

Multi-layer Perception

Initialization (Advanced)

Quang-Vinh Dinh
Ph.D. in Computer Science

Outline

- **Case Studies**
- **Gradient Vanishing**
- **Gradient Explosion**
- **Xavier Glorot Initialization**
- **Kaiming He Initialization**

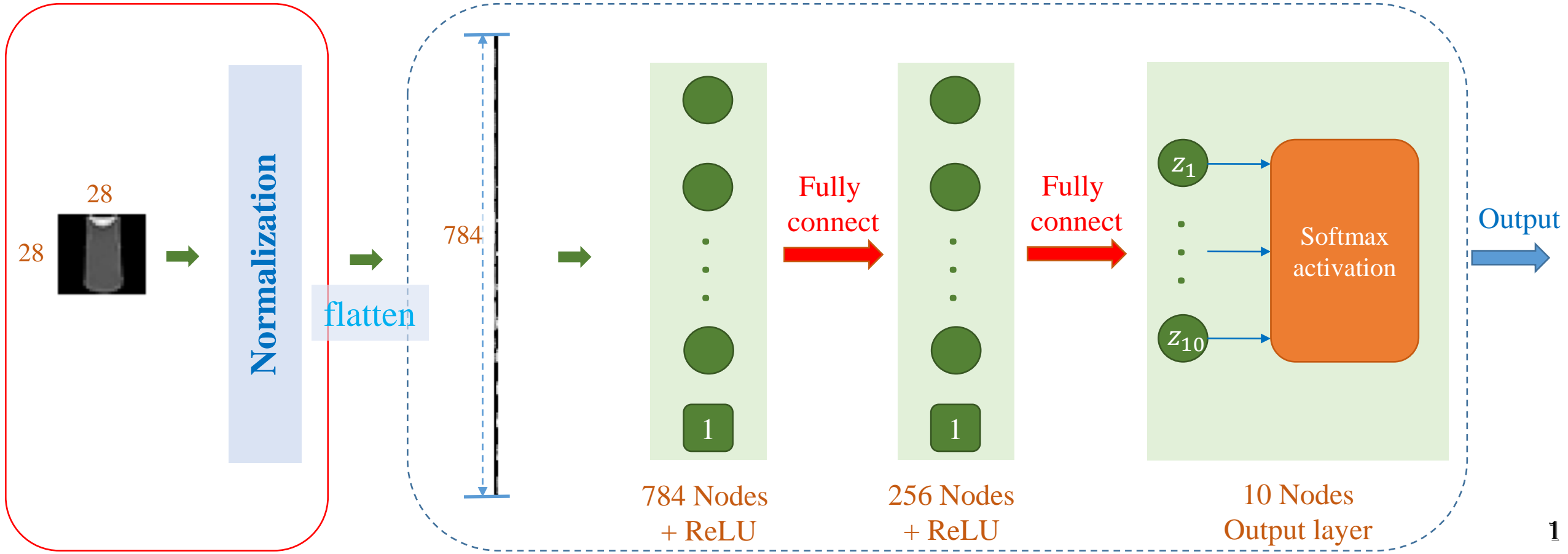
$$X \in [0, 255]$$

Normalize(*mean*, *std*)

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

```
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)
)
```



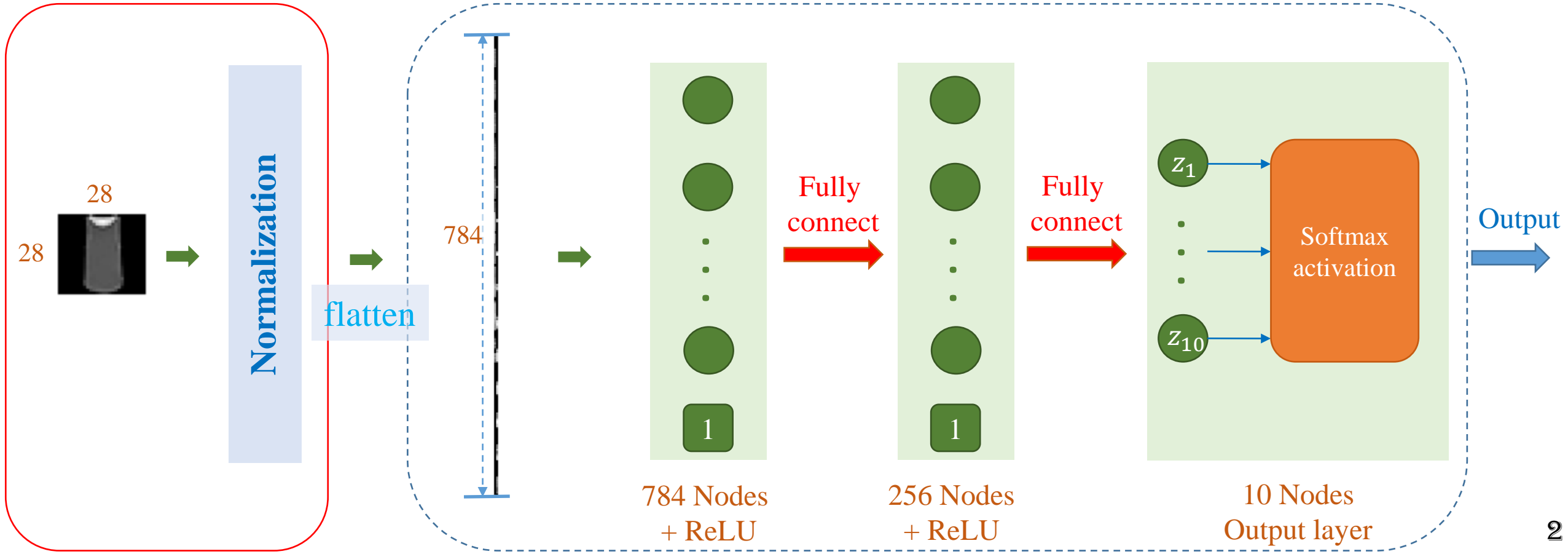
$$X \in [-1, 1]$$

Normalize(*mean*, *std*)

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

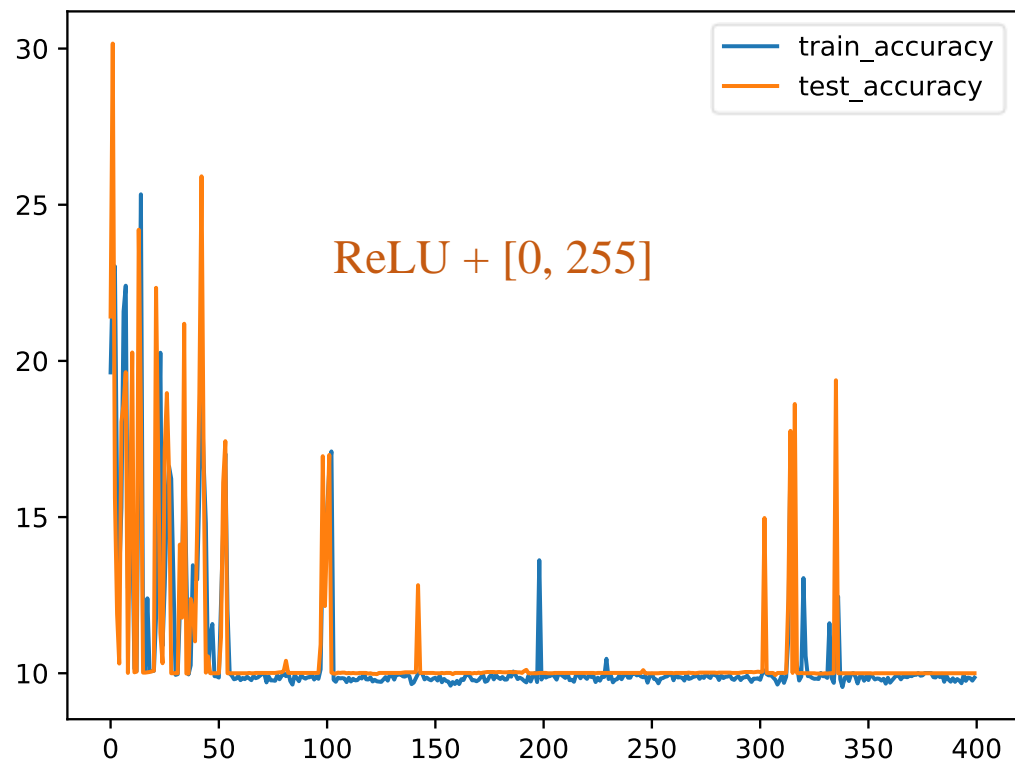
```
transform = transforms.Compose([transforms.ToTensor(),
                                transforms.Normalize((0.5,),
                                                    (0.5,))])

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

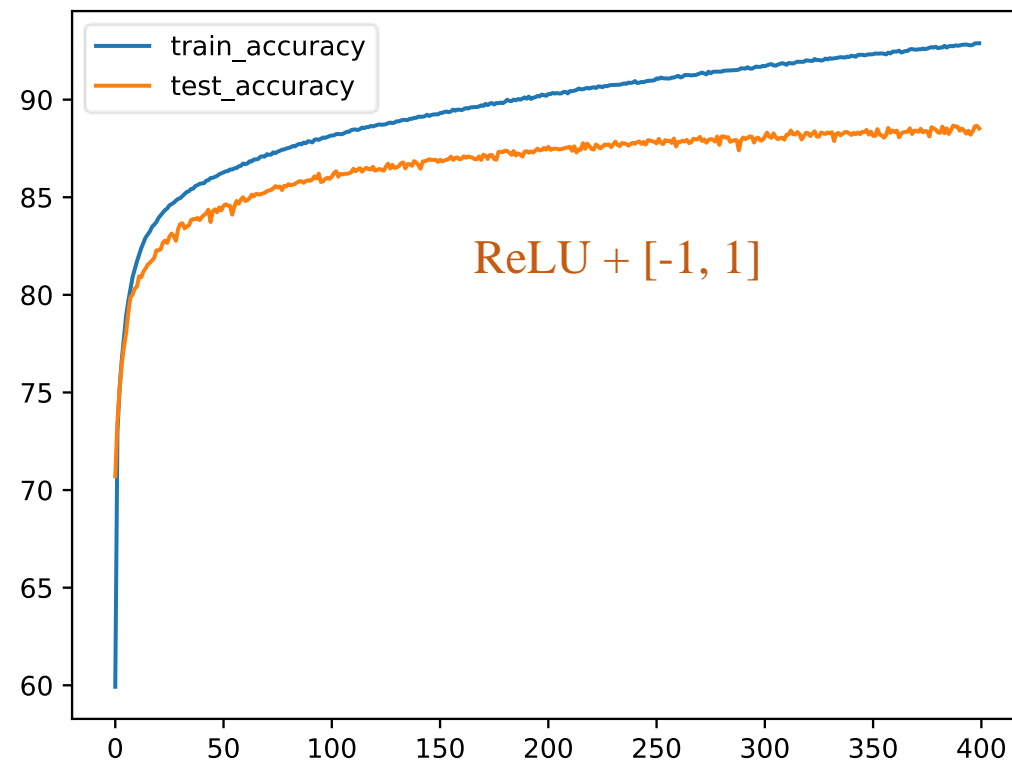


Experimental Results

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



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



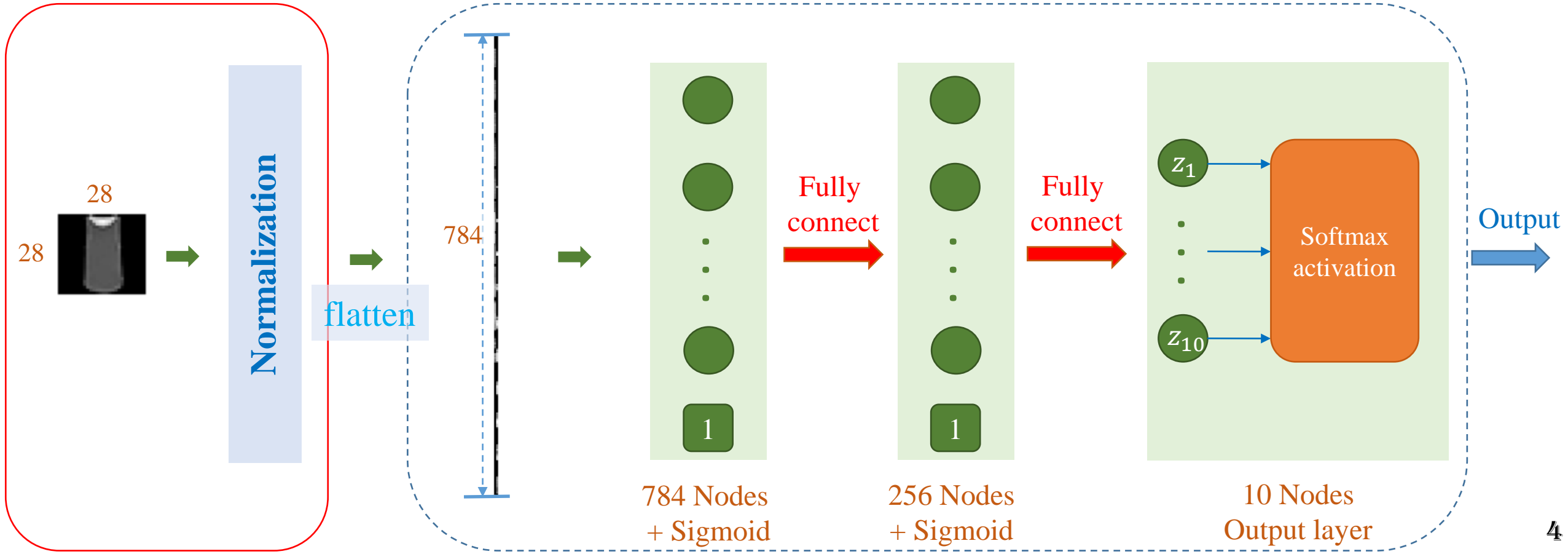
$$X \in [0, 255]$$

Normalize(*mean*, *std*)

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

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

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



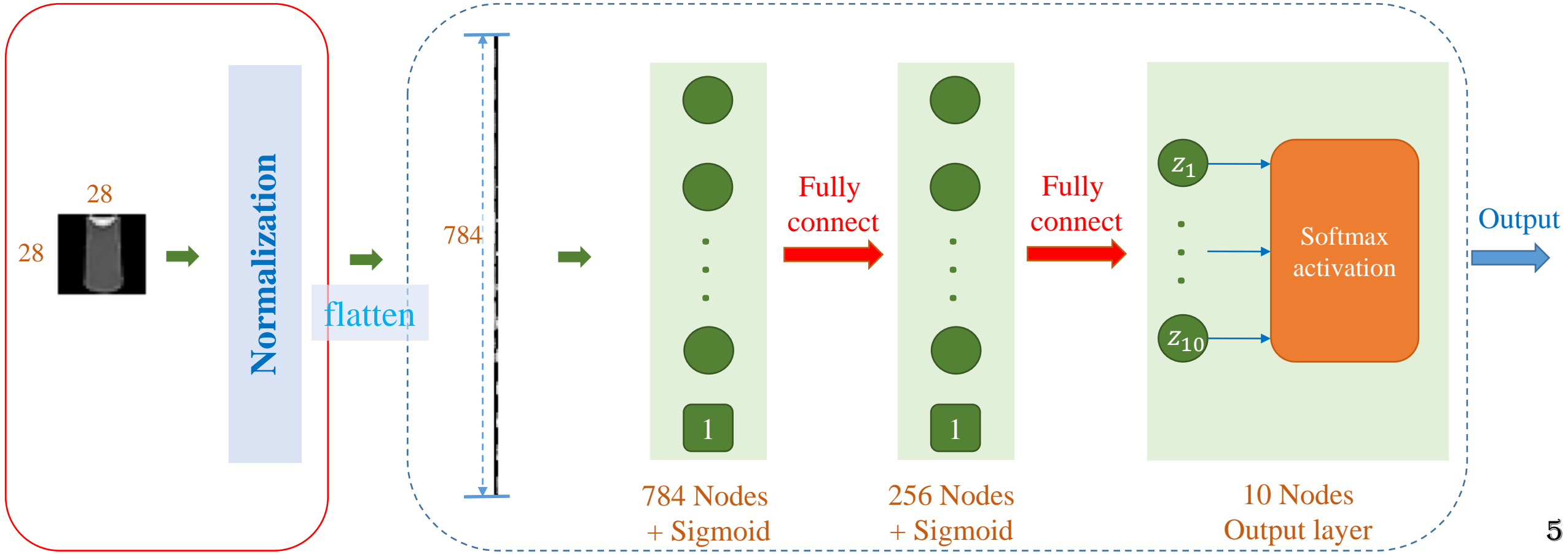
$$X \in [-1, 1]$$

Normalize(*mean*, *std*)

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

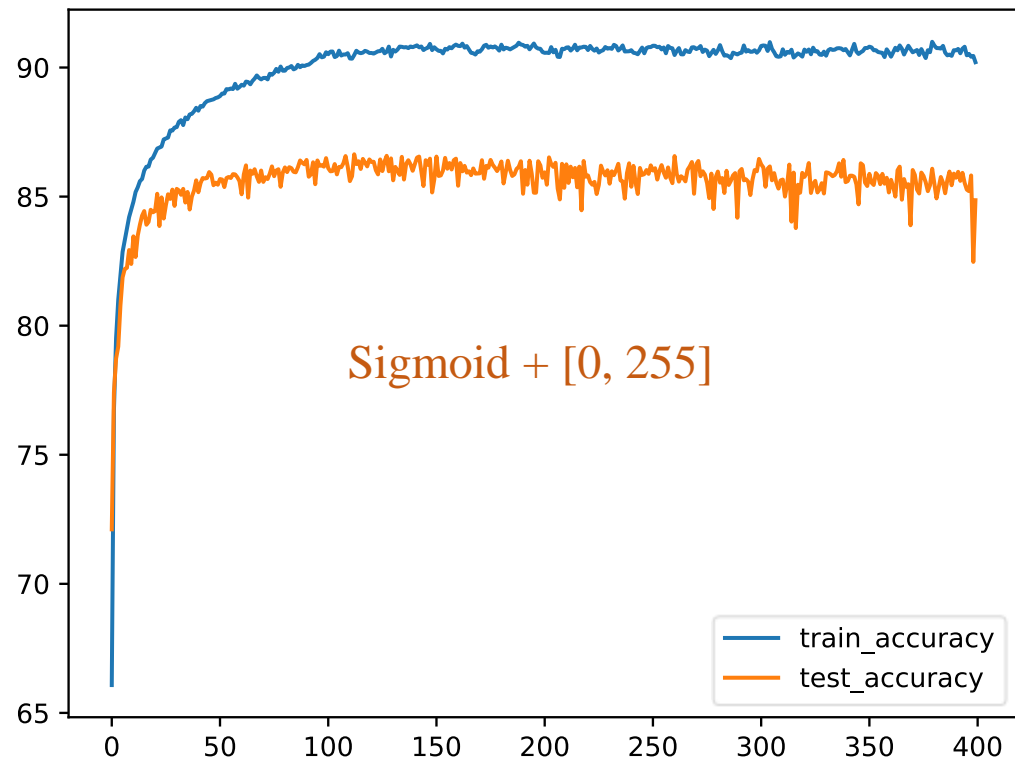
```
transform = Compose([ToTensor(),
                      Normalize((0.5,),
                                (0.5,))])
```

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

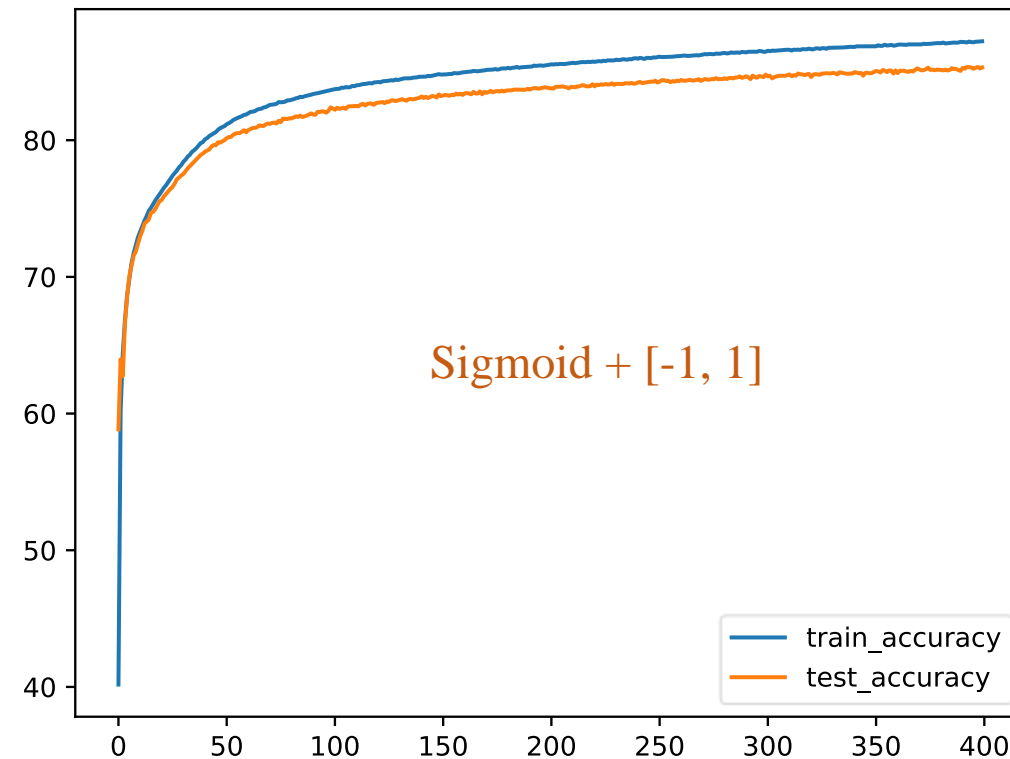


Experimental Results

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



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

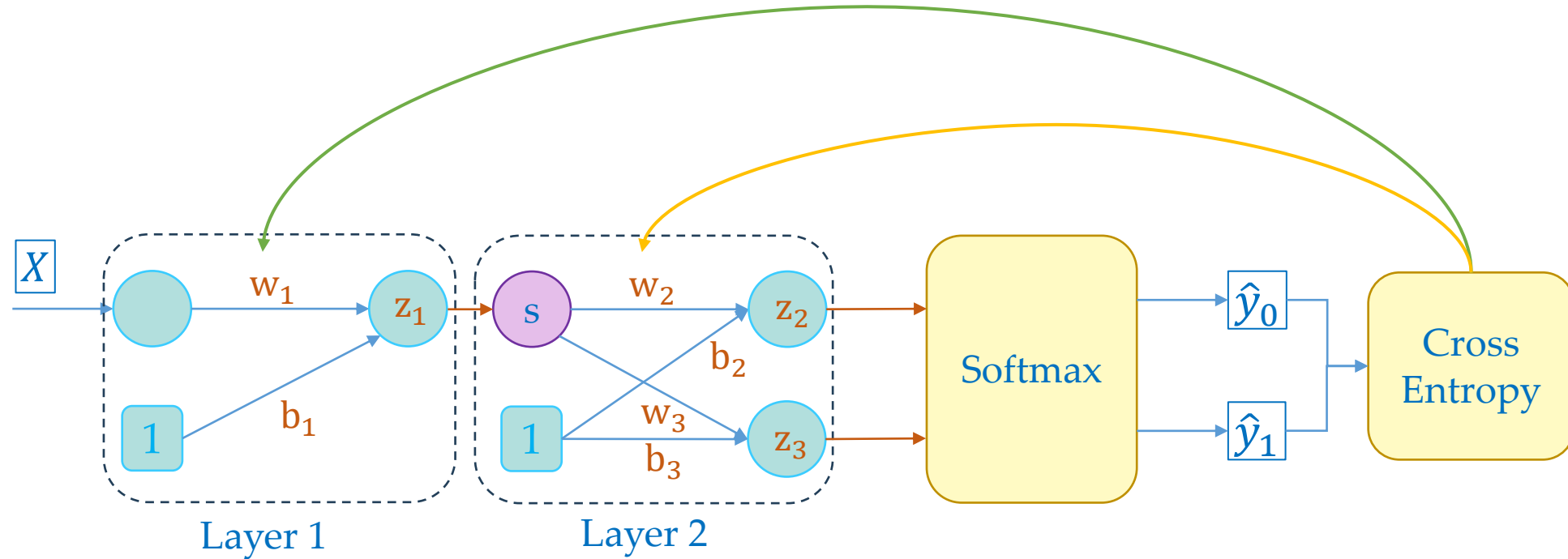


Outline

- **Case Studies**
- **Gradient Vanishing**
- **Gradient Explosion**
- **Xavier Glorot Initialization**
- **Kaiming He Initialization**

Gradient Vanishing

Large weight initialization



Gradient Vanishing

Large weight initialization

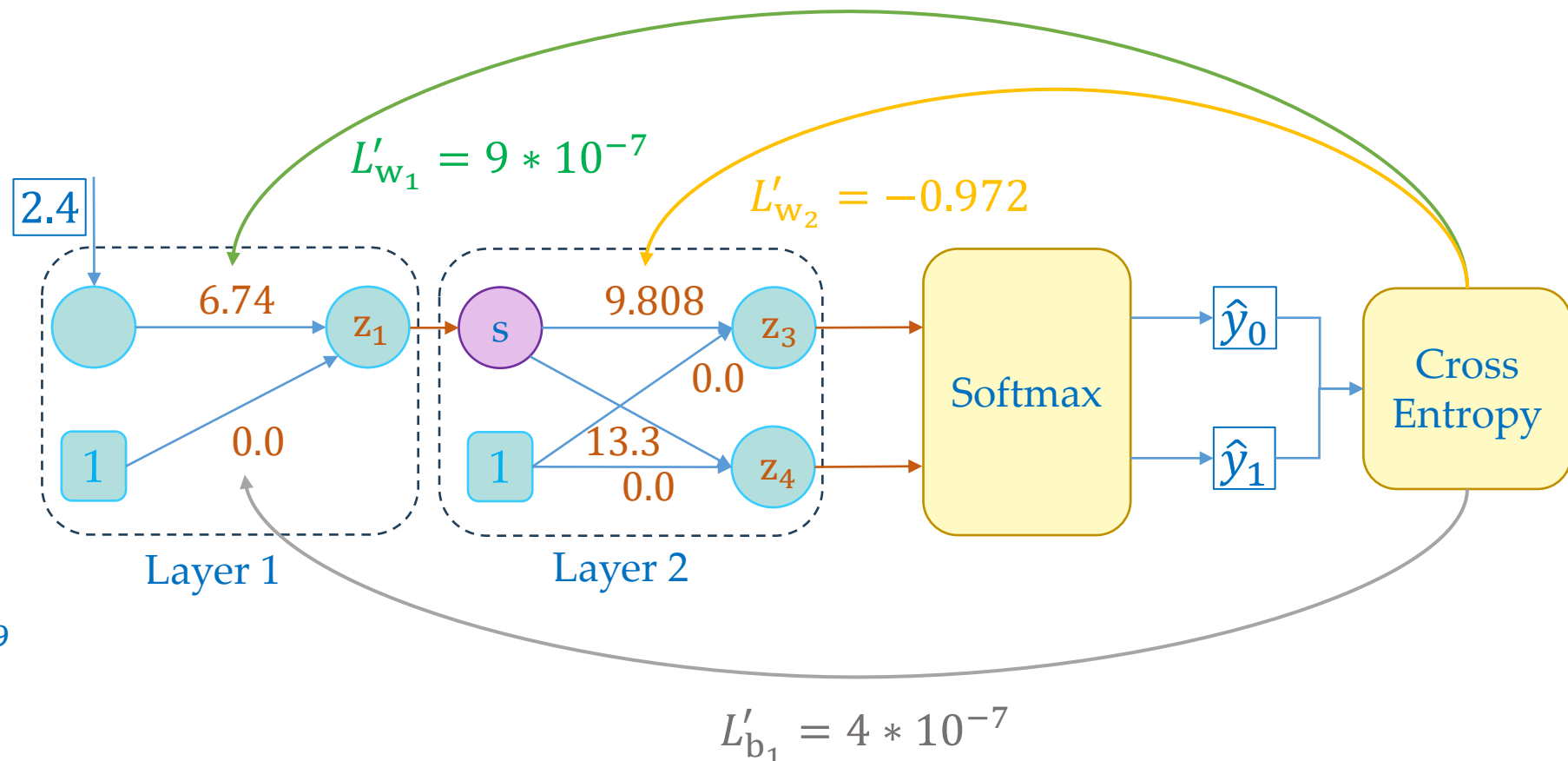
```
linear1 = nn.Linear(1, 1)
linear2 = nn.Linear(1, 2)

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

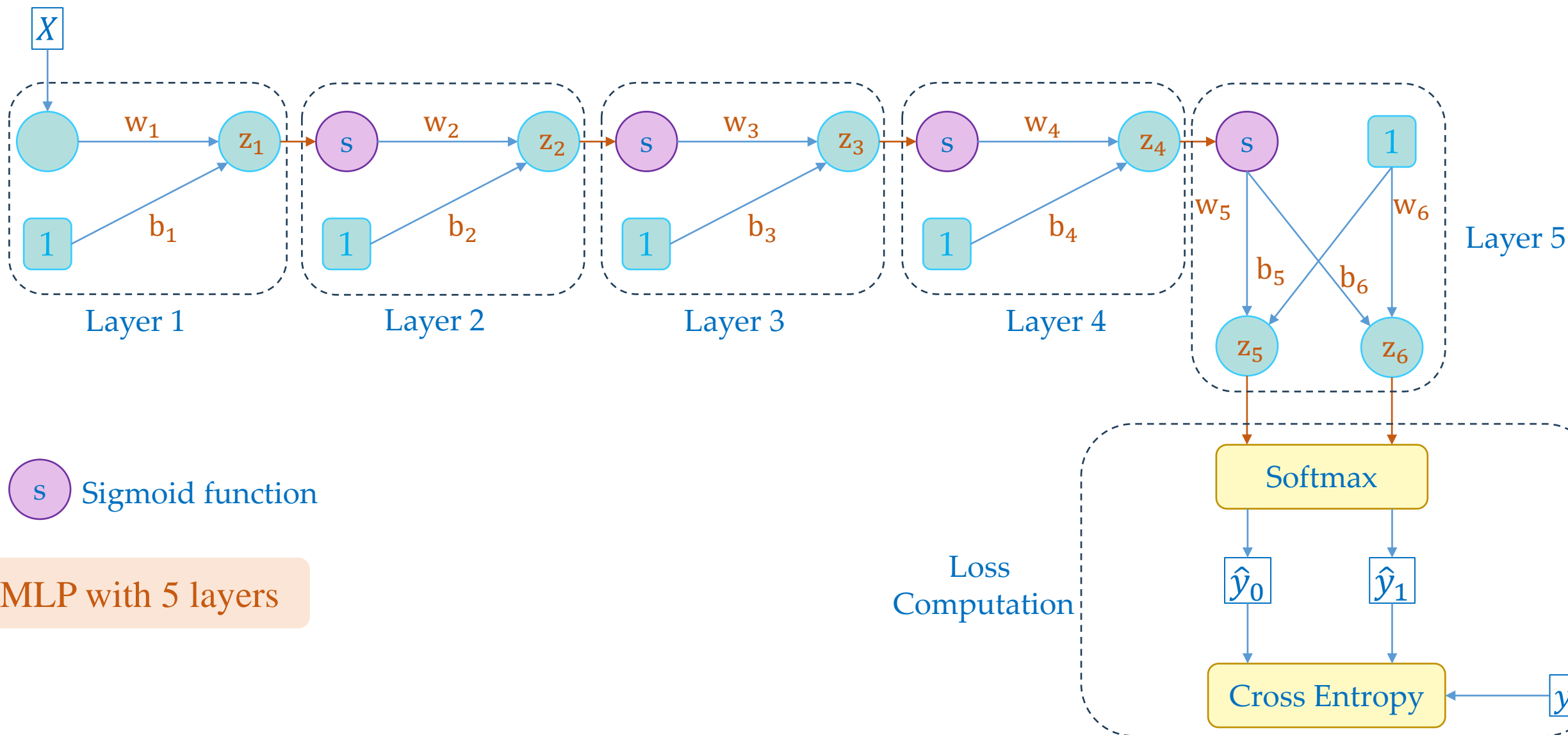
with $\eta = 0.01$

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

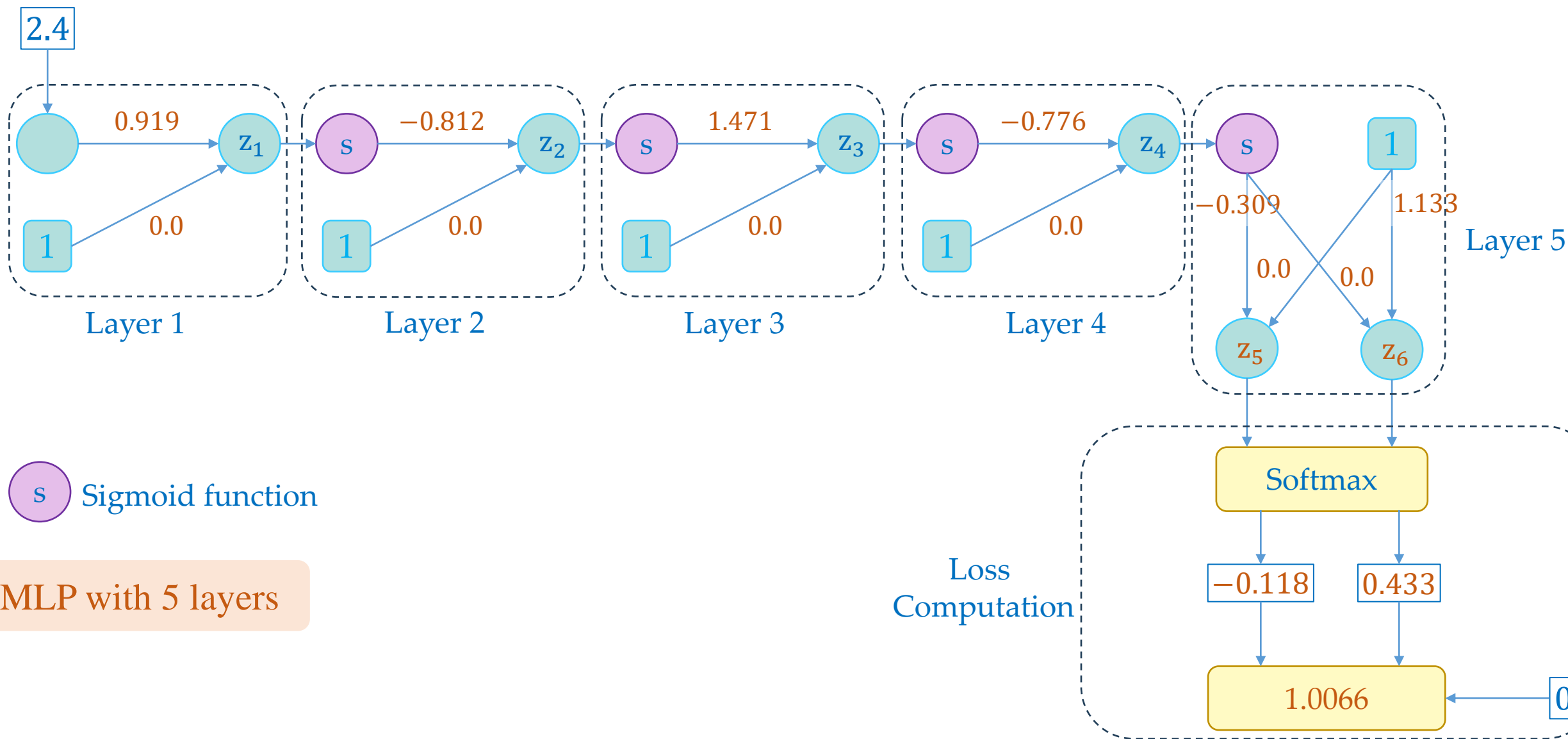
$$\eta L'_{b_1} = 4 * 10^{-9}$$



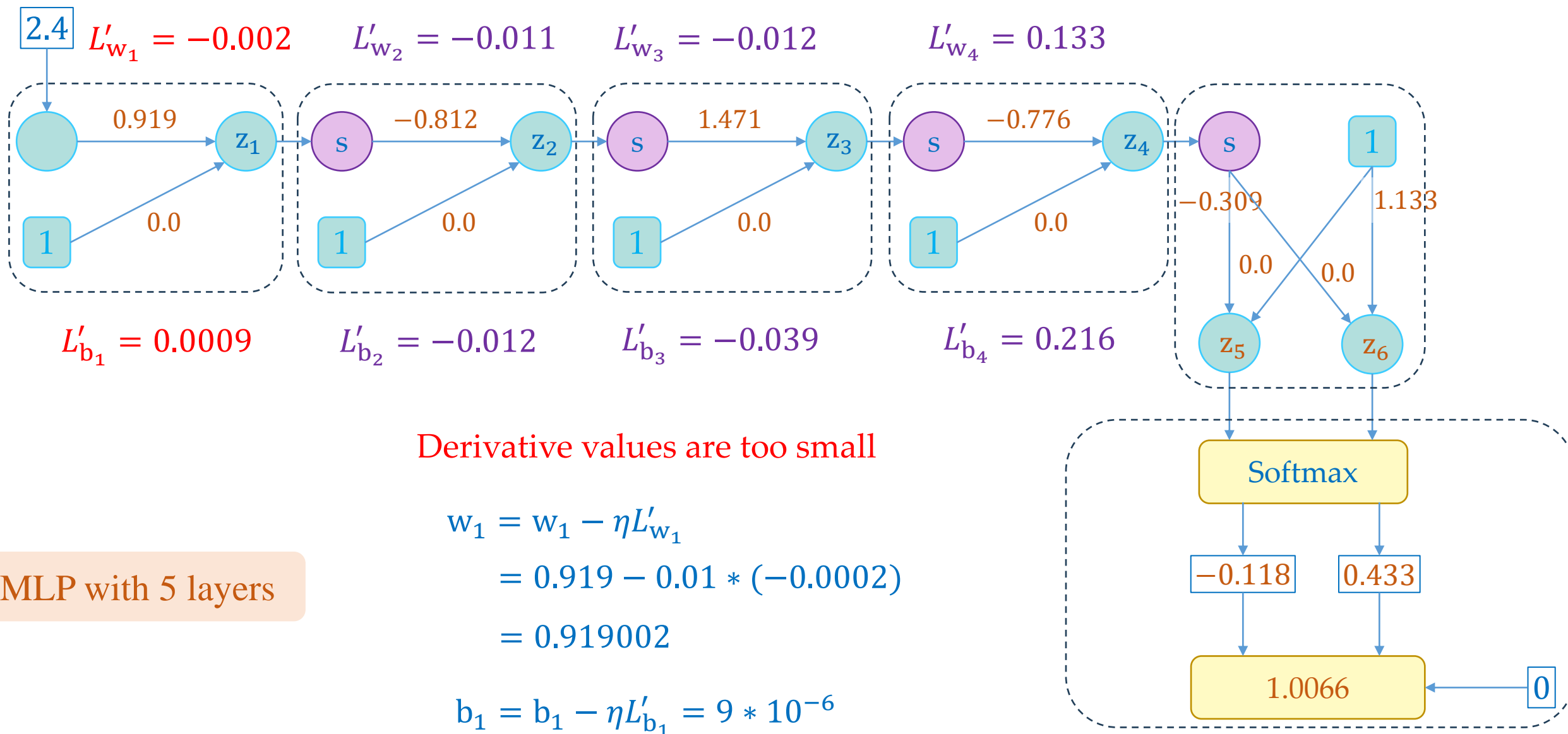
Gradient Vanishing



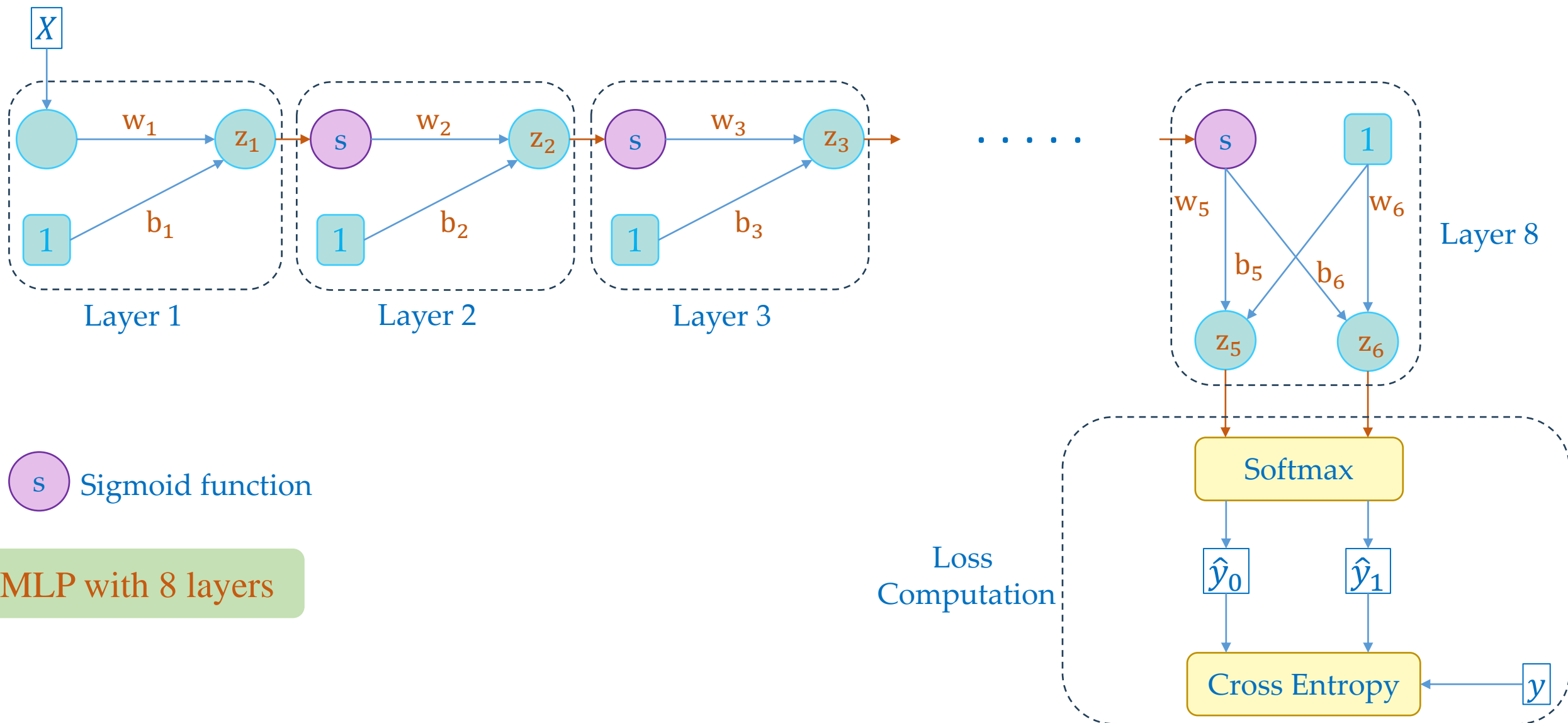
Gradient Vanishing



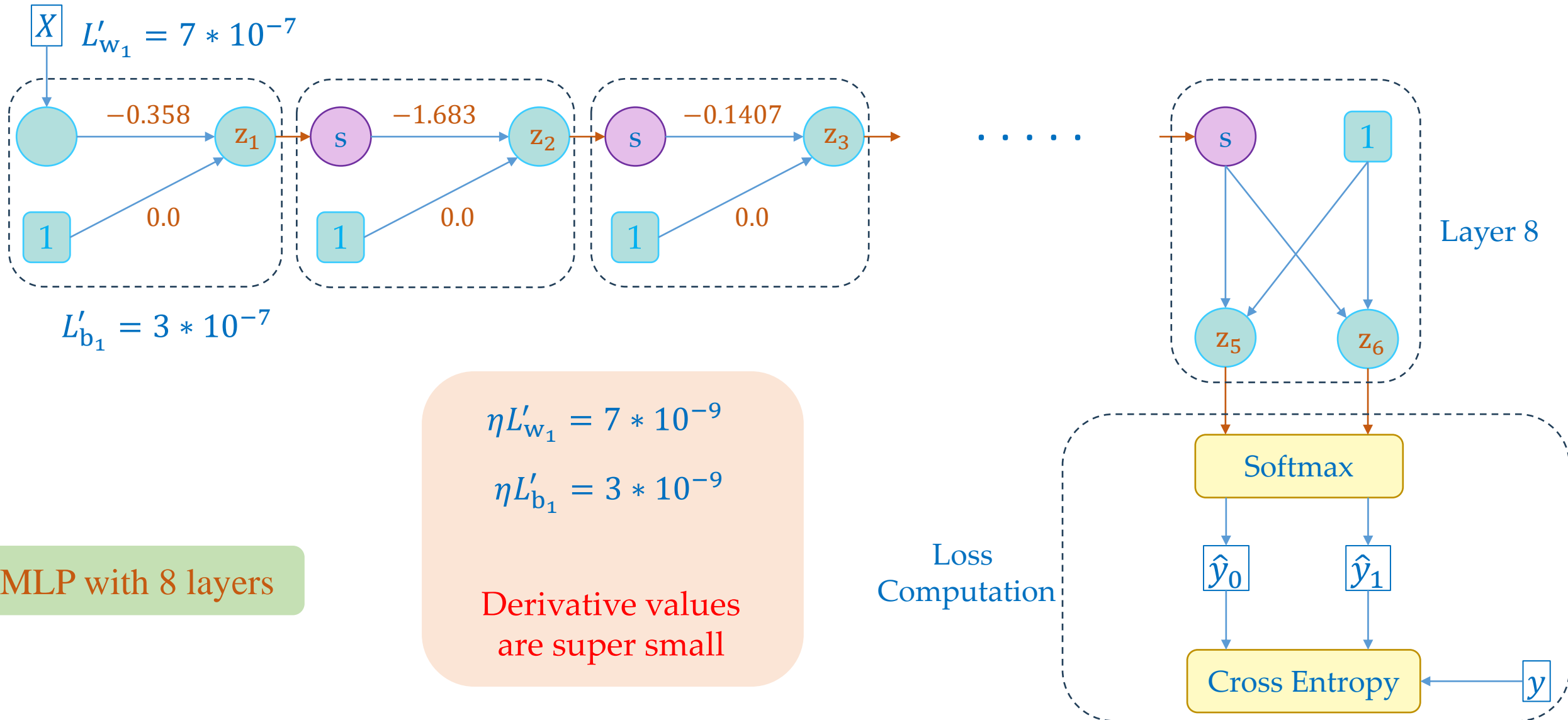
Gradient Vanishing



Gradient Vanishing



Gradient Vanishing



Gradient Explosion

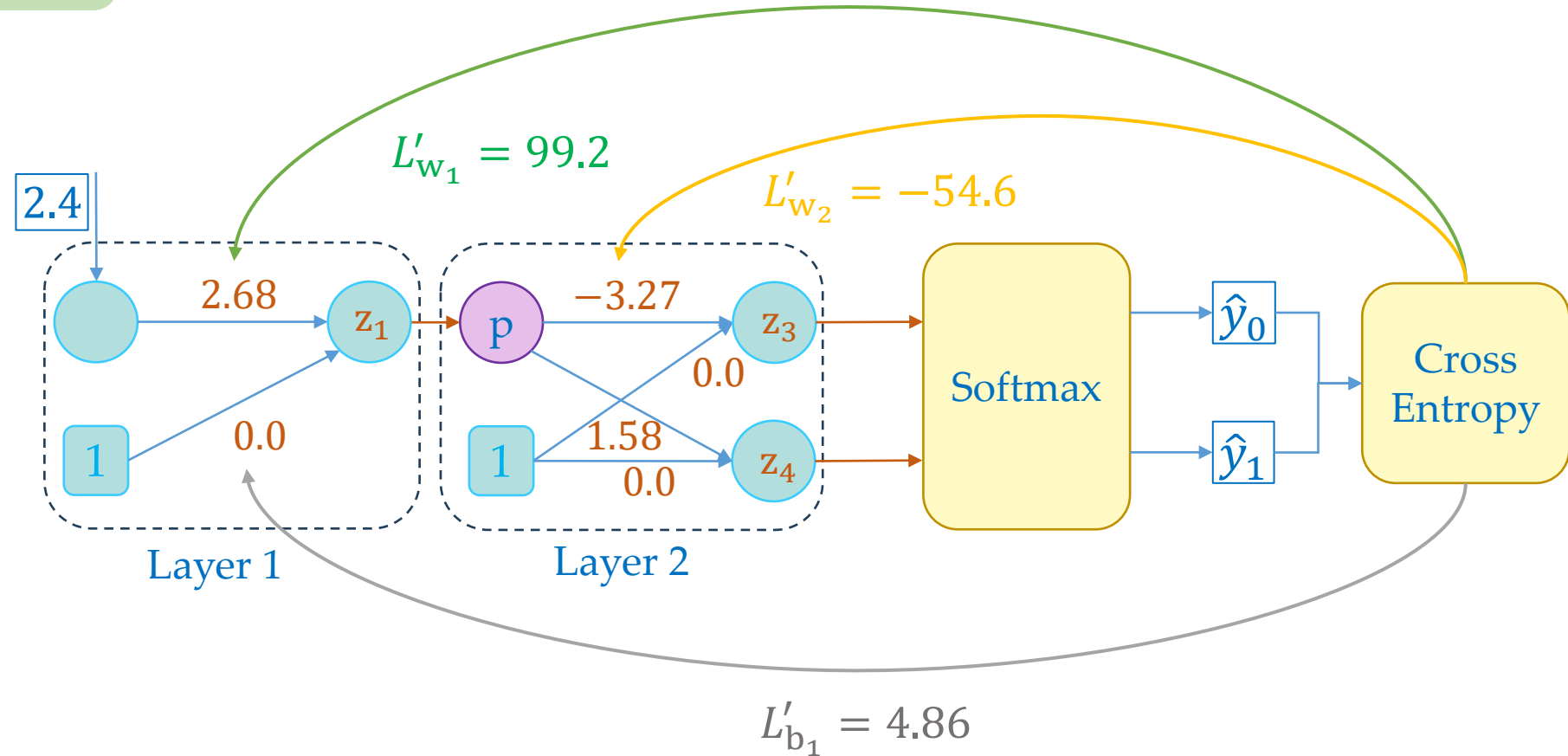
Large weight initialization
and large learning rate

s PReLU function

with $\eta = 10$

$$\eta L'_{w_1} = 99$$

$$\eta L'_{b_1} = 48.6$$



Outline

- **Case Studies**
- **Gradient Vanishing**
- **Gradient Explosion**
- **Xavier Glorot Initialization**
- **Kaiming He Initialization**

Mean

Data

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

Formula

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

Given the data

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

$$N = 6$$

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

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

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

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

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

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

Mean

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

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

Standard
deviation

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

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

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

$$\begin{aligned} \text{var}(X) &= \frac{1}{5} [(5 - 5)^2 + (3 - 5)^2 + (6 - 5)^2 + \\ &\quad (7 - 5)^2 + (4 - 5)^2] \\ &= \frac{1}{5} (0 + 4 + 1 + 4 + 1) = 2 \end{aligned}$$

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

Variance

Formula

mean

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

variance

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

Standard deviation

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

$$\begin{aligned} \text{var}(X) &= \sum_{i=1}^N (X_i - E(X))^2 P_X(X_i) \\ &= \sum_{i=1}^N (X_i^2 - 2X_i E(X) + E(X)^2) P_X(X_i) \\ &= \sum_{i=1}^N X_i^2 P_X(X_i) - \sum_{i=1}^N 2X_i E(X) P_X(X_i) \\ &\quad + \sum_{i=1}^N E(X)^2 P_X(X_i) \\ &= E(X^2) - 2E(X) \left[\sum_{i=1}^N X_i P_X(X_i) \right] + E(X)^2 \\ &= E(X^2) - (E(X))^2 \end{aligned}$$

Variance

$$\text{var}(X) = E(X^2) - (E(X))^2$$

$$\begin{aligned}\text{var}(XY) &= E(X^2Y^2) - (E(XY))^2 \\ &= E(X^2)E(Y^2) - (E(X)E(Y))^2 \\ &= \left[\text{var}(X) + (E(X))^2 \right] \left[\text{var}(Y) + (E(Y))^2 \right] - (E(X)E(Y))^2 \\ &= \text{var}(X)\text{var}(Y) + \text{var}(X)(E(Y))^2 + \text{var}(Y)(E(X))^2\end{aligned}$$

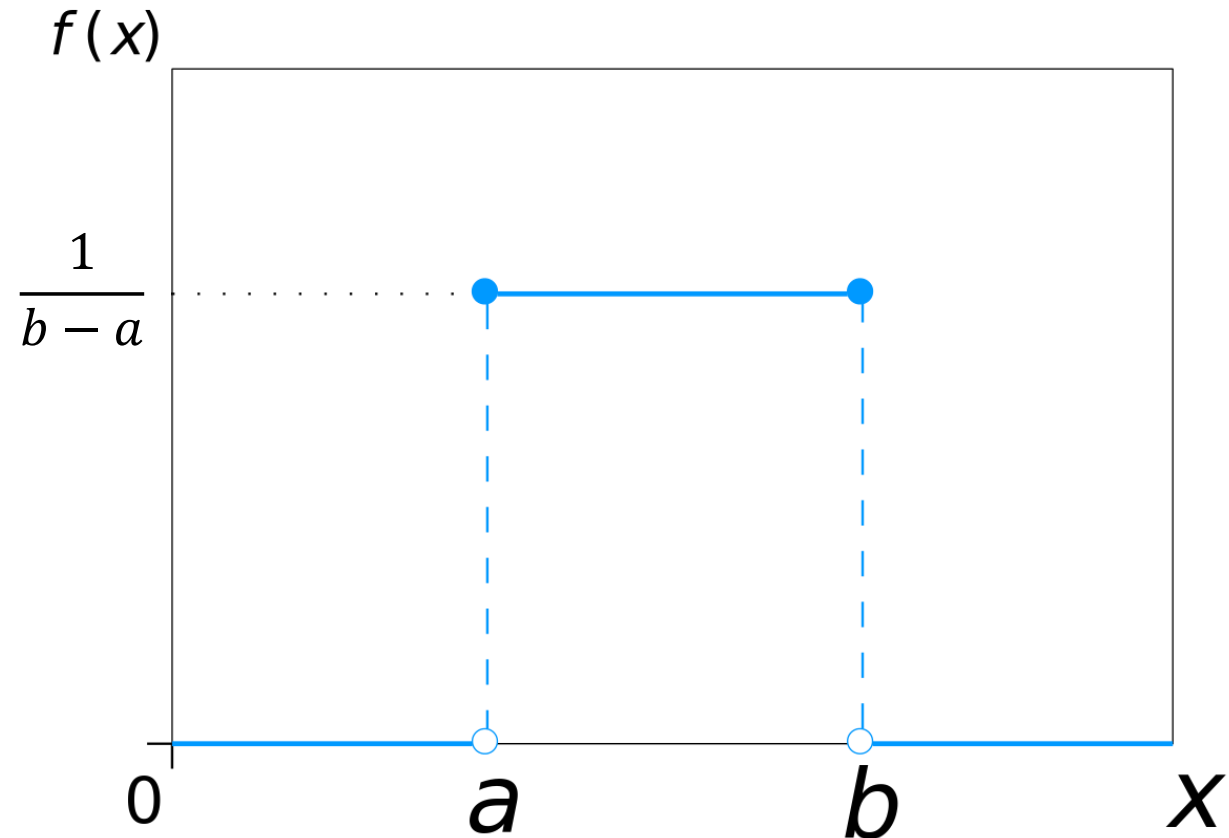
Initialization Methods

Xavier Initialization

Uniform Distribution

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

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

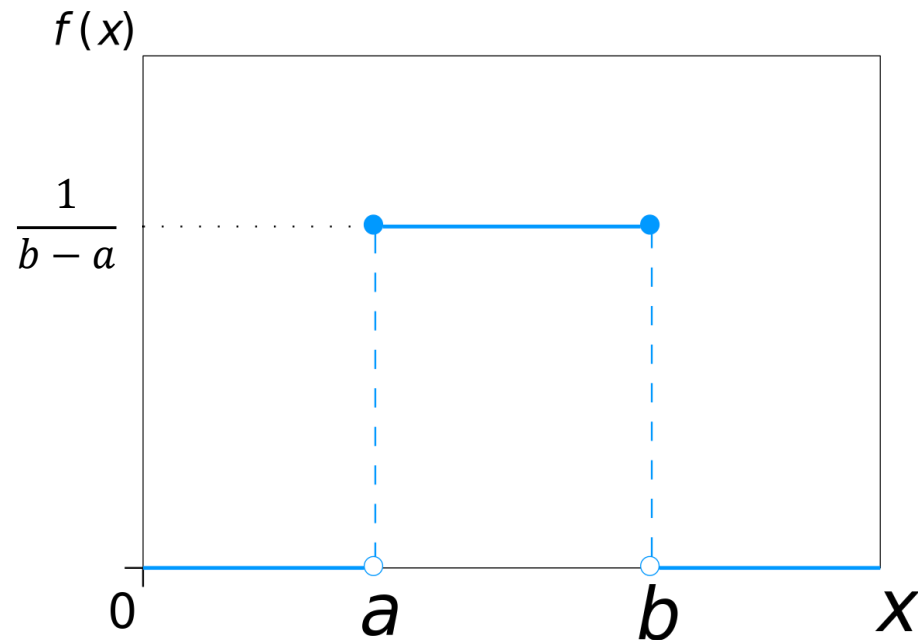


Initialization Methods

Uniform Distribution

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

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

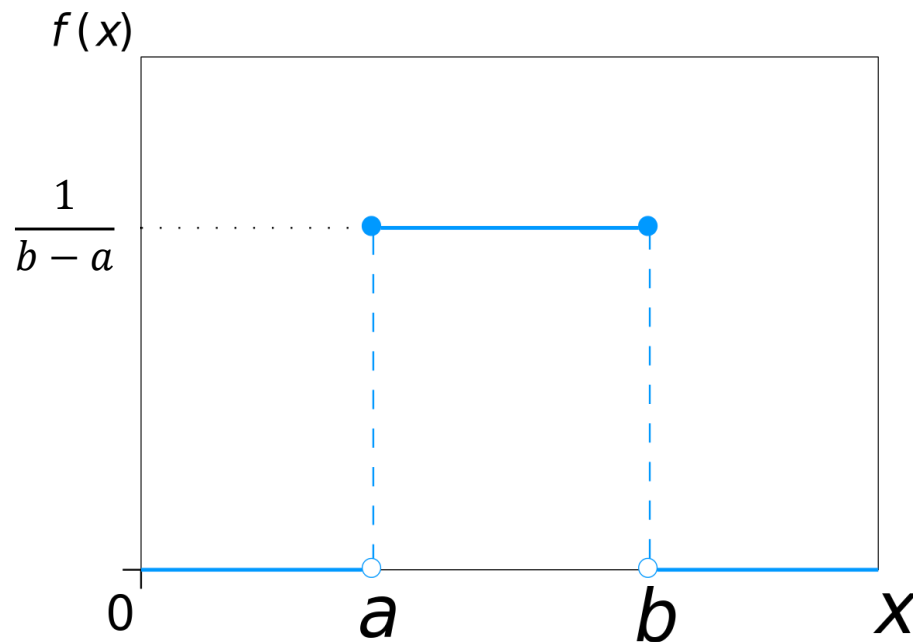


$$\begin{aligned} E[X] &= \int_{-\infty}^{\infty} x f(x) dx = \int_a^b x \frac{1}{b-a} dx \\ &= \frac{x^2}{2(b-a)} \Big|_a^b = \frac{b^2 - a^2}{2(b-a)} = \frac{a+b}{2} \end{aligned}$$

Initialization Methods

Uniform Distribution

$$\begin{aligned} X &\sim U(a, b) & E[X] &= \frac{a + b}{2} \\ f(x) &= \frac{1}{b - a} & \text{var}[X] &= \frac{(b - a)^2}{12} \end{aligned}$$



$$\begin{aligned} \text{var}[X] &= E\left((X - E(X))^2\right) = \int_{-\infty}^{\infty} (x - E(X))^2 f(x) dx \\ &= \int_a^b \left(x - \frac{a + b}{2}\right)^2 \frac{1}{b - a} dx \\ &= \frac{1}{b - a} \left[\int_a^b x^2 dx - \int_a^b 2x \frac{a + b}{2} dx \right] + \int_a^b \left(\frac{a + b}{2}\right)^2 dx \\ &= \frac{1}{b - a} \left[\frac{x^3}{3} \Big|_a^b - \frac{x^2(a + b)}{2} \Big|_a^b + \left(\frac{a + b}{2}\right)^2 x \Big|_a^b \right] \\ &= \frac{1}{b - a} \left[\frac{b^3 - a^3}{3} - \frac{(b^2 - a^2)(a + b)}{2} + \left(\frac{a + b}{2}\right)^2 (b - a) \right] \\ &= \frac{a^2 + ab + b^2}{3} - \frac{a^2 + 2ab + b^2}{2} + \frac{a^2 + 2ab + b^2}{4} \\ &= \frac{4(a^2 + ab + b^2) - 3(a^2 + 2ab + b^2)}{12} = \frac{(b - a)^2}{12} \end{aligned}$$

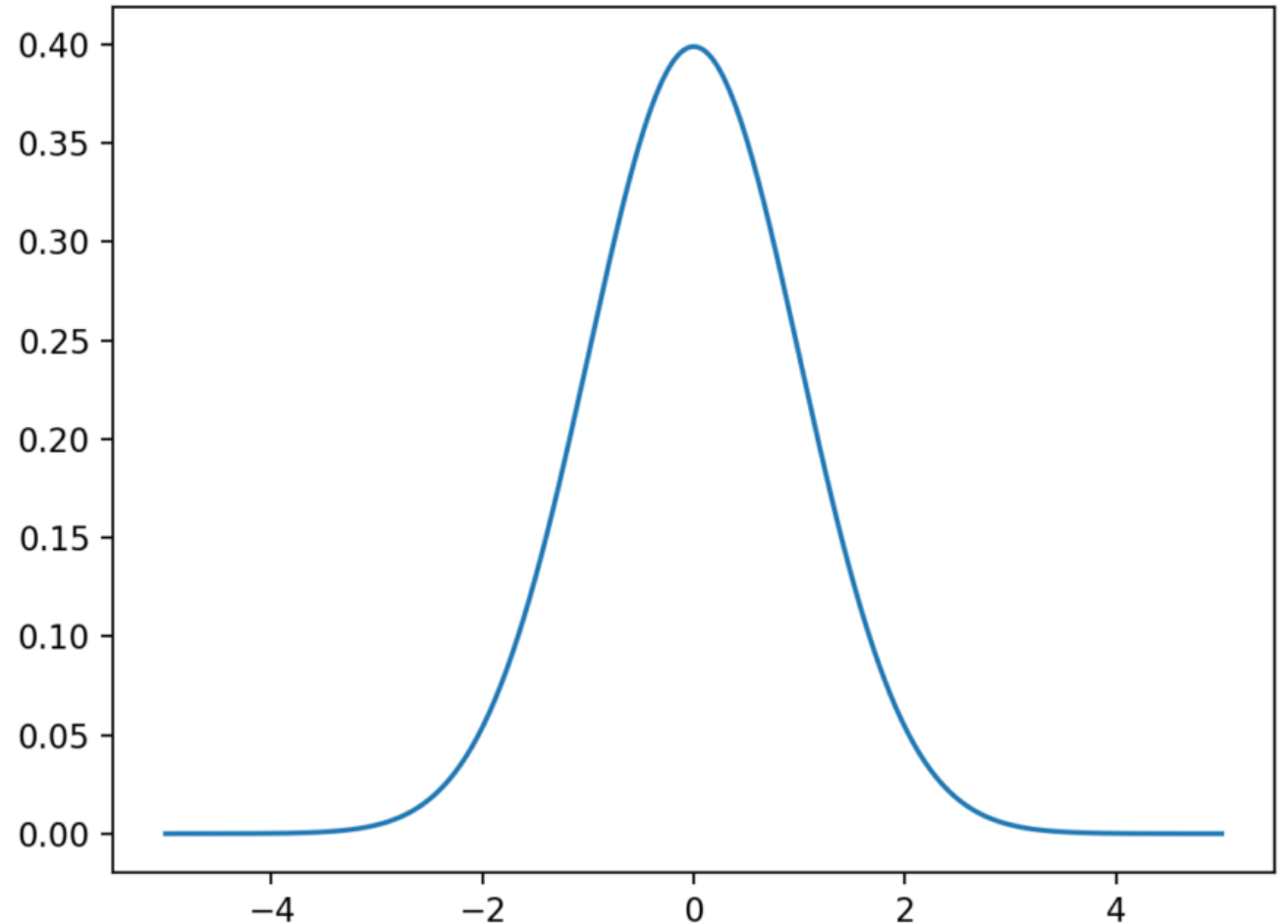
Initialization Methods

Xavier Initialization

Gaussian Distribution

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

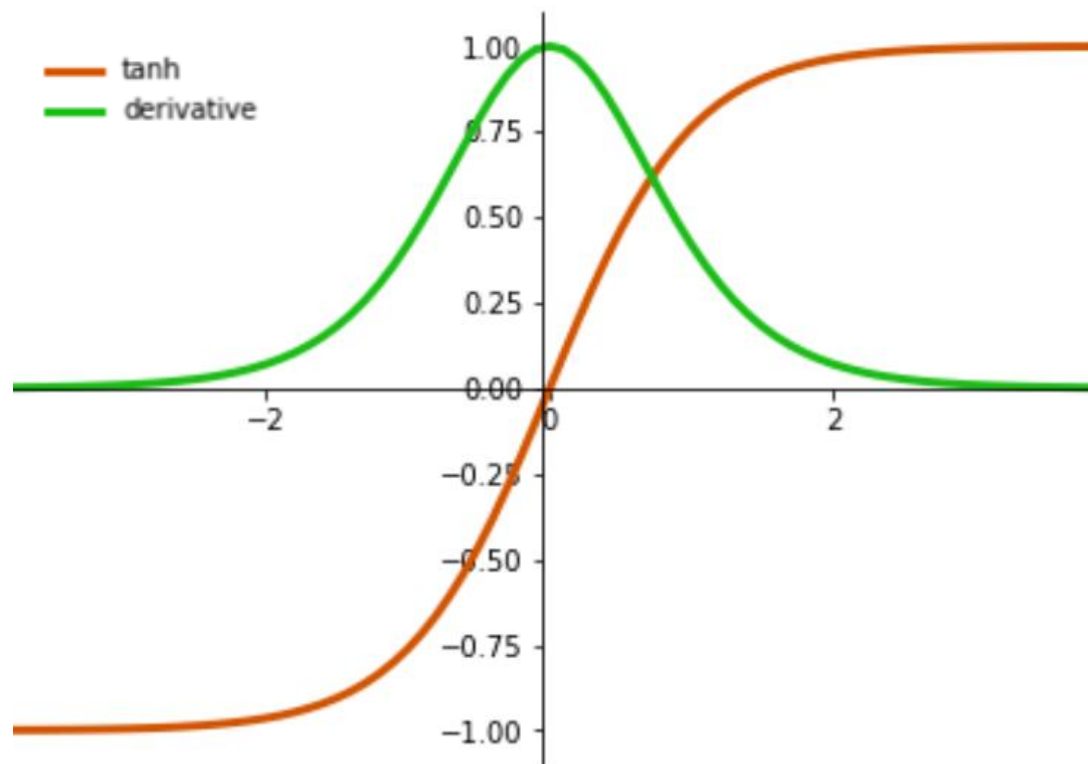
$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$



Maclaurin series

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

$$f(x) = \sum_{n=0}^{\infty} f^{(n)}(0) \frac{x^n}{n!}$$
$$= f(0) + f'(0)x + \frac{f''(0)}{2!}x^2 + \frac{f^{(3)}(0)}{3!}x^3 + \dots$$



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

$$\tanh(0) = 0$$

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

$$\begin{aligned}\tanh''(0) &= (1 - \tanh^2(0))' \\ &= -2\tanh(0)\tanh'(0) = 0\end{aligned}$$

$$\begin{aligned}\tanh^{(3)}(0) &= (-2\tanh(0)\tanh'(0))' \\ &= -2[\tanh'(0)\tanh'(0) + \tanh''(0)\tanh(0)] = -2\end{aligned}$$

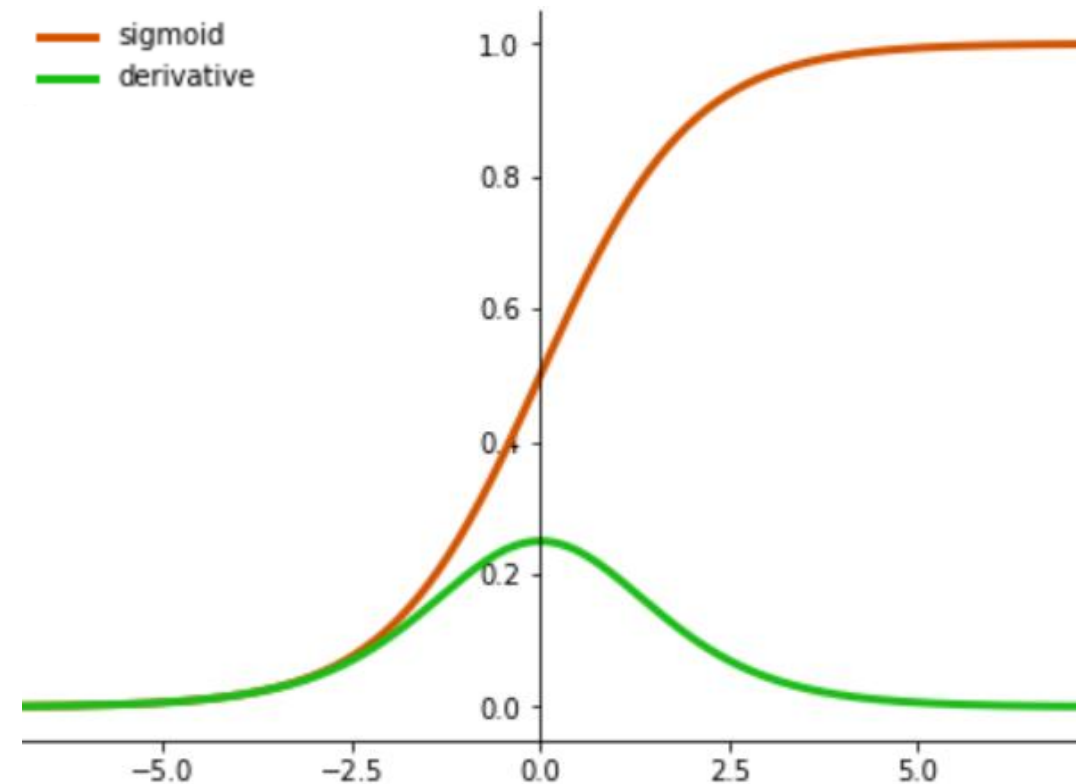
$$\begin{aligned}\tanh(x) &= f(0) + f'(0)x + \frac{f''(0)}{2!}x^2 + \frac{f^{(3)}(0)}{3!}x^3 + \dots \\ &= x - \frac{x^3}{3!} + \dots\end{aligned}$$

$$\Rightarrow \tanh(x) \approx x$$

Maclaurin series

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

$$f(x) = \sum_{n=0}^{\infty} f^{(n)}(0) \frac{x^n}{n!}$$
$$= f(0) + f'(0)x + \frac{f''(0)}{2!}x^2 + \frac{f^{(3)}(0)}{3!}x^3 + \dots$$



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

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

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

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

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

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

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

$$\Rightarrow \text{sigmoid}(x) \approx \frac{1}{2} + \frac{x}{4}$$

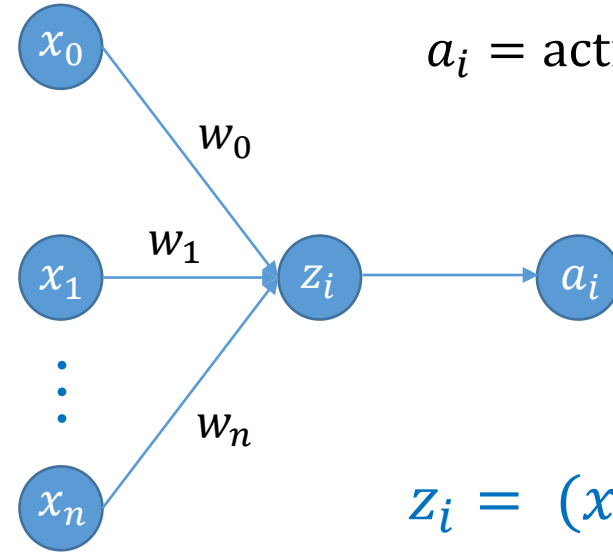
Initialization Methods

Xavier Initialization

$$E(XY) = E(X)E(Y)$$
$$\text{var}(XY) = \text{var}(X)\text{var}(Y) + \text{var}(X)(E(Y))^2 + \text{var}(Y)(E(X))^2$$

Uniform Distribution

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



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

$$E(X) = 0$$

$$E(W) = 0$$

$$b = 0$$

$$z_i = (x_1w_1 + \dots + x_nw_n + b)$$

$$\begin{aligned} \text{var}(z_i) &= \text{var}(x_1w_1 + \dots + x_nw_n + b) \\ &= n\text{var}(x_iw_i) = n\text{var}(x_i)\text{var}(w_i) \end{aligned}$$

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

$$\begin{aligned} \text{var}(X) \approx \text{var}(\mathbf{a}) &\xrightarrow{\text{iid}} \text{var}(x_i) \approx \text{var}(a_i) \rightarrow n\text{var}(w_i) \approx 1 \\ &\rightarrow \text{var}(w_i) \approx \frac{1}{n} \end{aligned}$$

Initialization Methods

Xavier Initialization

activation = tanh

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

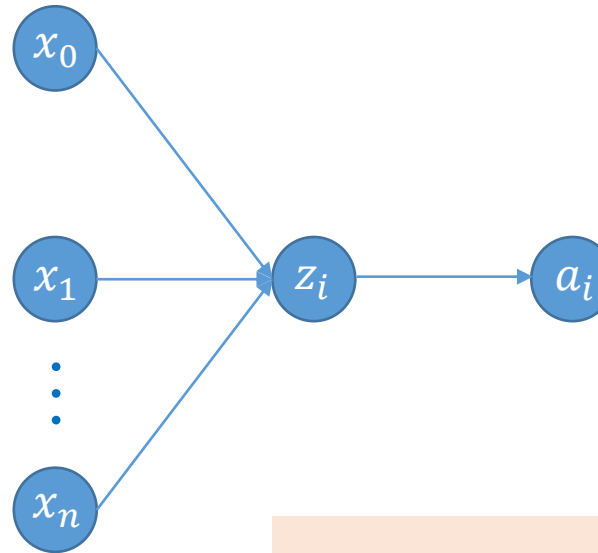
$$\begin{aligned} \text{var}(XY) = & \text{var}(X)\text{var}(Y) + \\ & \text{var}(X)(E(Y))^2 + \\ & \text{var}(Y)(E(X))^2 \end{aligned}$$

Uniform Distribution

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

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

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



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

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

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

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

Initialization Methods

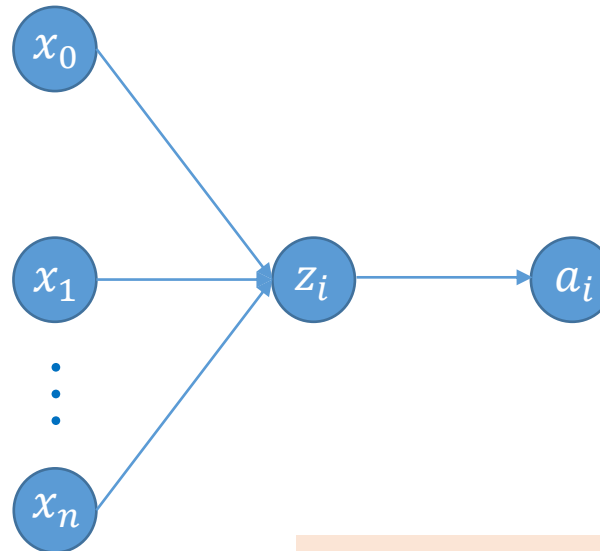
Xavier Initialization

activation = tanh

$$E(XY) = E(X)E(Y)$$
$$\text{var}(XY) = \text{var}(X)\text{var}(Y) + \text{var}(X)(E(Y))^2 + \text{var}(Y)(E(X))^2$$

Gaussian Distribution

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



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

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

$$\sigma^2 = \frac{1}{n} \quad \rightarrow \quad \sigma = \frac{1}{\sqrt{n}}$$

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

Initialization Methods

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

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

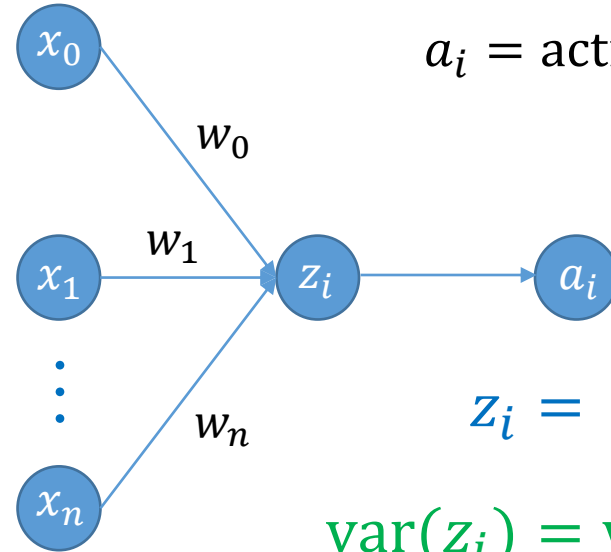
$$\begin{aligned} \text{var}(XY) = & \text{var}(X)\text{var}(Y) + \\ & \text{var}(X)(E(Y))^2 + \\ & \text{var}(Y)(E(X))^2 \end{aligned}$$

Uniform Distribution

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

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

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



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

$$E(X) = 0$$

$$E(W) = 0$$

$$b = 0$$

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

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

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

$$\rightarrow 16 \text{var}(a_i) = \text{var}(z_i)$$

$$\begin{aligned} \text{var}(X) \approx \text{var}(\mathbf{a}) &\xrightarrow{\text{iid}} \text{var}(x_i) \approx \text{var}(a_i) \rightarrow n \text{var}(w_i) \approx 16 \\ &\rightarrow \text{var}(w_i) \approx \frac{16}{n} \end{aligned}$$

Initialization Methods

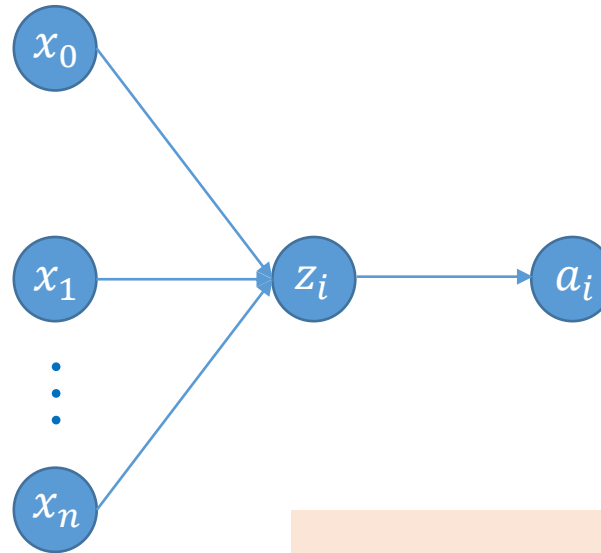
Xavier Initialization

activation = sigmoid

$$E(XY) = E(X)E(Y)$$
$$\text{var}(XY) = \text{var}(X)\text{var}(Y) + \text{var}(X)(E(Y))^2 + \text{var}(Y)(E(X))^2$$

Uniform Distribution

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



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

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

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

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

Initialization Methods

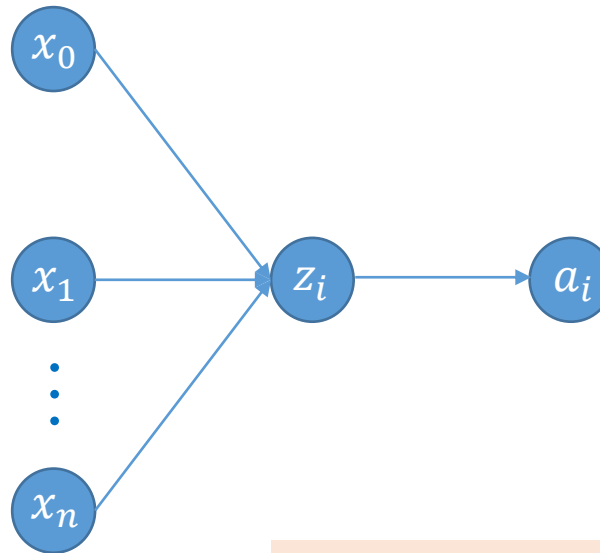
Xavier Initialization

activation = sigmoid

$$E(XY) = E(X)E(Y)$$
$$\text{var}(XY) = \text{var}(X)\text{var}(Y) + \text{var}(X)(E(Y))^2 + \text{var}(Y)(E(X))^2$$

Gaussian Distribution

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



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

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

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

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

Initialization Methods

Kaiming He Initialization

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

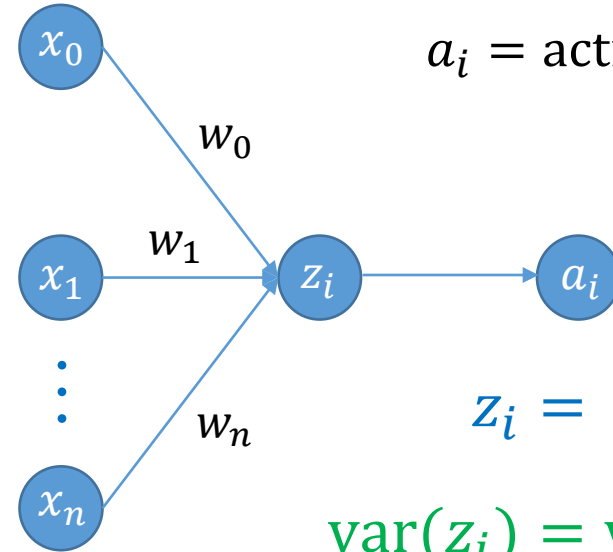
$$\text{var}(XY) = \text{var}(X)\text{var}(Y) + \text{var}(X)(E(Y))^2 + \text{var}(Y)(E(X))^2$$

Uniform Distribution

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

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

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



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

$$E(X) = 0$$

$$E(W) = 0$$

$$b = 0$$

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

$$\begin{aligned} \text{var}(z_i) &= \text{var}(x_1 w_1 + \dots + x_n w_n + b) \\ &= n \text{var}(x_i w_i) = n \text{var}(x_i) \text{var}(w_i) \end{aligned}$$

$$\text{activation} = \text{relu} \quad \rightarrow \quad a_i = \max(0, z_i)$$

$$\rightarrow 2\text{var}(a_i) = \text{var}(z_i)$$

$$\text{var}(X) \approx \text{var}(\mathbf{a}) \xrightarrow{\text{iid}} \text{var}(x_i) \approx \text{var}(a_i) \rightarrow n \text{var}(w_i) \approx 2$$

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

Initialization Methods

He Initialization

activation = he

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

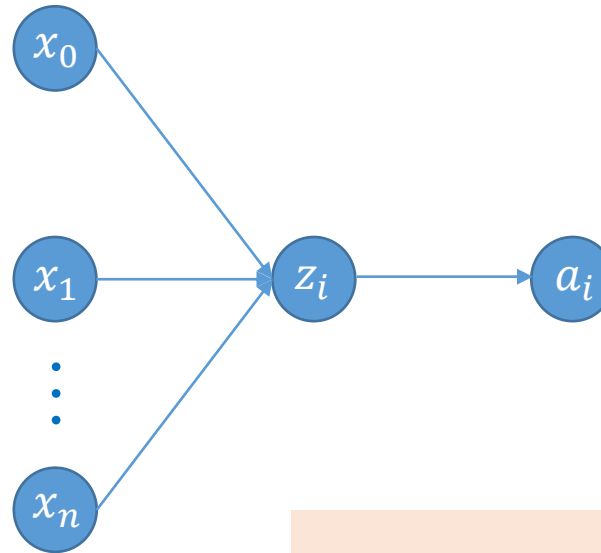
$$\begin{aligned} \text{var}(XY) = & \text{var}(X)\text{var}(Y) + \\ & \text{var}(X)(E(Y))^2 + \\ & \text{var}(Y)(E(X))^2 \end{aligned}$$

Uniform Distribution

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

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

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



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

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

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

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

Initialization Methods

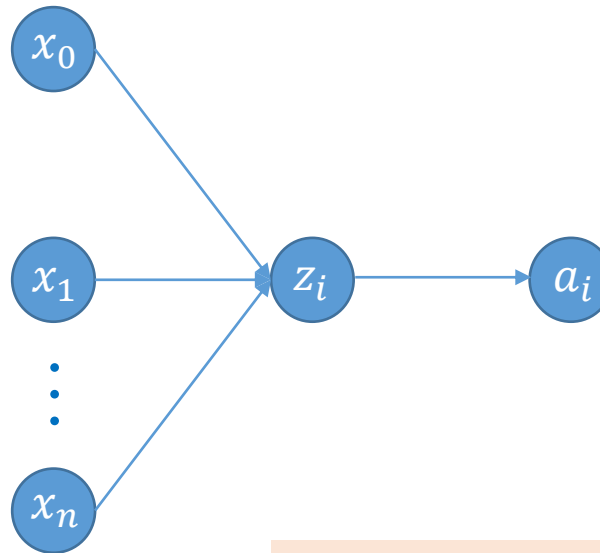
He Initialization

activation = he

$$E(XY) = E(X)E(Y)$$
$$\text{var}(XY) = \text{var}(X)\text{var}(Y) + \text{var}(X)(E(Y))^2 + \text{var}(Y)(E(X))^2$$

Gaussian Distribution

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



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

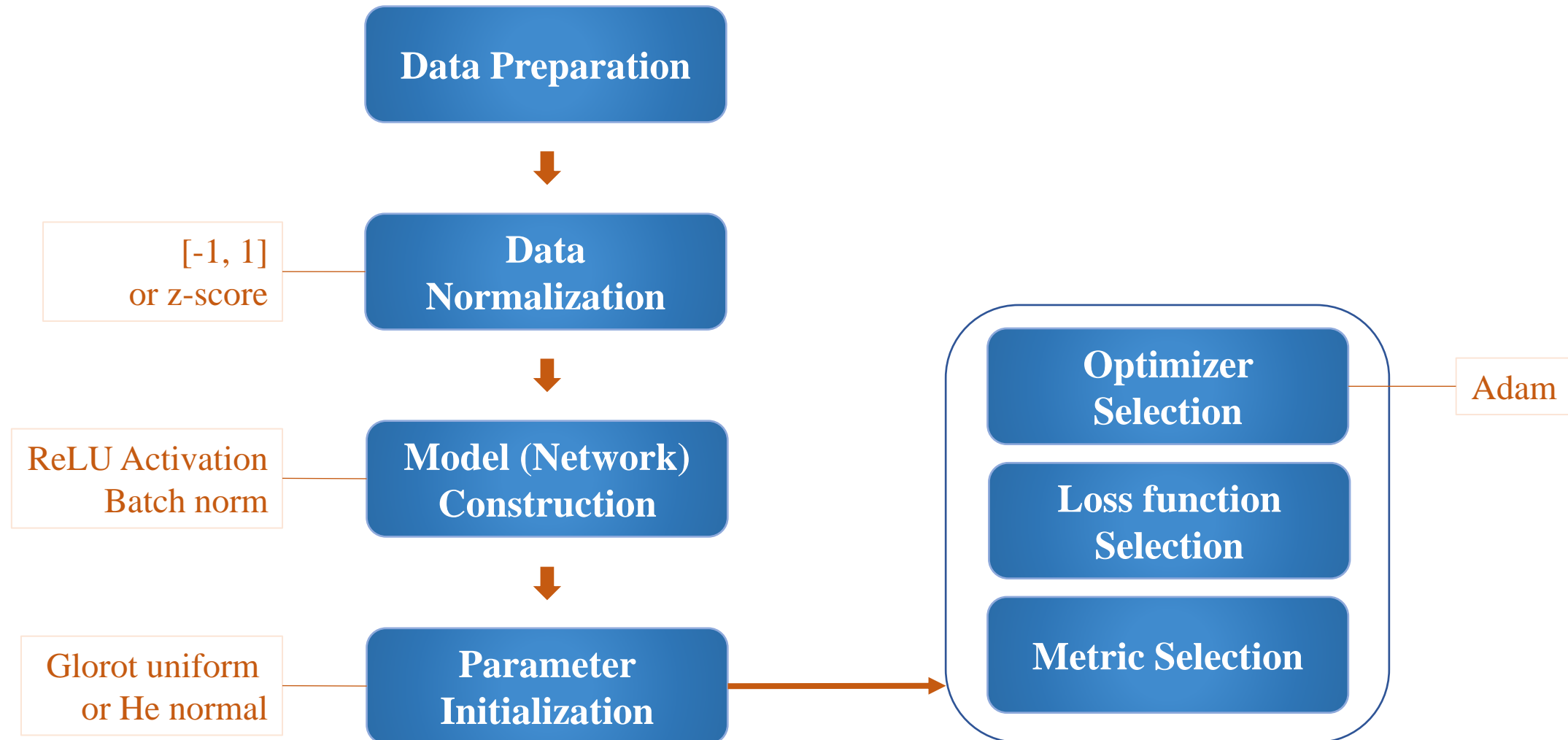
$$w_i \sim N(0, \sigma^2)$$

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

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

Summary

Recommendation



Further Reading

Dying ReLU

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

Initialization

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

