

Multi-layer Perception

Activation and Initialization

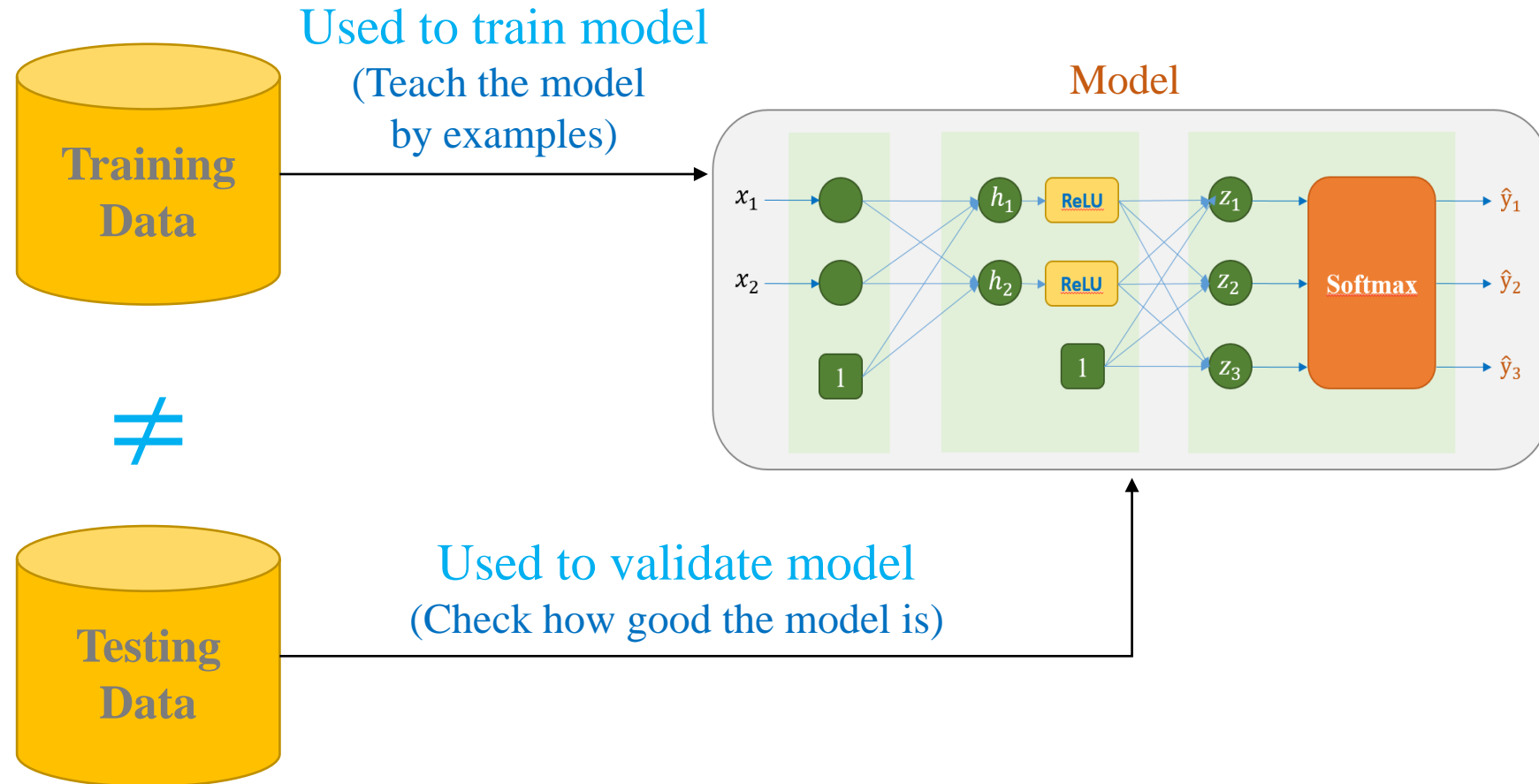
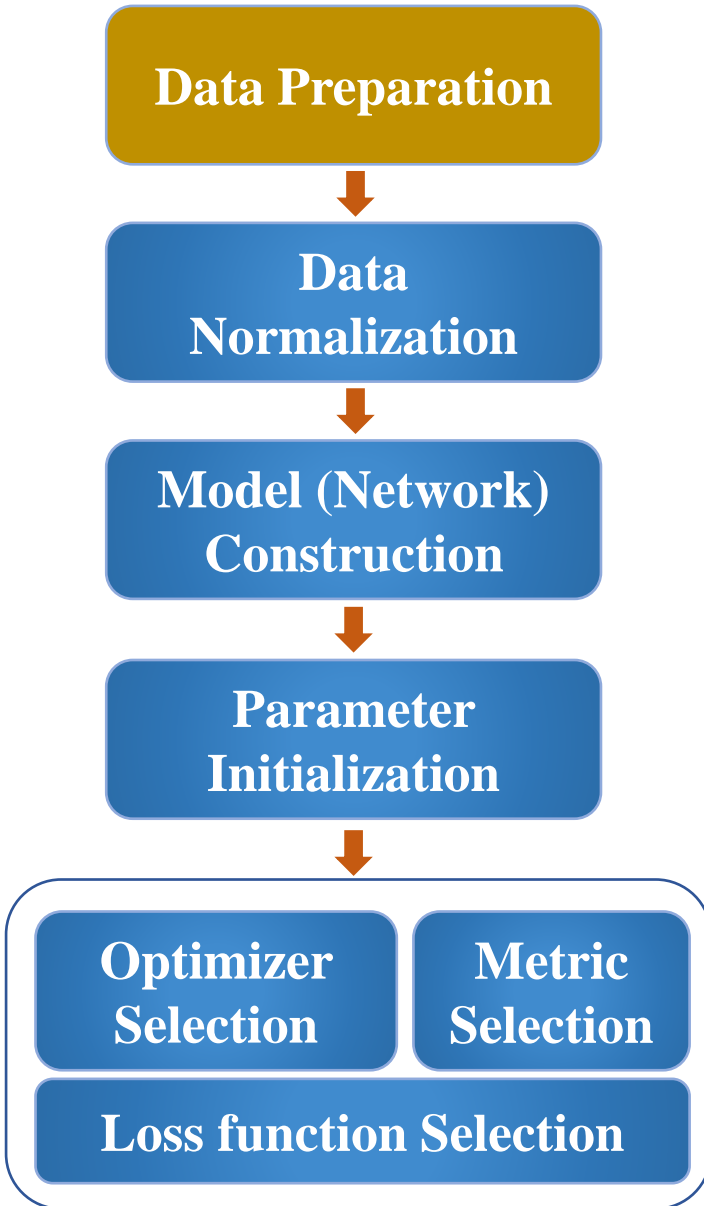
Quang-Vinh Dinh
Ph.D. in Computer Science

Outline

- **Pipeline Recommendation**
- **Data Normalization**
- **Activation Functions**
- **MLP Examples**
- **Initialization Methods**

To-do List for Training

Data Preparation



Data Normalization

Data Preparation



Data Normalization



Model (Network) Construction



Parameter Initialization



Optimizer Selection

Metric Selection

Loss function Selection

In Theory

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

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

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

In Pytorch

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

Normalize(*mean*, *std*)

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

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

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

Compute mean and std from data

```
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)

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)

# 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)
```

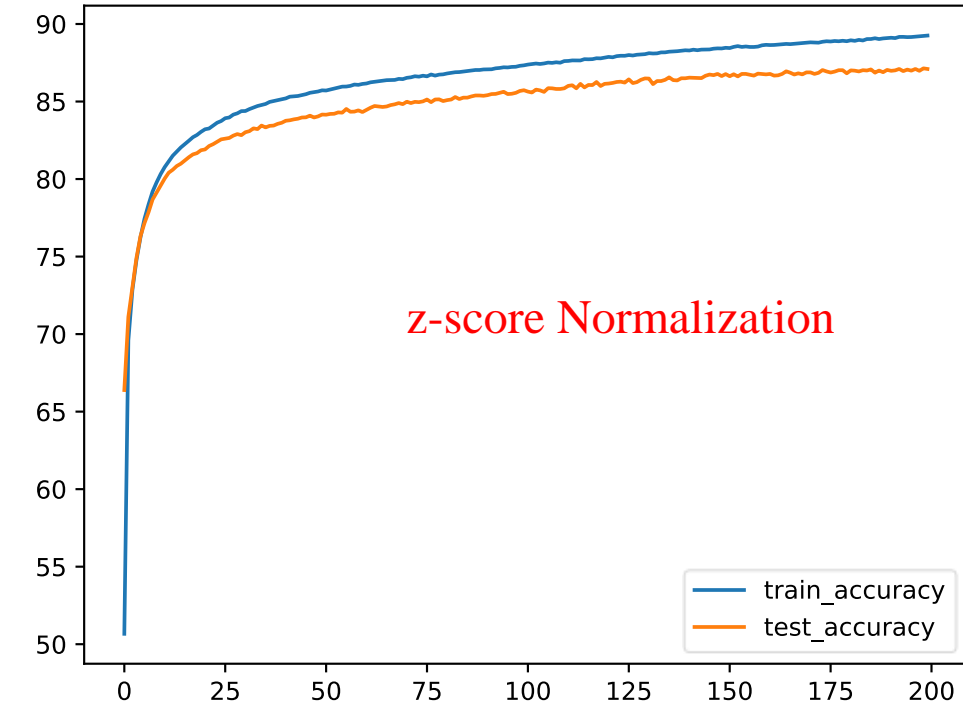
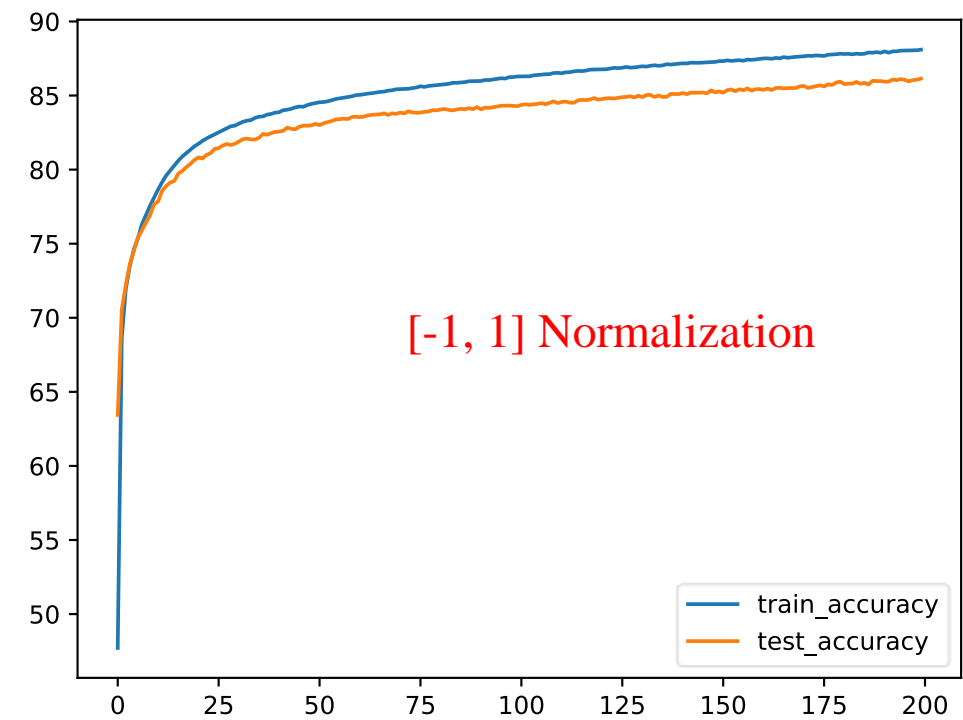
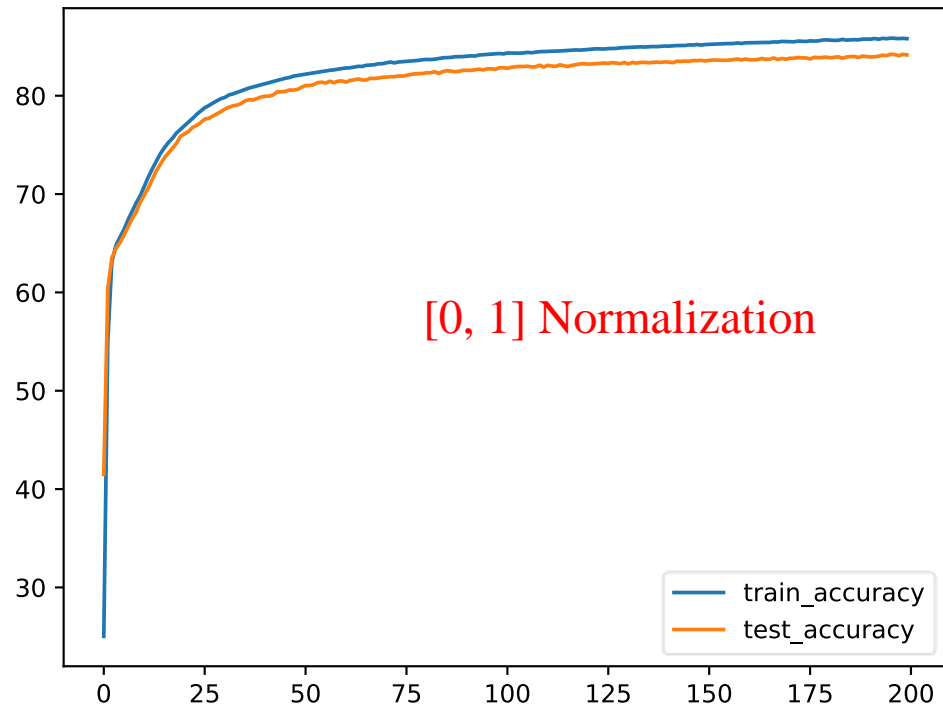
(a) [0, 1] Normalization

(b) [-1, 1] Normalization

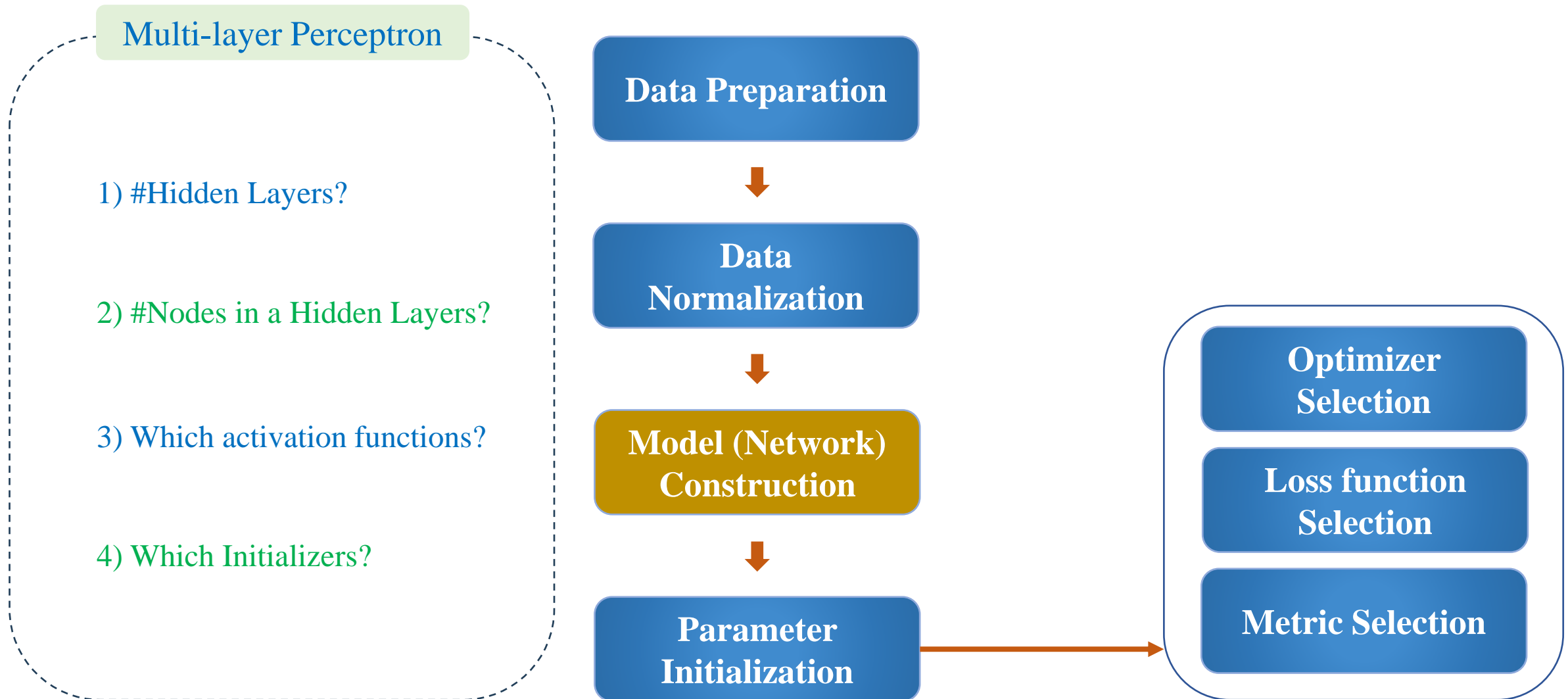
(c) z-score Normalization

Data Normalization

```
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)
```

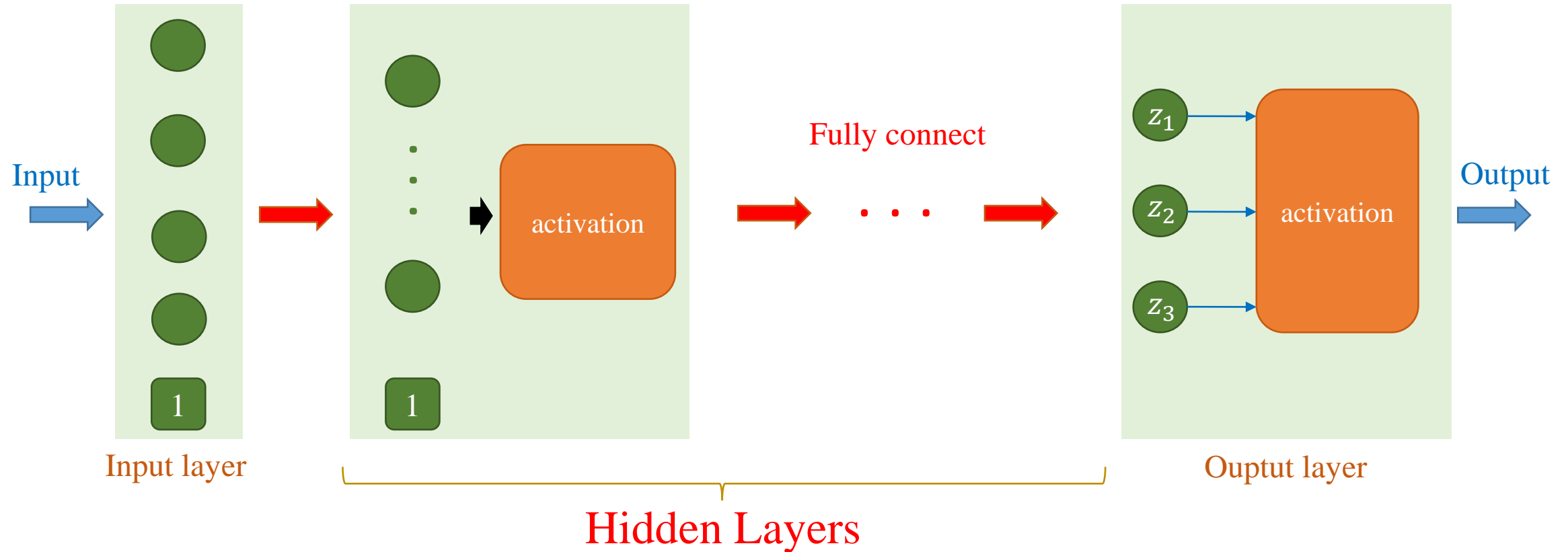


Training Pipeline



Training Pipeline

Model (Network) Construction

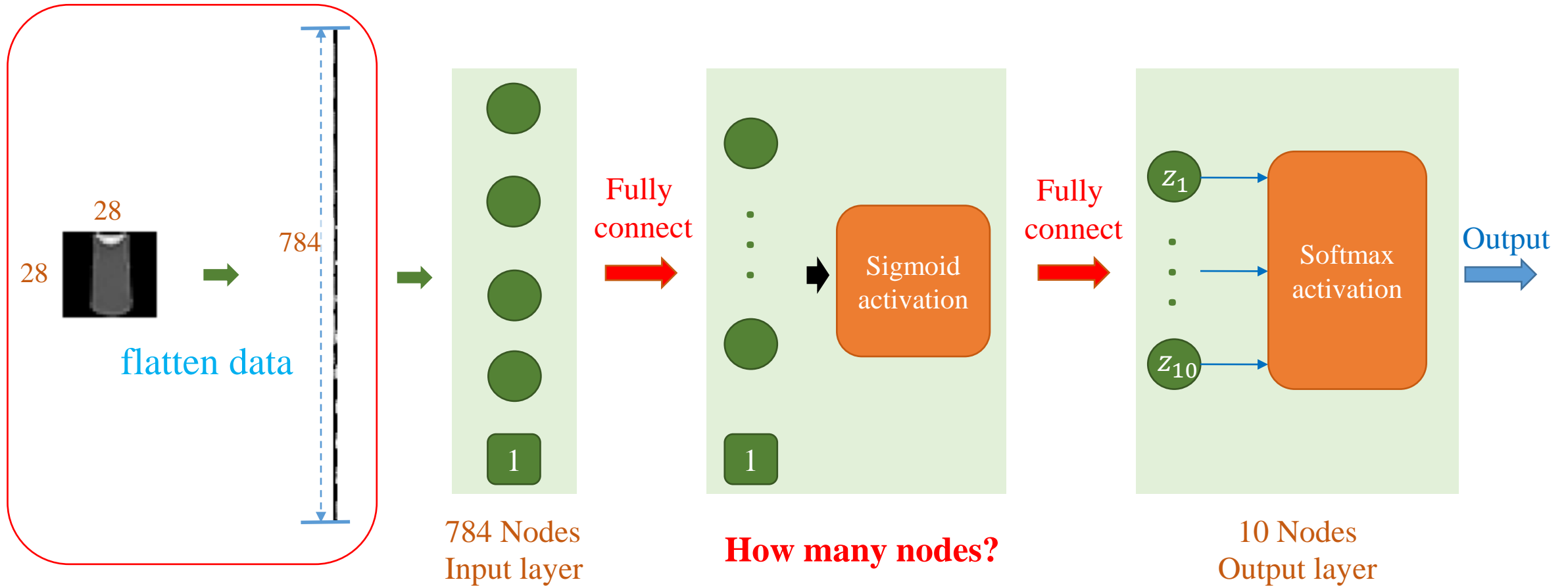


How many hidden layers?
How many nodes in a hidden layer?

Which activation function?
Which network components?

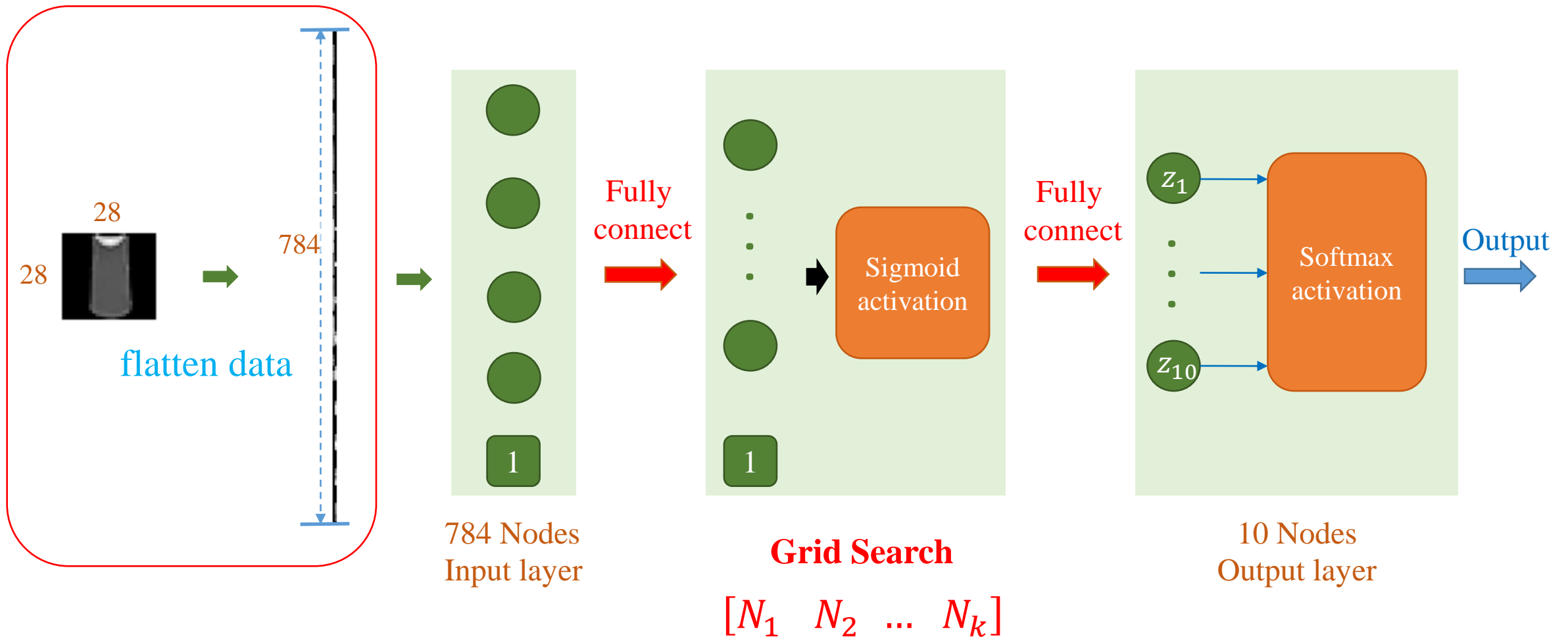
How many nodes?

Model (Network) Construction



How many nodes?

Model (Network) Construction

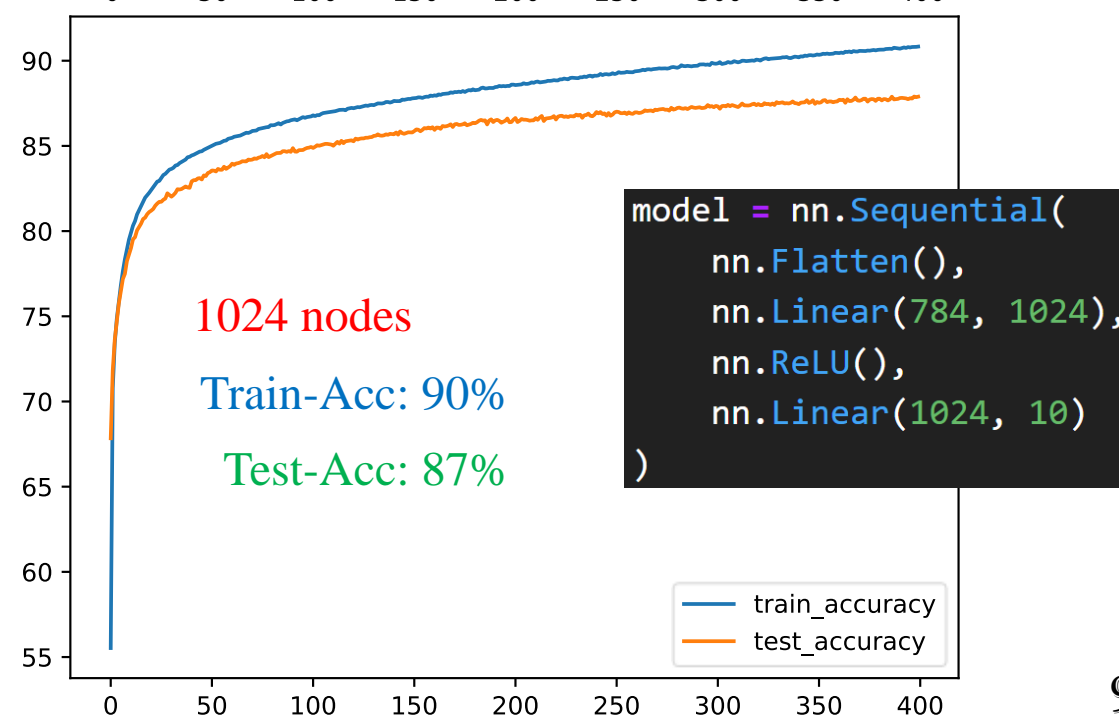
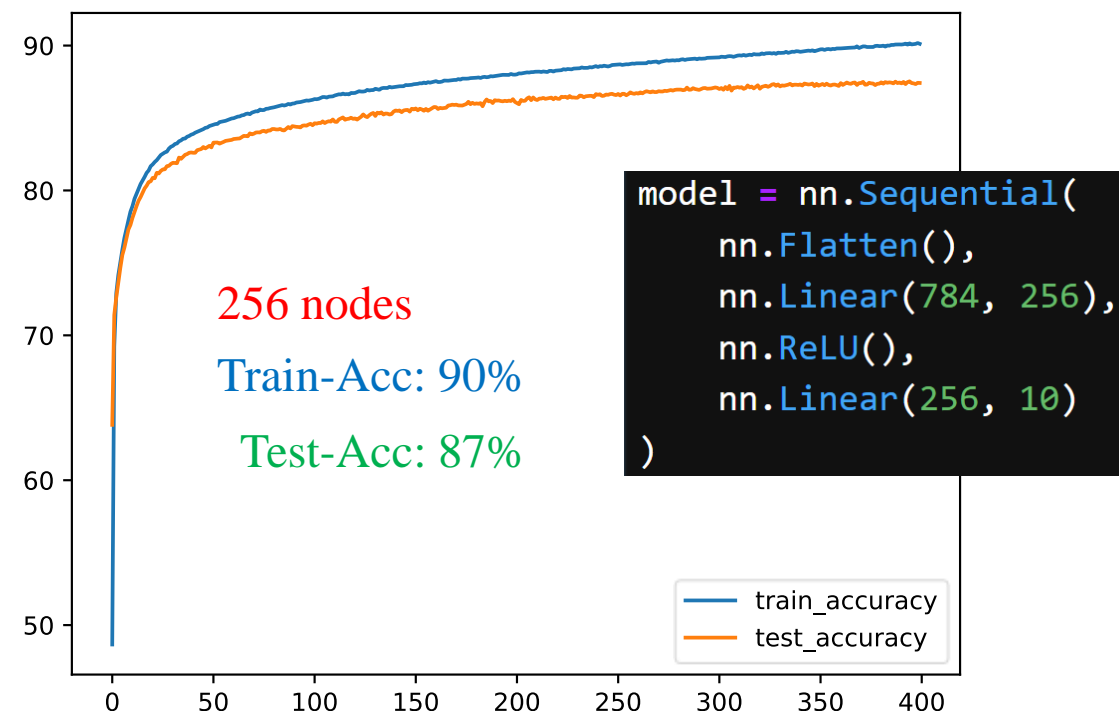
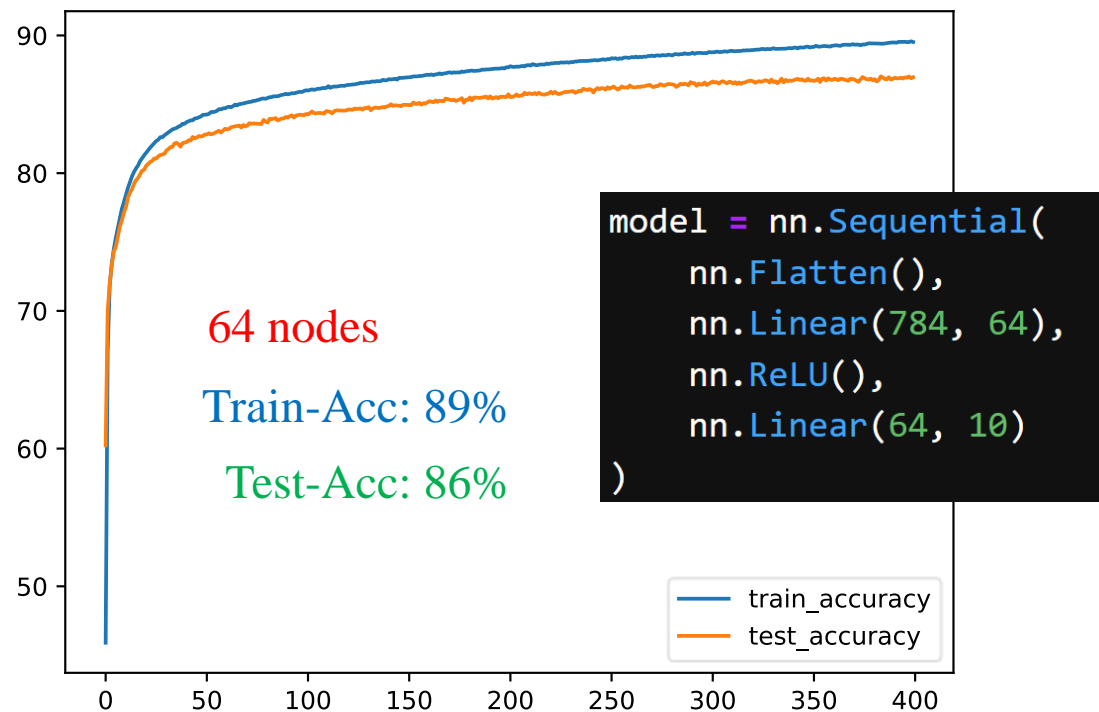


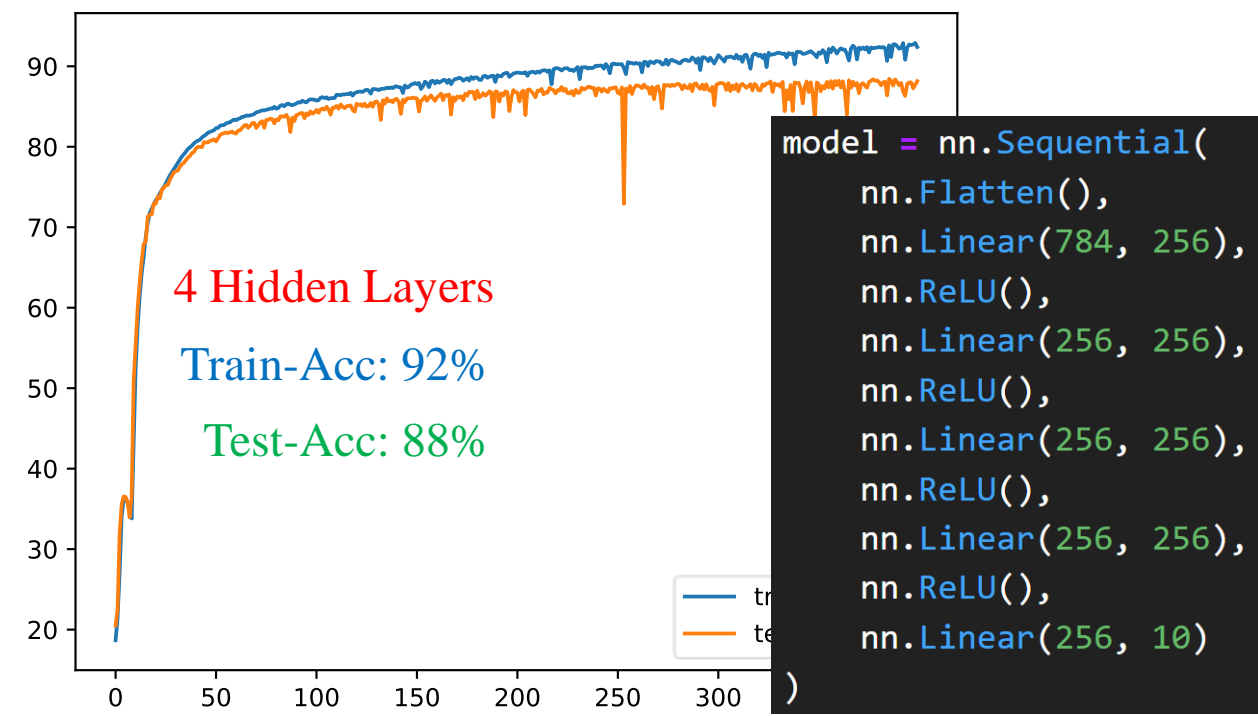
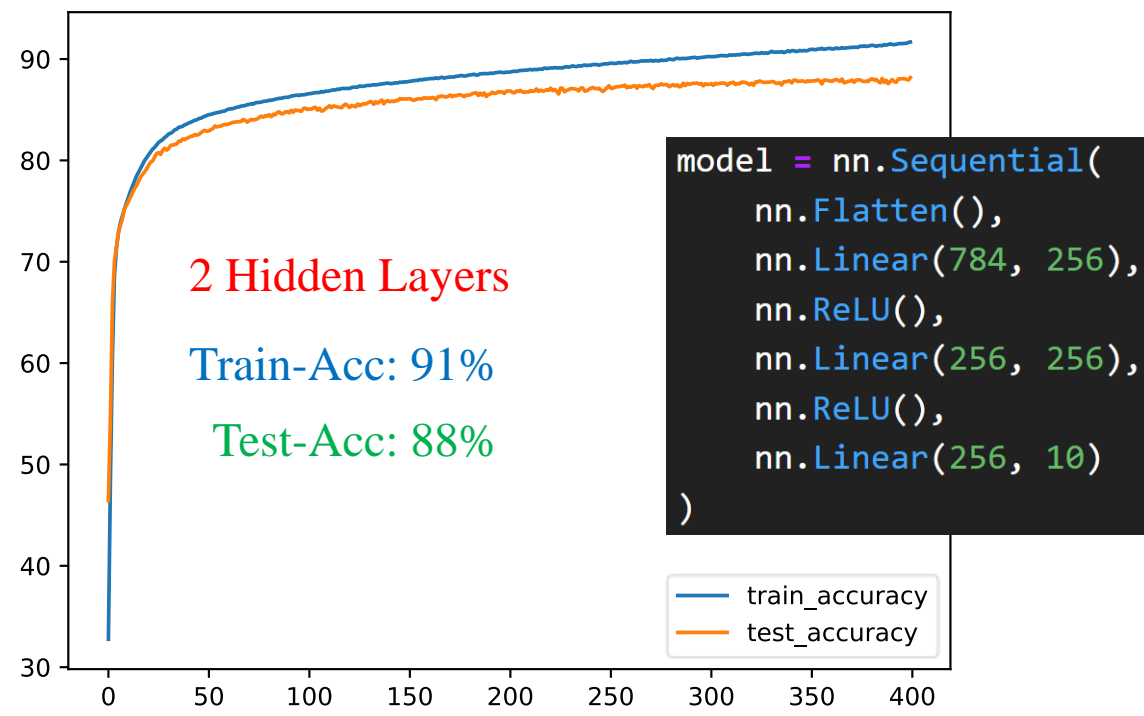
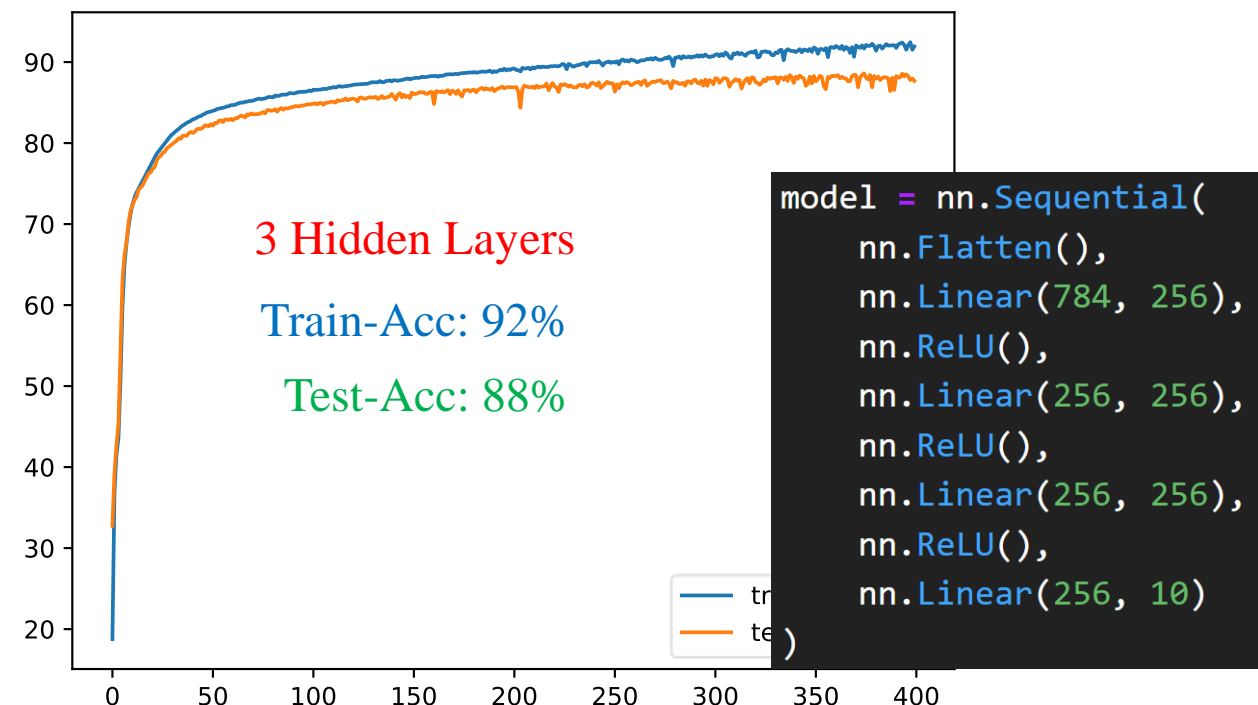
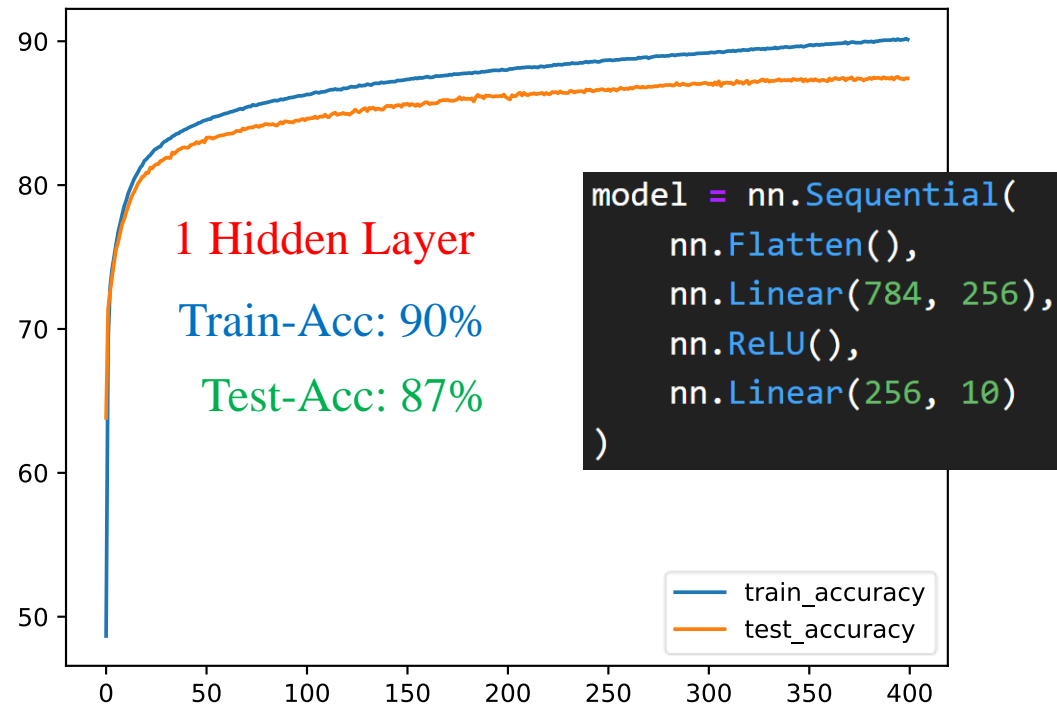
How many nodes?

[-1, 1] Normalization

Cross-entropy Loss

SGD with lr=0.01





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

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

$$\lambda \approx 1.0507$$

$$\alpha \approx 1.6733$$

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

2015

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

2001

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

2015

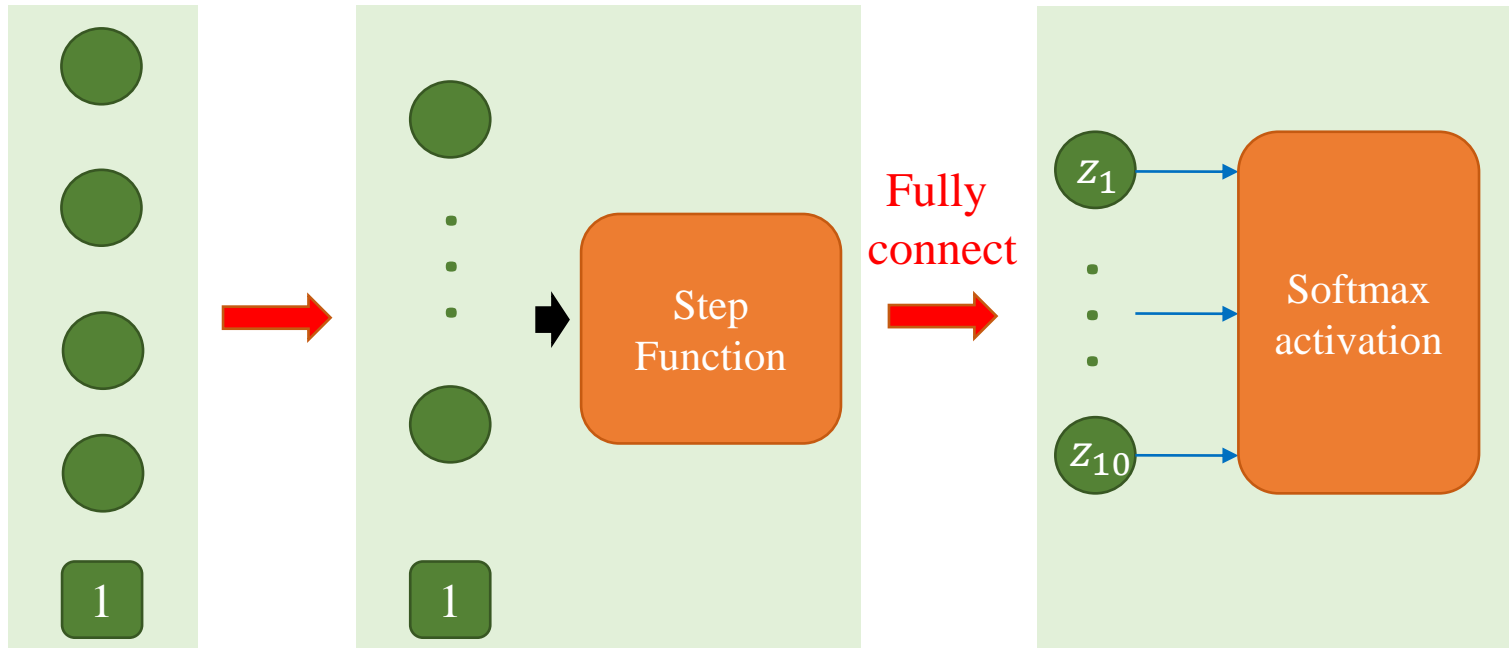
$$\text{PReLU}(x) = \begin{cases} \alpha x & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$$

2017

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

Activation Functions

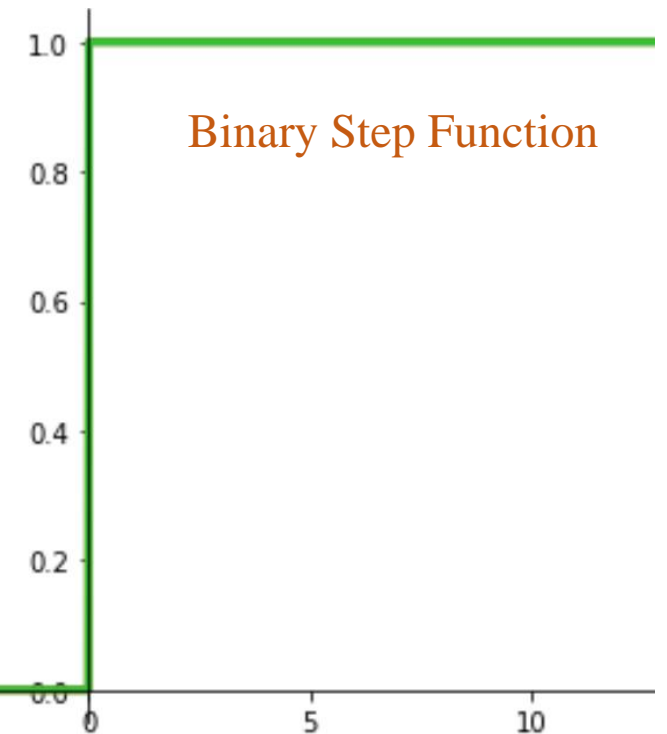
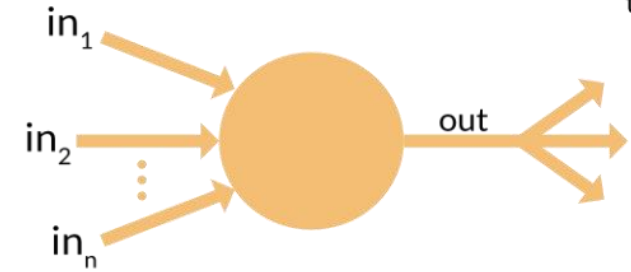
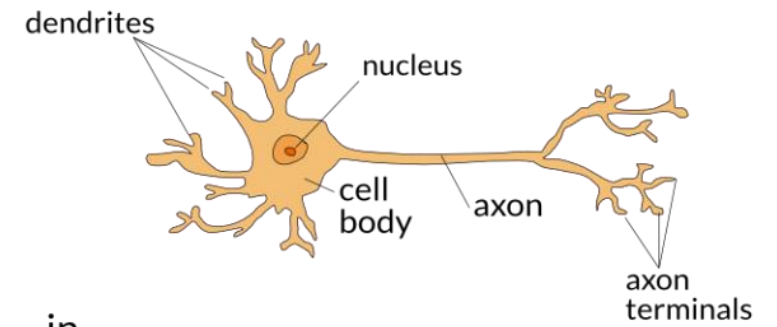
❖ Step function



Input layer

10 Nodes
Output layer

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



Activation Functions

❖ Sigmoid function

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

data =

1

5

-4

3

-2

data_a = sigmoid(data)

data_a =

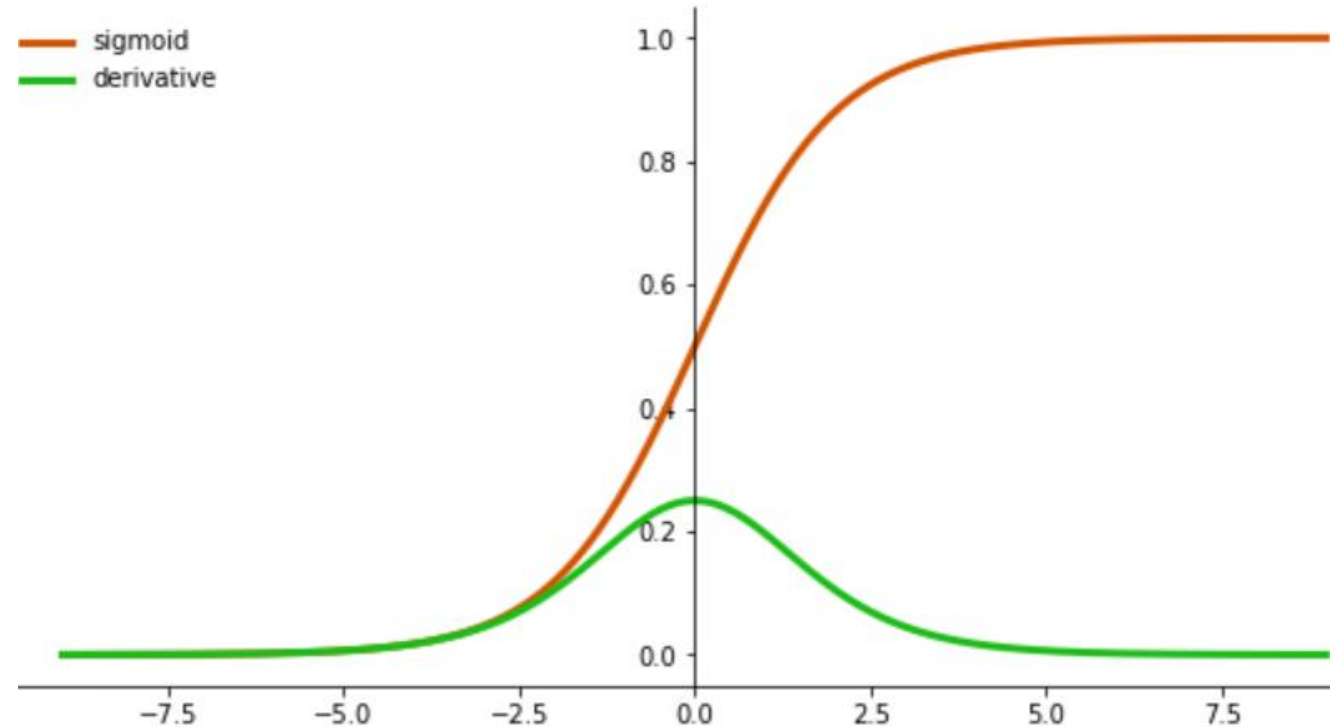
0.731

0.993

0.017

0.95

0.119



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

Activation Functions

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

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

Activation Functions

❖ Tanh function

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

data =

1

5

-4

3

-2

data_a = **tanh**(data)

data_a =

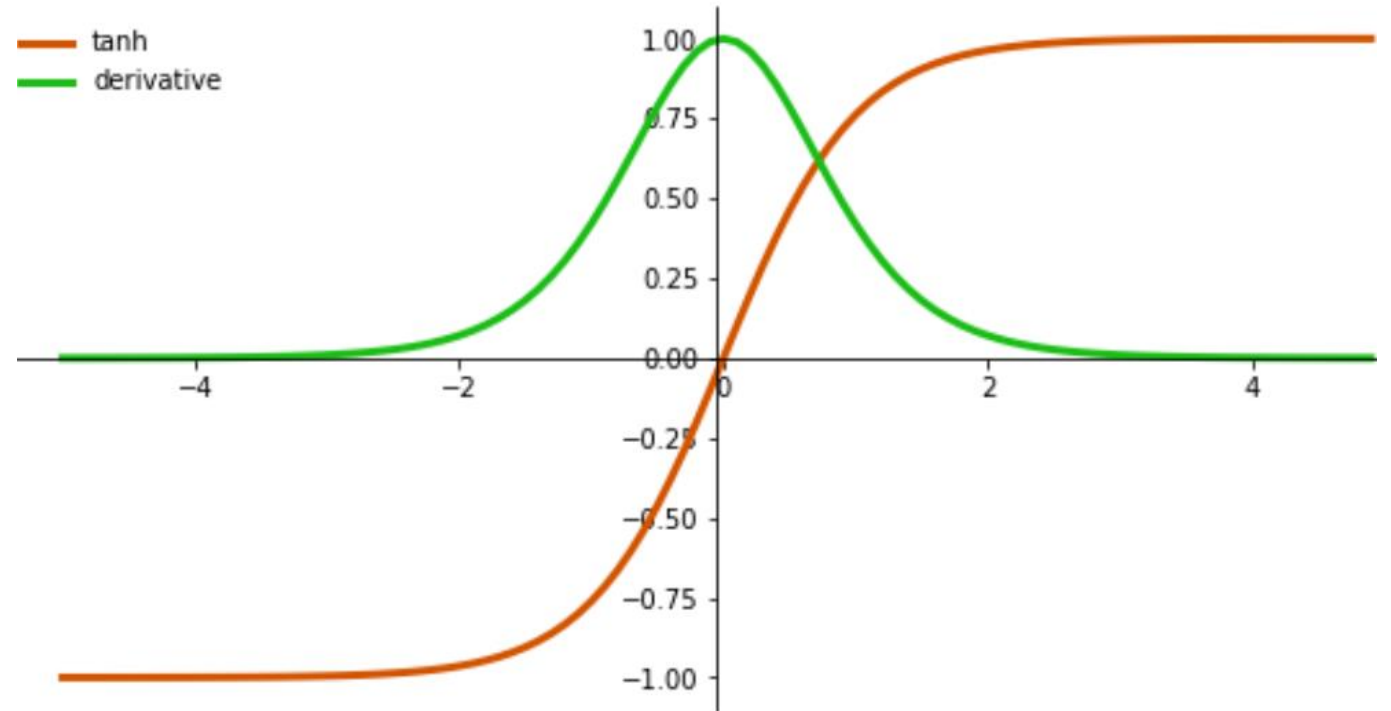
0.761

0.999

-0.999

0.995

-0.964



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

Activation Functions

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

$$\begin{aligned}\tanh'(x) &= \left(\frac{e^x - e^{-x}}{e^x + e^{-x}} \right)' = \frac{(e^x + e^{-x})(e^x + e^{-x}) - (e^x - e^{-x})(e^x - e^{-x})}{(e^x + e^{-x})^2} \\ &= \frac{(e^x + e^{-x})^2 - (e^x - e^{-x})^2}{(e^x + e^{-x})^2} \\ &= 1 - \left(\frac{e^x - e^{-x}}{e^x + e^{-x}} \right)^2 = 1 - \tanh^2(x)\end{aligned}$$

Activation Functions

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

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

Activation Functions

❖ Softplus function

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

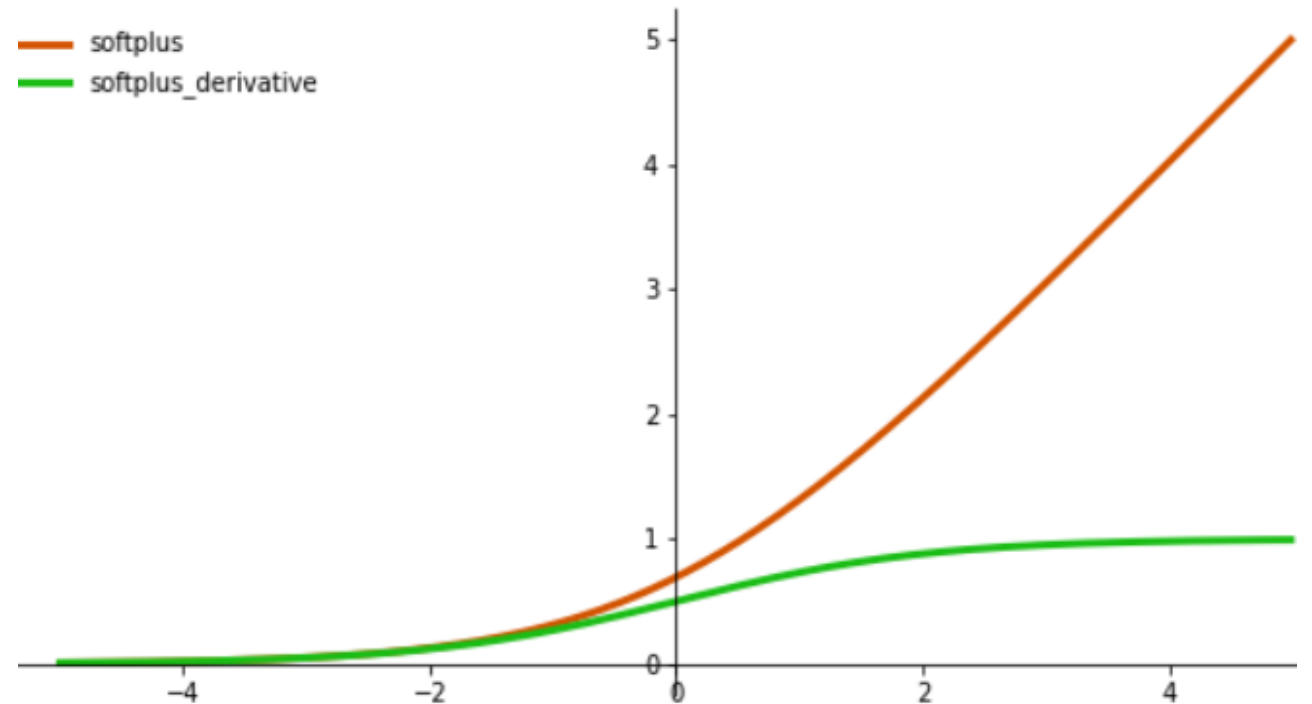
data =

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

data_a = softplus(data)

data_a =

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



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

Activation Functions

❖ ReLU function

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

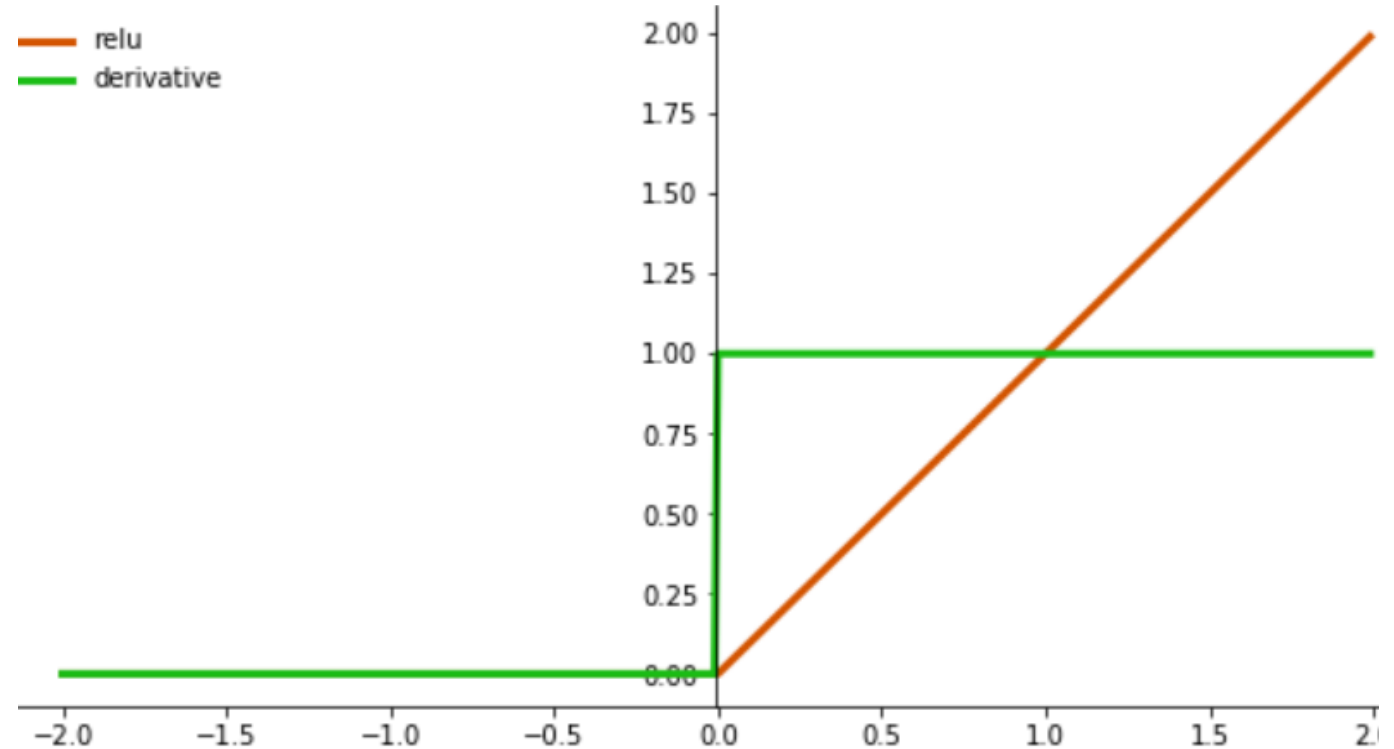
data =

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

data_a = **ReLU**(data)

data_a =

1	5	0	3	0
---	---	---	---	---



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

Activation Functions

❖ LeakyReLU function

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

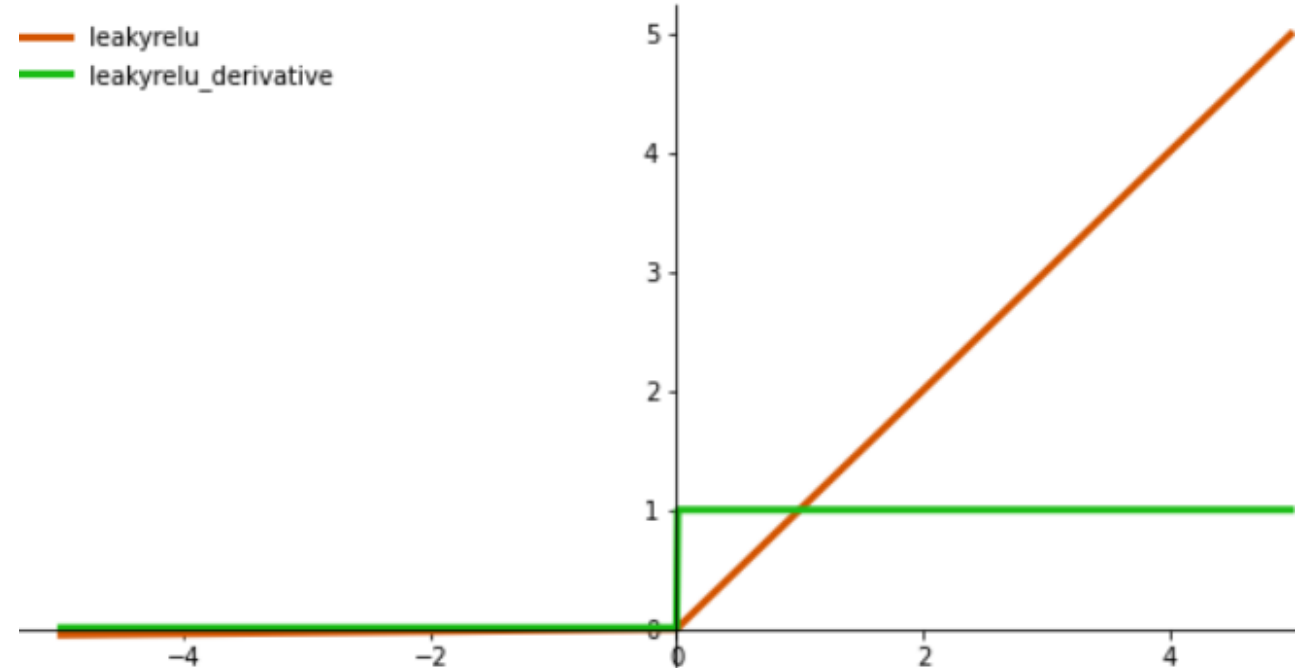
data =

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

data_a = leakyrelu(data)

data_a =

1	5	-0.04	3	-0.02
---	---	-------	---	-------



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

Activation Functions

❖ ELU function

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

$$\alpha = 0.1$$

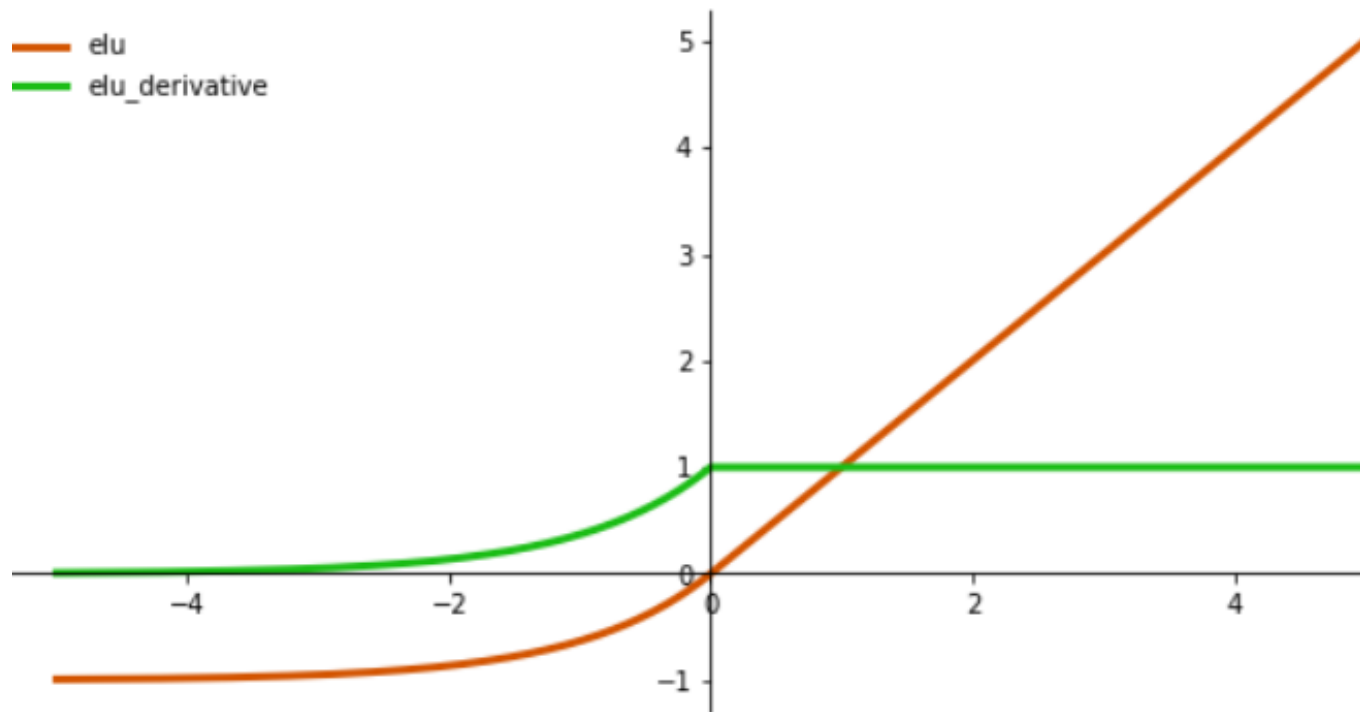
data =

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

data_a = ELU(data)

data_a =

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



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

Activation Functions

❖ PReLU function

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

$$\alpha = 0.1$$

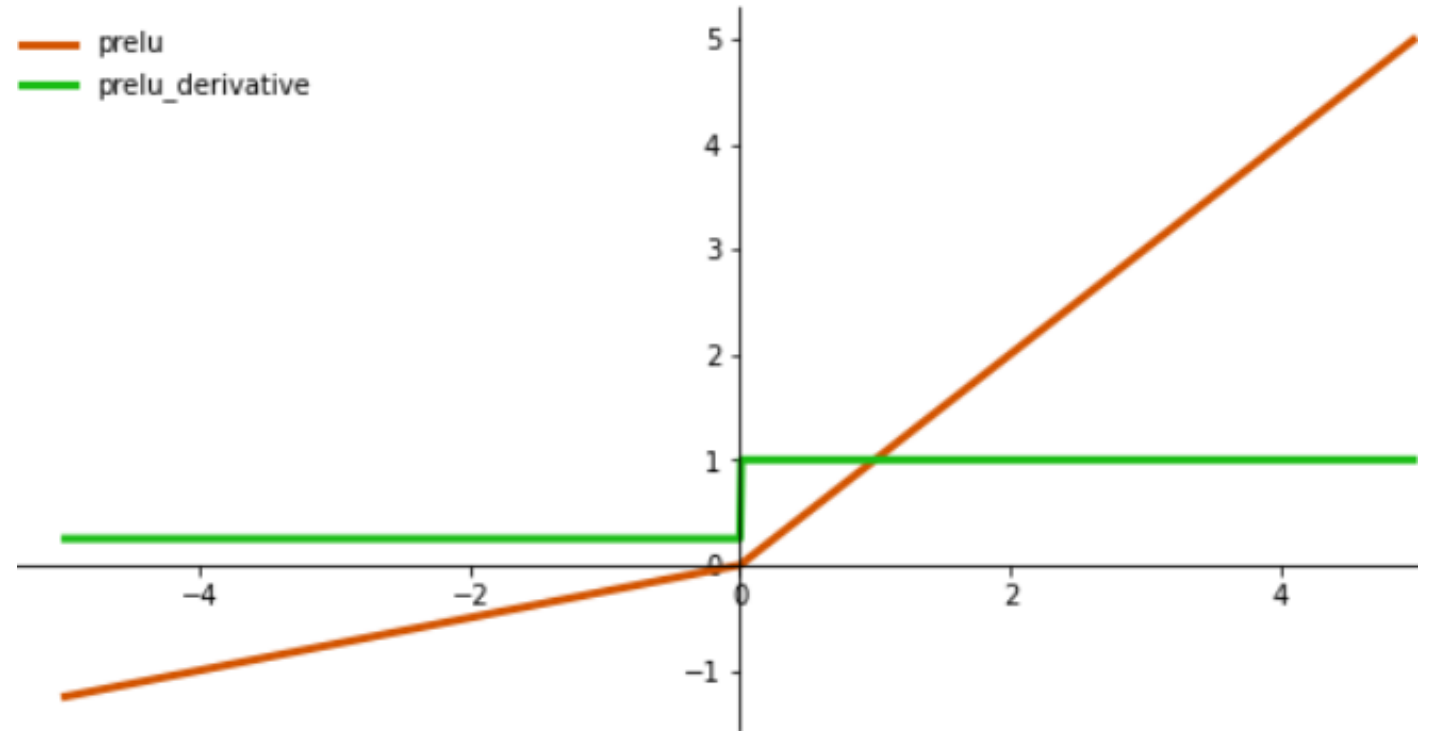
data =

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

data_a = PRELU(data)

data_a =

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



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

Activation Functions

❖ Swish function

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

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

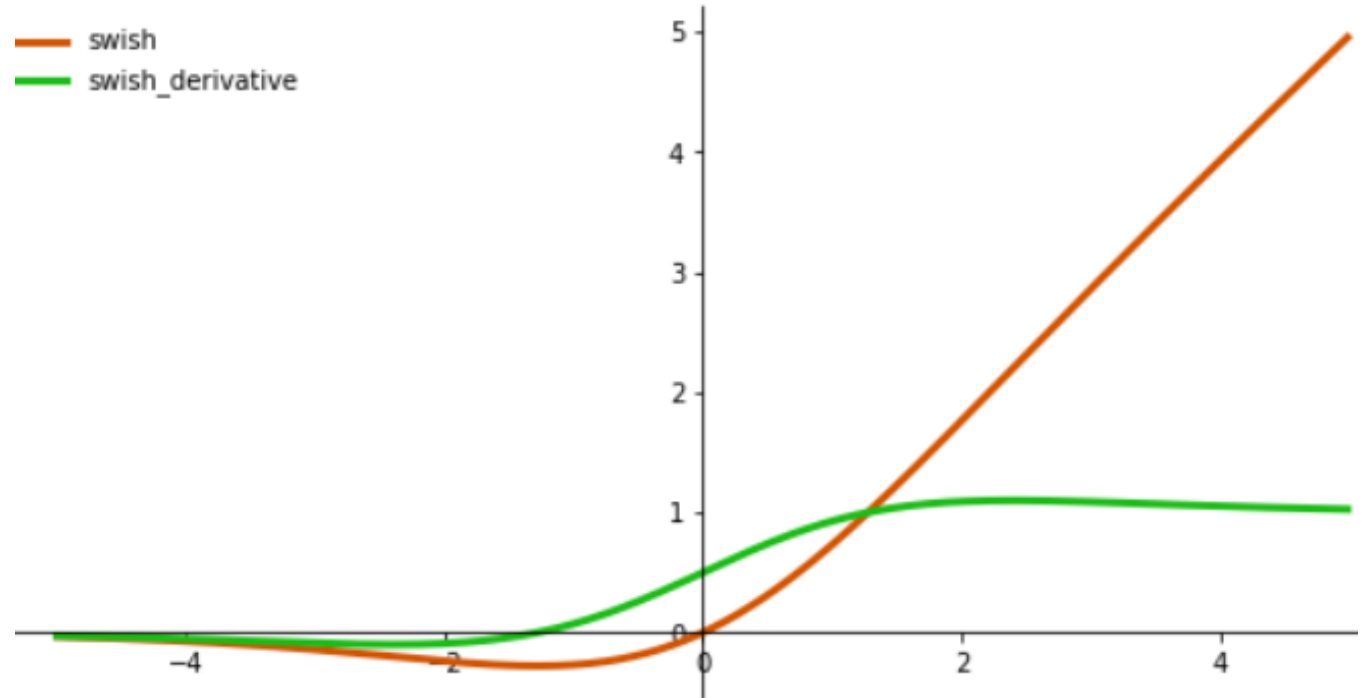
data =

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

data_a = swish(data)

data_a =

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



$$\text{swish}'(x) = \text{swish}(x) + \sigma(x) (1 - \text{swish}(x))$$

Activation Functions

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

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

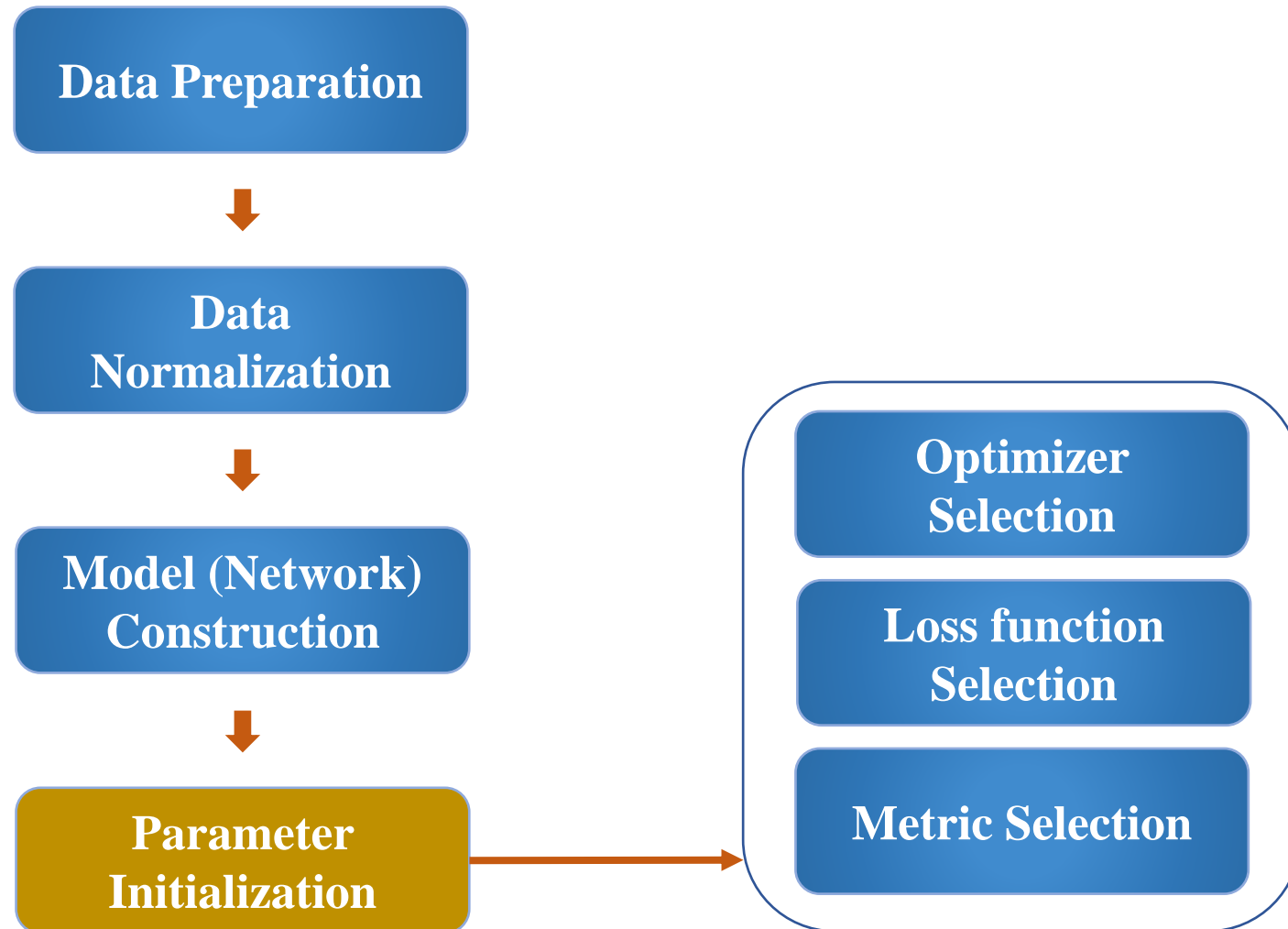
$$\begin{aligned}\text{swish}'(x) &= (x \sigma(x))' = (x)' \sigma(x) + x(\sigma(x))' \\ &= \sigma(x) + x \sigma(x) (1 - \sigma(x)) \\ &= \sigma(x) + x \sigma(x) - x \sigma(x)^2 \\ &= x \sigma(x) + \sigma(x)(1 - x \sigma(x)) \\ &= \text{swish}(x) + \sigma(x) (1 - \text{swish}(x))\end{aligned}$$

Outline

- **Pipeline Recommendation**
- **Data Normalization**
- **Activation Functions**
- **MLP Examples**
- **Initialization Methods**

To-do List for Training

Train a model

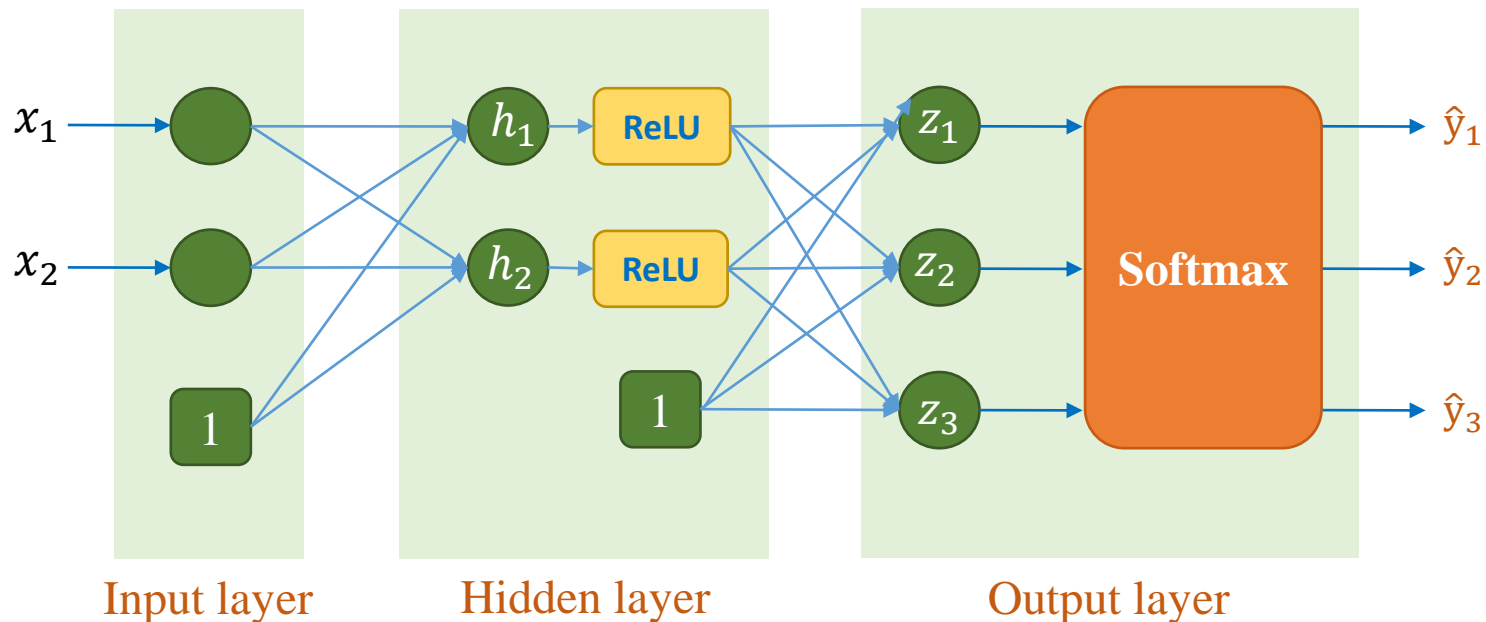


MLP Example 1

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

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

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



$$\mathbf{W}_h = [\mathbf{W}_{h1} \quad \mathbf{W}_{h2}]$$

$$= \begin{bmatrix} 0.0 & 0.0 \\ 0.86 & -1.04 \\ 0.41 & -0.65 \end{bmatrix}$$

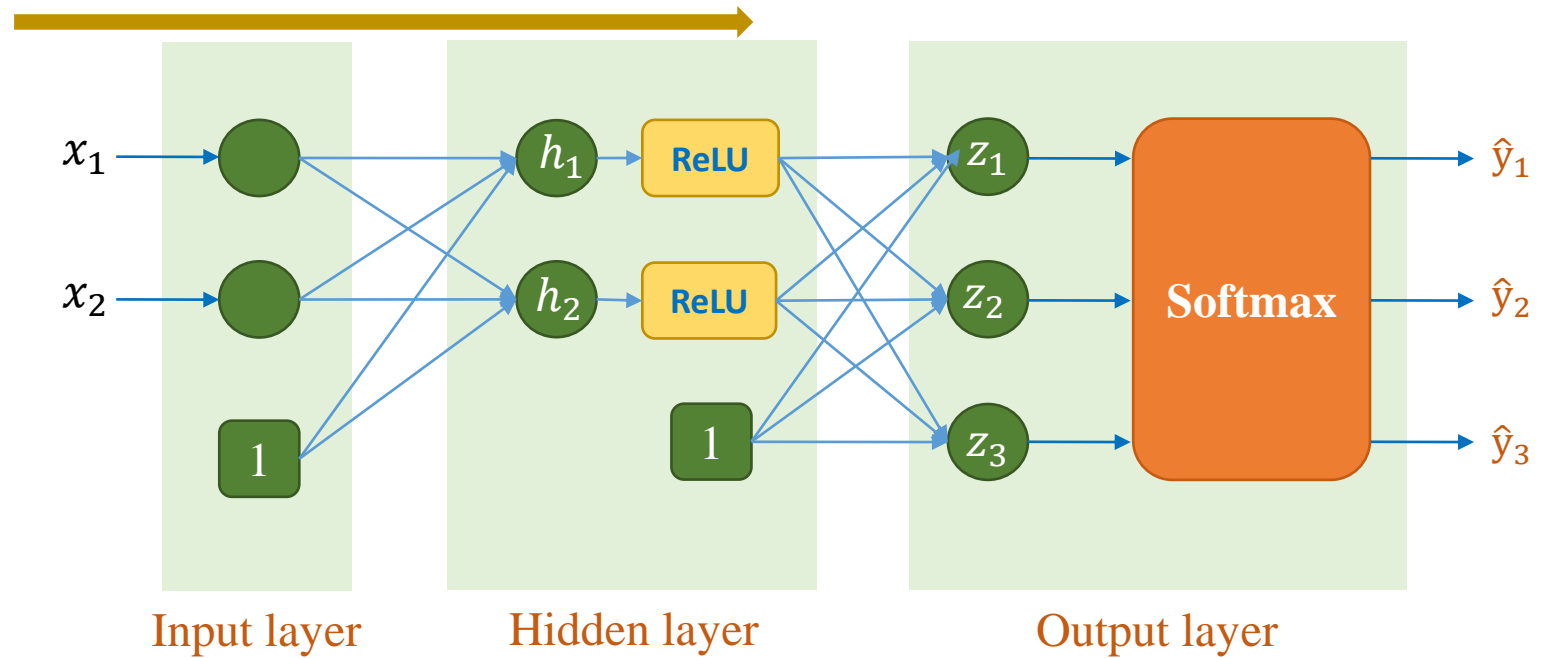
$$\mathbf{W}_z = [\mathbf{W}_{z1} \quad \mathbf{W}_{z2} \quad \mathbf{W}_{z3}]$$

$$= \begin{bmatrix} 0.0 & 0.0 & 0.0 \\ 0.32 & 0.25 & 0.14 \\ -0.47 & -1.06 & 0.063 \end{bmatrix}$$

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

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

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



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

$$\begin{aligned} \mathbf{W}_h &= [\mathbf{W}_{h1} \quad \mathbf{W}_{h2}] & \mathbf{W}_z &= [\mathbf{W}_{z1} \quad \mathbf{W}_{z2} \quad \mathbf{W}_{z3}] \\ &= \begin{bmatrix} 0.0 & 0.0 \\ 0.86 & -1.04 \\ 0.41 & -0.65 \end{bmatrix} & &= \begin{bmatrix} 0.0 & 0.0 & 0.0 \\ 0.32 & 0.25 & 0.14 \\ -0.47 & -1.06 & 0.063 \end{bmatrix} \end{aligned}$$

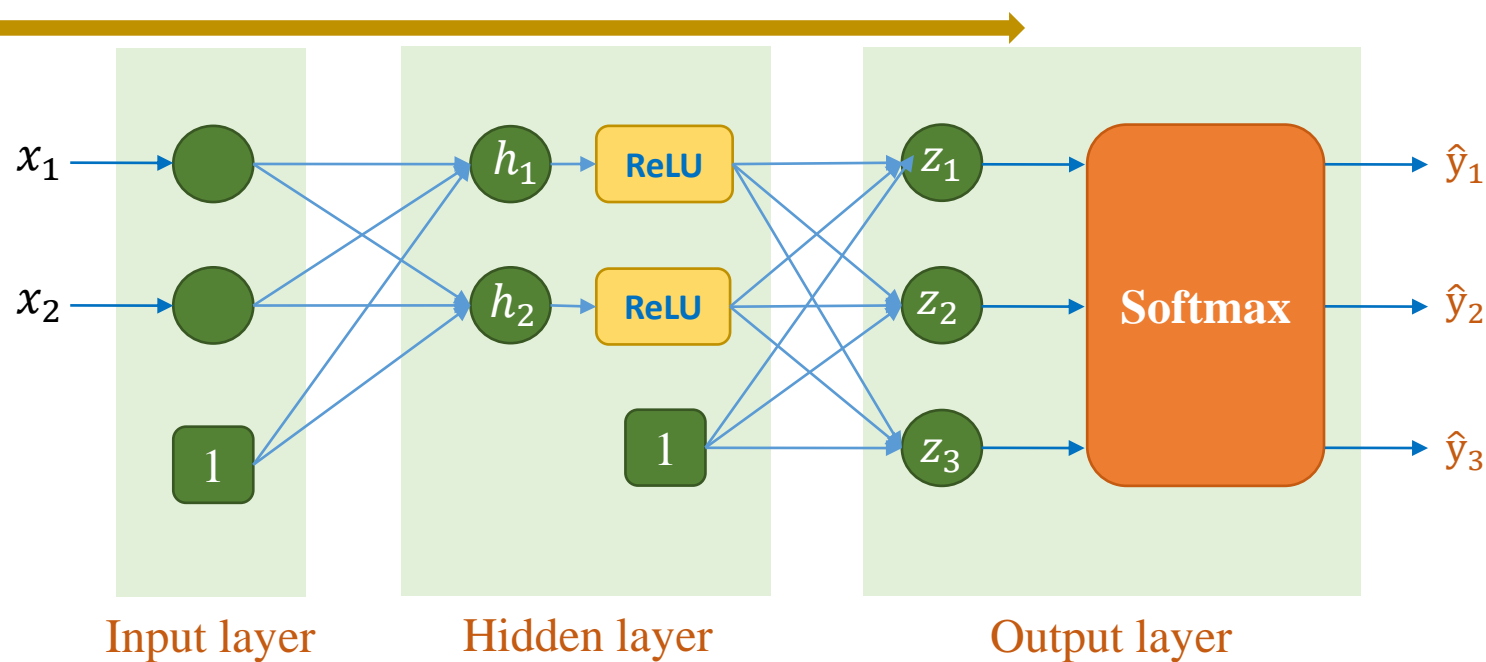
$$\text{ReLU}(\mathbf{h}) = \begin{bmatrix} 1.373 & 0 \\ 4.708 & 0 \\ 5.731 & 0 \end{bmatrix}$$

$$[\mathbf{1} \quad \text{ReLU}(\mathbf{h})] = \begin{bmatrix} 1 & 1.373 & 0 \\ 1 & 4.708 & 0 \\ 1 & 5.731 & 0 \end{bmatrix}$$

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

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

$$\begin{aligned} \mathbf{z} = [\mathbf{1} \quad \text{ReLU}(\mathbf{h})] \mathbf{W}_z &= \begin{bmatrix} 1 & 1.373 & 0 \\ 1 & 4.708 & 0 \\ 1 & 5.731 & 0 \end{bmatrix} \begin{bmatrix} 0.0 & 0.0 & 0.0 \\ 0.32 & 0.25 & 0.14 \\ -0.47 & -1.06 & 0.063 \end{bmatrix} \\ &= \begin{bmatrix} 0.439 & 0.356 & 0.195 \\ 1.507 & 1.220 & 0.670 \\ 1.835 & 1.485 & 0.816 \end{bmatrix} \end{aligned}$$



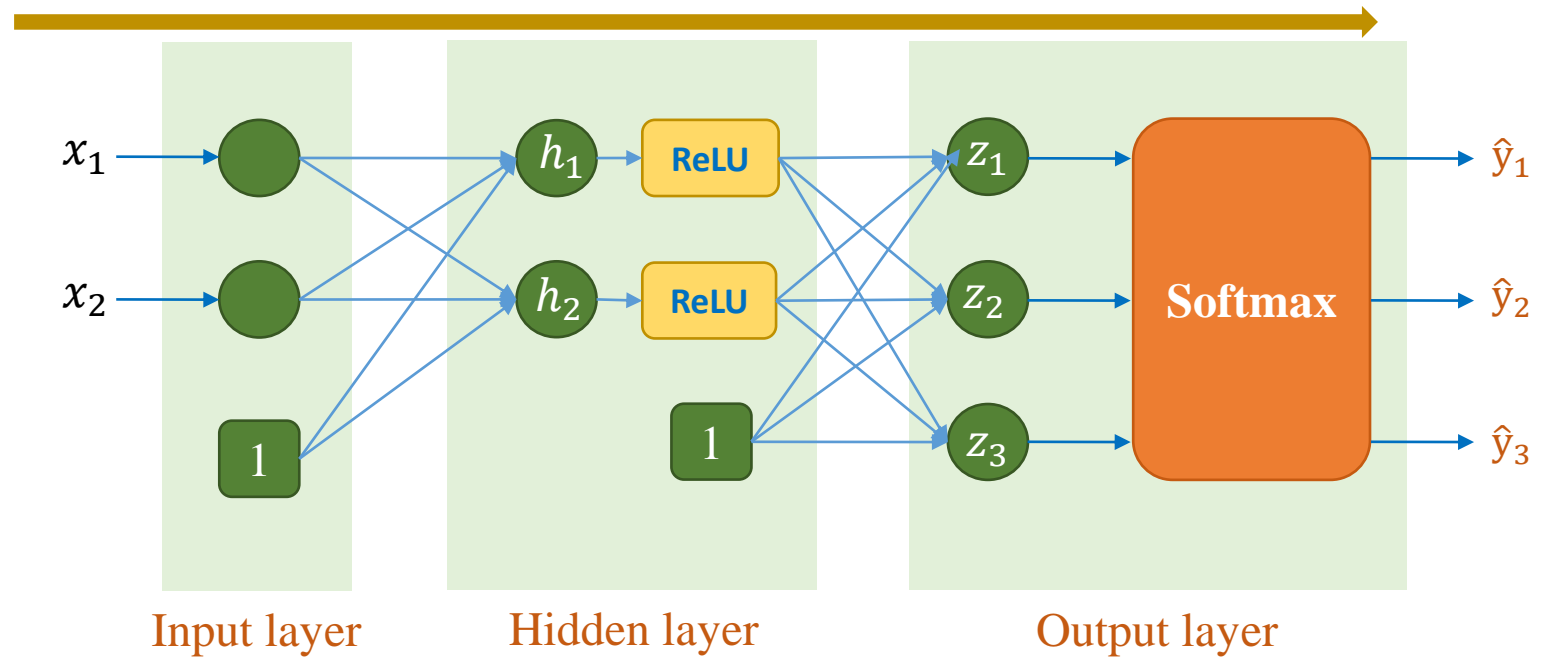
$$\begin{aligned} \mathbf{W}_h &= [\mathbf{W}_{h1} \quad \mathbf{W}_{h2}] & \mathbf{W}_z &= [\mathbf{W}_{z1} \quad \mathbf{W}_{z2} \quad \mathbf{W}_{z3}] \\ &= \begin{bmatrix} 0.0 & 0.0 \\ 0.86 & -1.04 \\ 0.41 & -0.65 \end{bmatrix} & &= \begin{bmatrix} 0.0 & 0.0 & 0.0 \\ 0.32 & 0.25 & 0.14 \\ -0.47 & -1.06 & 0.063 \end{bmatrix} \end{aligned}$$

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

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

$$\text{loss} = 1.269$$

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



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

$$\mathbf{W}_h = [\mathbf{W}_{h1} \quad \mathbf{W}_{h2}]$$

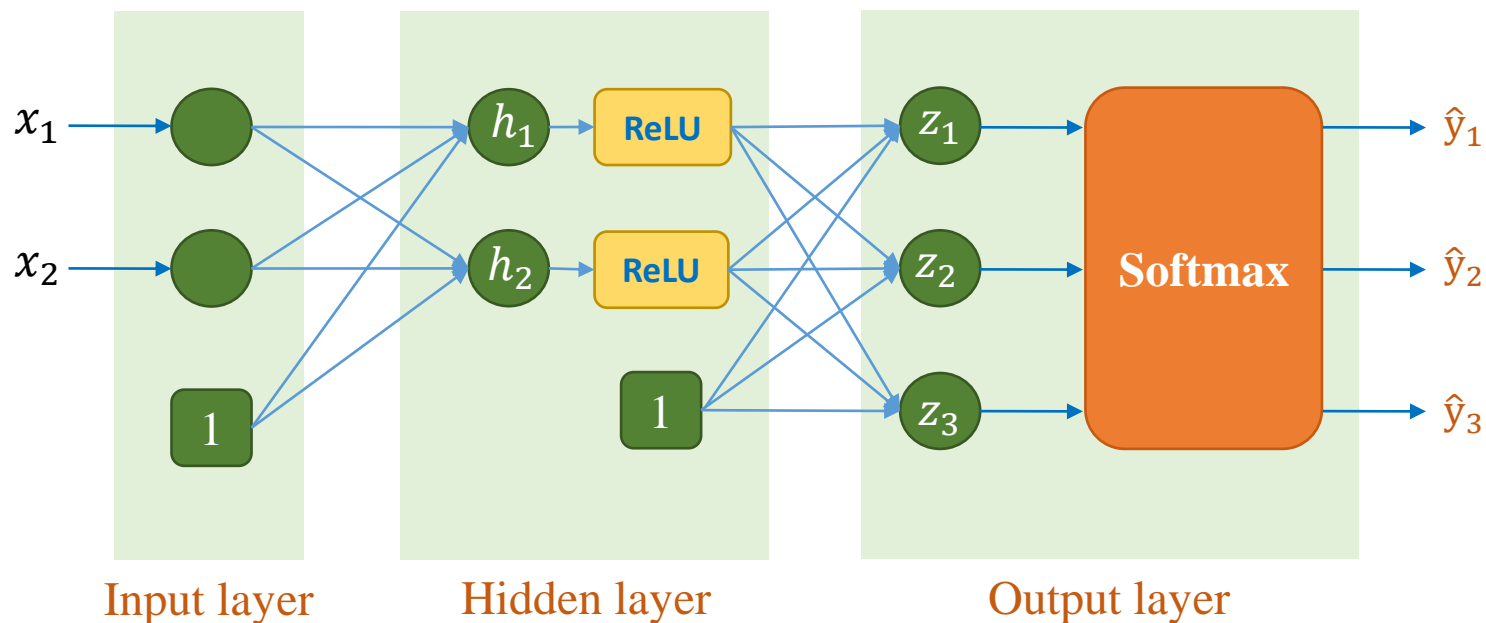
$$= \begin{bmatrix} 0.0 & 0.0 \\ 0.86 & -1.04 \\ 0.41 & -0.65 \end{bmatrix}$$

$$\mathbf{W}_z = [\mathbf{W}_{z1} \quad \mathbf{W}_{z2} \quad \mathbf{W}_{z3}]$$

$$= \begin{bmatrix} 0.0 & 0.0 & 0.0 \\ 0.32 & 0.25 & 0.14 \\ -0.47 & -1.06 & 0.063 \end{bmatrix}$$

Example 2 - Dying ReLU

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



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

$$m = [m_1 \quad m_2]$$

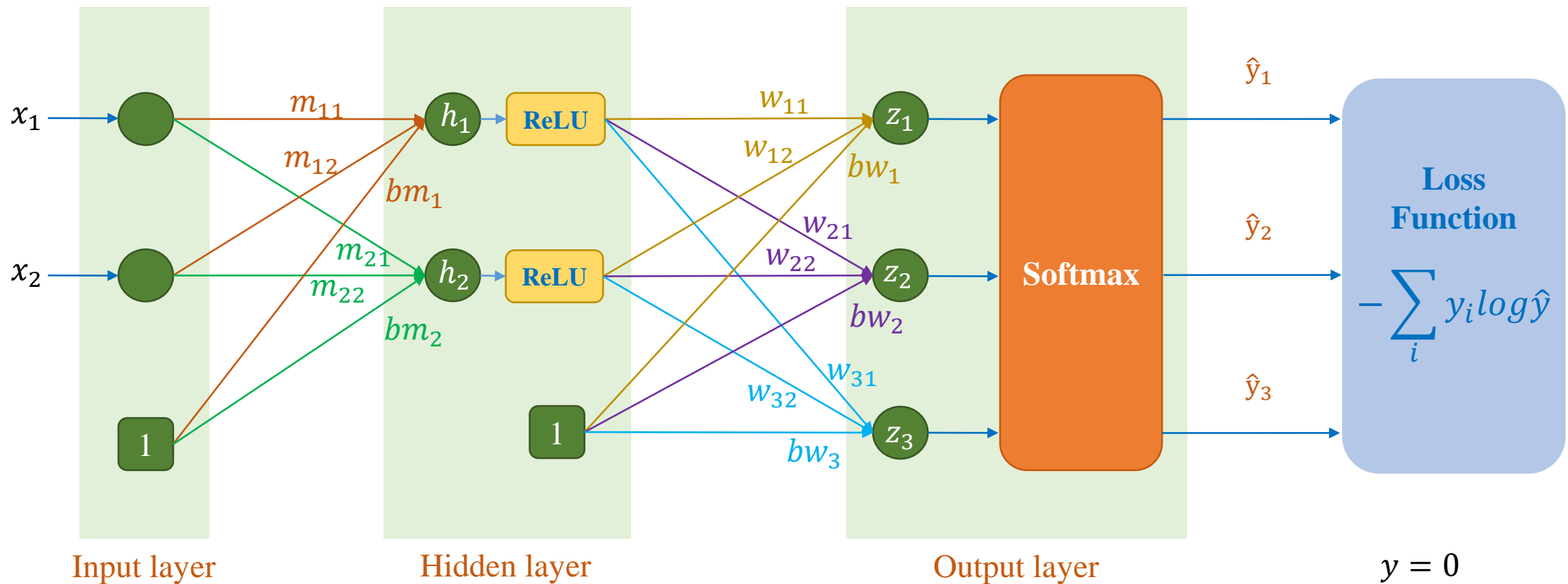
$$= \begin{bmatrix} 0.86 & -1.04 \\ 0.41 & -0.65 \end{bmatrix}$$

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

$$w = [w_1 \quad w_2 \quad w_3]$$

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

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



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

$$m = [\mathbf{m}_1 \quad \mathbf{m}_2]$$

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

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

$$w = [\mathbf{w}_1 \quad \mathbf{w}_2 \quad \mathbf{w}_3]$$

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

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

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

Forward pass

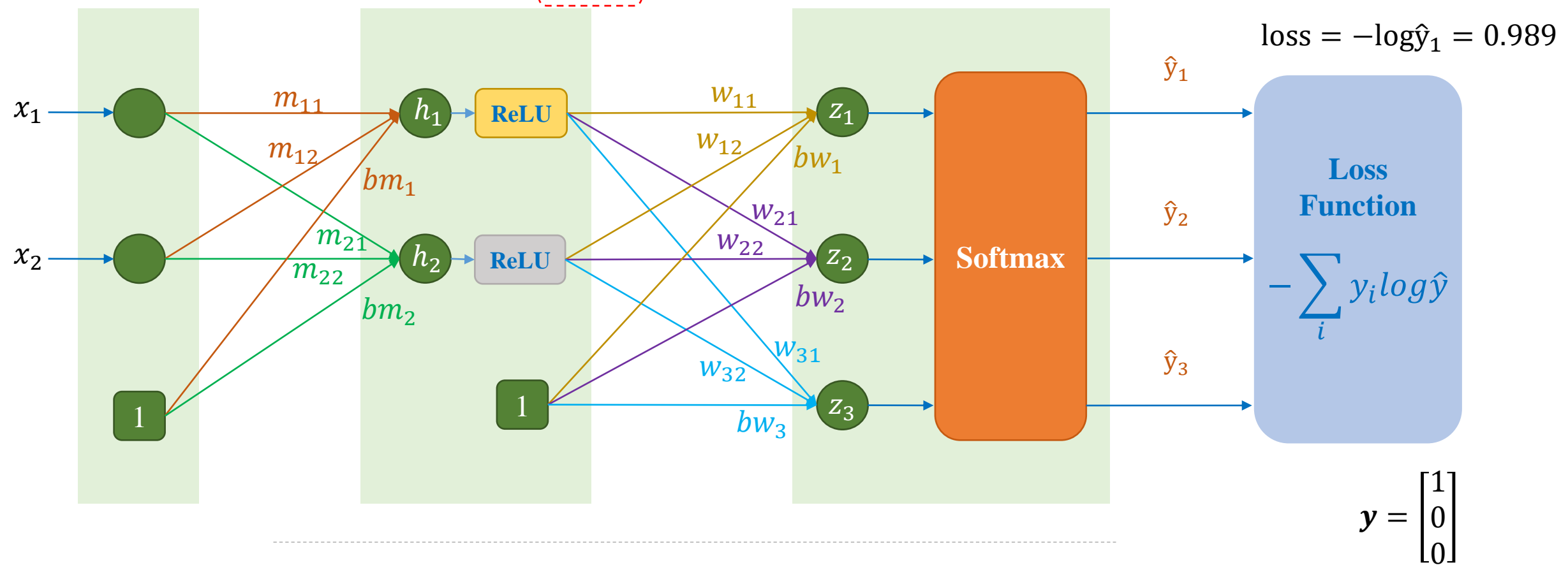
$$h = \begin{bmatrix} 1.372 \\ -1.68 \end{bmatrix}$$

zero value

$$\text{ReLU} = \begin{bmatrix} 1.372 \\ 0.0 \end{bmatrix}$$

$$z = \begin{bmatrix} 0.439 \\ 0.343 \\ 0.192 \end{bmatrix}$$

$$\hat{y} = \begin{bmatrix} 0.372 \\ 0.338 \\ 0.290 \end{bmatrix}$$

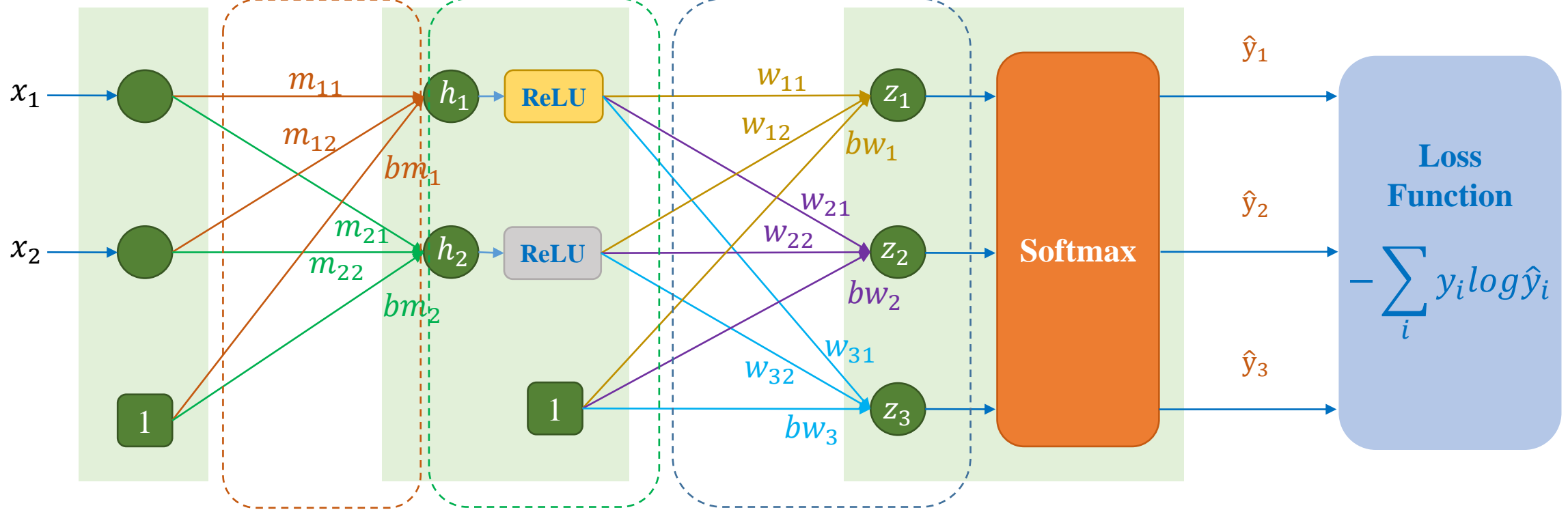


$$m = \begin{bmatrix} m_1 & m_2 \\ 0.86 & -1.04 \\ 0.41 & -0.65 \end{bmatrix}$$

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

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

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



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

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

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

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

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

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

$$\frac{\partial L}{\partial w_{ij}} = x_j \frac{\partial L}{\partial z_i}$$

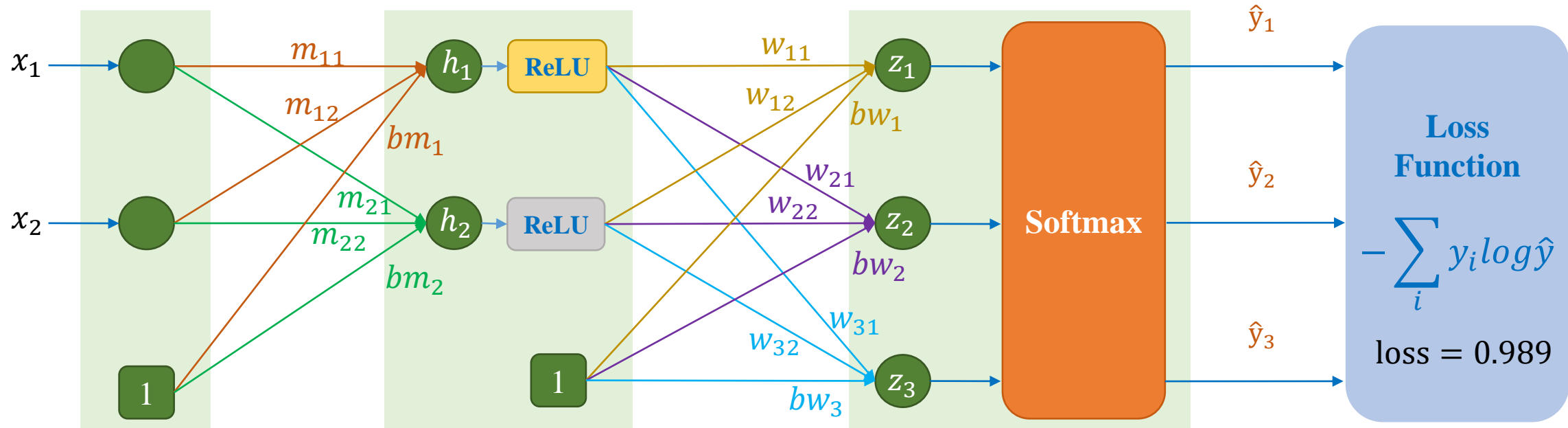
$$\frac{\partial L}{\partial bw_i} = \frac{\partial L}{\partial z_i}$$

Backward
pass

Backward
pass

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

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



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

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

$$\frac{\partial L}{\partial w_{ij}} = x_j \frac{\partial L}{\partial z_i}$$

$$\nabla_w L = \begin{bmatrix} -0.628 & 0.338 & 0.29 \\ 0.0 & 0.0 & 0.0 \end{bmatrix}$$

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

$$\nabla_{bw} L = \begin{bmatrix} -0.628 \\ 0.338 \\ 0.290 \end{bmatrix}$$

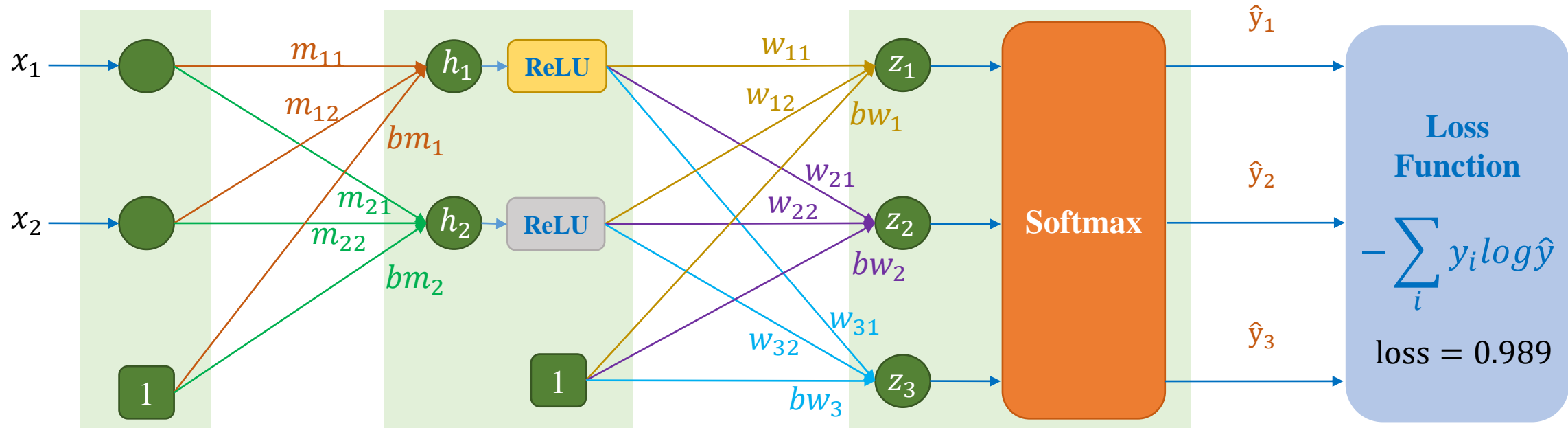
$$\frac{\partial L}{\partial z_i} = \hat{y}_i - y_i$$

$$\nabla_z L = \begin{bmatrix} -0.628 \\ 0.338 \\ 0.290 \end{bmatrix}$$

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

Backward
pass

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



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

$$\nabla_{\mathbf{m}} L = \begin{bmatrix} -0.114 & 0.0 \\ -0.015 & 0.0 \end{bmatrix}$$

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

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

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

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

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

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

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

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

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

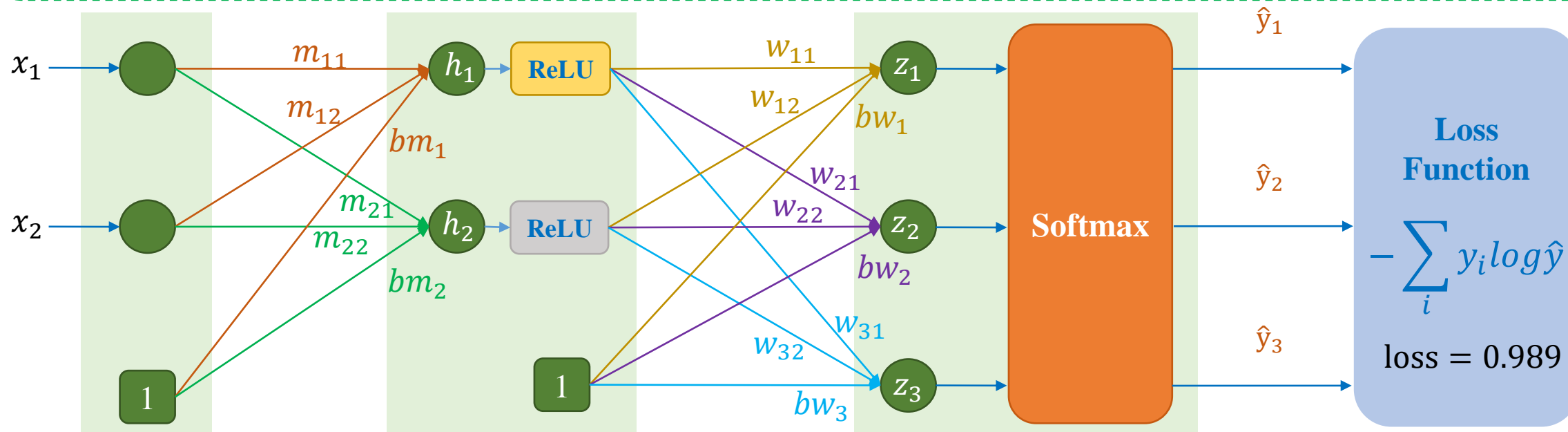
$$bw = \begin{bmatrix} 0.0 \\ 0.0 \\ 0.0 \end{bmatrix}$$

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

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

$$\nabla_w L = \begin{bmatrix} -0.628 & 0.338 & 0.29 \\ 0.0 & 0.0 & 0.0 \end{bmatrix}$$

$$\nabla_{bw} L = \begin{bmatrix} -0.628 \\ 0.338 \\ 0.290 \end{bmatrix}$$



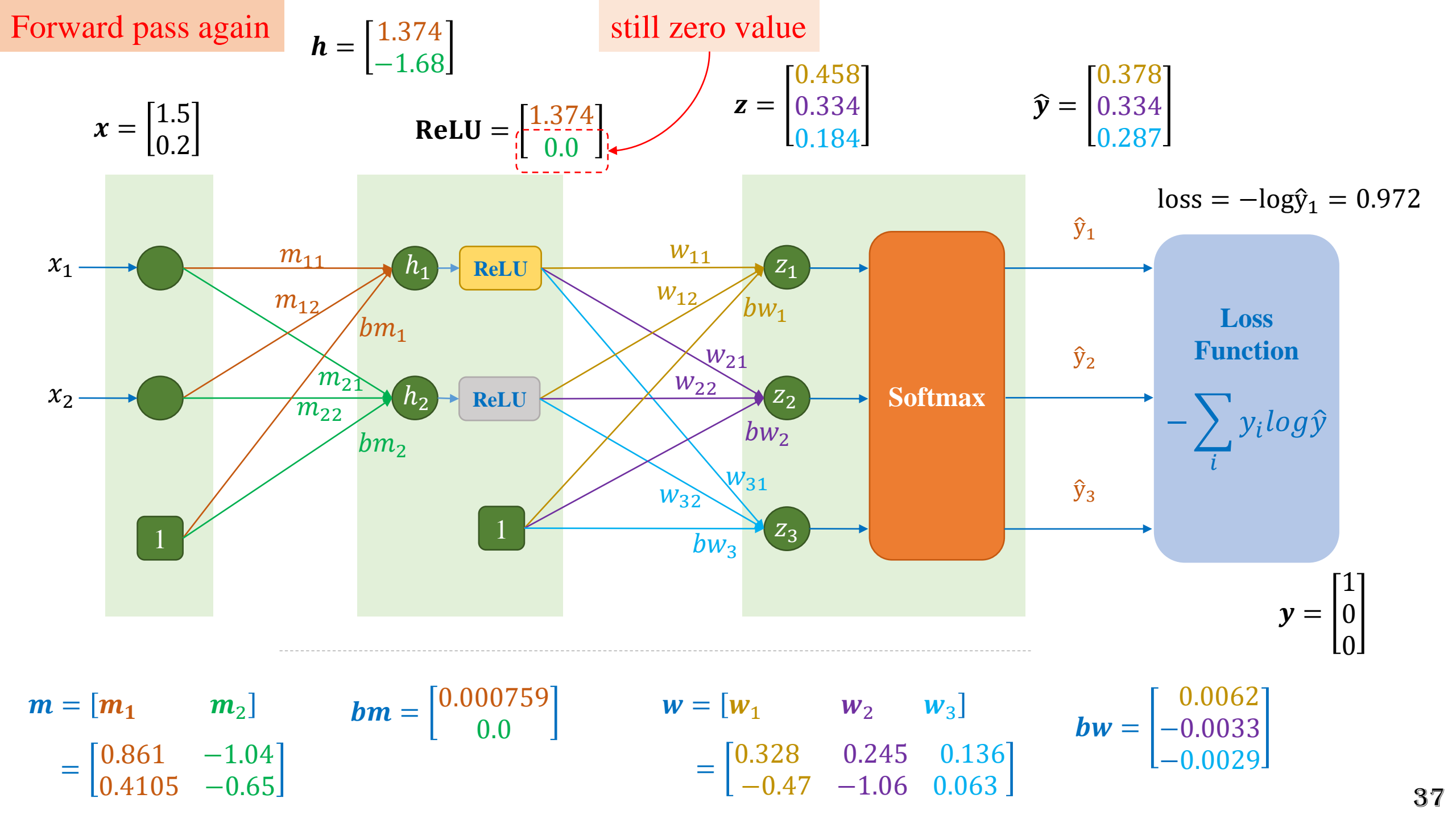
Update the parameters with $\eta = 0.01$

$$m = \begin{bmatrix} 0.861 & -1.04 \\ 0.4105 & -0.65 \end{bmatrix}$$

$$bm = \begin{bmatrix} 0.000759 \\ 0.0 \end{bmatrix}$$

$$w = \begin{bmatrix} 0.328 & 0.245 & 0.136 \\ -0.47 & -1.06 & 0.063 \end{bmatrix}$$

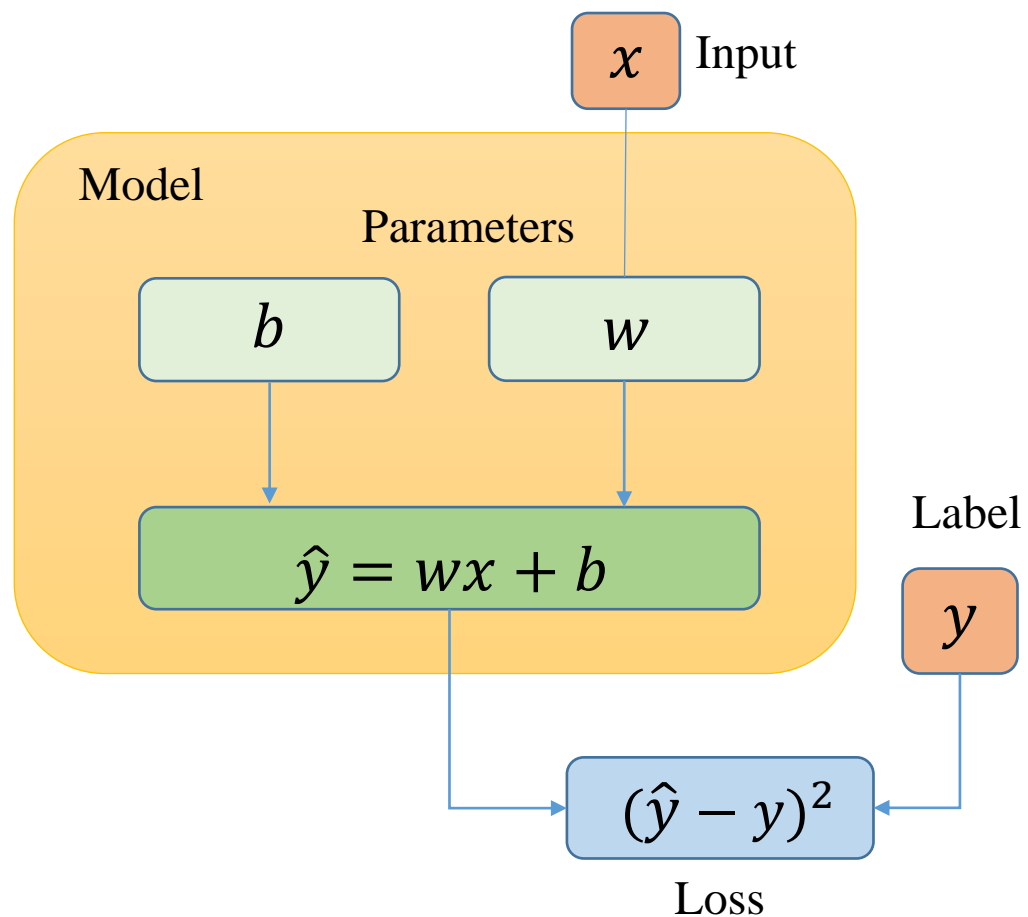
$$bw = \begin{bmatrix} 0.0062 \\ -0.0033 \\ -0.0029 \end{bmatrix}$$



Example 3 - Zero Initialization

❖ 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

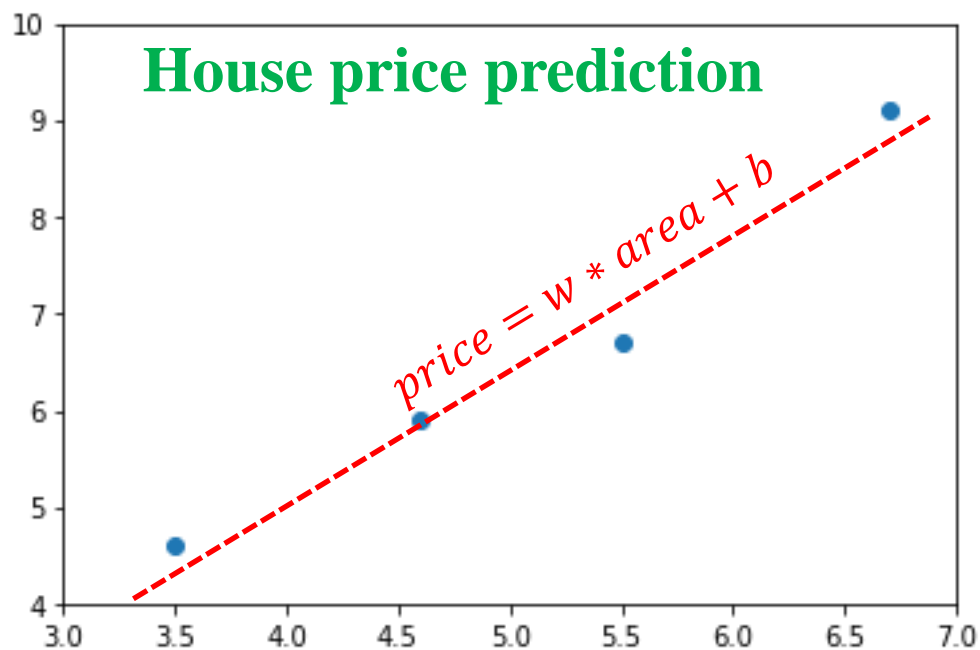
$$w = w - \eta L'_w$$

$$b = b - \eta L'_b$$

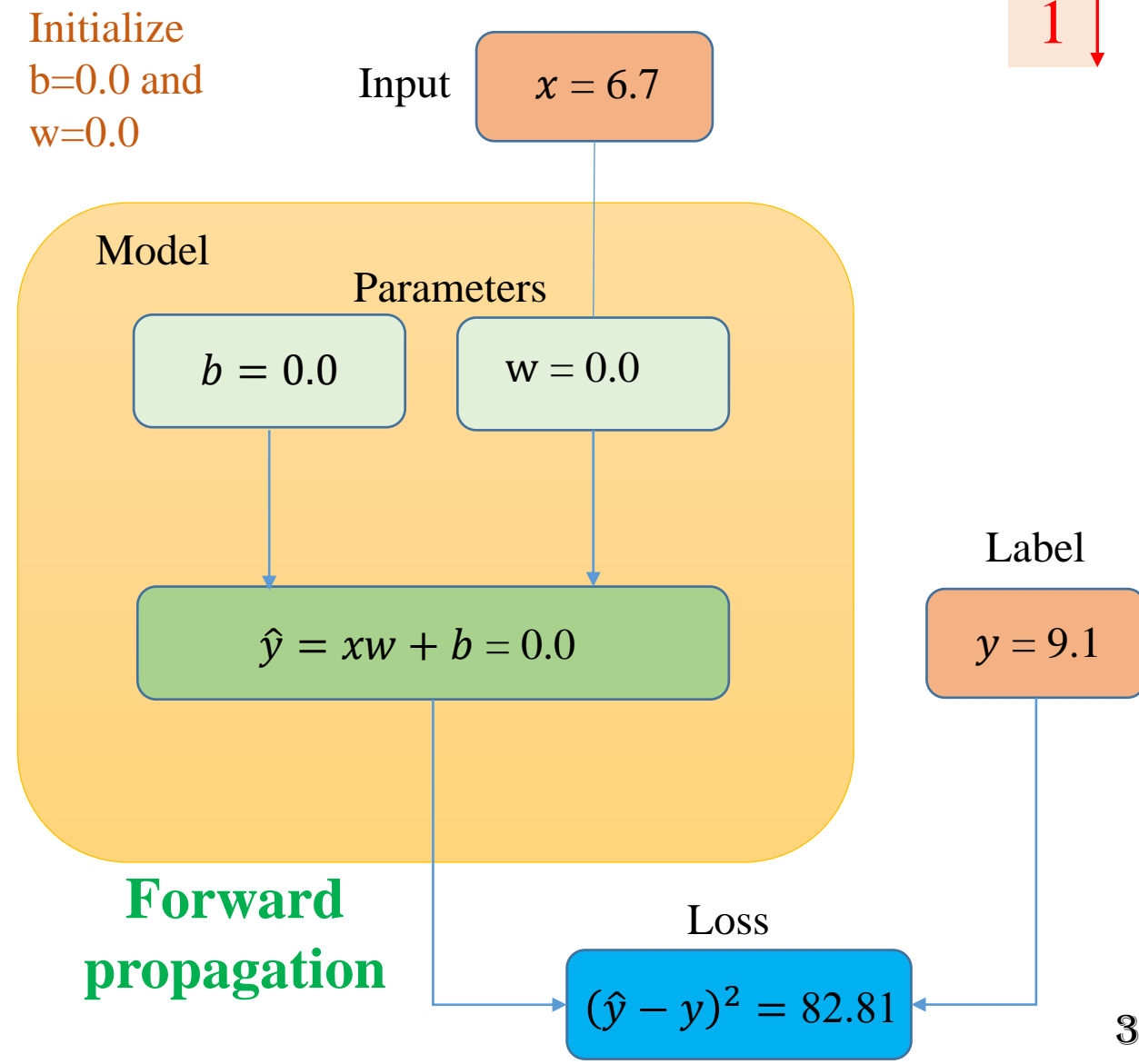
Example 3 - Zero Initialization

Given
sample
data

Feature	Label
area	price
6.7	9.1
4.6	5.9
3.5	4.6
5.5	6.7



Initialize
 $b=0.0$ and
 $w=0.0$



2

Input

 $x = 0.67$

Backpropagation

Model

Parameters

 $b = 0.0$ $w = 0.0$

$$b = b - \eta L'_b$$

$$w = w - \eta L'_w$$

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

Label

 $y = 9.1$

Loss

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

$$L'_w = 2x(\hat{y} - y) = -121.94$$

$$L'_b = 2(\hat{y} - y) = -18.2$$

$$b = b - \eta L'_b = 0.182$$

$$w = w - \eta L'_w = 1.2194$$

 $\eta = 0.01$

Input

 $x = 0.67$

Forward propagation

Model

Parameters

 $b = 0.182$ $w = 1.2194$

$$b = b - \eta L'_b$$

$$w = w - \eta L'_w$$

$$\hat{y} = xw + b = 8.351$$

Label

 $y = 9.1$

Loss

$$(\hat{y} - y)^2 = 0.559$$

New w and b help
the loss reduce

Example 4 - Zero Initialization

❖ Logistic regression

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

2) Compute output \hat{y}

$$z = \theta^T x$$

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

3) Compute loss

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

4) Compute derivative

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

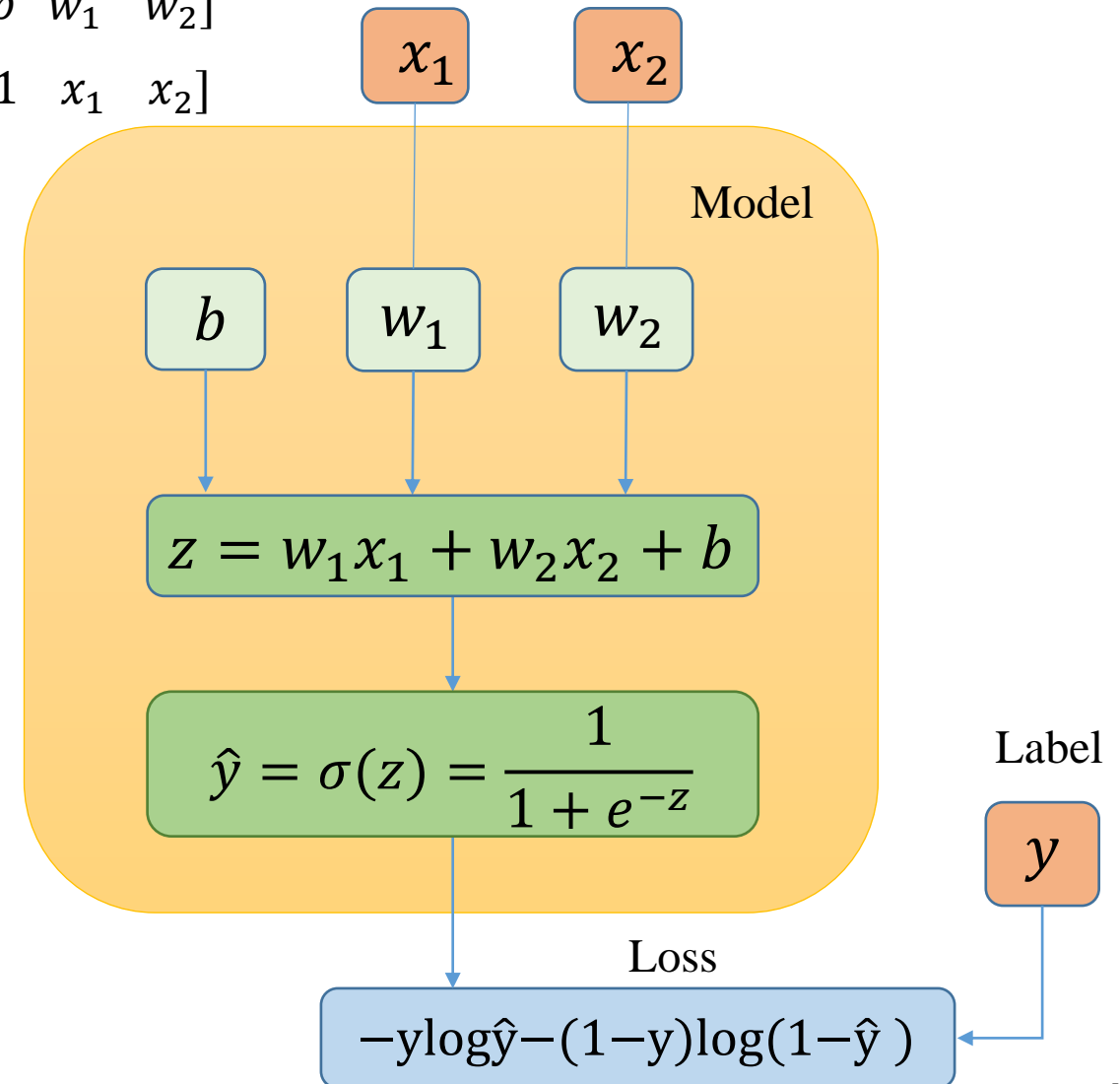
5) Update parameters

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

η is learning rate

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

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



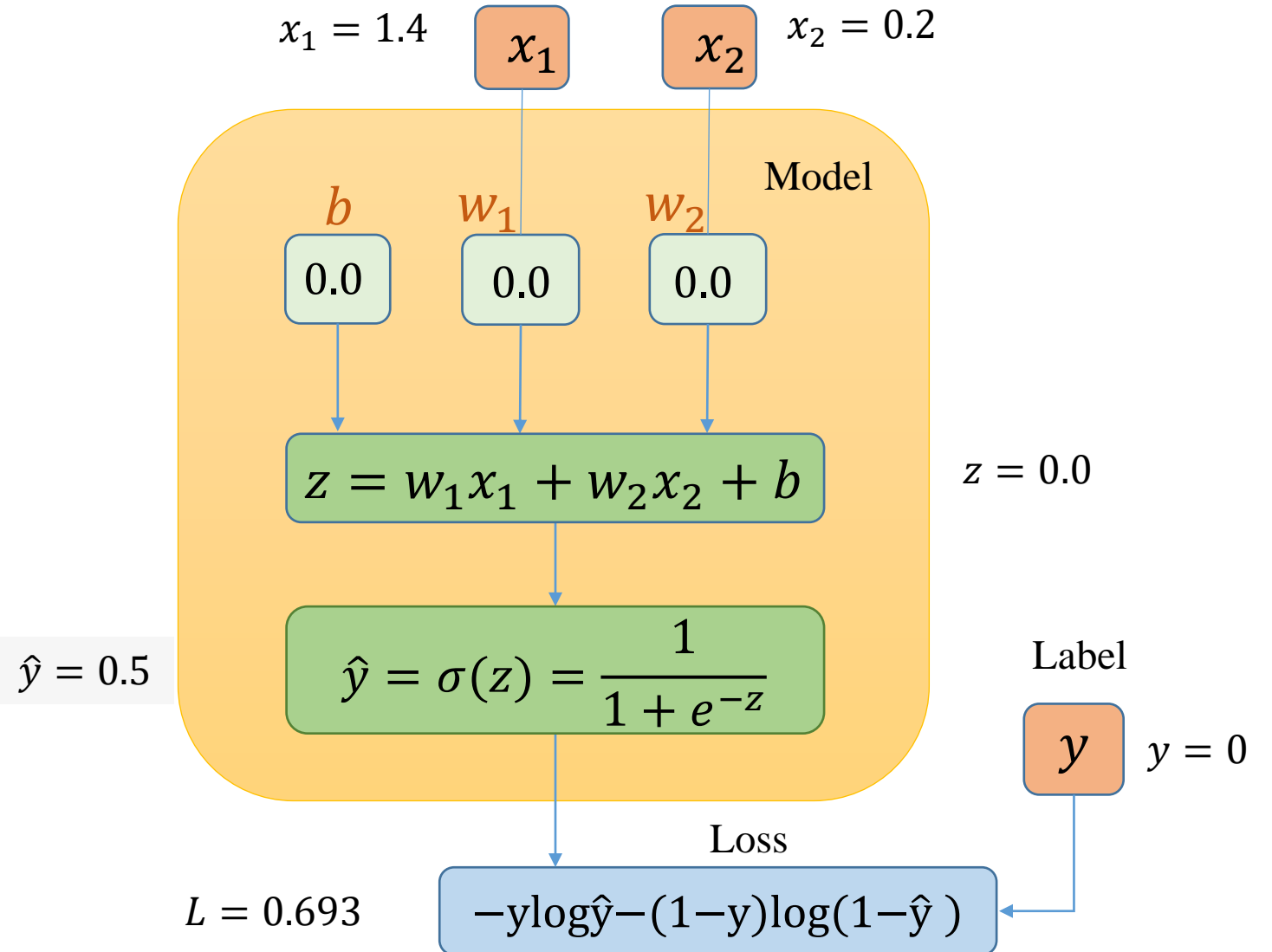
Example 4 - Zero Initialization

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

$$\mathbf{x} = \begin{bmatrix} 1 \\ 1.4 \\ 0.2 \end{bmatrix}$$

$$\mathbf{y} = [0]$$



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

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

$$\eta = 0.01$$

$$b = 0.005$$

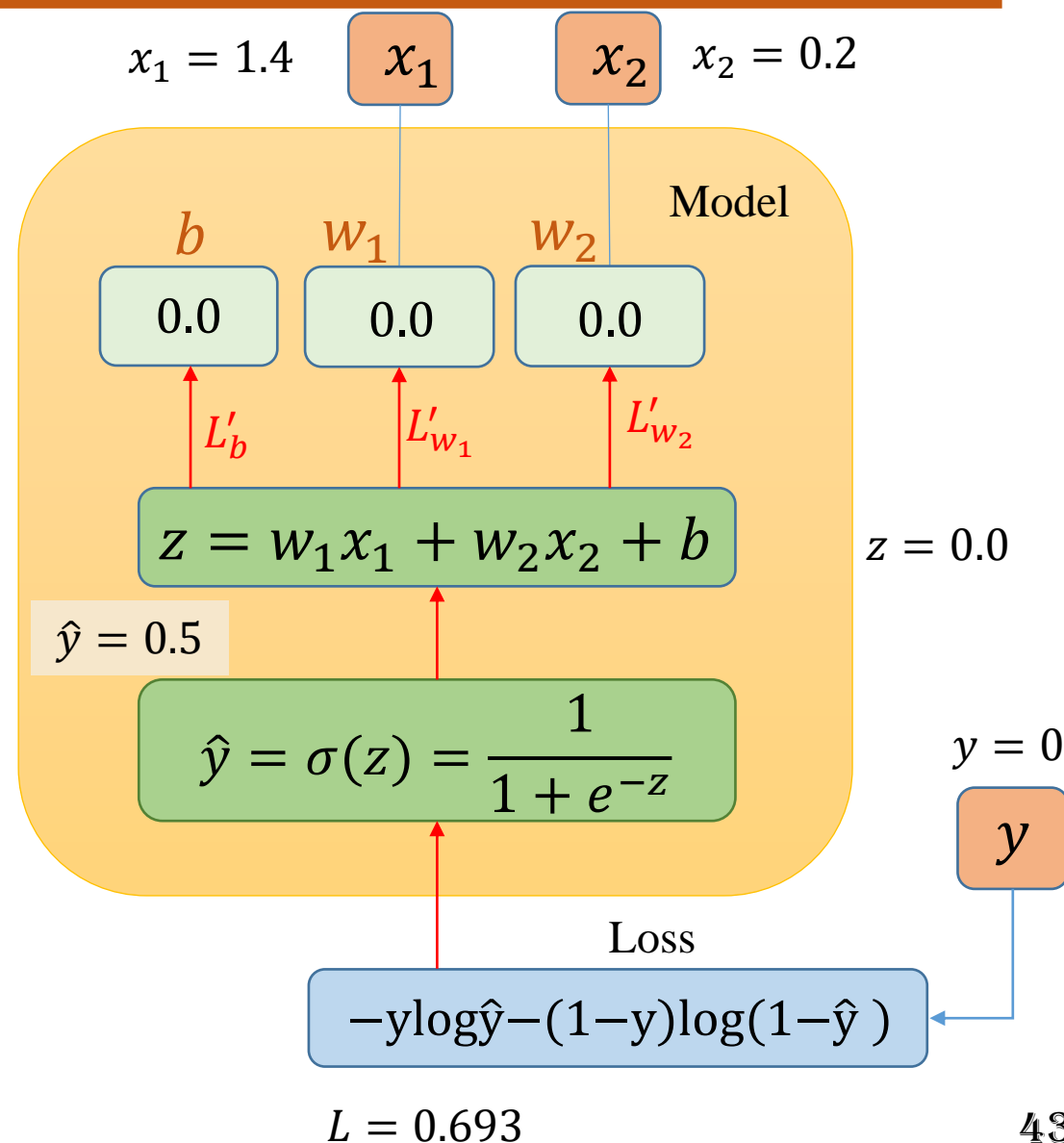
$$w_1 = 0.007$$

$$w_2 = 0.001$$

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

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

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



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

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

$$\eta = 0.01$$

$$b = -0.005$$

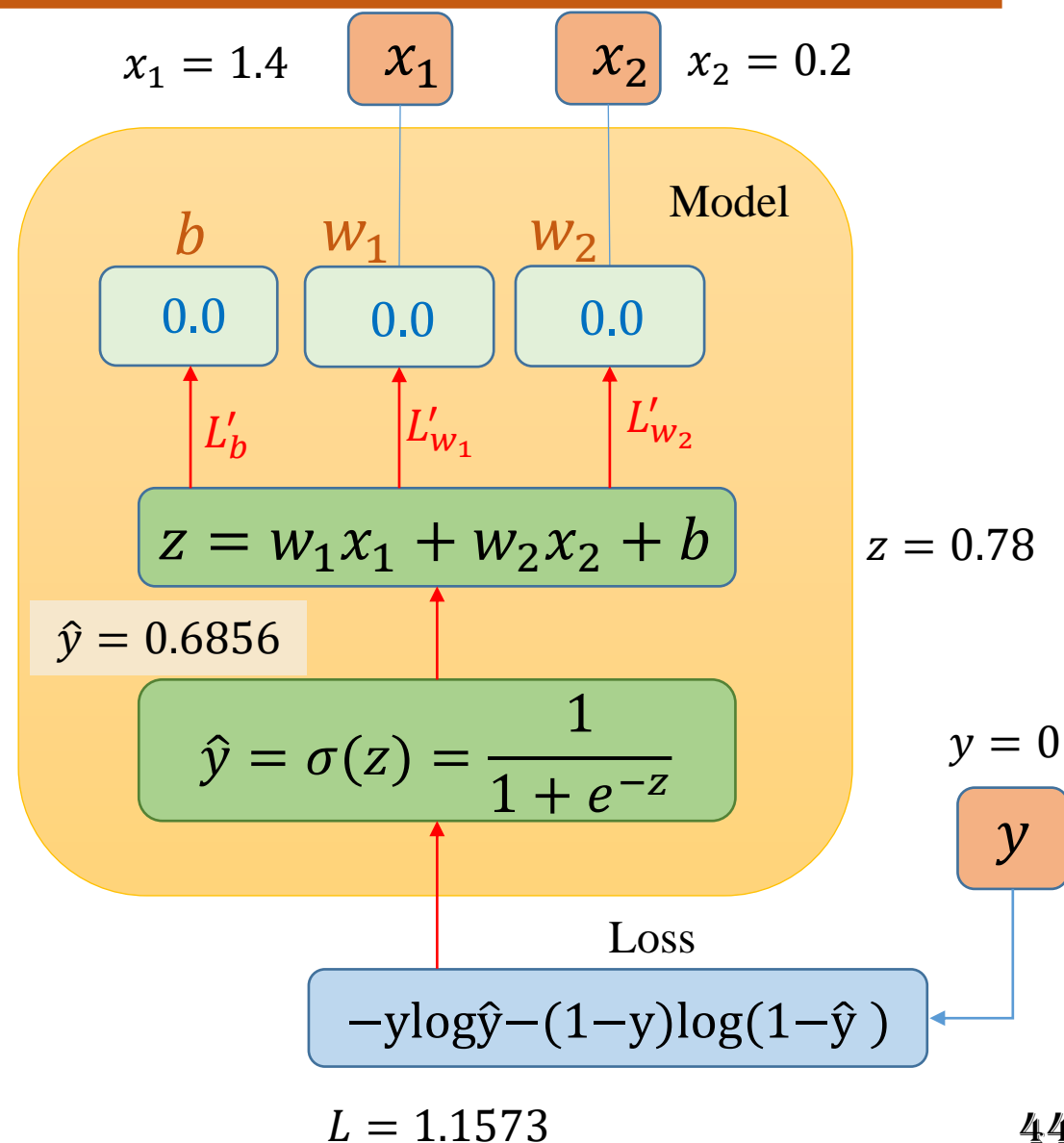
$$w_1 = -0.007$$

$$w_2 = -0.001$$

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

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

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



Example 4 - Zero Initialization

Dataset

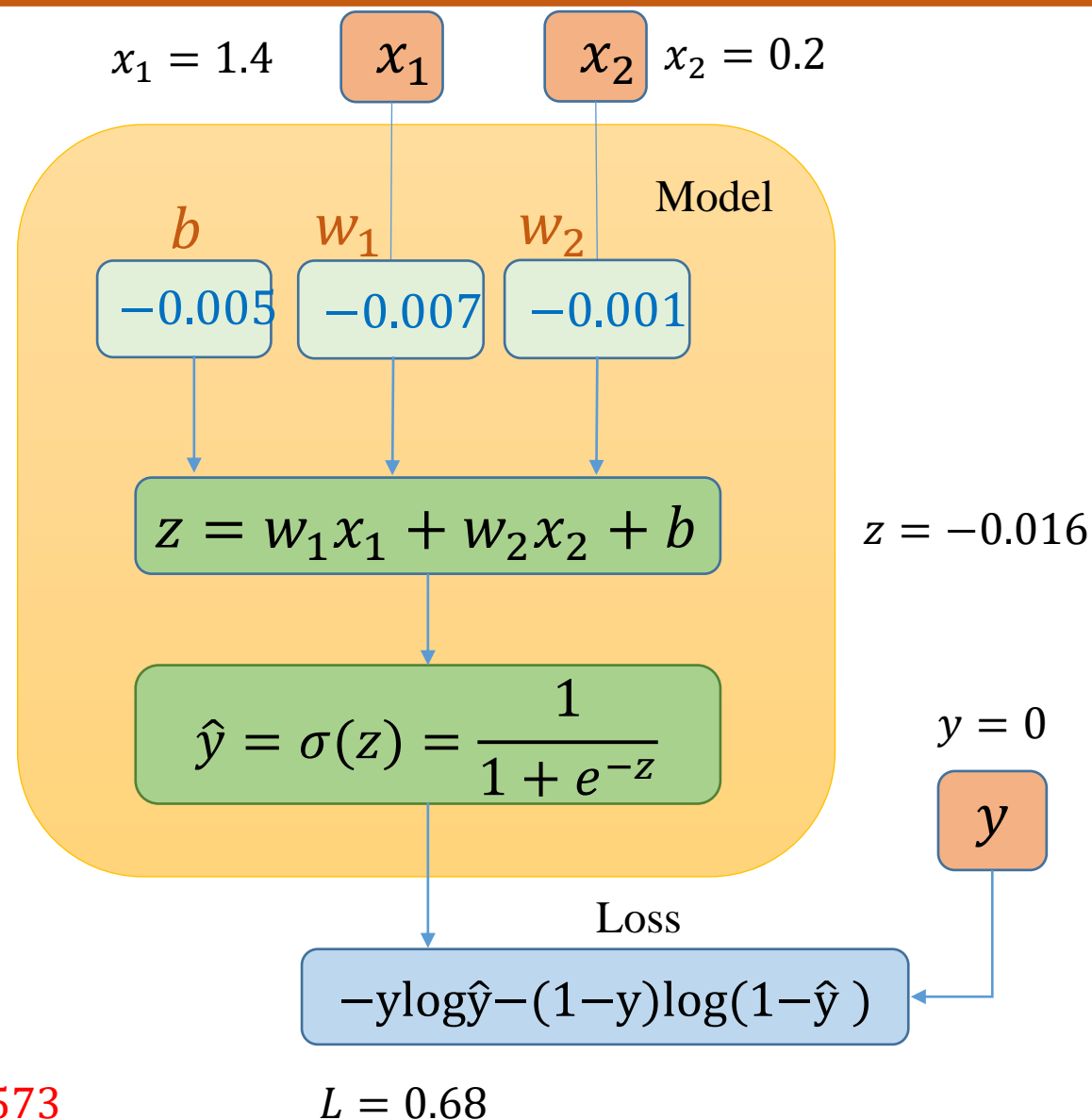
Petal_Length	Petal_Width	Label
1.4	0.2	0
1.5	0.2	0
3	1.1	1
4.1	1.3	1

$$\mathbf{x} = \begin{bmatrix} 1 \\ 1.4 \\ 0.2 \end{bmatrix}$$

$$\mathbf{y} = [0]$$

$$\hat{y} = 0.49$$

previous $L = 1.1573$



Example 5 - Zero Initialization

❖ Softmax regression

Training data

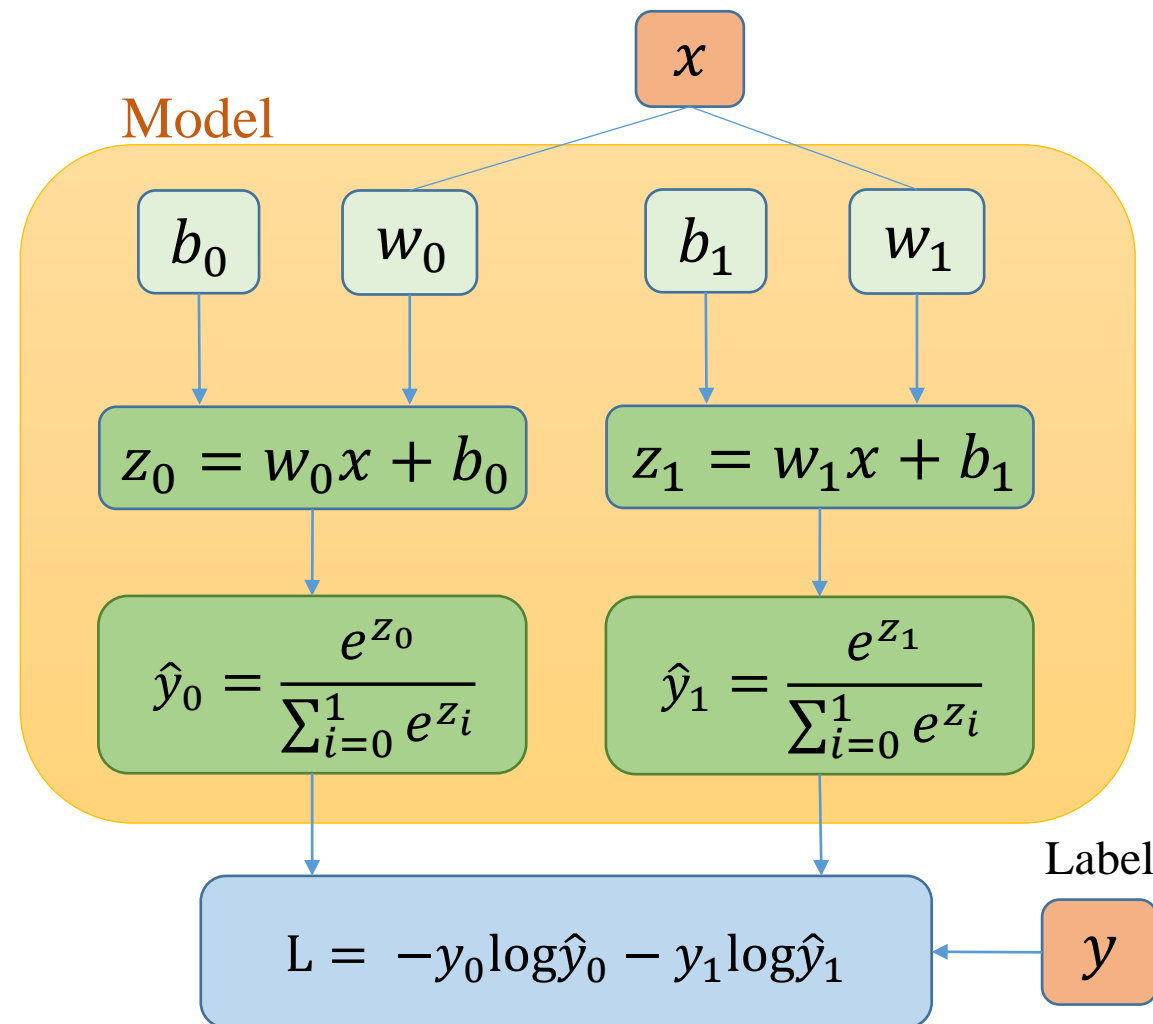
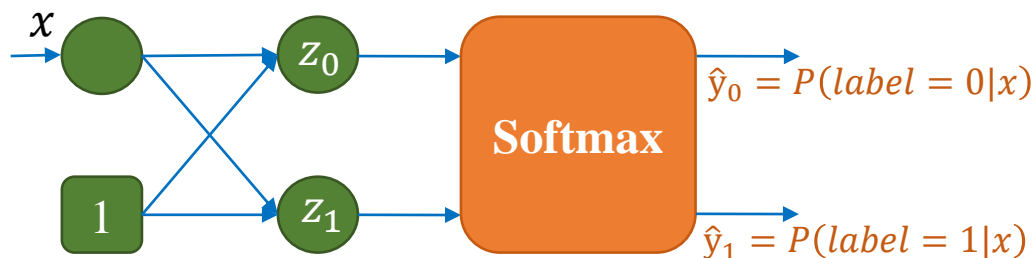
Feature	Label	
Petal_Length	Label	
1.4	0	Category A
1.3	0	
1.5	0	
4.5	1	Category B
4.1	1	
4.6	1	

One-hot encoding for labels

index 0 1

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

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



Example 5 - Zero Initialization

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

#class=2

#feature=1

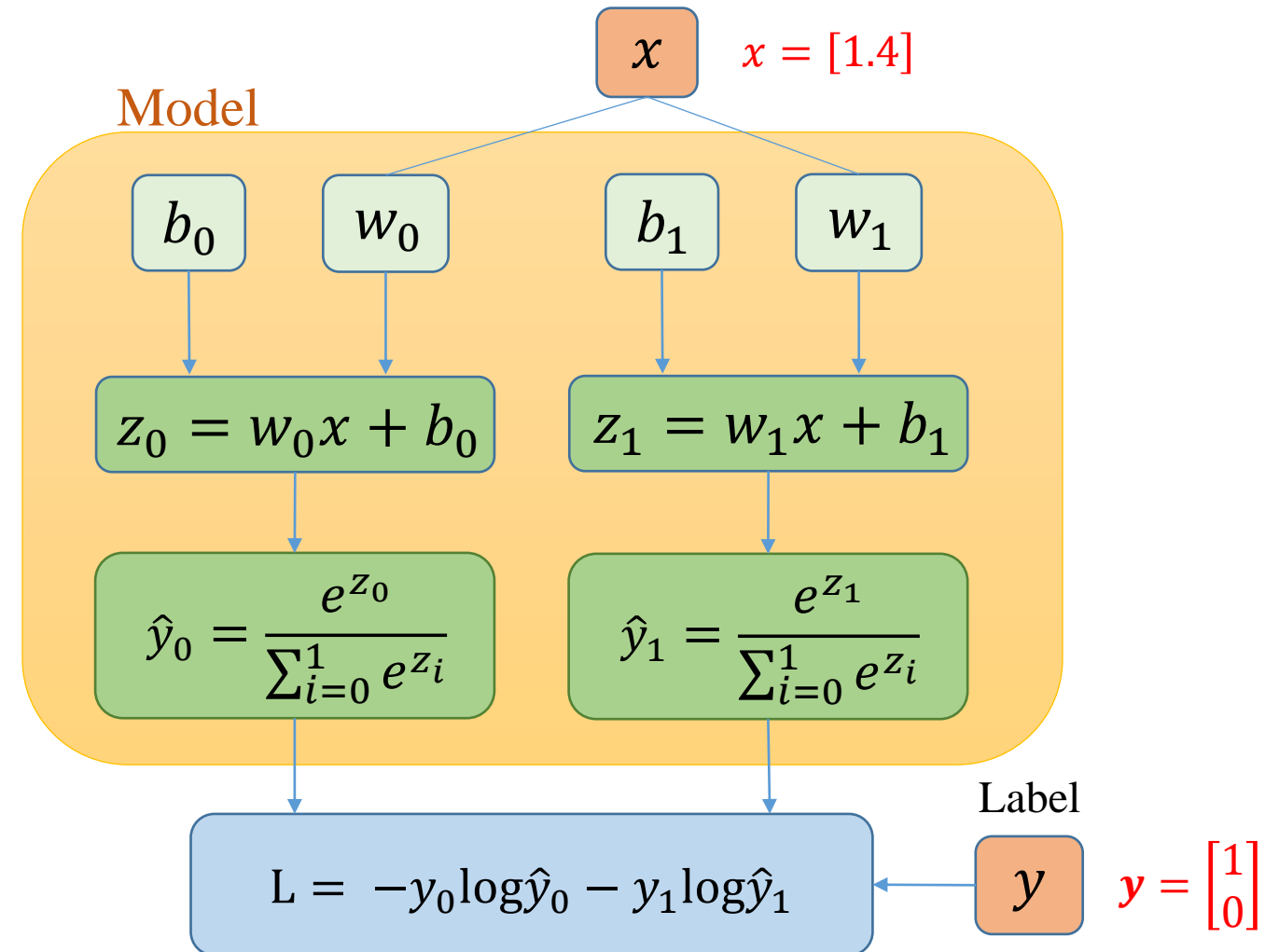
One-hot encoding for label

$$y = 0 \rightarrow \mathbf{y}^T = \begin{bmatrix} 1 & 0 \end{bmatrix}$$

$$y = 1 \rightarrow \mathbf{y}^T = \begin{bmatrix} 0 & 1 \end{bmatrix}$$

Training example

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



Example 5 - Zero Initialization

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

#class=2

#feature=1

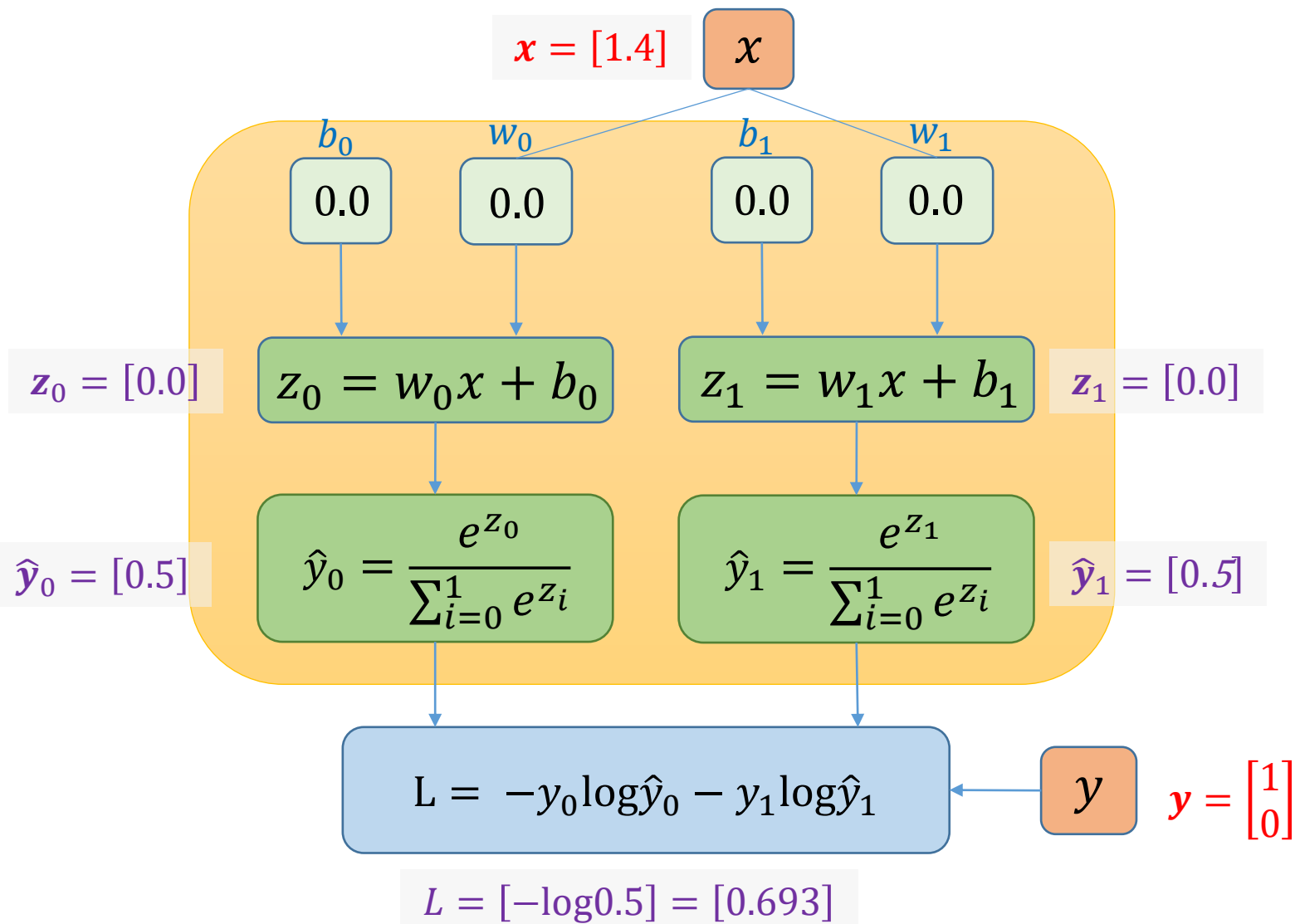
One-hot encoding for label

$$y = 0 \rightarrow \mathbf{y}^T = \begin{bmatrix} 1 & 0 \end{bmatrix}$$

$$y = 1 \rightarrow \mathbf{y}^T = \begin{bmatrix} 0 & 1 \end{bmatrix}$$

Training example

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



Example 5 - Zero Initialization

Derivative

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

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

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

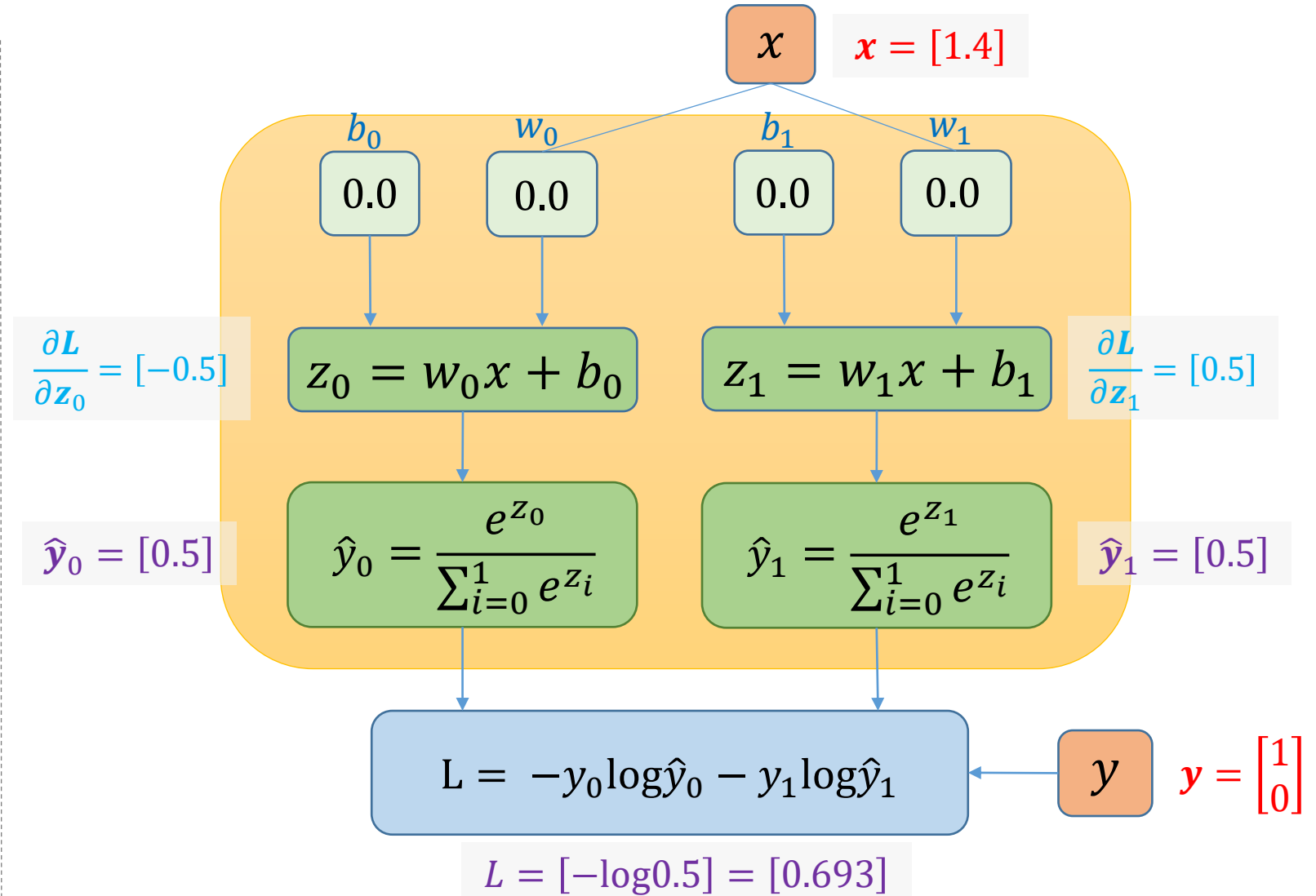
$$y = 0 \rightarrow \mathbf{y}^T = \begin{bmatrix} y_0 & y_1 \\ 1 & 0 \end{bmatrix}$$

$$y = 1 \rightarrow \mathbf{y}^T = \begin{bmatrix} 0 & 1 \end{bmatrix}$$

$$\frac{\partial L}{\partial z_0} = \hat{y}_0 - 1$$

$$= 0.5 - 1 = -0.5$$

$$\frac{\partial L}{\partial z_1} = \hat{y}_1 - 0 = 0.5$$



Example 5 - Zero Initialization

Derivative

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

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

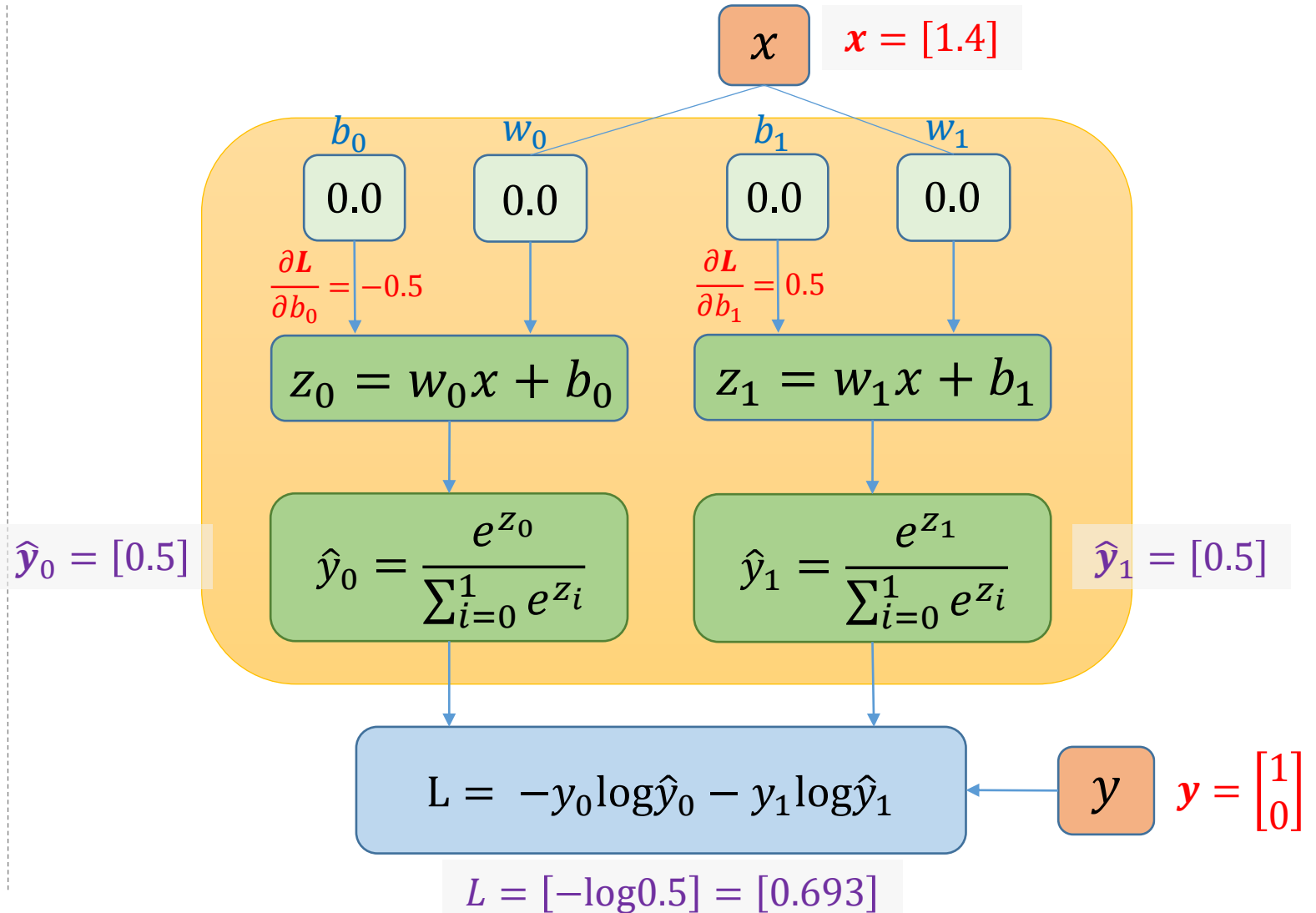
$$\frac{\partial L}{\partial b_i} = \hat{y}_i - y_i$$

$$y = 0 \rightarrow \mathbf{y}^T = \begin{bmatrix} y_0 & y_1 \\ 1 & 0 \end{bmatrix}$$

$$y = 1 \rightarrow \mathbf{y}^T = \begin{bmatrix} 0 & 1 \end{bmatrix}$$

$$\frac{\partial L}{\partial b_0} = (\hat{y}_0 - 1) = -0.5$$

$$\frac{\partial L}{\partial b_1} = (\hat{y}_1 - 0) = 0.5$$



Example 5 - Zero Initialization

Derivative

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

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

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

$$y = 0 \rightarrow \mathbf{y}^T = \begin{bmatrix} y_0 & y_1 \\ 1 & 0 \end{bmatrix}$$

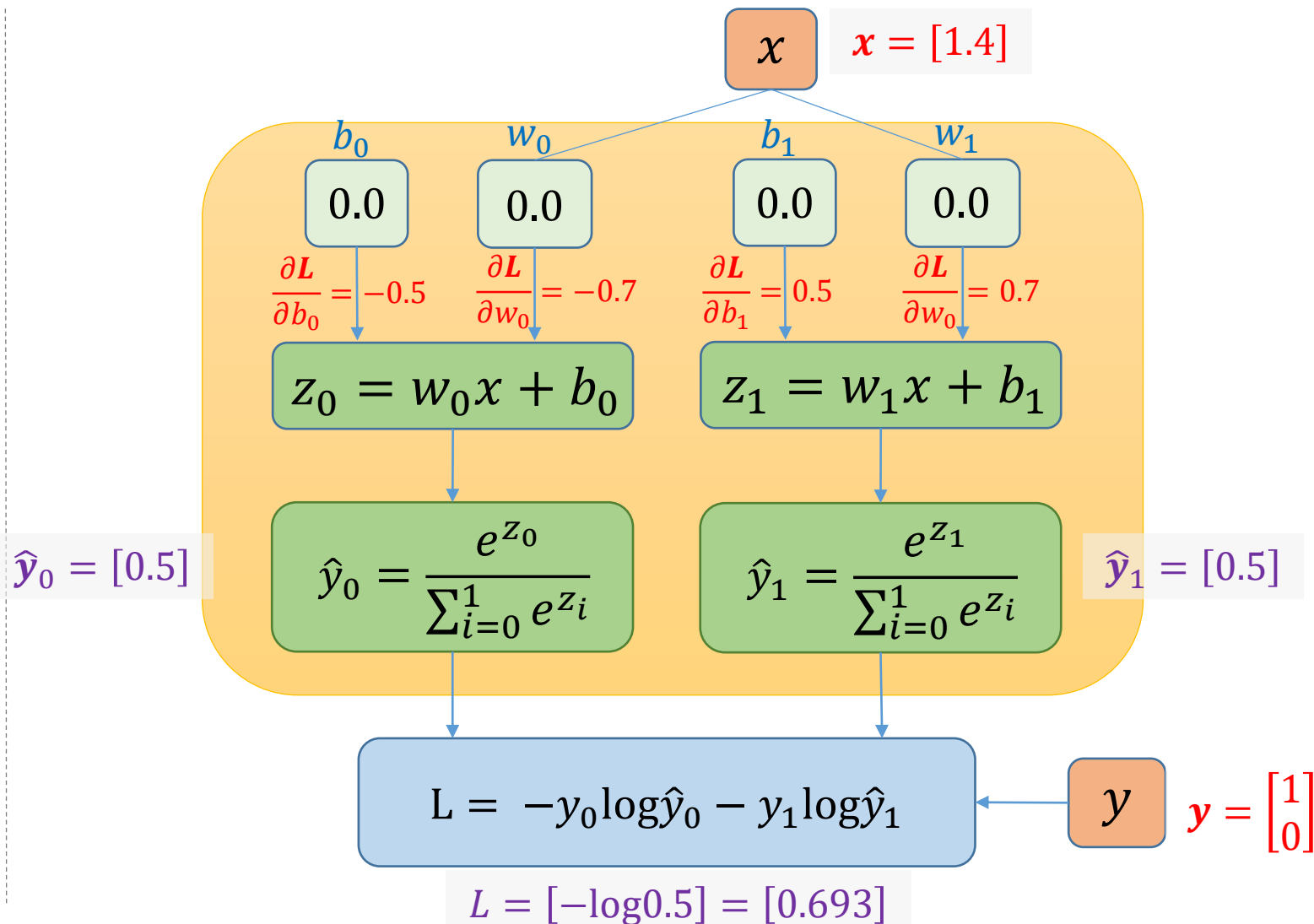
$$y = 1 \rightarrow \mathbf{y}^T = \begin{bmatrix} 0 & 1 \end{bmatrix}$$

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

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

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

$$= 0.5 * 1.4 = 0.7$$



Example 5 - Zero Initialization

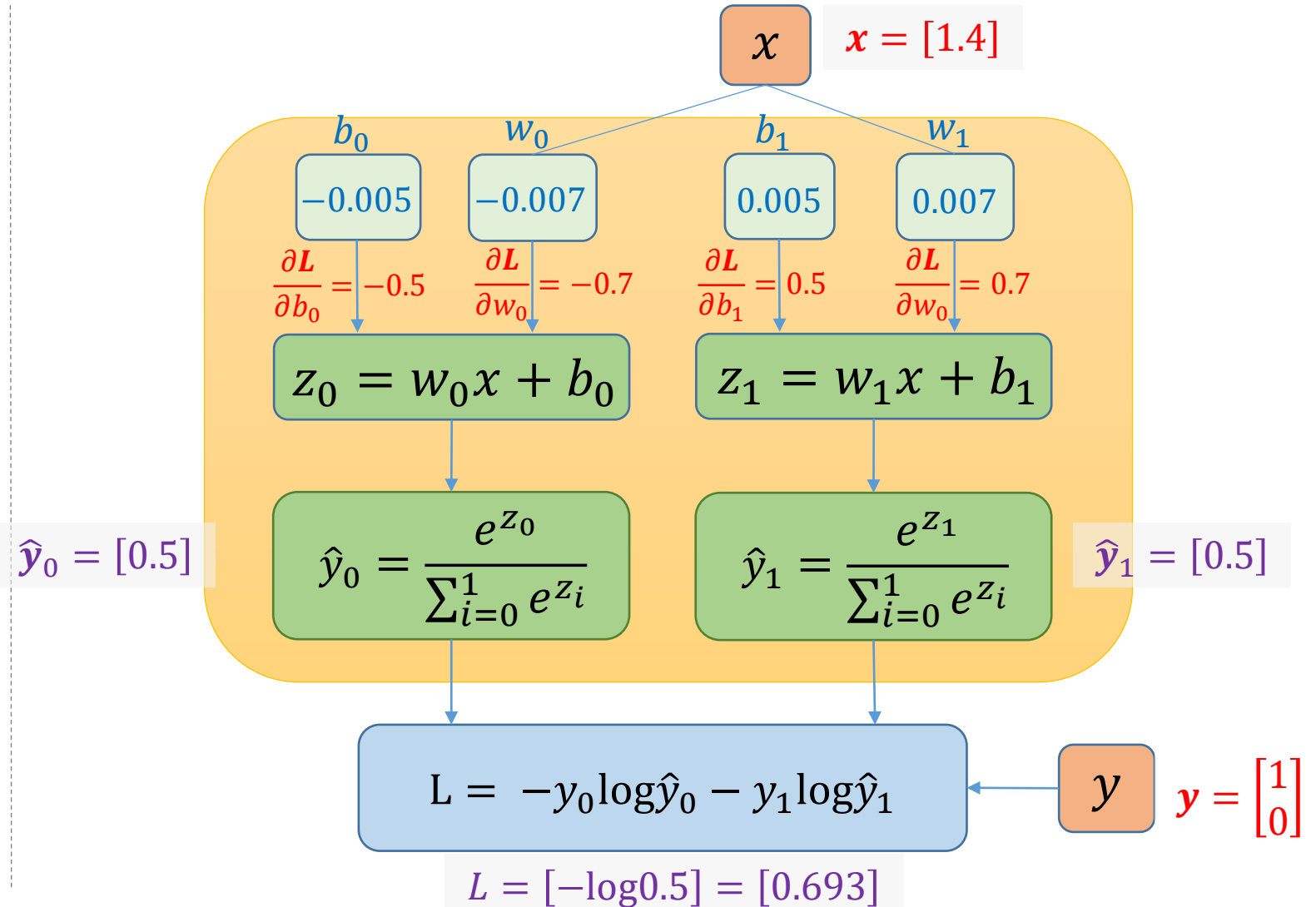
Update parameters

$$\theta = \theta - \eta L'_\theta$$

η is learning rate

$$\theta = \begin{bmatrix} b_0 & b_1 \\ w_0 & w_1 \end{bmatrix} \quad \eta = 0.1 \quad L'_\theta = \begin{bmatrix} \frac{\partial L}{\partial b_0} & \frac{\partial L}{\partial b_1} \\ \frac{\partial L}{\partial w_0} & \frac{\partial L}{\partial w_1} \end{bmatrix}$$

$$\theta = \begin{bmatrix} 0.0 & 0.0 \\ 0.0 & 0.0 \end{bmatrix} - 0.01 \begin{bmatrix} -0.5 & 0.5 \\ -0.7 & 0.7 \end{bmatrix} = \begin{bmatrix} -0.005 & 0.005 \\ -0.007 & 0.007 \end{bmatrix}$$



Example 5 - Zero Initialization

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

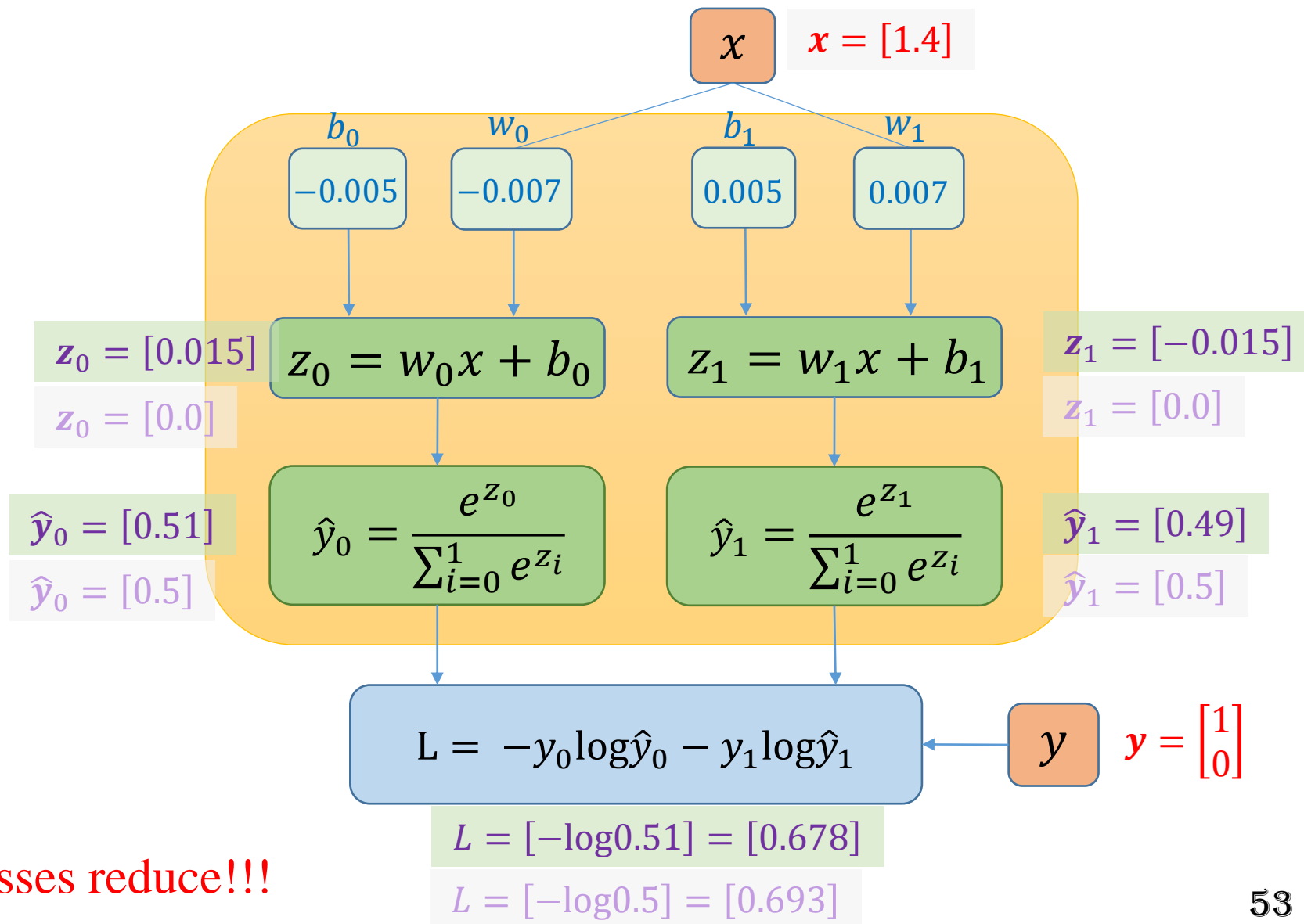
One-hot encoding for label

$$y = 0 \rightarrow \mathbf{y}^T = \begin{bmatrix} 1 & 0 \end{bmatrix}$$

$$y = 1 \rightarrow \mathbf{y}^T = \begin{bmatrix} 0 & 1 \end{bmatrix}$$

Training example

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

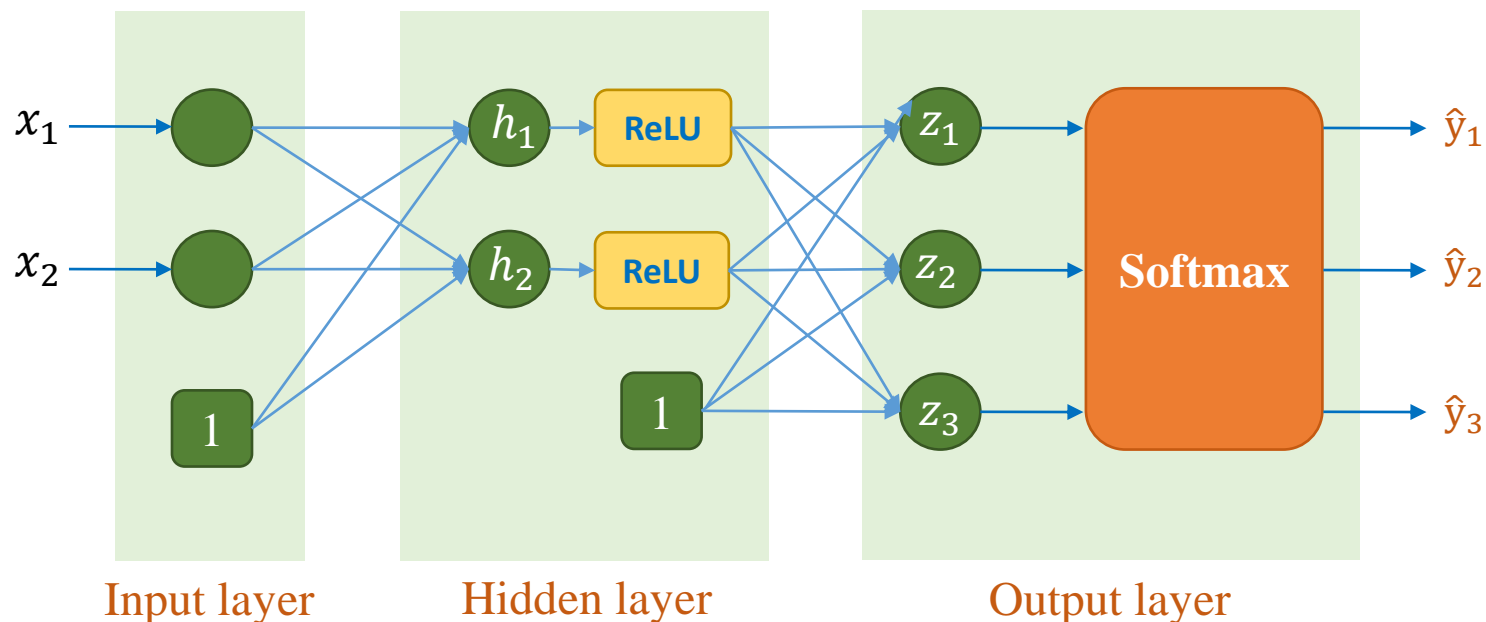


Example 6 - Zero Initialization

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

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

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



$$\mathbf{h} = [\mathbf{h}_1 \quad \mathbf{h}_2]$$

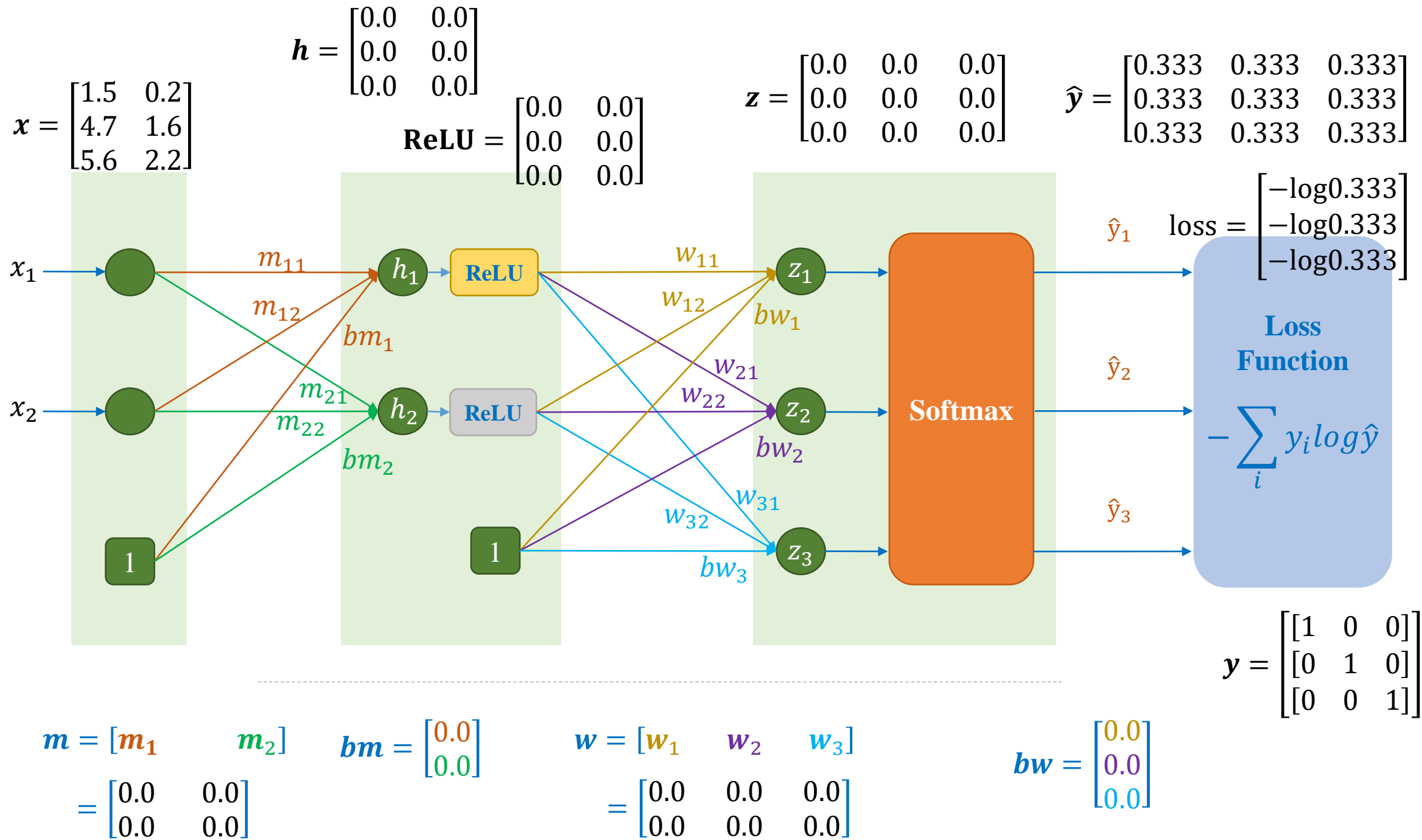
$$= \begin{bmatrix} 0.0 & 0.0 \\ 0.0 & 0.0 \end{bmatrix}$$

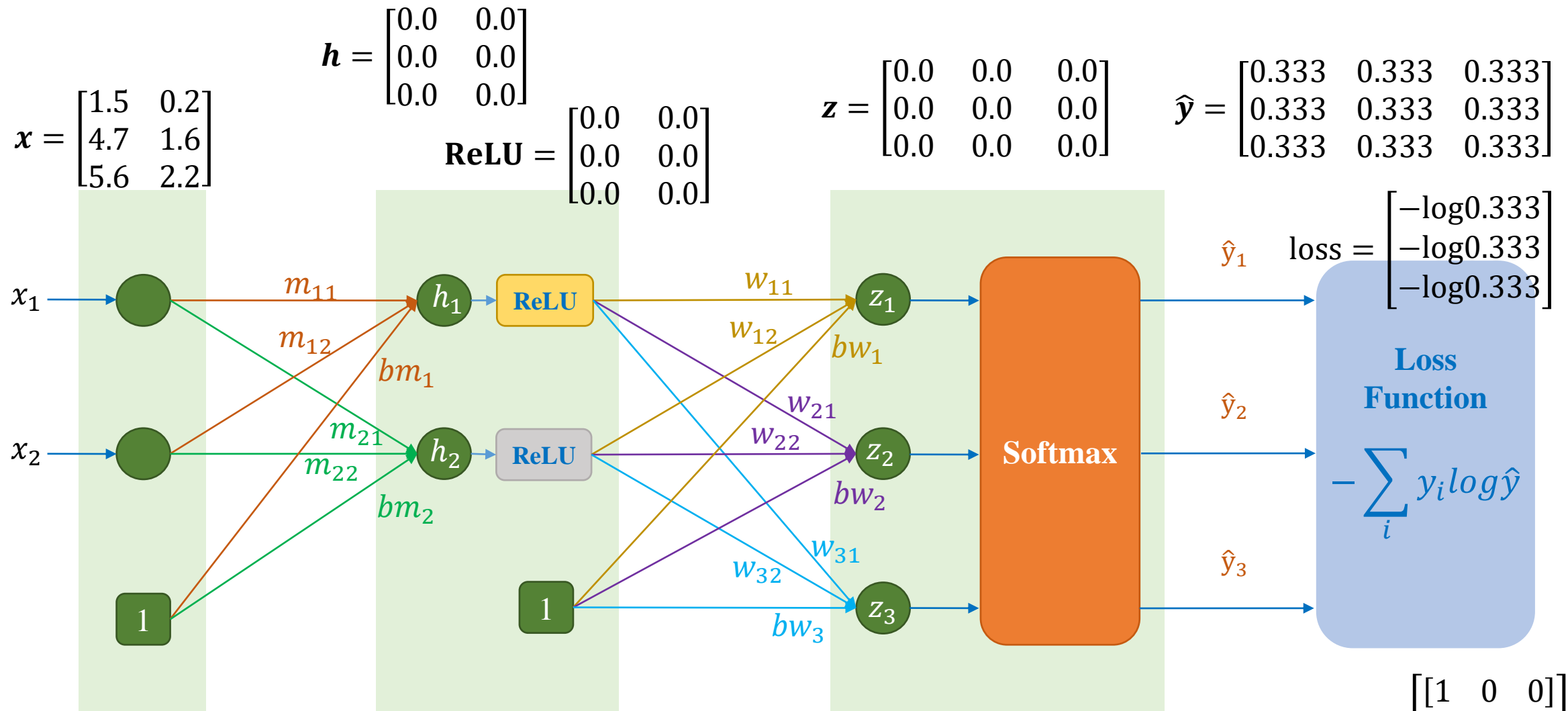
$$\mathbf{b}_h = \begin{bmatrix} 0.0 \\ 0.0 \end{bmatrix}$$

$$\mathbf{w} = [\mathbf{w}_1 \quad \mathbf{w}_2 \quad \mathbf{w}_3]$$

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

$$\mathbf{b}_w = \begin{bmatrix} 0.0 \\ 0.0 \\ 0.0 \end{bmatrix}$$





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

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

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

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

Optimizers

Optimizer Selection

Data Preparation



Data
Normalization



Model (Network)
Construction



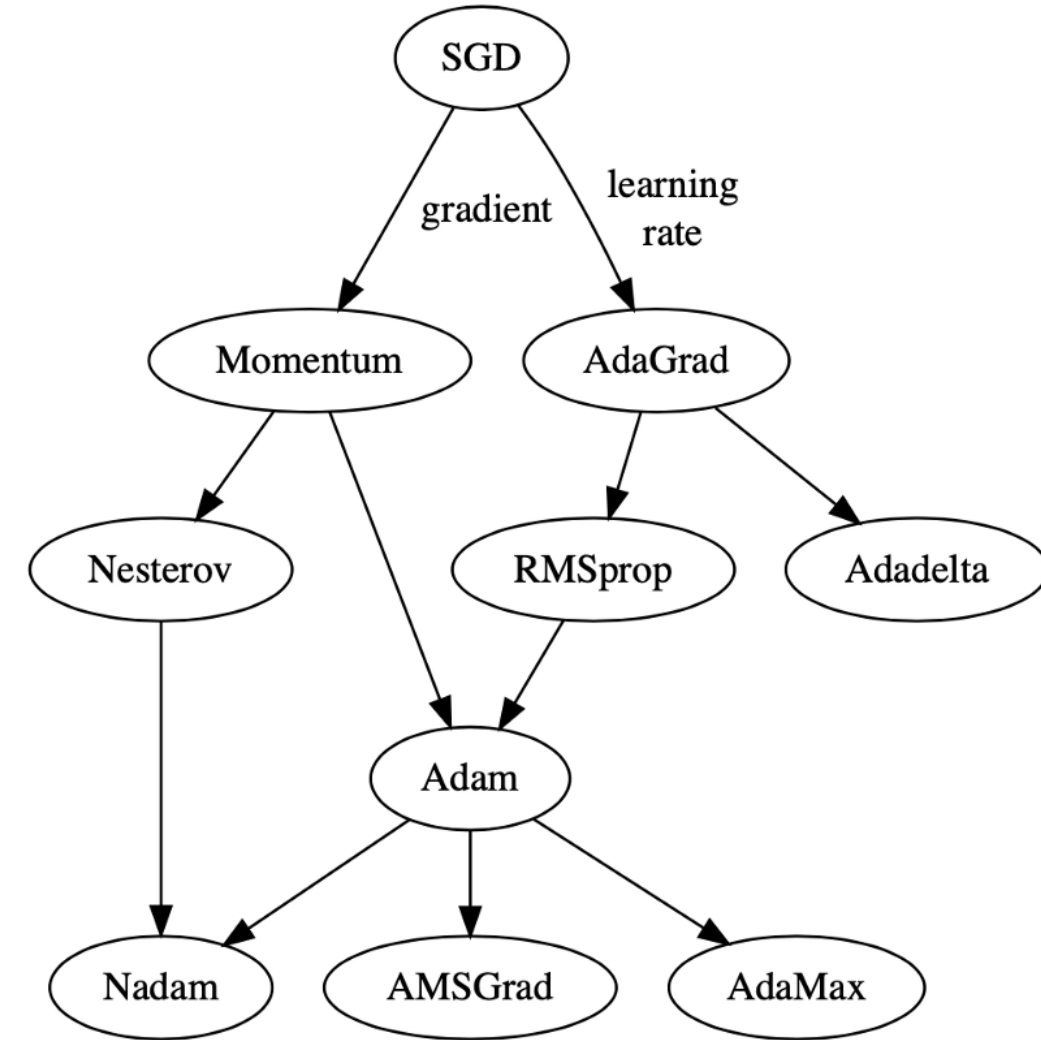
Parameter
Initialization

Define a way to update parameters

Optimizer
Selection

Loss function
Selection

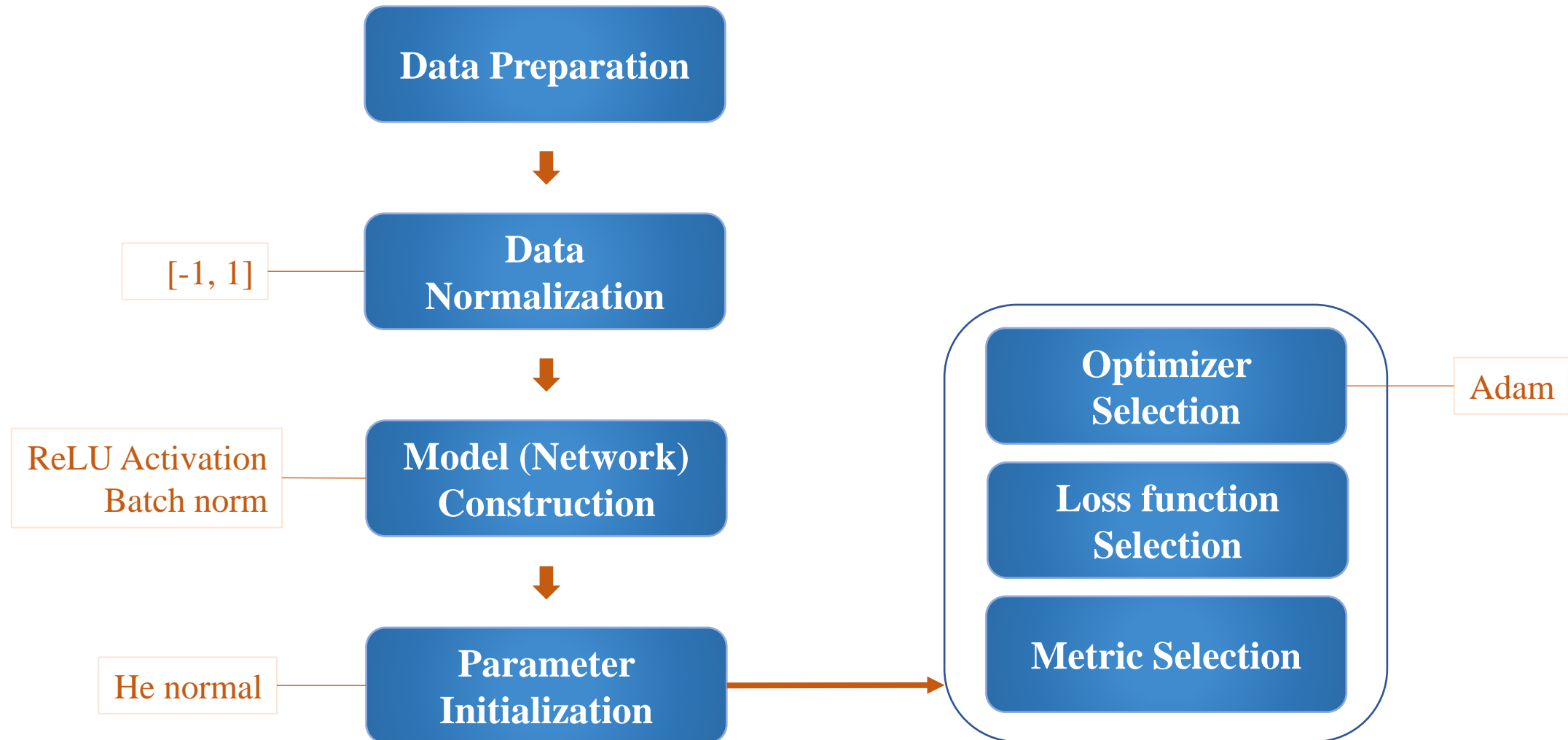
Metric Selection



<https://www.kdnuggets.com/2019/06/gradient-descent-algorithms-cheat-sheet.html>

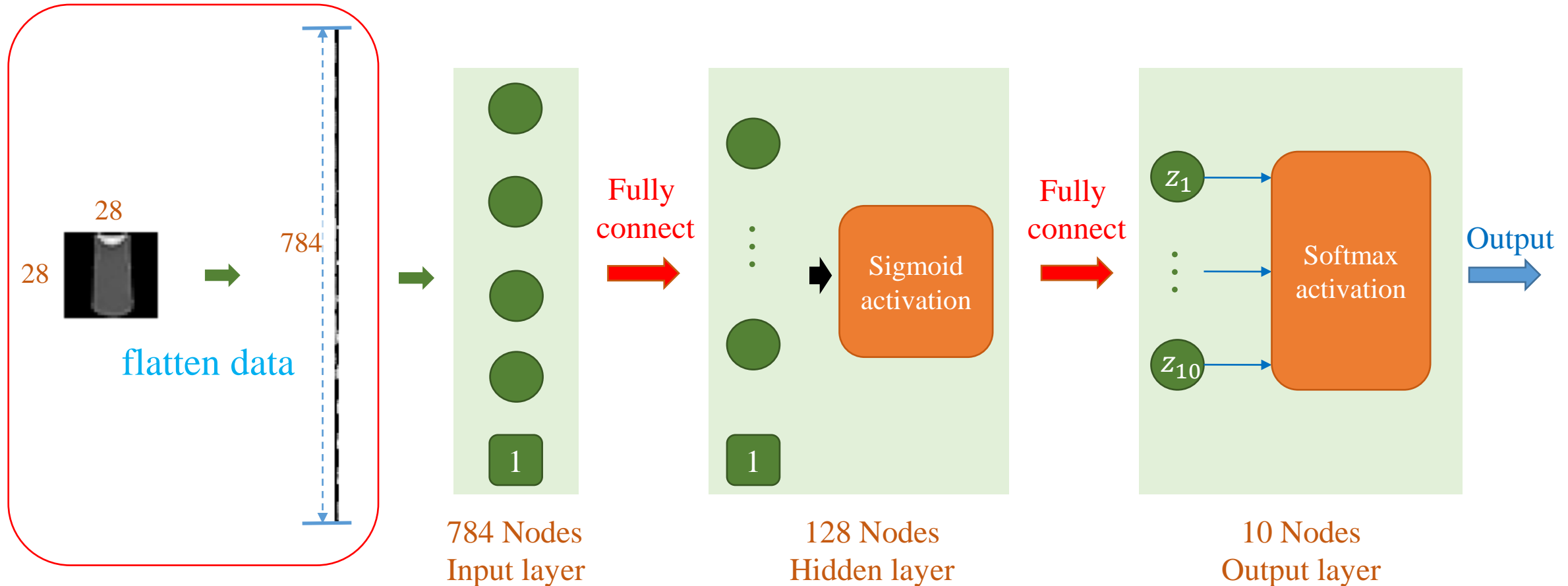
Summary

Recommendation



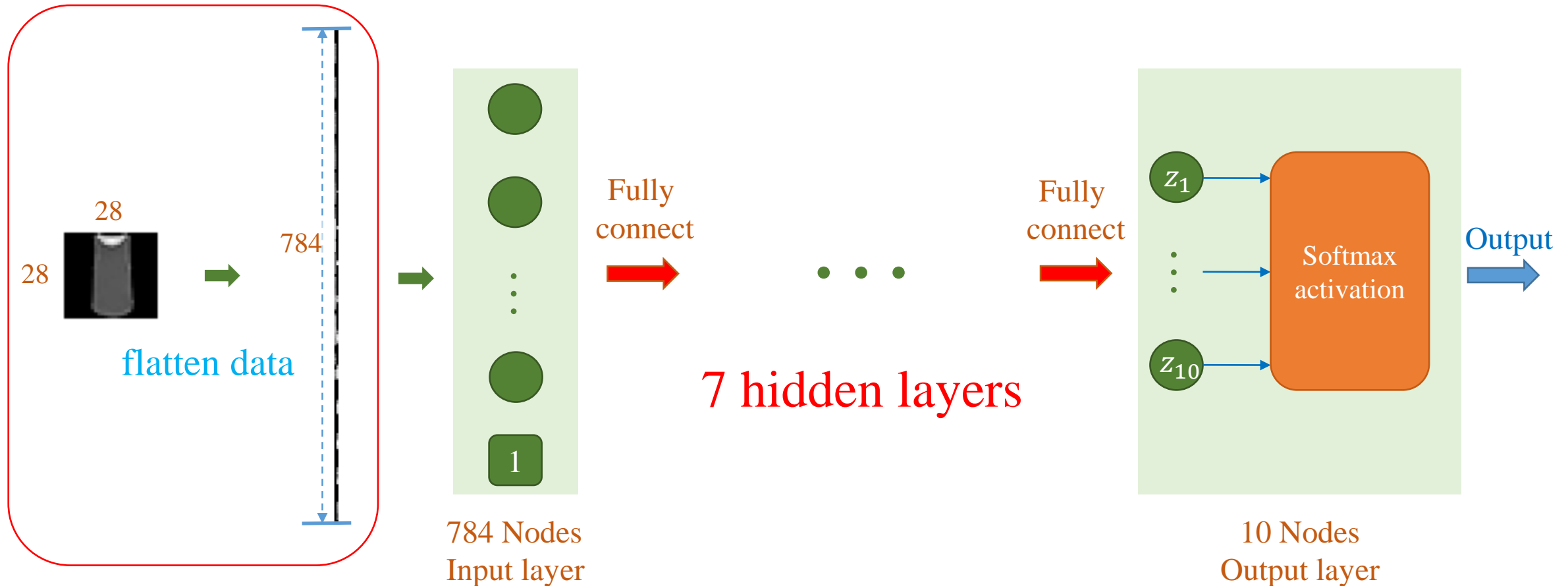
Discussion

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



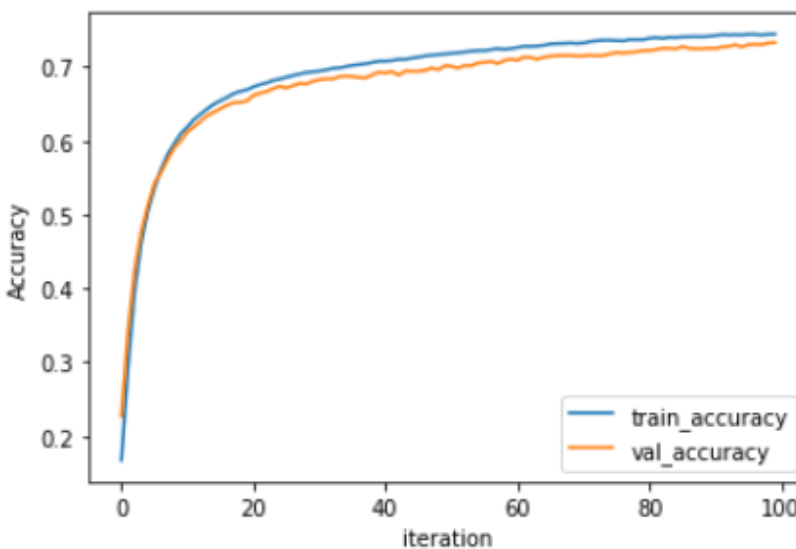
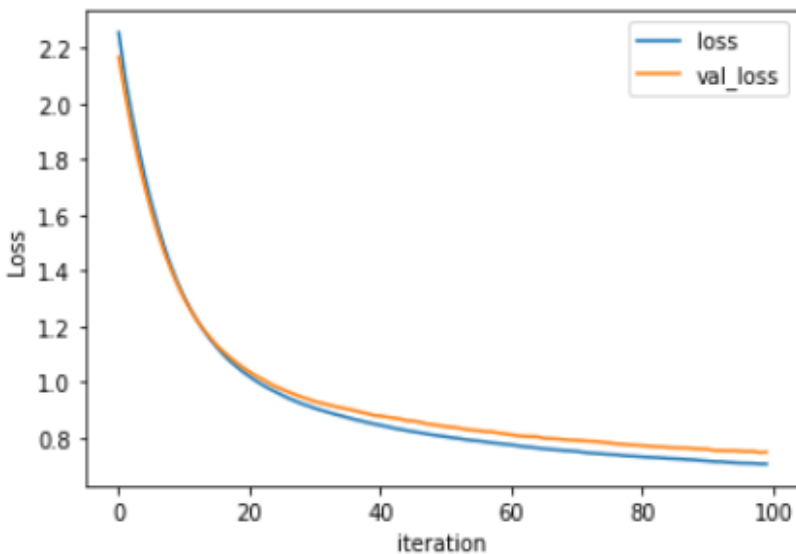
Discussion

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

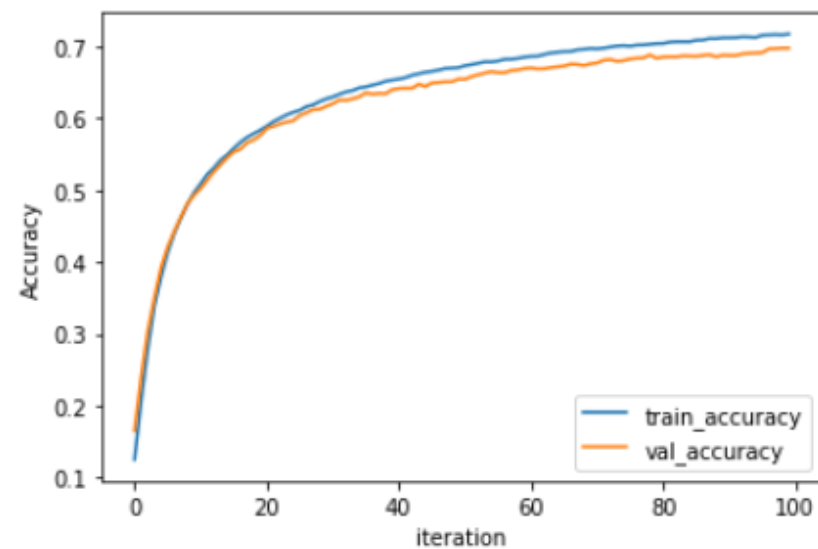
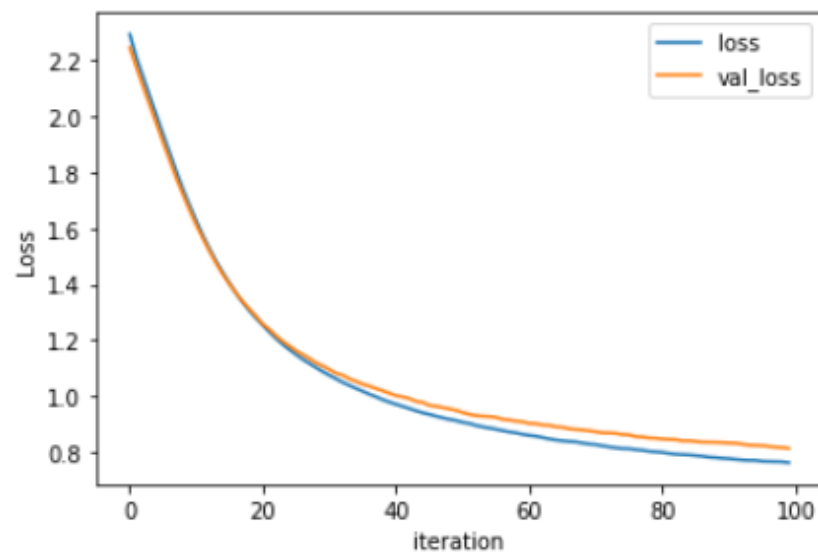


Discussion

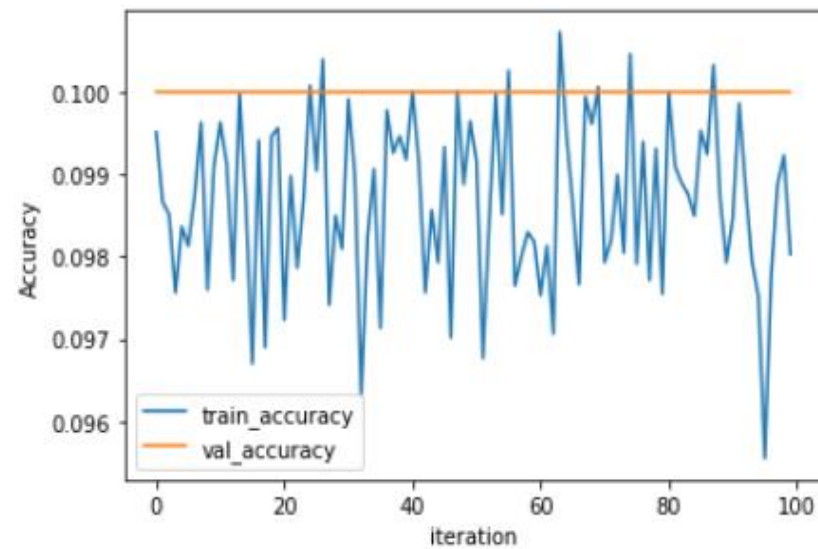
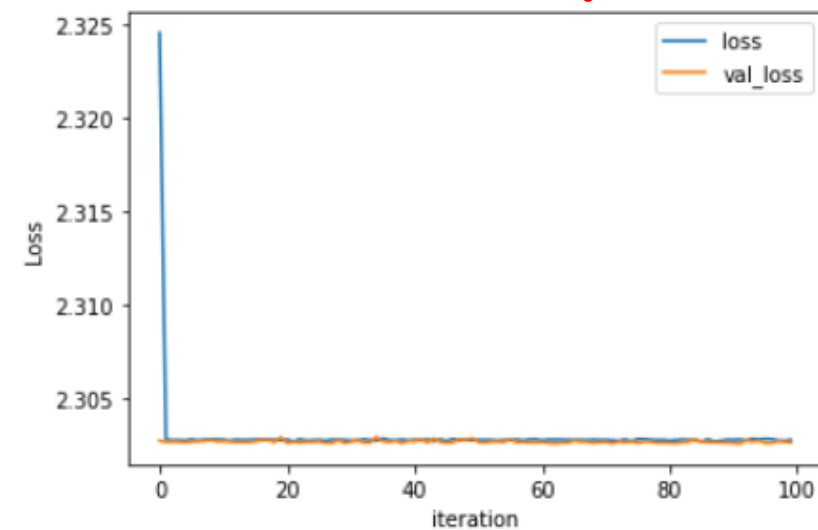
2 hidden layers



5 hidden layers (!)



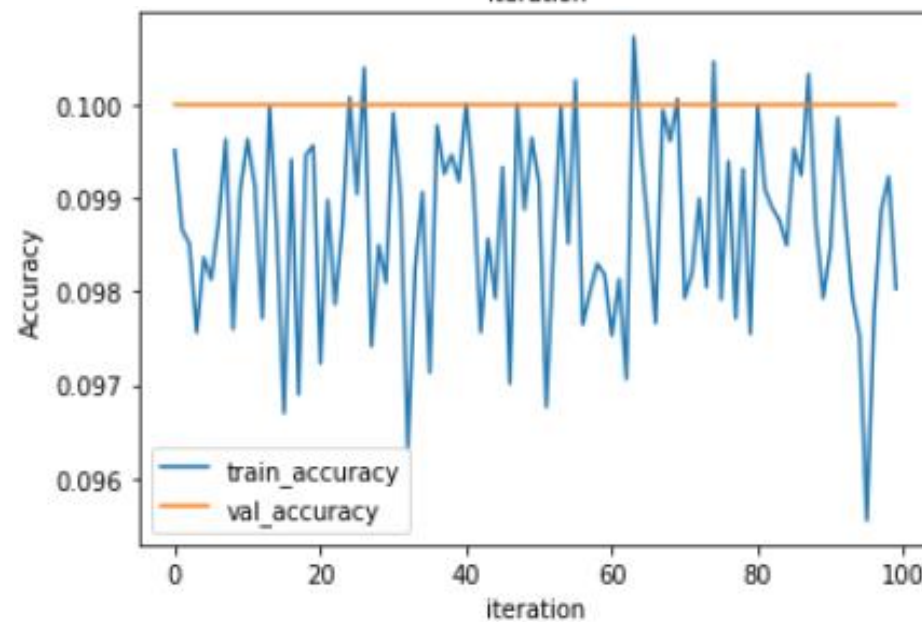
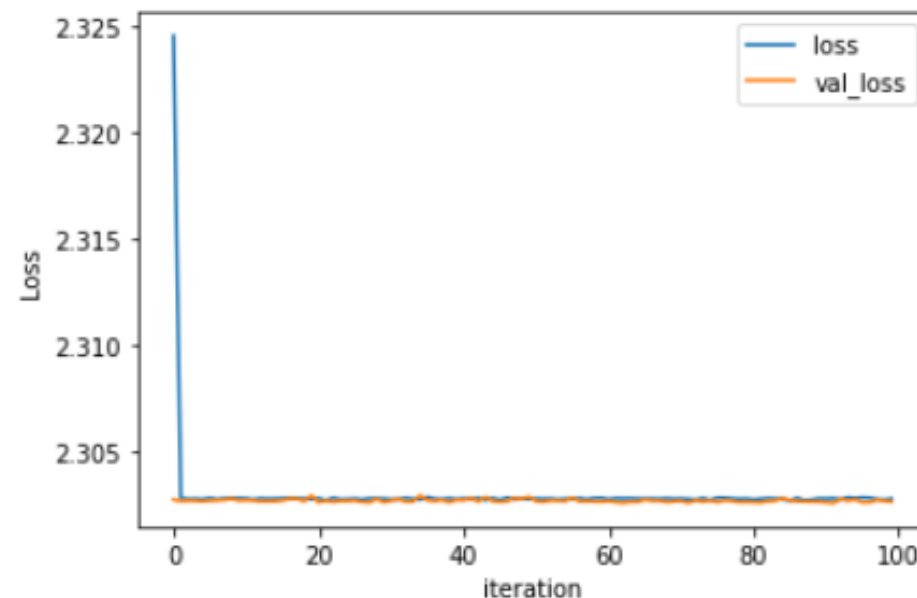
7 hidden layers



Discussion

Tensorflow

```
1 import tensorflow as tf
2 import tensorflow.keras as keras
3
4 # set seed
5 tf.random.set_seed(1234)
6 initializer = tf.keras.initializers.RandomNormal()
7
8 # create model
9 model = keras.Sequential()
10 model.add(keras.Input(shape=(784,)))
11 model.add(keras.layers.Dense(128, activation='sigmoid',
12                               kernel_initializer=initializer))
13 model.add(keras.layers.Dense(128, activation='sigmoid',
14                               kernel_initializer=initializer))
15 model.add(keras.layers.Dense(128, activation='sigmoid',
16                               kernel_initializer=initializer))
17 model.add(keras.layers.Dense(128, activation='sigmoid',
18                               kernel_initializer=initializer))
19 model.add(keras.layers.Dense(128, activation='sigmoid',
20                               kernel_initializer=initializer))
21 model.add(keras.layers.Dense(128, activation='sigmoid',
22                               kernel_initializer=initializer))
23 model.add(keras.layers.Dense(128, activation='sigmoid',
24                               kernel_initializer=initializer))
25 model.add(keras.layers.Dense(10, activation='softmax'))
```



Further Reading

Dying ReLU

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

Initialization

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

