



Module 05 – Extra Class

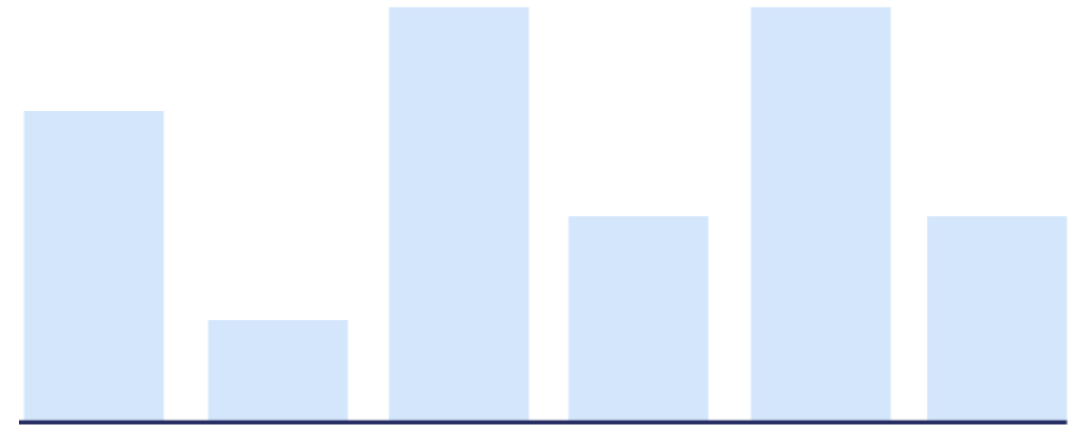
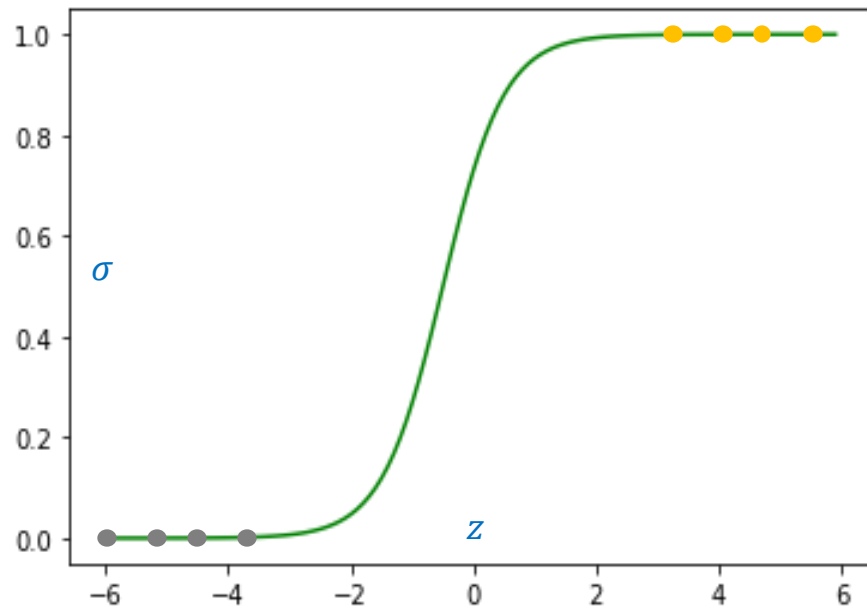
SOFTMAX REGRESSION

Nguyen Quoc Thai

Objectives

Logistic Regression (Review)

- ❖ Logistic Regression
- ❖ Sigmoid Function
- ❖ Gradient Descent



Softmax Regression

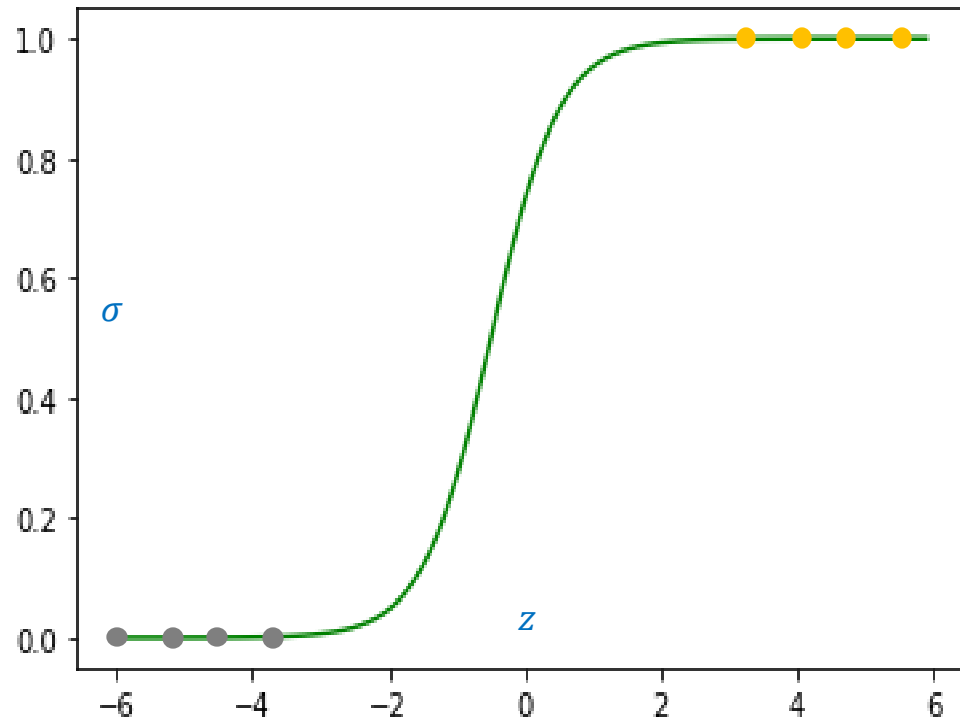
- ❖ Softmax Regression
- ❖ Softmax Function
- ❖ One Sample
- ❖ N Sample



Outline

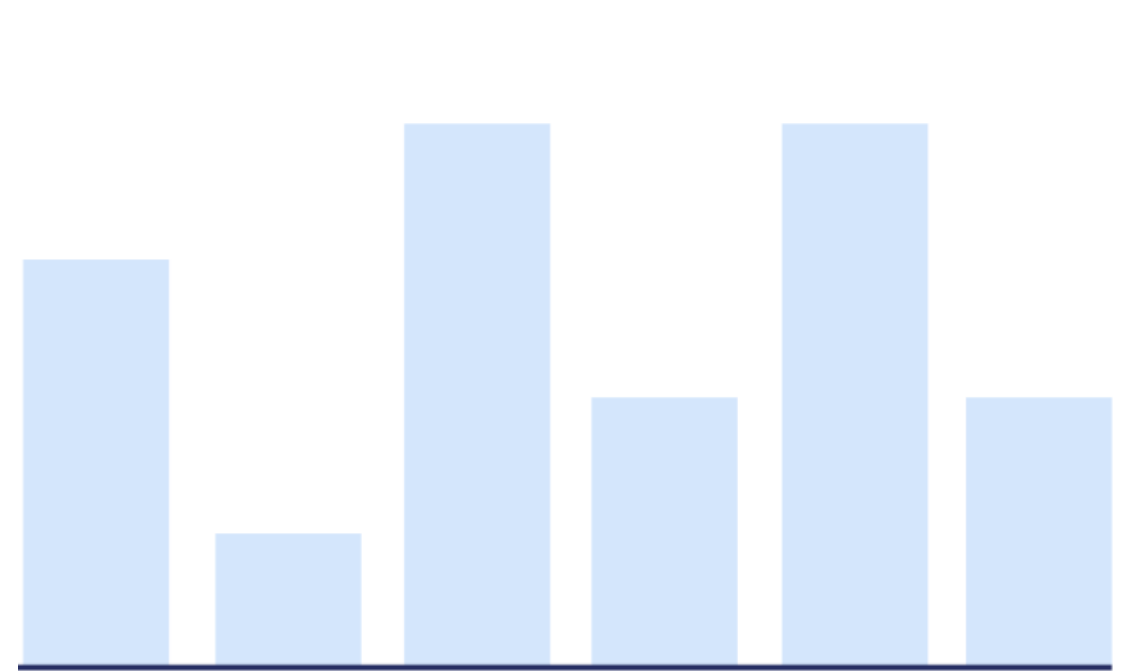
SECTION 1

Logistic Regression



SECTION 2

Softmax Regression

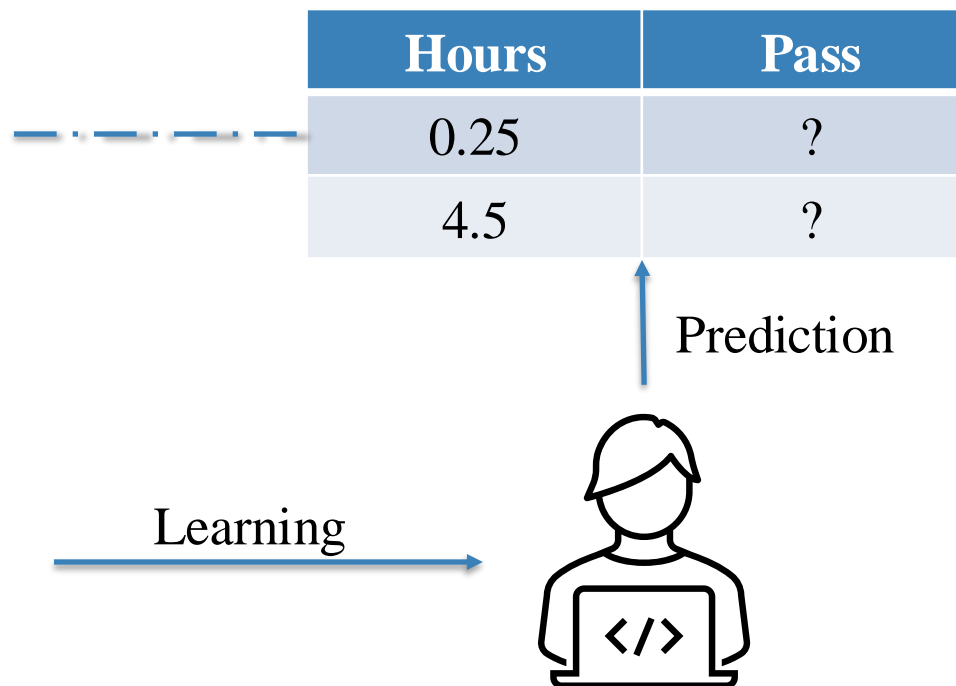


Logistic Regression



Classification Task

Hours	Pass
0.5	0
1.0	0
1.5	0
2.0	0
2.5	1
3.0	1
3.5	1
4.0	1



Logistic Regression



Classification Task

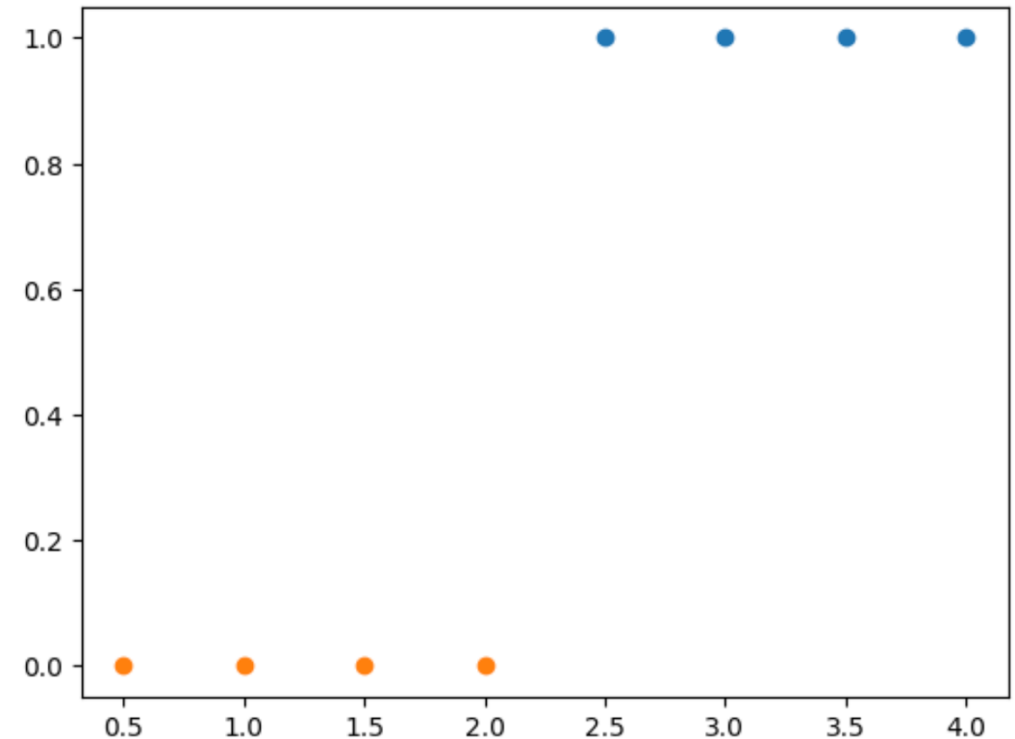
Hours	Pass
0.5	0
1.0	0
1.5	0
2.0	0
2.5	1
3.0	1
3.5	1
4.0	1

$$y = f(x)$$

Find a function to
fit the data

Sigmoid function

Visualization



Logistic Regression

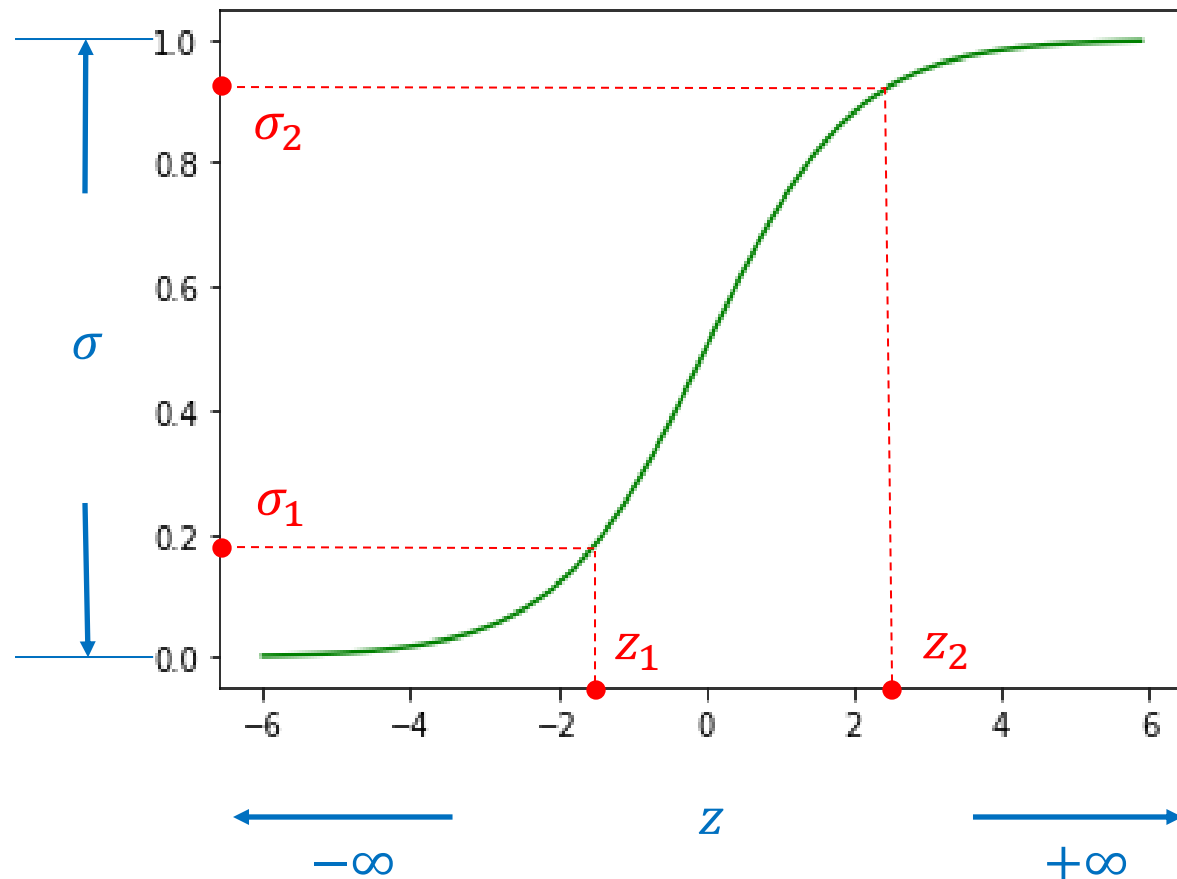


Sigmoid Function

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

$$z \in (-\infty, +\infty)$$

$$\sigma(z) \in (0, 1)$$



Logistic Regression

! Logistic Regression using Gradient Descent

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

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

$$\theta = \theta - \eta \nabla_{\theta} L$$

η is learning rate

$$\theta = \begin{bmatrix} b \\ w \end{bmatrix} = \begin{bmatrix} 0.1 \\ 0.1 \end{bmatrix}$$

$$x = \begin{bmatrix} 1.0 \\ 0.5 \end{bmatrix}$$

$$y = [0]$$

$$\eta = 0.1$$

Hours	Pass
0.5	0
1.0	0
1.5	0
2.0	0
2.5	1
3.0	1
3.5	1
4.0	1

Logistic Regression

! Logistic Regression using Gradient Descent

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

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

$$\theta = \theta - \eta \nabla_{\theta} L$$

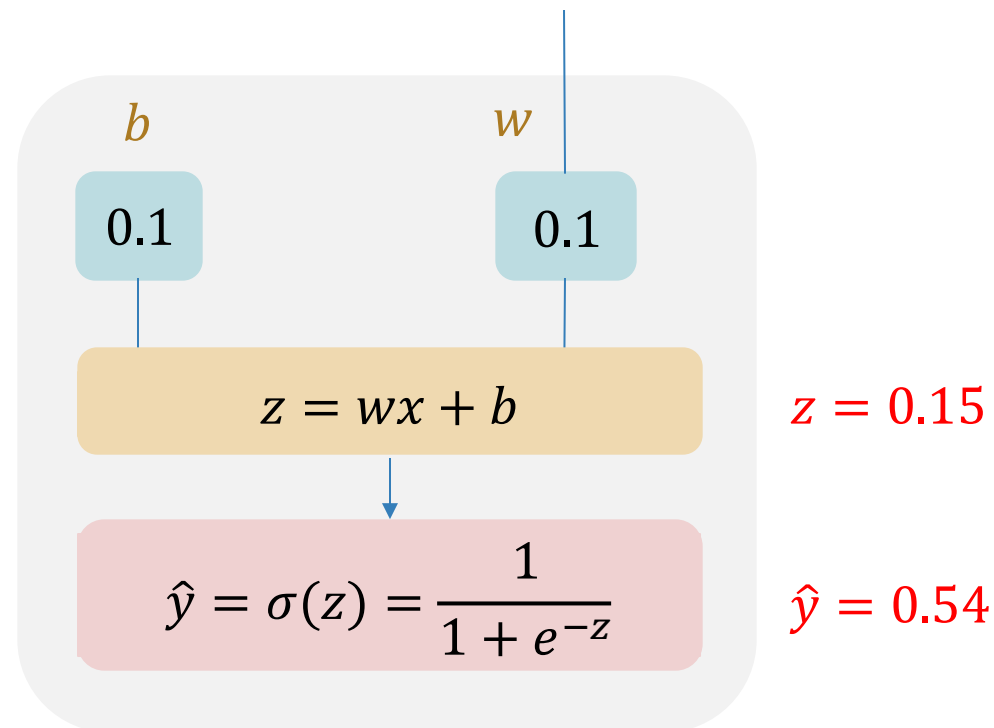
η is learning rate

$$\eta = 0.1$$

$$\theta = \begin{bmatrix} 0.1 \\ 0.1 \end{bmatrix}$$

$$x = \begin{bmatrix} 1.0 \\ 0.5 \end{bmatrix}$$

$$y = [0]$$



Logistic Regression

! Logistic Regression using Gradient Descent

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

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

$$\theta = \theta - \eta \nabla_{\theta} L$$

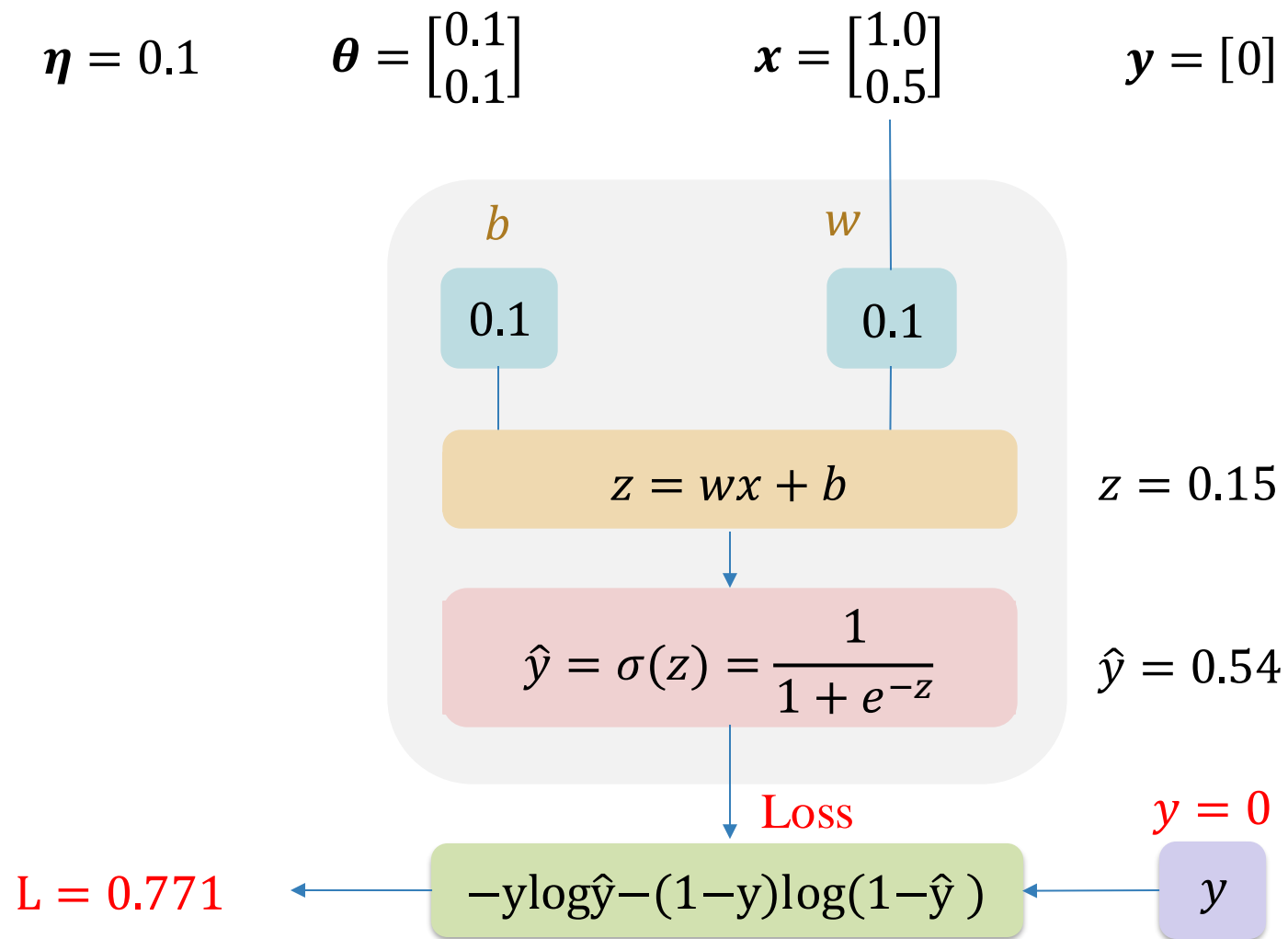
η is learning rate

$$\eta = 0.1$$

$$\theta = \begin{bmatrix} 0.1 \\ 0.1 \end{bmatrix}$$

$$x = \begin{bmatrix} 1.0 \\ 0.5 \end{bmatrix}$$

$$y = [0]$$



Logistic Regression

! Logistic Regression using Gradient Descent

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

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

$$\theta = \theta - \eta \nabla_{\theta} L$$

η is learning rate

$$\eta = 0.1$$

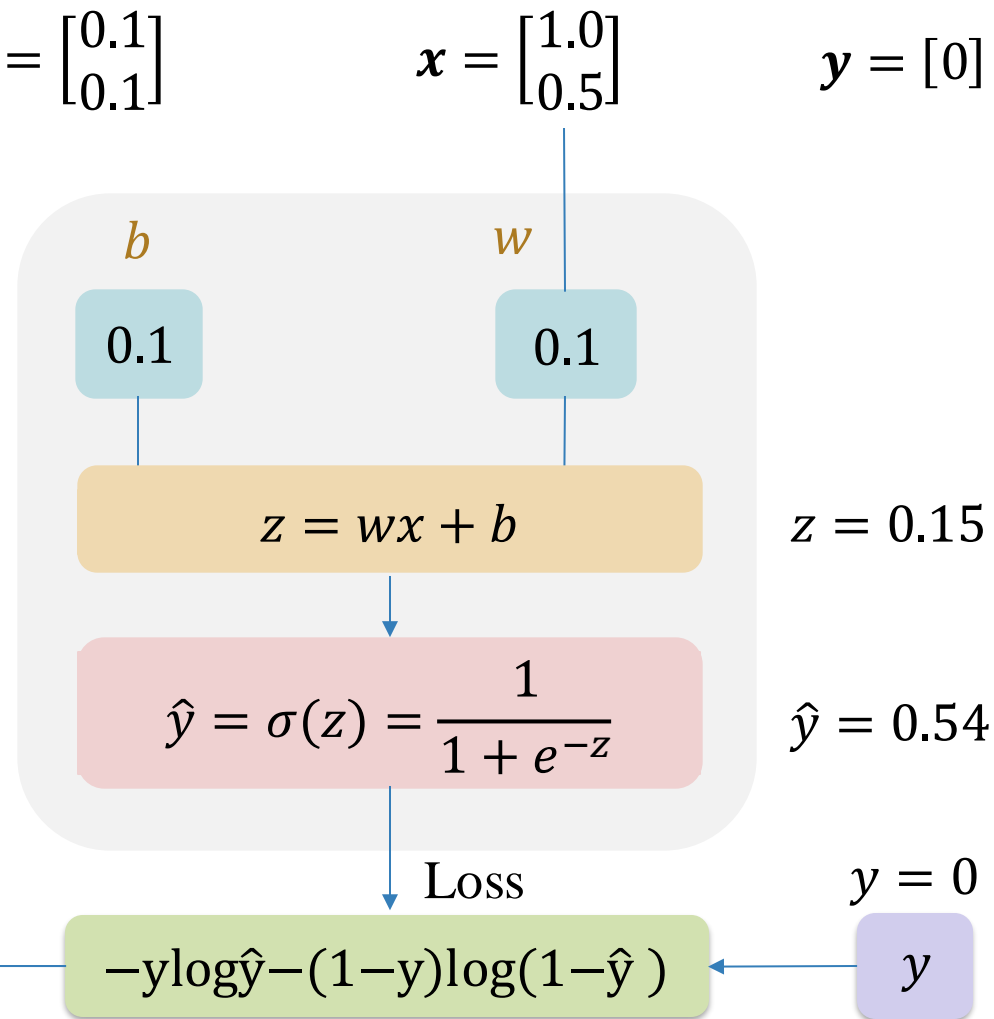
$$\theta = \begin{bmatrix} 0.1 \\ 0.1 \end{bmatrix}$$

$$x = \begin{bmatrix} 1.0 \\ 0.5 \end{bmatrix}$$

$$y = [0]$$

$$\begin{aligned} \nabla_{\theta} L &= x(\hat{y} - y) \\ &= \begin{bmatrix} 1.0 \\ 0.5 \end{bmatrix} * 0.54 \\ &= \begin{bmatrix} 0.54 \\ 0.27 \end{bmatrix} \end{aligned}$$

$$L = 0.771$$



Logistic Regression

! Logistic Regression using Gradient Descent

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

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

$$\theta = \theta - \eta \nabla_{\theta} L$$

η is learning rate

$$\eta = 0.1$$

$$\theta = \begin{bmatrix} 0.1 \\ 0.1 \end{bmatrix}$$

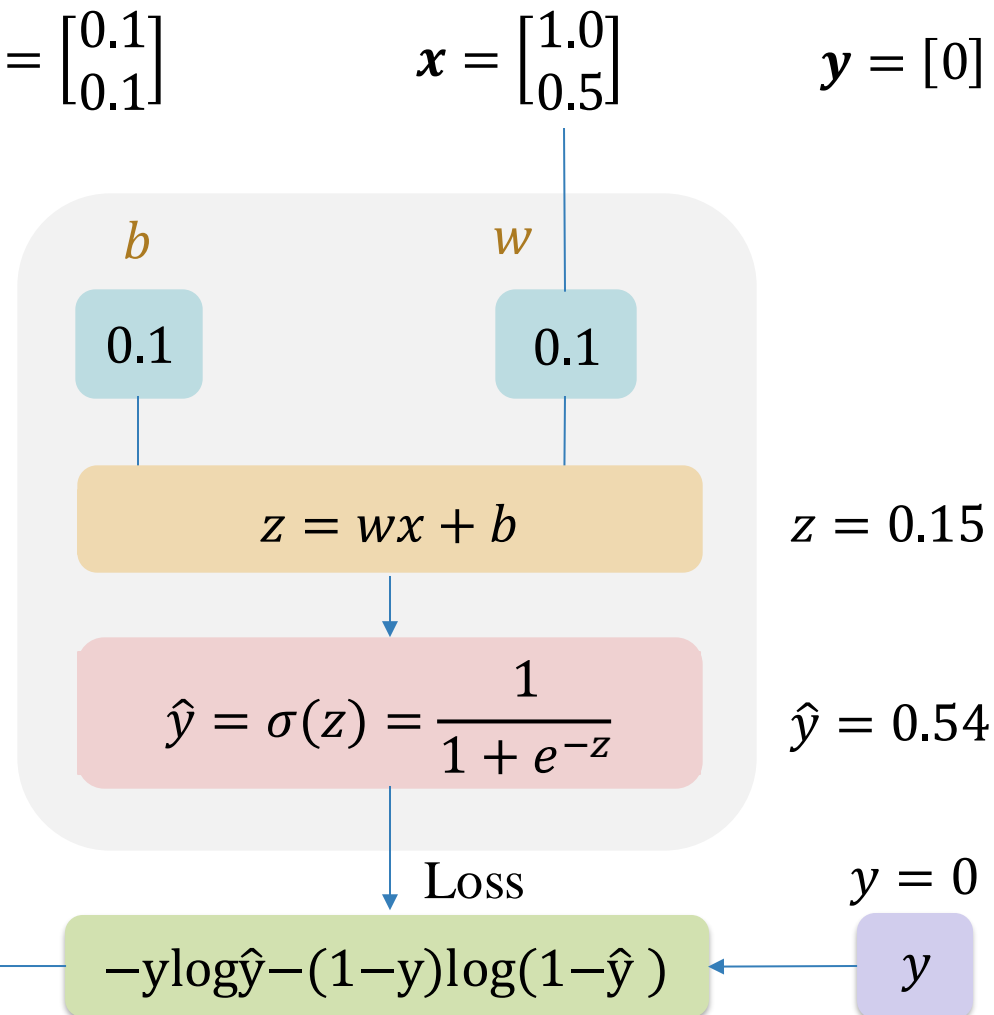
$$x = \begin{bmatrix} 1.0 \\ 0.5 \end{bmatrix}$$

$$y = [0]$$

$$\begin{aligned} \theta &= \theta - \eta \nabla_{\theta} L \\ &= \begin{bmatrix} 0.046 \\ 0.073 \end{bmatrix} \end{aligned}$$

$$\begin{aligned} \nabla_{\theta} L &= x(\hat{y} - y) \\ &= \begin{bmatrix} 1.0 \\ 0.5 \end{bmatrix} * 0.54 \\ &= \begin{bmatrix} 0.54 \\ 0.27 \end{bmatrix} \end{aligned}$$

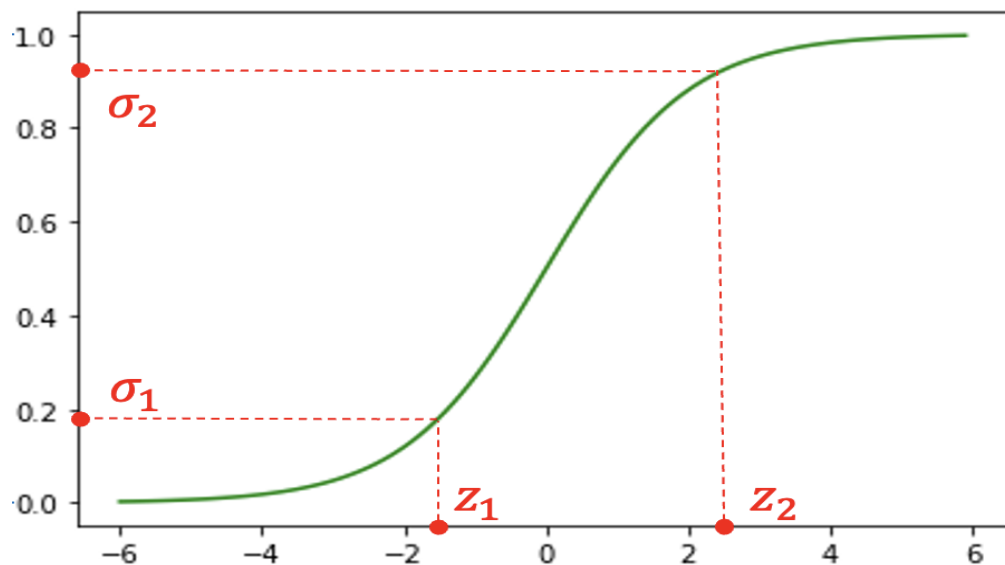
$$L = 0.771$$



Outline

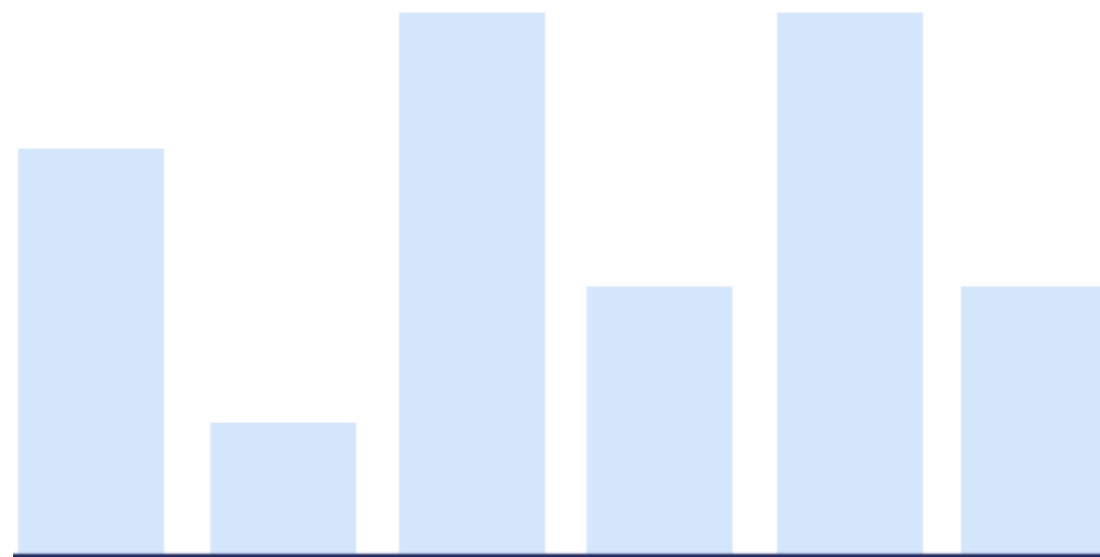
SECTION 1

Logistic Regression



SECTION 2

Softmax Regression



Softmax Regression



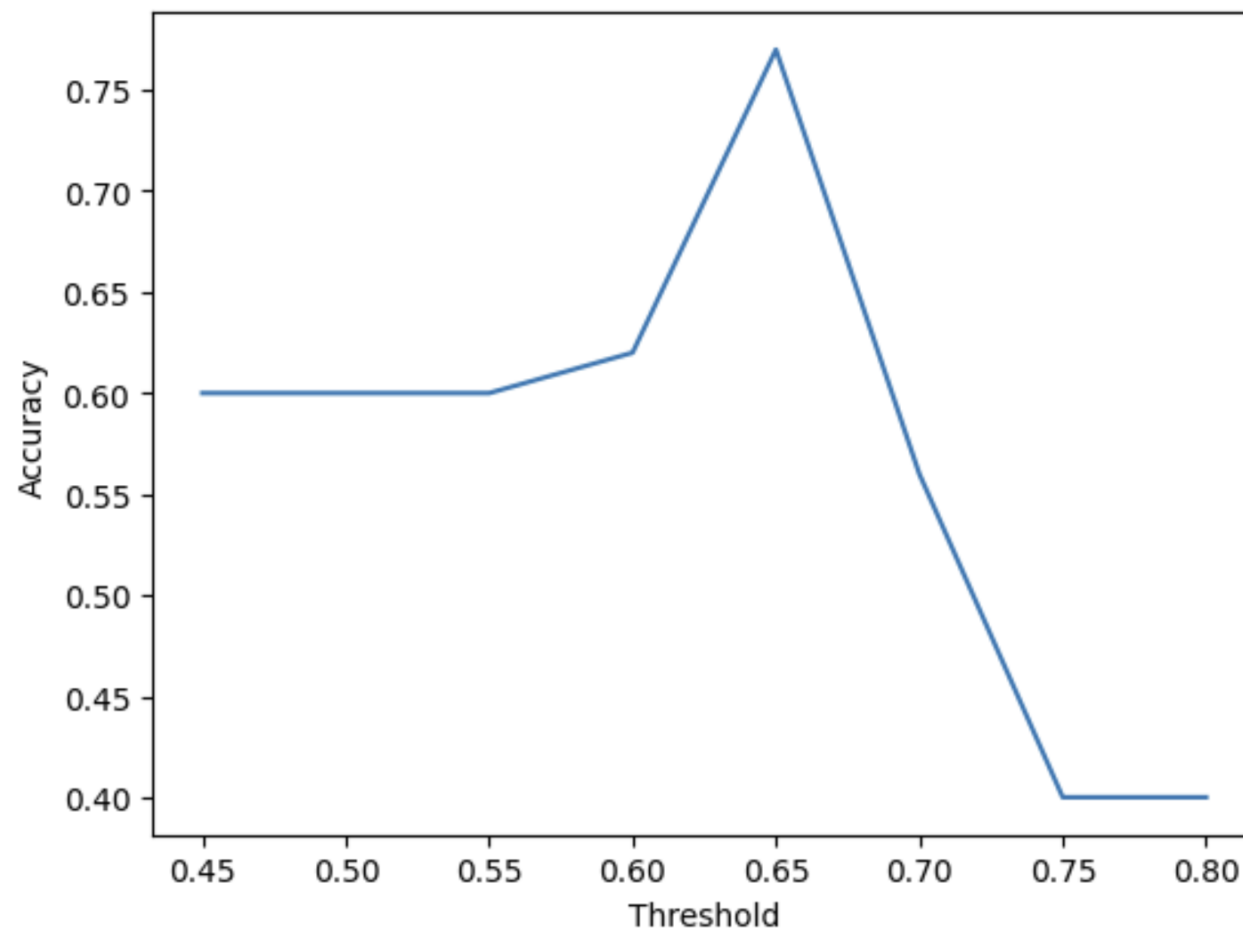
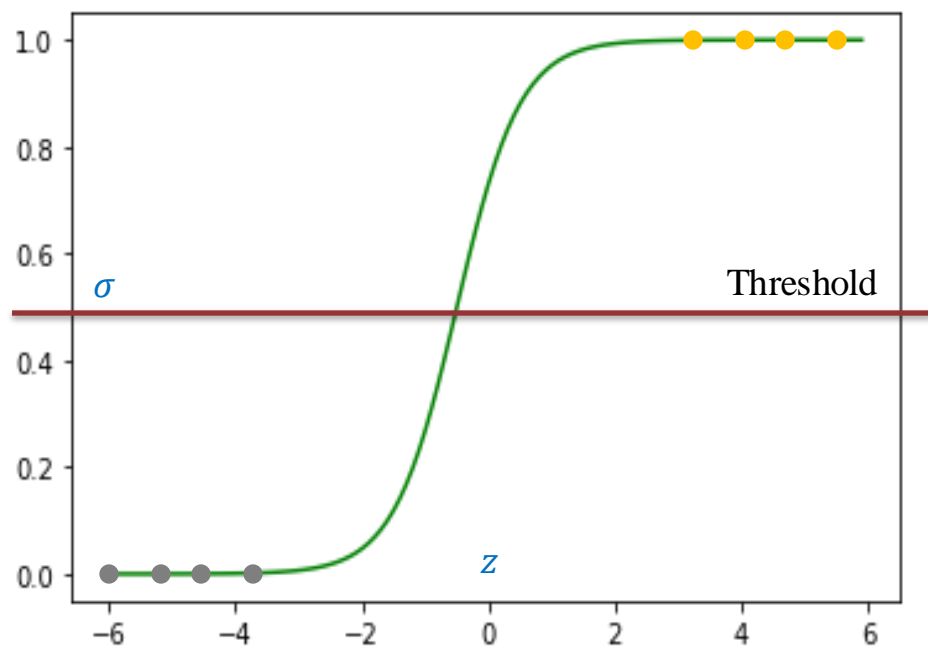
Problem

Sigmoid function

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

$$z \in (-\infty, +\infty)$$

$$\sigma(z) \in (0, 1)$$



Softmax Regression



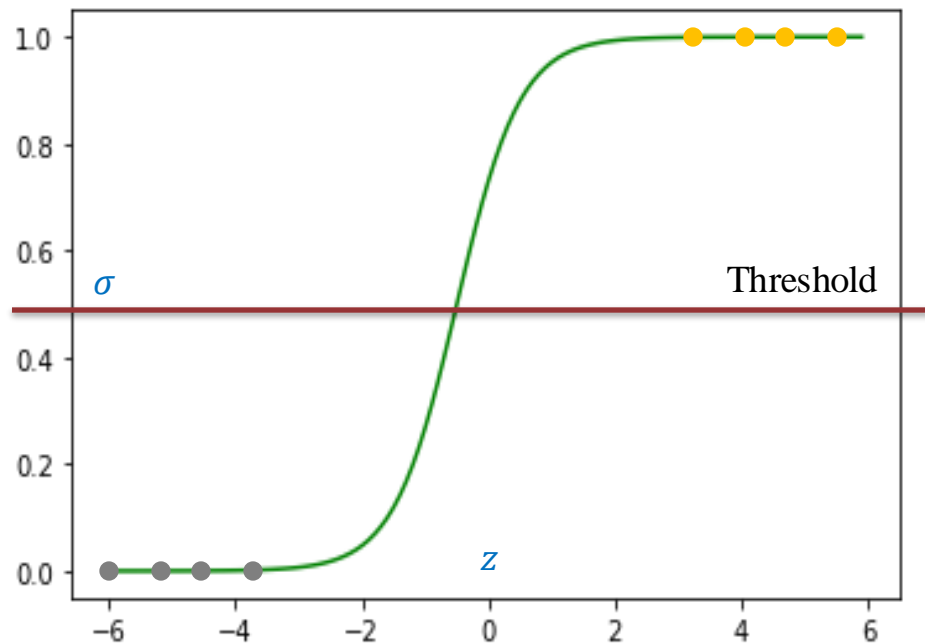
Problem

Sigmoid function

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

$$z \in (-\infty, +\infty)$$

$$\sigma(z) \in (0, 1)$$



Hours	Pass
0.5	0
1.0	0
1.5	1
2.0	1

Classes: {0, 1}
Binary Classification

Hours	Score
0.5	0
1.0	0
1.5	1
2.0	1
2.5	2
3.0	2
3.5	3
4.0	3

Classes: {0, 1, 2, 3}
Multi-class Classification

Softmax Regression



Problem

Classes: {0, 1}

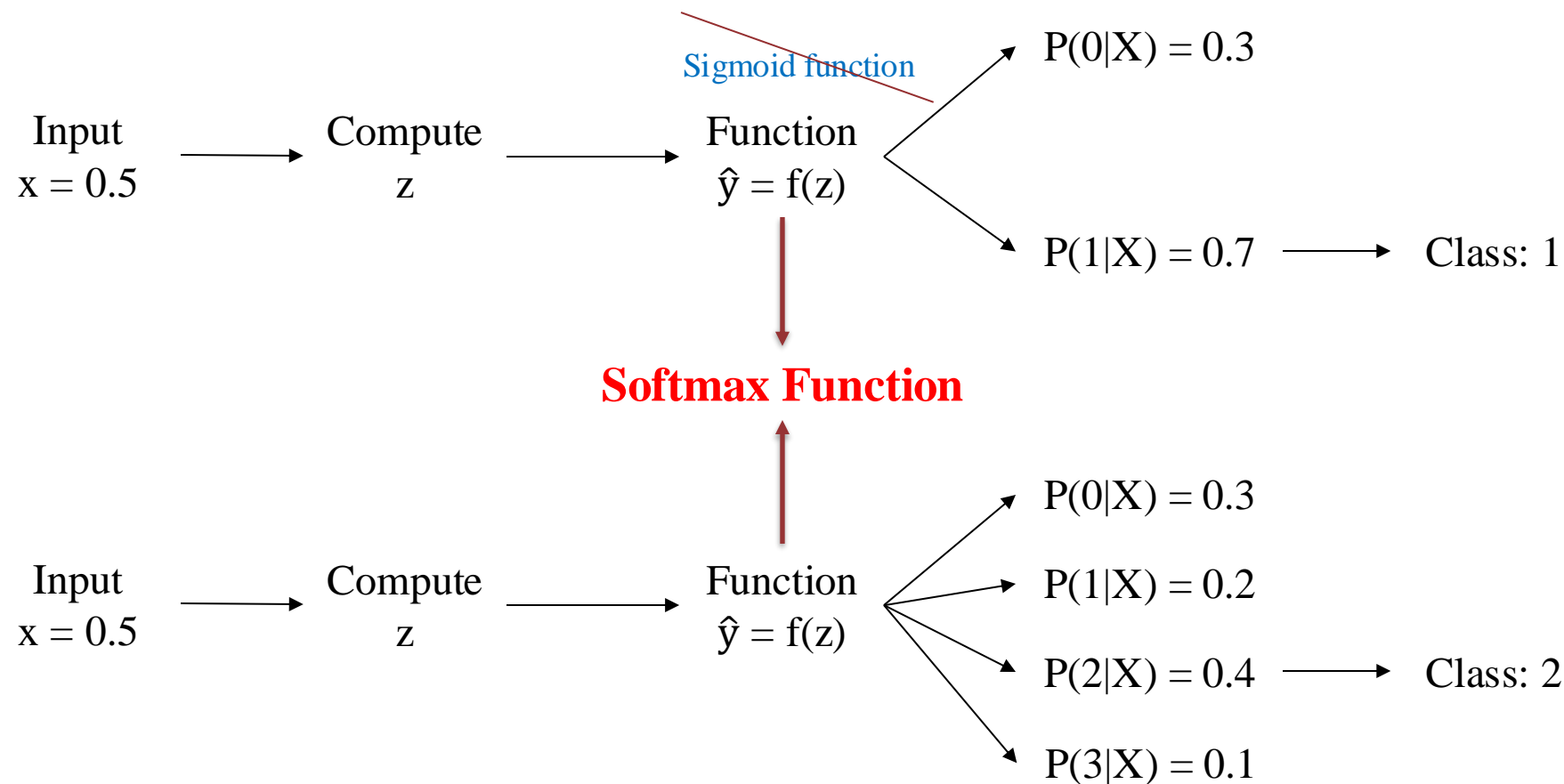
Binary Classification

Hours	Pass
0.5	0
2.0	1

Classes: {0, 1, 2, 3}

Multi-class Classification

Hours	Score
0.5	0
1.5	1
3.0	2
4.0	3



Softmax Regression

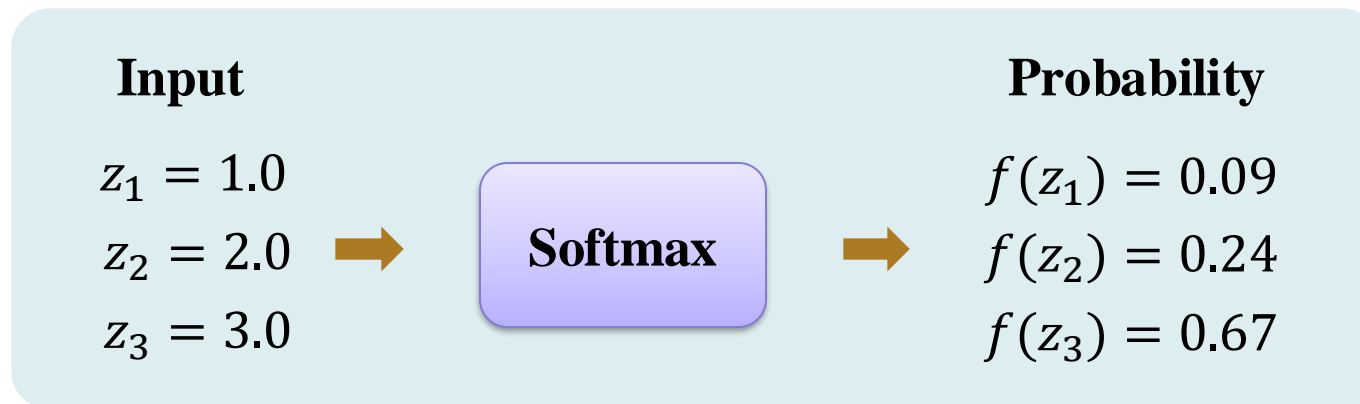
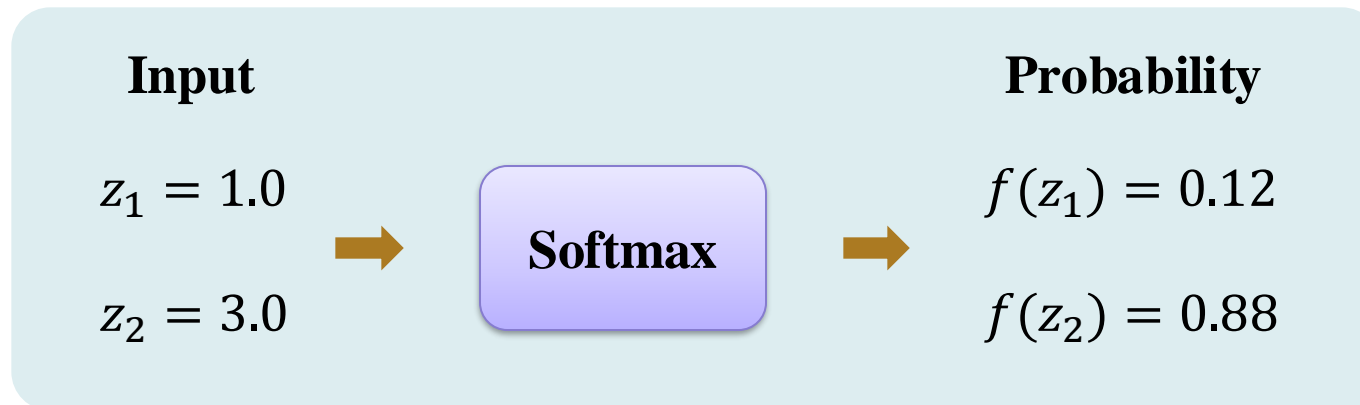


Softmax Function

$$P_i = f(z_i) = \frac{e^{z_i}}{\sum_j e^{z_j}}$$

$$0 \leq f(z_i) \leq 1$$

$$\sum_i f(z_i) = 1$$



Softmax Regression

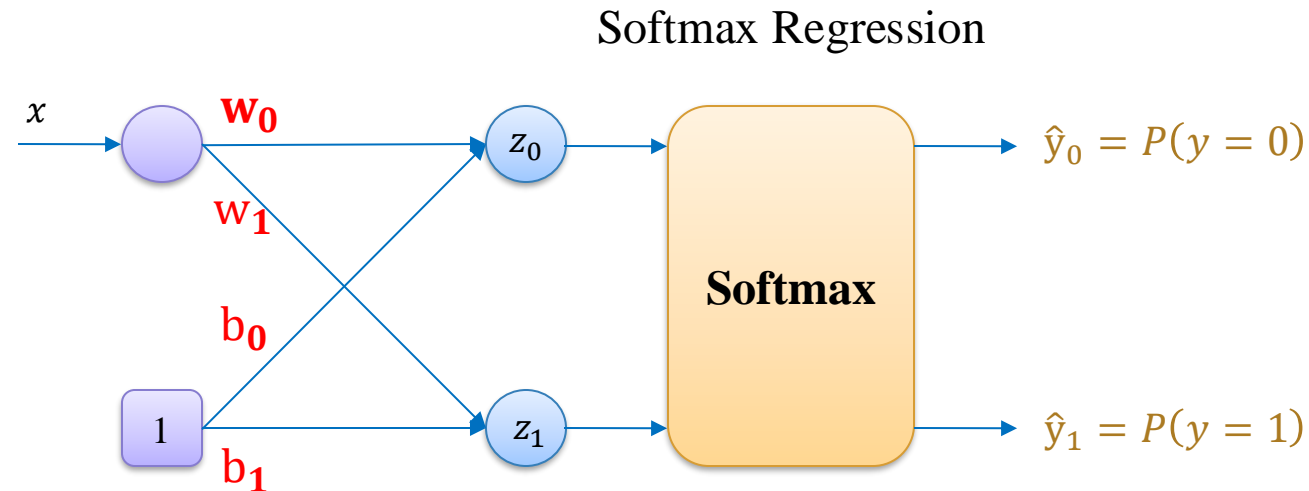
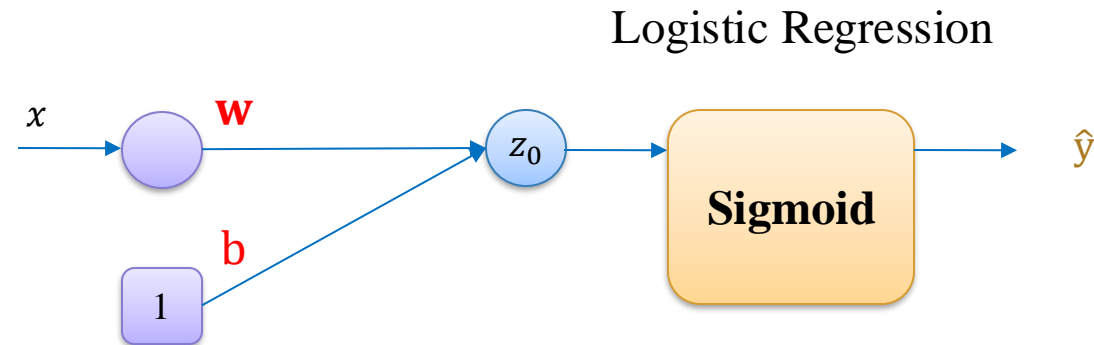


Parameters

Classes: $\{0, 1\}$
Binary Classification

Hours	Pass
0.5	0
2.0	1

#feature: 1
#class: 2



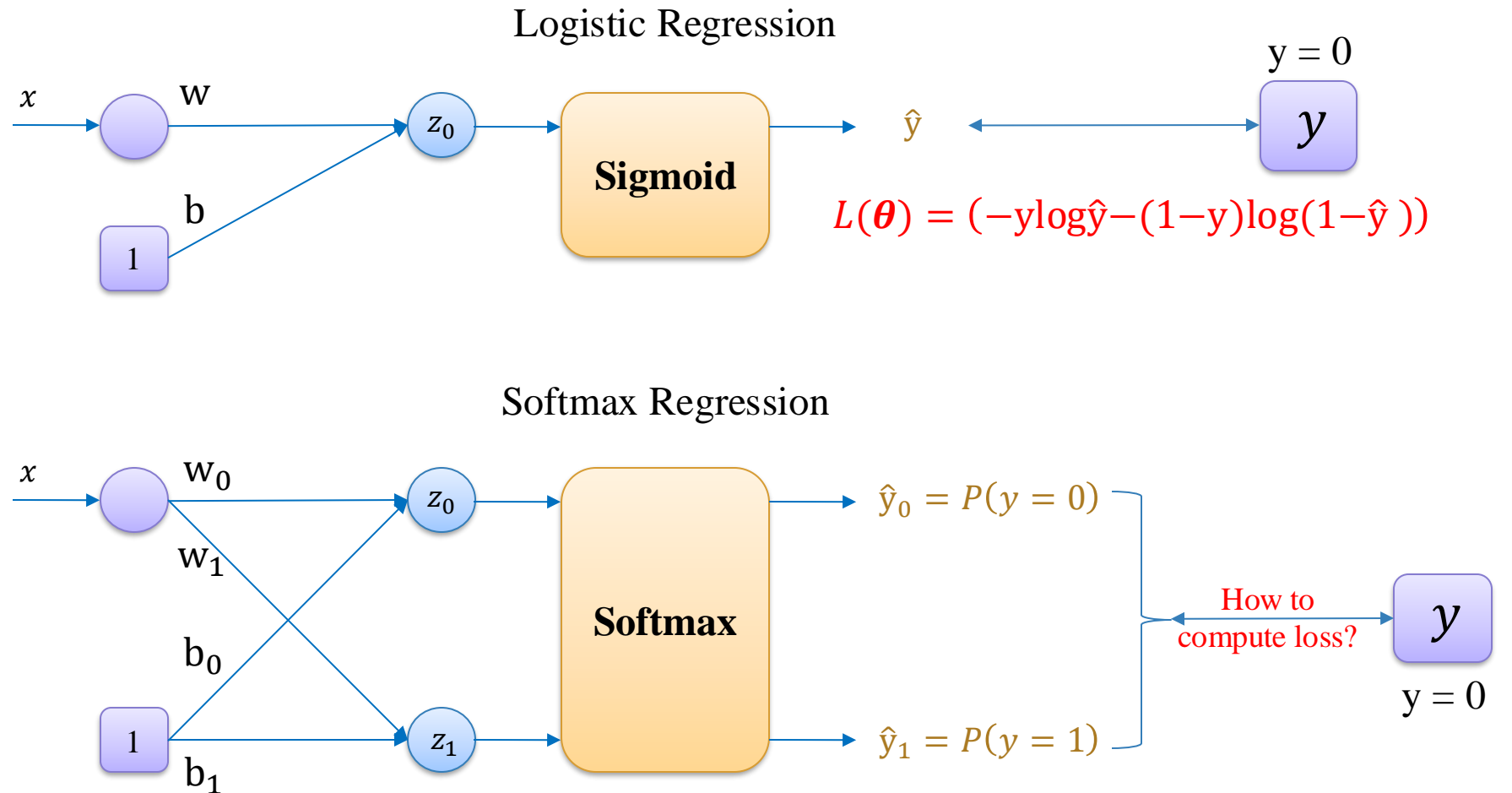
Softmax Regression

! Loss Function

Classes: {0, 1}
Binary Classification

Hours	Pass
0.5	0
2.0	1

#feature: 1
#class: 2



Softmax Regression



One-hot Encoding

$$\mathbf{y} = \begin{bmatrix} y_0 \\ \vdots \\ y_C \end{bmatrix}$$

$$y_i \in \{0,1\}$$

$$\sum_i y_i = 1$$

$$C = \text{\#classes}$$

Classes: {0, 1}
Binary Classification

Hours	Pass
0.5	0
2.0	1

#feature: 1
#class: 2

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

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

Classes: {0, 1, 2}
Multi-class Classification

Hours	Score
0.5	0
1.5	1
3.0	2

#feature: 1
#class: 3

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

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

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

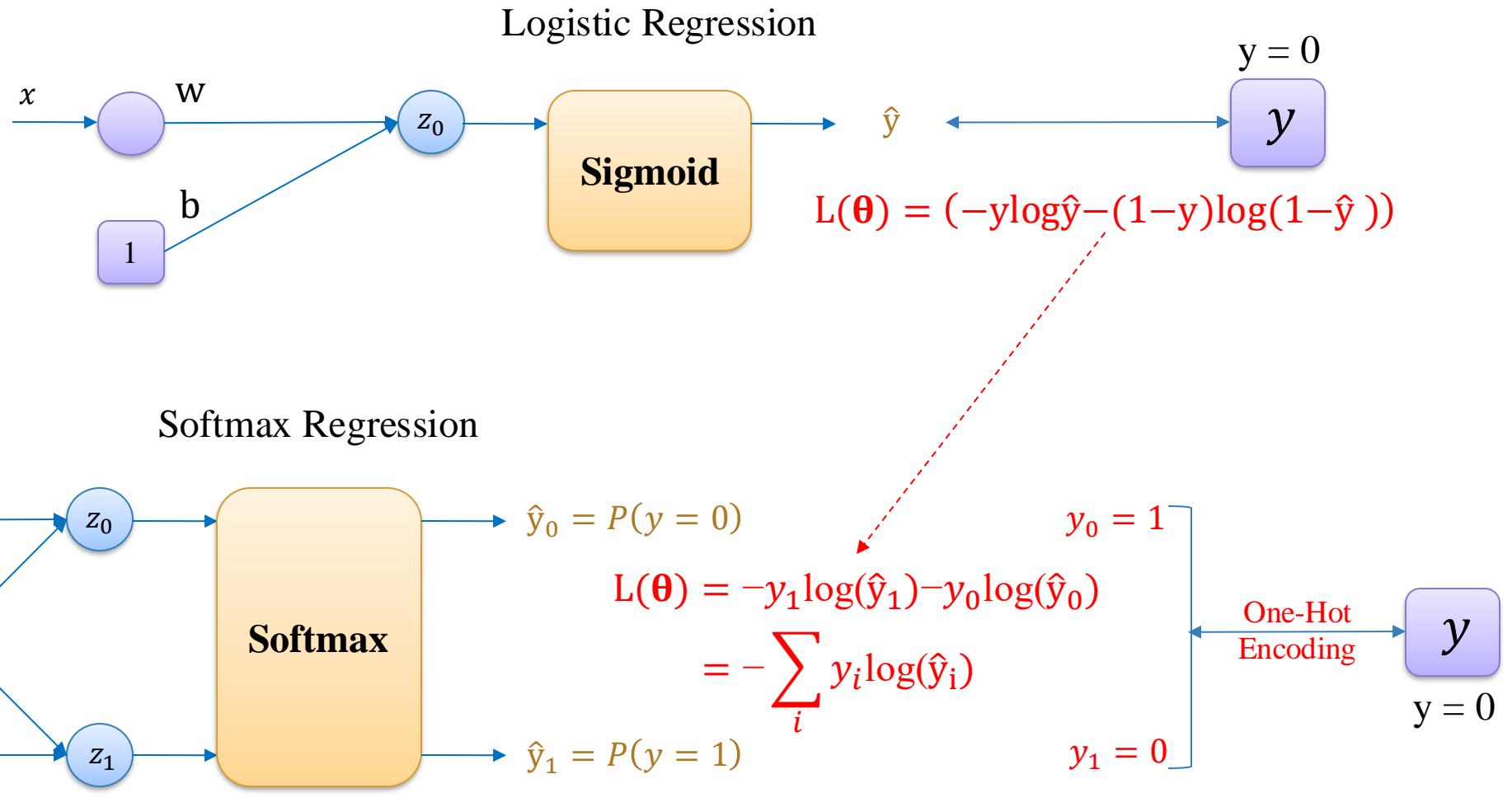
Softmax Regression

! Loss Function

Classes: {0, 1}
Binary Classification

Hours	Pass
0.5	0
2.0	1

#feature: 1
#class: 2



Softmax Regression

! Softmax Regression using Gradient Descent

1) Pick a sample from training data

2) Compute output \hat{y}

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

$$\mathbf{d} = [1 \dots 1] \mathbf{e}^{\mathbf{z}}$$

$$\hat{\mathbf{y}} = \mathbf{e}^{\mathbf{z}} \oslash \mathbf{d}$$

\oslash is
Hadamard
division

3) Compute loss (cross-entropy)

$$L(\boldsymbol{\theta}) = -\mathbf{y}^T \log \hat{\mathbf{y}}$$

4) Compute derivative

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

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} L$$

η is learning rate

$$\mathbf{x} = \begin{bmatrix} 1.0 \\ 0.5 \end{bmatrix}$$

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

One-hot encoding for label

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

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

$$\boldsymbol{\theta} = \begin{bmatrix} b_0 & b_1 \\ w_0 & w_1 \end{bmatrix} = \begin{bmatrix} 0.1 & 0.2 \\ 0.3 & 0.4 \end{bmatrix}$$

$$\eta = 0.1$$

Hours	Pass
0.5	0
1.0	0
1.5	1
2.0	1

Softmax Regression

! Softmax Regression using Gradient Descent

1) Pick a sample from training data

2) Compute output \hat{y}

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

$$\mathbf{d} = [1 \dots 1] \mathbf{e}^{\mathbf{z}}$$

$$\hat{\mathbf{y}} = \mathbf{e}^{\mathbf{z}} \oslash \mathbf{d}$$

\oslash is
Hadamard
division

3) Compute loss (cross-entropy)

$$L(\boldsymbol{\theta}) = -\mathbf{y}^T \log \hat{\mathbf{y}}$$

4) Compute derivative

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

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} L$$

η is learning rate

$$\eta = 0.1$$

$$\mathbf{x} = \begin{bmatrix} 1.0 \\ 0.5 \end{bmatrix}$$

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

$$\boldsymbol{\theta} = \begin{bmatrix} 0.1 & 0.2 \\ 0.3 & 0.4 \end{bmatrix}$$

Softmax Regression

! Softmax Regression using Gradient Descent

1) Pick a sample from training data

2) Compute output \hat{y}

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

$$\mathbf{d} = [1 \dots 1] \mathbf{e}^{\mathbf{z}}$$

$$\hat{\mathbf{y}} = \mathbf{e}^{\mathbf{z}} \oslash \mathbf{d}$$

\oslash is
Hadamard
division

3) Compute loss (cross-entropy)

$$L(\boldsymbol{\theta}) = -\mathbf{y}^T \log \hat{\mathbf{y}}$$

4) Compute derivative

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

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} L$$

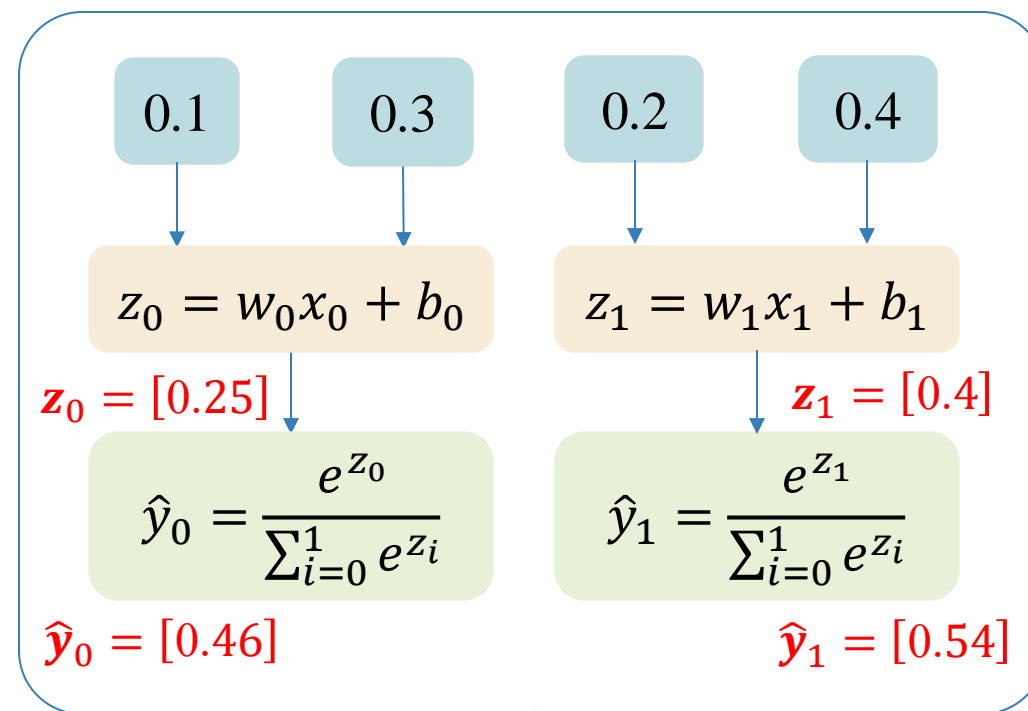
η is learning rate

$$\eta = 0.1$$

$$\mathbf{x} = \begin{bmatrix} 1.0 \\ 0.5 \end{bmatrix}$$

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

$$\boldsymbol{\theta} = \begin{bmatrix} 0.1 & 0.2 \\ 0.3 & 0.4 \end{bmatrix}$$



Softmax Regression

! Softmax Regression using Gradient Descent

1) Pick a sample from training data

2) Compute output \hat{y}

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

$$\mathbf{d} = [1 \dots 1] \mathbf{e}^{\mathbf{z}}$$

$$\hat{\mathbf{y}} = \mathbf{e}^{\mathbf{z}} \oslash \mathbf{d}$$

\oslash is
Hadamard
division

3) Compute loss (cross-entropy)

$$L(\boldsymbol{\theta}) = -\mathbf{y}^T \log \hat{\mathbf{y}}$$

4) Compute derivative

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

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} L$$

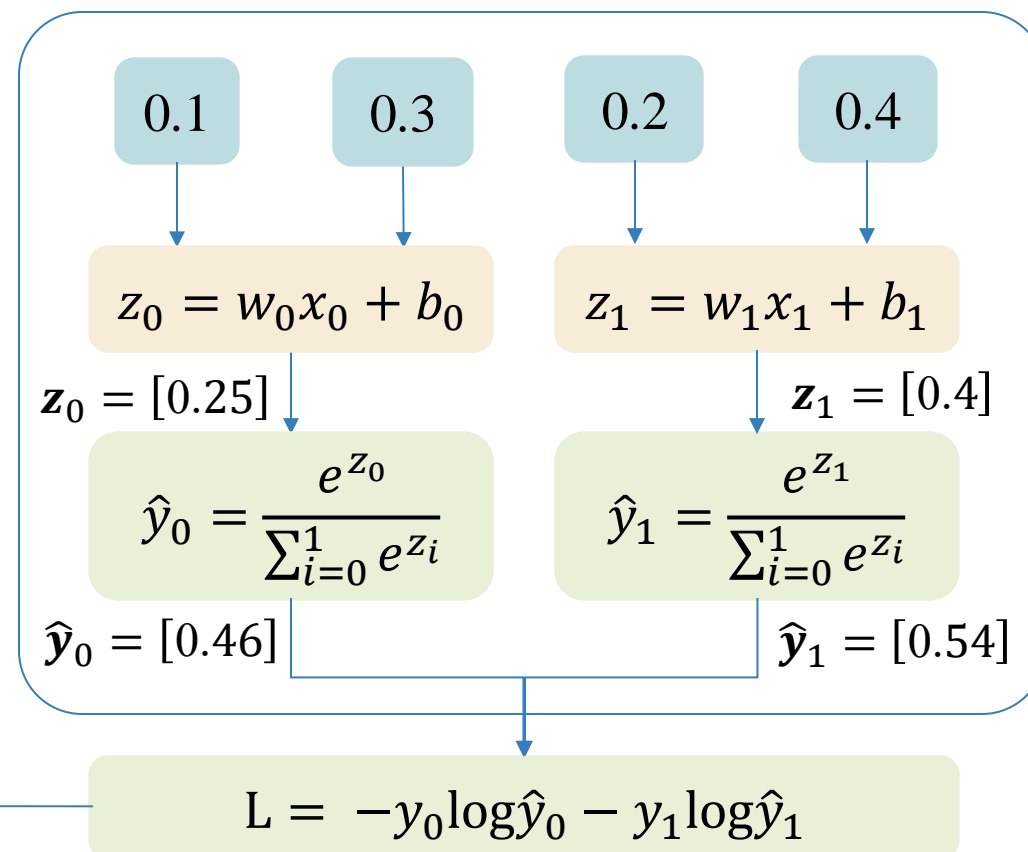
η is learning rate

$$\eta = 0.1$$

$$\mathbf{x} = \begin{bmatrix} 1.0 \\ 0.5 \end{bmatrix}$$

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

$$\boldsymbol{\theta} = \begin{bmatrix} 0.1 & 0.2 \\ 0.3 & 0.4 \end{bmatrix}$$



$$L = [0.77]$$

Softmax Regression

! Softmax Regression using Gradient Descent

1) Pick a sample from training data

2) Compute output \hat{y}

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

$$\mathbf{d} = [1 \dots 1] \mathbf{e}^{\mathbf{z}}$$

\emptyset is Hadamard division

$$\hat{\mathbf{y}} = \mathbf{e}^{\mathbf{z}} \emptyset \mathbf{d}$$

3) Compute loss (cross-entropy)

$$L(\boldsymbol{\theta}) = -\mathbf{y}^T \log \hat{\mathbf{y}}$$

4) Compute derivative

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

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} L$$

η is learning rate

$$\eta = 0.1$$

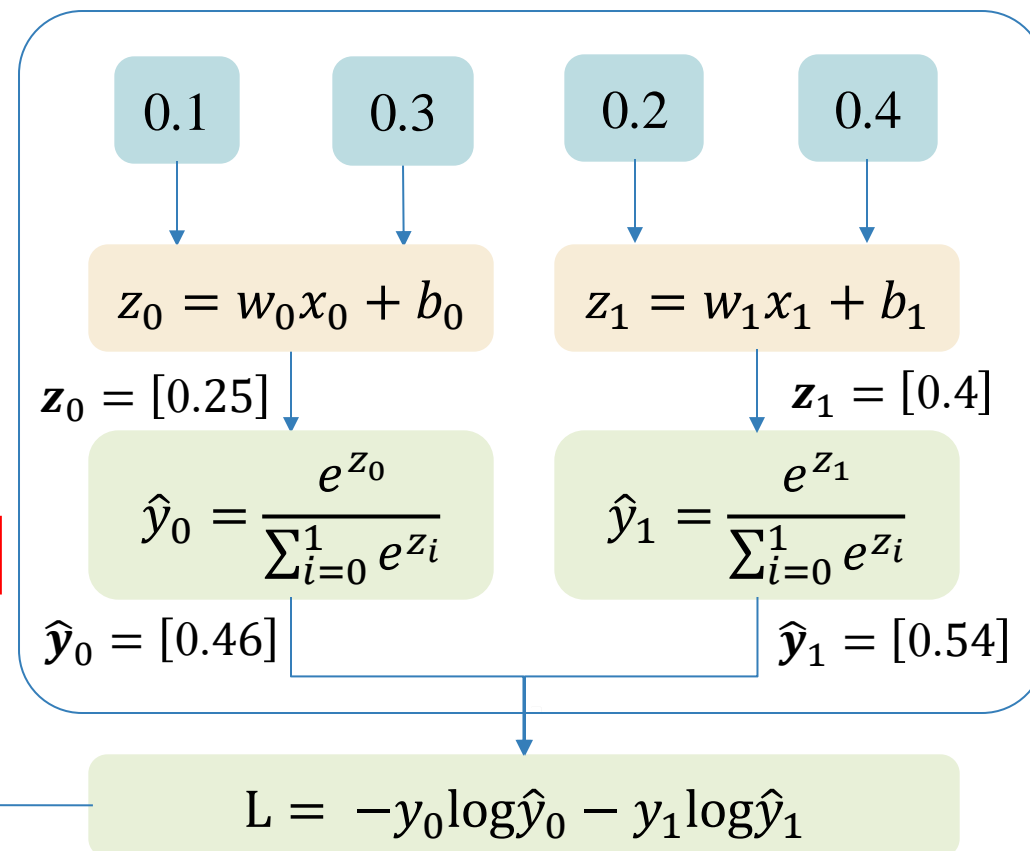
$$\mathbf{x} = \begin{bmatrix} 1.0 \\ 0.5 \end{bmatrix}$$

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

$$\boldsymbol{\theta} = \begin{bmatrix} 0.1 & 0.2 \\ 0.3 & 0.4 \end{bmatrix}$$

$$\begin{aligned} \nabla_{\boldsymbol{\theta}} L &= \mathbf{x}(\hat{\mathbf{y}} - \mathbf{y})^T \\ &= \begin{bmatrix} -0.54 & 0.54 \\ -0.27 & 0.27 \end{bmatrix} \end{aligned}$$

$$L = [0.77]$$



Softmax Regression

! Softmax Regression using Gradient Descent

1) Pick a sample from training data

2) Compute output \hat{y}

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

$$\mathbf{d} = [1 \dots 1] \mathbf{e}^{\mathbf{z}}$$

$$\hat{\mathbf{y}} = \mathbf{e}^{\mathbf{z}} \oslash \mathbf{d}$$

\oslash is
Hadamard
division

3) Compute loss (cross-entropy)

$$L(\boldsymbol{\theta}) = -\mathbf{y}^T \log \hat{\mathbf{y}}$$

4) Compute derivative

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

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} L$$

η is learning rate

$$\eta = 0.1$$

$$\mathbf{x} = \begin{bmatrix} 1.0 \\ 0.5 \end{bmatrix}$$

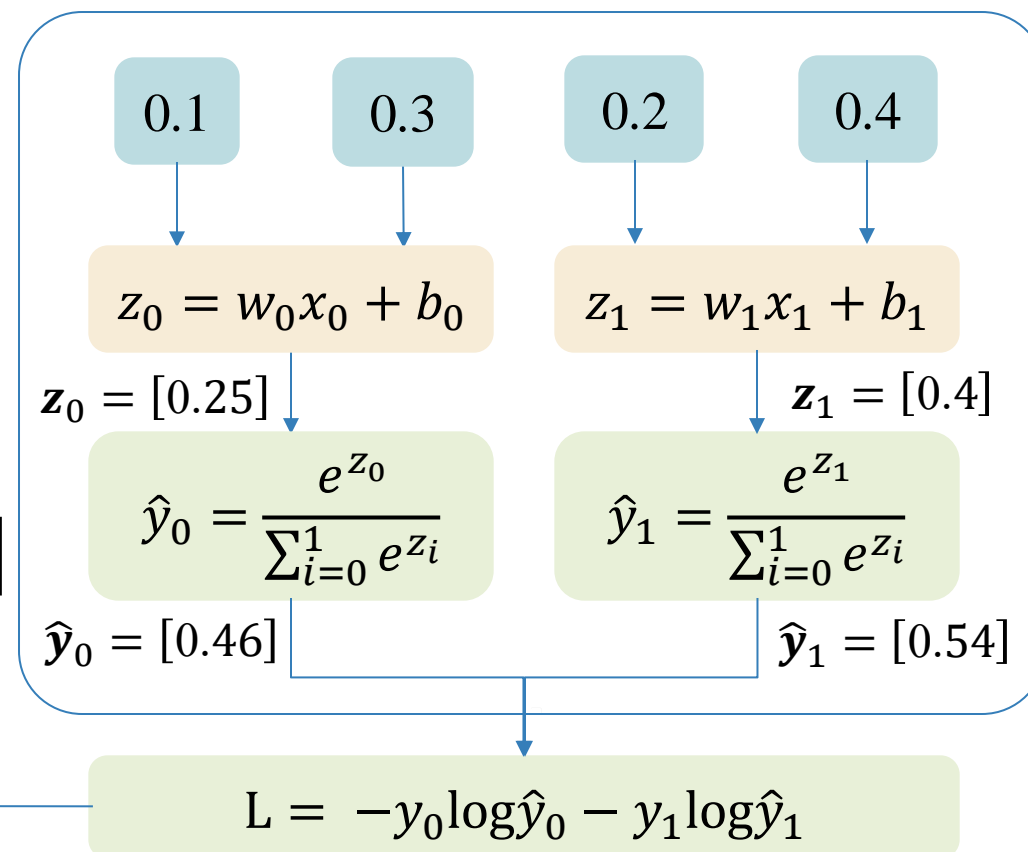
$$\mathbf{y} = 0 \rightarrow \mathbf{y} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$\boldsymbol{\theta} = \begin{bmatrix} 0.1 & 0.2 \\ 0.3 & 0.4 \end{bmatrix}$$

$$\begin{aligned} \boldsymbol{\theta} &= \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} L \\ &= \begin{bmatrix} 0.105 & 0.195 \\ 0.303 & 0.397 \end{bmatrix} \end{aligned}$$

$$\begin{aligned} \nabla_{\boldsymbol{\theta}} L &= \mathbf{x}(\hat{\mathbf{y}} - \mathbf{y})^T \\ &= \begin{bmatrix} -0.54 & 0.54 \\ -0.27 & 0.27 \end{bmatrix} \end{aligned}$$

$$L = [0.77]$$



Softmax Regression

! Softmax Regression using Gradient Descent

1) Pick a sample from training data

2) Compute output \hat{y}

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

$$\mathbf{d} = [1 \dots 1] \mathbf{e}^{\mathbf{z}}$$

\emptyset is
Hadamard
division

$$\hat{\mathbf{y}} = \mathbf{e}^{\mathbf{z}} \emptyset \mathbf{d}$$

3) Compute loss (cross-entropy)

$$L(\boldsymbol{\theta}) = -\mathbf{y}^T \log \hat{\mathbf{y}}$$

4) Compute derivative

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

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} L$$

η is learning rate

$$\mathbf{x} = \begin{bmatrix} 1.0 \\ 1.5 \end{bmatrix}$$

$$\mathbf{y} = [1]$$

One-hot encoding for label

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

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

$$\boldsymbol{\theta} = \begin{bmatrix} b_0 & b_1 \\ w_0 & w_1 \end{bmatrix} = \begin{bmatrix} 0.1 & 0.2 \\ 0.3 & 0.4 \end{bmatrix}$$

$$\eta = 0.1$$

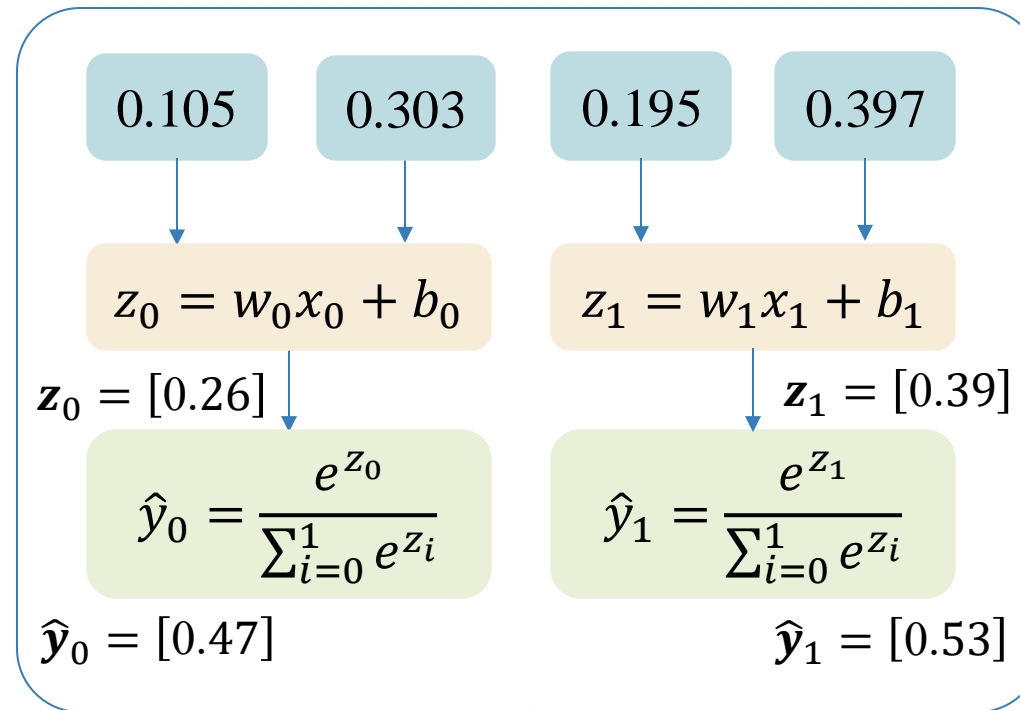
Hours	Pass
0.5	0
1.0	0
1.5	1
2.0	1

Softmax Regression



Softmax Regression using Gradient Descent

Hours	Pass
0.5	?



$y_{pred} = 1$

QUIZ TIME

Softmax Regression

! Softmax Regression using Gradient Descent

1) Pick a sample from training data

2) Compute output \hat{y}

$$z = \theta^T x$$

$$d = [1 \dots 1] e^z$$

\emptyset is Hadamard division

$$\hat{y} = e^z \emptyset d$$

3) Compute loss (cross-entropy)

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

4) Compute derivative

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

5) Update parameters

$$\theta = \theta - \eta \nabla_{\theta} L$$

η is learning rate

```
1 X = np.array([
2     [0.5],
3     [1.0],
4     [1.5],
5     [2.0],
6     [2.5],
7     [3.0],
8     [3.5],
9     [4.0],
10 ])
11 Y = np.array([0, 0, 1, 0, 0, 1, 1, 1])
```

```
1 def convert_one_hot(y, k):
2     one_hot = np.zeros((len(y), k))
3     one_hot[np.arange(len(y)), y] = 1
4     return one_hot
5
6 n_classes = 2
7 Y_onehot = convert_one_hot(Y, n_classes)
```

```
1 X = np.hstack([np.ones((X.shape[0], 1)), X])
```

Softmax Regression

! Softmax Regression using Gradient Descent

1) Pick a sample from training data

2) Compute output \hat{y}

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

$$\mathbf{d} = [1 \dots 1] \mathbf{e}^{\mathbf{z}}$$

$$\hat{\mathbf{y}} = \mathbf{e}^{\mathbf{z}} \oslash \mathbf{d}$$

\oslash is
Hadamard
division

3) Compute loss (cross-entropy)

$$L(\boldsymbol{\theta}) = -\mathbf{y}^T \log \hat{\mathbf{y}}$$

4) Compute derivative

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

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} L$$

η is learning rate

$$\eta = 0.1$$

$$\mathbf{x} = \begin{bmatrix} 1.0 \\ 0.5 \end{bmatrix}$$

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

$$\boldsymbol{\theta} = \begin{bmatrix} 0.1 & 0.2 \\ 0.3 & 0.4 \end{bmatrix}$$

```
1 x = X[0].reshape((2,1))
2 y = Y_onehot[0].reshape((2,1))
3 x, y
```

```
(array([[1. ],
        [0.5]]),
 array([[1.],
        [0.]])
```

```
1 theta = np.array([[0.1, 0.2], [0.3, 0.4]])
2 theta
```

```
array([[0.1, 0.2],
        [0.3, 0.4]])
```

Softmax Regression

! Softmax Regression using Gradient Descent

1) Pick a sample from training data

2) Compute output \hat{y}

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

$$\mathbf{d} = [1 \dots 1] \mathbf{e}^{\mathbf{z}}$$

\emptyset is Hadamard division

$$\hat{\mathbf{y}} = \mathbf{e}^{\mathbf{z}} \emptyset \mathbf{d}$$

3) Compute loss (cross-entropy)

$$L(\boldsymbol{\theta}) = -\mathbf{y}^T \log \hat{\mathbf{y}}$$

4) Compute derivative

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

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} L$$

η is learning rate

$$\eta = 0.1$$

$$\mathbf{x} = \begin{bmatrix} 1.0 \\ 0.5 \end{bmatrix}$$

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

$$\boldsymbol{\theta} = \begin{bmatrix} 0.1 & 0.2 \\ 0.3 & 0.4 \end{bmatrix}$$

```
1 # define softmax function
2 def softmax_function(z):
3     return np.exp(z) / np.sum(np.exp(z))
```

```
1 # compute y_hat
2 def predict(x, theta):
3     z = theta.T.dot(x)
4     y_hat = softmax_function(z)
5     return z, y_hat
6
7 z, y_hat = predict(x, theta)
8 z, y_hat
```

```
(array([[0.25],
        [0.4 ]]),
 array([[0.46257015],
        [0.53742985]]))
```


Softmax Regression

! Softmax Regression using Gradient Descent

1) Pick a sample from training data

2) Compute output \hat{y}

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

$$\mathbf{d} = [1 \dots 1] \mathbf{e}^{\mathbf{z}}$$

\emptyset is Hadamard division

$$\hat{\mathbf{y}} = \mathbf{e}^{\mathbf{z}} \emptyset \mathbf{d}$$

3) Compute loss (cross-entropy)

$$L(\boldsymbol{\theta}) = -\mathbf{y}^T \log \hat{\mathbf{y}}$$

4) Compute derivative

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

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} L$$

η is learning rate

$$\eta = 0.1$$

$$\mathbf{x} = \begin{bmatrix} 1.0 \\ 0.5 \end{bmatrix}$$

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

$$\boldsymbol{\theta} = \begin{bmatrix} 0.1 & 0.2 \\ 0.3 & 0.4 \end{bmatrix}$$

```
1 # compute loss
2 def compute_loss(y_hat, y):
3     loss = -np.log(y.T.dot(y_hat))
4     return loss
5
6 loss = compute_loss(y_hat, y)
7 loss
```

array([[0.77095705]])

Softmax Regression

! Softmax Regression using Gradient Descent

1) Pick a sample from training data

2) Compute output \hat{y}

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

$$\mathbf{d} = [1 \dots 1] \mathbf{e}^{\mathbf{z}}$$

\emptyset is Hadamard division

$$\hat{\mathbf{y}} = \mathbf{e}^{\mathbf{z}} \emptyset \mathbf{d}$$

3) Compute loss (cross-entropy)

$$L(\boldsymbol{\theta}) = -\mathbf{y}^T \log \hat{\mathbf{y}}$$

4) Compute derivative

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

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} L$$

η is learning rate

$$\eta = 0.1$$

$$\mathbf{x} = \begin{bmatrix} 1.0 \\ 0.5 \end{bmatrix}$$

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

$$\boldsymbol{\theta} = \begin{bmatrix} 0.1 & 0.2 \\ 0.3 & 0.4 \end{bmatrix}$$

```
1 # compute gradient
2 def compute_gradient(y_hat, y, x):
3     gradient = x.dot((y_hat - y).T)
4     return gradient
5
6 gradient = compute_gradient(y_hat, y, x)
7 gradient
```

```
array([[ -0.53742985,  0.53742985],
       [ -0.26871492,  0.26871492]])
```

Softmax Regression

! Softmax Regression using Gradient Descent

1) Pick a sample from training data

2) Compute output \hat{y}

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

$$\mathbf{d} = [1 \dots 1] \mathbf{e}^{\mathbf{z}}$$

\emptyset is Hadamard division

$$\hat{\mathbf{y}} = \mathbf{e}^{\mathbf{z}} \emptyset \mathbf{d}$$

3) Compute loss (cross-entropy)

$$L(\boldsymbol{\theta}) = -\mathbf{y}^T \log \hat{\mathbf{y}}$$

4) Compute derivative

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

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} L$$

η is learning rate

$$\eta = 0.1 \quad \mathbf{x} = \begin{bmatrix} 1.0 \\ 0.5 \end{bmatrix} \quad \mathbf{y} = 0 \rightarrow \mathbf{y} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad \boldsymbol{\theta} = \begin{bmatrix} 0.1 & 0.2 \\ 0.3 & 0.4 \end{bmatrix}$$

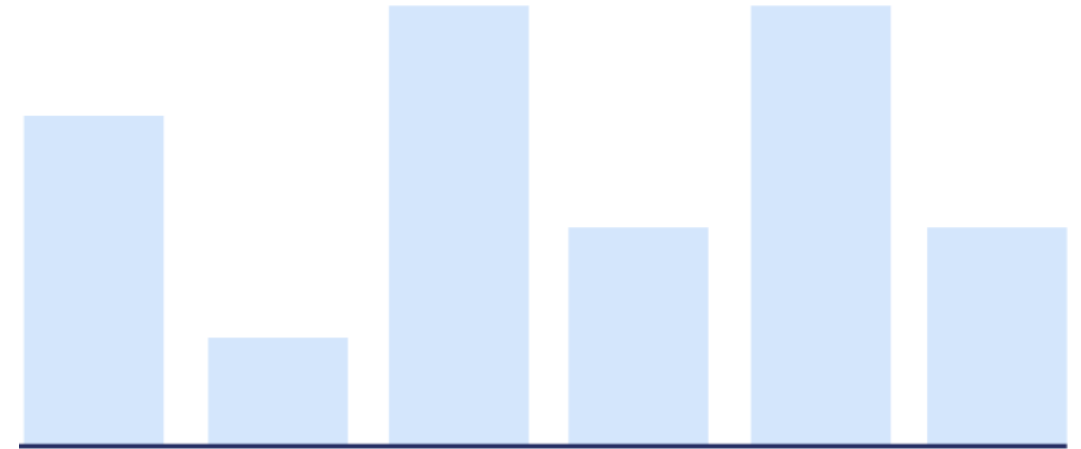
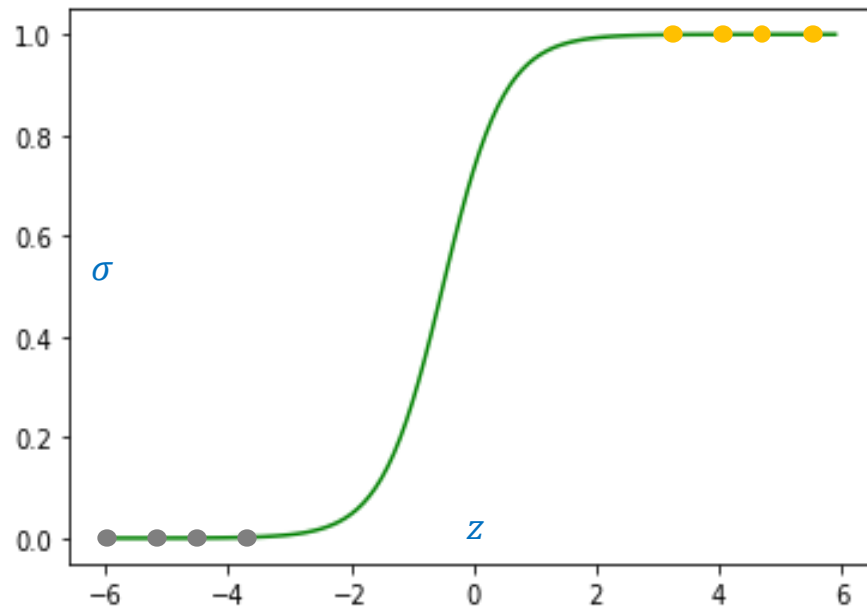
```
1 # update weights
2 learning_rate = 0.01
3
4 def update_weight(theta, gradient, learning_rate):
5     theta -= (learning_rate * gradient)
6     return theta
7
8 theta = update_weight(theta, gradient, learning_rate)
9 theta
```

```
array([[0.1053743 , 0.1946257 ],
       [0.30268715, 0.39731285]])
```

Objectives

Logistic Regression (Review)

- ❖ Logistic Regression
- ❖ Sigmoid Function
- ❖ Gradient Descent



Softmax Regression

- ❖ Softmax Regression
- ❖ Softmax Function
- ❖ One Sample
- ❖ N Sample



AI VIET NAM

@aivietnam.edu.vn

Thanks!

Any questions?