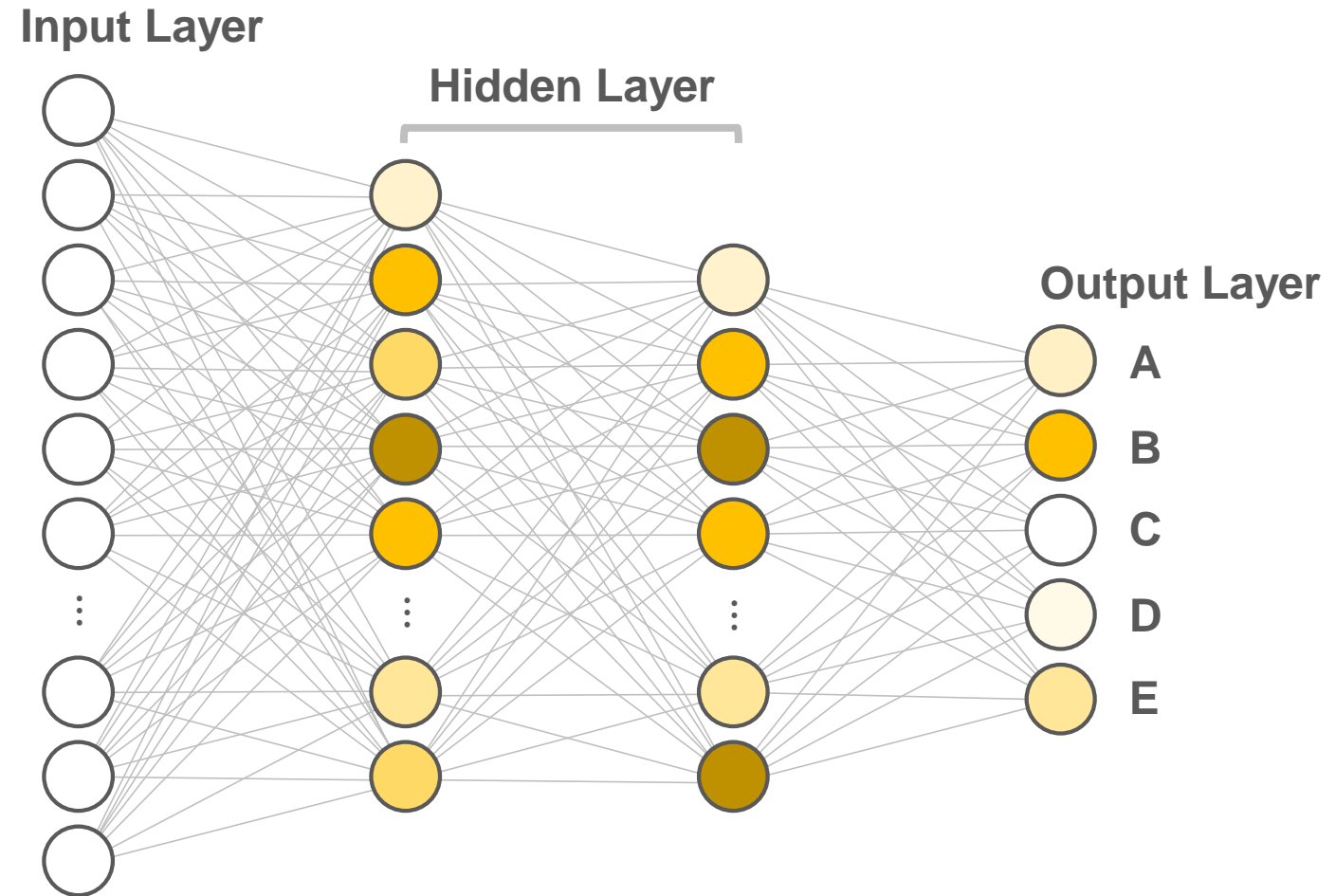- Back Propagation
- Summary

# Forward Propagation

- Forward propagation is the process of calculating the output of a neural network given an input. In other words, it is the process of feeding input data through the network's layers in a forward direction to produce an output.

- During forward propagation, each neuron in a layer receives input from the previous layer and applies an activation function to produce an output, which is then passed to the next layer. The output of the final layer is the predicted output of the neural network.

**Input Layer**

**Hidden Layer**

**Output Layer**

A

B

C

D

E

# Back Propagation

- Backpropagation computes the gradient of the loss function with respect to the weights of the network for a single input–output example, and does so efficiently, unlike a naive direct computation of the gradient with respect to each weight individually.



**Input Layer**

**Hidden Layer**

**Output Layer**

A
B
C
D
E

# Chain Rule

$$\Delta w \rightarrow \Delta x \rightarrow \Delta y \rightarrow \Delta z$$

$$\frac{\partial z}{\partial w} = \frac{\partial z}{\partial y}\frac{\partial y}{\partial x}\frac{\partial x}{\partial w}$$

$$= f'(y)f'(x)f'(w)$$

**Forward** propagation for computing the **cost**

$$= f'\left(f(f(w))\right)f'(f(w))f'(w)$$

**Back** propagation for computing the **gradient**

# Looking Back – Gradient Decsent

$$y = f(x) = \sigma(W^L \dots \sigma(W^2 \sigma(W^1 x + b^1) + b^2) \dots + b^L)$$

$$\theta = \{W^1, b^1, W^2, b^2, \dots, W^L, b^L\}$$

$$W^l = \begin{bmatrix} w^l_{11} & w^l_{12} & \cdots \\ w^l_{21} & w^l_{22} & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix}, b^l = \begin{bmatrix} \vdots \\ b^l_i \\ \vdots \end{bmatrix}$$

$$\nabla C(\theta) = \begin{bmatrix} \vdots \\ \dfrac{\partial C(\theta)}{\partial w^l_{ij}} \\ \vdots \\ \dfrac{\partial C(\theta)}{\partial b^l_i} \end{bmatrix}$$

**Algorithm**
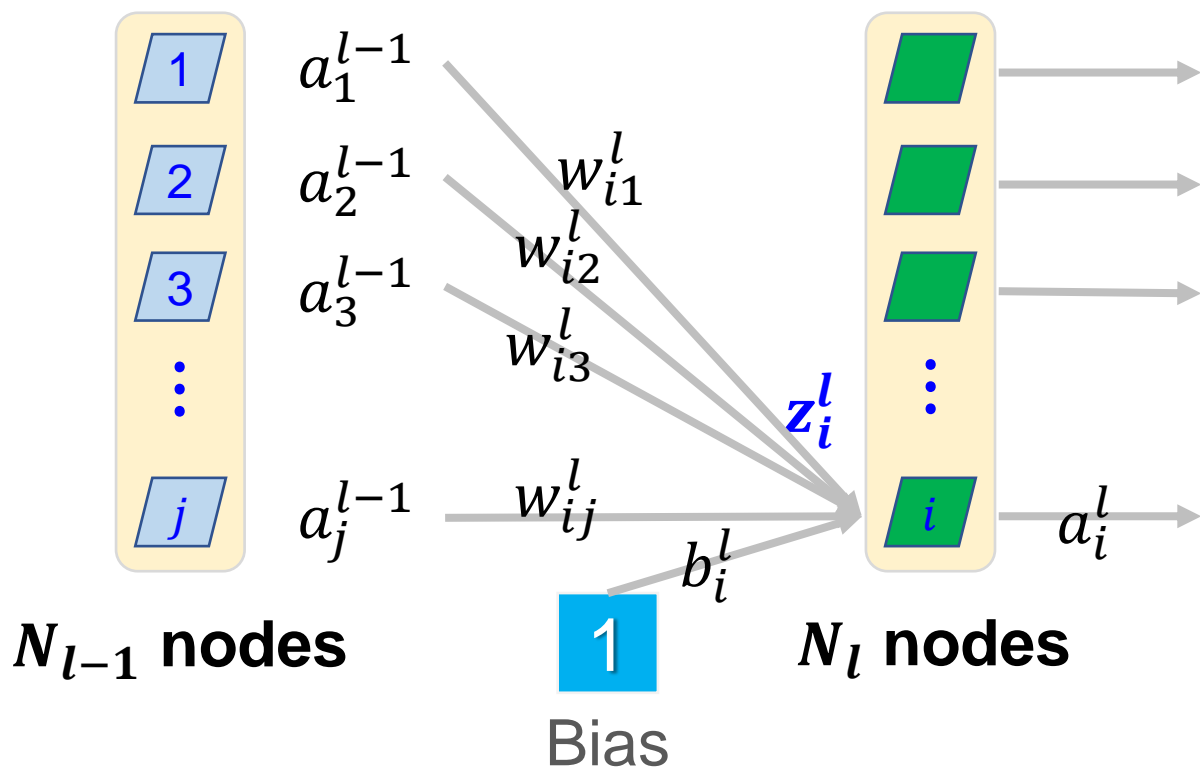Initialization: start at $\theta^0$
while ($\theta^{(i+1)} \neq \theta^i$)
{

    compute the gradients at $\theta^i$
    update parameters
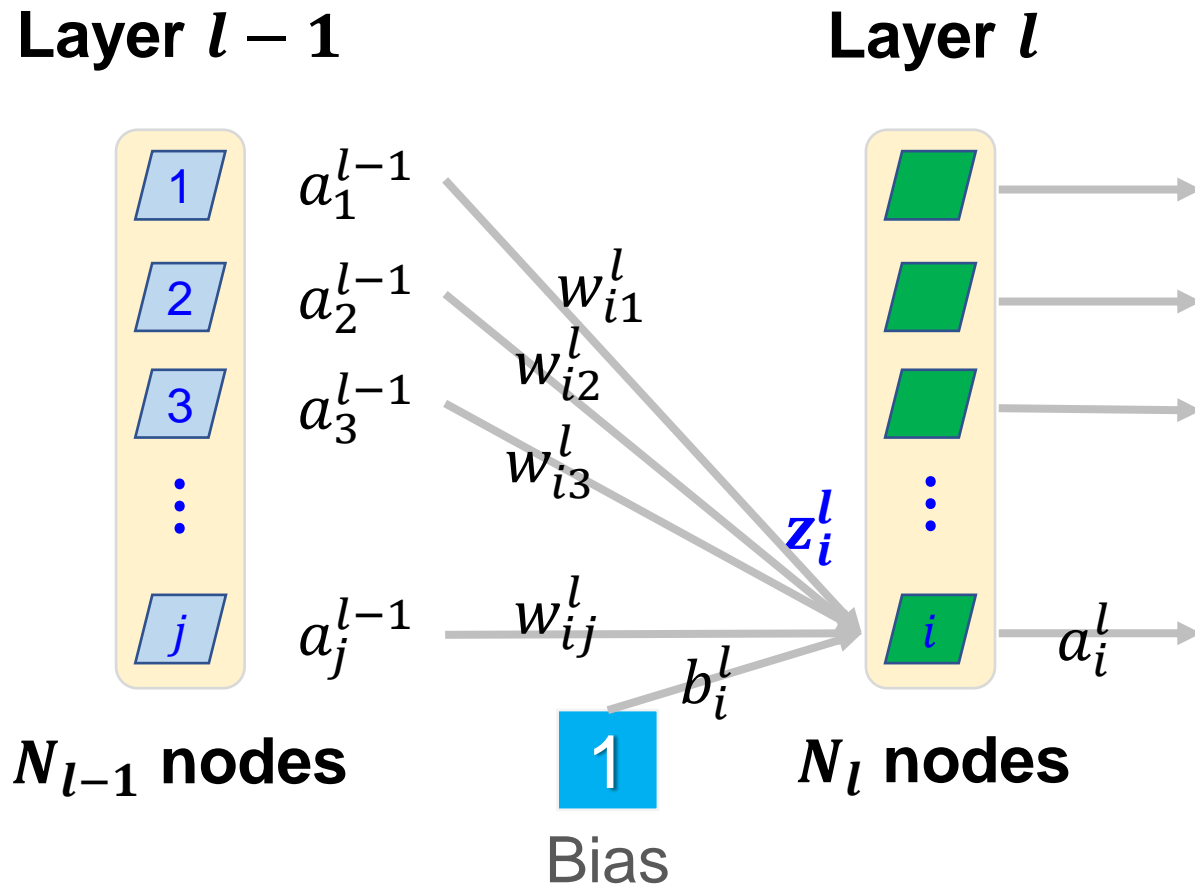    $\theta^{i+1} \leftarrow \theta^i - \eta \nabla_\theta C(\theta^i)$

}

# $\partial C(\theta)/\partial w_{ij}^l$

**Layer $l-1$**

**Layer $l$**



$a_1^{l-1}$

$a_2^{l-1}$

$w_{i1}^l$

$a_3^{l-1}$

$w_{i2}^l$

$w_{i3}^l$

$z_i^l$

$a_j^{l-1}$  $w_{ij}^l$

$b_i^l$

$a_i^l$

$N_{l-1}$ **nodes**

**1**

Bias

$N_l$ **nodes**

$$\frac{\partial C(\theta)}{\partial w_{ij}^l} = \frac{\partial C(\theta)}{\partial z_i^l} \boxed{\frac{\partial z_i^l}{\partial w_{ij}^l}}$$

# $\partial z_i^l / \partial w_{ij}^l$ $(l > 1)$

**Layer $l-1$**

**Layer $l$**



$1$

$a_1^{l-1}$

$2$

$a_2^{l-1}$

$w_{i1}^l$

$3$

$a_3^{l-1}$

$w_{i2}^l$

$w_{i3}^l$

$z_i^l$

$j$

$a_j^{l-1}$

$w_{ij}^l$

$i$

$a_i^l$

$b_i^l$
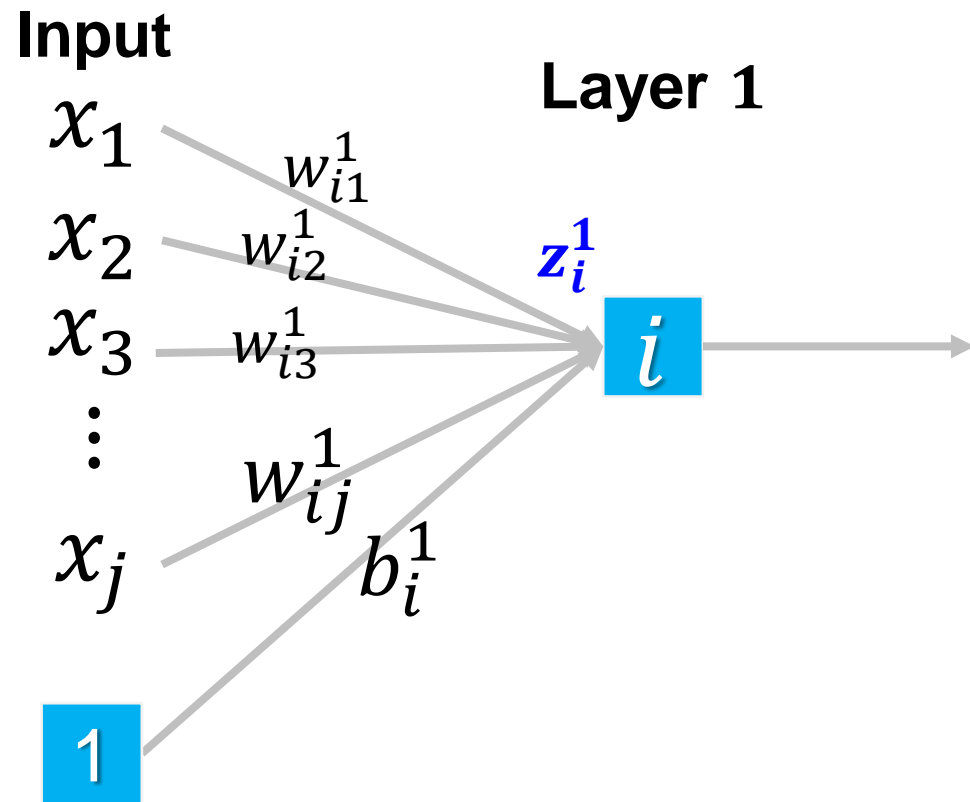
**$N_{l-1}$ nodes**

$1$

**$N_l$ nodes**

Bias

$$z^l = W^l a^{l-1} + b^l$$

$$z_i^l = \sum_j w_{ij}^l a_j^{l-1} + b_i^l$$

$$\frac{\partial z_i^l}{\partial w_{ij}^l} = a_j^{l-1}$$

$$\partial z_i^l / \partial w_{ij}^l \ (l = 1)$$

**Input**

**Layer 1**

$x_1$

$\quad w_{i1}^1$

$x_2$

$\quad w_{i2}^1$

$\quad\quad z_i^1$

$x_3$

$\quad w_{i3}^1$

$\qquad i$

$\vdots$

$\quad w_{ij}^1$
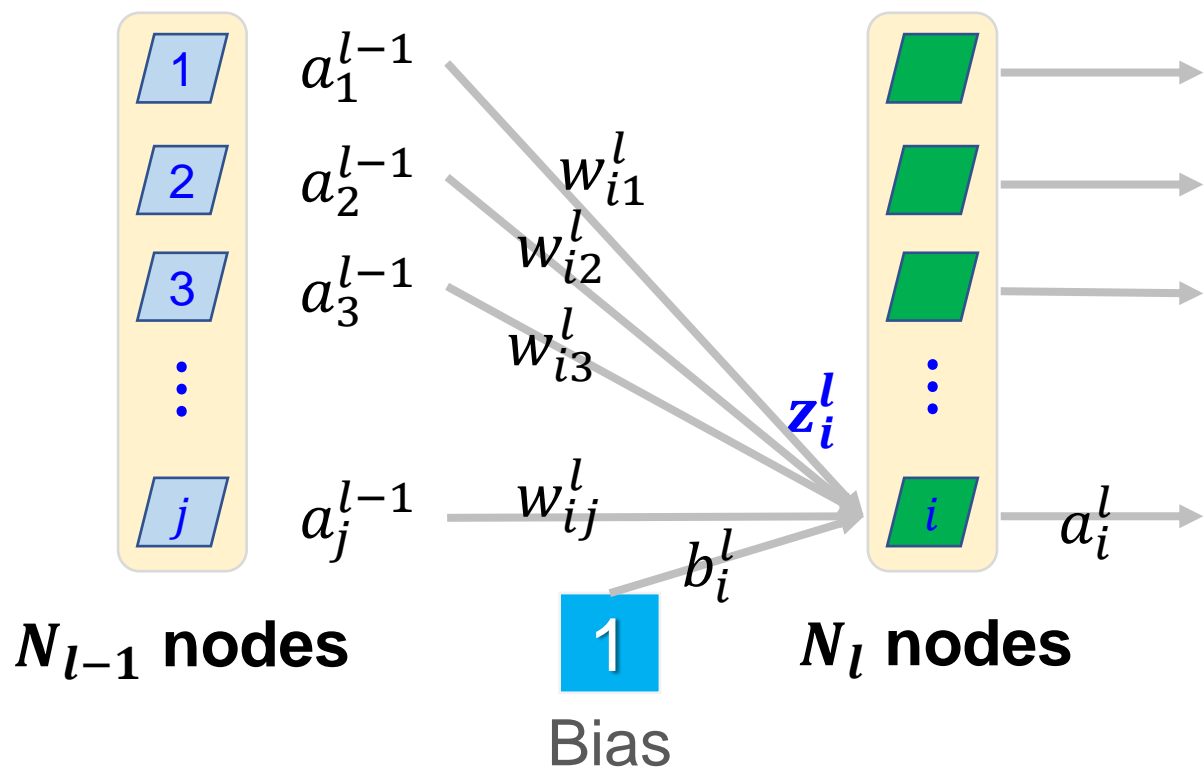
$x_j$

$\quad\quad b_i^1$

$1$

$$z^1 = W^1 x + b^1$$

$$z_i^1 = \sum_j w_{ij}^1 x_j + b_i^1$$

$$\frac{\partial z_i^1}{\partial w_{ij}^1} = x_j$$

# $\partial C(\theta)/\partial w_{ij}^l$

**Layer $l-1$**                    **Layer $l$**



$a_1^{l-1}$

$a_2^{l-1}$        $w_{i1}^l$

$a_3^{l-1}$        $w_{i2}^l$

        $w_{i3}^l$

$\vdots$            $z_i^l$

$j$    $a_j^{l-1}$    $w_{ij}^l$        $i$    $a_i^l$

            $b_i^l$

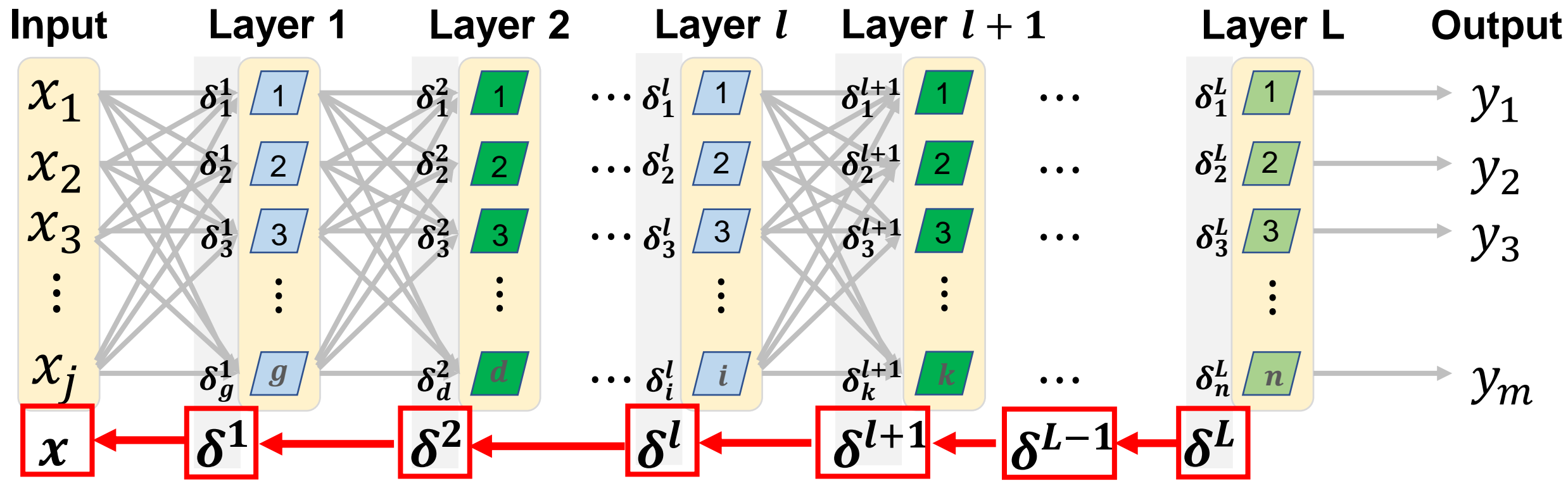$N_{l-1}$ **nodes**        **1**        $N_l$ **nodes**

Bias

$$\frac{\partial C(\theta)}{\partial w_{ij}^l} = \frac{\partial C(\theta)}{\partial z_i^l} \frac{\partial z_i^l}{\partial w_{ij}^l}$$

$$\frac{\partial z_i^l}{\partial w_{ij}^l} = \begin{cases} a_j^{l-1} & , l > 1 \\ x_j & , l = 1 \end{cases}$$

$$\partial C(\pmb{\theta})/\partial z_i^l \qquad \frac{\partial C(\theta)}{\partial w_{ij}^l} = \frac{\partial C(\theta)}{\partial z_i^l} \frac{\partial z_i^l}{\partial w_{ij}^l}$$

$\delta_i^l: the\ propagated\ gradient\ corresponding\ to\ the\ l-th\ layer$

$$\partial C(\theta)/\partial z_i^l = \delta_i^l$$

- **Procedure: from layer $L$ to layer 1**

1. Initialization: compute $\delta^L$

$$\delta_i^L = \frac{\partial C}{\partial z_i^L} = \frac{\partial C}{\partial y_i}\frac{\partial y_i}{\partial z_i^L}, \quad \boxed{\frac{\partial C}{\partial y_i} \; depends \; on \; the \; loss \; function}$$

$$\Delta z_i^L \rightarrow a_i^L = \Delta y_i \rightarrow \Delta C$$

2. Compute $\delta^L$ based on $\delta^{l+1}$

$$\partial C(\pmb{\theta})/\partial z_i^l = \pmb{\delta}_i^l$$

- **Procedure: from layer $L$ to layer 1**

1. Initialization: compute $\delta^L$

2. Compute $\delta^L$ based on $\delta^{l+1}$

$$\Delta z_i^L \rightarrow \Delta a_i^L = \Delta y_i \rightarrow \Delta C$$

$$\delta_i^L = \frac{\partial C}{\partial z_i^L}$$

$$\delta_i^L = \frac{\partial C}{\partial y_i}\frac{\partial y_i}{\partial z_i^L} = \frac{\partial C}{\partial y_i}\sigma'(z_i^L),$$
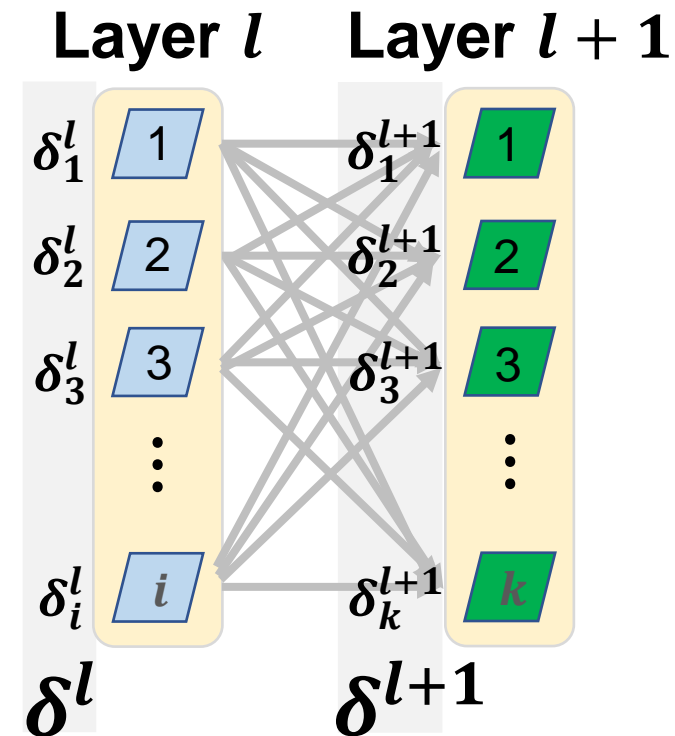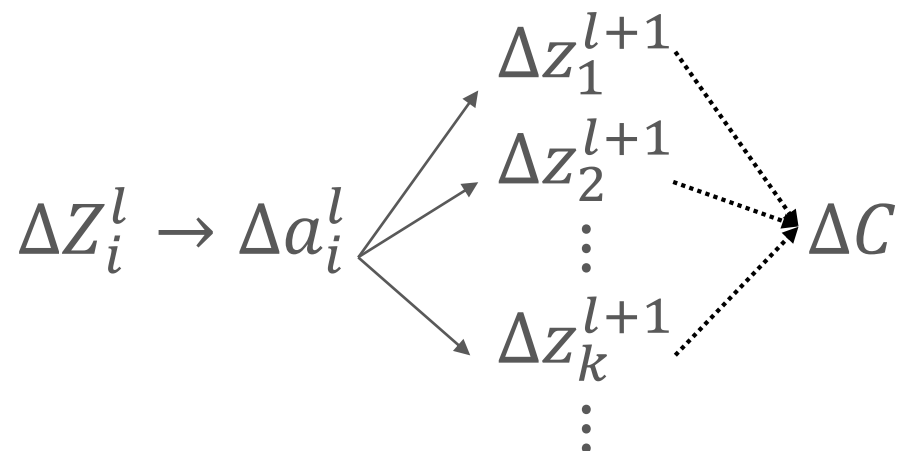
$$\text{where } y_i = a_i^L = \sigma(z_i^L)$$

$$\sigma'(z^L) = \begin{bmatrix} \sigma'(z_1^L) \\ \sigma'(z_2^L) \\ \vdots \\ \sigma'(z_i^L) \\ \vdots \end{bmatrix}, \nabla C(y) = \begin{bmatrix} \dfrac{\partial C}{\partial y_1} \\ \dfrac{\partial C}{\partial y_2} \\ \vdots \\ \dfrac{\partial C}{\partial y_i} \\ \vdots \end{bmatrix}$$

$$\pmb{\delta^L = \sigma'(z^L) \odot \nabla C(y)}$$

$$\partial C(\boldsymbol{\theta})/\partial z_i^l = \delta_i^l$$

- **Procedure: from layer $L$ to layer $1$**

1. Initialization: compute $\delta^L$

2. Compute $\delta^L$ based on $\delta^{l+1}$



$$\delta_i^l = \frac{\partial C}{\partial z_i^l} = \sum_k \left( \frac{\partial C}{\partial z_k^{l+1}} \frac{\partial z_k^{l+1}}{\partial a_i^l} \frac{\partial a_i^l}{\partial z_i^l} \right) = \frac{\partial a_i^l}{\partial z_i^l} \sum_k \left( \boxed{\frac{\partial C}{\partial z_k^{l+1}}} \frac{\partial z_k^{l+1}}{\partial a_i^l} \right)$$

$$\partial C(\theta)/\partial z_i^l = \delta_i^l$$
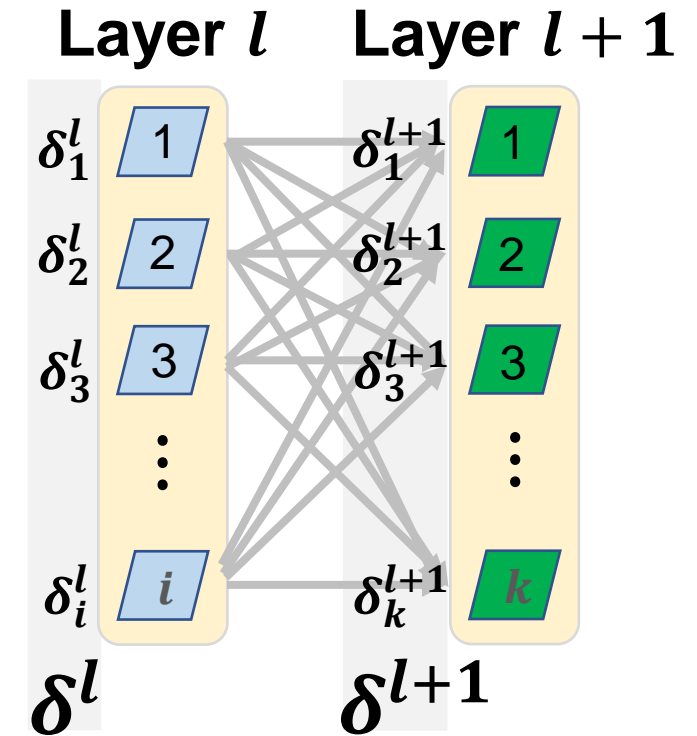
- **Procedure: from layer $L$ to layer $1$**

1. Initialization: compute $\delta^L$

2. Compute $\delta^L$ based on $\delta^{l+1}$

$$\delta_i^l = \frac{\partial a_i^l}{\partial z_i^l} \sum_k \frac{\partial z_k^{l+1}}{\partial a_i^l} \delta_k^{l+1}$$

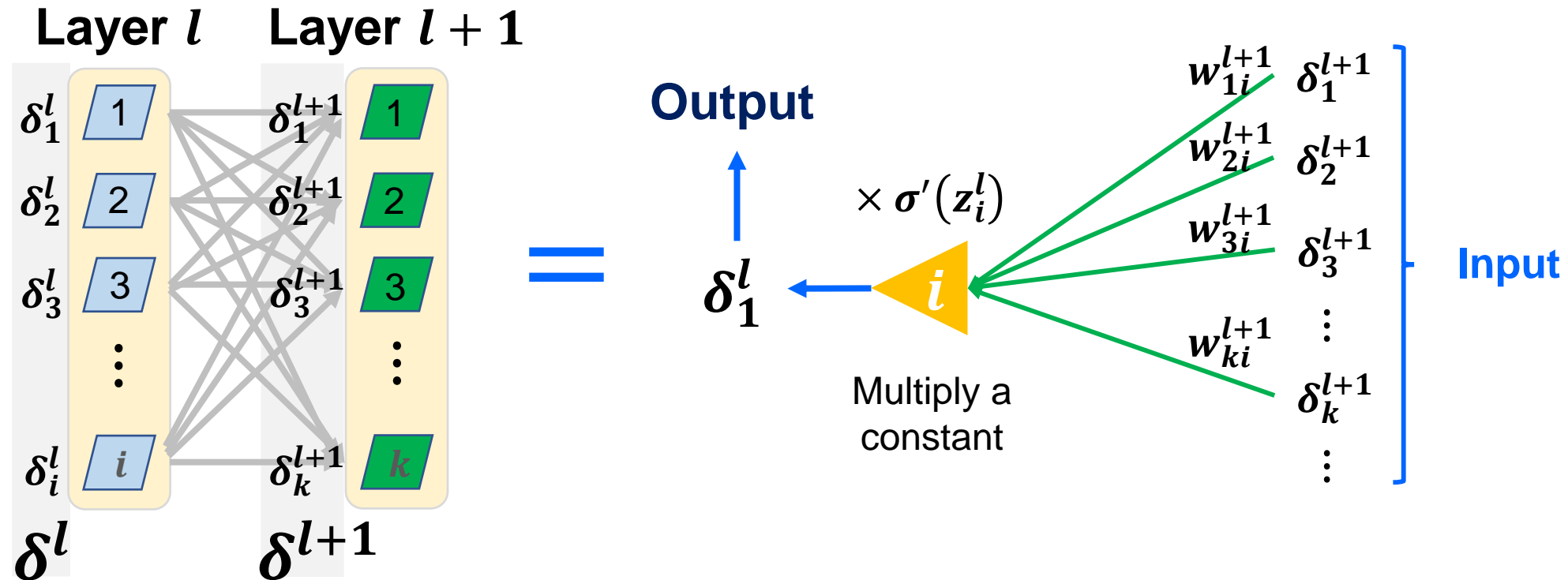$$\delta_i^l = \sigma'(z_i) \sum_k \boxed{\frac{\partial z_k^{l+1}}{\partial a_i^l}} \delta_k^{l+1} \quad \longrightarrow \quad \boxed{= \sum_i w_{ki}^{l+1} a_i^l + b_k^{l+1}}$$

$$\delta_i^l = \sigma'(z_i) \sum_k w_{ki}^{l+1} \delta_k^{l+1}$$

**Layer $l$**　　**Layer $l+1$**

$\delta_1^l$ — 1 ⟶ $\delta_1^{l+1}$ — 1

$\delta_2^l$ — 2 ⟶ $\delta_2^{l+1}$ — 2

$\delta_3^l$ — 3 ⟶ $\delta_3^{l+1}$ — 3

⋮ ⋮

$\delta_i^l$ — $i$ ⟶ $\delta_k^{l+1}$ — $k$

$\boldsymbol{\delta^l}$　　$\boldsymbol{\delta^{l+1}}$

$$\partial C(\theta)/\partial z_i^l = \delta_i^l$$

- The propagation: $\delta_i^l = \sigma'(z_i) \sum_k w_{ki}^{l+1} \delta_k^{l+1}$



**Layer** $l$     **Layer** $l+1$

$\delta_1^l$ [1]     $\delta_1^{l+1}$ [1]

$\delta_2^l$ [2]     $\delta_2^{l+1}$ [2]

$\delta_3^l$ [3]     $\delta_3^{l+1}$ [3]

$\delta_i^l$ [i]     $\delta_k^{l+1}$ [k]

$\delta^l$     $\delta^{l+1}$

$=$

**Output**

$\times \sigma'(z_i^l)$

$\delta_1^l \leftarrow$ [i]

Multiply a constant

$w_{1i}^{l+1}$  $\delta_1^{l+1}$
$w_{2i}^{l+1}$  $\delta_2^{l+1}$
$w_{3i}^{l+1}$  $\delta_3^{l+1}$
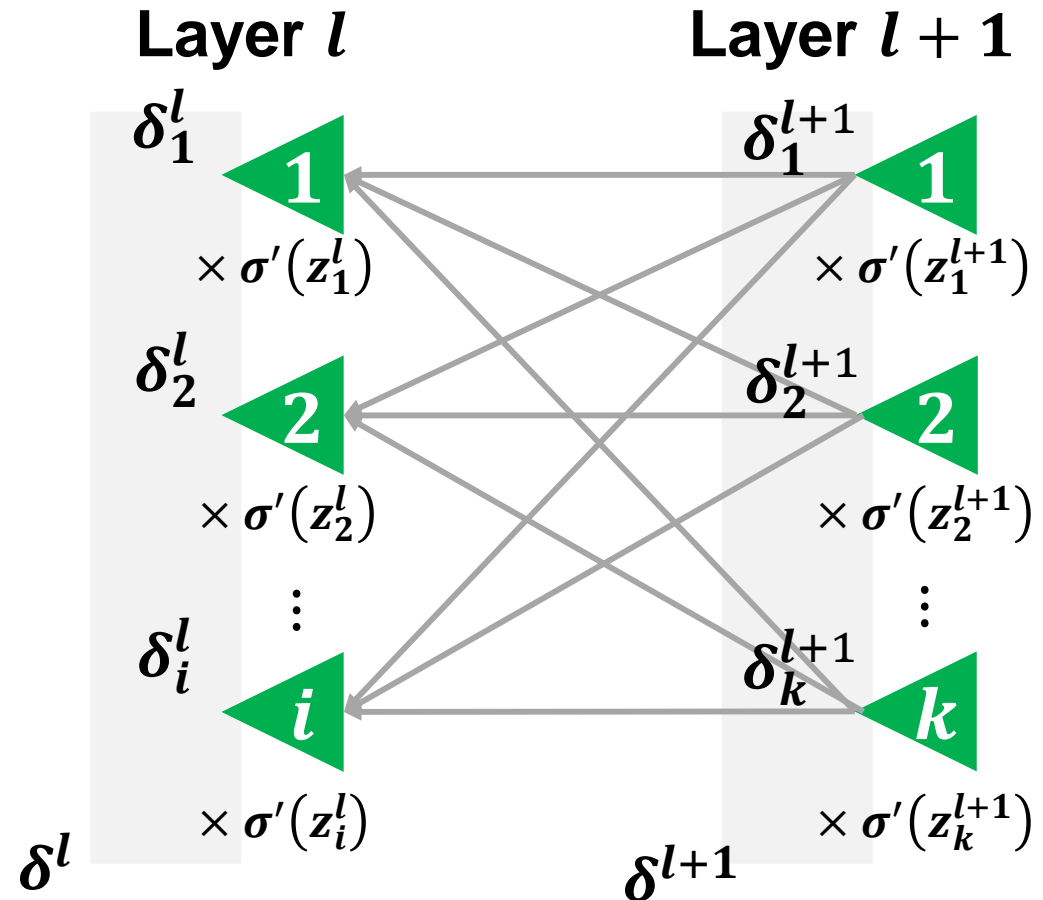$w_{ki}^{l+1}$  $\delta_k^{l+1}$

**Input**

$$\partial C(\theta)/\partial z_i^l = \delta_i^l$$

$$\delta_i^l = \sigma'(z_i)\sum_k w_{ki}^{l+1}\delta_k^{l+1}$$

$$\sigma'(z^l) = \begin{bmatrix} \sigma'(z_1^l) \\ \sigma'(z_2^l) \\ \vdots \\ \sigma'(z_i^l) \\ \vdots \end{bmatrix}$$

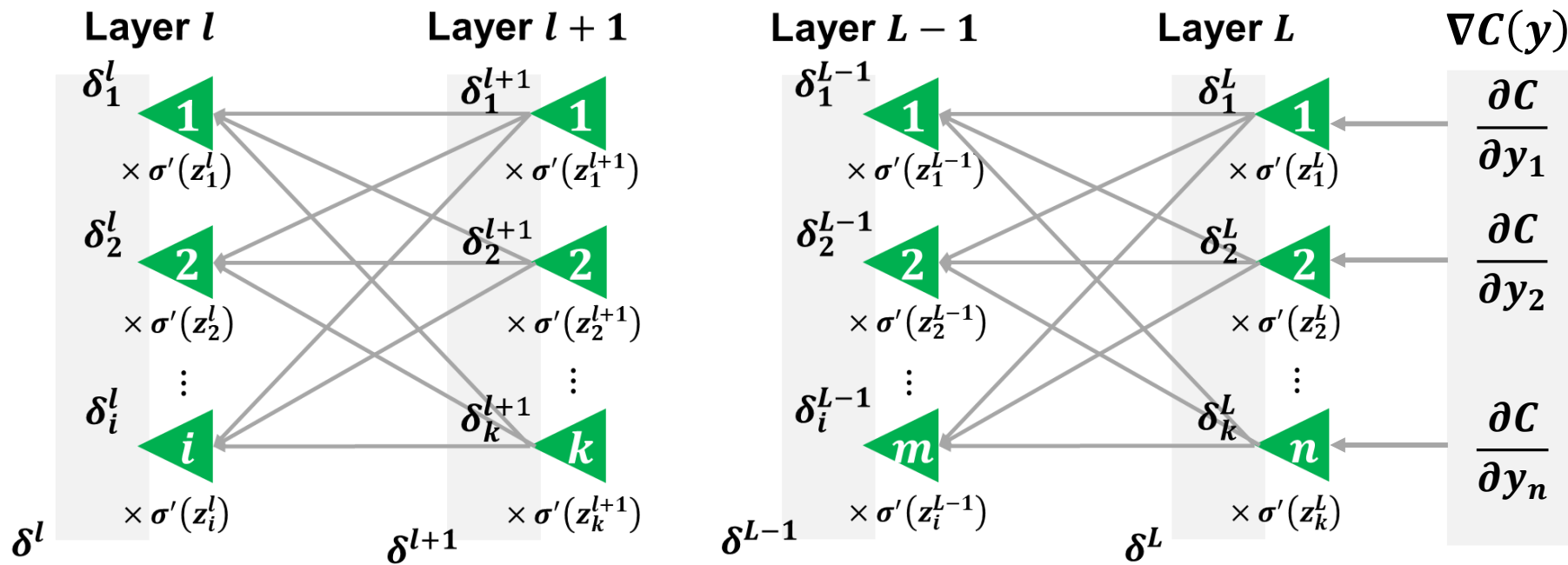$$\delta^l = \sigma'(z^l) \odot (W^{l+1})^T \delta^{l+1}$$

$$\partial C(\theta)/\partial z_i^l = \delta_i^l$$

$$\frac{\partial C(\theta)}{\partial w_{ij}^l} = \frac{\partial C(\theta)}{\partial z_i^l}\frac{\partial z_i^l}{\partial w_{ij}^l}$$

- **Procedure: from layer $L$ to layer $1$**
1. Initialization: compute $\delta^L$     $\delta^L = \sigma'(z^L) \odot \nabla C(y)$
2. Compute $\delta^{l-1}$ based on $\delta^l$    $\delta^l = \sigma'(z^l) \odot (W^{l+1})^T \delta^{l+1}$

# Back Propagation

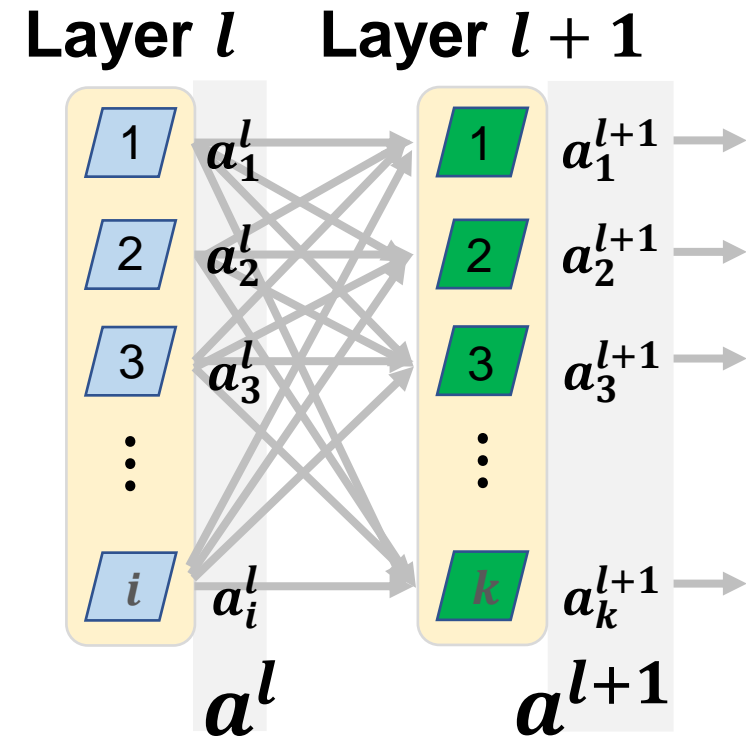$$\frac{\partial C(\theta)}{\partial w_{ij}^l} = \frac{\partial C(\theta)}{\partial z_i^l}\frac{\partial z_i^l}{\partial w_{ij}^l}$$

$$\frac{\delta z_i^i}{\partial w_{ij}^l} = \begin{cases} a_j^{l-1} & , l > 1 \\ x_j & , l = 1 \end{cases}$$

**Forward Direction**

$$\boldsymbol{z^1 = W^1 x + b^1} \qquad \boldsymbol{a^1 = \sigma(z^1)}$$

$$\vdots \qquad\qquad \vdots$$

$$\boldsymbol{z^l = W^l x^{l-1} + b^l} \qquad \boldsymbol{a^l = \sigma(z^l)}$$

$$\vdots \qquad\qquad \vdots$$

**Layer $l$**   **Layer $l+1$**

# Back Propagation

$$\frac{\partial C(\theta)}{\partial w_{ij}^l} = \frac{\partial C(\theta)}{\partial z_i^l} \frac{\partial z_i^l}{\partial w_{ij}^l}$$

$$\frac{\partial C(\theta)}{\partial z_i^l} = \delta_i^l$$

**Backward Direction**

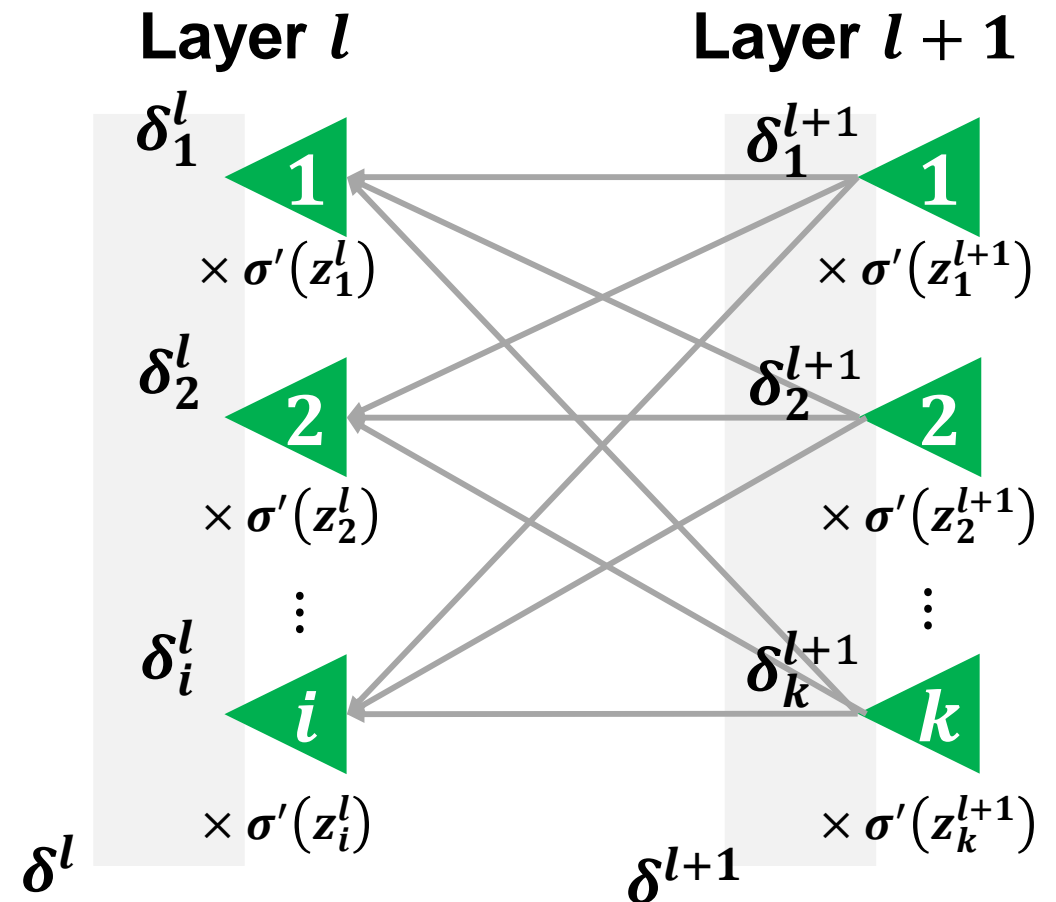$$\boldsymbol{\delta^L = \sigma'(z^L) \odot \nabla C(y)}$$

$$\boldsymbol{\delta^{L-1} = \sigma'(z^{L-1}) \odot (W^L)^T \delta^L}$$

$$\vdots$$

$$\boldsymbol{\delta^l = \sigma'(z^l) \odot (W^{l+1})^T \delta^{l+1}}$$

$$\vdots$$

**Layer $l$**                    **Layer $l+1$**

# Gradient Descent

$$y = f(x) = \sigma(W^L \dots \sigma(W^2 \sigma(W^1 x + b^1) + b^2) \dots + b^L)$$

$$\theta = \{W^1, b^1, W^2, b^2, \dots, W^L, b^L\}$$

$$W^l = \begin{bmatrix} w_{11}^l & w_{12}^l & \cdots \\ w_{21}^l & w_{22}^l & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix}, b^l = \begin{bmatrix} \vdots \\ b_i^l \\ \vdots \end{bmatrix}$$

$$\nabla C(\theta) = \begin{bmatrix} \vdots \\ \dfrac{\partial C(\theta)}{\partial w_{ij}^l} \\ \vdots \\ \dfrac{\partial C(\theta)}{\partial b_i^l} \end{bmatrix}$$
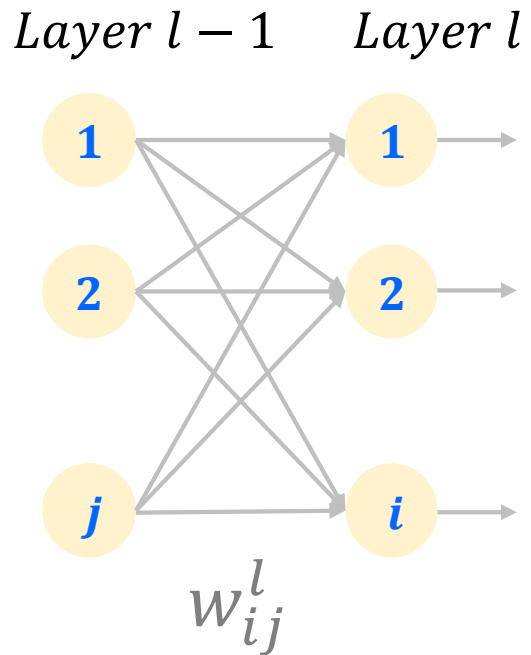
**Algorithm**
Initialization: start at $\theta^0$
while $(\theta^{(i+1)} \neq \theta^i)$
{
    compute the gradients at $\theta^i$
    update parameters
    $\theta^{i+1} \leftarrow \theta^i - \eta \nabla_\theta C(\theta^i)$
}

Gradient descent in neural network requires to update millions of parameters ...
However, it is time-consuming and an inefficient process…
Therefore, we will introduce another solution – **backpropagation**.

# Summary

$$\frac{\partial C(\theta)}{\partial w_{ij}^l} = \frac{\partial C(\theta)}{\partial z_i^l} \frac{\partial z_i^l}{\partial w_{ij}^l}$$

*Layer l − 1*   *Layer l*



$\delta_i^l$

$$\begin{cases} a_j^{l-1} & , l > 1 \\ x_j & , l = 1 \end{cases}$$

$w_{ij}^l$

**Forward Direction**

$$z^1 = W^1 x + b^1 \qquad a^1 = \sigma(z^1)$$
$$\vdots \qquad\qquad \vdots$$
$$z^l = W^l x^{l-1} + b^l \qquad a^l = \sigma(z^l)$$
$$\vdots \qquad\qquad \vdots$$

**Backward Direction**

$$\delta^L = \sigma'(z^L) \odot \nabla C(y)$$
$$\delta^{L-1} = \sigma'(z^{L-1}) \odot (W^L)^T \delta^L$$
$$\vdots$$
$$\delta^l = \sigma'(z^l) \odot (W^{l+1})^T \delta^{l+1}$$
$$\vdots$$

# **References**

- 臺大電機系李宏毅教授講義
- 臺大資工系陳縕儂教授講義
- NVIDIA Deep Learning Tutorial

# Thank you for your attention!