# Lec07: Artificial Neural Network (Part I)

충북대학교
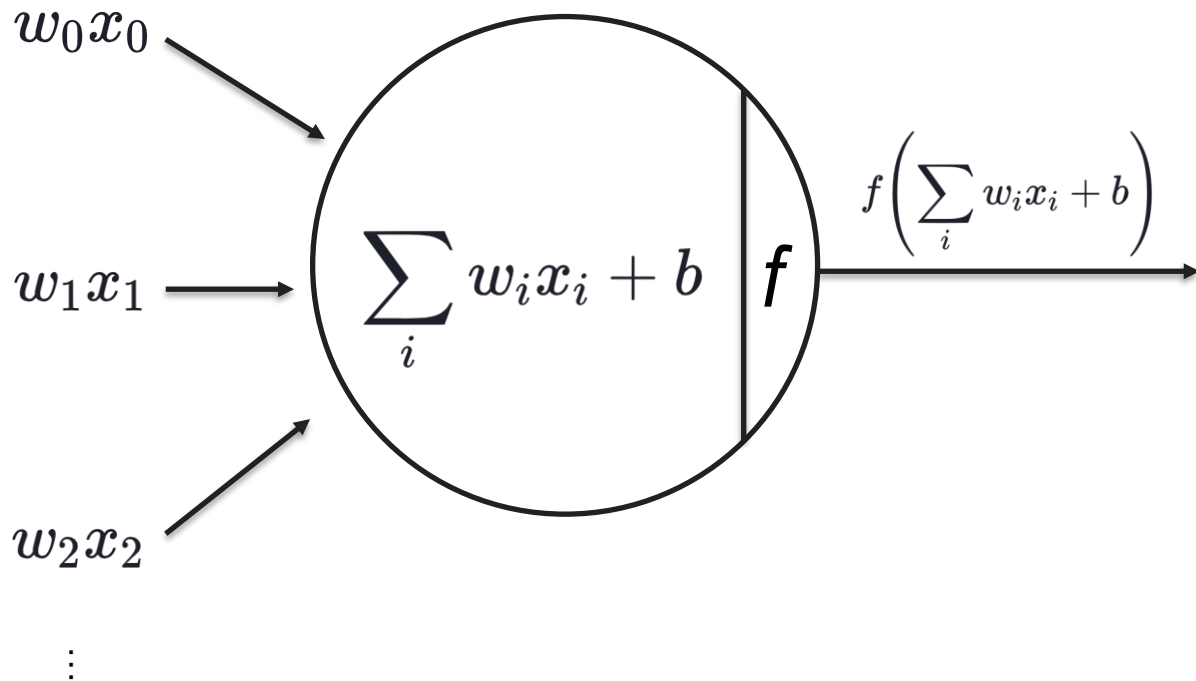
문성태 (지능로봇공학과)

stmoon@cbnu.ac.kr
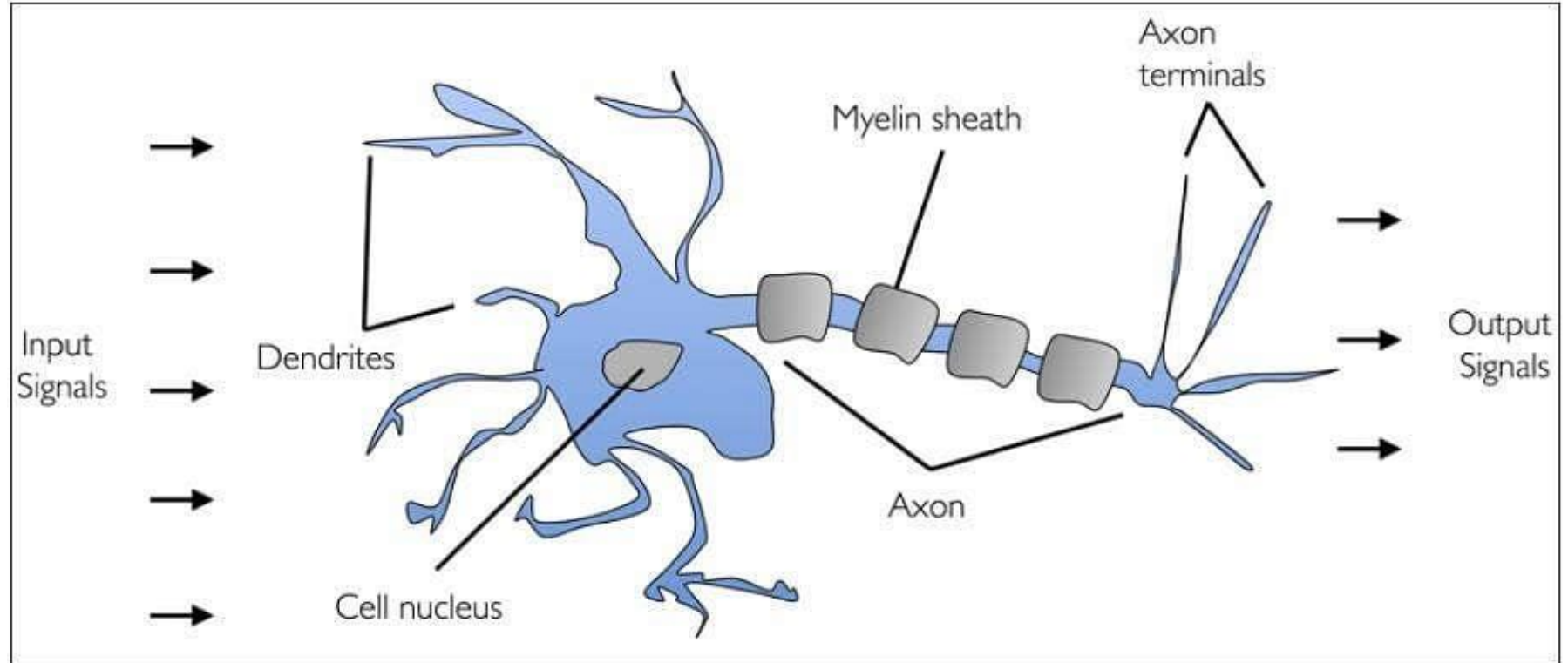
# 01

## Perceptron

# Recap: Logistic Regression



$w_0 x_0$

$w_1 x_1$

$w_2 x_2$

$$\sum_i w_i x_i + b$$

$f$

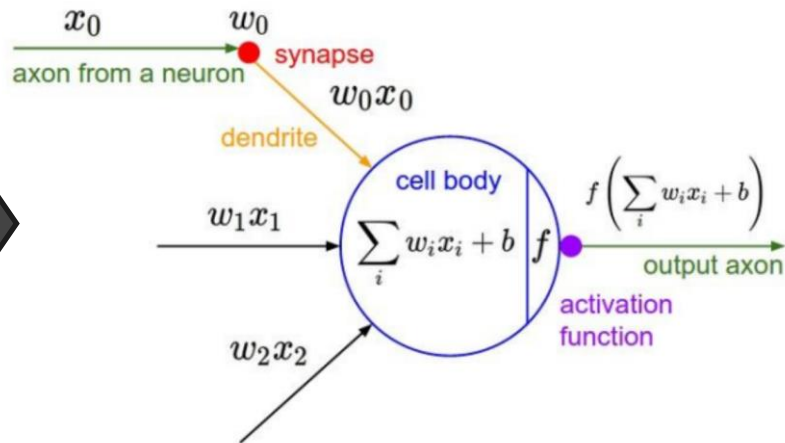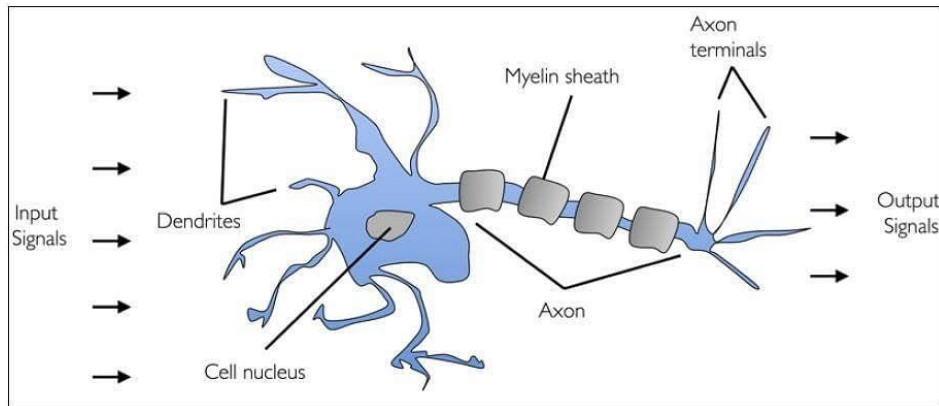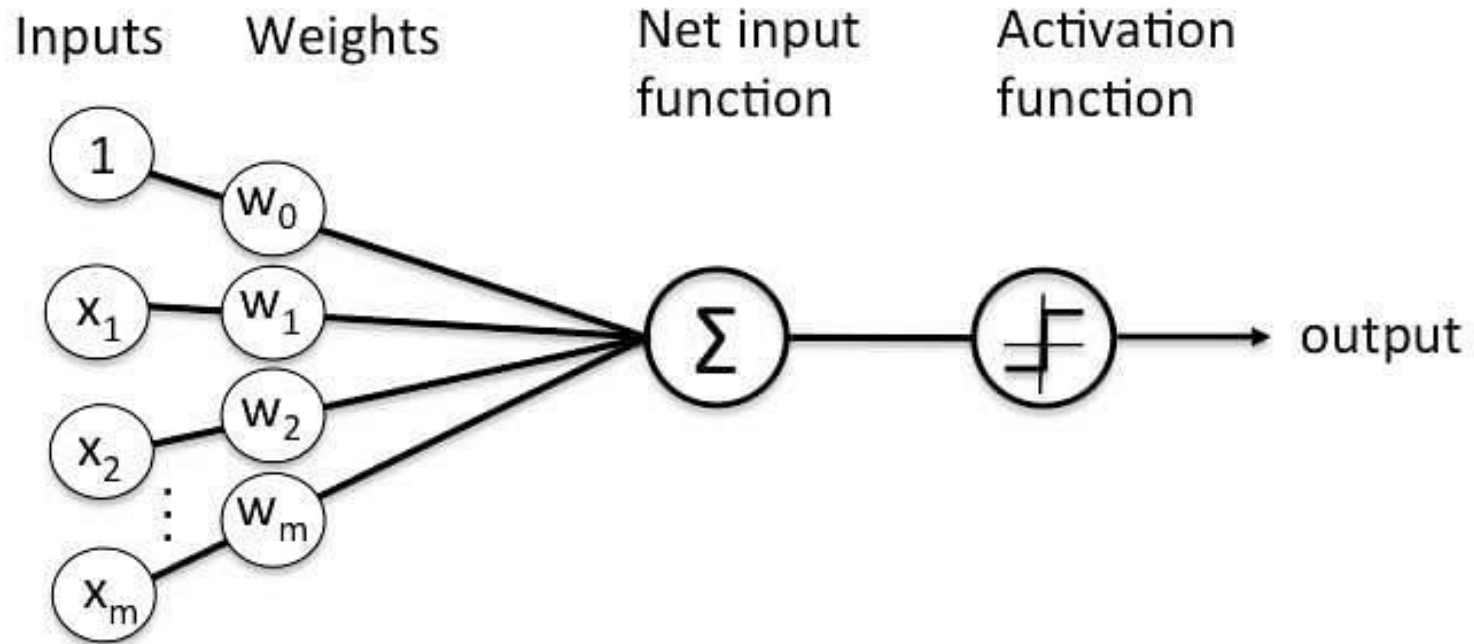$$f\left(\sum_i w_i x_i + b\right)$$

# Biological Neuron

# Perceptron

- 인공 신경망(Artificial Neural Network, ANN)의 구성 요소(unit)로서 다수의 값을 입력 받아 하나의 값으로 출력하는 알고리즘

# Perceptron

- Perceptron was introduced by Frank Rosenblatt in 1957

Inputs    Weights        Net input        Activation
                         function         function

$1$ ── $w_0$

$x_1$ ── $w_1$

$x_2$ ── $w_2$                 $\Sigma$ ──── $\boxed{\neq}$ ──→ output

$x_m$ ── $w_m$

# Perceptron

- Frank Rosenblatt

# Perceptron





FIG. 1 — Organization of a biological brain. (Red areas indicate active cells, responding to the letter X.)

FIG. 2 — Organization of a perceptron.

# Perceptron Hardware

- Frank Rosenblatt with his Mark I perceptron



Figure I ORGANIZATION OF THE MARK I PERCEPTRON

# False Promises

"The Navy revealed the embryo of an electronic computer today that _it expects will be able to walk, talk, see, write, reproduce itself_ an be conscious of its existence … Dr. Frank Rosenblatt, a research psychologist at the Cornell Aeronautical Laboratory, Buffalo, said Perceptrons might be fired to the planets as mechanical space explorers" The New York Times   July 08, 1958

# AND/OR Problem

- perceptron can separate its input space with a <span style="color:red">hyperplane</span>

# XOR Problem

- Linearly separable?

| XOR | | |
|-----|-----|-----|
| $I_1$ | $I_2$ | out |
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

$I_1$

(1, 0)

(1, 1)

?

(0, 0)

(0, 1)

$I_2$

# Perceptrons (1969)

- Perceptrons (1969) by Marvin Minsky, founder of the MIT AI Lab

## "No one on earth had found a viable way to train"

We need to use MLP, multilayer perceptrons (multilayer neural nets)

No one on earth had found a viable way to train MLPs good enough to learn such simple functions

# AI Winter I

AI HAS A LONG HISTORY OF BEING "THE NEXT BIG THING"...



**Timeline of AI Development**

- **1950s-1960s**: First AI boom - the age of reasoning, prototype AI developed
- **1970s**: AI winter I
- **1980s-1990s**: Second AI boom: the age of Knowledge representation (appearance of expert systems capable of reproducing human decision-making)
- **1990s**: AI winter II
- **1997**: Deep Blue beats Gary Kasparov
- **2006**: University of Toronto develops Deep Learning
- **2011**: IBM's Watson won Jeopardy
- **2016**: Go software based on Deep Learning beats world's champions

# 02

# Multi-Layer Perceptron

# XOR Problem

| XOR | | |
|-----|-----|-----|
| $I_1$ | $I_2$ | out |
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |



$x_1 \rightarrow$  [ $\Sigma$ / sig ] $\rightarrow y_1$

$x_2 \rightarrow$

# Quiz

- 성냥개비 6개로 정삼각형 4개만 만드세요

# Quiz Solution

# XOR Solution

# XOR Problem Solution

# XOR Problem Solution



$w = \begin{bmatrix} 5 \\ 5 \end{bmatrix}, b = -8$

Perceptron
Wx+b

$\Sigma$ sig

$y_1$

$w = \begin{bmatrix} -15 \\ -15 \end{bmatrix}, b = 6$

$x_1$

$x_2$

Perceptron
Wx+b

$\Sigma$ sig

$\hat{y}$

Perceptron
Wx+b

$\Sigma$ sig

$y_2$

$w = \begin{bmatrix} -7 \\ -7 \end{bmatrix}, b = 3$

How can we learn W and b from training data?

# 03

**Backpropagation**

# Backpropagation

- 1974, 1982 by Paul Werbos, 1986 by Hinton
  - Paul Werbos, based on his 1974 Ph.D. thesis, publicly proposes the use of Backpropagation for propagating errors during the training of Neural Networks

Before learning backpropagation…

# Basic derivative

$$\frac{df(x)}{dx} = \lim_{h \to 0} \frac{f(x+h) - f(x)}{h}$$

- $f(x) = 3$

- $f(x) = 2x$

- $f(x) = x + 3$

$\cancel{\emptyset}$

$$\frac{2 + (2x)}{dx} = 2$$

$$= 1$$

# Basic derivative (Chain Rule)

$$F(x) = f(g(x))$$
$$F'(x) = f'(g(x)) \cdot g'(x)$$

하나 스타

# Basic derivative (Sigmoid)

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

# Computing gradients of weights in neural network

- Simplest example: two-layer neural network with one hidden node

$$\hat{y} = f(x; \mathbf{W})$$

$$\{w_1, w_2\}$$

# Computing gradients of weights in neural network

- Simplest example: two-layer neural network with one hidden node

$$\hat{y} = f(x; \mathbf{W})$$



$$\frac{\partial \mathcal{J}(\mathbf{W})}{\partial w_1} =$$

?

# Computing gradients of weights in neural network

- Simplest example: two-layer neural network with one hidden node

$$\hat{y} = f(x; \mathbf{W})$$



Chain rule: propagating the gradient across the layers

$$\frac{\partial \mathcal{J}(\mathbf{W})}{\partial w_1} = \frac{\partial \mathcal{J}(\mathbf{W})}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial w_1}$$

# Computing gradients of weights in neural network

- Simplest example: two-layer neural network with one hidden node

$$\hat{y} = f(x; \mathbf{W})$$



Chain rule: propagating the gradient across the layers

$$\frac{\partial \mathcal{J}(\mathbf{W})}{\partial w_1} = \frac{\partial \mathcal{J}(\mathbf{W})}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial w_1}$$

# Computing gradients of weights in neural network

- Simplest example: two-layer neural network with one hidden node

$$\hat{y} = f(x; \mathbf{W})$$



Chain rule: propagating the gradient across the layers

$$\frac{\partial \mathcal{J}(\mathbf{W})}{\partial w_1} = \frac{\partial \mathcal{J}(\mathbf{W})}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial w_1}$$

# Computing gradients of weights in neural network
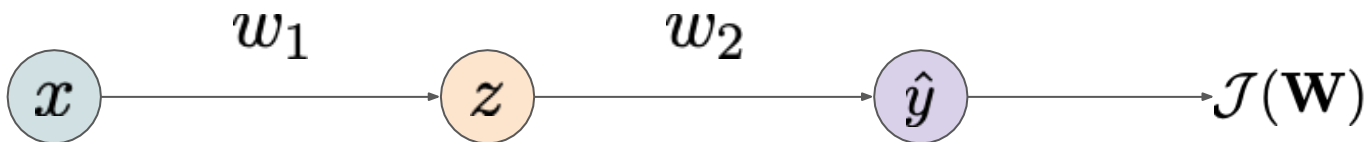
- Simplest example: two-layer neural network with one hidden node

$$\hat{y} = f(x; \mathbf{W})$$



Chain rule: propagating the gradient across the layers

$$\frac{\partial \mathcal{J}(\mathbf{W})}{\partial w_1} = \frac{\partial \mathcal{J}(\mathbf{W})}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial w_1}$$

# 03

**Backpropagation for XOR**

# XOR neural network



$$W^{[1]} = \begin{bmatrix} w_{00}^{[1]} & w_{01}^{[1]} \\ w_{10}^{[1]} & w_{11}^{[1]} \end{bmatrix} \quad B^{[1]} = \begin{bmatrix} b_0^{[1]} \\ b_1^{[1]} \end{bmatrix} \quad W^{[2]} = \begin{bmatrix} w_{00}^{[2]} \\ w_{10}^{[2]} \end{bmatrix} \quad B^{[2]} = \begin{bmatrix} b_0^{[2]} \end{bmatrix}$$

# Forward propagation



$$z_0^{[1]} = w_{00}^{[1]} x_o + w_{10}^{[1]} x_1 + b_o^{[1])}$$
$$z_1^{[1]} = w_{01}^{[1]} x_o + w_{11}^{[1]} x_1 + b_1^{[1]}$$

➡ $Z^{[1]} = W^{[1]^T} X + B^{[1]}$

$$W^{(1)} = \begin{bmatrix} w_{00}^{(1)} & w_{01}^{(1)} \\ w_{10}^{(1)} & w_{11}^{(1)} \end{bmatrix} \quad B^{(1)} = \begin{bmatrix} b_0^{(1)} \\ b_1^{(1)} \end{bmatrix}$$

$$h_0 = \sigma(z_0^{(1)})$$
$$h_1 = \sigma(z_1^{(1)})$$

➡ $H^{(1)} = \sigma(Z^{[1]})$

# Forward propagation



$$z_0^{[2]} = w_0^{[2]} h_o + w_1^{[2]} h_1 + b_o^{[2]} \quad \blacktriangleright \quad z^{[2]} = W^{[2]^T} H + b_0^{[2]} \qquad W^{(2)} = \begin{bmatrix} w_{00}^{(2)} \\ w_{10}^{(2)} \end{bmatrix}$$

$$\hat{y} = \sigma(z^{[2]})$$

# Forward propagation



$$\mathcal{L} = -\frac{1}{m} \sum_i^m \{ y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i) \}$$

# Backward Propagation



$$\frac{\partial \mathcal{L}}{\partial W^{(1)}}, \ \frac{\partial \mathcal{L}}{\partial B^{(1)}}, \ \frac{\partial \mathcal{L}}{\partial W^{(2)}}, \ \frac{\partial \mathcal{L}}{\partial B^{(2)}}$$

# Backward Propagation



$$\frac{\partial \mathcal{L}}{\partial w_{00}^{[2]}} = \frac{\partial \mathcal{L}}{\partial z_0^{[2]}} \frac{\partial z_0^{[2]}}{\partial w_{00}^{[2]}} = \frac{\partial \mathcal{L}}{\partial z_0^{[2]}} h_0$$

$$\frac{\partial \mathcal{L}}{\partial w_{10}^{[2]}} = \frac{\partial \mathcal{L}}{\partial z_0^{[2]}} \frac{\partial z_0^{[2]}}{\partial w_{10}^{[2]}} = \frac{\partial \mathcal{L}}{\partial z_0^{[2]}} h_1$$

$$\frac{\partial \mathcal{L}}{\partial b^{[2]}} = \frac{\partial \mathcal{L}}{\partial z_0^{[2]}}$$

$$\frac{\partial \mathcal{L}}{\partial z_0^{[2]}} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_0^{[2]}}$$

$$= \frac{\partial \mathcal{L}}{\partial \hat{y}} \hat{y}(1 - \hat{y})$$

$$= \left( -\frac{y}{\hat{y}} + \frac{1-y}{1-\hat{y}} \right) \hat{y}(1 - \hat{y})$$

$$= \hat{y} - y$$

$$\frac{\partial \mathcal{L}}{\partial \hat{y}} = \frac{1}{\partial \hat{y}} (-y \log \hat{y} + (1-y) \log(1 - \hat{y})$$

$$= -\frac{y}{\hat{y}} + \frac{1-y}{1-\hat{y}}$$

# Backward **Propagation**



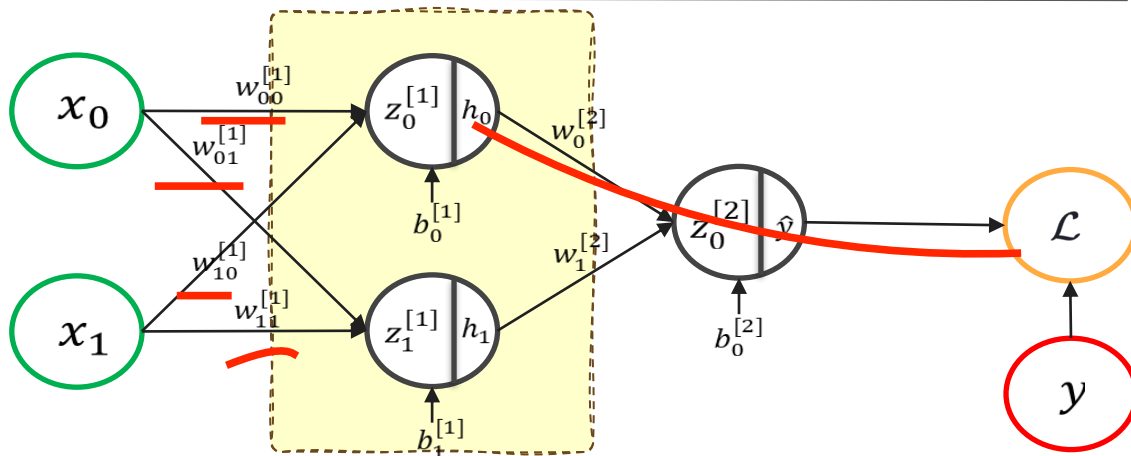$$\frac{\partial \mathcal{L}}{\partial W^{[1]}} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial w_{00}^{[1]}} & \frac{\partial \mathcal{L}}{\partial w_{01}^{[1]}} \\ \frac{\partial \mathcal{L}}{\partial w_{10}^{[1]}} & \frac{\partial \mathcal{L}}{\partial w_{11}^{[1]}} \end{bmatrix} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial z_0^{[1]}} x_0 & \frac{\partial \mathcal{L}}{\partial z_1^{[1]}} x_0 \\ \frac{\partial \mathcal{L}}{\partial z_0^{[1]}} x_1 & \frac{\partial \mathcal{L}}{\partial z_1^{[1]}} x_1 \end{bmatrix}$$

$$\frac{\partial \mathcal{L}}{\partial B^{[1]}} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial b_0^{[1]}} \\ \frac{\partial \mathcal{L}}{\partial b_1^{[1]}} \end{bmatrix}$$

$$\frac{\partial \mathcal{L}}{\partial w_{00}^{[1]}} = \frac{\partial \mathcal{L}}{\partial z_0^{[1]}} \frac{\partial z_0^{[1]}}{\partial w_{00}^{[1]}} = \frac{\partial \mathcal{L}}{\partial z_0^{[1]}} x_0$$

$$\frac{\partial \mathcal{L}}{\partial w_{01}^{[1]}} = \frac{\partial \mathcal{L}}{\partial z_0^{[1]}} \frac{\partial z_0^{[1]}}{\partial w_{01}^{[1]}} = \frac{\partial \mathcal{L}}{\partial z_1^{[1]}} x_0$$

$$\frac{\partial \mathcal{L}}{\partial z_0^{[1]}} = \frac{\partial \mathcal{L}}{\partial h_0} \frac{\partial h_0}{\partial z_0^{[1]}}$$

$$= \frac{\partial \mathcal{L}}{\partial h_0} h_0 (1 - h_0)$$

$$\frac{\partial \mathcal{L}}{\partial h_0} = \frac{\partial \mathcal{L}}{\partial z_0^{[2]}} \frac{\partial z_0^{[2]}}{\partial h_0} = \frac{\partial \mathcal{L}}{\partial z_0^{[2]}} w_{00}^{[2]}$$

$$\frac{\partial \mathcal{L}}{\partial b_0^{[1]}} = \frac{\partial \mathcal{L}}{\partial z_0^{[1]}} \frac{\partial z_0^{[1]}}{\partial b_0^{[1]}} = \frac{\partial \mathcal{L}}{\partial z_0^{[1]}}$$

$$\frac{\partial \mathcal{L}}{\partial w_{10}^{[1]}} = \frac{\partial \mathcal{L}}{\partial z_0^{[1]}} \frac{\partial z_0^{[1]}}{\partial w_{10}^{[1]}} = \frac{\partial \mathcal{L}}{\partial z_0^{[1]}} x_1$$

$$\frac{\partial \mathcal{L}}{\partial z_1^{[1]}} = \frac{\partial \mathcal{L}}{\partial h_1} \frac{\partial h_1}{\partial z_1^{[1]}}$$

$$\frac{\partial \mathcal{L}}{\partial h_1} = \frac{\partial \mathcal{L}}{\partial z_0^{[2]}} \frac{\partial z_0^{[2]}}{\partial h_1} = \frac{\partial \mathcal{L}}{\partial z_0^{[2]}} w_{10}^{[2]}$$

$$\frac{\partial \mathcal{L}}{\partial b_1^{[1]}} = \frac{\partial \mathcal{L}}{\partial z_1^{[1]}} \frac{\partial z_1^{[1]}}{\partial b_1^{[1]}} = \frac{\partial \mathcal{L}}{\partial z_1^{[1]}}$$

$$\frac{\partial \mathcal{L}}{\partial w_{11}^{[1]}} = \frac{\partial \mathcal{L}}{\partial z_1^{[1]}} \frac{\partial z_1^{[1]}}{\partial w_{11}^{[1]}} = \frac{\partial \mathcal{L}}{\partial z_1^{[1]}} x_1$$

$$= \frac{\partial \mathcal{L}}{\partial h_1} h_1 (1 - h_1)$$