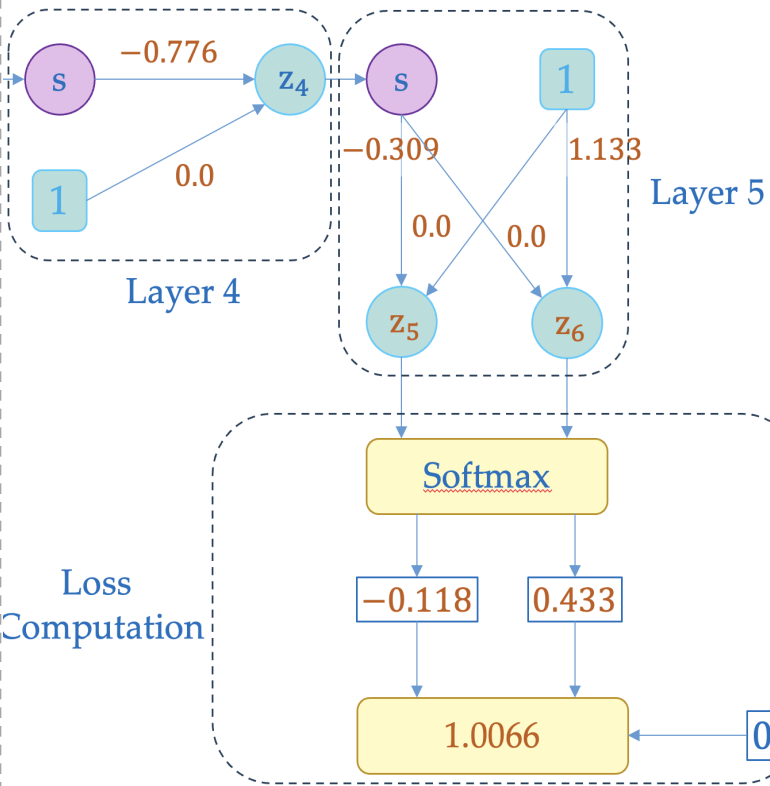# Multi-layer Perception

## Model Initialization

Quang–Vinh Dinh
Ph.D. in Computer Science

Year 2024        code&data

# Objectives



## Case Studies

## Xavier Glorot Init.

## Kaiming He Init.

$$W_i \sim U\left(-\frac{4\sqrt{3}}{\sqrt{n}}, \frac{4\sqrt{3}}{\sqrt{n}}\right)$$

$$W_i \sim U\left(-\frac{\sqrt{6}}{\sqrt{n}}, \frac{\sqrt{6}}{\sqrt{n}}\right)$$
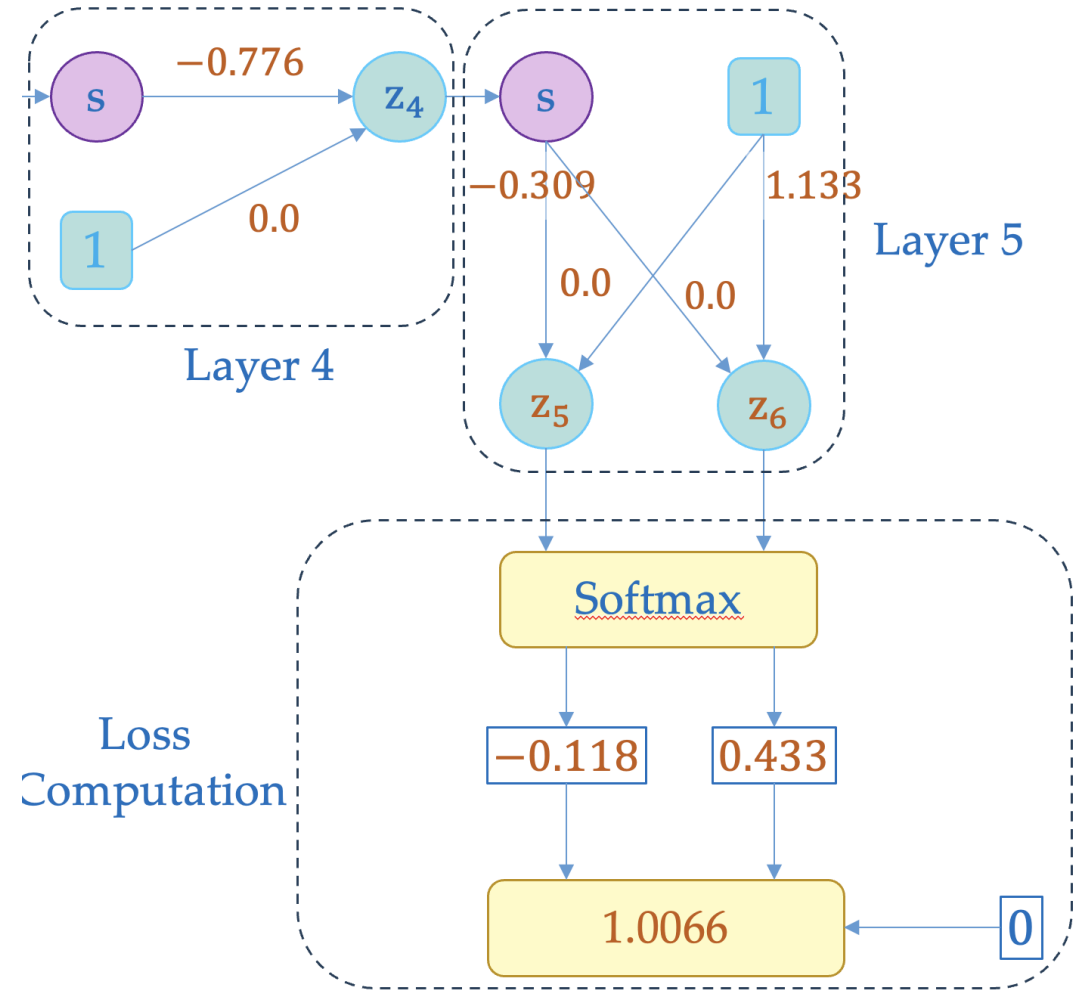
$$W_i \sim N\left(0, \frac{1}{n}\right)$$

$$W_i \sim N\left(0, \frac{2}{n}\right)$$

# Outline

$$X \in [0, 255]$$

$$\text{Normalize}(mean, \text{std})$$

$$\text{Image} = \frac{\text{Image} - mean}{\text{std}}$$

```python
transform = transforms.Compose([transforms.ToTensor(),
                                transforms.Normalize((0,),
                                                     (1.0/255,))])


model = nn.Sequential(
    nn.Flatten(), nn.Linear(784, 256),
    nn.ReLU(), nn.Linear(256, 10)
)
```



28

28

Normalization

flatten

784

Fully connect

Fully connect

$z_1$

$z_{10}$

Softmax activation

Output

784 Nodes + ReLU

256 Nodes + ReLU

10 Nodes Output layer

1

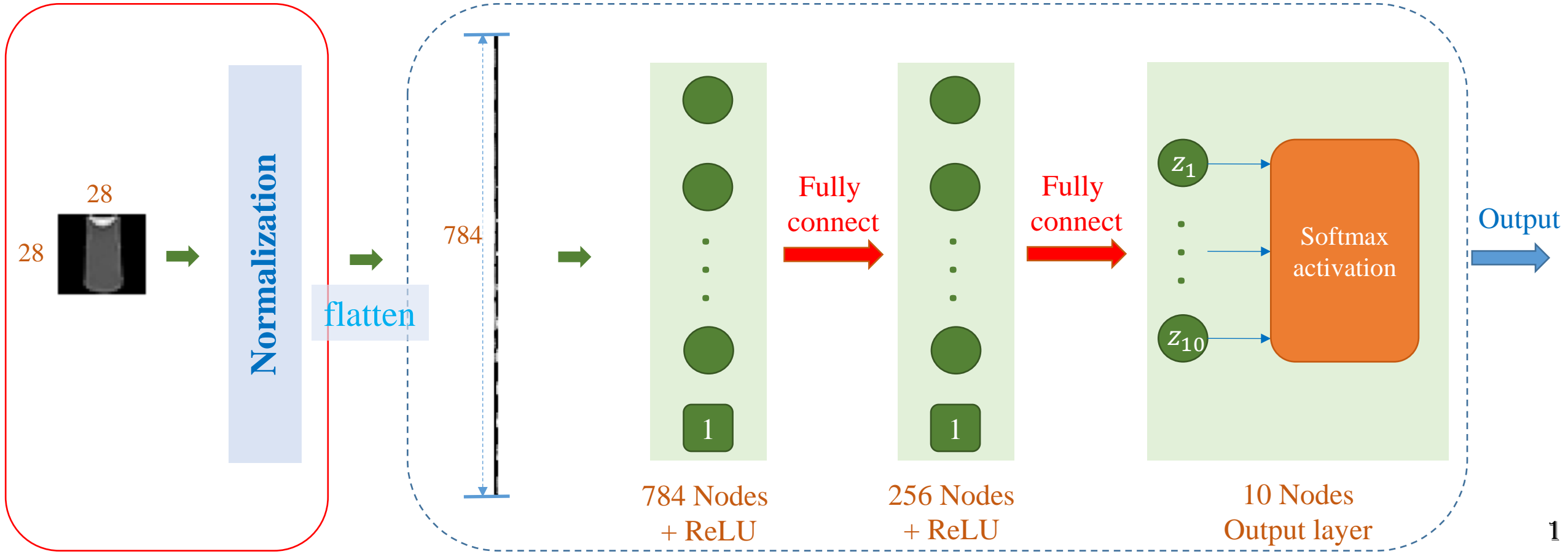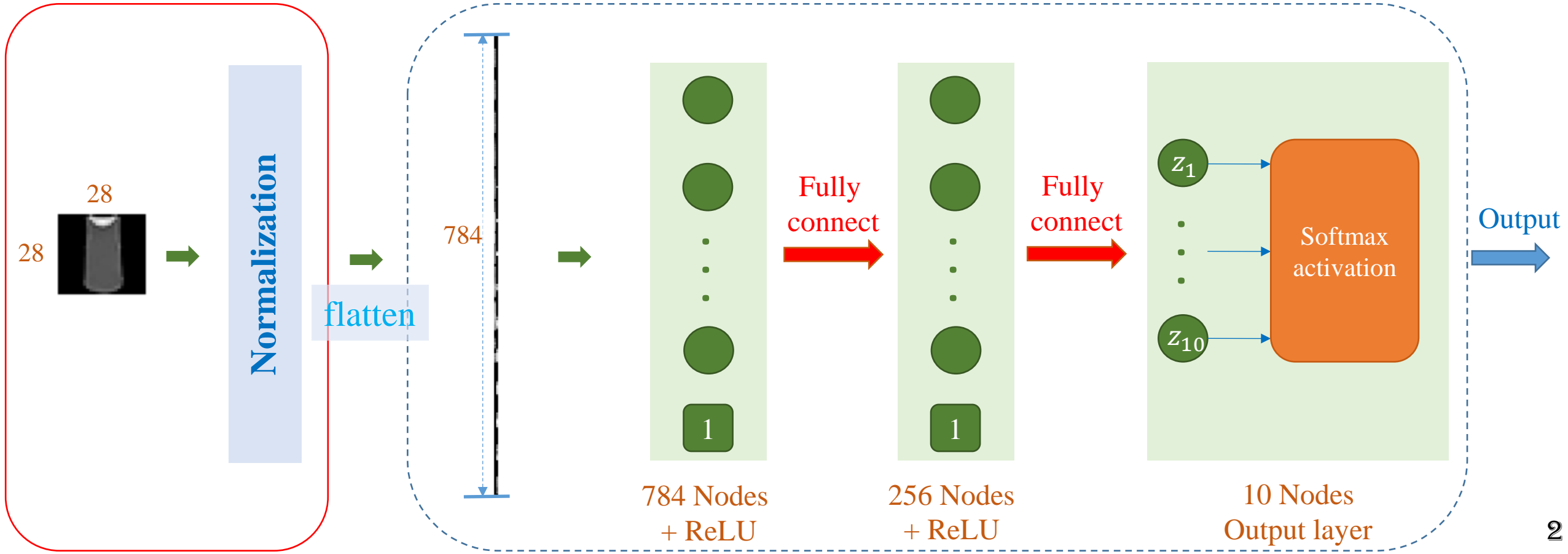$X \in [-1, 1]$

Normalize($mean$, std)

$$\text{Image} = \frac{\text{Image} - mean}{\text{std}}$$

```python
transform = transforms.Compose([transforms.ToTensor(),
                                transforms.Normalize((0.5,),
                                                     (0.5,))])


model = nn.Sequential(
    nn.Flatten(), nn.Linear(784, 256),
    nn.ReLU(), nn.Linear(256, 10)
)
```
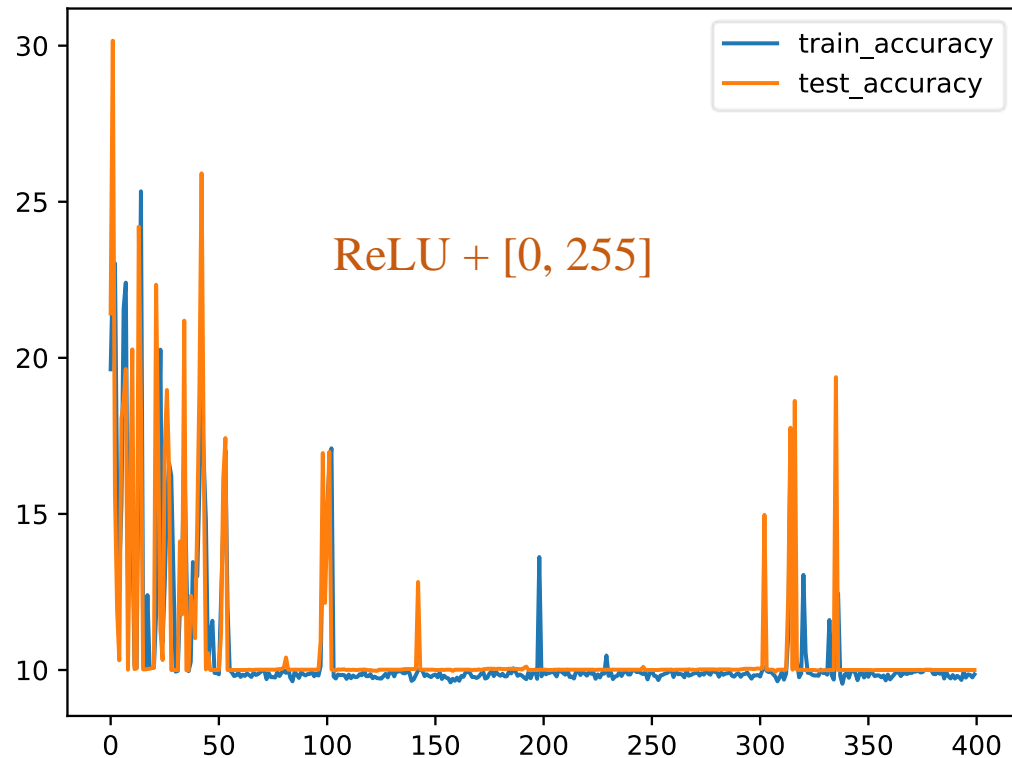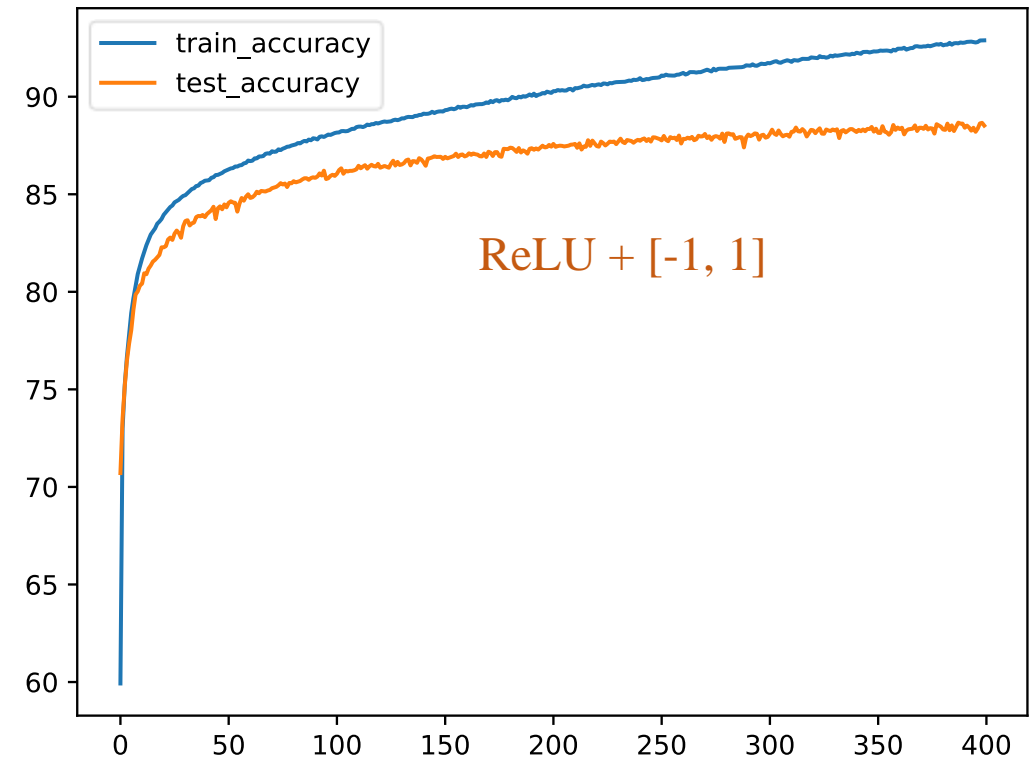
28
28

Normalization

flatten

784

Fully connect

Fully connect

$z_1$
$z_{10}$

Softmax activation

Output

784 Nodes + ReLU

256 Nodes + ReLU

10 Nodes Output layer

2

```
Compose([transforms.ToTensor(),
         transforms.Normalize((0,),
                              (1.0/255,))])
```

```
Compose([transforms.ToTensor(),
         transforms.Normalize((0.5,),
                              (0.5,))])
```
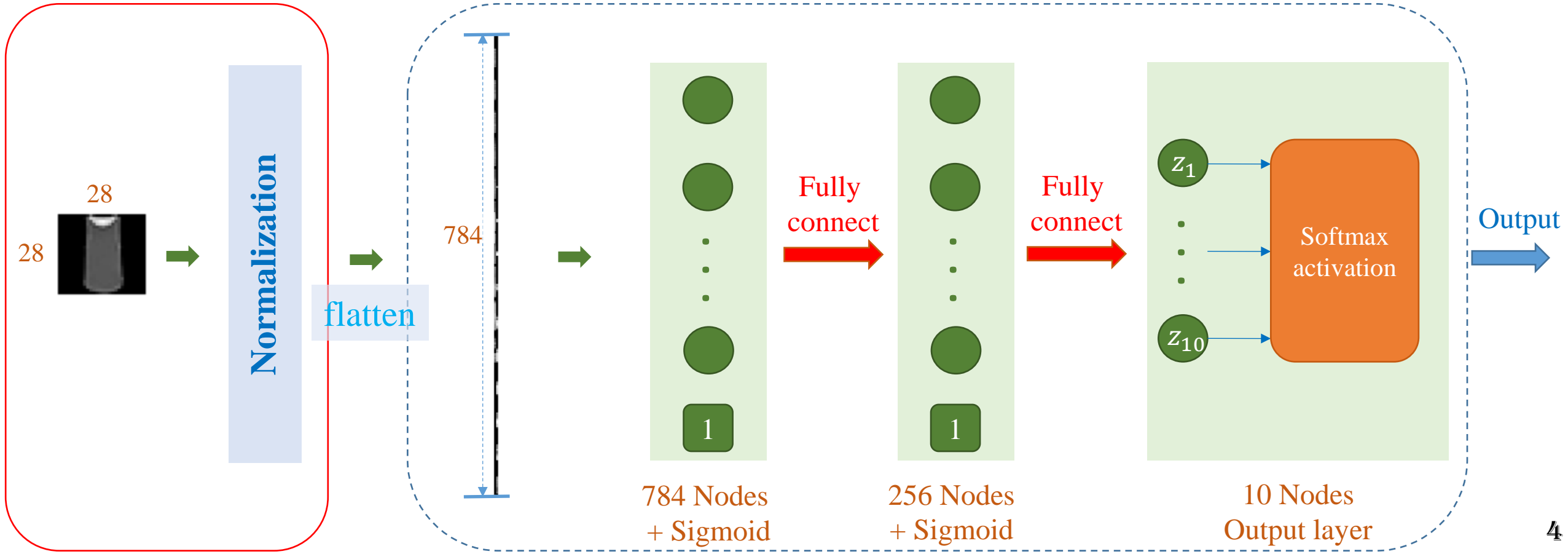
ReLU + [0, 255]

ReLU + [-1, 1]

$$X \in [0, 255]$$

$$\text{Normalize}(mean, \text{std})$$

$$\text{Image} = \frac{\text{Image} - mean}{\text{std}}$$

```python
transform = Compose([ToTensor(),
                     Normalize((0,),
                               (1.0/255,))])


model = nn.Sequential(
    nn.Flatten(), nn.Linear(784, 256),
    nn.Sigmoid(), nn.Linear(256, 10)
)
```



28

28

Normalization

flatten

784

Fully connect

Fully connect

$z_1$

$z_{10}$

Softmax activation

Output

1

1

784 Nodes + Sigmoid

256 Nodes + Sigmoid

10 Nodes Output layer

4

$$X \in [-1, 1]$$

$$\text{Normalize}(mean, \text{std})$$

$$\text{Image} = \frac{\text{Image} - mean}{\text{std}}$$

```python
transform = Compose([ToTensor(),
                     Normalize((0.5,),
                               (0.5,))])


model = nn.Sequential(
    nn.Flatten(), nn.Linear(784, 256),
    nn.Sigmoid(), nn.Linear(256, 10)
)
```
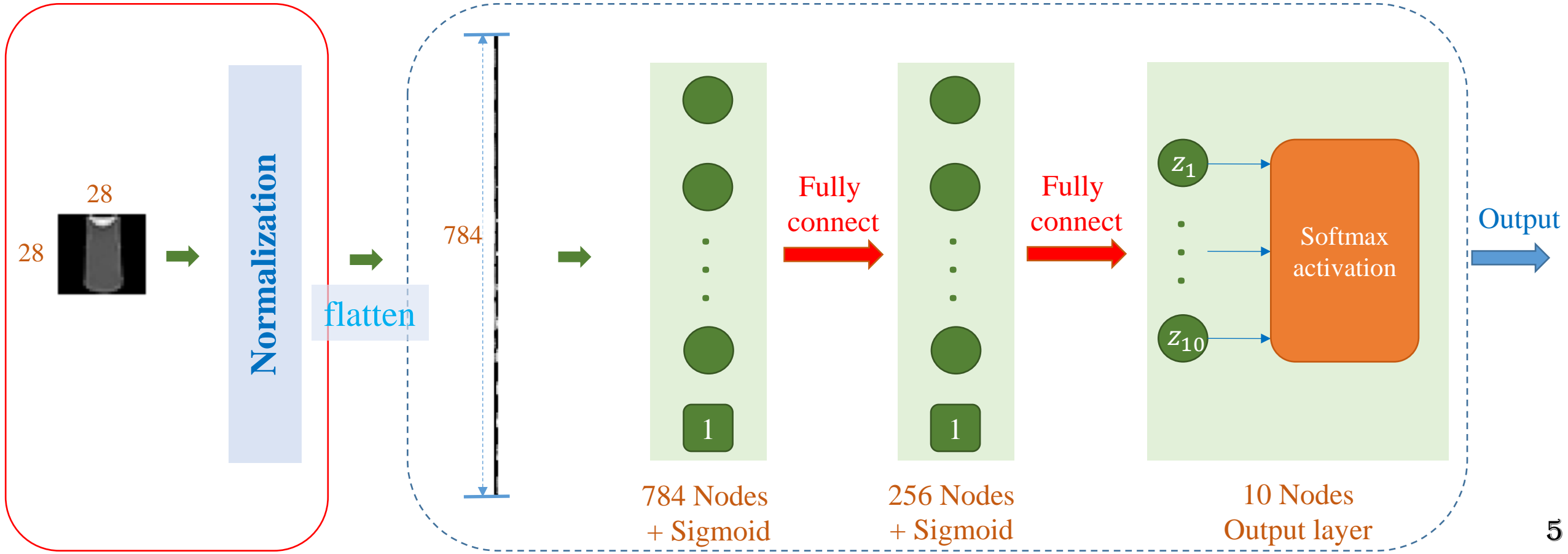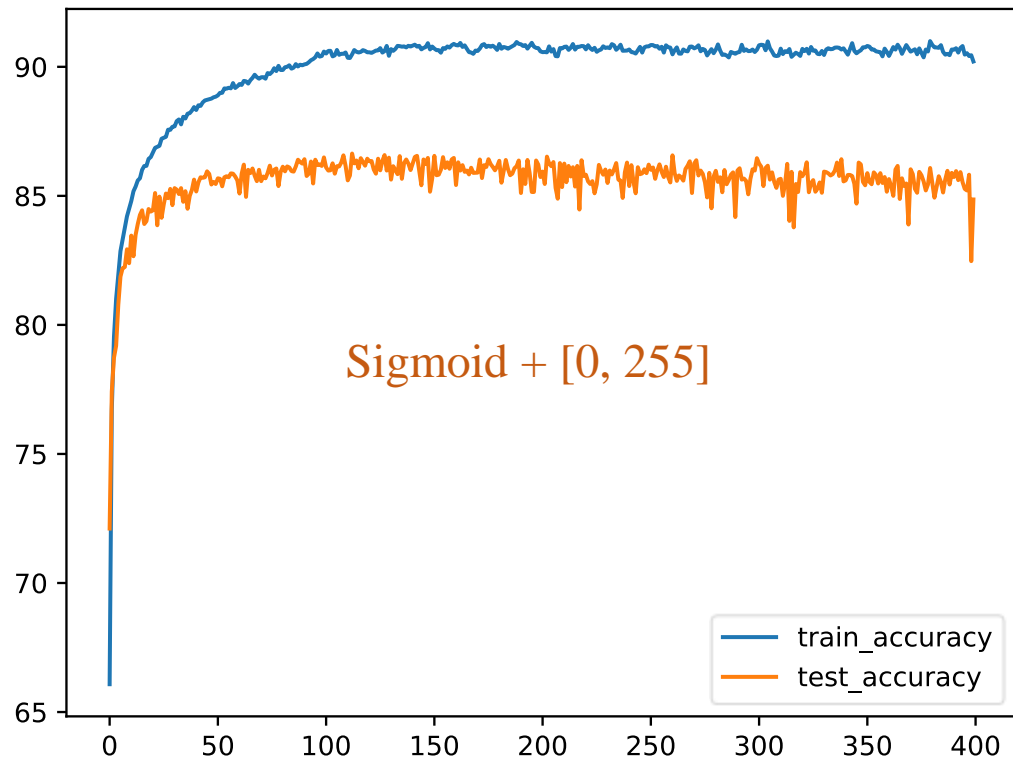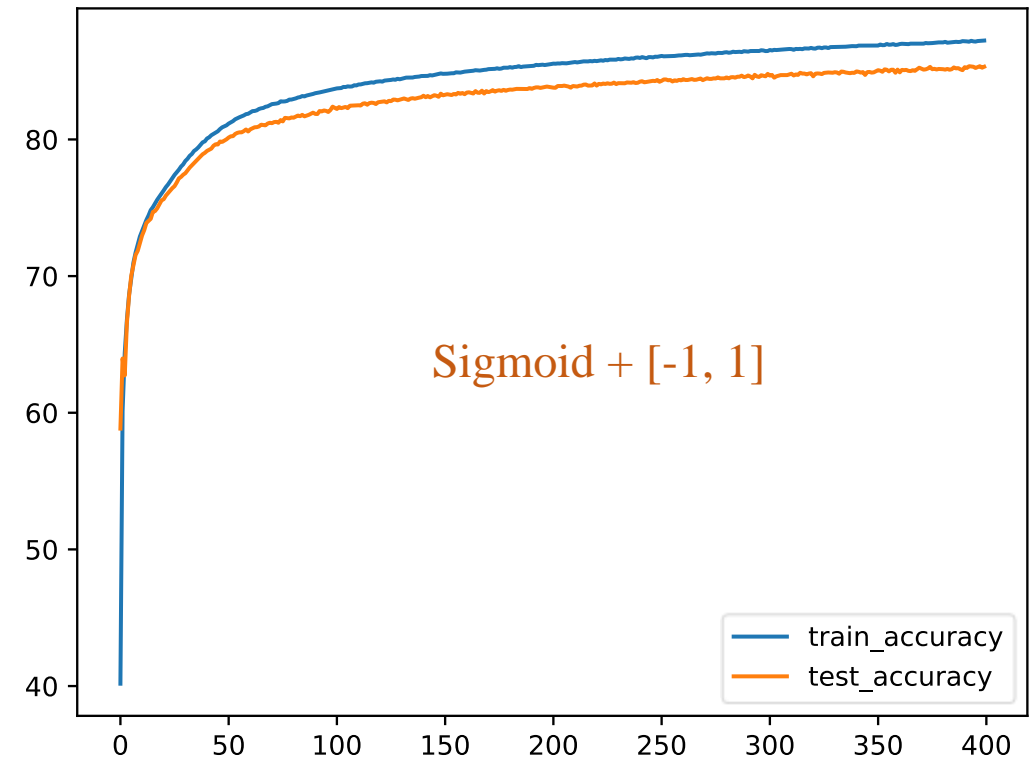
28
28

**Normalization**

flatten

784

Fully connect

Fully connect

$z_1$

$z_{10}$

Softmax activation

Output

1

1

784 Nodes + Sigmoid

256 Nodes + Sigmoid

10 Nodes Output layer

5

```
Compose([transforms.ToTensor(),
         transforms.Normalize((0,),
                              (1.0/255,))])
```

Sigmoid + [0, 255]

```
Compose([transforms.ToTensor(),
         transforms.Normalize((0.5,),
                              (0.5,))])
```

Sigmoid + [-1, 1]

**AI VIET NAM**
@aivietnam.edu.vn

Large weight initialization



X

$w_1$   $z_1$   s   $w_2$   $z_2$

$b_1$   1   $w_3$   $b_2$   $z_3$

1   $b_3$

Softmax

$\hat{y}_0$

$\hat{y}_1$

Cross Entropy

Layer 1   Layer 2

s   Sigmoid function

Problem???

# Activation Functions

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

data = 

| 1 | 5 | -4 | 3 | -2 |
|---|---|----|---|----|

data_a = **sigmoid(data)**

data_a = 

| 0.731 | 0.993 | 0.017 | 0.95 | 0.119 |
|-------|-------|-------|------|-------|



— sigmoid
— derivative

$$
\begin{aligned}
\text{sigmoid}'(x) &= \left(\frac{1}{1 + e^{-x}}\right)' = \frac{-1}{(1 + e^{-x})^2}(-e^{-x}) \\[2mm]
&= \frac{e^{-x}}{(1 + e^{-x})^2} = \frac{e^{-x} + 1 - 1}{(1 + e^{-x})^2} \\[2mm]
&= \frac{1}{1 + e^{-x}} - \frac{1}{(1 + e^{-x})^2} \\[2mm]
&= \frac{1}{1 + e^{-x}}\left(1 - \frac{1}{1 + e^{-x}}\right) \\[2mm]
&= \text{sigmoid}(x)\,(1 - \text{sigmoid}(x))
\end{aligned}
$$

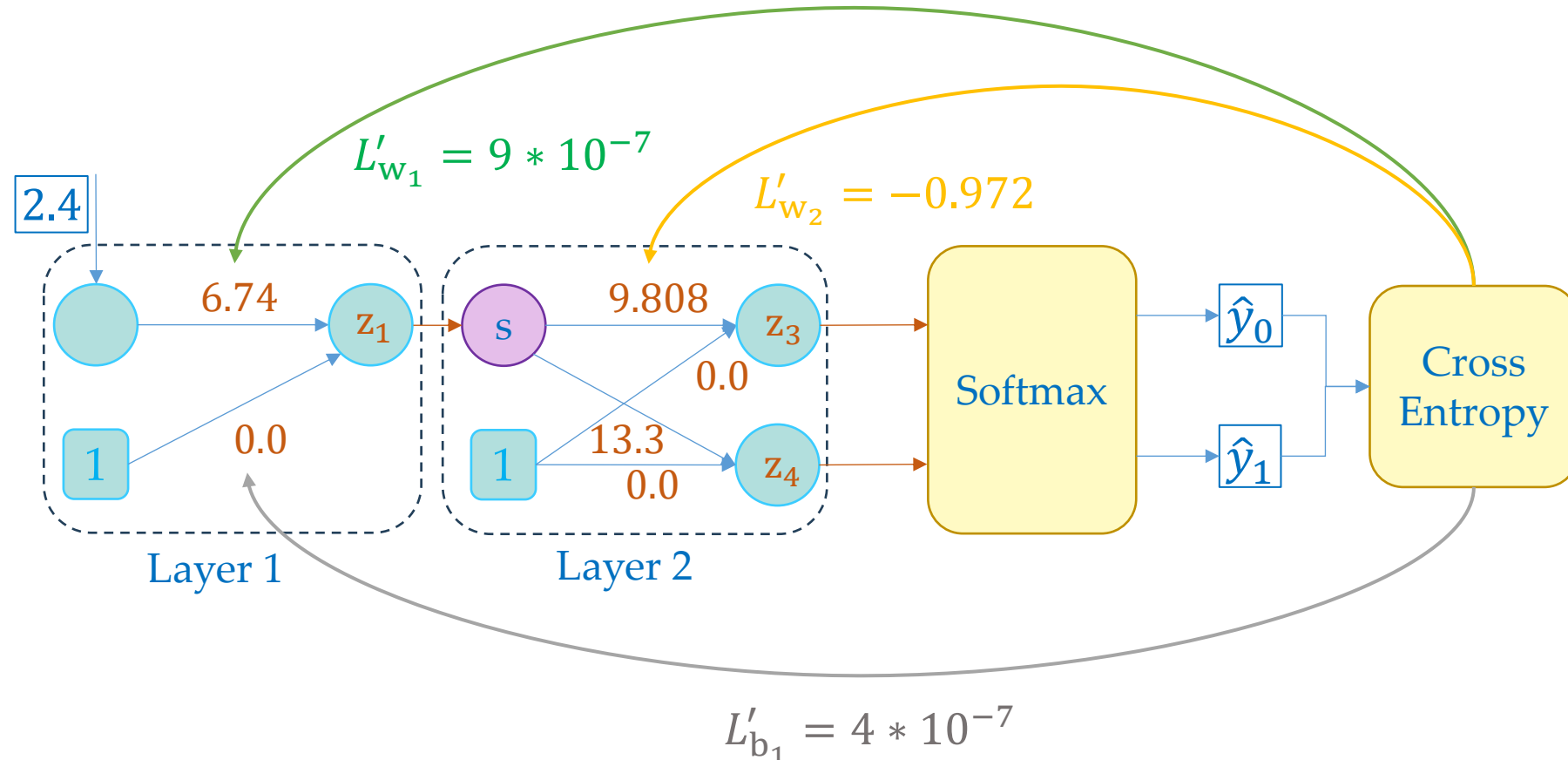Large weight initialization

```python
linear1 = nn.Linear(1, 1)
linear2 = nn.Linear(1, 2)

init.normal_(linear1.weight,
             mean=0, std=10)
init.normal_(linear2.weight,
             mean=0, std=10)
```

$$L'_{w_1} = 9 * 10^{-7}$$
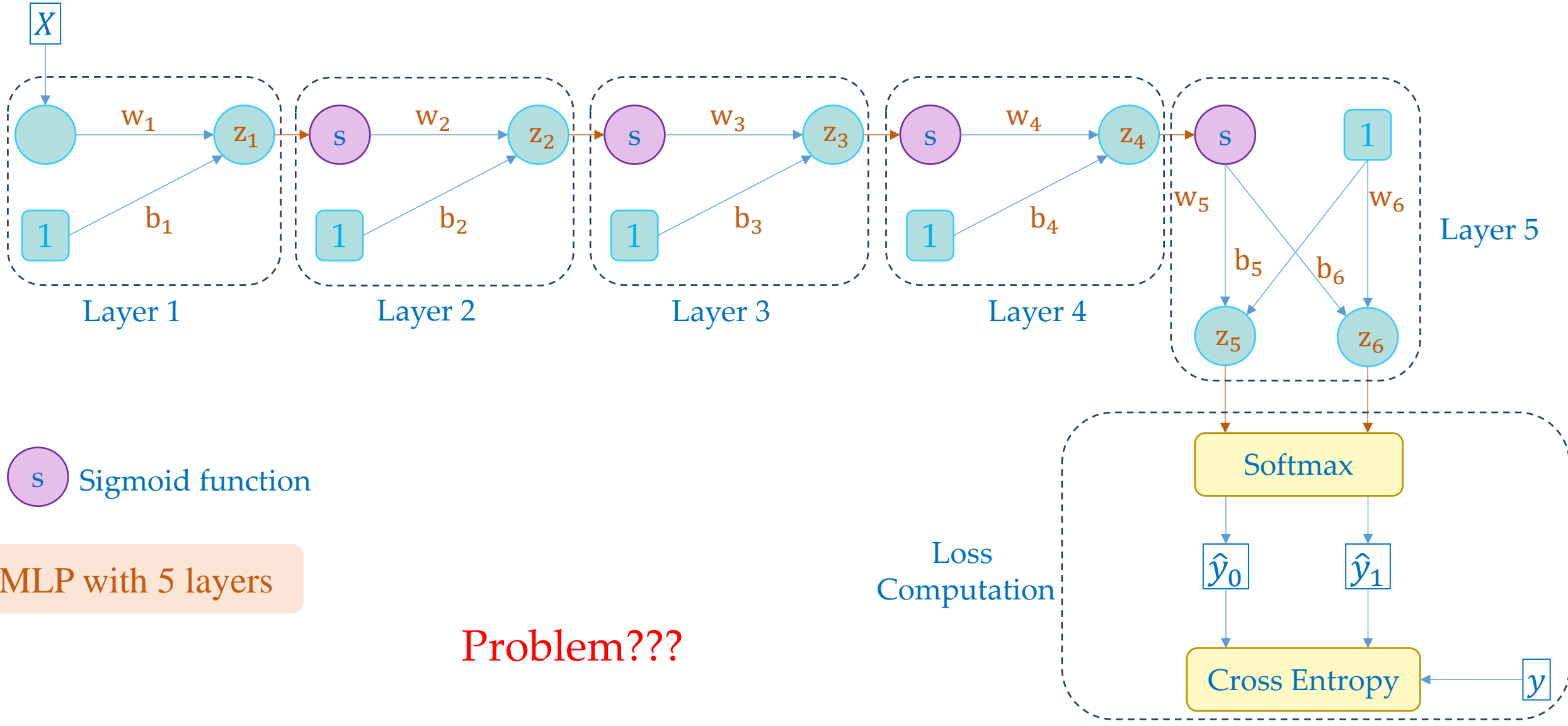
$$L'_{w_2} = -0.972$$

2.4

6.74

$z_1$

s

9.808

$z_3$

0.0

1

0.0

1

13.3

$z_4$

0.0

Softmax

$\hat{y}_0$

$\hat{y}_1$

Cross Entropy

Layer 1

Layer 2

with $\eta = 0.01$

$$\eta L'_{w_1} = 9 * 10^{-9}$$

$$\eta L'_{b_1} = 4 * 10^{-9}$$

$$L'_{b_1} = 4 * 10^{-7}$$

Derivative values are too small

8

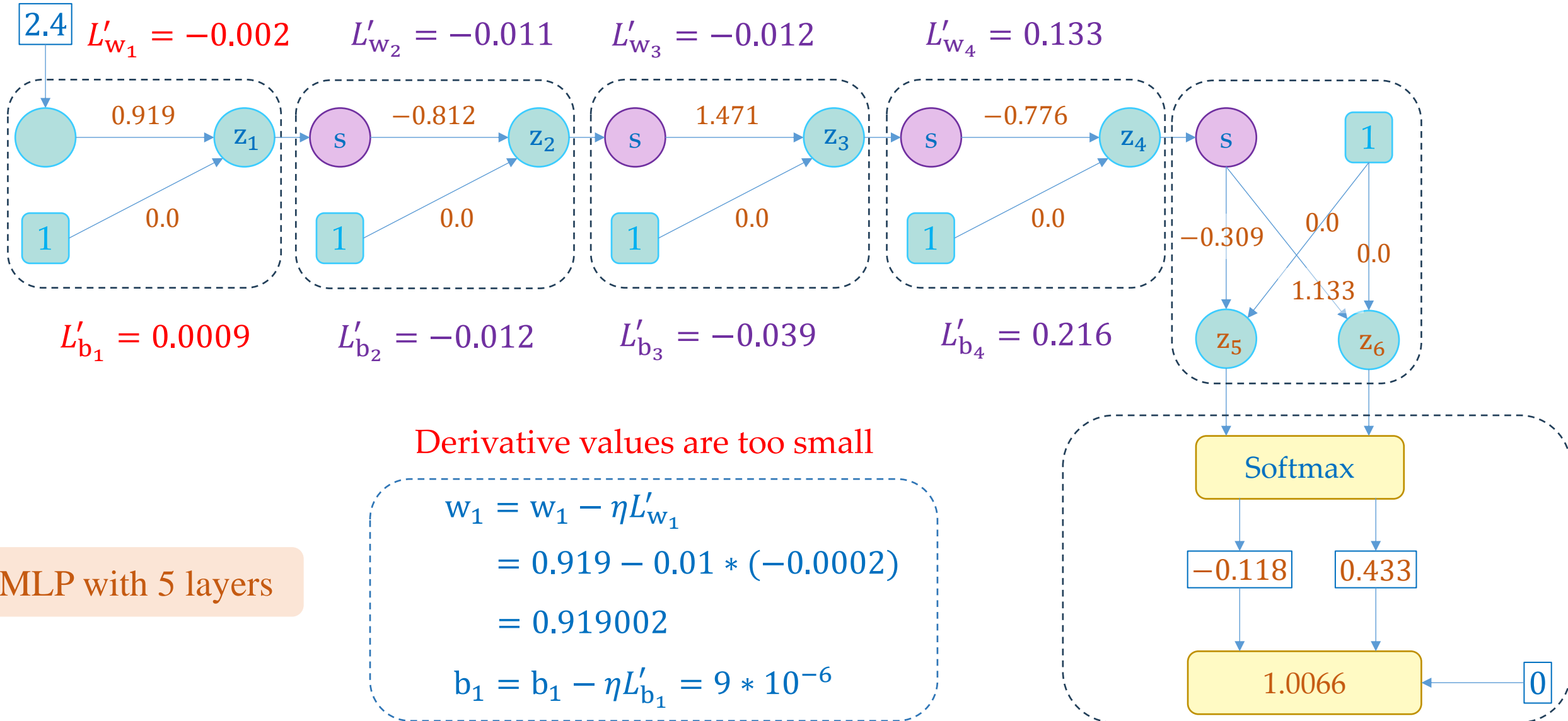# Gradient Vanishing

Using appropriate weight initialization



$X$

Layer 1

Layer 2

Layer 3

Layer 4

Layer 5

$w_1$ $z_1$ $b_1$ $w_2$ $s$ $z_2$ $b_2$ $w_3$ $s$ $z_3$ $b_3$ $w_4$ $s$ $z_4$ $b_4$ $s$ $w_5$ $w_6$ $b_5$ $b_6$ $z_5$ $z_6$

$s$ Sigmoid function

MLP with 5 layers

Problem???

Loss Computation

Softmax

$\hat{y}_0$ $\hat{y}_1$

Cross Entropy $\leftarrow$ $y$

# Gradient Vanishing

Using appropriate weight initialization



**Layer 1**

**Layer 2**

**Layer 3**

**Layer 4**

**Layer 5**

$s$ Sigmoid function

MLP with 5 layers

Loss Computation

Softmax

$\hat{y}_0$   $\hat{y}_1$

1.0066

# Gradient Vanishing

## Using appropriate weight initialization

$2.4$

$L'_{w_1} = -0.002$     $L'_{w_2} = -0.011$     $L'_{w_3} = -0.012$     $L'_{w_4} = 0.133$



$0.919$   $z_1$    s   $-0.812$   $z_2$    s   $1.471$   $z_3$    s   $-0.776$   $z_4$    s   $1$

$1$   $0.0$    $1$   $0.0$    $1$   $0.0$    $1$   $0.0$

$-0.309$   $0.0$   $0.0$   $1.133$

$z_5$    $z_6$

$L'_{b_1} = 0.0009$     $L'_{b_2} = -0.012$     $L'_{b_3} = -0.039$     $L'_{b_4} = 0.216$

### Derivative values are too small

$$w_1 = w_1 - \eta L'_{w_1}$$
$$= 0.919 - 0.01 * (-0.0002)$$
$$= 0.919002$$
$$b_1 = b_1 - \eta L'_{b_1} = 9 * 10^{-6}$$

MLP with 5 layers

Softmax

$-0.118$    $0.433$

$1.0066$    $0$

# Gradient Vanishing

## MLP with 8 layers

$X$

$w_1$ ... $z_1$ ... $b_1$ ... Layer 1

$s$ ... $w_2$ ... $z_2$ ... $b_2$ ... Layer 2

$s$ ... $w_3$ ... $z_3$ ... $b_3$ ... Layer 3

$s$ ... $1$ ... $w_5$ ... $w_6$ ... $b_5$ ... $b_6$ ... $z_5$ ... $z_6$ ... Layer 8

$s$ — Sigmoid function

Softmax

$\hat{y}_0$ ... $\hat{y}_1$

Loss Computation

Cross Entropy ← $y$

**❖ PReLU function**

$$PReLU(x) = \begin{cases} \alpha x & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$$



prelu
prelu_derivative

data =

| 1 | 5 | -4 | 3 | -2 |
|---|---|----|---|----|

data_a = **PRELU**(data)

data_a =

| 1 | 5 | -0.4 | 3 | -0.2 |
|---|---|------|---|------|

$$PReLU'(x) = \begin{cases} \alpha & \text{if } x \leq 0 \\ 1 & \text{if } x > 0 \end{cases}$$

# Outline

$$W_i \sim U\left(-\frac{4\sqrt{3}}{\sqrt{n}}, \frac{4\sqrt{3}}{\sqrt{n}}\right)$$

$$W_i \sim N\left(0, \frac{1}{n}\right)$$

**Data**

$$X = \{X_1, \ldots, X_N\}$$

**Formula**

$$E(X) = \sum_{i=1}^{N} X_i P_X(X_i)$$

**Given the data**

$$X = \{2, 8, 5, 4, 1, 4\}$$

$$N = 6$$

$$P_X(X = 2) = \frac{1}{6} \qquad P_X(X = 4) = \frac{2}{6}$$

$$P_X(X = 8) = \frac{1}{6} \qquad P_X(X = 1) = \frac{1}{6}$$

$$P_X(X = 5) = \frac{1}{6}$$

$$E(X) = 2 \times \frac{1}{6} + 8 \times \frac{1}{6} + 5 \times \frac{1}{6} + 4 \times \frac{2}{6} + 1 \times \frac{1}{6}$$

$$= \frac{2}{6} + \frac{8}{6} + \frac{5}{6} + \frac{8}{6} + \frac{1}{6} = 4$$

15

**Data**

$$X = \{X_1, \dots, X_N\}$$

**Formula**

$$E(X) = \sum_{i=1}^{N} X_i P_X(X_i)$$

$$E(XY) = \sum_{i=1}^{N} \sum_{j=1}^{N} X_i Y_j P(X_i, Y_j)$$

$$= \sum_{i=1}^{N} \sum_{j=1}^{N} X_i Y_j P(X_i) P(Y_j)$$

$$= \sum_{i=1}^{N} X_i P(X_i) \sum_{j=1}^{N} Y_j P(Y_j)$$

$$= E(X) E(Y)$$

# Variance

### Formula

**mean**

$$E(X) = \sum_{i=1}^{N} X_i P_X(X_i)$$

**variance**

$$var(X) = E\left(\left(X - E(X)\right)^2\right)$$

$$= \sum_{i=1}^{N} (X_i - E(X))^2 P_X(X_i)$$

**Standard deviation**

$$\sigma = \sqrt{var(X)}$$

### Example:   $X = \{5, 3\ 6, 7, 4\}$

$$E(X) = 5 \times \frac{1}{5} + 3 \times \frac{1}{5} + 6 \times \frac{1}{5} + 7 \times \frac{1}{5} + 4 \times \frac{1}{5}$$

$$= 5$$

$$var(X) = \frac{1}{5}[(5-5)^2 + (3-5)^2 + (6-5)^2 +$$
$$(7-5)^2 + (4-5)^2]$$

$$= \frac{1}{5}(0+4+1+4+1)=2$$

$$\sigma = \sqrt{var(X)} = 1.41$$

17

# Variance

## Formula

**mean**

$$E(X) = \sum_{i=1}^{N} X_i P_X(X_i)$$

**variance**

$$var(X) = E\left(\left(X - E(X)\right)^2\right)$$

$$= \sum_{i=1}^{N} \left(X_i - E(X)\right)^2 P_X(X_i)$$

**Standard deviation**

$$\sigma = \sqrt{var(X)}$$

$$var(X) = \sum_{i=1}^{N} \left(X_i - E(X)\right)^2 P_X(X_i)$$

$$= \sum_{i=1}^{N} \left(X_i^2 - 2X_i E(X) + E(X)^2\right) P_X(X_i)$$

$$= \sum_{i=1}^{N} X_i^2 P_X(X_i) - \sum_{i=1}^{N} 2X_i E(X) P_X(X_i)$$

$$+ \sum_{i=1}^{N} E(X)^2 P_X(X_i)$$

$$= E(X^2) - 2E(X)\left[\sum_{i=1}^{N} X_i P_X(X_i)\right] + E(X)^2$$

$$= E(X^2) - \left(E(X)\right)^2$$

$$var(X) = E(X^2) - \left(E(X)\right)^2$$

$$var(XY) = E(X^2Y^2) - \left(E(XY)\right)^2$$

$$= E(X^2)E(Y^2) - \left(E(X)E(Y)\right)^2$$

$$= \left[var(X) + \left(E(X)\right)^2\right]\left[var(Y) + \left(E(Y)\right)^2\right] - \left(E(X)E(Y)\right)^2$$

$$= var(X)var(Y) + var(X)\left(E(Y)\right)^2 + var(Y)\left(E(X)\right)^2$$

❖ **Xavier Glorot Initialization**

Uniform Distribution

$X \sim U(a, b)$    $E[X] = \dfrac{a + b}{2}$

$f(x) = \dfrac{1}{b - a}$    $var[X] = \dfrac{(b - a)^2}{12}$

## Uniform Distribution

$$X \sim U(a, b) \qquad E[X] = \frac{a+b}{2}$$

$$f(x) = \frac{1}{b-a} \qquad var[X] = \frac{(b-a)^2}{12}$$



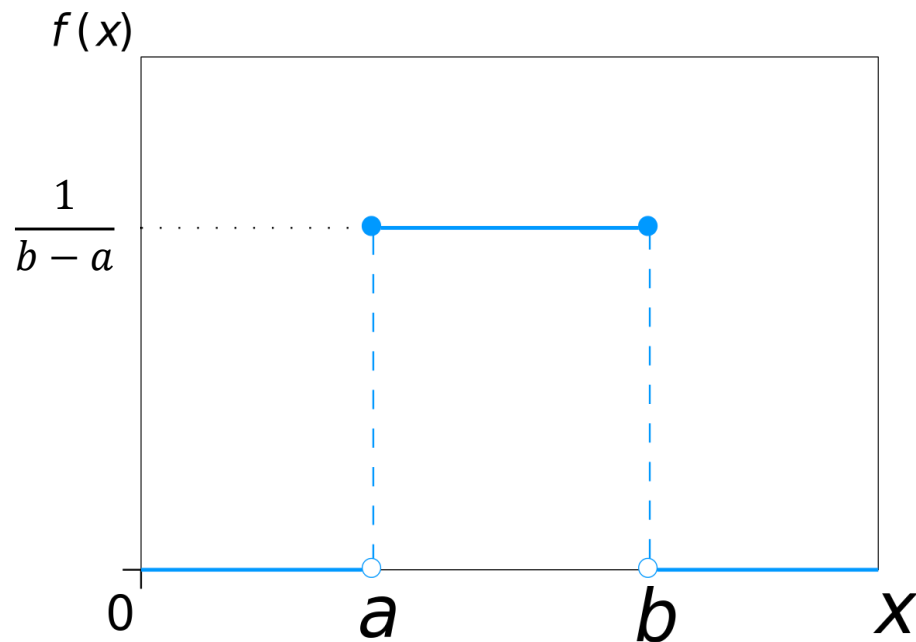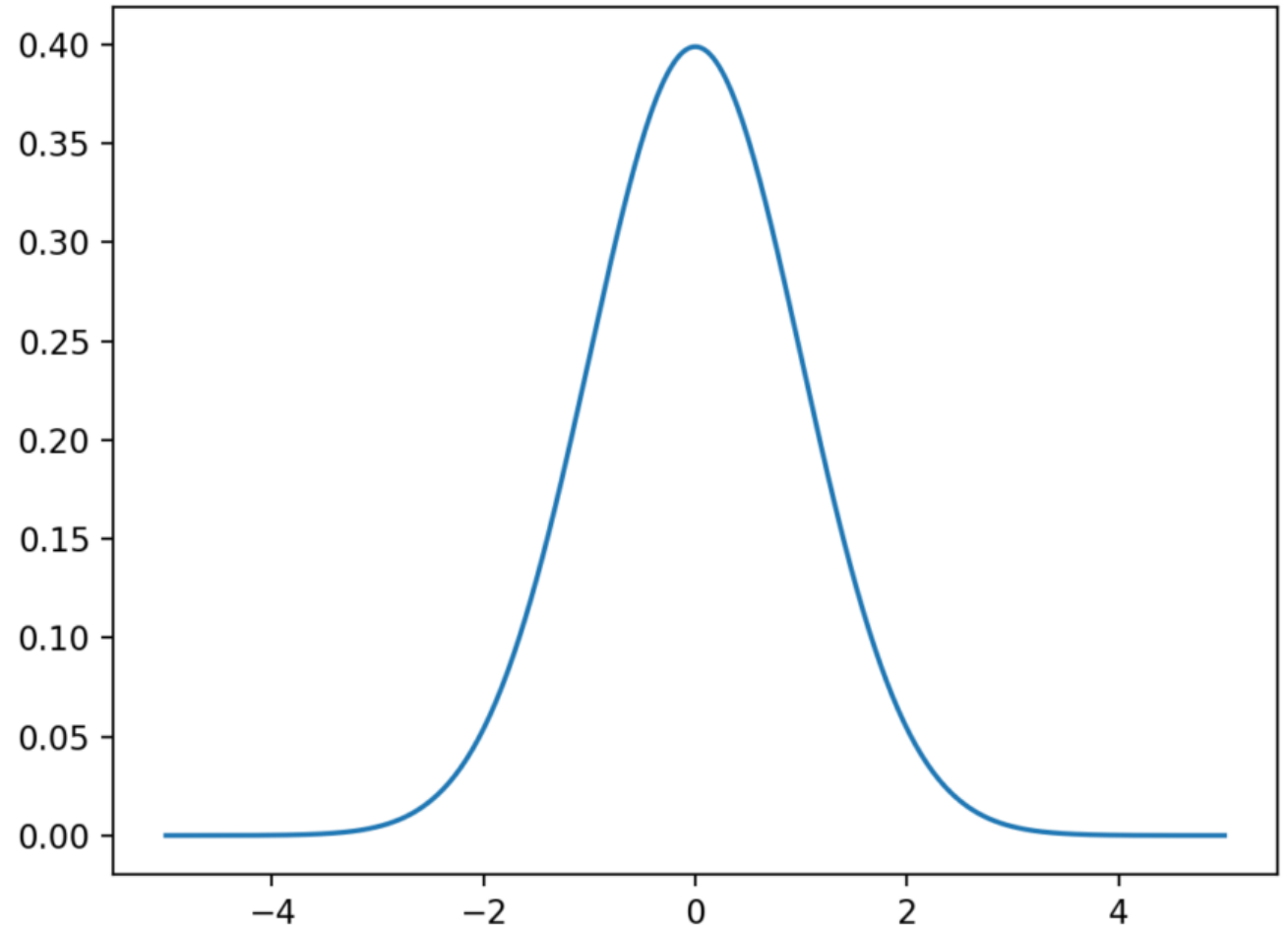$$E[X] = \int_{-\infty}^{\infty} x f(x) dx = \int_{a}^{b} x \frac{1}{b-a} dx$$

$$= \frac{x^2}{2(b-a)} \Big|_{a}^{b} = \frac{b^2 - a^2}{2(b-a)} = \frac{a+b}{2}$$

## Uniform Distribution

$$X \sim U(a, b) \qquad E[X] = \frac{a+b}{2}$$

$$f(x) = \frac{1}{b-a} \qquad var[X] = \frac{(b-a)^2}{12}$$



$$var[X] = E\left((X - E(X))^2\right) = \int_{-\infty}^{\infty} (x - E(X))^2 f(x) dx$$

$$= \int_a^b \left(x - \frac{a+b}{2}\right)^2 \frac{1}{b-a} dx$$

$$= \frac{1}{b-a}\left[\int_a^b x^2 dx - \int_a^b 2x \frac{a+b}{2} dx + \int_a^b \left(\frac{a+b}{2}\right)^2 dx\right]$$

$$= \frac{1}{b-a}\left[\frac{x^3}{3}\Big|_a^b - \frac{x^2(a+b)}{2}\Big|_a^b + \left(\frac{a+b}{2}\right)^2 x\Big|_a^b\right]$$

$$= \frac{1}{b-a}\left[\frac{b^3 - a^3}{3} - \frac{(b^2 - a^2)(a+b)}{2} + \left(\frac{a+b}{2}\right)^2 (b-a)\right]$$

$$= \frac{a^2 + ab + b^2}{3} - \frac{a^2 + 2ab + b^2}{2} + \frac{a^2 + 2ab + b^2}{4}$$

$$= \frac{4(a^2 + ab + b^2) - 3(a^2 + 2ab + b^2)}{12} = \frac{(b-a)^2}{12}$$

❖ **Xavier Initialization**
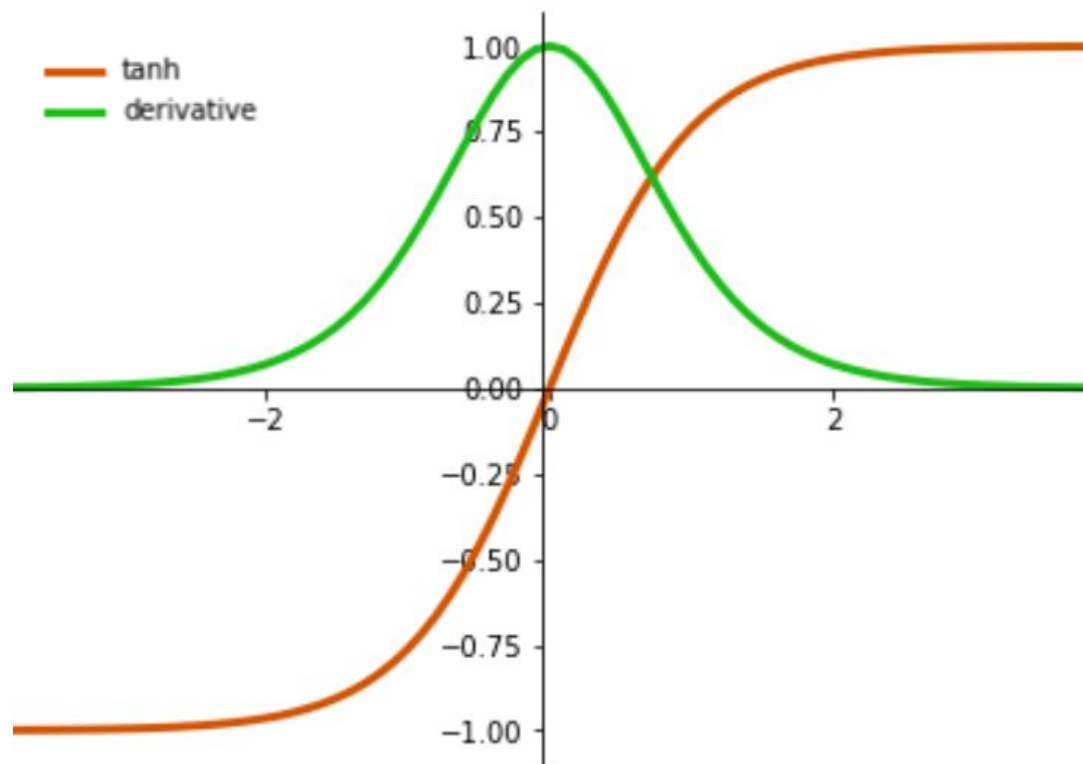
Gaussian Distribution

$$X \sim N(\mu, \sigma^2)$$

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

# Maclaurin series

Tính giá trị xấp xỉ hàm f(x) cho những giá trị $x \approx 0$

$$f(x) = \sum_{n=0}^{\infty} f^{(n)}(0) \frac{x^n}{n!}$$

$$= f(0) + f'(0)x + \frac{f''(0)}{2!}x^2 + \frac{f^{(3)}(0)}{3!}x^3 + \cdots$$



$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = 1 - \frac{2}{e^{2x} + 1} = \frac{2}{e^{-2x} + 1} - 1$$

$$\tanh(0) = 0$$

$$\tanh'(0) = 1 - tanh^2(0) = 1$$

$$\tanh''(0) = \left(1 - tanh^2(0)\right)'$$

$$= -2tanh(0)\tanh'(0) = 0$$

$$\tanh^{(3)}(0) = \left(-2tanh(0)\tanh'(0)\right)'$$

$$= -2[\tanh'(0)\tanh'(0) + \tanh''(0)tanh(0)] = -2$$
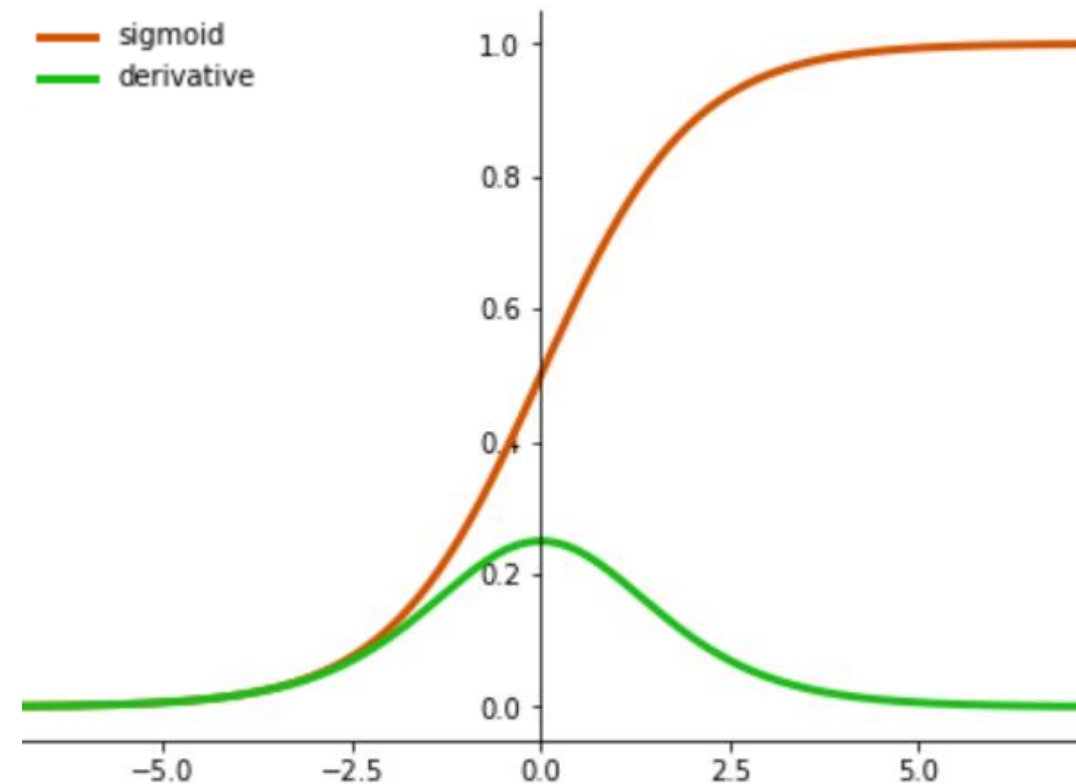
$$\tanh(x) = f(0) + f'(0)x + \frac{f''(0)}{2!}x^2 + \frac{f^{(3)}(0)}{3!}x^3 + \cdots$$

$$= x - \frac{2x^3}{3!} + \cdots$$

$$\Longrightarrow \quad \tanh(x) \approx x$$

# Maclaurin series

Tính giá trị xấp xỉ hàm f(x) cho những giá trị $x \approx 0$

$$f(x) = \sum_{n=0}^{\infty} f^{(n)}(0) \frac{x^n}{n!}$$

$$= f(0) + f'(0)x + \frac{f''(0)}{2!}x^2 + \frac{f^{(3)}(0)}{3!}x^3 + \cdots$$



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

$$\text{sigmoid}(0) = \frac{1}{2}$$

$$\text{sigmoid}'(0) = \text{sigmoid}(0)\,(1 - \text{sigmoid}(0)) = \frac{1}{4}$$

$$\text{sigmoid}''(0) = [\text{sigmoid}(0)\,(1 - \text{sigmoid}(0))]'$$

$$= \text{sigmoid}'(0) - 2\,\text{sigmoid}(0)\text{sigmoid}'(0) = 0$$

$$\text{sigmoid}(x) = f(0) + f'(0)x + \frac{f''(0)}{2!}x^2 + \frac{f^{(3)}(0)}{3!}x^3 + \cdots$$

$$= \frac{1}{2} + \frac{x}{4} + \cdots$$

$$\Longrightarrow \quad \text{sigmoid}(x) \approx \frac{1}{2} + \frac{x}{4}$$
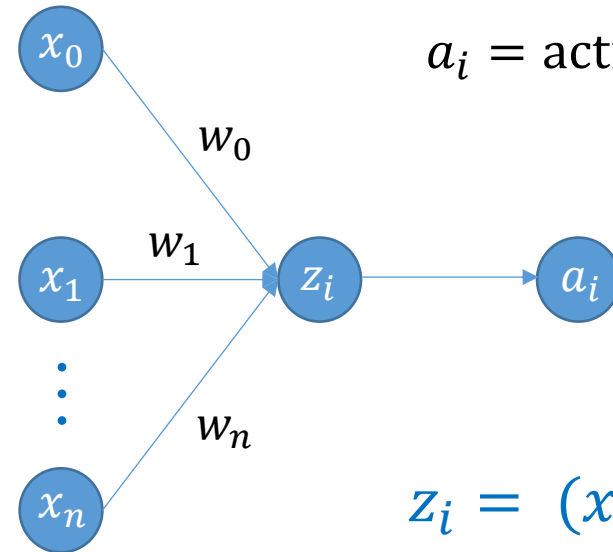
## ❖ Xavier Initialization

$$E(XY) = E(X)E(Y)$$

$$var(XY) = var(X)var(Y) +$$
$$var(X)\big(E(Y)\big)^2 +$$
$$var(y)\big(E(X)\big)^2$$

Uniform Distribution

$$X \sim U(a, b)$$

$$f(x) = \frac{1}{b-a}$$

$$var[X] = \frac{(b-a)^2}{12}$$

$a_i = \text{activation}(z_i)$

$E(X) = 0$

$E(W) = 0$

$b = 0$

$$z_i = (x_1 w_1 + \cdots + x_n w_n + b)$$

$$\text{var}(z_i) = \text{var}(x_1 w_1 + \cdots + x_n w_n + b)$$
$$= n\text{var}(x_i w_i) = n\text{var}(x_i)\text{var}(w_i)$$

$\text{activation} = \tanh \Rightarrow a_i = \tanh(z_i) \approx z_i \Rightarrow \text{var}(a_i) = \text{var}(z_i)$

$\text{var}(X) = \text{var}(\mathbf{a}) \xrightarrow{\text{iid}} \text{var}(x_i) = \text{var}(a_i) \Rightarrow n\text{var}(w_i) = 1$

$$\Rightarrow \text{var}(w_i) = \frac{1}{n}$$

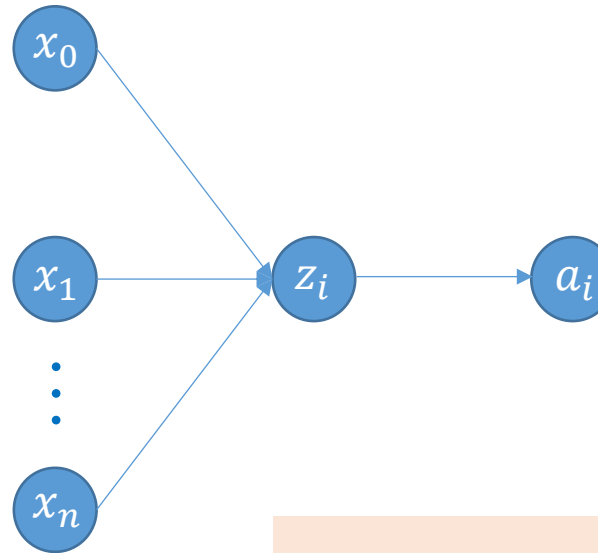❖ **Xavier Initialization**

activation = tanh

$$E(XY) = E(X)E(Y)$$

$$var(XY) = var(X)var(Y) +$$
$$var(X)\big(E(Y)\big)^2 +$$
$$var(y)\big(E(X)\big)^2$$

Uniform Distribution

$$X \sim U(a, b)$$

$$f(x) = \frac{1}{b-a}$$

$$var[X] = \frac{(b-a)^2}{12}$$

$$var(w_i) \approx \frac{1}{n}$$

$$w_i \sim U(-r, r)$$

$$var[w_i] = \frac{r^2}{3}$$

$$W_i \sim U\left(-\frac{\sqrt{3}}{\sqrt{n}}, \frac{\sqrt{3}}{\sqrt{n}}\right)$$
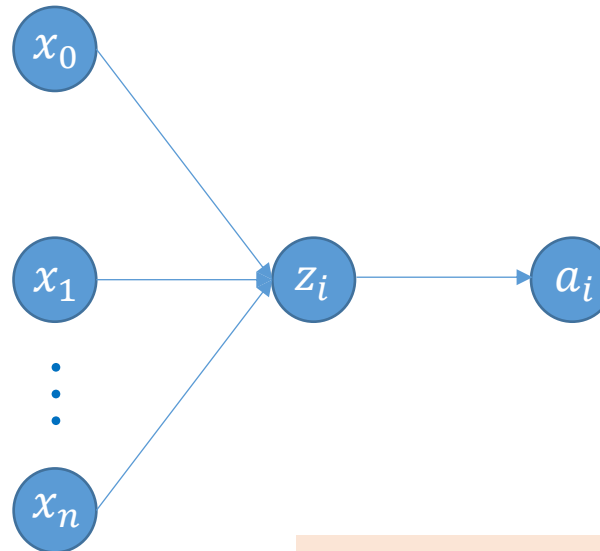
**❖ Xavier Initialization**

activation = tanh

$E(XY) = E(X)E(Y)$

$var(XY) = var(X)var(Y) +$

$\qquad var(X)\big(E(Y)\big)^2 +$

$\qquad var(y)\big(E(X)\big)^2$

Gaussian Distribution

$X \sim N(0, \sigma^2)$

$x_0$

$x_1$

$\vdots$

$x_n$

$z_i$

$a_i$

$var(w_i) \approx \dfrac{1}{n}$

$w_i \sim N(0, \sigma^2)$

$\sigma^2 = \dfrac{1}{n} \quad \Rightarrow \quad \sigma = \dfrac{1}{\sqrt{n}}$

$$W_i \sim N\left(0, \frac{1}{n}\right)$$

❖ **Xavier Initialization**

activation = tanh

**Uniform Distribution**

$$W_{ij} \sim U\left(-\frac{\sqrt{3}}{\sqrt{n}}, \frac{\sqrt{3}}{\sqrt{n}}\right)$$

**Gaussian Distribution**

$$W_{ij} \sim N\left(0, \frac{1}{n}\right)$$

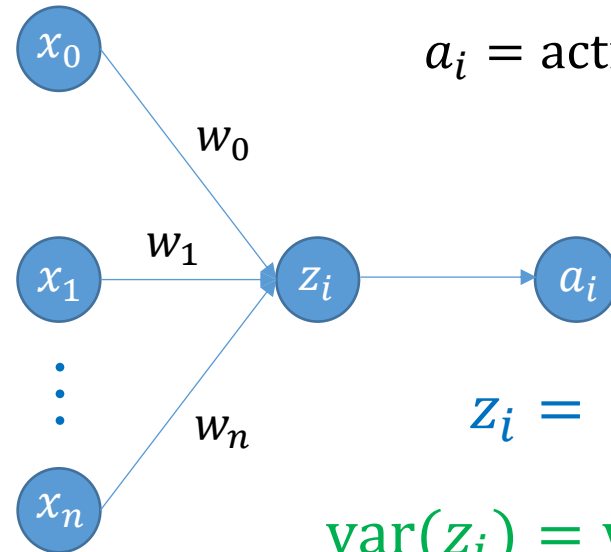# Initialization Methods

❖ **Xavier Initialization**



$a_i = \text{activation}(z_i)$

$E(X) = 0$

$E(W) = 0$

$b = 0$

$$E(XY) = E(X)E(Y)$$

$$var(XY) = var(X)var(Y) +$$
$$var(X)\big(E(Y)\big)^2 +$$
$$var(y)\big(E(X)\big)^2$$

$$z_i = (x_1 w_1 + \cdots + x_n w_n + b)$$

$$var(z_i) = var(x_1 w_1 + \cdots + x_n w_n + b)$$

$$= n\,var(x_i w_i) = n\,var(x_i)var(w_i)$$

## Uniform Distribution

$$X \sim U(a, b)$$

$$f(x) = \frac{1}{b - a}$$

$$var[X] = \frac{(b - a)^2}{12}$$

$\text{activation} = \text{sigmoid} \;\Rightarrow\; a_i = \text{sigmoid}(z_i) \approx \frac{1}{2} + \frac{z_i}{4}$

$\Rightarrow\; 16\,var(a_i) = var(z_i)$

$var(X) = var(\mathbf{a}) \;\xrightarrow{\text{iid}}\; var(x_i) = var(a_i) \;\Rightarrow\; n\,var(w_i) = 16$

$\Rightarrow\; var(w_i) = \frac{16}{n}$

❖ **Xavier Initialization**

activation = sigmoid

$$E(XY) = E(X)E(Y)$$

$$var(XY) = var(X)var(Y) + var(X)\big(E(Y)\big)^2 + var(y)\big(E(X)\big)^2$$

Uniform Distribution

$$X \sim U(a, b)$$

$$f(x) = \frac{1}{b-a}$$

$$var[X] = \frac{(b-a)^2}{12}$$

$x_0$

$x_1$

$\vdots$

$x_n$

$z_i$

$a_i$

$$var(w_i) \approx \frac{16}{n}$$

$$w_i \sim U(-r, r)$$

$$var[w_i] = \frac{r^2}{3}$$

$$W_i \sim U\left(-\frac{4\sqrt{3}}{\sqrt{n}}, \frac{4\sqrt{3}}{\sqrt{n}}\right)$$
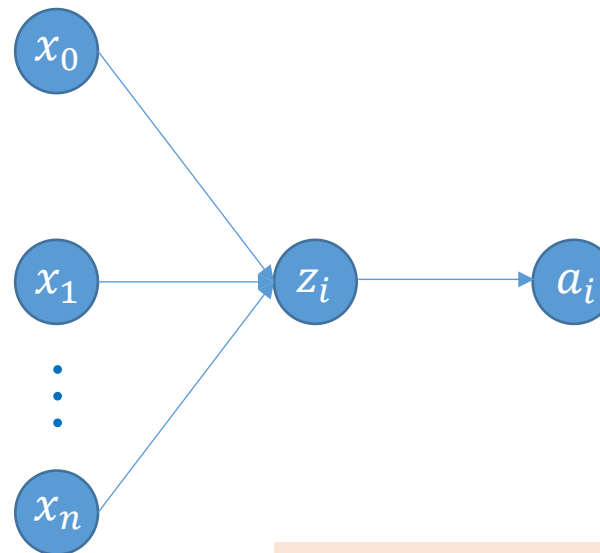
31

❖ **Xavier Initialization**

activation = sigmoid

$$E(XY) = E(X)E(Y)$$

$$var(XY) = var(X)var(Y) +$$
$$var(X)\big(E(Y)\big)^2 +$$
$$var(y)\big(E(X)\big)^2$$

Gaussian Distribution

$$X \sim N(0, \sigma^2)$$

$x_0$

$x_1$

$\vdots$

$x_n$

$z_i$

$a_i$
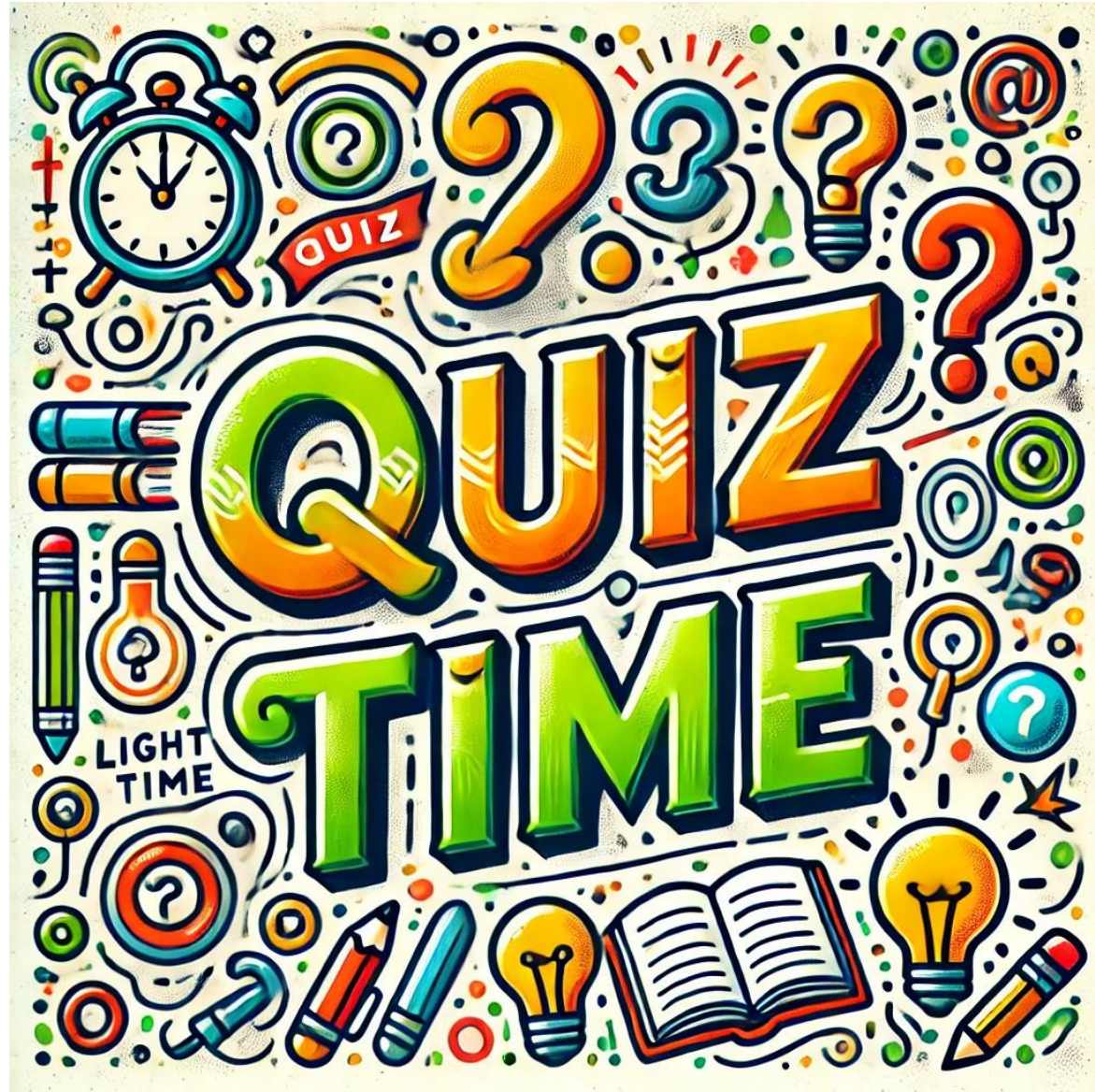
$$var(w_i) \approx \frac{16}{n}$$

$$w_i \sim N(0, \sigma^2)$$

$$\sigma^2 = \frac{1}{n}$$

$$W_i \sim N\left(0, \frac{16}{n}\right)$$

# Question 1

❖ Chuẩn hóa dữ liệu nào nên dùng cho Glorot Initialization (chọn nhiều đáp án)?

a) Sau chuẩn hóa có range là [0, 255]

b) Có range là [0, 1]

c) Có range là [-1, 1]

d) Dạng z-score

# Question 2

❖ Glorot Initialization giả định activation đang dùng là gì (chọn nhiều đáp án)?

a) Sigmoid                                    b) Tanh

c) ReLU                                       d) PReLU

# Question 3

❖ Code nào nên dùng khi sử dụng với Xavier Init.?



```
Compose([transforms.ToTensor(), transforms.Normalize((0.5,), (0.5,))])    1

Compose([transforms.ToTensor(), transforms.Normalize((0,), (1.0,))])    2

Compose([transforms.ToTensor(), transforms.Normalize((mean,), (std,))])    3

transforms.Compose([transforms.ToTensor(),
                    transforms.Normalize((0,),
                                        (1.0/255,))])    4
```
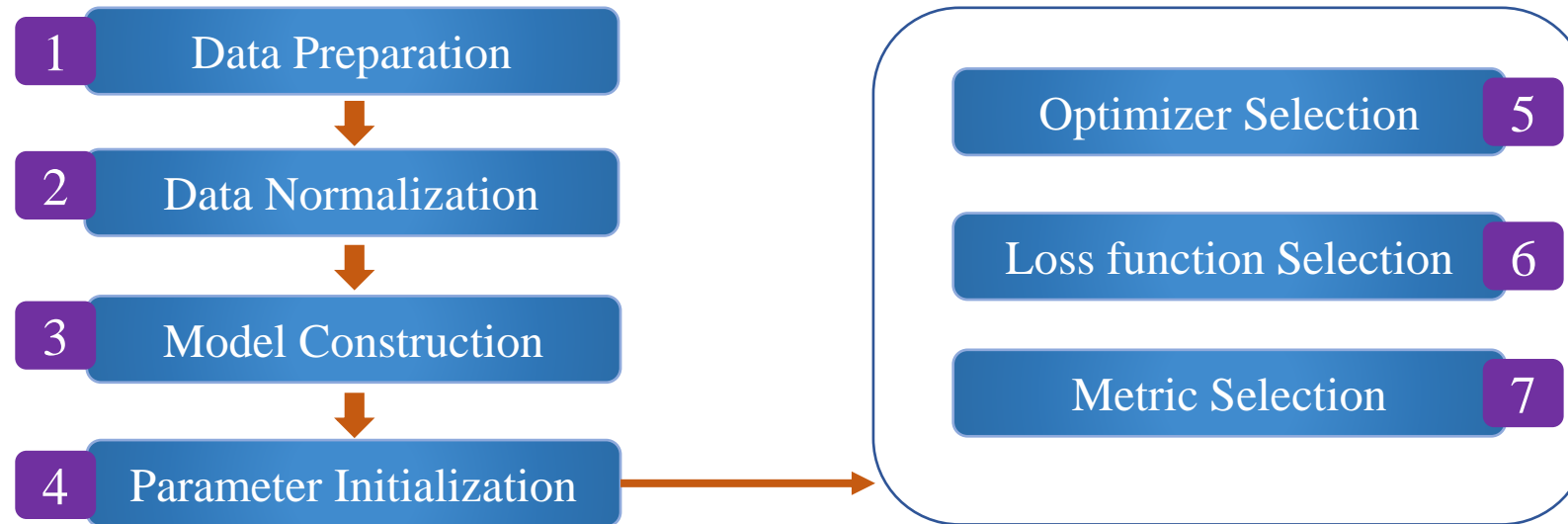
a) Code 1                    b) Code 2

c) Code 3                    d) Code 4

❖ Dựa vào kiến thức AIO tới thời điểm này, nếu chọn bước (4) là Glorot init., thì bước nào cần có hành động (gì đó) tương ứng?

| 1 | Data Preparation |
| 2 | Data Normalization |
| 3 | Model Construction |
| 4 | Parameter Initialization |

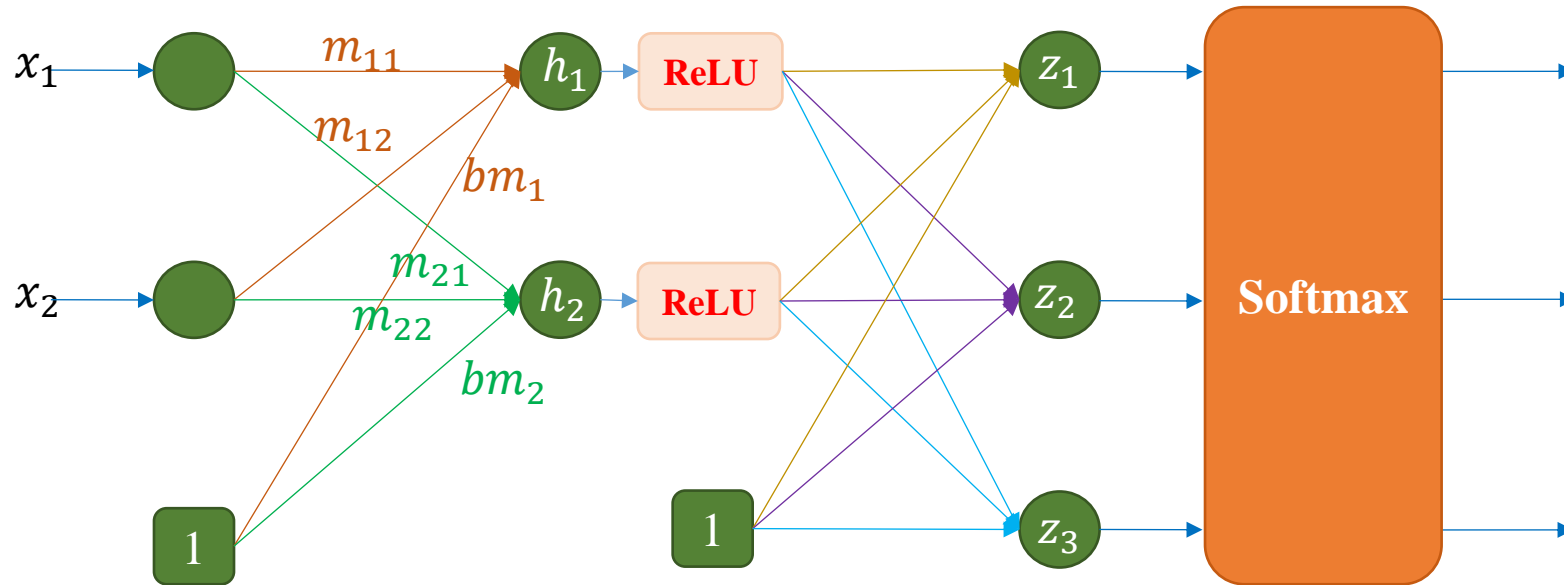| Optimizer Selection | 5 |
| Loss function Selection | 6 |
| Metric Selection | 7 |

a) Bước (1)

b) Bước (2)

c) Bước (3)

d) Bước (5)

# Question 5

❖ Hãy chọn 1 giải pháp hay nhất để khắc phục vấn đề dying relu?



a) Thay relu bằng sigmoid

b) Thay relu bằng tanh

c) Thay relu bằng prelu

d) Giữ relu và dùng He init.

❖ Glorot Init. có những giả định (điều kiện cho trước) nào?

a) Activation là sigmoid

b) Activation là tanh

c) Activation là relu

d) Data input có mean=0

# Outline

$$W_i \sim U\left(-\frac{\sqrt{6}}{\sqrt{n}}, \frac{\sqrt{6}}{\sqrt{n}}\right)$$

$$W_i \sim N\left(0, \frac{2}{n}\right)$$

$x_i \sim a_i$

❖ **Kaiming He Initialization**



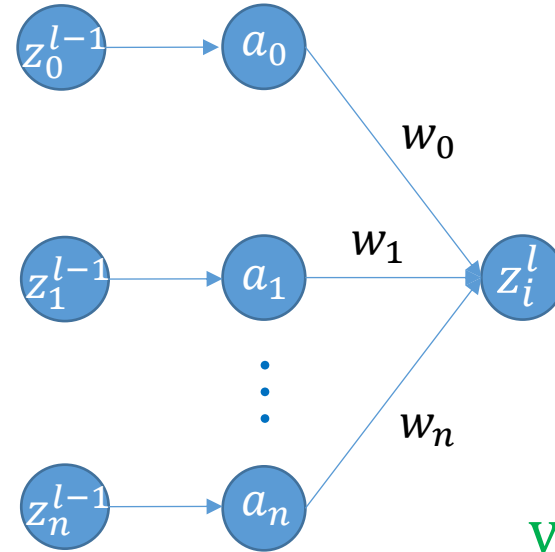$a_i = \text{activation}(z_i)$

$E(W) = 0$

$b = 0$

$E(XY) = E(X)E(Y)$

$var(XY) = var(X)var(Y) +$
$var(X)\big(E(Y)\big)^2 +$
$var(y)\big(E(X)\big)^2$

$z_i = (a_1 w_1 + \cdots + a_n w_n + b)$

$var(z_i) = var(a_1 w_1 + \cdots + a_n w_n + b)$

Uniform Distribution

$X \sim U(a, b)$

$f(x) = \dfrac{1}{b - a}$

$var[X] = \dfrac{(b-a)^2}{12}$

activation = relu ➡ $a_i = max(0, z_i)$

$var(z^{l-1}) = var(z^l)$ $\xrightarrow{\text{iid}}$ $var(z_i^{l-1}) = var(z_i^l)$ ➡ $nvar(w_i) = 2$

➡ $var(w_i) = \dfrac{2}{n}$

# Initialization Methods

❖ **He Initialization**

activation = relu

$$E(XY) = E(X)E(Y)$$

$$var(XY) = var(X)var(Y) +$$
$$var(X)\big(E(Y)\big)^2 +$$
$$var(y)\big(E(X)\big)^2$$



$$var(w_i) \approx \frac{2}{n}$$

$$w_i \sim U(-r, r)$$

$$var[w_i] = \frac{r^2}{3}$$

## Uniform Distribution

$$X \sim U(a, b)$$

$$f(x) = \frac{1}{b - a}$$

$$var[X] = \frac{(b - a)^2}{12}$$

$$W_i \sim U\left(-\frac{\sqrt{6}}{\sqrt{n}}, \frac{\sqrt{6}}{\sqrt{n}}\right)$$
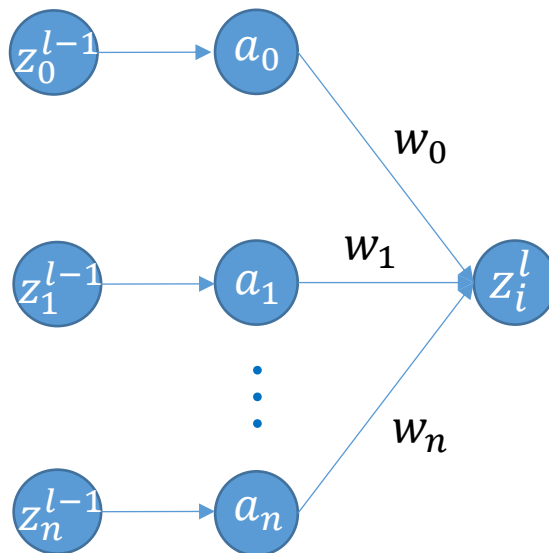
34

# Initialization Methods

❖ **He Initialization**

activation = he

$$E(XY) = E(X)E(Y)$$

$$var(XY) = var(X)var(Y) +$$
$$var(X)\big(E(Y)\big)^2 +$$
$$var(y)\big(E(X)\big)^2$$

Gaussian Distribution

$$X \sim N(0, \sigma^2)$$



$z_0^{l-1} \rightarrow a_0$

$z_1^{l-1} \rightarrow a_1$ $\quad w_0$

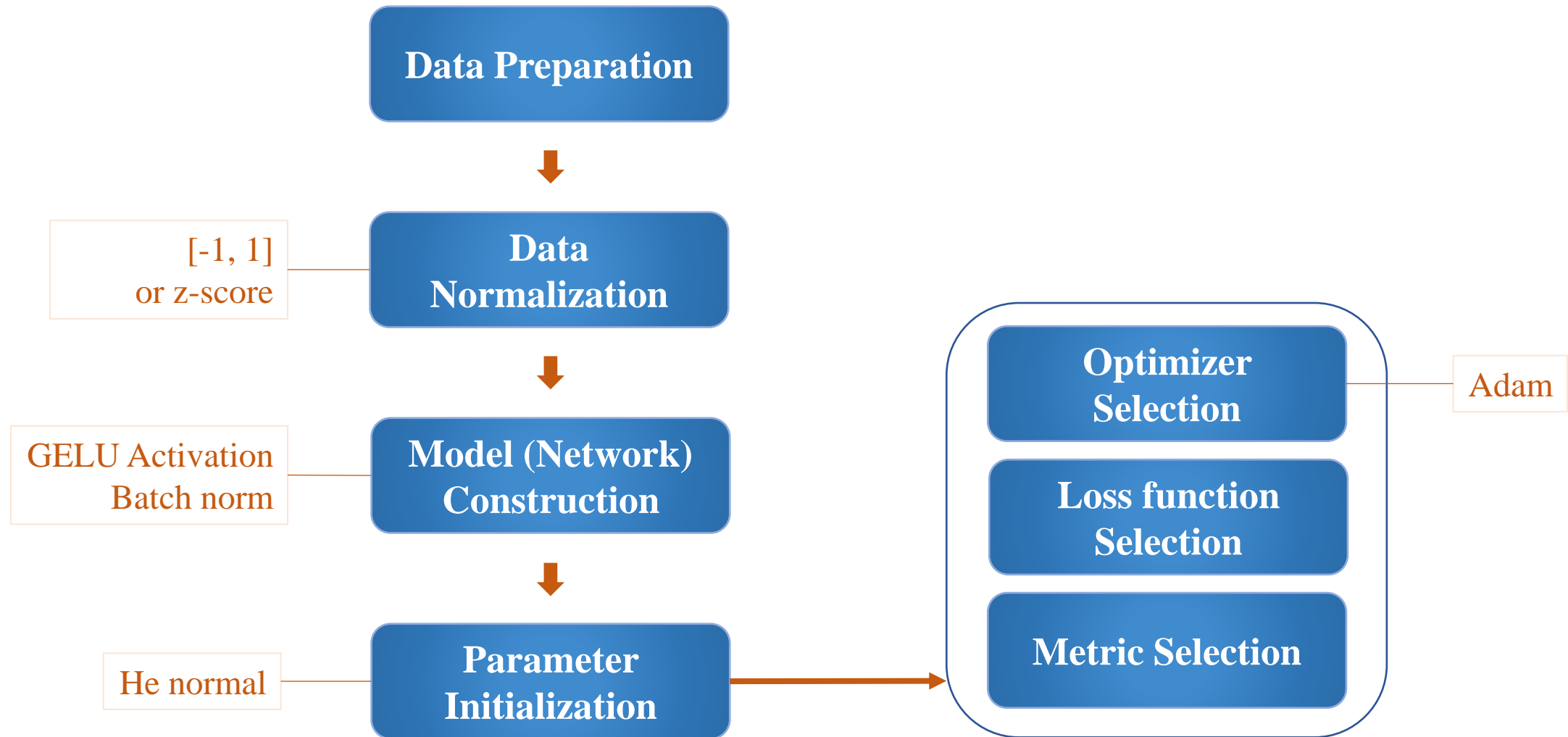$\qquad w_1 \quad z_i^l$

$z_n^{l-1} \rightarrow a_n$ $\quad w_n$

$$\text{var}(w_i) \approx \frac{2}{n}$$

$$w_i \sim N(0, \sigma^2)$$

$$\sigma^2 = \frac{1}{n}$$

$$W_i \sim N\left(0, \frac{2}{n}\right)$$

35

# **Further Reading**

## Dying ReLU

https://towardsdatascience.com/the-dying-relu-problem-clearly-explained-42d0c54e0d24

## Initialization

https://www.deeplearning.ai/ai-notes/initialization/index.html