# Insight into Multi-layer Perceptron

Quang-Vinh Dinh
Ph.D. in Computer Science

*Year 2024*  code&data

# Objectives



## MLP Insight

Data Normalization

Model (Network) Construction

Parameter Initialization

Optimizer Selection | Metric Selection

Loss function Selection

## MLP Examples

$h_1$ — ReLU — $w_{11}$ — $z_1$

$w_{12}$ — $bw_1$

$h_2$ — ReLU — $w_{21}$ — $z_2$

$w_{22}$ — $bw_2$

1 — $w_{31}$ — $w_{32}$ — $z_3$

$bw_3$

## Init. Examples

$\boldsymbol{x} = [1.4]$    $x$

$b_0$   $w_0$   $b_1$   $w_1$

0.0   0.0   0.0   0.0

$z_0 = w_0 x + b_0$    $z_1 = w_1 x + b_1$    $\boldsymbol{z_1} = [0.0]$

$\hat{y}_0 = \dfrac{e^{z_0}}{\sum_{i=0}^{1} e^{z_i}}$    $\hat{y}_1 = \dfrac{e^{z_1}}{\sum_{i=0}^{1} e^{z_i}}$    $\boldsymbol{\hat{y}_1} = [0.5]$

$\mathrm{L} = -y_0 \log \hat{y}_0 - y_1 \log \hat{y}_1$    $y$    $\boldsymbol{y} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$

$L = [-\log 0.5] = [0.693]$

# Outline

**Data Normalization**

↓

**Model (Network) Construction**

↓

**Parameter Initialization**

↓

**Optimizer Selection**

**Metric Selection**

**Loss function Selection**

# Data Normalization

**Data Preparation**

↓

**Data Normalization**

↓

**Model (Network) Construction**

↓

**Parameter Initialization**

↓

**Optimizer Selection** | **Metric Selection**

**Loss function Selection**

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

$$X \in [0, 1]$$

**Convert to the range [0,1]**

$$\text{Image} = \frac{\text{Image}}{255}$$

**Convert to the range [-1,1]**

$$\text{Image} = \frac{\text{Image}}{127.5} - 1$$

**Z-score normalization**

$$\text{Image} = \frac{\text{Image} - \mu}{\sigma}$$

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

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

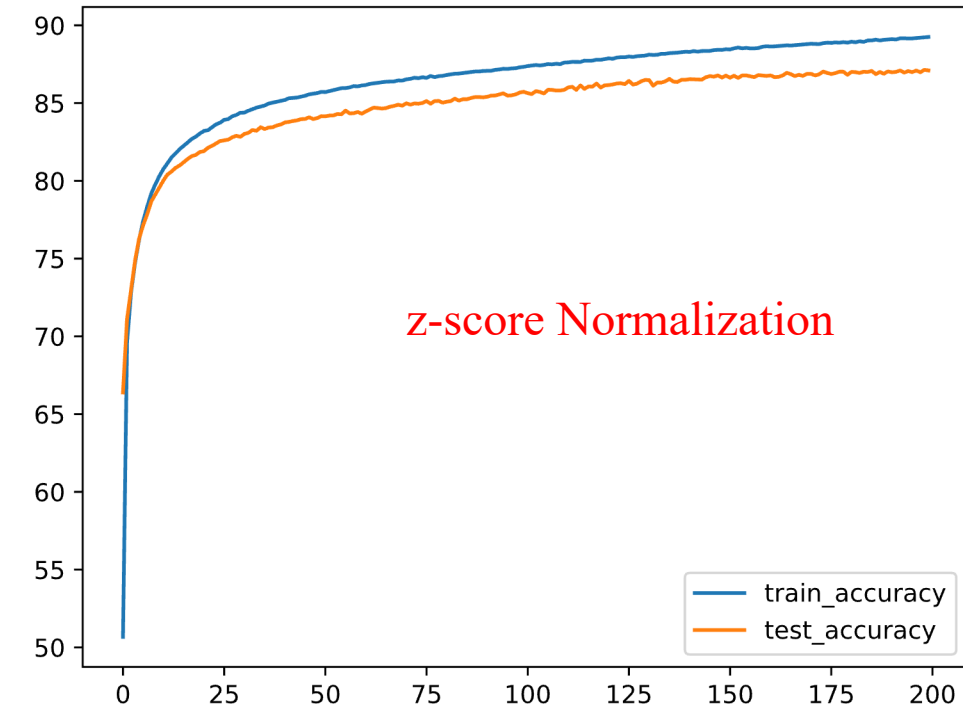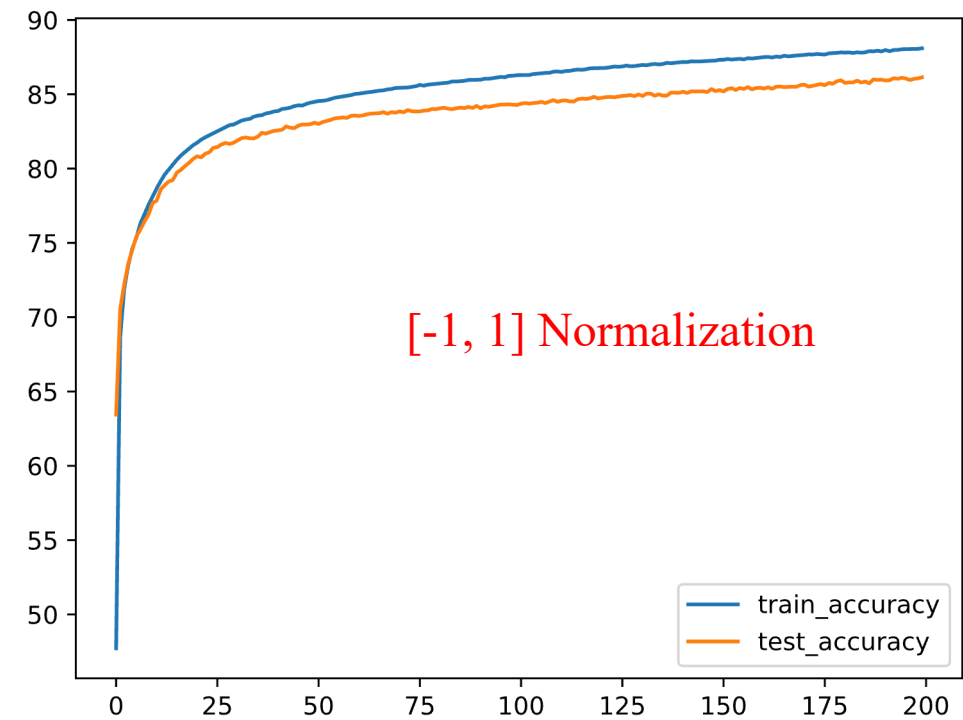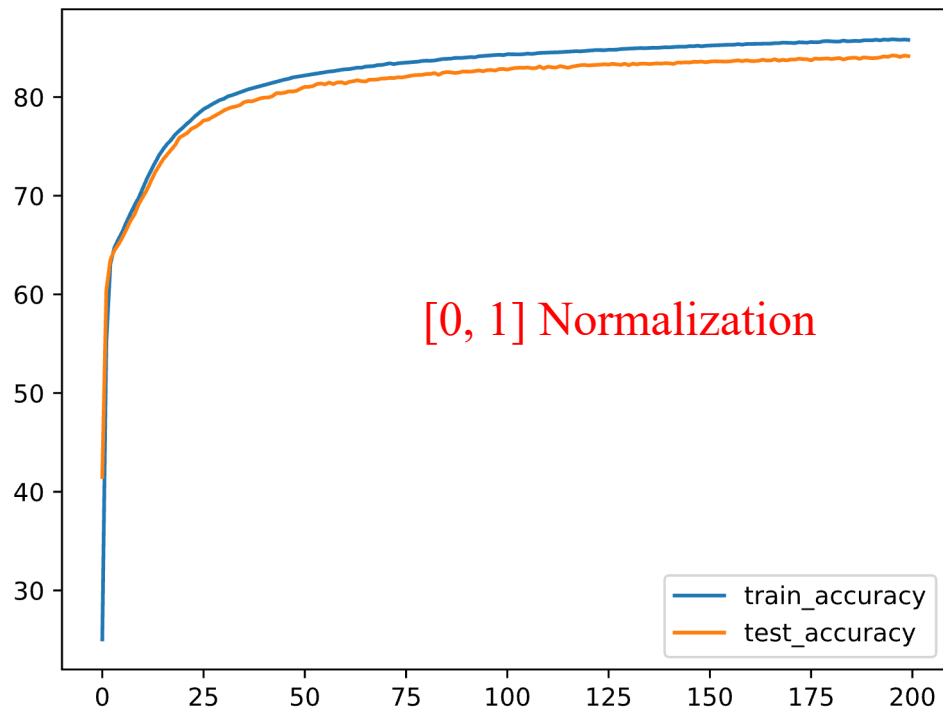| | |
|---|---|
| [0,1] | mean = 0 ; std = 1 |
| [-1,1] | mean = 0.5; std = 0.5 |

Compute mean and std from data

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

trainset = torchvision.datasets.FashionMNIST(root='data', train=True, download=True, transform=transform)
trainloader = torch.utils.data.DataLoader(trainset, batch_size=1024, num_workers=10, shuffle=True)

testset = torchvision.datasets.FashionMNIST(root='data', train=False, download=True, transform=transform)
testloader = torch.utils.data.DataLoader(testset, batch_size=1024, num_workers=10, shuffle=False)
```

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

trainset = torchvision.datasets.FashionMNIST(root='data', train=True, download=True, transform=transform)
trainloader = torch.utils.data.DataLoader(trainset, batch_size=1024, num_workers=10, shuffle=True)

testset = torchvision.datasets.FashionMNIST(root='data', train=False, download=True, transform=transform)
testloader = torch.utils.data.DataLoader(testset, batch_size=1024, num_workers=10, shuffle=False)
```
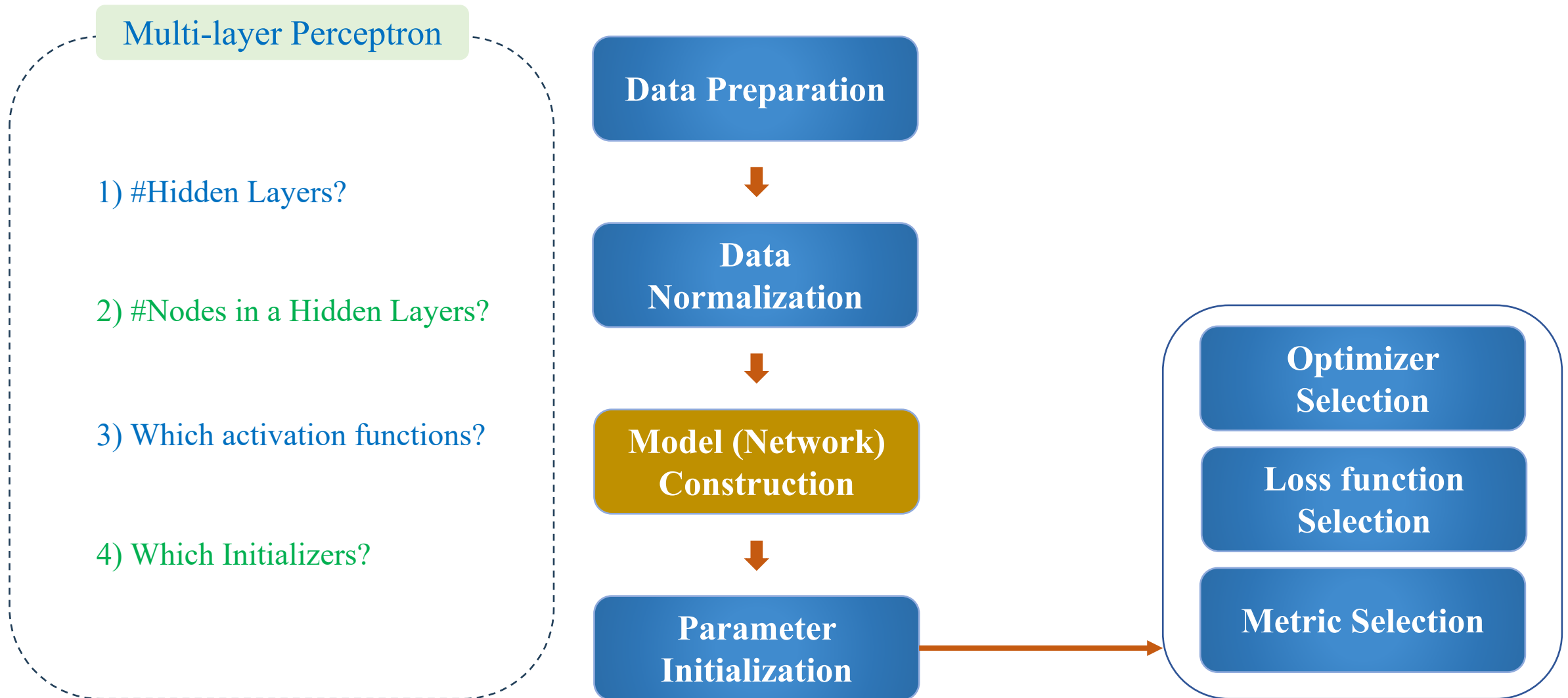
```python
# computed mean and std in advance
transform = transforms.Compose([transforms.ToTensor(), transforms.Normalize((mean,), (std,))])

trainset = torchvision.datasets.FashionMNIST(root='data', train=True, download=True, transform=transform)
trainloader = torch.utils.data.DataLoader(trainset, batch_size=1024, num_workers=10, shuffle=True)

testset = torchvision.datasets.FashionMNIST(root='data', train=False, download=True, transform=transform)
testloader = torch.utils.data.DataLoader(testset, batch_size=1024, num_workers=10, shuffle=False)
```

(a) [0, 1] Normalization     (b) [-1, 1] Normalization     (c) z-score Normalization

# Data Normalization

```python
model = nn.Sequential(
    nn.Flatten(), nn.Linear(784, 256),
    nn.ReLU(), nn.Linear(256, 10)
)
criterion = nn.CrossEntropyLoss()
optimizer = optim.SGD(model.parameters(),
                      lr=0.01)
```
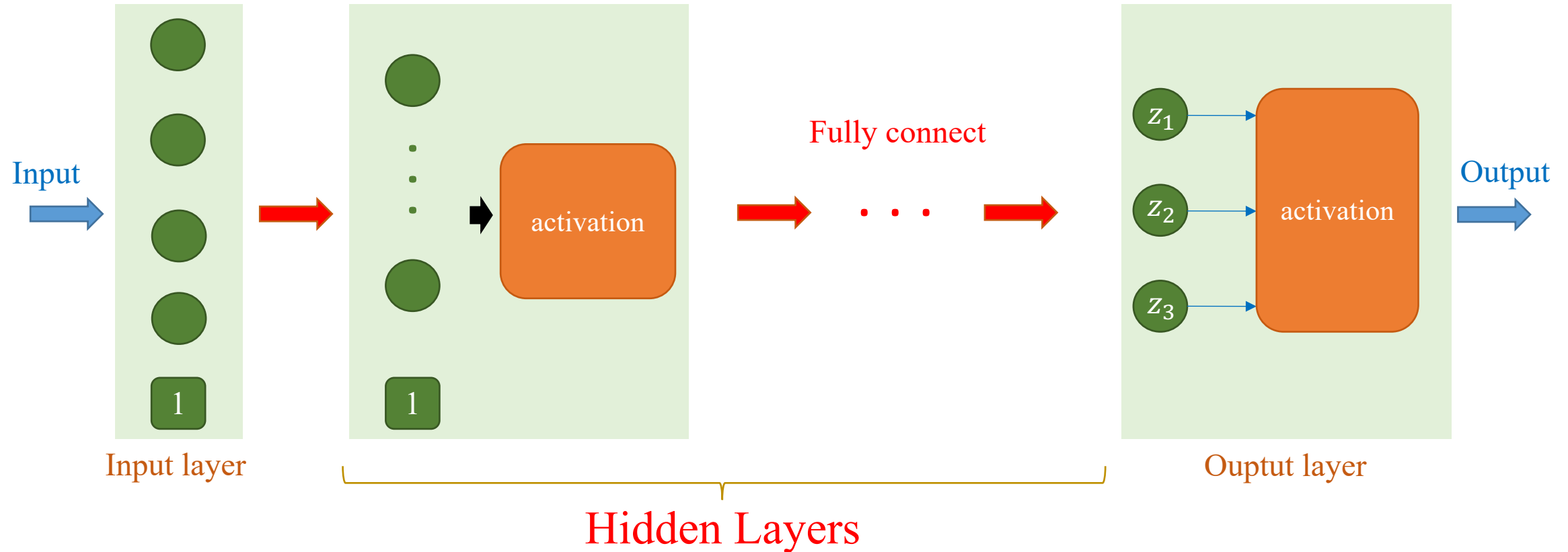


[-1, 1] Normalization

[0, 1] Normalization

z-score Normalization

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

❖ **Model (Network) Construction**



Input

Input layer

activation

Fully connect

$z_1$
$z_2$
$z_3$
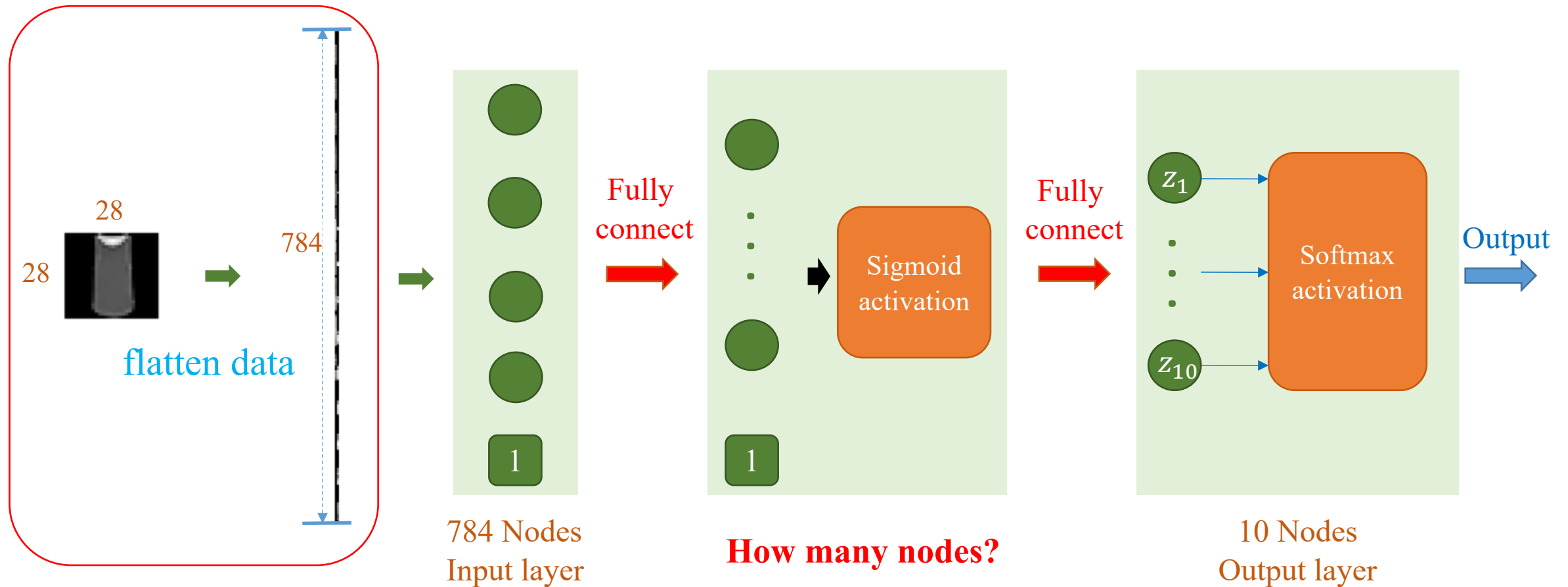activation

Output

Ouptut layer

Hidden Layers

How many hidden layers?
How many nodes in a hidden layer?

Which activation function?
Which network components?

6

# How many nodes?

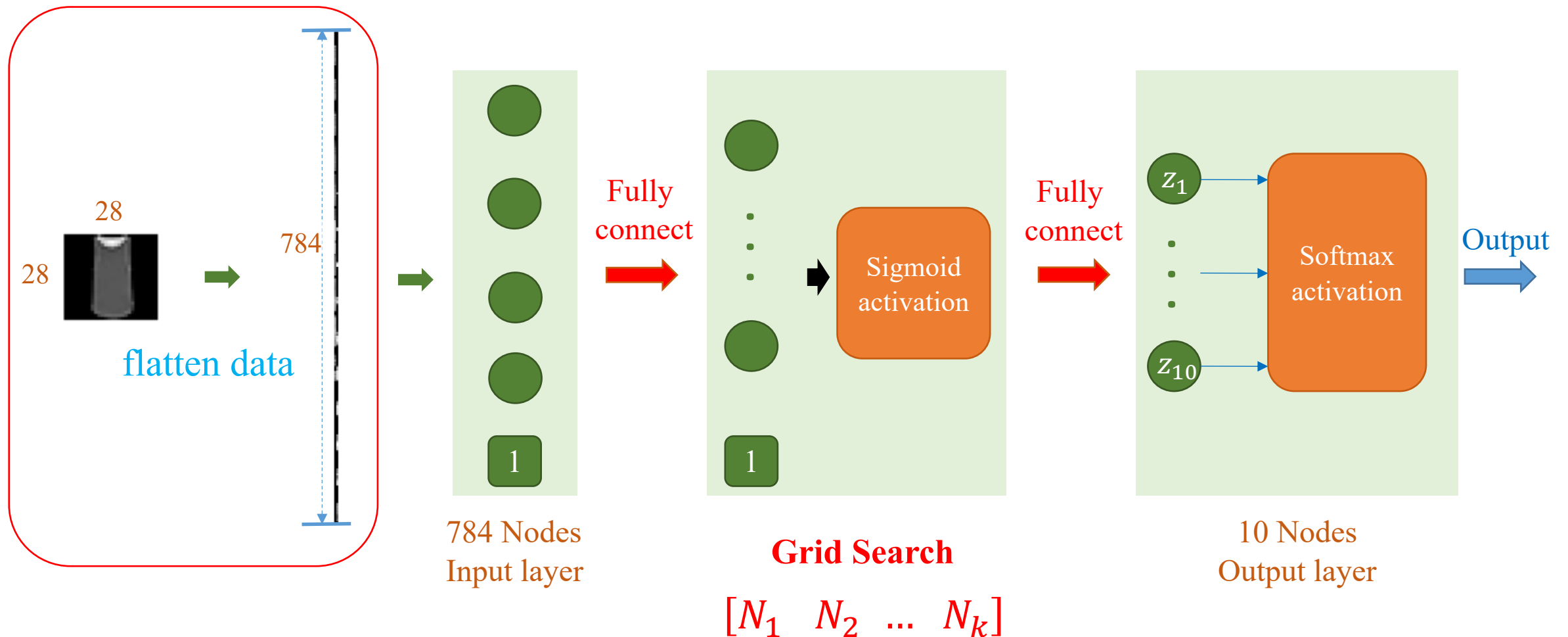❖ **Model (Network) Construction**

784 Nodes
Input layer

How many nodes?

10 Nodes
Output layer

❖ **Model (Network) Construction**



784 Nodes
Input layer

**Grid Search**

$$[N_1 \quad N_2 \quad ... \quad N_k]$$

10 Nodes
Output layer

# How many nodes?

[-1, 1] Normalization

Cross-entropy Loss

SGD with lr=0.01



256 nodes

Train-Acc: 90%

Test-Acc: 87%

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



64 nodes

Train-Acc: 89%

Test-Acc: 86%

```python
model = nn.Sequential(
    nn.Flatten(),
    nn.Linear(784, 64),
    nn.ReLU(),
    nn.Linear(64, 10)
)
```
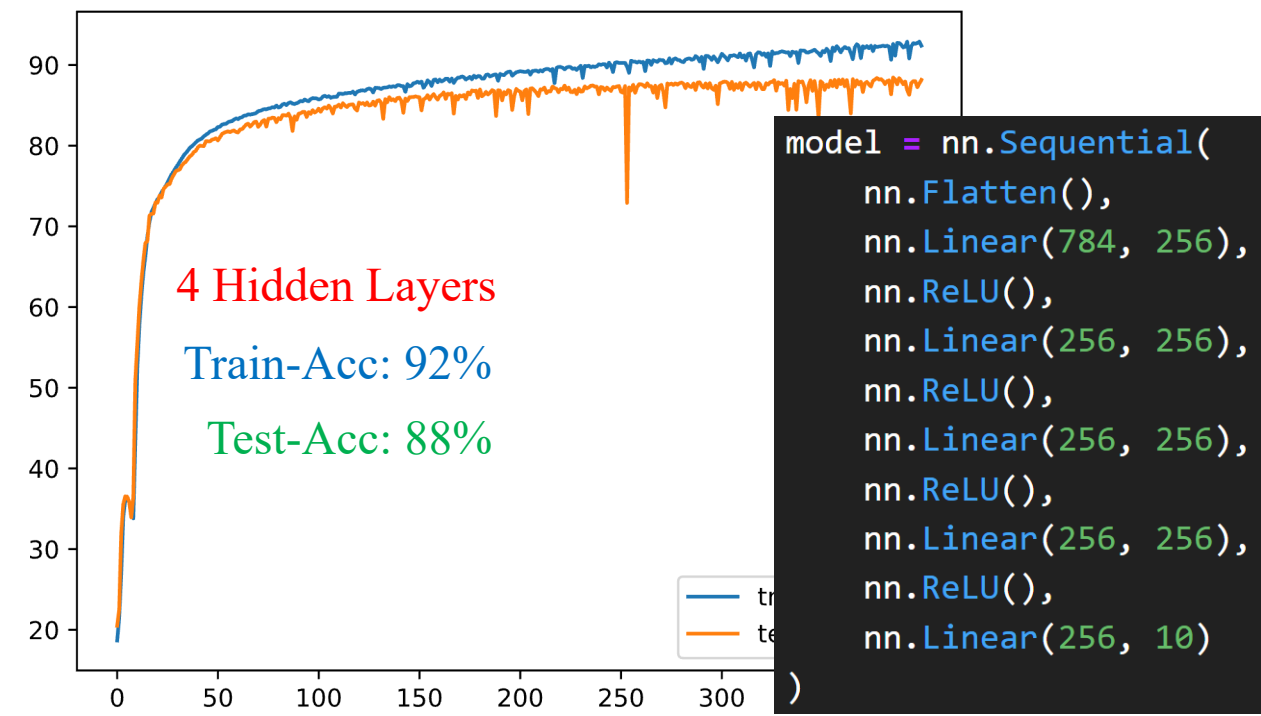


1024 nodes

Train-Acc: 90%

Test-Acc: 87%

```python
model = nn.Sequential(
    nn.Flatten(),
    nn.Linear(784, 1024),
    nn.ReLU(),
    nn.Linear(1024, 10)
)
```

**1 Hidden Layer**

Train-Acc: 90%

Test-Acc: 87%

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

train_accuracy
test_accuracy

**3 Hidden Layers**

Train-Acc: 92%

Test-Acc: 88%

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

**2 Hidden Layers**

Train-Acc: 91%

Test-Acc: 88%

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

train_accuracy
test_accuracy

**4 Hidden Layers**

Train-Acc: 92%

Test-Acc: 88%

```python
model = nn.Sequential(
    nn.Flatten(),
    nn.Linear(784, 256),
    nn.ReLU(),
    nn.Linear(256, 256),
    nn.ReLU(),
    nn.Linear(256, 256),
    nn.ReLU(),
    nn.Linear(256, 256),
    nn.ReLU(),
    nn.Linear(256, 10)
)
```
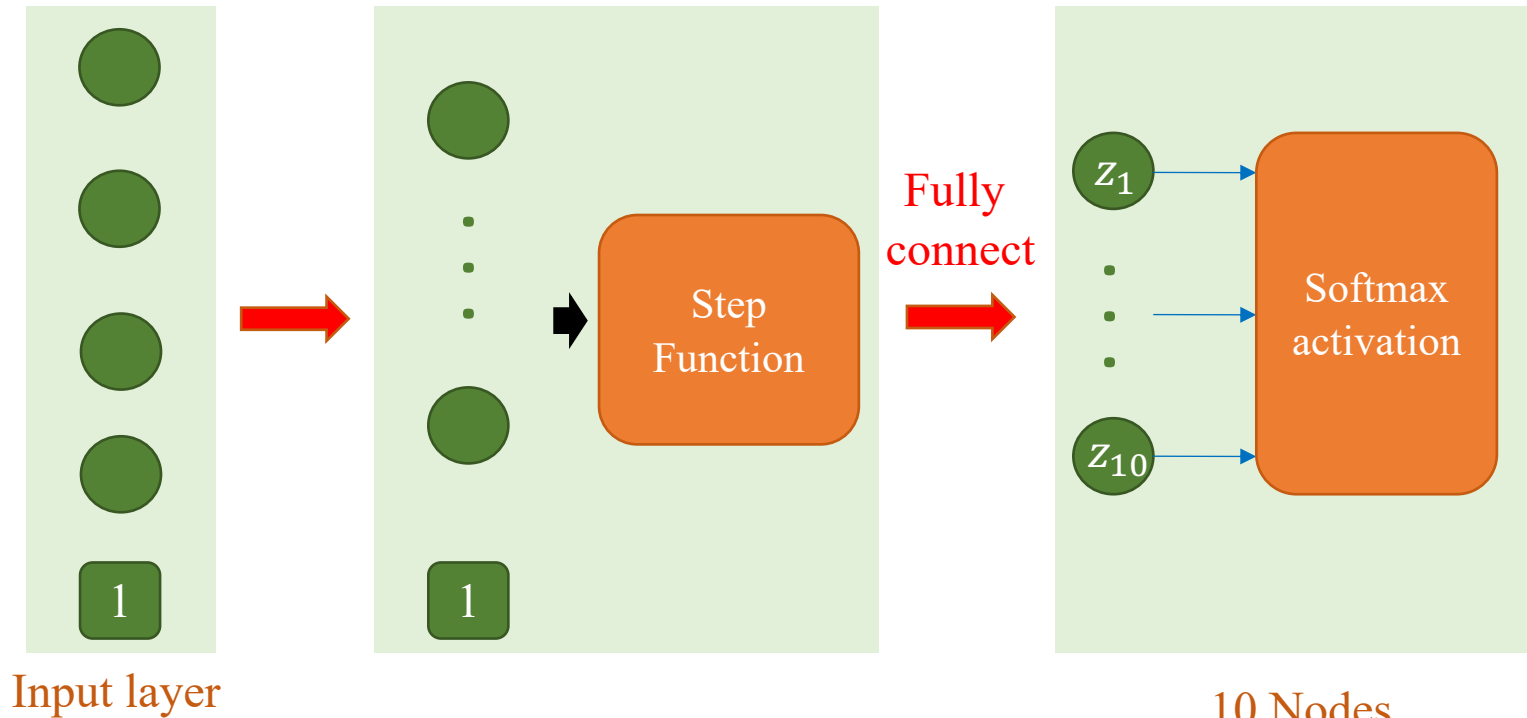
# Activation Functions

❖ **Model (Network) Construction**
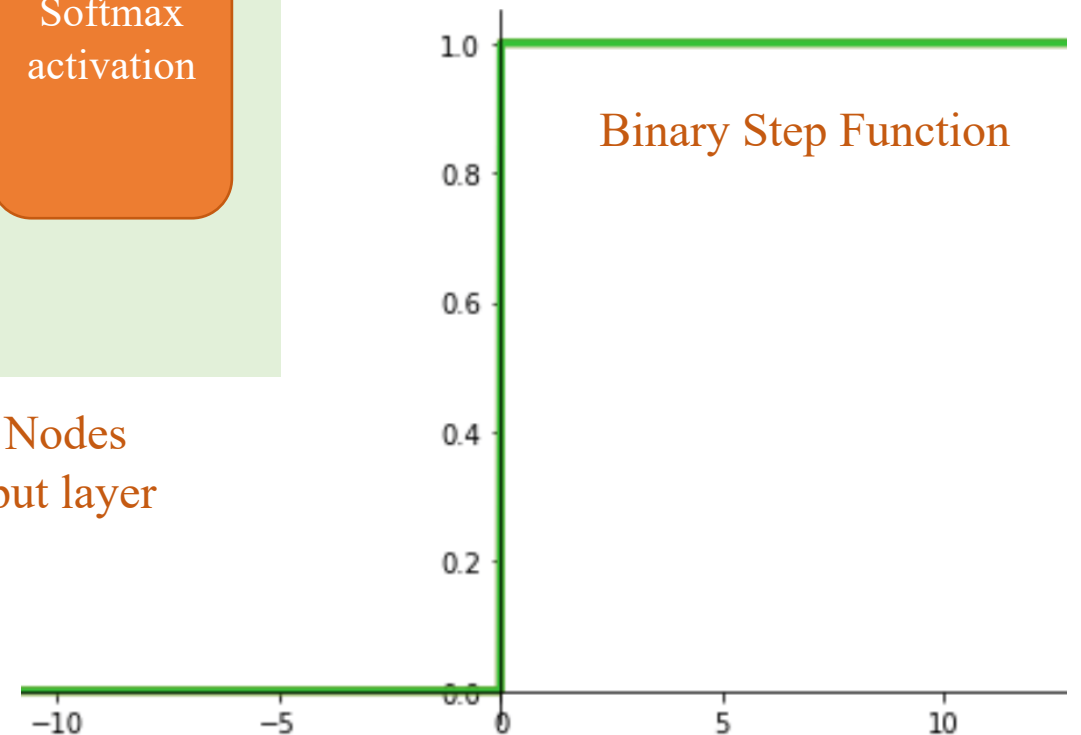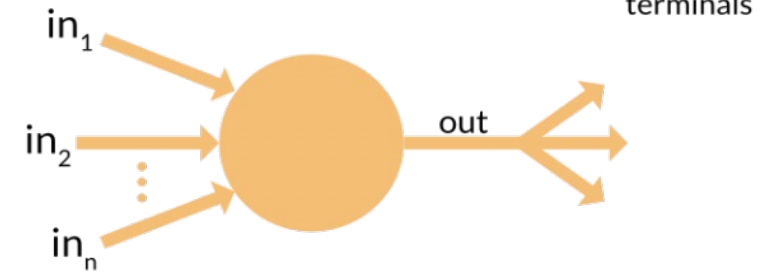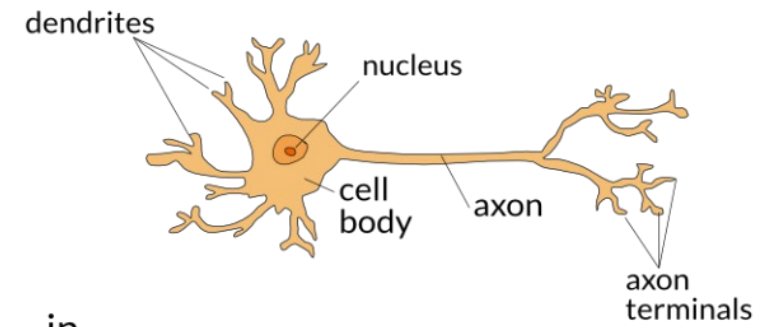
Which activation function?

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

2010
$$\text{ReLU}(x) = \begin{cases} 0 & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$$

2017
$$SELU(x) = \begin{cases} \lambda x & \text{if } x \geq 0 \\ \lambda \alpha(e^x - 1) & \text{if } x < 0 \end{cases}$$

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

2015
$$\text{ELU}(x) = \begin{cases} \alpha(e^x - 1) & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$$

$$\lambda \approx 1.0507$$
$$\alpha \approx 1.6733$$

2001
$$\text{softplus}(x) = \log(1 + e^x)$$

2015
$$\text{PReLU}(x) = \begin{cases} \alpha x & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$$

2017
$$swish1(x) = x * \frac{1}{1 + e^{-x}}$$

2023
$$\text{GELU}(x) = x\phi(x) \approx x * \text{sigmoid}(1.702x)$$



11

# Activation Functions

❖ **Step function**



Input layer

Step Function

Fully connect

$z_1$

$z_{10}$

Softmax activation

10 Nodes
Output layer

dendrites
nucleus
cell body
axon
axon terminals

$in_1$
$in_2$
$in_n$
out

Binary Step Function

$$f(x) = \begin{cases} 0 & if \ \ x < 0 \\ 1 & if \ \ x \geq 0 \end{cases}$$

# Activation Functions

❖ **Sigmoid function**

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

**data** =

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

**data_a** = **sigmoid**(data)

**data_a** =

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



$$\text{sigmoid}'(x) = \text{sigmoid}(x)\,(1 - \text{sigmoid}(x))$$

13

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

$$\text{sigmoid}'(x) = \left(\frac{1}{1 + e^{-x}}\right)' = \frac{-1}{(1 + e^{-x})^2}(-e^{-x})$$

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

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

$$= \frac{1}{1 + e^{-x}}\left(1 - \frac{1}{1 + e^{-x}}\right)$$

$$= \text{sigmoid}(x)\,(1 - \text{sigmoid}(x))$$

14

# Activation Functions

❖ **Tanh function**

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

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

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



data = | 1 | 5 | -4 | 3 | -2 |

data_a = **tanh**(data)

data_a = | 0.761 | 0.999 | -0.999 | 0.995 | -0.964 |

$$tanh'(x) = 1 - tanh^2(x)$$

15

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

$$\tanh'(x) = \left(\frac{e^x - e^{-x}}{e^x + e^{-x}}\right)' = \frac{(e^x + e^{-x})(e^x + e^{-x}) - (e^x - e^{-x})(e^x - e^{-x})}{(e^x + e^{-x})^2}$$

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

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

16

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

$$tanh'(x) = \left(\frac{2}{e^{-2x} + 1} - 1\right)' = \frac{4e^{-2x}}{(e^{-2x} + 1)^2} = 4\left(\frac{e^{-2x} + 1 - 1}{(e^{-2x} + 1)^2}\right)$$

$$= 4\left(\frac{1}{e^{-2x} + 1} - \frac{1}{(e^{-2x} + 1)^2}\right) = -\left(\frac{4}{(e^{-2x} + 1)^2} - \frac{4}{e^{-2x} + 1}\right)$$

$$= -\left(\frac{4}{(e^{-2x} + 1)^2} - \frac{4}{e^{-2x} + 1} + 1 - 1\right) = 1 - \left(\frac{2}{e^{-2x} + 1} - 1\right)^2 = 1 - tanh^2(x)$$

# Activation Functions

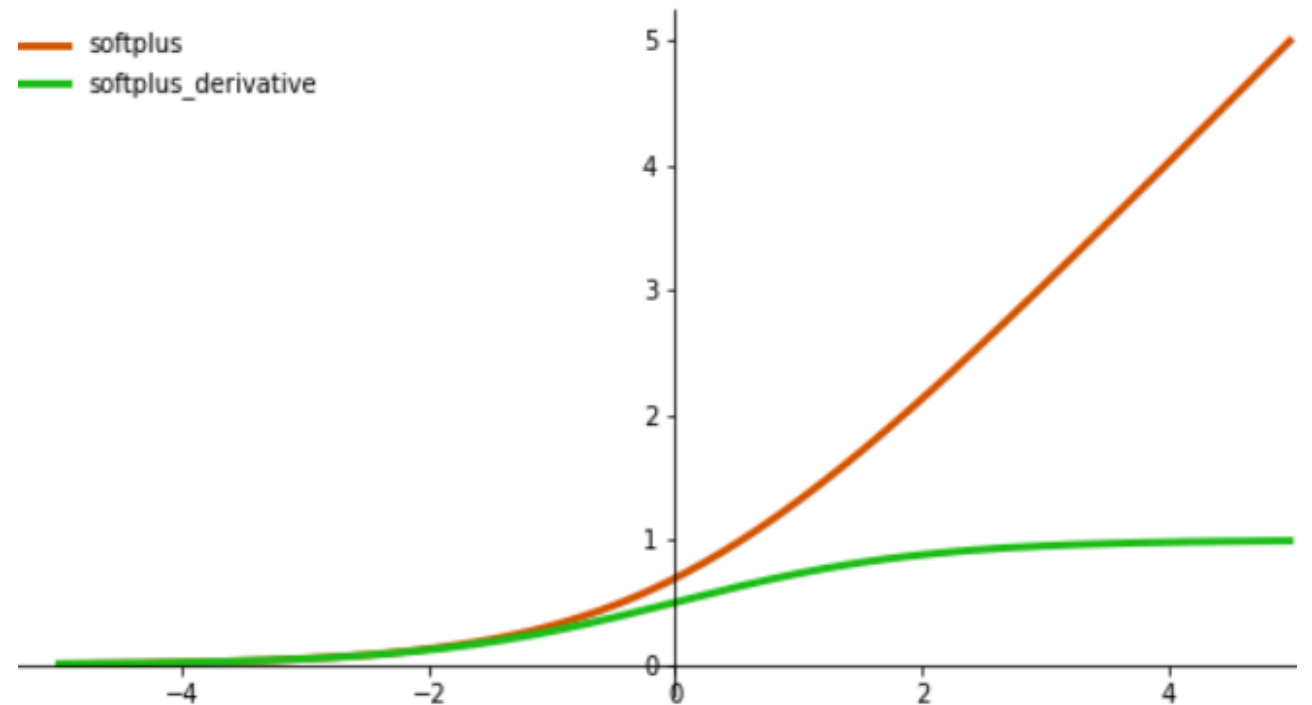❖ **Softplus function**

$$\text{softplus}(x) = \log(1 + e^x)$$

**data =**

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

**data_a = softplus(data)**

**data_a =**

| 1.313 | 5.006 | 0.018 | 3.048 | 0.126 |
|-------|-------|-------|-------|-------|



- softplus
- softplus_derivative

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

18

❖ **ReLU function**

$$\text{ReLU}(x) = \begin{cases} 0 & \text{if } x \leq 0 \\ x & \text{if } x > 0 \end{cases}$$

data = | 1 | 5 | -4 | 3 | -2 |

data_a = **ReLU(data)**

data_a = | 1 | 5 | 0 | 3 | 0 |



— relu
— derivative
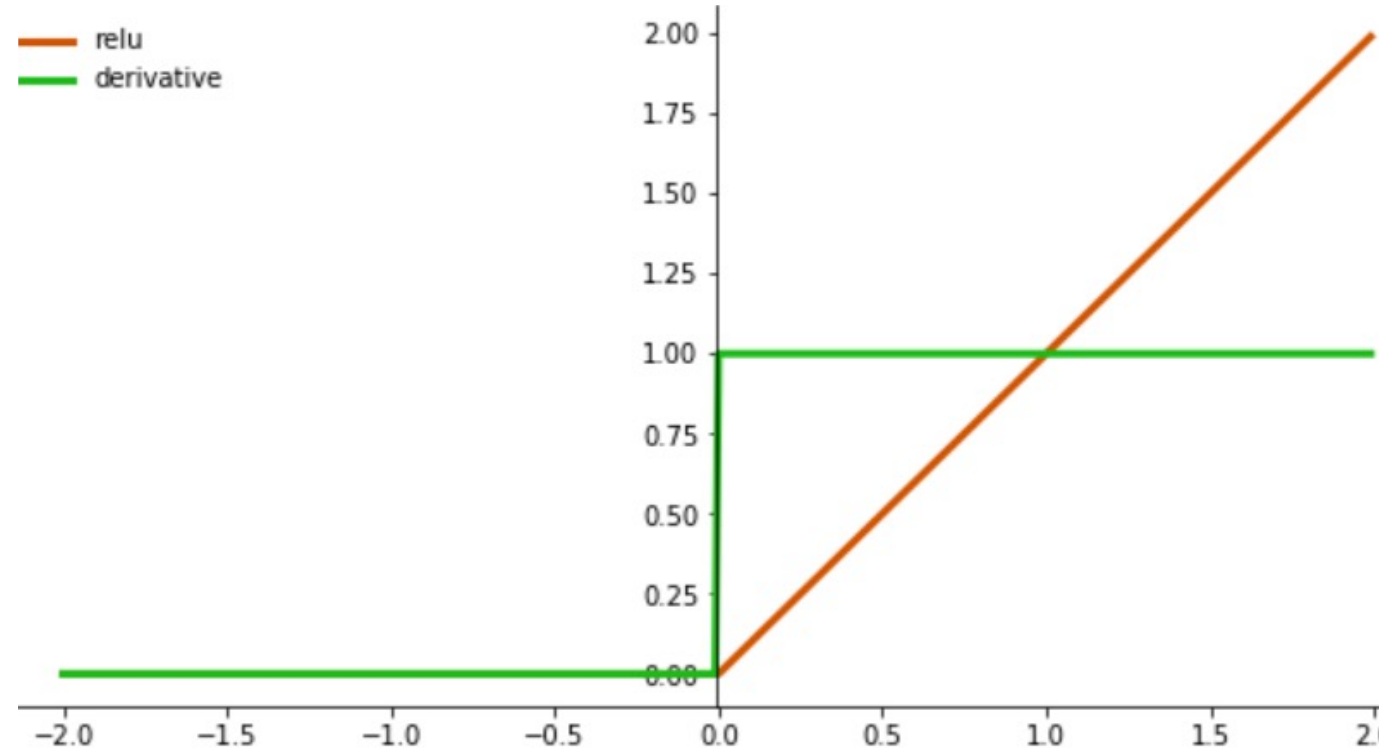
$$\text{ReLU}'(x) = \begin{cases} 0 & \text{if } x \leq 0 \\ 1 & \text{if } x > 0 \end{cases}$$

19

# Activation Functions

❖ **LeakyReLU function**

$$\text{LeakyReLU}(x) = \begin{cases} 0.01x & \text{if } x \leq 0 \\ x & \text{if } x > 0 \end{cases}$$
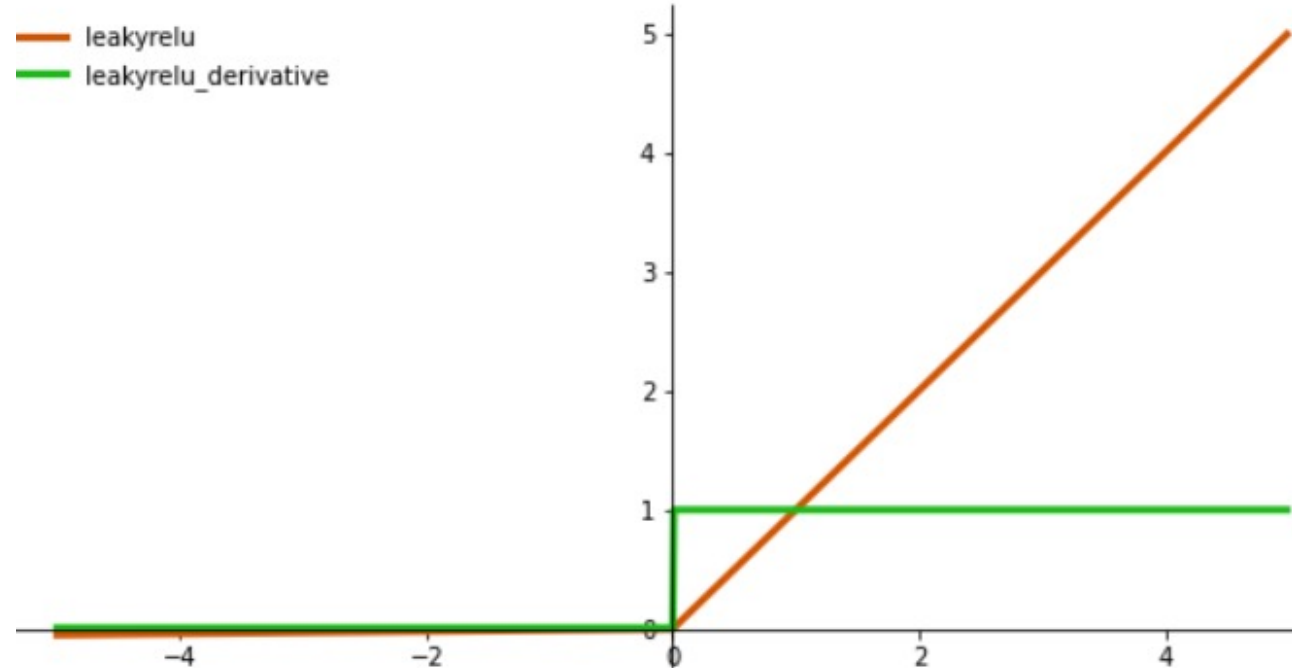
data =
| 1 | 5 | -4 | 3 | -2 |

data_a = **leakyrelu**(data)

data_a =
| 1 | 5 | -0.04 | 3 | -0.02 |


— leakyrelu
— leakyrelu_derivative

$$\text{LeakyReLU}'(x) = \begin{cases} 0.01 & \text{if } x \leq 0 \\ 1 & \text{if } x > 0 \end{cases}$$

20

# Activation Functions

❖ **ELU function**

$$ELU(x) = \begin{cases} \alpha(e^x - 1) & \text{if } x \leq 0 \\ x & \text{if } x > 0 \end{cases}$$



— elu
— elu_derivative
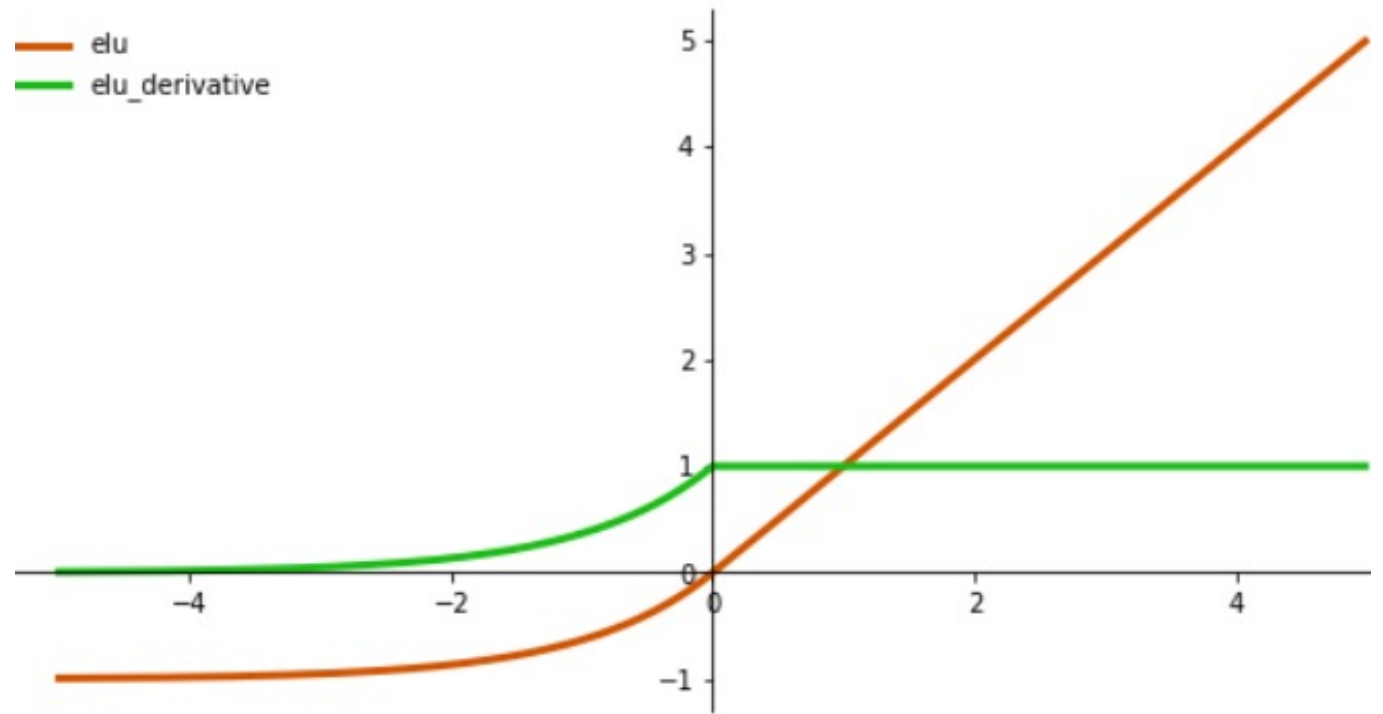
$$\alpha = 0.1$$

**data** =

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

**data_a** = **ELU(data)**

**data_a** =

| 1 | 5 | -0.098 | 3 | -0.086 |
|---|---|--------|---|--------|

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

21

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

❖ **PReLU function**

$$\text{PReLU}(x) = \begin{cases} \alpha x & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$$
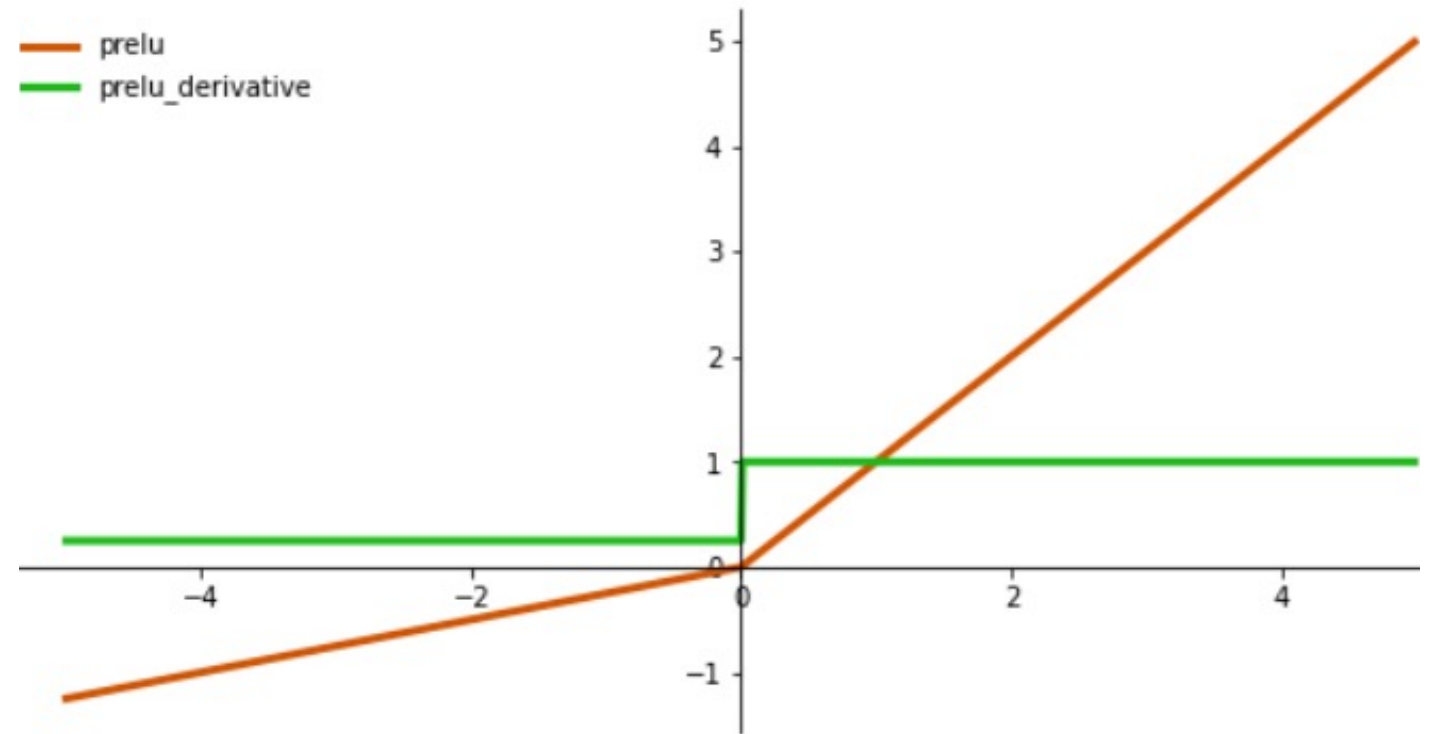
$\alpha = 0.1$



data =

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

data_a = **PRELU**(data)

data_a =

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

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

22

# Activation Functions

## ❖ Swish function

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$
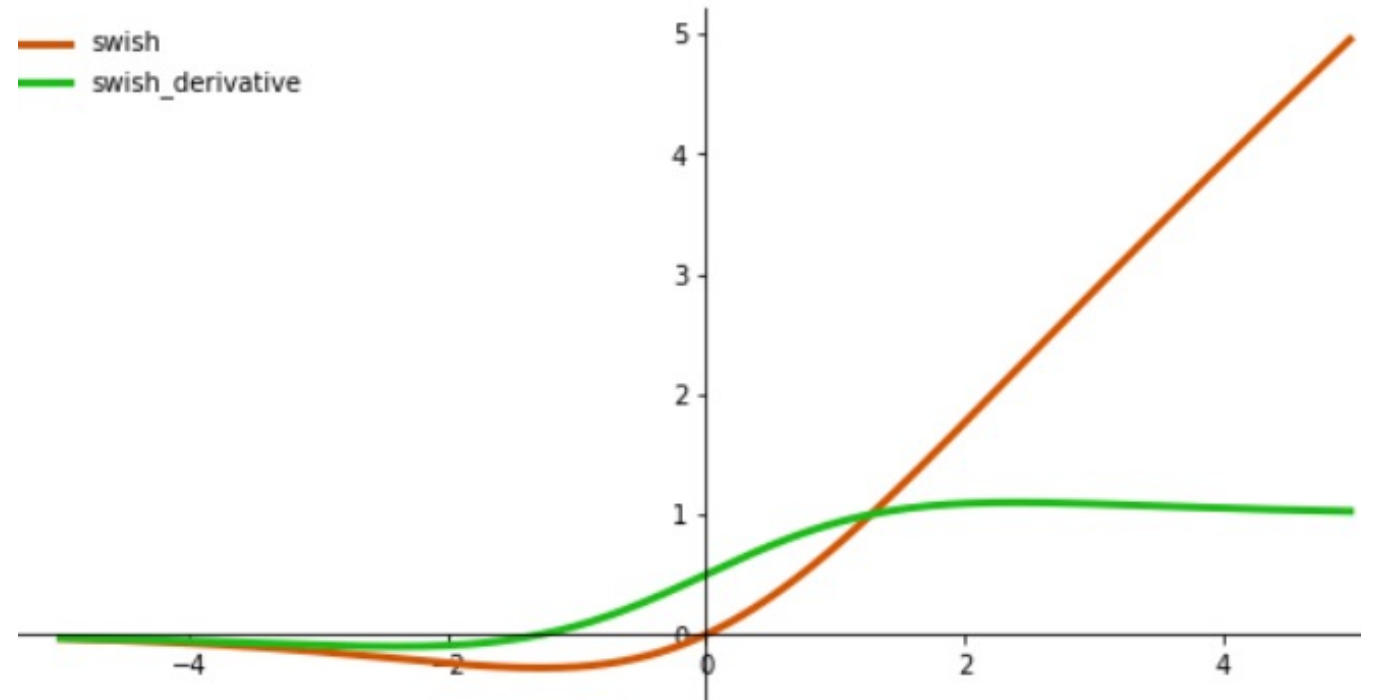
$$swish(x) = \frac{x}{1 + e^{-x}} = x\,\sigma(x)$$



swish
swish_derivative

**data =**

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

**data_a = swish(data)**

**data_a =**

| 0.731 | 4.966 | -0.071 | 2.857 | -0.238 |
|-------|-------|--------|-------|--------|

$$swish'(x) = swish(x) + \sigma(x)\,(1 - swish(x))$$

23

❖ **Swish function**

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

$$swish(x) = \frac{x}{1 + e^{-x}} = x\,\sigma(x)$$

$$swish'(x) = (x\,\sigma(x))' = (x)'\,\sigma(x) + x(\sigma(x))'$$

$$= \sigma(x) + x\,\sigma(x)\,(1 - \sigma(x))$$

$$= \sigma(x) + x\,\sigma(x) - x\,\sigma(x)^2$$

$$= x\,\sigma(x) + \sigma(x)(1 - x\,\sigma(x))$$

$$= swish(x) + \sigma(x)\,(1 - swish(x))$$

# Outline

# MLP Example 1

$$h = xW_h = \begin{bmatrix} 1 & 1.5 & 0.2 \\ 1 & 4.7 & 1.6 \\ 1 & 5.6 & 2.2 \end{bmatrix} \begin{bmatrix} 0.0 & 0.0 \\ 0.86 & -1.04 \\ 0.41 & -0.65 \end{bmatrix} = \begin{bmatrix} 1.373 & -1.696 \\ 4.708 & -5.951 \\ 5.731 & -7.281 \end{bmatrix}$$

$$\text{ReLU}(h) = \begin{bmatrix} 1.373 & 0 \\ 4.708 & 0 \\ 5.731 & 0 \end{bmatrix}$$



| Feature | | Label |
|---|---|---|
| Petal Length | Petal Width | Label |
| 1.5 | 0.2 | 0 |
| 1.4 | 0.2 | 0 |
| 1.6 | 0.2 | 0 |
| 4.7 | 1.6 | 1 |
| 3.3 | 1.1 | 1 |
| 4.6 | 1.3 | 1 |
| 5.6 | 2.2 | 2 |
| 5.1 | 1.5 | 2 |
| 5.6 | 1.4 | 2 |

Input layer    Hidden layer    Output layer

$$x = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ x^{(3)} \end{bmatrix} = \begin{bmatrix} 1 & 1.5 & 0.2 \\ 1 & 4.7 & 1.6 \\ 1 & 5.6 & 2.2 \end{bmatrix} \quad y = \begin{bmatrix} 0 \\ 1 \\ 2 \end{bmatrix}$$

$$W_h = \begin{bmatrix} W_{h1} & W_{h2} \end{bmatrix} \qquad W_z = \begin{bmatrix} W_{z1} & W_{z2} & W_{z3} \end{bmatrix}$$

$$= \begin{bmatrix} 0.0 & 0.0 \\ 0.86 & -1.04 \\ 0.41 & -0.65 \end{bmatrix} \qquad = \begin{bmatrix} 0.0 & 0.0 & 0.0 \\ 0.32 & 0.25 & 0.14 \\ -0.47 & -1.06 & 0.063 \end{bmatrix}$$

27

$$\text{ReLU}(\boldsymbol{h}) = \begin{bmatrix} 1.373 & 0 \\ 4.708 & 0 \\ 5.731 & 0 \end{bmatrix}$$

$$\boldsymbol{z} = \begin{bmatrix} 1 & \text{ReLU}(\boldsymbol{h}) \end{bmatrix} \boldsymbol{W_z} = \begin{bmatrix} 1 & 1.373 & 0 \\ 1 & 4.708 & 0 \\ 1 & 5.731 & 0 \end{bmatrix} \begin{bmatrix} 0.0 & 0.0 & 0.0 \\ 0.32 & 0.25 & 0.14 \\ -0.47 & -1.06 & 0.063 \end{bmatrix}$$

$$\begin{bmatrix} 1 & \text{ReLU}(\boldsymbol{h}) \end{bmatrix} = \begin{bmatrix} 1 & 1.373 & 0 \\ 1 & 4.708 & 0 \\ 1 & 5.731 & 0 \end{bmatrix}$$

$$= \begin{bmatrix} 0.439 & 0.356 & 0.195 \\ 1.507 & 1.220 & 0.670 \\ 1.835 & 1.485 & 0.816 \end{bmatrix}$$

| Feature | | Label |
|---|---|---|
| **Petal Length** | **Petal Width** | **Label** |
| 1.5 | 0.2 | 0 |
| 1.4 | 0.2 | 0 |
| 1.6 | 0.2 | 0 |
| 4.7 | 1.6 | 1 |
| 3.3 | 1.1 | 1 |
| 4.6 | 1.3 | 1 |
| 5.6 | 2.2 | 2 |
| 5.1 | 1.5 | 2 |
| 5.6 | 1.4 | 2 |



Input layer  Hidden layer  Output layer

$$\boldsymbol{x} = \begin{bmatrix} \boldsymbol{x}^{(1)} \\ \boldsymbol{x}^{(2)} \\ \boldsymbol{x}^{(3)} \end{bmatrix} = \begin{bmatrix} 1 & 1.5 & 0.2 \\ 1 & 4.7 & 1.6 \\ 1 & 5.6 & 2.2 \end{bmatrix} \quad \boldsymbol{y} = \begin{bmatrix} 0 \\ 1 \\ 2 \end{bmatrix}$$

$$\boldsymbol{W_h} = \begin{bmatrix} \boldsymbol{W_{h1}} & \boldsymbol{W_{h2}} \end{bmatrix} \qquad \boldsymbol{W_z} = \begin{bmatrix} \boldsymbol{W_{z1}} & \boldsymbol{W_{z2}} & \boldsymbol{W_{z3}} \end{bmatrix}$$

$$= \begin{bmatrix} 0.0 & 0.0 \\ 0.86 & -1.04 \\ 0.41 & -0.65 \end{bmatrix} \qquad = \begin{bmatrix} 0.0 & 0.0 & 0.0 \\ 0.32 & 0.25 & 0.14 \\ -0.47 & -1.06 & 0.063 \end{bmatrix}$$
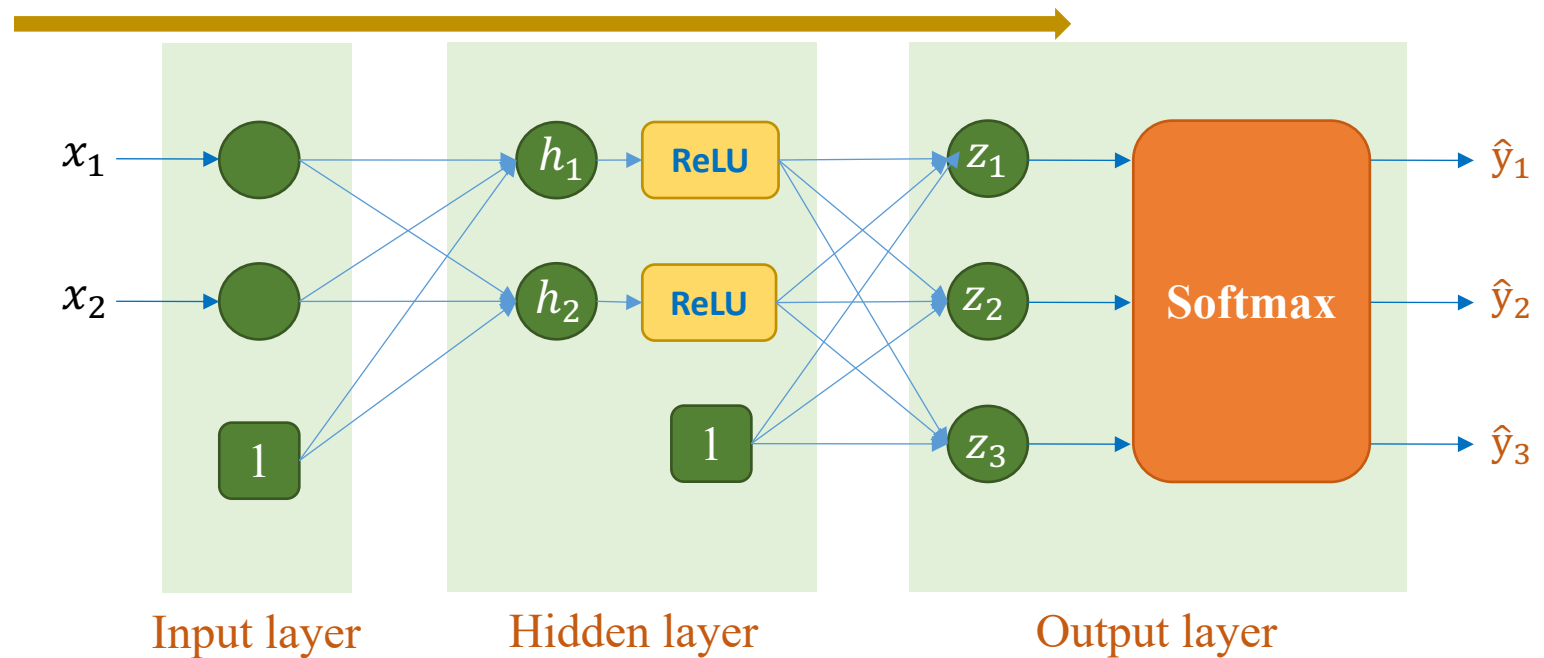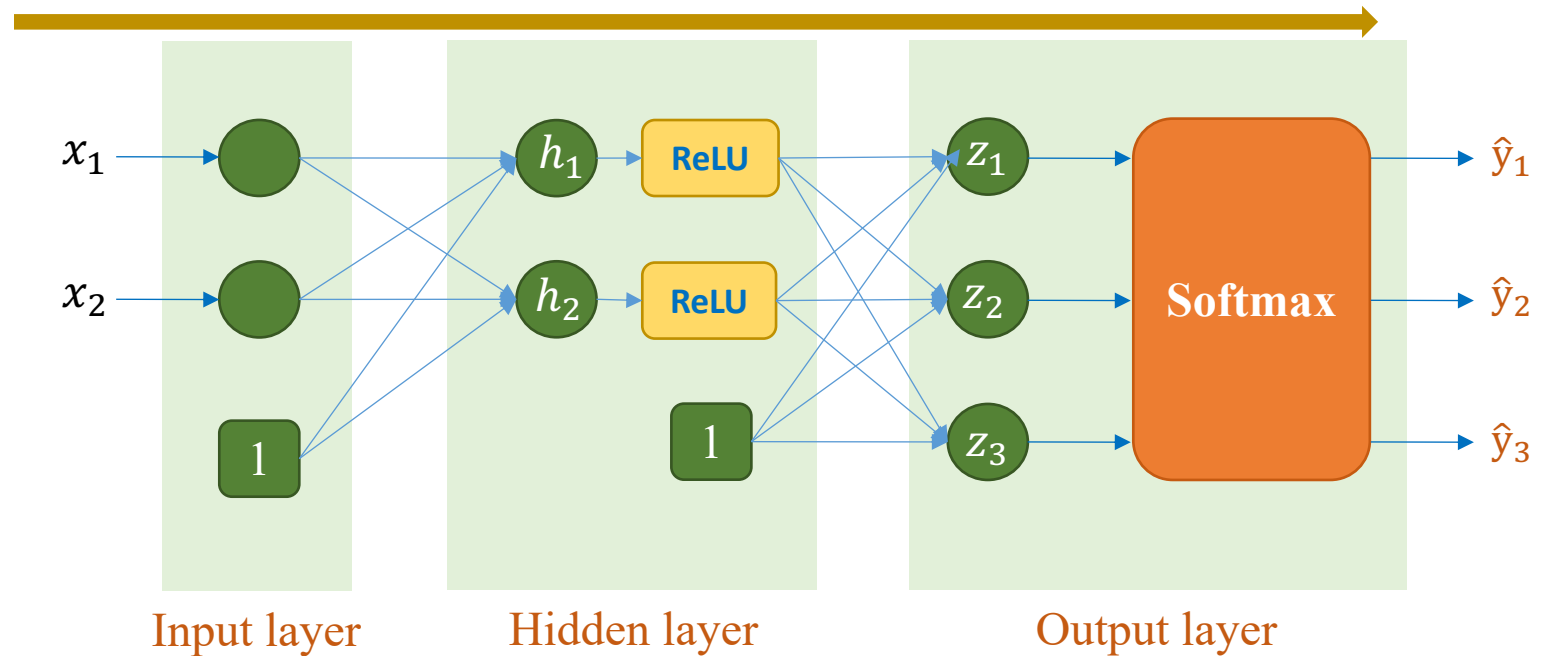
$$z = \begin{bmatrix} 0.439 & 0.356 & 0.195 \\ 1.507 & 1.220 & 0.670 \\ 1.835 & 1.485 & 0.816 \end{bmatrix}$$

$$\hat{y} = \text{softmax}(z) = \begin{bmatrix} \hat{y}^{(1)} \\ \hat{y}^{(2)} \\ \hat{y}^{(3)} \end{bmatrix} = \begin{bmatrix} 0.369 & 0.340 & 0.289 \\ 0.458 & 0.343 & 0.198 \\ 0.484 & 0.341 & 0.174 \end{bmatrix}$$

loss = 1.269

| Feature | | Label |
|---|---|---|
| Petal Length | Petal Width | Label |
| 1.5 | 0.2 | 0 |
| 1.4 | 0.2 | 0 |
| 1.6 | 0.2 | 0 |
| 4.7 | 1.6 | 1 |
| 3.3 | 1.1 | 1 |
| 4.6 | 1.3 | 1 |
| 5.6 | 2.2 | 2 |
| 5.1 | 1.5 | 2 |
| 5.6 | 1.4 | 2 |



Input layer      Hidden layer      Output layer

$$x = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ x^{(3)} \end{bmatrix} = \begin{bmatrix} 1 & 1.5 & 0.2 \\ 1 & 4.7 & 1.6 \\ 1 & 5.6 & 2.2 \end{bmatrix} \quad y = \begin{bmatrix} 0 \\ 1 \\ 2 \end{bmatrix}$$

$$W_h = \begin{bmatrix} W_{h1} & W_{h2} \end{bmatrix} \qquad W_z = \begin{bmatrix} W_{z1} & W_{z2} & W_{z3} \end{bmatrix}$$

$$= \begin{bmatrix} 0.0 & 0.0 \\ 0.86 & -1.04 \\ 0.41 & -0.65 \end{bmatrix} \qquad = \begin{bmatrix} 0.0 & 0.0 & 0.0 \\ 0.32 & 0.25 & 0.14 \\ -0.47 & -1.06 & 0.063 \end{bmatrix}$$

**Example 2 - Dying ReLU**

AI VIET NAM
@aivietnam.edu.vn

| Feature | | Label |
|---|---|---|
| **Petal Length** | **Petal Width** | **Label** |
| 1.5 | 0.2 | 0 |
| 1.4 | 0.2 | 0 |
| 1.6 | 0.2 | 0 |
| 4.7 | 1.6 | 1 |
| 3.3 | 1.1 | 1 |
| 4.6 | 1.3 | 1 |
| 5.6 | 2.2 | 2 |
| 5.1 | 1.5 | 2 |
| 5.6 | 1.4 | 2 |

Input layer        Hidden layer        Output layer

$$x = \begin{bmatrix} 1.5 \\ 0.2 \end{bmatrix} \qquad y = 0$$

$$m = \begin{bmatrix} m_1 & m_2 \end{bmatrix}$$

$$= \begin{bmatrix} 0.86 & -1.04 \\ 0.41 & -0.65 \end{bmatrix}$$

$$w = \begin{bmatrix} w_1 & w_2 & w_3 \end{bmatrix}$$
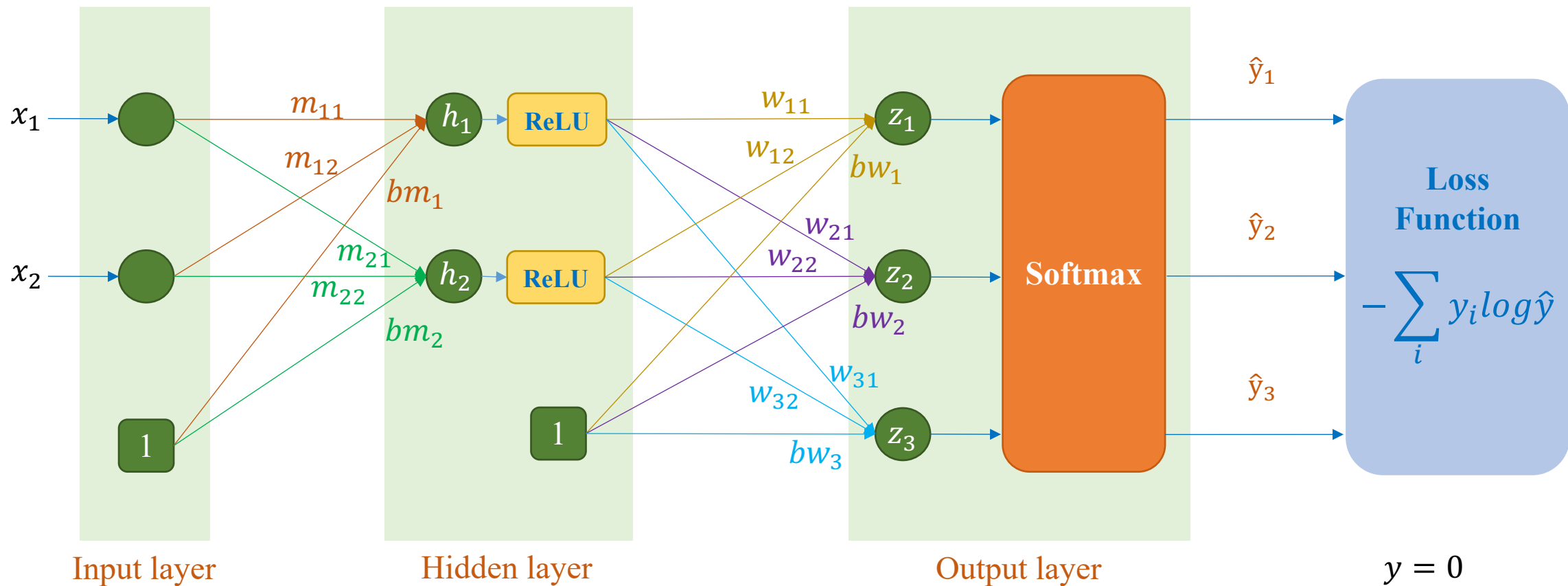
$$= \begin{bmatrix} 0.32 & 0.25 & 0.14 \\ -0.47 & -1.06 & 0.063 \end{bmatrix}$$

$$bm = \begin{bmatrix} 0.0 \\ 0.0 \end{bmatrix}$$

$$bw = \begin{bmatrix} 0.0 \\ 0.0 \\ 0.0 \end{bmatrix}$$

30

$$\boldsymbol{x} = \begin{bmatrix} 1.5 \\ 0.2 \end{bmatrix}$$

$$\boldsymbol{m} = \begin{bmatrix} \boldsymbol{m_1} & \boldsymbol{m_2} \end{bmatrix}$$

$$= \begin{bmatrix} 0.86 & -1.04 \\ 0.41 & -0.65 \end{bmatrix}$$

$$\boldsymbol{w} = \begin{bmatrix} \boldsymbol{w_1} & \boldsymbol{w_2} & \boldsymbol{w_3} \end{bmatrix}$$

$$= \begin{bmatrix} 0.32 & 0.25 & 0.14 \\ -0.47 & -1.06 & 0.063 \end{bmatrix}$$

$$\rightarrow \boldsymbol{y} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

$$\boldsymbol{bm} = \begin{bmatrix} 0.0 \\ 0.0 \end{bmatrix}$$

$$\boldsymbol{bw} = \begin{bmatrix} 0.0 \\ 0.0 \\ 0.0 \end{bmatrix}$$

31

Backward pass

$$\frac{\partial L}{\partial m_{jk}} = x_k \frac{\partial L}{\partial h_j}$$

$$\frac{\partial L}{\partial bm_j} = \frac{\partial L}{\partial h_j}$$

$$\frac{\partial L}{\partial relu_j} = \sum_i w_{ij} \frac{\partial L}{\partial z_i}$$

$$\text{ReLU}'(h_j) = \begin{cases} 0 & \text{if } h_j \leq 0 \\ 1 & \text{if } h_j > 0 \end{cases}$$

$$\frac{\partial L}{\partial h_j} = \begin{cases} 0 & \text{if } h_j \leq 0 \\ \dfrac{\partial L}{\partial relu_j} & \text{if } h_j > 0 \end{cases}$$

$$\frac{\partial L}{\partial z_i} = \hat{y}_i - y_i$$

$$\frac{\partial L}{\partial w_{ij}} = \text{ReLU}_j \frac{\partial L}{\partial z_i}$$

$$\frac{\partial L}{\partial bw_i} = \frac{\partial L}{\partial z_i}$$
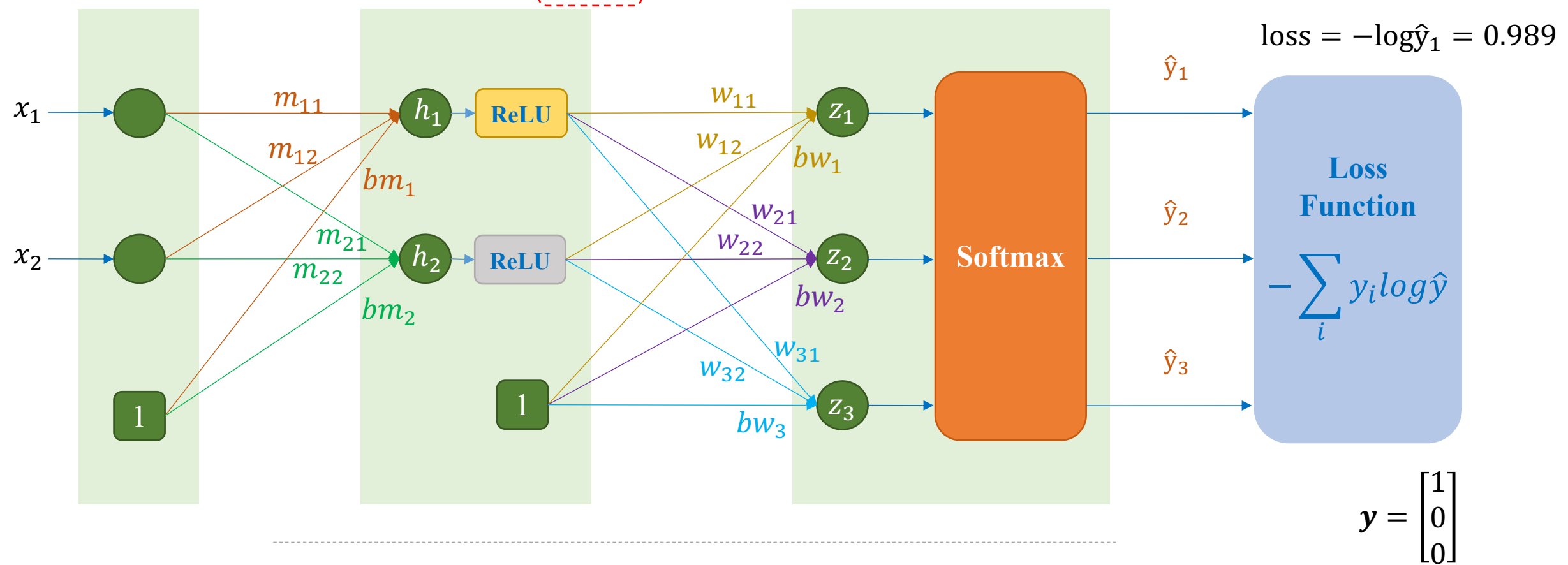
33

$$x = \begin{bmatrix} 1.5 \\ 0.2 \end{bmatrix} \qquad h = \begin{bmatrix} 1.372 \\ -1.68 \end{bmatrix} \qquad \text{ReLU} = \begin{bmatrix} 1.372 \\ 0.0 \end{bmatrix} \qquad z = \begin{bmatrix} 0.439 \\ 0.343 \\ 0.192 \end{bmatrix} \qquad \hat{y} = \begin{bmatrix} 0.372 \\ 0.338 \\ 0.290 \end{bmatrix}$$
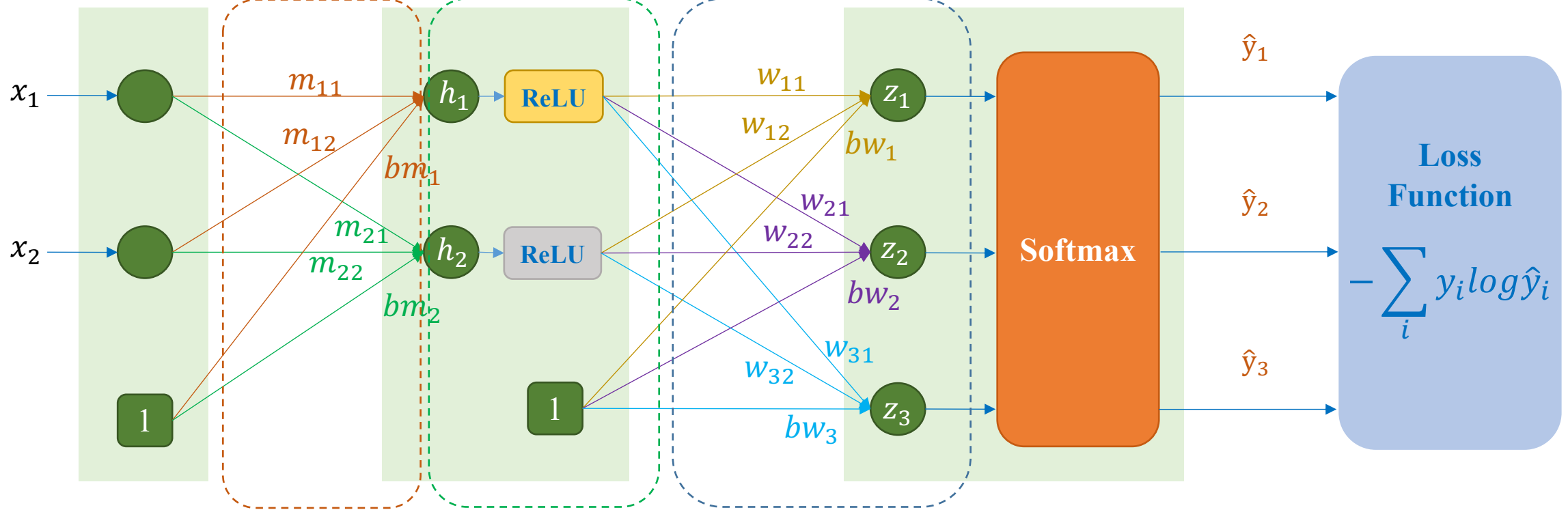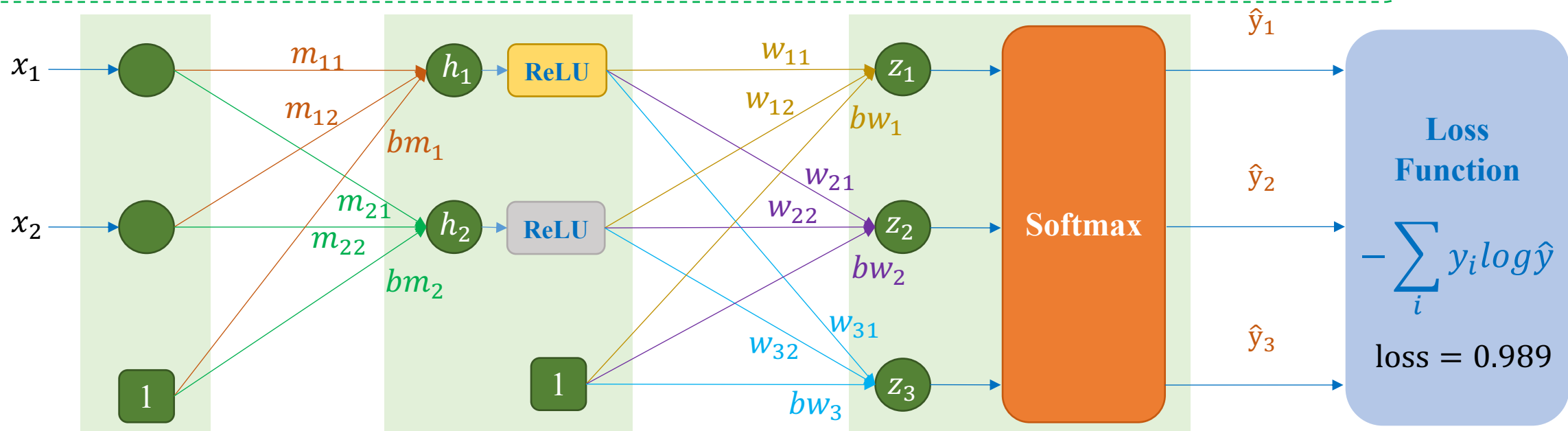
$$m = \begin{bmatrix} 0.86 & -1.04 \\ 0.41 & -0.65 \end{bmatrix} \quad bm = \begin{bmatrix} 0.0 \\ 0.0 \end{bmatrix} \quad w = \begin{bmatrix} 0.32 & 0.25 & 0.14 \\ -0.47 & -1.06 & 0.063 \end{bmatrix} \quad bw = \begin{bmatrix} 0.0 \\ 0.0 \\ 0.0 \end{bmatrix}$$

$x_1$ · $m_{11}$ · $m_{12}$ · $h_1$ · ReLU · $w_{11}$ · $w_{12}$ · $z_1$ · $bw_1$ · $\hat{y}_1$

$bm_1$

$x_2$ · $m_{21}$ · $m_{22}$ · $h_2$ · ReLU · $w_{21}$ · $w_{22}$ · $z_2$ · $bw_2$ · $\hat{y}_2$

$bm_2$

$1$ · $w_{32}$ · $w_{31}$ · $z_3$ · $bw_3$ · $\hat{y}_3$

**Softmax**

**Loss Function**

$$-\sum_i y_i \log \hat{y}$$

$$\text{loss} = 0.989$$

$$\frac{\partial L}{\partial relu_j} = \sum_i w_{ij} \frac{\partial L}{\partial z_i} \qquad \frac{\partial L}{\partial w_{ij}} = \text{ReLU}_j \frac{\partial L}{\partial z_i} \qquad \frac{\partial L}{\partial bw_i} = \frac{\partial L}{\partial z_i} \qquad \frac{\partial L}{\partial z_i} = \hat{y}_i - y_i$$

$$\nabla_{\text{ReLU}} L = \begin{bmatrix} -0.0759 \\ -0.0445 \end{bmatrix} \qquad \nabla_w L = \begin{bmatrix} -0.861 & 0.463 & 0.398 \\ 0.0 & 0.0 & 0.0 \end{bmatrix} \qquad \nabla_{bw} L = \begin{bmatrix} -0.628 \\ 0.338 \\ 0.290 \end{bmatrix} \qquad \nabla_z L = \begin{bmatrix} -0.628 \\ 0.338 \\ 0.290 \end{bmatrix}$$

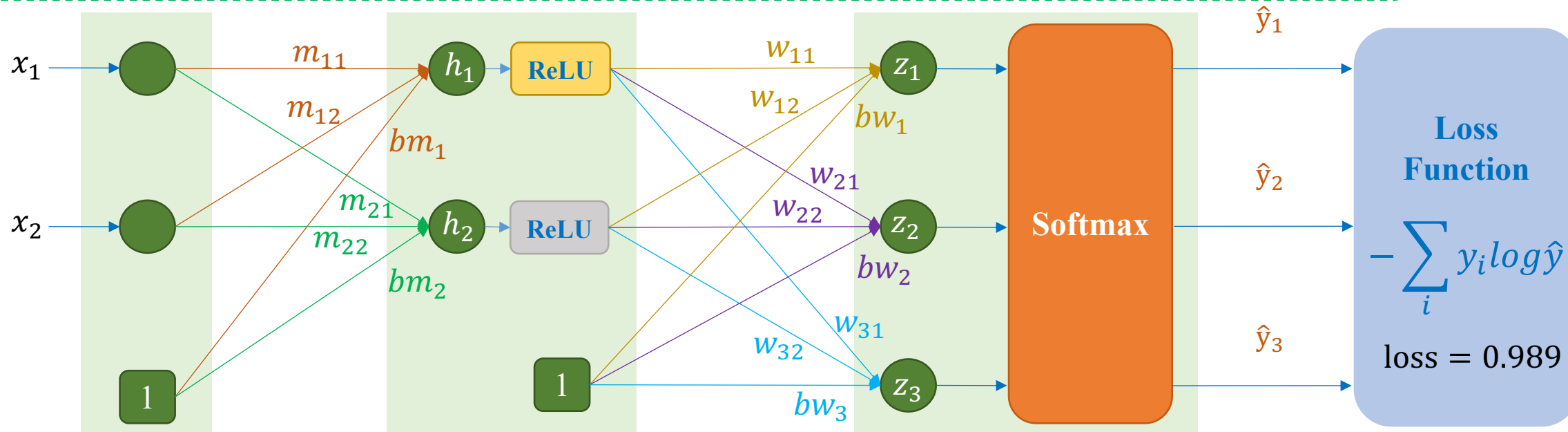$$x = \begin{bmatrix} 1.5 \\ 0.2 \end{bmatrix} \qquad h = \begin{bmatrix} 1.372 \\ -1.68 \end{bmatrix} \qquad \text{ReLU} = \begin{bmatrix} 1.372 \\ 0.0 \end{bmatrix} \qquad z = \begin{bmatrix} 0.439 \\ 0.343 \\ 0.192 \end{bmatrix} \qquad \hat{y} = \begin{bmatrix} 0.372 \\ 0.338 \\ 0.290 \end{bmatrix}$$

Backward pass

$$m = \begin{bmatrix} 0.86 & -1.04 \\ 0.41 & -0.65 \end{bmatrix} \qquad bm = \begin{bmatrix} 0.0 \\ 0.0 \end{bmatrix} \qquad w = \begin{bmatrix} 0.32 & 0.25 & 0.14 \\ -0.47 & -1.06 & 0.063 \end{bmatrix} \qquad bw = \begin{bmatrix} 0.0 \\ 0.0 \\ 0.0 \end{bmatrix}$$

Loss Function

$$-\sum_i y_i \log \hat{y}$$

loss = 0.989

$$\frac{\partial L}{\partial m_{jk}} = x_k \frac{\partial L}{\partial h_j}$$

$$\nabla_m L = \begin{bmatrix} -0.114 & 0.0 \\ -0.015 & 0.0 \end{bmatrix}$$

$$\frac{\partial L}{\partial bm_j} = \frac{\partial L}{\partial h_j}$$

$$\nabla_{bm} L = \begin{bmatrix} -0.0759 \\ 0.0 \end{bmatrix}$$

$$\frac{\partial L}{\partial h_j} = \begin{cases} 0 & \text{if } h_j \leq 0 \\ \dfrac{\partial L}{\partial relu_j} & \text{if } h_j > 0 \end{cases}$$

$$\nabla_h L = \begin{bmatrix} -0.0759 \\ 0.0 \end{bmatrix}$$

$$\frac{\partial L}{\partial relu_j} = \sum_i w_{ij} \frac{\partial L}{\partial z_i}$$

$$\nabla_{\text{ReLU}} L = \begin{bmatrix} -0.0759 \\ -0.0445 \end{bmatrix}$$

$$m = \begin{bmatrix} 0.86 & -1.04 \\ 0.41 & -0.65 \end{bmatrix} \qquad bm = \begin{bmatrix} 0.0 \\ 0.0 \end{bmatrix} \qquad w = \begin{bmatrix} 0.32 & 0.25 & 0.14 \\ -0.47 & -1.06 & 0.063 \end{bmatrix} \qquad bw = \begin{bmatrix} 0.0 \\ 0.0 \\ 0.0 \end{bmatrix}$$
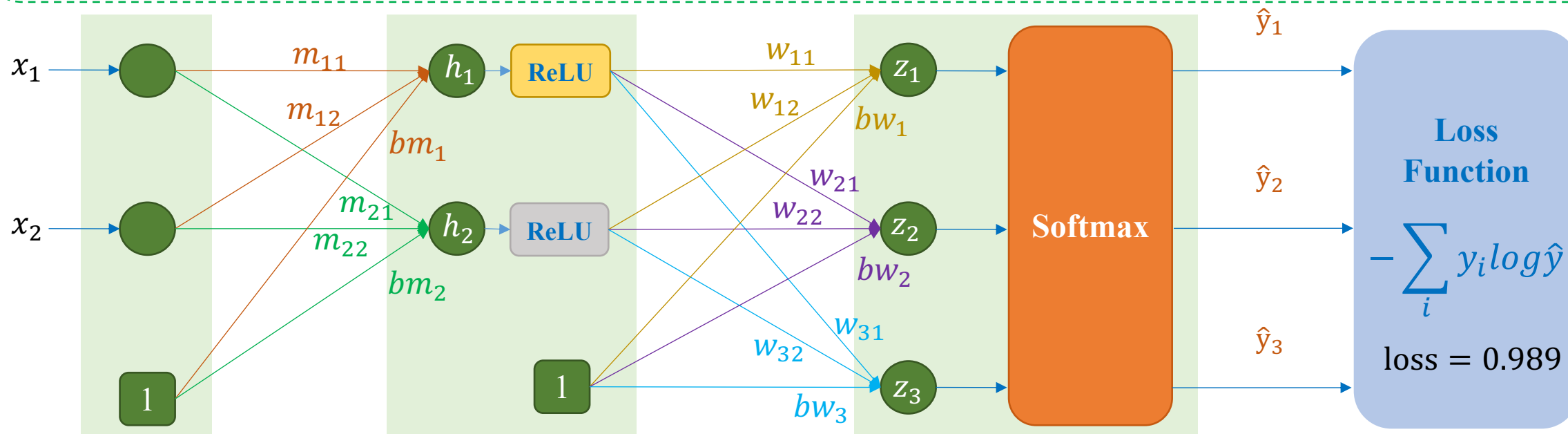
$$\nabla_m L = \begin{bmatrix} -0.114 & 0.0 \\ -0.015 & 0.0 \end{bmatrix} \qquad \nabla_{bm} L = \begin{bmatrix} -0.0759 \\ 0.0 \end{bmatrix} \qquad \nabla_w L = \begin{bmatrix} -0.628 & 0.338 & 0.29 \\ 0.0 & 0.0 & 0.0 \end{bmatrix} \qquad \nabla_{bw} L = \begin{bmatrix} -0.628 \\ 0.338 \\ 0.290 \end{bmatrix}$$

$x_1$

$x_2$

$m_{11}$

$m_{12}$

$bm_1$

$m_{21}$

$m_{22}$

$bm_2$

1

$h_1$

$h_2$

**ReLU**

**ReLU**

1

$w_{11}$

$w_{12}$

$bw_1$

$w_{21}$

$w_{22}$

$bw_2$

$w_{31}$

$w_{32}$

$bw_3$

$z_1$

$z_2$

$z_3$

**Softmax**

$\hat{y}_1$

$\hat{y}_2$

$\hat{y}_3$

**Loss Function**

$$-\sum_i y_i \log \hat{y}$$

loss = 0.989

Update the parameters with $\eta = 0.01$

$$m = \begin{bmatrix} 0.861 & -1.04 \\ 0.4105 & -0.65 \end{bmatrix} \qquad bm = \begin{bmatrix} 0.000759 \\ 0.0 \end{bmatrix} \qquad w = \begin{bmatrix} 0.328 & 0.245 & 0.136 \\ -0.47 & -1.06 & 0.063 \end{bmatrix} \qquad bw = \begin{bmatrix} 0.0062 \\ -0.0033 \\ -0.0029 \end{bmatrix}$$

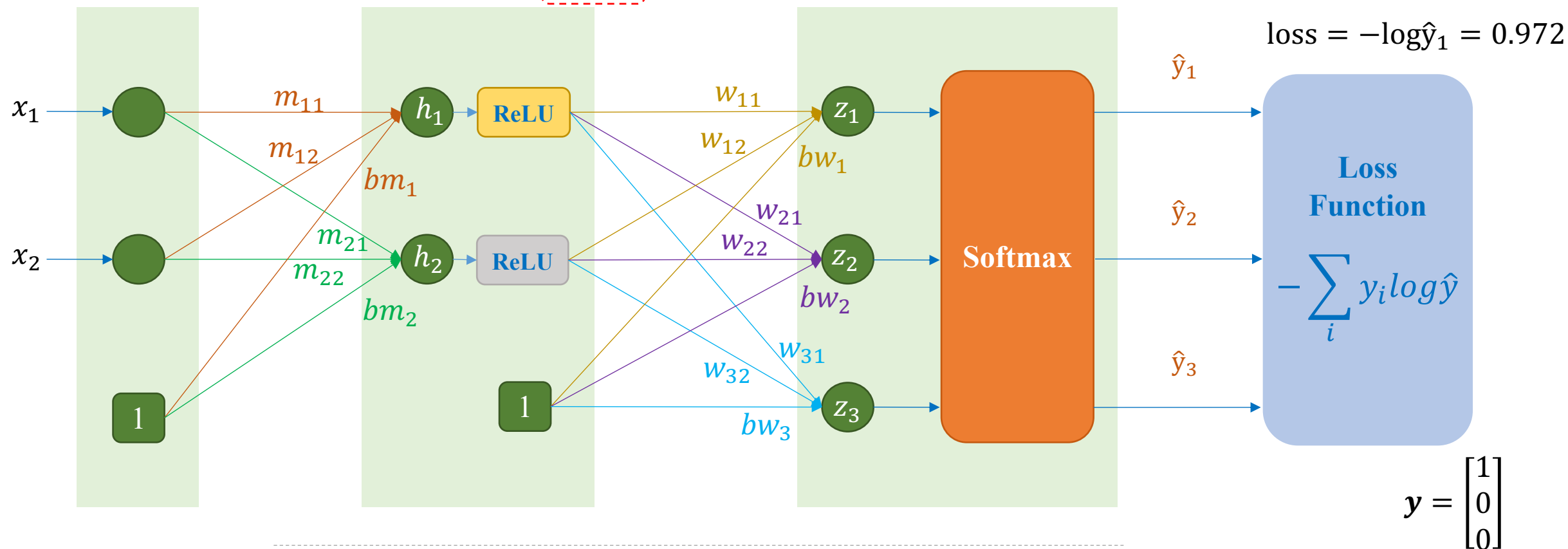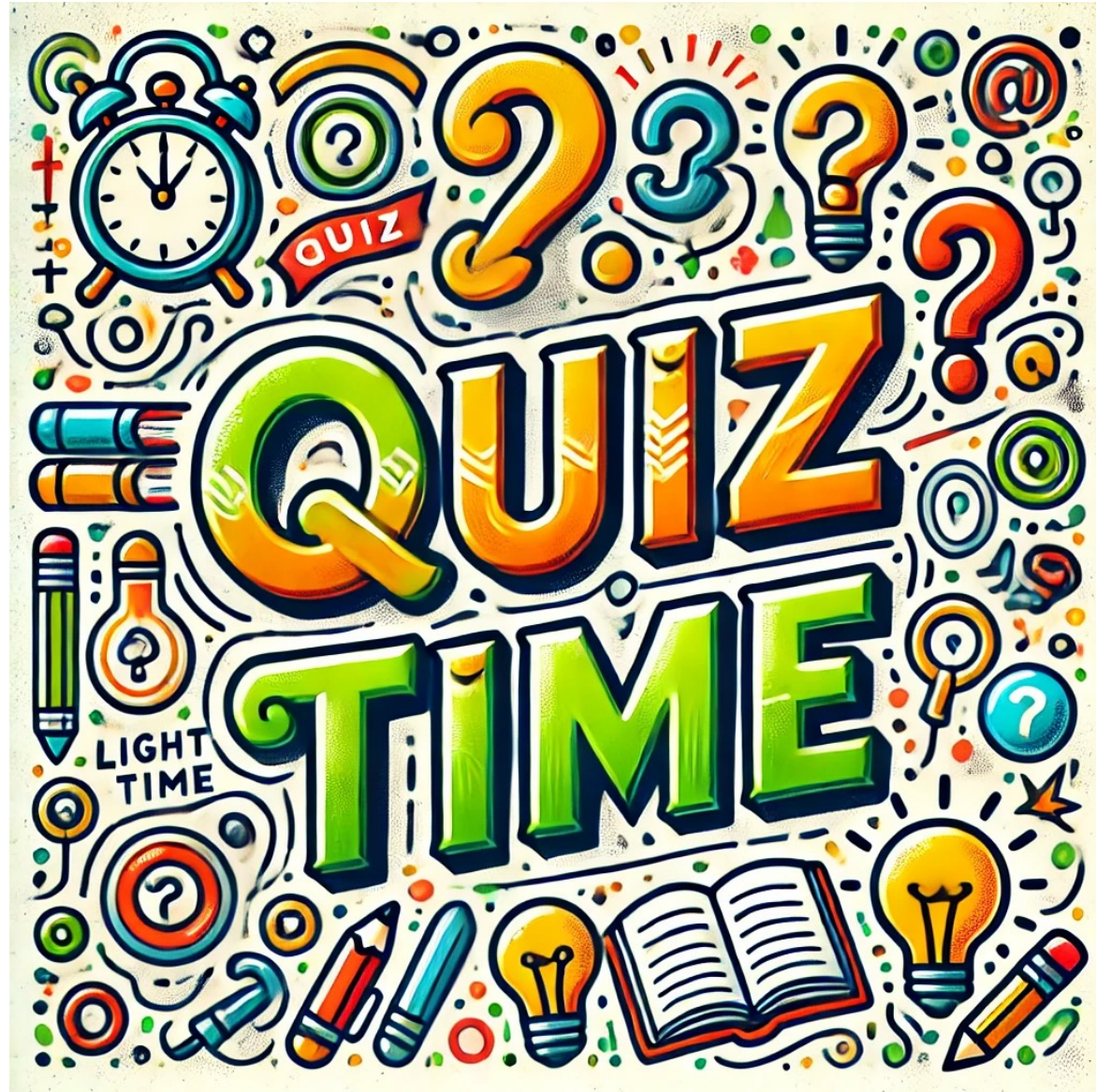❖ Chuẩn hóa dữ liệu ảnh nào có giá trị trung bình của data bằng 0 (chọn nhiều đáp án)?

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

❖ Code nào chuẩn hóa data và kết quả thuộc đoạn [0, 255]?

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

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

**3**
```python
Compose([transforms.ToTensor(), transforms.Normalize((mean,), (std,))])
```
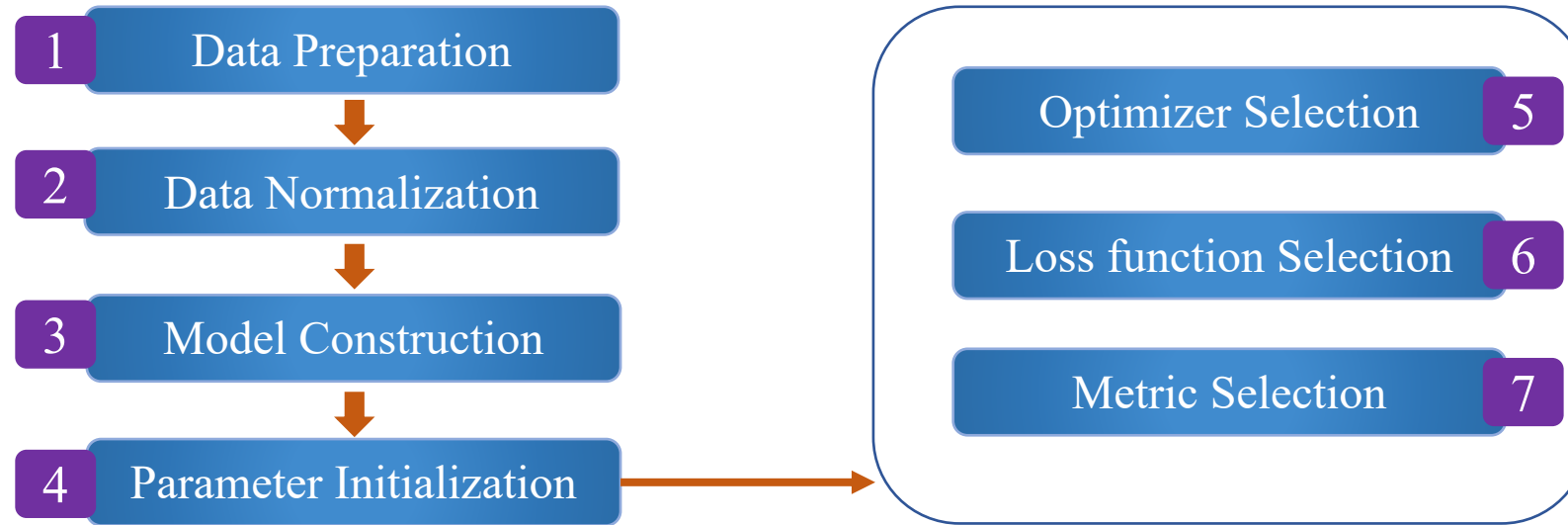
a) Code 1

b) Code 2

c) Code 3

d) Không code nào ở trên

# Question 3

❖ Chọn 2 thành phần ít quan trọng nhất từ hình pipeline huấn luyện sau?



a) Thành phần (1) hoặc (2)

b) Thành phần (3) hoặc (4)

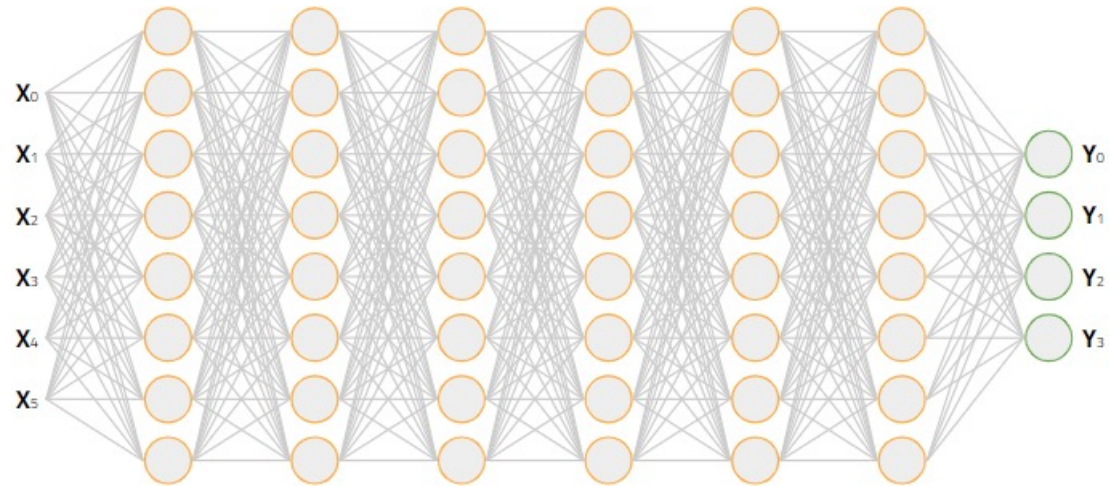c) Thành phần (5) hoặc (6)

d) Thành phần (7)

❖ Activation nào không nên dùng cho mô hình sau?

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

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

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

$$\text{GELU}(x) = x\phi(x)$$

$$\approx x * \text{sigmoid}(1.702x)$$



a) Sigmoid(.)

b) Tanh(.)

c) ReLU(.)

d) GELU(.)

❖ Khởi tạo tất cả tham số của model sau đều bằng 0. Việc huấn luyện mô hình sẽ như thế nào?



a) Vẫn huấn luyện được

b) Không huấn luyện được

c) Không xác định được

d) Các câu trả lời trên đều sai
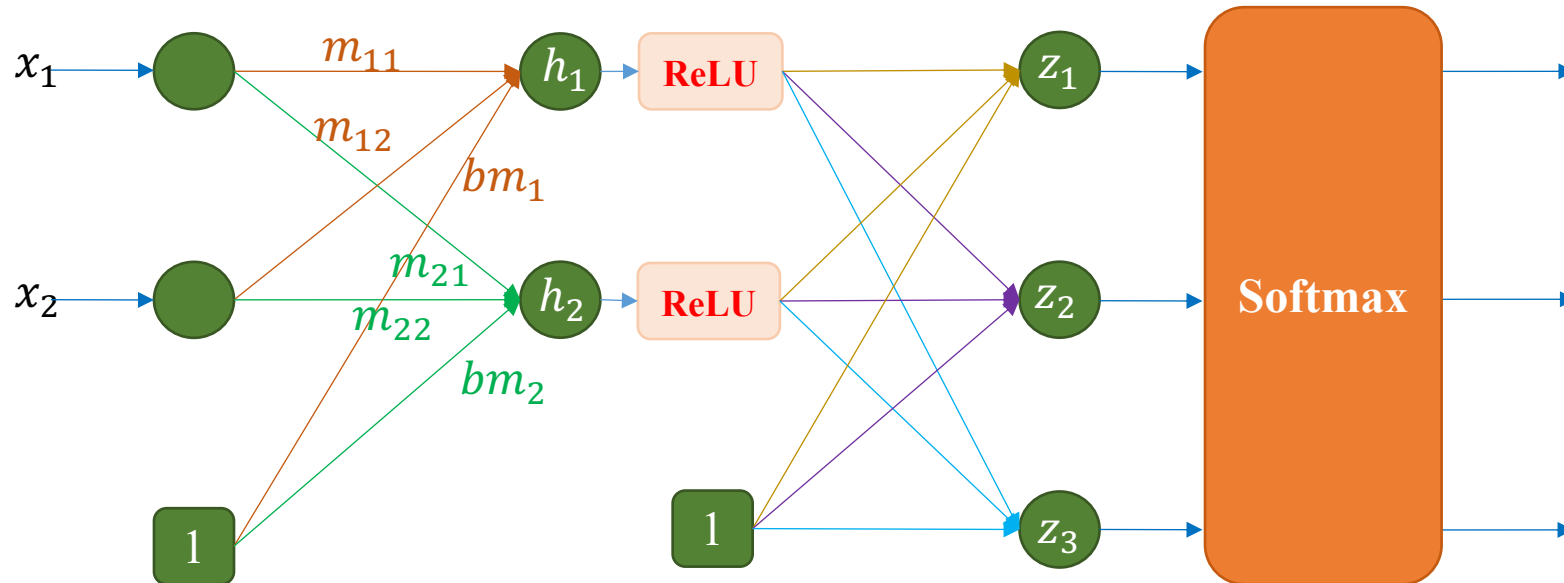
❖ Khởi tạo tất cả tham số của model sau đều bằng 0. Việc huấn luyện mô hình sẽ như thế nào?



a) Vẫn huấn luyện được

b) Không huấn luyện được

c) Không xác định được

d) Các câu trả lời trên đều sai

# Outline

$x = [1.4]$    $x$

$b_0$    $w_0$      $b_1$    $w_1$

0.0    0.0      0.0    0.0

$z_0 = w_0 x + b_0$      $z_1 = w_1 x + b_1$    $z_1 = [0.0]$

$$\hat{y}_0 = \frac{e^{z_0}}{\sum_{i=0}^{1} e^{z_i}}$$      $$\hat{y}_1 = \frac{e^{z_1}}{\sum_{i=0}^{1} e^{z_i}}$$    $\hat{y}_1 = [0.5]$

$$\text{L} = -y_0 \log \hat{y}_0 - y_1 \log \hat{y}_1$$   $y$   $y = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$

$L = [-\log 0.5] = [0.693]$

# Example 3 - Zero Initialization

AI VIET NAM
@aivietnam.edu.vn

❖ **Linear regression**

Diagram



Cheat sheet

Compute the output $\hat{y}$
$$\hat{y} = wx + b$$

Compute the loss
$$L = (\hat{y} - y)^2$$

Compute derivative
$$L'_w = 2x(\hat{y} - y)$$
$$L'_b = 2(\hat{y} - y)$$

Update parameters
$$w = w - \eta L'_w$$
$$b = b - \eta L'_b$$

38

Example 3 - Zero Initialization

Feature | Label

Given sample data

| area | price |
|------|-------|
| 6.7 | 9.1 |
| 4.6 | 5.9 |
| 3.5 | 4.6 |
| 5.5 | 6.7 |

House price prediction

$price = w * area + b$

Initialize b=0.0 and w=0.0

1

Input $x = 6.7$

Model

Parameters

$b = 0.0$   w = 0.0

$\hat{y} = xw + b = 0.0$

Label $y = 9.1$

Forward propagation

Loss $(\hat{y} - y)^2 = 82.81$

39

# Example 4 - Zero Initialization

❖ **Logistic regression**

$$\boldsymbol{\theta}^T = [b \quad w_1 \quad w_2]$$

$$\boldsymbol{x}^T = [1 \quad x_1 \quad x_2]$$

1) Pick a sample $(x, y)$ from training data

2) Compute output $\hat{y}$

$$z = \boldsymbol{\theta}^T \boldsymbol{x}$$

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

3) Compute loss

$$L(\boldsymbol{\theta}) = (-y\log\hat{y} - (1-y)\log(1-\hat{y}))$$

4) Compute derivative

$$\nabla_{\boldsymbol{\theta}} L = \mathbf{x}(\hat{y} - y)$$

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta L'_{\boldsymbol{\theta}}$$

$\eta$ is learning rate



$x_1$  $x_2$

Model

$b$  $w_1$  $w_2$

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

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

Label

$y$

Loss

$$-y\log\hat{y} - (1-y)\log(1-\hat{y})$$

41

# Example 4 - Zero Initialization



**Dataset**

| Petal_Length | Petal_Width | Label |
|---|---|---|
| 1.4 | 0.2 | 0 |
| 1.5 | 0.2 | 0 |
| 3 | 1.1 | 1 |
| 4.1 | 1.3 | 1 |

$$x = \begin{bmatrix} 1 \\ 1.4 \\ 0.2 \end{bmatrix} \qquad y = [0]$$

$$\eta = 0.01$$

$$b = -0.005$$
$$w_1 = -0.007$$
$$w_2 = -0.001$$

$$L'_{\boldsymbol{\theta}} = \mathbf{x}(\hat{y} - y)$$

$$= \begin{bmatrix} 1 \\ 1.4 \\ 0.2 \end{bmatrix} [0.5]$$

$$= \begin{bmatrix} 0.5 \\ 0.7 \\ 0.1 \end{bmatrix} = \begin{bmatrix} L'_b \\ L'_{w_1} \\ L'_{w_2} \end{bmatrix}$$

$$x_1 = 1.4 \qquad x_1 \qquad x_2 \qquad x_2 = 0.2$$

Model

$$b \qquad w_1 \qquad w_2$$

$$0.0 \qquad 0.0 \qquad 0.0$$

$$L'_b \qquad L'_{w_1} \qquad L'_{w_2}$$

$$z = w_1 x_1 + w_2 x_2 + b \qquad z = 0.0$$

$$\hat{y} = 0.5$$

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

$$y = 0$$

$$y$$

Loss

$$-y \log \hat{y} - (1-y) \log(1 - \hat{y})$$

$$L = 0.693$$

43

# Example 4 - Zero Initialization

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

**Dataset**

| Petal_Length | Petal_Width | Label |
|---|---|---|
| 1.4 | 0.2 | 0 |
| 1.5 | 0.2 | 0 |
| 3 | 1.1 | 1 |
| 4.1 | 1.3 | 1 |

$$x = \begin{bmatrix} 1 \\ 1.4 \\ 0.2 \end{bmatrix} \qquad y = [0]$$

$x_1 = 1.4$  $x_1$  $x_2$  $x_2 = 0.2$

Model

$b$  $w_1$  $w_2$

$-0.005$  $-0.007$  $-0.001$

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

$z = -0.016$

$\hat{y} = 0.49$

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

$y = 0$

$y$

Loss

$$-y\log\hat{y} - (1-y)\log(1-\hat{y})$$

previous $L = 1.1573$

$L = 0.68$

44

# Example 5 - Zero Initialization

## ❖ Softmax regression

**Training data**

| Feature | Label |
|---------|-------|
| Petal_Length | Label |
| 1.4 | 0 |
| 1.3 | 0 |
| 1.5 | 0 |
| 4.5 | 1 |
| 4.1 | 1 |
| 4.6 | 1 |

Category A

Category B

One-hot encoding for labels

index 0  1

$y = 0 \rightarrow \boldsymbol{y}^T = [1, 0]$

$y = 1 \rightarrow \boldsymbol{y}^T = [0, 1]$



$\hat{y}_0 = P(label = 0|x)$

$\hat{y}_1 = P(label = 1|x)$

**Model**

Model

$x$

$b_0$   $w_0$   $b_1$   $w_1$

$z_0 = w_0 x + b_0$   $z_1 = w_1 x + b_1$

$\hat{y}_0 = \dfrac{e^{z_0}}{\sum_{i=0}^1 e^{z_i}}$   $\hat{y}_1 = \dfrac{e^{z_1}}{\sum_{i=0}^1 e^{z_i}}$

Label

$\mathrm{L} = -y_0 \log \hat{y}_0 - y_1 \log \hat{y}_1$   $y$

45

**Example 5 - Zero Initialization**

AI VIET NAM
@aivietnam.edu.vn

❖ **Softmax regression**

| Feature | Label |
|---|---|
| Petal_Length | Label |
| 1.4 | 0 |
| 1.3 | 0 |
| 1.5 | 0 |
| 4.5 | 1 |
| 4.1 | 1 |
| 4.6 | 1 |

#class=2

#feature=1

One-hot encoding for label

$$y = 0 \rightarrow \boldsymbol{y}^T = [1 \quad 0]$$
$$y = 1 \rightarrow \boldsymbol{y}^T = [0 \quad 1]$$

Training example

$$(x, y) = (1.4, 0)$$



$x = [1.4]$

Model

$b_0$ $w_0$ $b_1$ $w_1$

$$z_0 = w_0 x + b_0$$

$$z_1 = w_1 x + b_1$$

$$\hat{y}_0 = \frac{e^{z_0}}{\sum_{i=0}^{1} e^{z_i}}$$

$$\hat{y}_1 = \frac{e^{z_1}}{\sum_{i=0}^{1} e^{z_i}}$$

$$L = -y_0 \log \hat{y}_0 - y_1 \log \hat{y}_1$$

Label

$y$

$\boldsymbol{y} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$

46

# Example 5 - Zero Initialization

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

❖ **Softmax regression**

$$\frac{\partial L}{\partial z_i} = \hat{y}_i - y_i$$

$$\frac{\partial L}{\partial w_i} = x(\hat{y}_i - y_i)$$

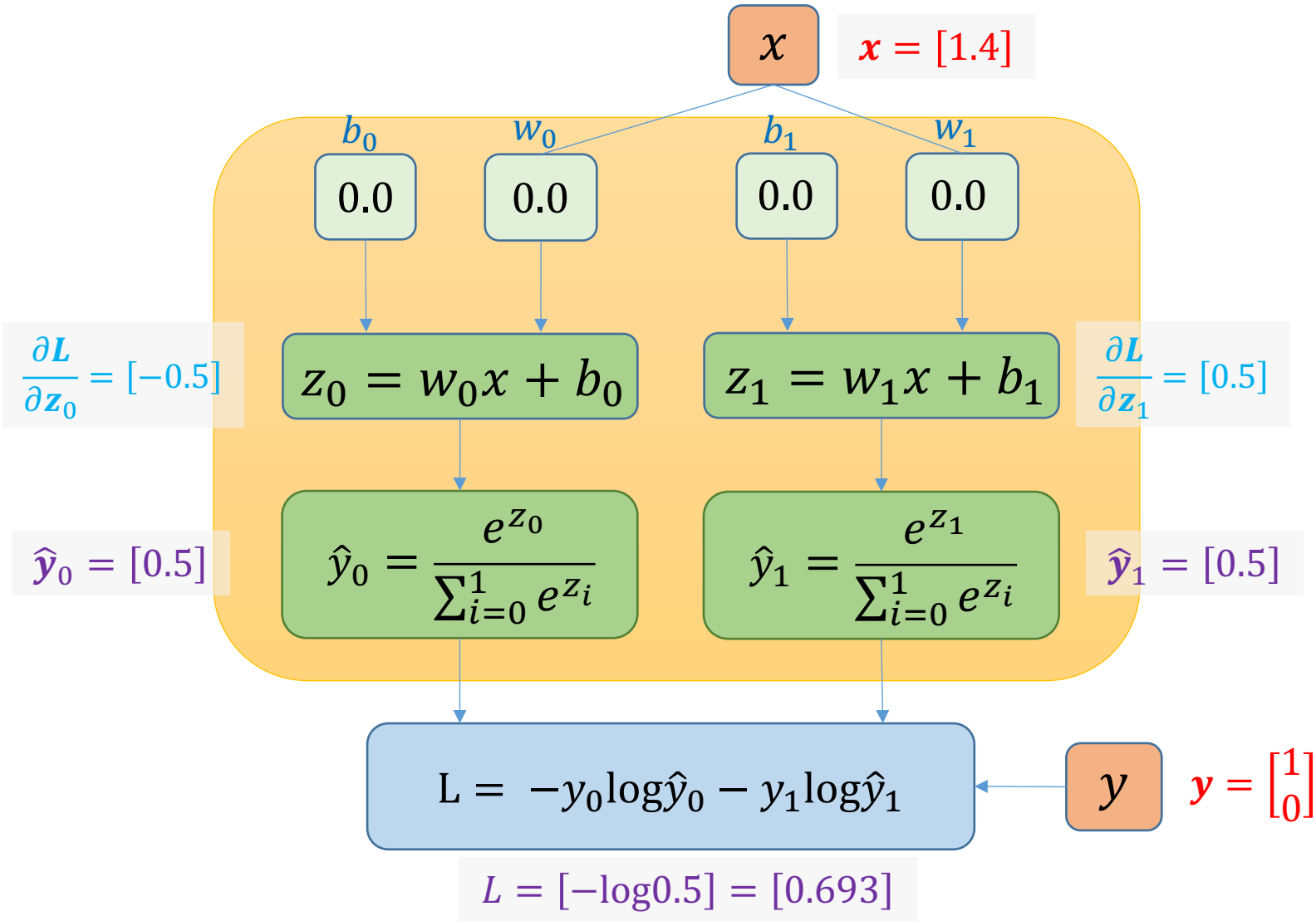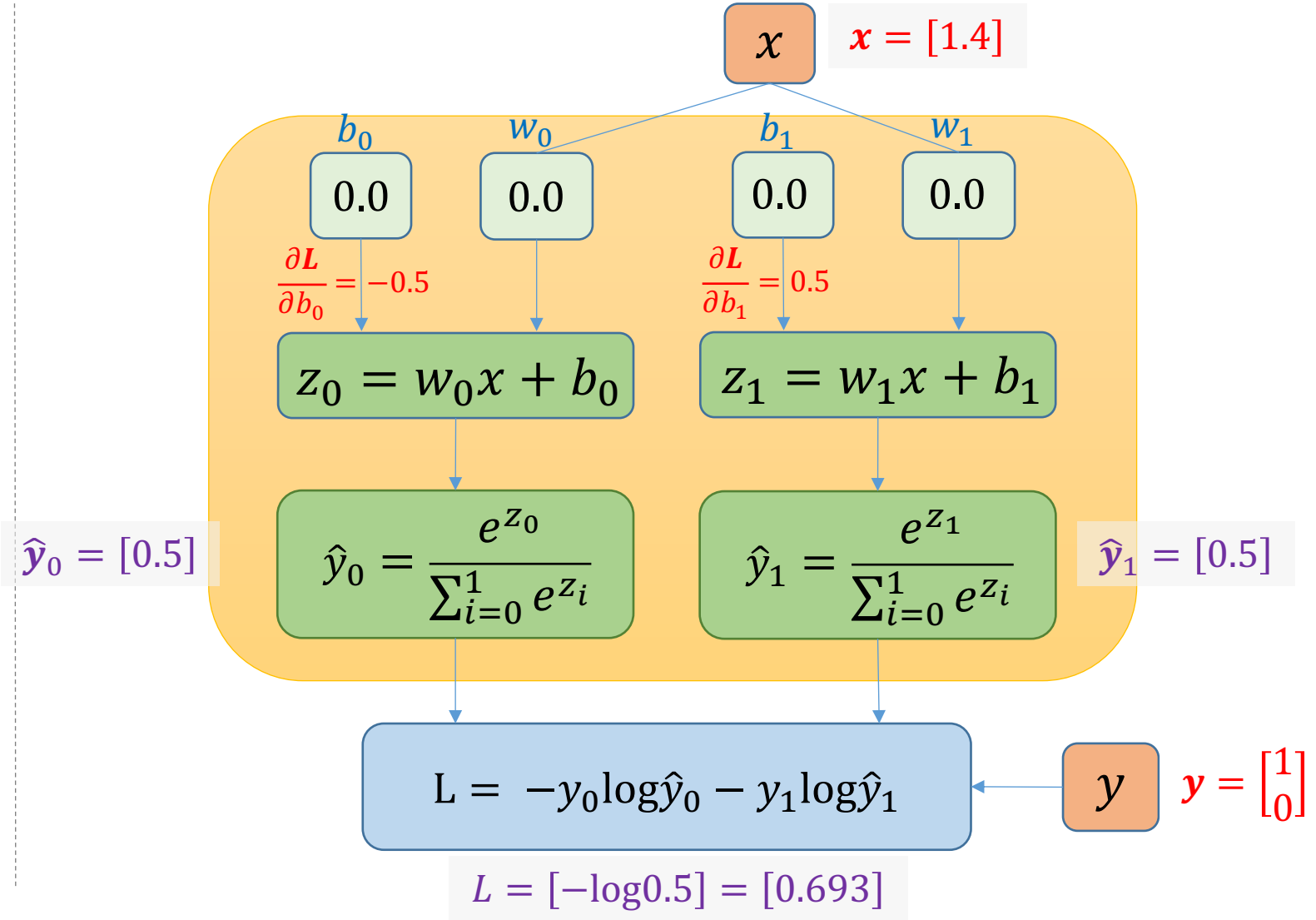$$\frac{\partial L}{\partial b_i} = \hat{y}_i - y_i$$

$$\begin{array}{cc} & y_0 \quad y_1 \\ y = 0 \rightarrow \boldsymbol{y}^T = [1 \quad 0] \\ y = 1 \rightarrow \boldsymbol{y}^T = [0 \quad 1] \end{array}$$

$$\frac{\partial L}{\partial w_0} = x(\hat{y}_0 - 1)$$

$$= -0.5 * 1.4 = -0.7$$

$$\frac{\partial L}{\partial w_1} = x(\hat{y}_1 - 0)$$

$$= 0.5 * 1.4 = 0.7$$

$x$   $\boldsymbol{x} = [1.4]$

$b_0$ : 0.0     $w_0$ : 0.0     $b_1$ : 0.0     $w_1$ : 0.0

$\frac{\partial \boldsymbol{L}}{\partial b_0} = -0.5$   $\frac{\partial \boldsymbol{L}}{\partial w_0} = -0.7$   $\frac{\partial \boldsymbol{L}}{\partial b_1} = 0.5$   $\frac{\partial \boldsymbol{L}}{\partial w_0} = 0.7$

$$z_0 = w_0 x + b_0 \qquad z_1 = w_1 x + b_1$$

$\hat{\boldsymbol{y}}_0 = [0.5]$

$$\hat{y}_0 = \frac{e^{z_0}}{\sum_{i=0}^{1} e^{z_i}} \qquad \hat{y}_1 = \frac{e^{z_1}}{\sum_{i=0}^{1} e^{z_i}}$$

$\hat{\boldsymbol{y}}_1 = [0.5]$

$$L = -y_0 \log \hat{y}_0 - y_1 \log \hat{y}_1$$

$y$   $\boldsymbol{y} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$

$L = [-\log 0.5] = [0.693]$

50

# Example 5 - Zero Initialization

❖ **Softmax regression**

**Update parameters**

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta L'_{\boldsymbol{\theta}}$$

$\eta$ is learning rate

$$\boldsymbol{\theta} = \begin{bmatrix} b_0 & b_1 \\ w_0 & w_1 \end{bmatrix} \qquad L'_{\boldsymbol{\theta}} = \begin{bmatrix} \dfrac{\partial L}{\partial b_0} & \dfrac{\partial L}{\partial b_1} \\ \dfrac{\partial L}{\partial w_0} & \dfrac{\partial L}{\partial w_1} \end{bmatrix}$$

$\eta = 0.1$

$$\boldsymbol{\theta} = \begin{bmatrix} 0.0 & 0.0 \\ 0.0 & 0.0 \end{bmatrix} - 0.01 \begin{bmatrix} -0.5 & 0.5 \\ -0.7 & 0.7 \end{bmatrix}$$

$$= \begin{bmatrix} -0.005 & 0.005 \\ -0.007 & 0.007 \end{bmatrix}$$

$x$ $\quad \boldsymbol{x} = [1.4]$

$b_0$ $\qquad w_0$ $\qquad\qquad b_1$ $\qquad w_1$

$-0.005$ $\qquad -0.007$ $\qquad\qquad 0.005$ $\qquad 0.007$

$\dfrac{\partial \boldsymbol{L}}{\partial b_0} = -0.5 \qquad \dfrac{\partial \boldsymbol{L}}{\partial w_0} = -0.7 \qquad \dfrac{\partial \boldsymbol{L}}{\partial b_1} = 0.5 \qquad \dfrac{\partial \boldsymbol{L}}{\partial w_0} = 0.7$

$$z_0 = w_0 x + b_0 \qquad\qquad z_1 = w_1 x + b_1$$

$\hat{\boldsymbol{y}}_0 = [0.5]$

$$\hat{y}_0 = \frac{e^{z_0}}{\sum_{i=0}^{1} e^{z_i}} \qquad\qquad \hat{y}_1 = \frac{e^{z_1}}{\sum_{i=0}^{1} e^{z_i}}$$

$\hat{\boldsymbol{y}}_1 = [0.5]$

$$L = -y_0 \log \hat{y}_0 - y_1 \log \hat{y}_1$$

$y \qquad \boldsymbol{y} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$

$$L = [-\log 0.5] = [0.693]$$

# Example 5 - Zero Initialization

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

❖ **Softmax regression**

| Feature | Label |
|---------|-------|
| Petal_Length | Label |
| 1.4 | 0 |
| 1.3 | 0 |
| 1.5 | 0 |
| 4.5 | 1 |
| 4.1 | 1 |
| 4.6 | 1 |

One-hot encoding for label

$$y = 0 \rightarrow \boldsymbol{y}^T = [1 \quad 0]$$
$$y = 1 \rightarrow \boldsymbol{y}^T = [0 \quad 1]$$

Training example

$$(x, y) = (1.4, 0)$$

$$x \qquad \boldsymbol{x} = [1.4]$$

$b_0$ : $-0.005$
$w_0$ : $-0.007$
$b_1$ : $0.005$
$w_1$ : $0.007$

$\boldsymbol{z}_0 = [0.015]$
$$z_0 = w_0 x + b_0$$
$\boldsymbol{z}_0 = [0.0]$

$\boldsymbol{z}_1 = [-0.015]$
$$z_1 = w_1 x + b_1$$
$\boldsymbol{z}_1 = [0.0]$

$\hat{\boldsymbol{y}}_0 = [0.51]$
$$\hat{y}_0 = \frac{e^{z_0}}{\sum_{i=0}^{1} e^{z_i}}$$
$\hat{\boldsymbol{y}}_0 = [0.5]$

$\hat{\boldsymbol{y}}_1 = [0.49]$
$$\hat{y}_1 = \frac{e^{z_1}}{\sum_{i=0}^{1} e^{z_i}}$$
$\hat{\boldsymbol{y}}_1 = [0.5]$

$$L = -y_0 \log \hat{y}_0 - y_1 \log \hat{y}_1$$

$$y \qquad \boldsymbol{y} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$L = [-\log 0.51] = [0.678]$
$L = [-\log 0.5] = [0.693]$

losses reduce!!!

52

# Example 6 - Zero Initialization

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

❖ **MLP**

| Feature | | Label |
|---|---|---|
| **Petal Length** | **Petal Width** | **Label** |
| 1.5 | 0.2 | 0 |
| 1.4 | 0.2 | 0 |
| 1.6 | 0.2 | 0 |
| 4.7 | 1.6 | 1 |
| 3.3 | 1.1 | 1 |
| 4.6 | 1.3 | 1 |
| 5.6 | 2.2 | 2 |
| 5.1 | 1.5 | 2 |
| 5.6 | 1.4 | 2 |



Input layer          Hidden layer          Output layer

$$x = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ x^{(3)} \end{bmatrix} = \begin{bmatrix} 1.5 & 0.2 \\ 4.7 & 1.6 \\ 5.6 & 2.2 \end{bmatrix}$$

$$y = \begin{bmatrix} 0 \\ 1 \\ 2 \end{bmatrix}$$

$$h = \begin{bmatrix} h_1 & h_2 \end{bmatrix}$$
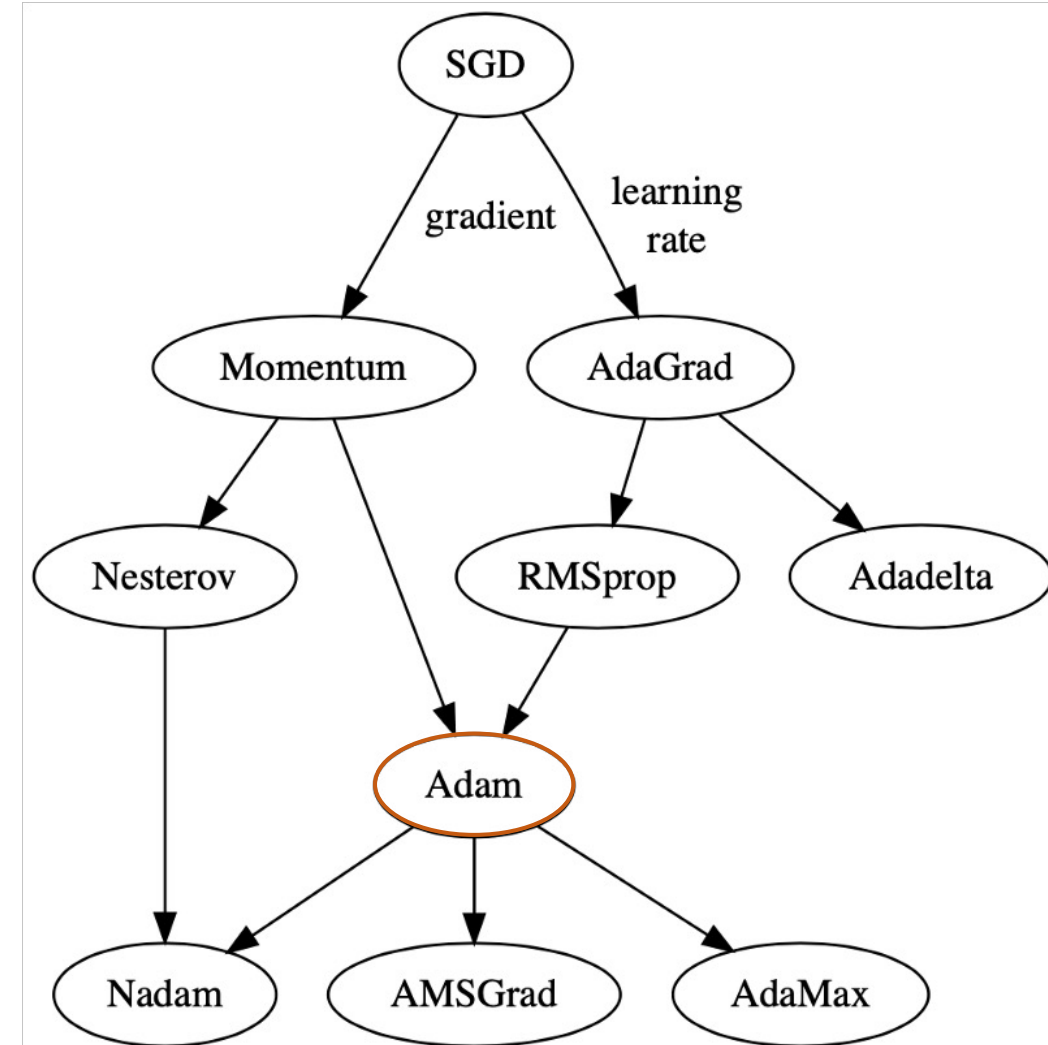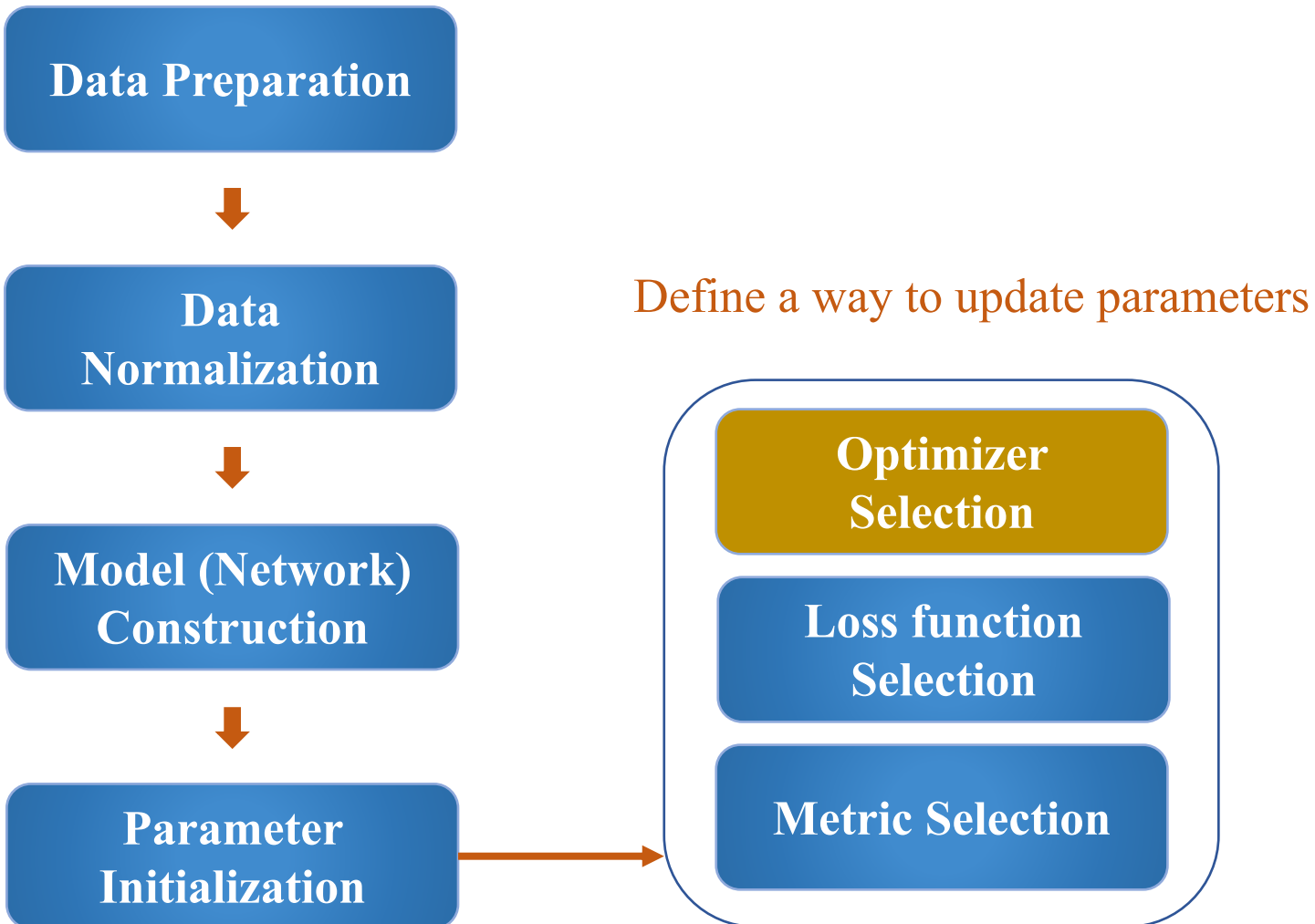$$= \begin{bmatrix} 0.0 & 0.0 \\ 0.0 & 0.0 \end{bmatrix}$$

$$b_h = \begin{bmatrix} 0.0 \\ 0.0 \end{bmatrix}$$

$$w = \begin{bmatrix} w_1 & w_2 & w_3 \end{bmatrix}$$
$$= \begin{bmatrix} 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \end{bmatrix}$$
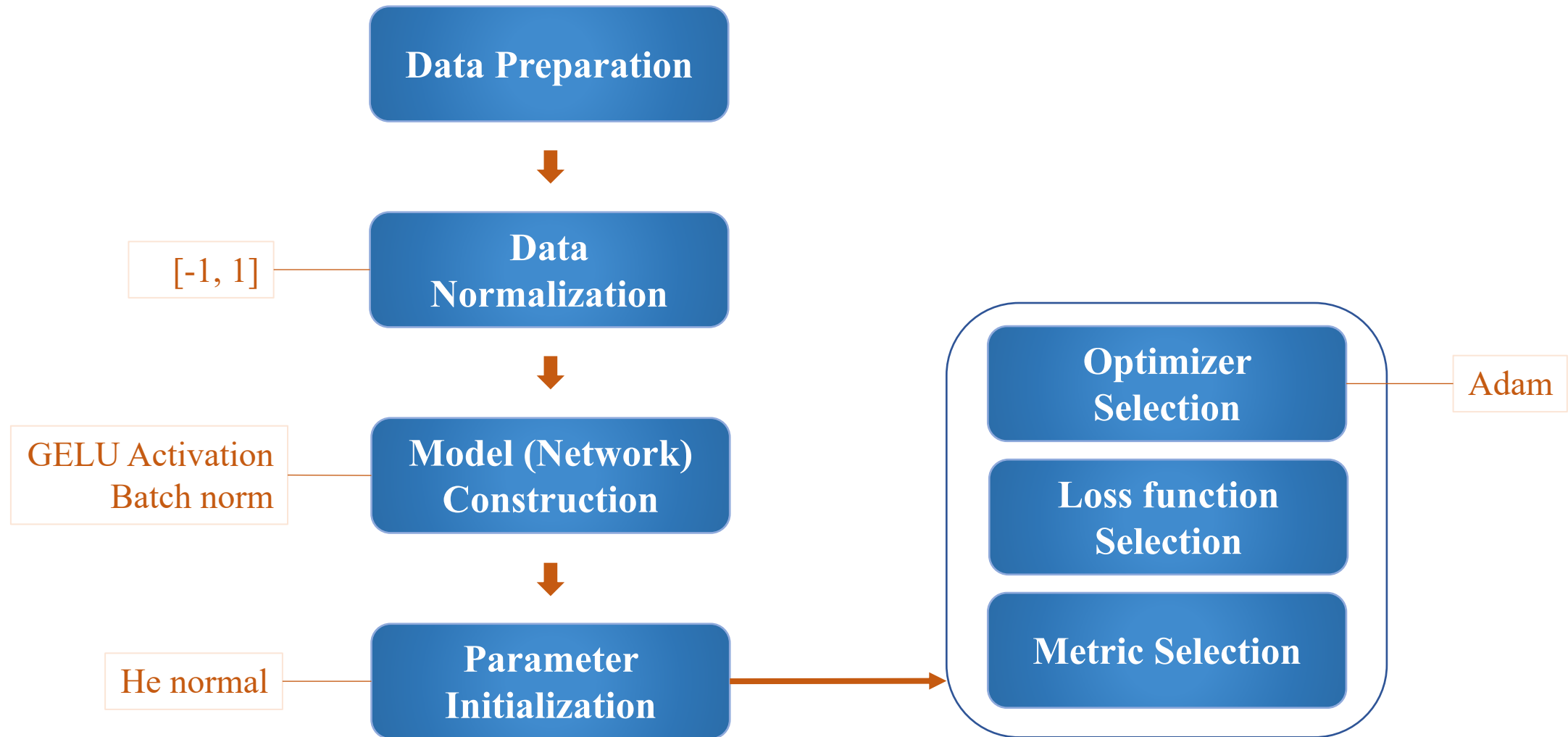
$$b_w = \begin{bmatrix} 0.0 \\ 0.0 \\ 0.0 \end{bmatrix}$$

53

$$x = \begin{bmatrix} 1.5 & 0.2 \\ 4.7 & 1.6 \\ 5.6 & 2.2 \end{bmatrix}$$

$$h = \begin{bmatrix} 0.0 & 0.0 \\ 0.0 & 0.0 \\ 0.0 & 0.0 \end{bmatrix}$$

$$\text{ReLU} = \begin{bmatrix} 0.0 & 0.0 \\ 0.0 & 0.0 \\ 0.0 & 0.0 \end{bmatrix}$$

$$z = \begin{bmatrix} 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \end{bmatrix}$$

$$\hat{y} = \begin{bmatrix} 0.333 & 0.333 & 0.333 \\ 0.333 & 0.333 & 0.333 \\ 0.333 & 0.333 & 0.333 \end{bmatrix}$$

$$\text{loss} = \begin{bmatrix} -\log 0.333 \\ -\log 0.333 \\ -\log 0.333 \end{bmatrix}$$

**Loss Function**

$$-\sum_i y_i \log \hat{y}$$

$$y = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$m = \begin{bmatrix} m_1 & m_2 \end{bmatrix} = \begin{bmatrix} 0.0 & 0.0 \\ 0.0 & 0.0 \end{bmatrix}$$

$$bm = \begin{bmatrix} 0.0 \\ 0.0 \end{bmatrix}$$

$$w = \begin{bmatrix} w_1 & w_2 & w_3 \end{bmatrix} = \begin{bmatrix} 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \end{bmatrix}$$

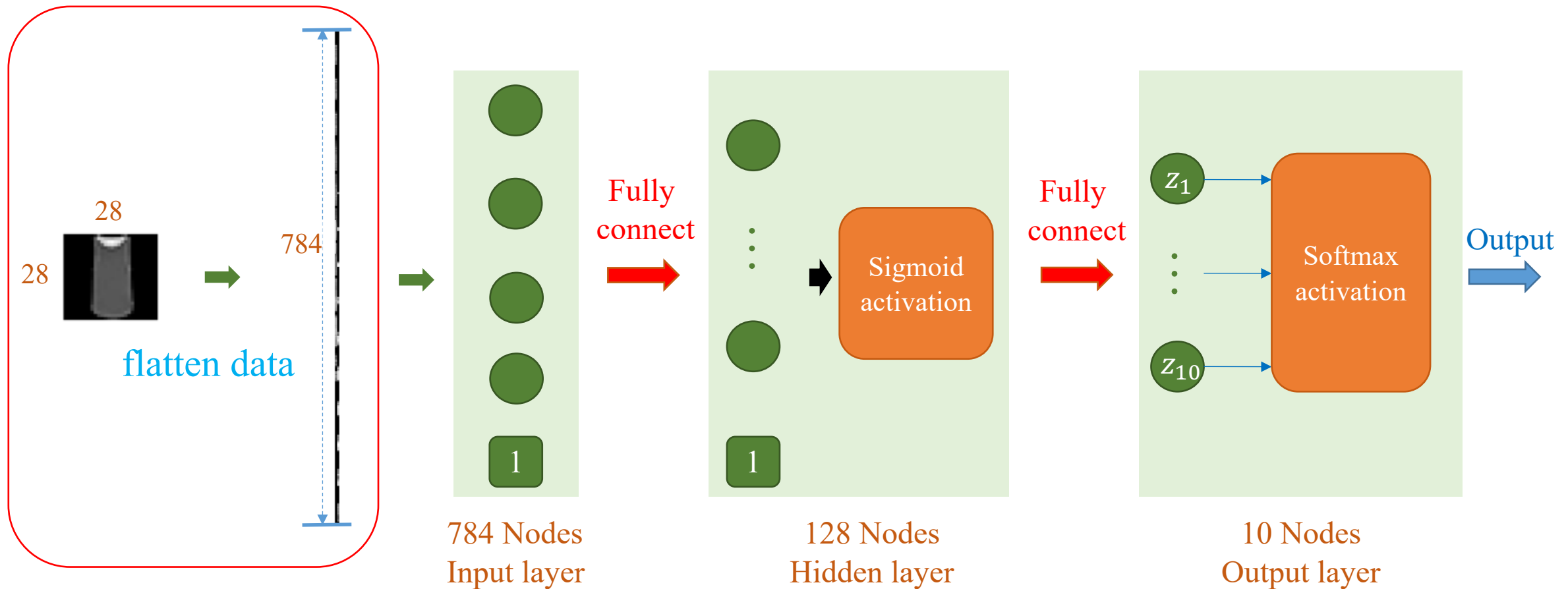$$bw = \begin{bmatrix} 0.0 \\ 0.0 \\ 0.0 \end{bmatrix}$$

54

55

# Optimizers

❖ **Optimizer Selection**

**Data Preparation**

↓

**Data Normalization**

↓

**Model (Network) Construction**

↓

**Parameter Initialization**

Define a way to update parameters

**Optimizer Selection**

**Loss function Selection**
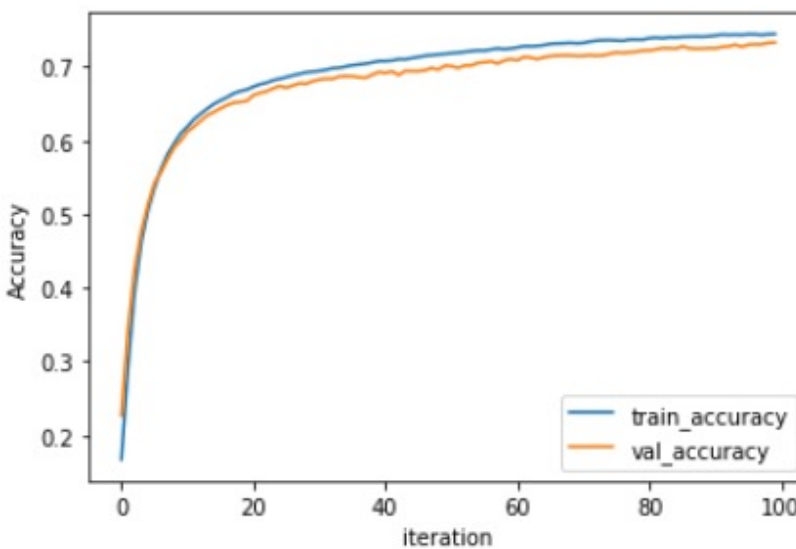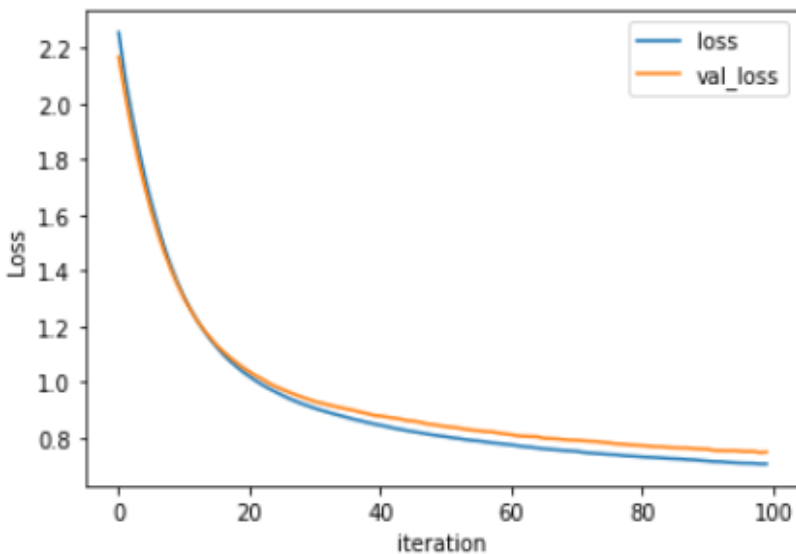
**Metric Selection**



56

# Summary and Discussion

- **Sigmoid and SGD**
- **W/o using normalization**

58

# Discussion

❖ **Sigmoid and SGD**

❖ **W/o using normalization**



28

28

784

flatten data

784 Nodes
Input layer

Fully
connect

• • •

7 hidden layers

Fully
connect

$z_1$

$z_{10}$

Softmax
activation

Output

10 Nodes
Output layer

59

# Discussion

## Dying ReLU

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

## Initialization

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