

Backpropagation

기계학습/딥러닝

임 경 태

Week	Chapter	Contents
1	1, 2장	강의 소개, 파이썬 복습
2	1, 3장	파이썬 복습, Numpy, Pandas
3	1, 4장	딥러닝을 위한 미분
4	5장	회귀
5	5장	분류
6	6장	XOR문제
7	7장	딥러닝
8	1~7장	중간고사
9	8장	MNIST 필기체 구현 및 팀 프로젝트 소개
10	9장	오차역전파
11	11장	Jetbot 자율주행 (Collision Avoidance, Transfer Learning)
12	12장	특강 (인공지능 활용 연구) + 팀프로젝트 자율 실습
13	10장	Jetbot 자율주행 (Road Following)
14	11장	합성곱 신경망(CNN), 순환 신경망(RNN) + 팀프로젝트 자율 실습
15	8~12장	기말고사 (or 프로젝트 발표)

CONTENTS

- 1 오차역전파 개념
- 2 출력층에서의 오차역전파
- 3 은닉층에서의 오차역전파
- 4 오차역전파 이용 MNIST 검증



목적 : 오차역전파를 이용하는 이유와 오차역전파에 대한 이해



목표 : 각 층에서의 오차역전파 진행 과정 이해



내용 : 수치 미분 문제점, 오차역전파 원리, MNIST 검증

CONTENTS

- 1 오차역전파 개념**
- 2 출력층에서의 오차역전파**
- 3 은닉층에서의 오차역전파**
- 4 오차역전파 이용 MNIST 검증**

수치 미분 문제점

수치 미분을 사용한 가중치 업데이트 과정

1. Training Data 입력

2. Feed Forward

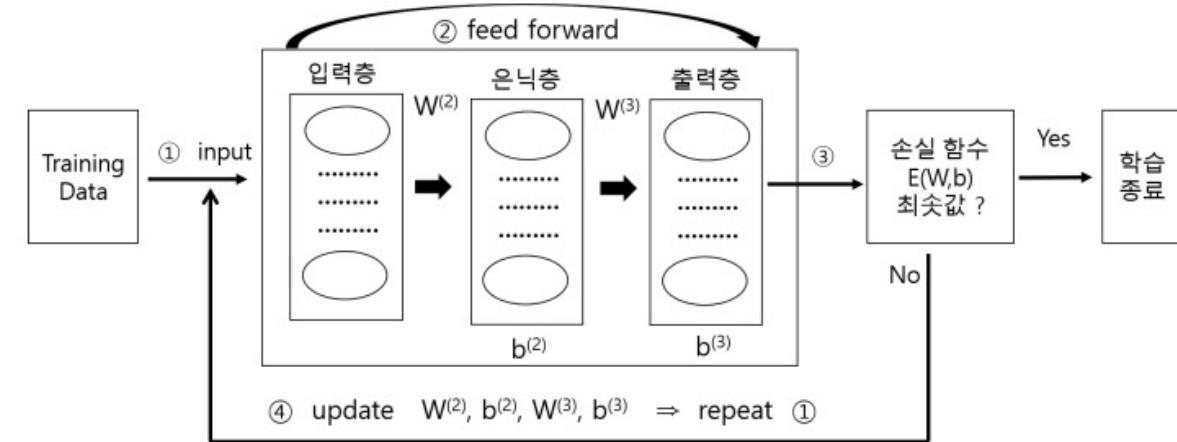
3. 손실함수 계산

4. 수치 미분을 통하여 가중치 업데이트

5. 2~4 반복

- 많은 시간 소요
- 수치 미분에서의 오차 포함

→ Chain Rule을 사용하여 분해, 오차역전파



$$W^{(2)} = W^{(2)} - \alpha \frac{\partial E(W,b)}{\partial W^{(2)}}$$

$$b^{(2)} = b^{(2)} - \alpha \frac{\partial E(W,b)}{\partial b^{(2)}}$$

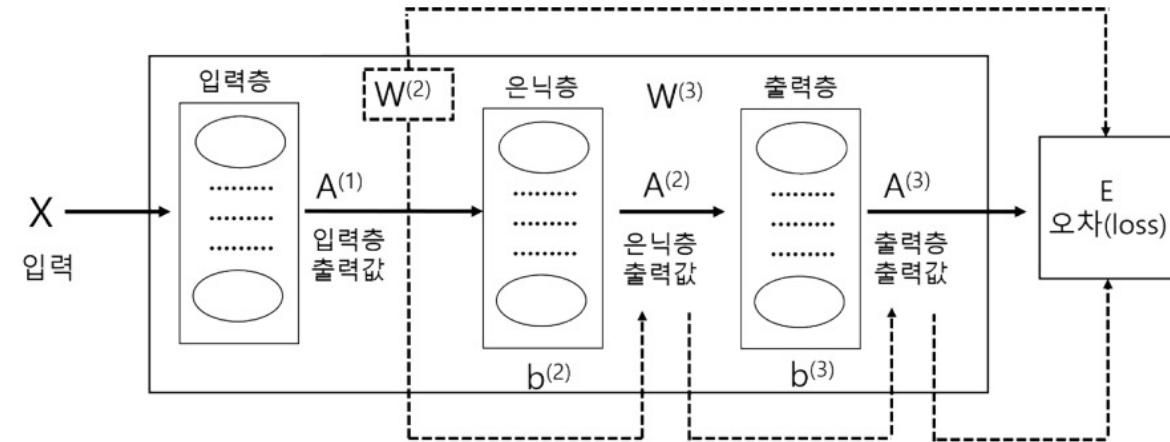
$$W^{(3)} = W^{(3)} - \alpha \frac{\partial E(W,b)}{\partial W^{(3)}}$$

$$b^{(3)} = b^{(3)} - \alpha \frac{\partial E(W,b)}{\partial b^{(3)}}$$

오차역전파 개념 및 동작 원리

Chain Rule을 적용하여 수식 분해

$A^{(1)}$: 입력층 출력 값
 $A^{(2)}$: 은닉층 출력 값
 $A^{(3)}$: 출력층 출력 값



$\frac{\partial E}{\partial W^{(2)}}$ 를 chain Rule을 사용하여 $A^{(2)}, A^{(3)}$ 를 포함하도록 분해하면?

$$\rightarrow \frac{\partial E}{\partial W^{(2)}} = \frac{\partial E}{\partial A^{(3)}} \cdot \frac{\partial A^{(3)}}{\partial W^{(2)}} = \frac{\partial E}{\partial A^{(3)}} \cdot \frac{\partial A^{(3)}}{\partial A^{(2)}} \cdot \frac{\partial A^{(2)}}{\partial W^{(2)}}$$

$$\frac{\partial A^{(2)}}{\partial W^{(2)}} \quad \frac{\partial A^{(3)}}{\partial A^{(2)}} \quad \frac{\partial A^{(3)}}{\partial W^{(2)}} \quad \frac{\partial E}{\partial A^{(3)}}$$

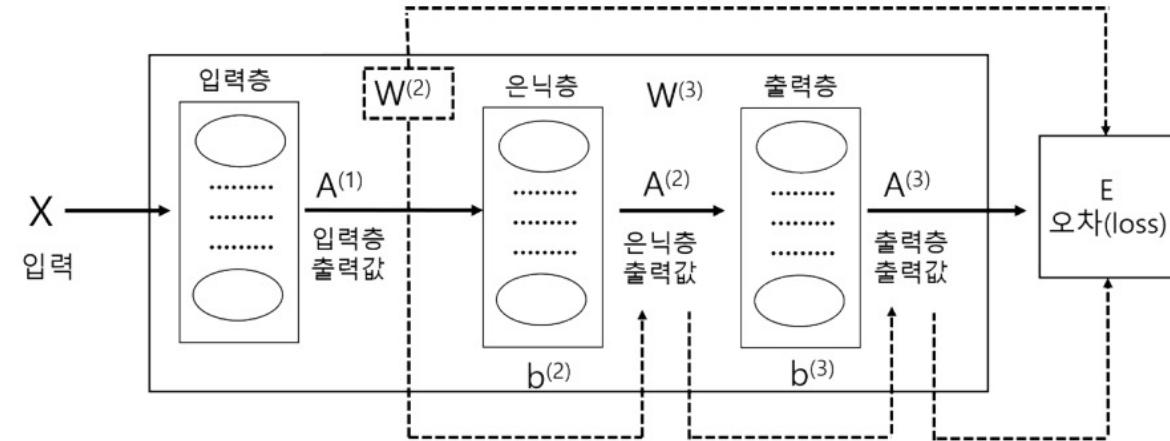
Chain Rule을 사용하여 $A^{(2)}, A^{(3)}$ 를 포함하도록 가중치 업데이트 식을 분해하면?

$$\rightarrow W^{(2)} := W^{(2)} - \alpha \frac{\partial E}{\partial W^{(2)}} = W^{(2)} - \alpha \frac{\partial E}{\partial A^{(3)}} \cdot \frac{\partial A^{(3)}}{\partial A^{(2)}} \cdot \frac{\partial A^{(2)}}{\partial W^{(2)}}$$

오차역전파 개념 및 동작 원리

Chain Rule을 적용하여 수식 분해

$A^{(1)}$: 입력층 출력 값
 $A^{(2)}$: 은닉층 출력 값
 $A^{(3)}$: 출력층 출력 값



$\frac{\partial E}{\partial W^{(2)}}$ 은 chain Rule에 의해 국소적인 미분의 곱셈으로 나타낼 수 있다!

✓ 국소적인 미분?

- 이전 layer의 출력 값 변화에 따른 다음 layer의 출력 값 변화율을 나타낼 수 있다.
- 계산 전체에서 어떤 일이 벌어지는 상관 없이 자신과 직접 관계된 정보만으로 결과를 출력한다.
- $\frac{\partial E}{\partial A^{(3)}}, \frac{\partial A^{(3)}}{\partial A^{(2)}}, \frac{\partial A^{(2)}}{\partial W^{(2)}}$ 와 같은 계산 결과를 모두 보관할 수 있다.

$\frac{\partial A^{(2)}}{\partial W^{(2)}}$: $W^{(2)}$ 에 대한 $A^{(2)}$ 의 변화율

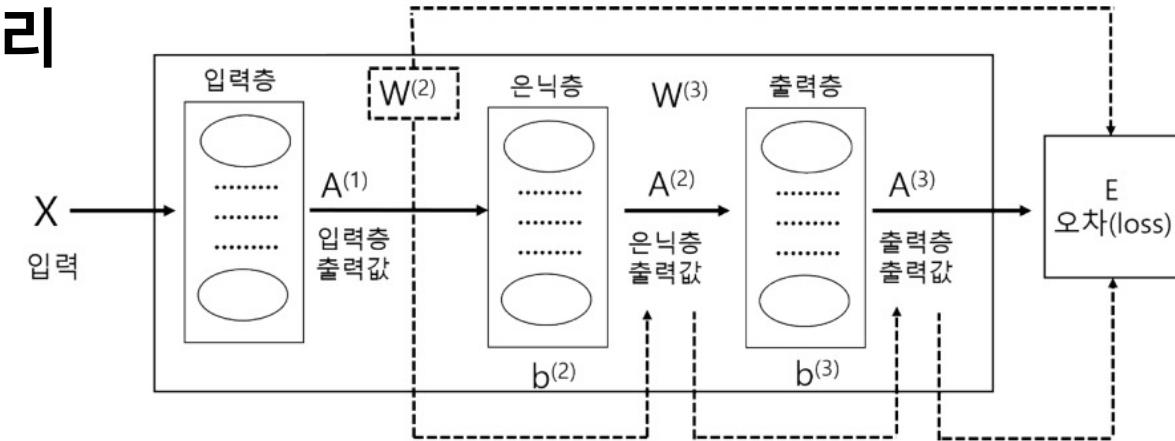
$\frac{\partial A^{(3)}}{\partial A^{(2)}}$: $A^{(2)}$ 에 대한 $A^{(3)}$ 의 변화율

$\frac{\partial E}{\partial A^{(3)}}$: $A^{(3)}$ 에 대한 E 의 변화율

오차역전파 개념 및 동작 원리

Chain Rule을 적용하여 수식 분해 - 정리

$A^{(1)}$: 입력층 출력 값
 $A^{(2)}$: 은닉층 출력 값
 $A^{(3)}$: 출력층 출력 값



- 가중치, 편향의 변화에 따른 오차가 얼마나 변하는지 나타내는 식을 Chain Rule을 통하여 다음과 같이 국소 미분으로 분리할 수 있다.

$$\frac{\partial E}{\partial W^{(2)}} = \frac{\partial E}{\partial A^{(3)}} \cdot \frac{\partial A^{(3)}}{\partial A^{(2)}} \cdot \frac{\partial A^{(2)}}{\partial W^{(2)}}$$

$$\frac{\partial E}{\partial W^{(3)}} = \frac{\partial E}{\partial A^{(3)}} \cdot \frac{\partial A^{(3)}}{\partial W^{(3)}}$$

$$\frac{\partial E}{\partial b^{(2)}} = \frac{\partial E}{\partial A^{(3)}} \cdot \frac{\partial A^{(3)}}{\partial A^{(2)}} \cdot \frac{\partial A^{(2)}}{\partial b^{(2)}}$$

$$\frac{\partial E}{\partial b^{(3)}} = \frac{\partial E}{\partial A^{(3)}} \cdot \frac{\partial A^{(3)}}{\partial b^{(3)}}$$

- 이 국소 미분으로 분리된 식을 이용하여 최종적으로 수치 미분이 아닌 곱셈 형태의 산술 식으로 계산하여 오차역전파에 적용한다.

$$\frac{\partial A^{(2)}}{\partial W^{(2)}} \quad \frac{\partial A^{(3)}}{\partial A^{(2)}} \quad \frac{\partial E}{\partial A^{(3)}}$$

$\frac{\partial A^{(2)}}{\partial W^{(2)}}$: $W^{(2)}$ 에 대한 $A^{(2)}$ 의 변화율

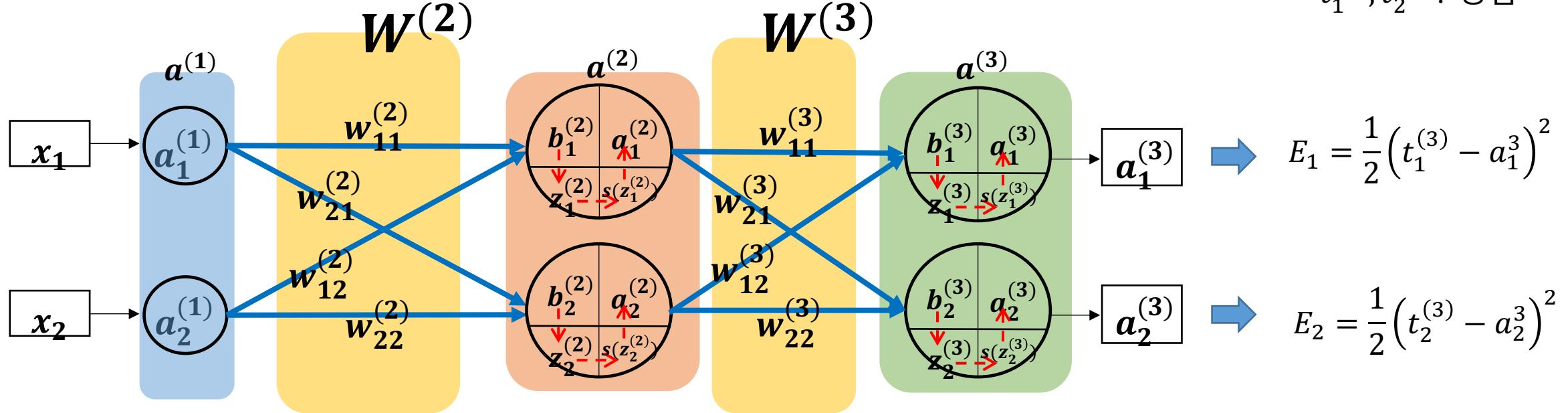
$\frac{\partial A^{(3)}}{\partial A^{(2)}}$: $A^{(2)}$ 에 대한 $A^{(3)}$ 의 변화율

$\frac{\partial E}{\partial A^{(3)}}$: $A^{(3)}$ 에 대한 E 의 변화율

오차역전파 개념 및 동작 원리

각 층에서의 가중치(W), 편향(b)과 오차(E)

$t_1^{(3)}, t_2^{(3)}$: 정답



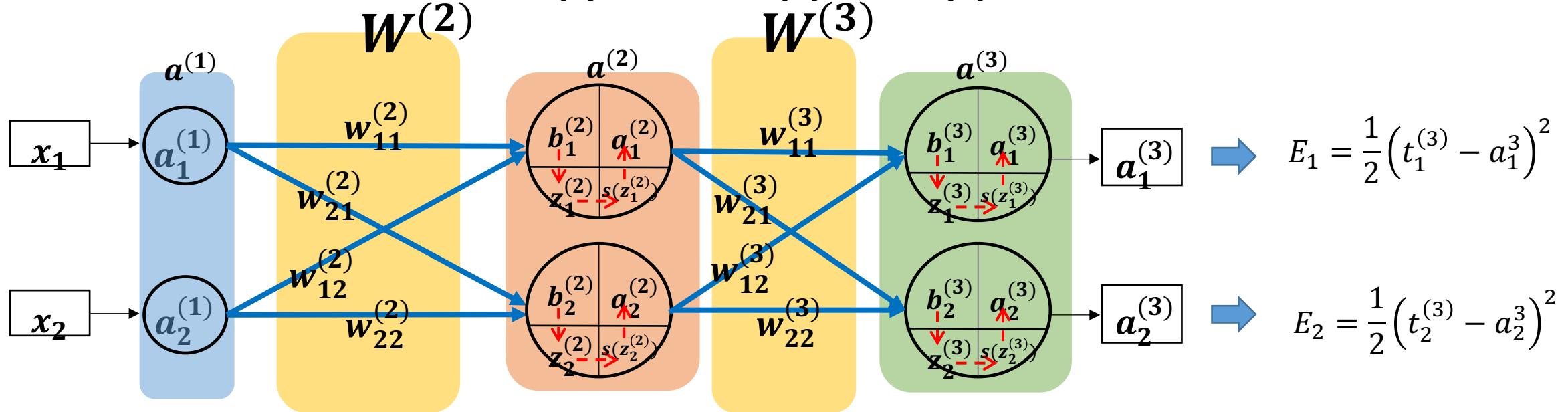
$$W^{(2)} = \begin{bmatrix} w_{11}^{(2)} & w_{21}^{(2)} \\ w_{12}^{(2)} & w_{22}^{(2)} \end{bmatrix} \quad b^{(2)} = [b_1^{(2)} \quad b_2^{(2)}]$$

$$W^{(3)} = \begin{bmatrix} w_{11}^{(3)} & w_{21}^{(3)} \\ w_{12}^{(3)} & w_{22}^{(3)} \end{bmatrix} \quad b^{(3)} = [b_1^{(3)} \quad b_2^{(3)}]$$

$$E = \frac{1}{2} \sum_{i=1}^2 (t_i^{(3)} - a_i^{(3)})^2 = \frac{1}{2} \left\{ (t_1^{(3)} - a_1^{(3)})^2 + (t_2^{(3)} - a_2^{(3)})^2 \right\} = E_1 + E_2$$

오차역전파 개념 및 동작 원리

☞ 각 층에서의 선형 회귀 값(z), 출력 값(a), 오차(E)



입력 층

$$\begin{aligned} a_1^{(1)} &= x_1 \\ a_2^{(1)} &= x_2 \end{aligned}$$

은닉 층

$$\begin{aligned} z_1^{(2)} &= a_1^{(1)}w_{11}^{(2)} + a_2^{(1)}w_{12}^{(2)} + b_1^{(2)} \\ z_2^{(2)} &= a_2^{(1)}w_{21}^{(2)} + a_2^{(1)}w_{22}^{(2)} + b_2^{(2)} \\ a_1^{(2)} &= \text{sigmoid}(z_1^{(2)}) \\ a_2^{(2)} &= \text{sigmoid}(z_2^{(2)}) \end{aligned}$$

출력 층

$$\begin{aligned} z_1^{(3)} &= a_1^{(2)}w_{11}^{(3)} + a_2^{(2)}w_{12}^{(3)} + b_1^{(3)} \\ z_2^{(3)} &= a_2^{(2)}w_{21}^{(3)} + a_2^{(2)}w_{22}^{(3)} + b_2^{(3)} \\ a_1^{(3)} &= \text{sigmoid}(z_1^{(3)}) \\ a_2^{(3)} &= \text{sigmoid}(z_2^{(3)}) \end{aligned}$$

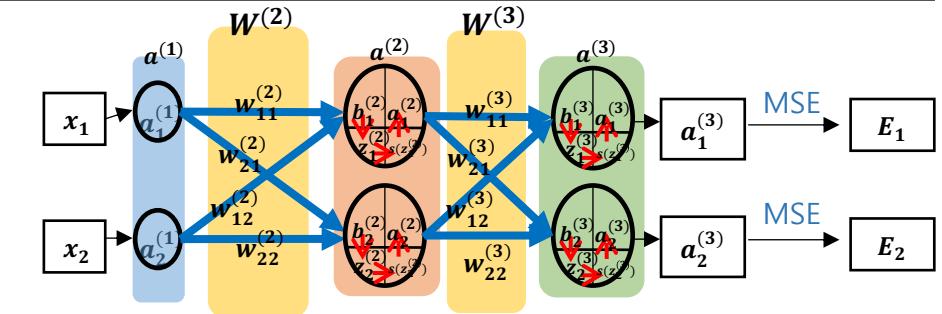
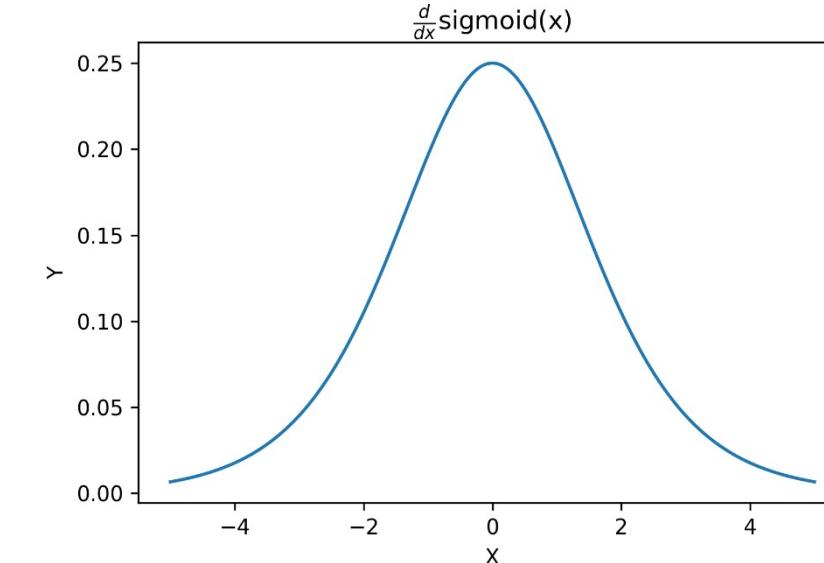
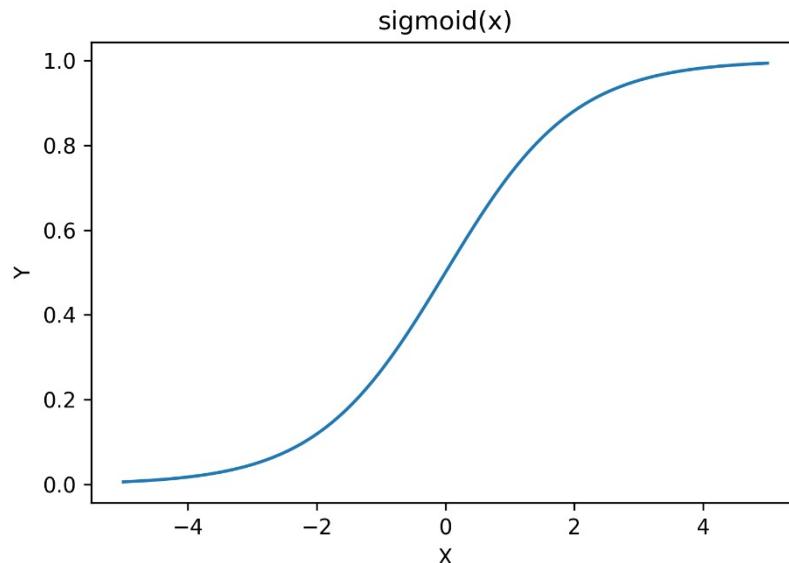
오차역전파 개념 및 동작 원리

Sigmoid 함수 미분

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

$$\frac{d}{dx} S(x) = \frac{d}{dx} \left(\frac{1}{1+e^{-x}} \right) = \frac{e^{-x}}{(1+e^{-x})^2} = \frac{1}{1+e^{-x}} \cdot \frac{e^{-x}}{1+e^{-x}} = \frac{1}{1+e^{-x}} \left(1 - \frac{1}{1+e^{-x}} \right)$$

$$= S(x)(1 - S(x))$$



CONTENTS

- 1 오차역전파 개념
- 2 출력층에서의 오차역전파
- 3 은닉층에서의 오차역전파
- 4 오차역전파 이용 MNIST 검증

출력 층에서의 오차역전파

$w_{11}^{(3)}$ 에 대한 E 의 변화율 - $\frac{\partial E}{\partial w_{11}^{(3)}}$

$$W^{(3)} := W^{(3)} - \alpha \frac{\partial E}{\partial W^{(3)}}$$

$$b^{(3)} := b^{(3)} - \alpha \frac{\partial E}{\partial b^{(3)}}$$

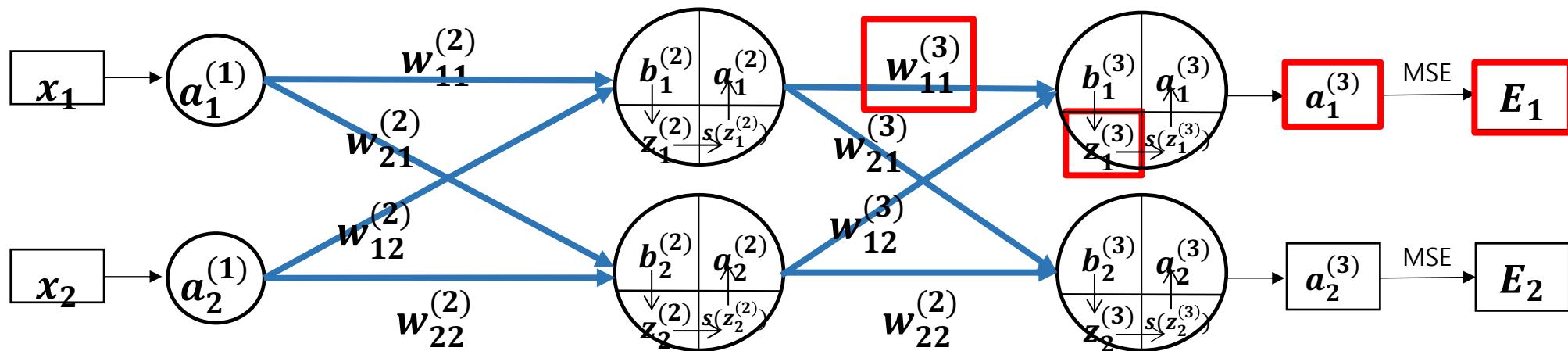
$$W^{(3)} = \begin{bmatrix} w_{11}^{(3)} & w_{21}^{(3)} \\ w_{12}^{(3)} & w_{22}^{(3)} \end{bmatrix} \quad b^{(3)} = [b_1^{(3)} \quad b_2^{(3)}]$$

$$\frac{\partial E}{\partial W^{(3)}} \rightarrow \frac{\partial E}{\partial b^{(3)}}$$

$$\frac{\partial E}{\partial W^{(3)}} = \begin{bmatrix} \frac{\partial E}{\partial w_{11}^{(3)}} & \frac{\partial E}{\partial w_{21}^{(3)}} \\ \frac{\partial E}{\partial w_{12}^{(3)}} & \frac{\partial E}{\partial w_{22}^{(3)}} \end{bmatrix}$$

$$\frac{\partial E}{\partial b^{(3)}} = \begin{bmatrix} \frac{\partial E}{\partial b_1^{(3)}} & \frac{\partial E}{\partial b_2^{(3)}} \end{bmatrix}$$

$$\frac{\partial E}{\partial w_{11}^{(3)}} = ?$$



출력 층에서의 오차역전파

$w_{11}^{(3)}$ 에 대한 E 의 변화율 - $\frac{\partial E}{\partial w_{11}^{(3)}}$

$$\frac{\partial E}{\partial W^{(3)}} = \begin{bmatrix} \frac{\partial E}{\partial w_{11}^{(3)}} & \frac{\partial E}{\partial w_{21}^{(3)}} \\ \frac{\partial E}{\partial w_{12}^{(3)}} & \frac{\partial E}{\partial w_{22}^{(3)}} \end{bmatrix} \quad \frac{\partial E}{\partial b^{(3)}} = \begin{bmatrix} \frac{\partial E}{\partial b_1^{(3)}} & \frac{\partial E}{\partial b_2^{(3)}} \end{bmatrix}$$

➤ $\frac{\partial E}{\partial w_{11}^{(3)}} = ?$

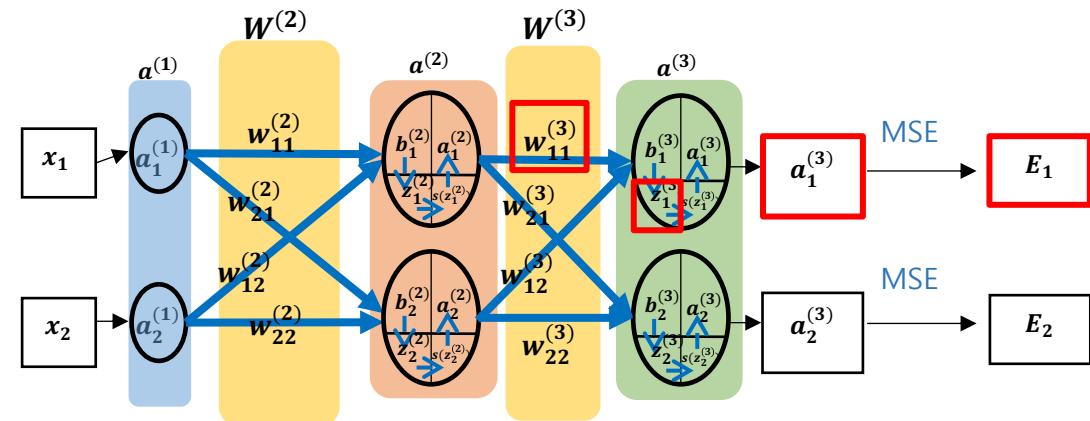
$$\frac{\partial E}{\partial w_{11}^{(3)}} = \frac{\partial E_1}{\partial w_{11}^{(3)}} + \frac{\partial E_2}{\partial w_{11}^{(3)}} \quad (\because E = E_1 + E_2)$$

$$= \frac{\partial E_1}{\partial w_{11}^{(3)}} \quad \left(\because \frac{\partial E_2}{\partial w_{11}^{(3)}} = 0 \right)$$

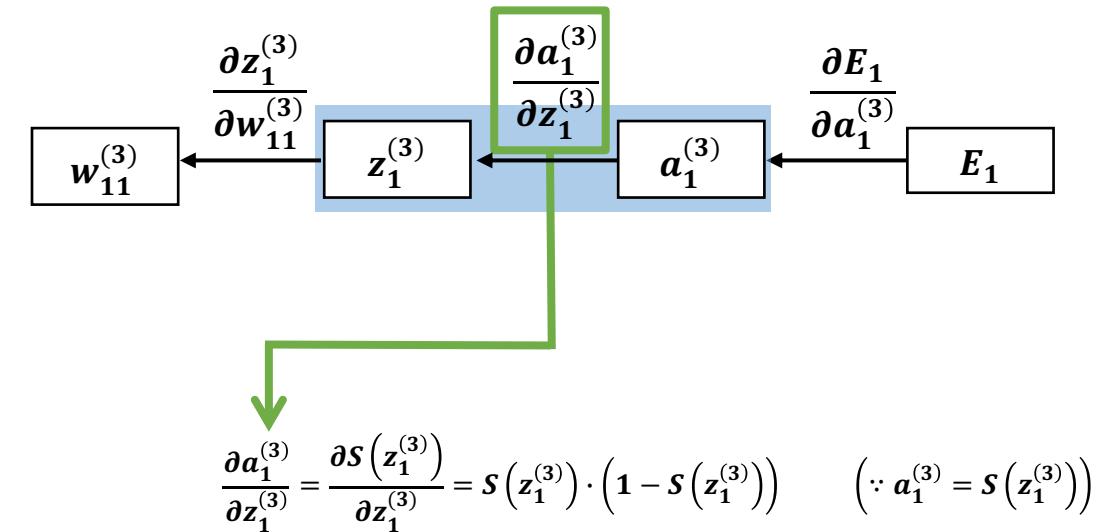
$$= \frac{\partial E_1}{\partial a_1^{(3)}} \cdot \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} \cdot \frac{\partial z_1^{(3)}}{\partial w_{11}^{(3)}}$$

$$= \frac{\partial \left(\frac{1}{2} (t_1^{(3)} - a_1^{(3)})^2 \right)}{\partial a_1^{(3)}} \cdot \frac{\partial S(z_1^{(3)})}{\partial z_1^{(3)}} \cdot \frac{\partial (a_1^{(2)} w_{11}^{(3)} + a_2^{(2)} w_{12}^{(3)} + b_1^{(3)})}{\partial w_{11}^{(3)}}$$

- ... ①
- ... ②
- ... ③
- ... ④



Back Propagation



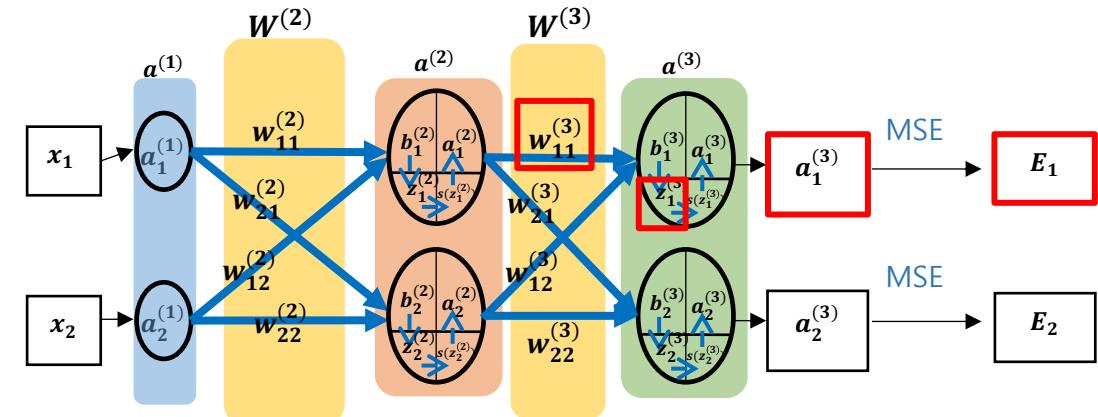
출력 층에서의 오차역전파

$w_{11}^{(3)}$ 에 대한 E 의 변화율 - $\frac{\partial E}{\partial w_{11}^{(3)}}$

$$\triangleright \frac{\partial E}{\partial w_{11}^{(3)}} = ?$$

$$= \frac{\partial \left(\frac{1}{2} (t_1^{(3)} - a_1^{(3)})^2 \right)}{\partial a_1^{(3)}} \cdot \frac{\partial S(z_1^{(3)})}{\partial z_1^{(3)}} \cdot \frac{\partial (a_1^{(2)} w_{11}^{(3)} + a_2^{(2)} w_{12}^{(3)} + b_1^{(3)})}{\partial w_{11}^{(3)}}$$

... ④



$$= (a_1^{(3)} - t_1^{(3)}) \cdot \left(S(z_1^{(3)}) \cdot (1 - S(z_1^{(3)})) \right) \cdot a_1^{(2)} \quad \dots \quad ⑤$$

$$= (a_1^{(3)} - t_1^{(3)}) \cdot a_1^{(3)} \cdot (1 - a_1^{(3)}) \cdot a_1^{(2)} \quad \dots \quad ⑥$$

$$\left(\because a_1^{(3)} = S(z_1^{(3)}) \right)$$

$$\therefore \frac{\partial E}{\partial w_{11}^{(3)}} = (a_1^{(3)} - t_1^{(3)}) \cdot a_1^{(3)} \cdot (1 - a_1^{(3)}) \cdot a_1^{(2)} \quad \dots \quad ⑦$$

- $\frac{\partial \left(\frac{1}{2} (t_1^{(3)} - a_1^{(3)})^2 \right)}{\partial a_1^{(3)}} = 2 \times \left(\frac{1}{2} (t_1^{(3)} - a_1^{(3)}) \right) \times (-1) = a_1^{(3)} - t_1^{(3)}$ 속미분
- $\frac{\partial S(z_1^{(3)})}{\partial z_1^{(3)}} = S(z_1^{(3)}) \cdot (1 - S(z_1^{(3)})) \quad (\because a_1^{(3)} = S(z_1^{(3)}))$
- $\frac{\partial (a_1^{(2)} w_{11}^{(3)} + a_2^{(2)} w_{12}^{(3)} + b_1^{(3)})}{\partial w_{11}^{(3)}} = a_1^{(2)}$

출력 층에서의 오차역전파

$w_{21}^{(3)}$ 에 대한 E 의 변화율 - $\frac{\partial E}{\partial w_{21}^{(3)}}$

$$W^{(3)} := W^{(3)} - \alpha \frac{\partial E}{\partial W^{(3)}}$$

$$b^{(3)} := b^{(3)} - \alpha \frac{\partial E}{\partial b^{(3)}}$$

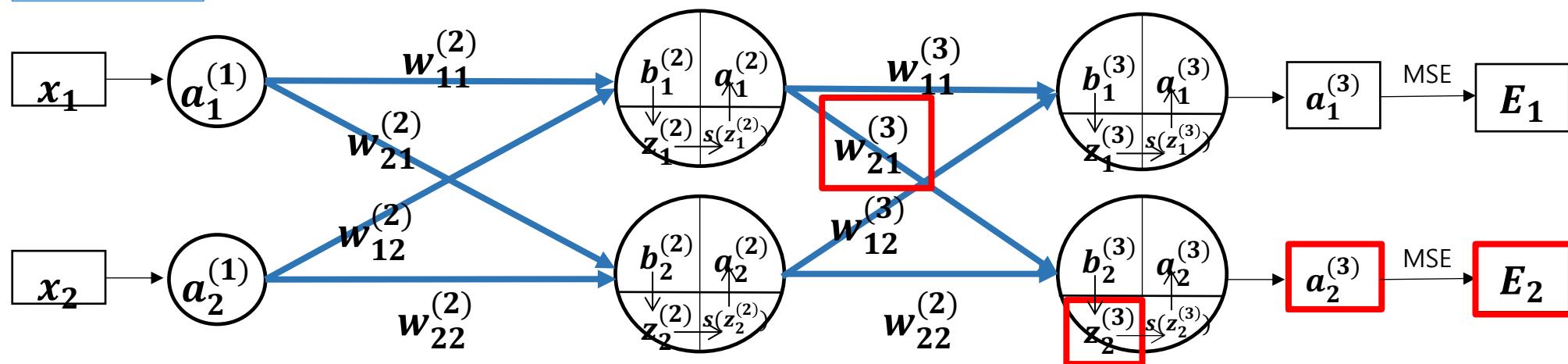
$$W^{(3)} = \begin{bmatrix} w_{11}^{(3)} & w_{21}^{(3)} \\ w_{12}^{(3)} & w_{22}^{(3)} \end{bmatrix} \quad b^{(3)} = [b_1^{(3)} \quad b_2^{(3)}]$$

$$\frac{\partial E}{\partial W^{(3)}} \rightarrow \frac{\partial E}{\partial b^{(3)}}$$

$$\frac{\partial E}{\partial W^{(3)}} = \begin{bmatrix} \frac{\partial E}{\partial w_{11}^{(3)}} & \frac{\partial E}{\partial w_{21}^{(3)}} \\ \frac{\partial E}{\partial w_{12}^{(3)}} & \frac{\partial E}{\partial w_{22}^{(3)}} \end{bmatrix}$$

$$\frac{\partial E}{\partial b^{(3)}} = \begin{bmatrix} \frac{\partial E}{\partial b_1^{(3)}} & \frac{\partial E}{\partial b_2^{(3)}} \end{bmatrix}$$

$$\frac{\partial E}{\partial w_{21}^{(3)}} = ?$$



출력 층에서의 오차역전파

$w_{21}^{(3)}$ 에 대한 E 의 변화율 - $\frac{\partial E}{\partial w_{21}^{(3)}}$

$$\frac{\partial E}{\partial W^{(3)}} = \begin{bmatrix} \frac{\partial E}{\partial w_{11}^{(3)}} & \frac{\partial E}{\partial w_{21}^{(3)}} \\ \frac{\partial E}{\partial w_{12}^{(3)}} & \frac{\partial E}{\partial w_{22}^{(3)}} \end{bmatrix} \quad \frac{\partial E}{\partial b^{(3)}} = \begin{bmatrix} \frac{\partial E}{\partial b_1^{(3)}} & \frac{\partial E}{\partial b_2^{(3)}} \end{bmatrix}$$

➤ $\frac{\partial E}{\partial w_{21}^{(3)}} = ?$

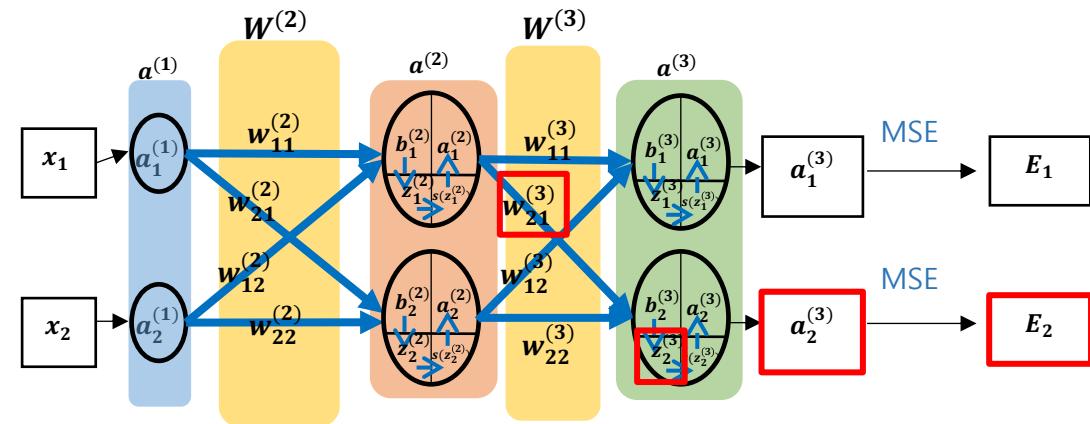
$$\frac{\partial E}{\partial w_{21}^{(3)}} = \frac{\partial E_1}{\partial w_{21}^{(3)}} + \frac{\partial E_2}{\partial w_{21}^{(3)}} \quad (\because E = E_1 + E_2)$$

$$= \frac{\partial E_2}{\partial w_{21}^{(3)}} \quad \left(\because \frac{\partial E_1}{\partial w_{21}^{(3)}} = 0 \right)$$

$$= \frac{\partial E_2}{\partial a_2^{(3)}} \cdot \frac{\partial a_2^{(3)}}{\partial z_2^{(3)}} \cdot \frac{\partial z_2^{(3)}}{\partial w_{21}^{(3)}}$$

$$= \frac{\partial \left(\frac{1}{2} (t_2^{(3)} - a_2^{(3)})^2 \right)}{\partial a_2^{(3)}} \cdot \frac{\partial S(z_2^{(3)})}{\partial z_2^{(3)}} \cdot \frac{\partial (a_1^{(2)} w_{21}^{(3)} + a_2^{(2)} w_{22}^{(3)} + b_2^{(3)})}{\partial w_{21}^{(3)}}$$

- ... ①
- ... ②
- ... ③
- ... ④



Back Propagation

$$\begin{aligned} \frac{\partial z_2^{(3)}}{\partial w_{21}^{(3)}} &\rightarrow \frac{\partial a_2^{(3)}}{\partial z_2^{(3)}} \rightarrow \frac{\partial E_2}{\partial a_2^{(3)}} \rightarrow E_2 \\ \frac{\partial a_2^{(3)}}{\partial z_2^{(3)}} &= \frac{\partial S(z_2^{(3)})}{\partial z_2^{(3)}} = S(z_2^{(3)}) \cdot (1 - S(z_2^{(3)})) \quad (\because a_2^{(3)} = S(z_2^{(3)})) \end{aligned}$$

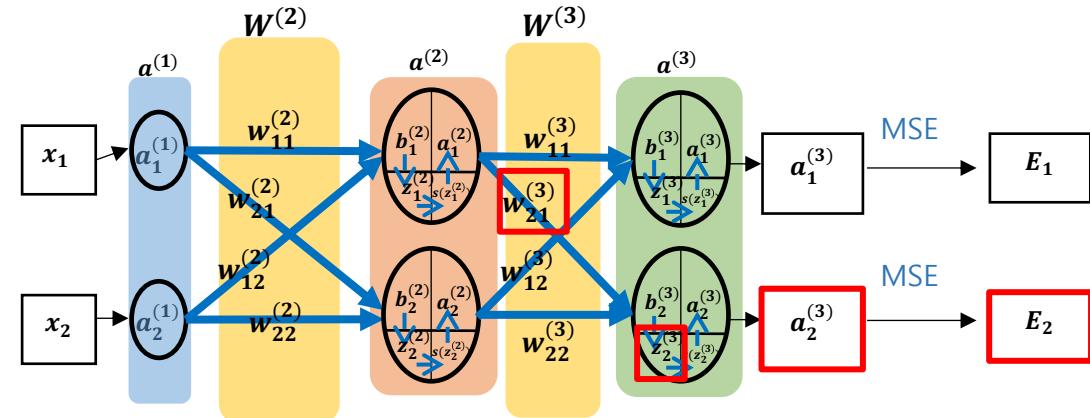
출력 층에서의 오차역전파

$w_{21}^{(3)}$ 에 대한 E 의 변화율 - $\frac{\partial E}{\partial w_{21}^{(3)}}$

$$\triangleright \frac{\partial E}{\partial w_{21}^{(3)}} = ?$$

$$= \frac{\partial \left(\frac{1}{2} (t_2^{(3)} - a_2^{(3)})^2 \right)}{\partial a_2^{(3)}} \cdot \frac{\partial S(z_2^{(3)})}{\partial z_2^{(3)}} \cdot \frac{\partial (a_1^{(2)} w_{21}^{(3)} + a_2^{(2)} w_{22}^{(3)} + b_2^{(3)})}{\partial w_{21}^{(3)}}$$

... ④



$$= (a_2^{(3)} - t_2^{(3)}) \cdot \left(S(z_2^{(3)}) \cdot (1 - S(z_2^{(3)})) \right) \cdot a_1^{(2)} \quad \dots \quad ⑤$$

$$= (a_1^{(3)} - t_1^{(3)}) \cdot a_2^{(3)} \cdot (1 - a_2^{(3)}) \cdot a_1^{(2)} \quad \dots \quad ⑥$$

$$\left(\because a_2^{(3)} = S(z_2^{(3)}) \right)$$

$$\therefore \frac{\partial E}{\partial w_{21}^{(3)}} = (a_2^{(3)} - t_2^{(3)}) \cdot a_2^{(3)} \cdot (1 - a_2^{(3)}) \cdot a_1^{(2)} \quad \dots \quad ⑦$$

<ul style="list-style-type: none"> $\frac{\partial \left(\frac{1}{2} (t_2^{(3)} - a_2^{(3)})^2 \right)}{\partial a_2^{(3)}} = 2 \times \left(\frac{1}{2} (t_2^{(3)} - a_2^{(3)}) \right) \times (-1) = a_2^{(3)} - t_2^{(3)}$ <small>속미분</small> $\frac{\partial S(z_2^{(3)})}{\partial z_2^{(3)}} = S(z_2^{(3)}) \cdot (1 - S(z_2^{(3)})) \quad (\because a_1^{(3)} = S(z_1^{(3)}))$ $\frac{\partial (a_1^{(2)} w_{21}^{(3)} + a_2^{(2)} w_{22}^{(3)} + b_2^{(3)})}{\partial w_{21}^{(3)}} = a_1^{(2)}$

출력 층에서의 오차역전파

$w_{12}^{(3)}$ 에 대한 E 의 변화율 - $\frac{\partial E}{\partial w_{12}^{(3)}}$

$$W^{(3)} := W^{(3)} - \alpha \frac{\partial E}{\partial W^{(3)}}$$

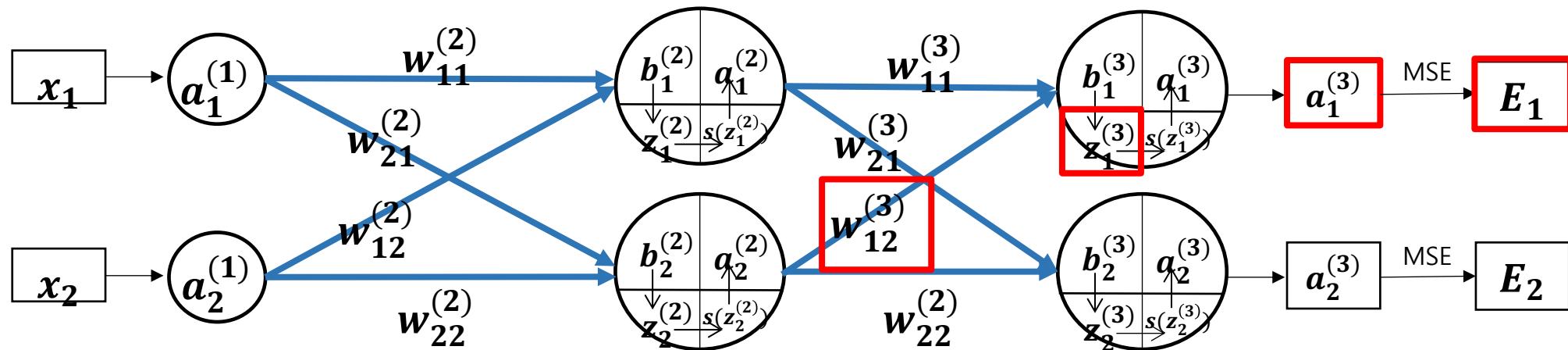
$$b^{(3)} := b^{(3)} - \alpha \frac{\partial E}{\partial b^{(3)}}$$

$$W^{(3)} = \begin{bmatrix} w_{11}^{(3)} & w_{21}^{(3)} \\ w_{12}^{(3)} & w_{22}^{(3)} \end{bmatrix} \quad b^{(3)} = [b_1^{(3)} \quad b_2^{(3)}]$$

$$\frac{\partial E}{\partial W^{(3)}} \rightarrow \frac{\partial E}{\partial b^{(3)}}$$

$$\frac{\partial E}{\partial W^{(3)}} = \begin{bmatrix} \frac{\partial E}{\partial w_{11}^{(3)}} & \frac{\partial E}{\partial w_{21}^{(3)}} \\ \frac{\partial E}{\partial w_{12}^{(3)}} & \frac{\partial E}{\partial w_{22}^{(3)}} \end{bmatrix} \quad \frac{\partial E}{\partial b^{(3)}} = \begin{bmatrix} \frac{\partial E}{\partial b_1^{(3)}} & \frac{\partial E}{\partial b_2^{(3)}} \end{bmatrix}$$

$$\frac{\partial E}{\partial w_{12}^{(3)}} = ?$$



출력 층에서의 오차역전파

$w_{12}^{(3)}$ 에 대한 E 의 변화율 - $\frac{\partial E}{\partial w_{12}^{(3)}}$

$$\frac{\partial E}{\partial W^{(3)}} = \begin{bmatrix} \frac{\partial E}{\partial w_{11}^{(3)}} & \frac{\partial E}{\partial w_{21}^{(3)}} \\ \frac{\partial E}{\partial w_{12}^{(3)}} & \frac{\partial E}{\partial w_{22}^{(3)}} \end{bmatrix} \quad \frac{\partial E}{\partial b^{(3)}} = \begin{bmatrix} \frac{\partial E}{\partial b_1^{(3)}} & \frac{\partial E}{\partial b_2^{(3)}} \end{bmatrix}$$

➤ $\frac{\partial E}{\partial w_{12}^{(3)}} = ?$

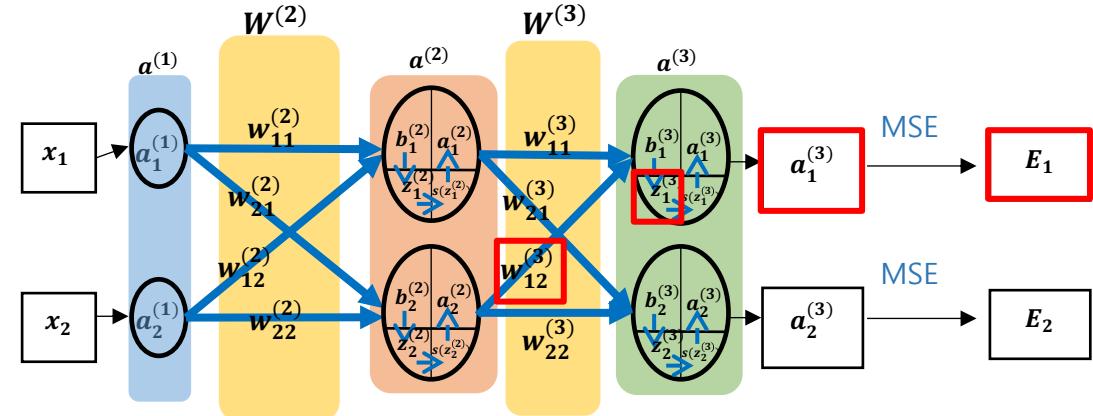
$$\frac{\partial E}{\partial w_{12}^{(3)}} = \frac{\partial E_1}{\partial w_{12}^{(3)}} + \frac{\partial E_2}{\partial w_{12}^{(3)}} \quad (\because E = E_1 + E_2)$$

$$= \frac{\partial E_1}{\partial w_{12}^{(3)}} \quad \left(\because \frac{\partial E_2}{\partial w_{12}^{(3)}} = 0 \right)$$

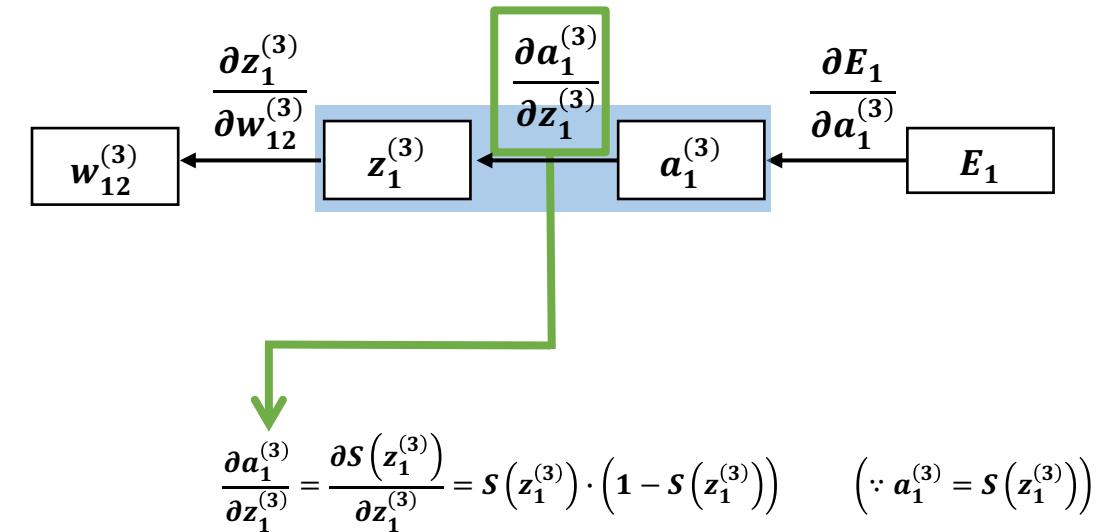
$$= \frac{\partial E_1}{\partial a_1^{(3)}} \cdot \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} \cdot \frac{\partial z_1^{(3)}}{\partial w_{12}^{(3)}}$$

$$= \frac{\partial \left(\frac{1}{2} (t_1^{(3)} - a_1^{(3)})^2 \right)}{\partial a_1^{(3)}} \cdot \frac{\partial S(z_1^{(3)})}{\partial z_1^{(3)}} \cdot \frac{\partial (a_1^{(2)} w_{11}^{(3)} + a_2^{(2)} w_{12}^{(3)} + b_1^{(3)})}{\partial w_{12}^{(3)}}$$

- ... ①
- ... ②
- ... ③
- ... ④



Back Propagation



출력 층에서의 오차역전파

$w_{12}^{(3)}$ 에 대한 E 의 변화율 - $\frac{\partial E}{\partial w_{12}^{(3)}}$

$$\triangleright \frac{\partial E}{\partial w_{12}^{(3)}} = ?$$

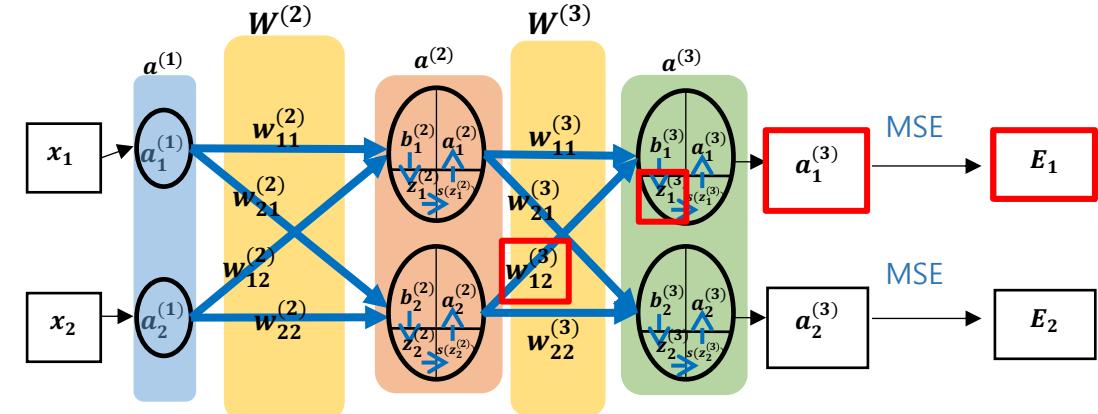
$$= \frac{\partial \left(\frac{1}{2} (t_1^{(3)} - a_1^{(3)})^2 \right)}{\partial a_1^{(3)}} \cdot \frac{\partial S(z_1^{(3)})}{\partial z_1^{(3)}} \cdot \frac{\partial (a_1^{(2)} w_{11}^{(3)} + a_2^{(2)} w_{12}^{(3)} + b_1^{(3)})}{\partial w_{12}^{(3)}} \quad \dots \quad ④$$

$$= (a_1^{(3)} - t_1^{(3)}) \cdot \left(S(z_1^{(3)}) \cdot (1 - S(z_1^{(3)})) \right) \cdot a_2^{(2)} \quad \dots \quad ⑤$$

$$= (a_1^{(3)} - t_1^{(3)}) \cdot a_1^{(3)} \cdot (1 - a_1^{(3)}) \cdot a_2^{(2)} \quad \dots \quad ⑥$$

$$\left(\because a_1^{(3)} = S(z_1^{(3)}) \right)$$

$$\therefore \frac{\partial E}{\partial w_{12}^{(3)}} = (a_1^{(3)} - t_1^{(3)}) \cdot a_1^{(3)} \cdot (1 - a_1^{(3)}) \cdot a_2^{(2)} \quad \dots \quad ⑦$$



- $\frac{\partial \left(\frac{1}{2} (t_1^{(3)} - a_1^{(3)})^2 \right)}{\partial a_1^{(3)}} = 2 \times \left(\frac{1}{2} (t_1^{(3)} - a_1^{(3)}) \right) \times (-1) = a_1^{(3)} - t_1^{(3)}$ 속미분
- $\frac{\partial S(z_1^{(3)})}{\partial z_1^{(3)}} = S(z_1^{(3)}) \cdot (1 - S(z_1^{(3)})) \quad (\because a_1^{(3)} = S(z_1^{(3)}))$
- $\frac{\partial (a_1^{(2)} w_{11}^{(3)} + a_2^{(2)} w_{12}^{(3)} + b_1^{(3)})}{\partial w_{12}^{(3)}} = a_2^{(2)}$

출력 층에서의 오차역전파

$w_{22}^{(3)}$ 에 대한 E 의 변화율 - $\frac{\partial E}{\partial w_{22}^{(3)}}$

$$W^{(3)} := W^{(3)} - \alpha \frac{\partial E}{\partial W^{(3)}}$$

$$b^{(3)} := b^{(3)} - \alpha \frac{\partial E}{\partial b^{(3)}}$$

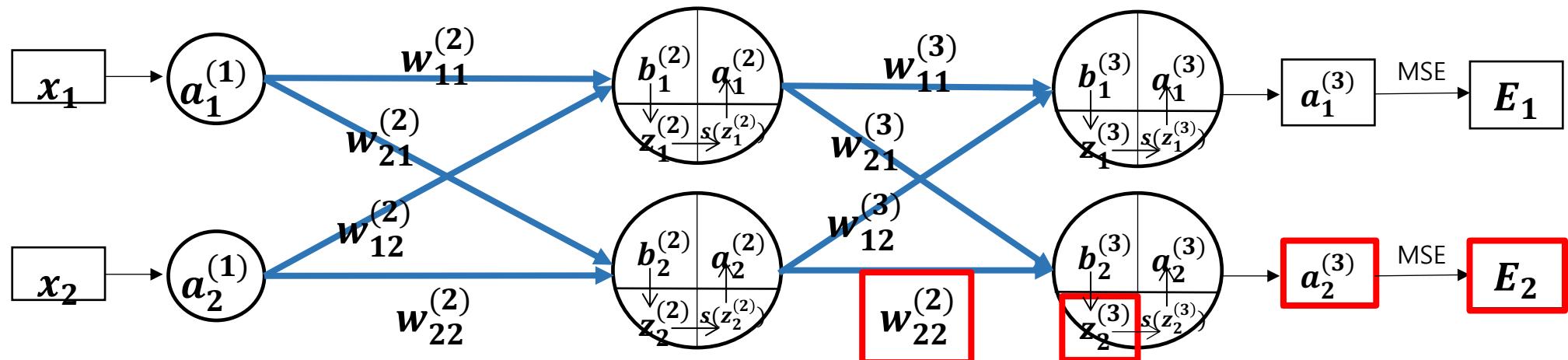
$$W^{(3)} = \begin{bmatrix} w_{11}^{(3)} & w_{21}^{(3)} \\ w_{12}^{(3)} & w_{22}^{(3)} \end{bmatrix} \quad b^{(3)} = [b_1^{(3)} \quad b_2^{(3)}]$$

$$\frac{\partial E}{\partial W^{(3)}} \rightarrow \frac{\partial E}{\partial b^{(3)}}$$

$$\frac{\partial E}{\partial W^{(3)}} = \begin{bmatrix} \frac{\partial E}{\partial w_{11}^{(3)}} & \frac{\partial E}{\partial w_{21}^{(3)}} \\ \frac{\partial E}{\partial w_{12}^{(3)}} & \frac{\partial E}{\partial w_{22}^{(3)}} \end{bmatrix}$$

$$\frac{\partial E}{\partial b^{(3)}} = \begin{bmatrix} \frac{\partial E}{\partial b_1^{(3)}} & \frac{\partial E}{\partial b_2^{(3)}} \end{bmatrix}$$

$$\frac{\partial E}{\partial w_{22}^{(3)}} = ?$$



출력 층에서의 오차역전파

$w_{22}^{(3)}$ 에 대한 E 의 변화율 - $\frac{\partial E}{\partial w_{22}^{(3)}}$

$$\frac{\partial E}{\partial W^{(3)}} = \begin{bmatrix} \frac{\partial E}{\partial w_{11}^{(3)}} & \frac{\partial E}{\partial w_{21}^{(3)}} \\ \frac{\partial E}{\partial w_{12}^{(3)}} & \frac{\partial E}{\partial w_{22}^{(3)}} \end{bmatrix} \quad \frac{\partial E}{\partial b^{(3)}} = \begin{bmatrix} \frac{\partial E}{\partial b_1^{(3)}} & \frac{\partial E}{\partial b_2^{(3)}} \end{bmatrix}$$

➤ $\frac{\partial E}{\partial w_{22}^{(3)}} = ?$

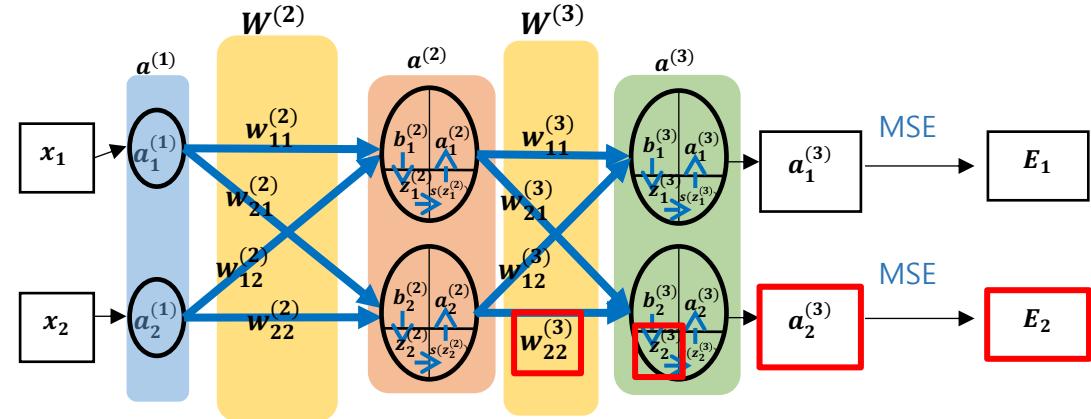
$$\frac{\partial E}{\partial w_{22}^{(3)}} = \frac{\partial E_1}{\partial w_{22}^{(3)}} + \frac{\partial E_2}{\partial w_{22}^{(3)}} \quad (\because E = E_1 + E_2)$$

$$= \frac{\partial E_2}{\partial w_{22}^{(3)}} \quad \left(\because \frac{\partial E_1}{\partial w_{22}^{(3)}} = 0 \right)$$

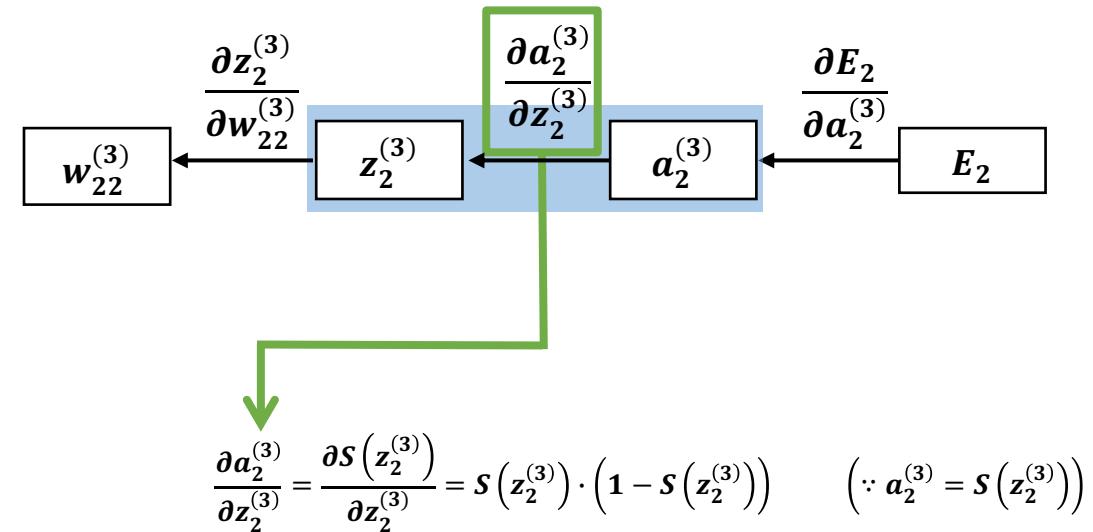
$$= \frac{\partial E_2}{\partial a_2^{(3)}} \cdot \frac{\partial a_2^{(3)}}{\partial z_2^{(3)}} \cdot \frac{\partial z_2^{(3)}}{\partial w_{22}^{(3)}}$$

$$= \frac{\partial \left(\frac{1}{2} (t_2^{(3)} - a_2^{(3)})^2 \right)}{\partial a_2^{(3)}} \cdot \frac{\partial S(z_2^{(3)})}{\partial z_2^{(3)}} \cdot \frac{\partial (a_1^{(2)} w_{21}^{(3)} + a_2^{(2)} w_{22}^{(3)} + b_2^{(3)})}{\partial w_{22}^{(3)}}$$

- ... ①
- ... ②
- ... ③
- ... ④



Back Propagation



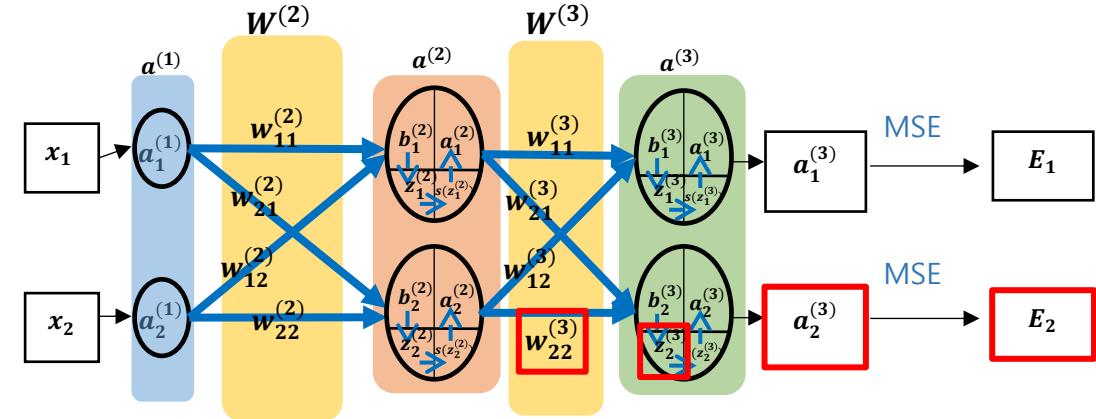
출력 층에서의 오차역전파

$w_{22}^{(3)}$ 에 대한 E 의 변화율 - $\frac{\partial E}{\partial w_{22}^{(3)}}$

$$\triangleright \frac{\partial E}{\partial w_{22}^{(3)}} = ?$$

$$= \frac{\partial \left(\frac{1}{2} (t_2^{(3)} - a_2^{(3)})^2 \right)}{\partial a_2^{(3)}} \cdot \frac{\partial S(z_2^{(3)})}{\partial z_2^{(3)}} \cdot \frac{\partial (a_1^{(2)} w_{21}^{(3)} + a_2^{(2)} w_{22}^{(3)} + b_2^{(3)})}{\partial w_{21}^{(3)}}$$

... ④



$$= (a_2^{(3)} - t_2^{(3)}) \cdot \left(S(z_2^{(3)}) \cdot (1 - S(z_2^{(3)})) \right) \cdot a_2^{(2)} \quad \dots \quad ⑤$$

$$= (a_1^{(3)} - t_1^{(3)}) \cdot a_2^{(3)} \cdot (1 - a_2^{(3)}) \cdot a_2^{(2)} \quad \dots \quad ⑥$$

$$\left(\because a_2^{(3)} = S(z_2^{(3)}) \right)$$

$$\therefore \frac{\partial E}{\partial w_{22}^{(3)}} = (a_2^{(3)} - t_2^{(3)}) \cdot a_2^{(3)} \cdot (1 - a_2^{(3)}) \cdot a_2^{(2)} \quad \dots \quad ⑦$$

<ul style="list-style-type: none"> $\frac{\partial \left(\frac{1}{2} (t_2^{(3)} - a_2^{(3)})^2 \right)}{\partial a_2^{(3)}} = 2 \times \left(\frac{1}{2} (t_2^{(3)} - a_2^{(3)}) \right) \times (-1) = a_2^{(3)} - t_2^{(3)}$ <small>속미분</small> $\frac{\partial S(z_2^{(3)})}{\partial z_2^{(3)}} = S(z_2^{(3)}) \cdot (1 - S(z_2^{(3)})) \quad (\because a_1^{(3)} = S(z_1^{(3)}))$ $\frac{\partial (a_1^{(2)} w_{21}^{(3)} + a_2^{(2)} w_{22}^{(3)} + b_2^{(3)})}{\partial w_{22}^{(3)}} = a_2^{(2)}$

출력 층에서의 오차역전파

$b_1^{(3)}$ 에 대한 E 의 변화율 - $\frac{\partial E}{\partial b_1^{(3)}}$

$$W^{(3)} := W^{(3)} - \alpha \frac{\partial E}{\partial W^{(3)}}$$

$$b^{(3)} := b^{(3)} - \alpha \frac{\partial E}{\partial b^{(3)}}$$

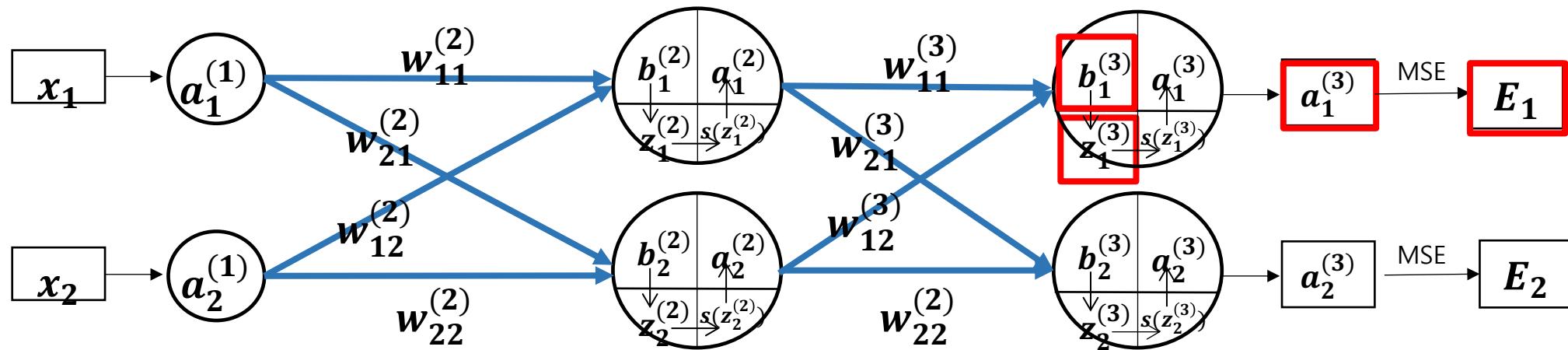
$$W^{(3)} = \begin{bmatrix} w_{11}^{(3)} & w_{21}^{(3)} \\ w_{12}^{(3)} & w_{22}^{(3)} \end{bmatrix} \quad b^{(3)} = [b_1^{(3)} \quad b_2^{(3)}]$$

$$\frac{\partial E}{\partial b^{(3)}} \rightarrow \frac{\partial E}{\partial b_1^{(3)}}$$

$$\frac{\partial E}{\partial W^{(3)}} = \begin{bmatrix} \frac{\partial E}{\partial w_{11}^{(3)}} & \frac{\partial E}{\partial w_{21}^{(3)}} \\ \frac{\partial E}{\partial w_{12}^{(3)}} & \frac{\partial E}{\partial w_{22}^{(3)}} \end{bmatrix}$$

$$\frac{\partial E}{\partial b^{(3)}} = \begin{bmatrix} \frac{\partial E}{\partial b_1^{(3)}} & \frac{\partial E}{\partial b_2^{(3)}} \end{bmatrix}$$

$$\frac{\partial E}{\partial w_{21}^{(3)}} = ?$$



출력 층에서의 오차역전파

$b_1^{(3)}$ 에 대한 E 의 변화율 - $\frac{\partial E}{\partial b_1^{(3)}}$

$$\frac{\partial E}{\partial W^{(3)}} = \begin{bmatrix} \frac{\partial E}{\partial w_{11}^{(3)}} & \frac{\partial E}{\partial w_{21}^{(3)}} \\ \frac{\partial E}{\partial w_{12}^{(3)}} & \frac{\partial E}{\partial w_{22}^{(3)}} \end{bmatrix} \quad \frac{\partial E}{\partial b^{(3)}} = \begin{bmatrix} \frac{\partial E}{\partial b_1^{(3)}} & \frac{\partial E}{\partial b_2^{(3)}} \end{bmatrix}$$

➤ $\frac{\partial E}{\partial b_1^{(3)}} = ?$

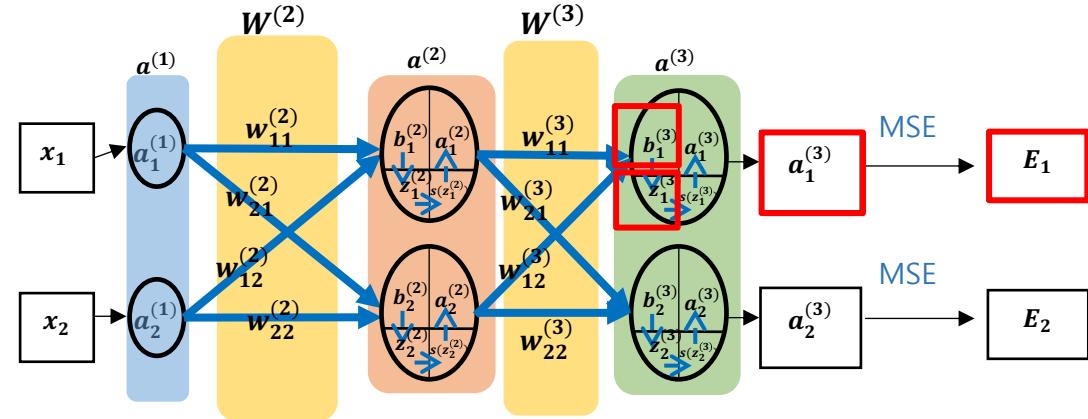
$$\frac{\partial E}{\partial b_1^{(3)}} = \frac{\partial E_1}{\partial b_1^{(3)}} + \frac{\partial E_2}{\partial b_1^{(3)}} \quad (\because E = E_1 + E_2)$$

$$= \frac{\partial E_1}{\partial b_1^{(3)}} \quad \left(\because \frac{\partial E_2}{\partial b_1^{(3)}} = 0 \right)$$

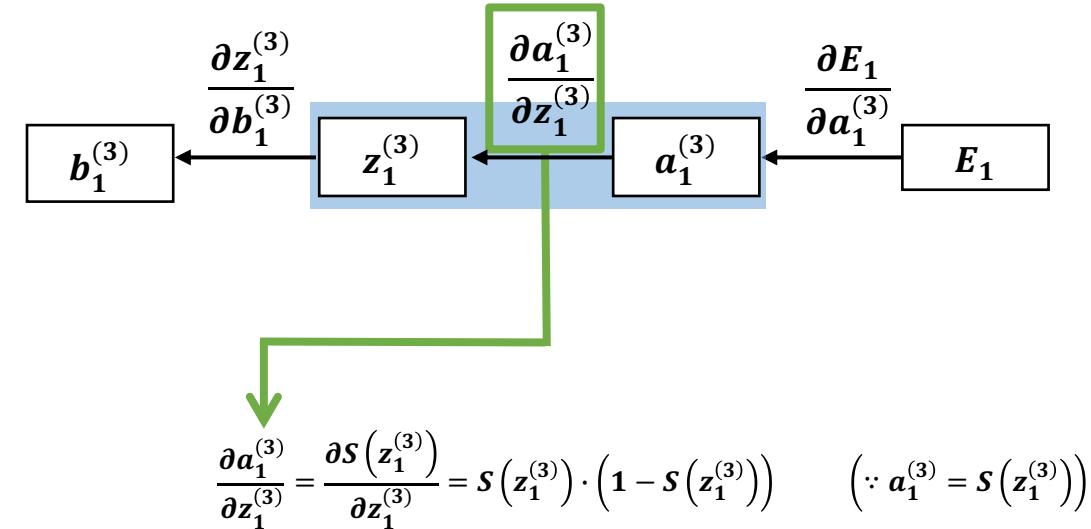
$$= \frac{\partial E_1}{\partial a_1^{(3)}} \cdot \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} \cdot \frac{\partial z_1^{(3)}}{\partial b_1^{(3)}}$$

$$= \frac{\partial \left(\frac{1}{2} (t_1^{(3)} - a_1^{(3)})^2 \right)}{\partial a_1^{(3)}} \cdot \frac{\partial s(z_1^{(3)})}{\partial z_1^{(3)}} \cdot \frac{\partial (a_1^{(2)} w_{11}^{(3)} + a_2^{(2)} w_{12}^{(3)} + b_1^{(3)})}{\partial b_1^{(3)}}$$

- ... ①
- ... ②
- ... ③
- ... ④



Back Propagation



출력 층에서의 오차역전파

$b_1^{(3)}$ 에 대한 E 의 변화율 - $\frac{\partial E}{\partial b_1^{(3)}}$

$$\triangleright \frac{\partial E}{\partial b_1^{(3)}} = ?$$

$$= \frac{\partial \left(\frac{1}{2} (t_1^{(3)} - a_1^{(3)})^2 \right)}{\partial a_1^{(3)}} \cdot \frac{\partial S(z_1^{(3)})}{\partial z_1^{(3)}} \cdot \frac{\partial (a_1^{(2)} w_{11}^{(3)} + a_2^{(2)} w_{12}^{(3)} + b_1^{(3)})}{\partial b_1^{(3)}}$$

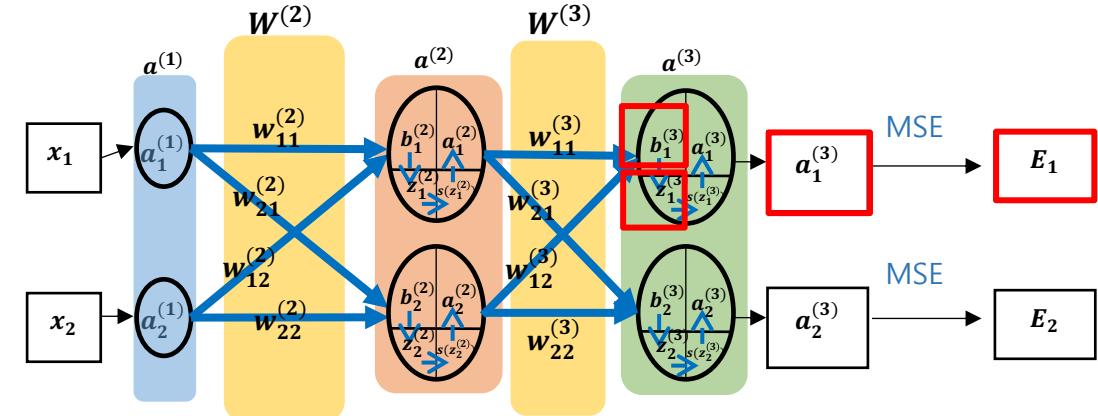
... ④

$$= (a_1^{(3)} - t_1^{(3)}) \cdot \left(S(z_1^{(3)}) \cdot (1 - S(z_1^{(3)})) \right) \cdot 1 \quad \dots \quad ⑤$$

$$= (a_1^{(3)} - t_1^{(3)}) \cdot a_1^{(3)} \cdot (1 - a_1^{(3)}) \cdot 1 \quad \dots \quad ⑥$$

$$\left(\because a_1^{(3)} = S(z_1^{(3)}) \right)$$

$$\therefore \frac{\partial E}{\partial b_1^{(3)}} = (a_1^{(3)} - t_1^{(3)}) \cdot a_1^{(3)} \cdot (1 - a_1^{(3)}) \quad \dots \quad ⑦$$



- $\frac{\partial \left(\frac{1}{2} (t_1^{(3)} - a_1^{(3)})^2 \right)}{\partial a_1^{(3)}} = 2 \times \left(\frac{1}{2} (t_1^{(3)} - a_1^{(3)}) \right) \times (-1) = a_1^{(3)} - t_1^{(3)}$ 속미분
- $\frac{\partial S(z_1^{(3)})}{\partial z_1^{(3)}} = S(z_1^{(3)}) \cdot (1 - S(z_1^{(3)})) \quad (\because a_1^{(3)} = S(z_1^{(3)}))$
- $\frac{\partial (a_1^{(2)} w_{11}^{(3)} + a_2^{(2)} w_{12}^{(3)} + b_1^{(3)})}{\partial b_1^{(3)}} = 1$

출력 층에서의 오차역전파

$b_2^{(3)}$ 에 대한 E 의 변화율 - $\frac{\partial E}{\partial b_2^{(3)}}$

$$W^{(3)} := W^{(3)} - \alpha \frac{\partial E}{\partial W^{(3)}}$$

$$b^{(3)} := b^{(3)} - \alpha \frac{\partial E}{\partial b^{(3)}}$$

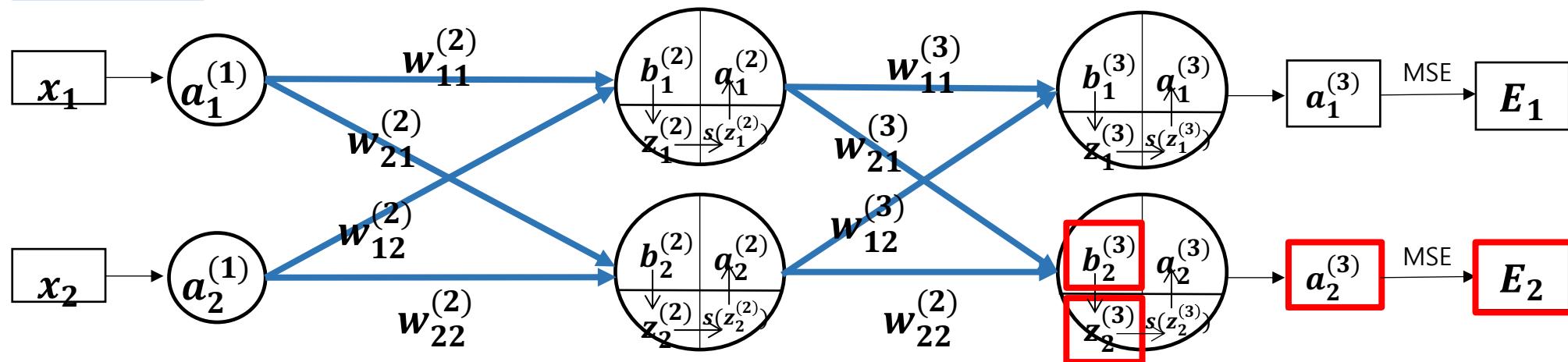
$$W^{(3)} = \begin{bmatrix} w_{11}^{(3)} & w_{21}^{(3)} \\ w_{12}^{(3)} & w_{22}^{(3)} \end{bmatrix} \quad b^{(3)} = [b_1^{(3)} \quad b_2^{(3)}]$$

$$\frac{\partial E}{\partial W^{(3)}} \quad \frac{\partial E}{\partial b^{(3)}}$$

$$\frac{\partial E}{\partial W^{(3)}} = \begin{bmatrix} \frac{\partial E}{\partial w_{11}^{(3)}} & \frac{\partial E}{\partial w_{21}^{(3)}} \\ \frac{\partial E}{\partial w_{12}^{(3)}} & \frac{\partial E}{\partial w_{22}^{(3)}} \end{bmatrix}$$

$$\frac{\partial E}{\partial b^{(3)}} = \left[\frac{\partial E}{\partial b_1^{(3)}} \quad \frac{\partial E}{\partial b_2^{(3)}} \right]$$

$$\frac{\partial E}{\partial w_{21}^{(3)}} = ?$$



출력 층에서의 오차역전파

$b_1^{(3)}$ 에 대한 E 의 변화율 - $\frac{\partial E}{\partial b_1^{(3)}}$

$$\frac{\partial E}{\partial W^{(3)}} = \begin{bmatrix} \frac{\partial E}{\partial w_{11}^{(3)}} & \frac{\partial E}{\partial w_{21}^{(3)}} \\ \frac{\partial E}{\partial w_{12}^{(3)}} & \frac{\partial E}{\partial w_{22}^{(3)}} \end{bmatrix} \quad \frac{\partial E}{\partial b^{(3)}} = \left[\frac{\partial E}{\partial b_1^{(3)}} \quad \boxed{\frac{\partial E}{\partial b_2^{(3)}}} \right]$$

➤ $\frac{\partial E}{\partial b_2^{(3)}} = ?$

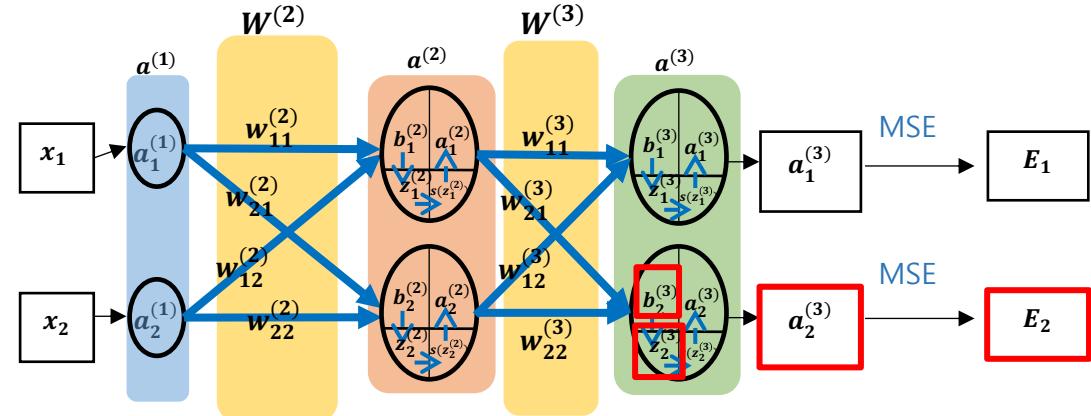
$$\frac{\partial E}{\partial b_2^{(3)}} = \frac{\partial E_1}{\partial b_2^{(3)}} + \frac{\partial E_2}{\partial b_2^{(3)}} \quad (\because E = E_1 + E_2)$$

$$= \frac{\partial E_2}{\partial b_2^{(3)}} \quad \left(\because \frac{\partial E_1}{\partial b_2^{(3)}} = 0 \right)$$

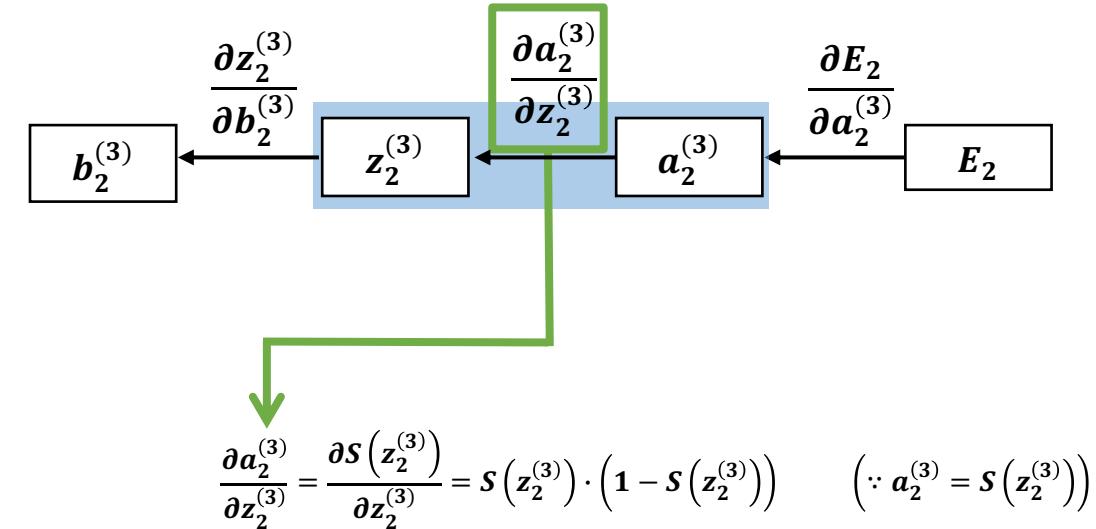
$$= \frac{\partial E_2}{\partial a_2^{(3)}} \cdot \frac{\partial a_2^{(3)}}{\partial z_2^{(3)}} \cdot \frac{\partial z_2^{(3)}}{\partial b_2^{(3)}}$$

$$= \frac{\partial \left(\frac{1}{2} (t_2^{(3)} - a_2^{(3)})^2 \right)}{\partial a_2^{(3)}} \cdot \frac{\partial S(z_2^{(3)})}{\partial z_2^{(3)}} \cdot \frac{\partial (a_1^{(2)} w_{21}^{(3)} + a_2^{(2)} w_{22}^{(3)} + b_2^{(3)})}{\partial b_2^{(3)}}$$

- ... ①
- ... ②
- ... ③
- ... ④



Back Propagation



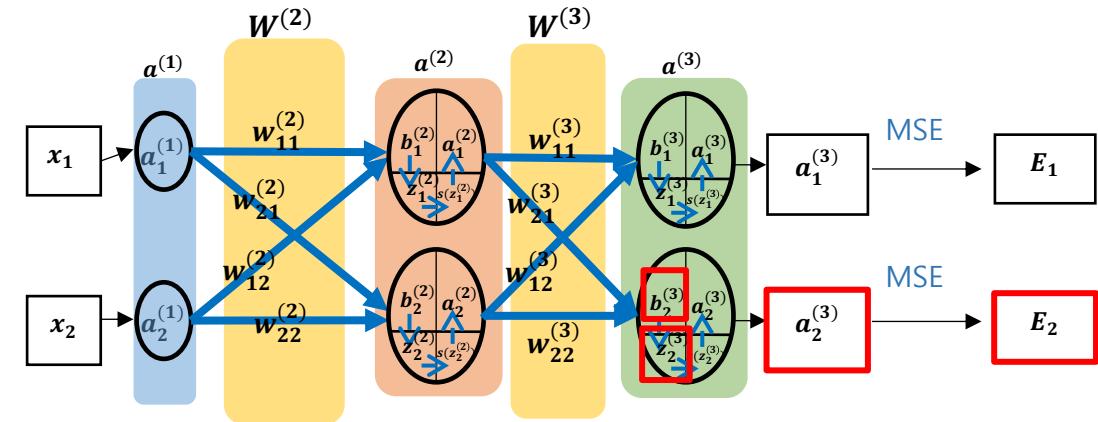
출력 층에서의 오차역전파

$b_1^{(3)}$ 에 대한 E 의 변화율 - $\frac{\partial E}{\partial b_1^{(3)}}$

$$\triangleright \frac{\partial E}{\partial b_2^{(3)}} = ?$$

$$= \frac{\partial \left(\frac{1}{2} (t_2^{(3)} - a_2^{(3)})^2 \right)}{\partial a_2^{(3)}} \cdot \frac{\partial S(z_2^{(3)})}{\partial z_2^{(3)}} \cdot \frac{\partial (a_1^{(2)} w_{21}^{(3)} + a_2^{(2)} w_{22}^{(3)} + b_2^{(3)})}{\partial b_2^{(3)}}$$

... ④



$$= (a_2^{(3)} - t_2^{(3)}) \cdot \left(S(z_2^{(3)}) \cdot (1 - S(z_2^{(3)})) \right) \cdot 1 \quad \dots \quad ⑤$$

$$= (a_1^{(3)} - t_1^{(3)}) \cdot a_2^{(3)} \cdot (1 - a_2^{(3)}) \cdot 1 \quad \dots \quad ⑥$$

$$\left(\because a_2^{(3)} = S(z_2^{(3)}) \right)$$

$$\therefore \frac{\partial E}{\partial b_2^{(3)}} = (a_2^{(3)} - t_2^{(3)}) \cdot a_2^{(3)} \cdot (1 - a_2^{(3)}) \quad \dots \quad ⑦$$

<ul style="list-style-type: none"> $\frac{\partial \left(\frac{1}{2} (t_2^{(3)} - a_2^{(3)})^2 \right)}{\partial a_2^{(3)}} = 2 \times \left(\frac{1}{2} (t_2^{(3)} - a_2^{(3)}) \right) \times (-1) = a_2^{(3)} - t_2^{(3)}$ 속미분 $\frac{\partial S(z_2^{(3)})}{\partial z_2^{(3)}} = S(z_2^{(3)}) \cdot (1 - S(z_2^{(3)})) \quad (\because a_1^{(3)} = S(z_1^{(3)}))$ $\frac{\partial (a_1^{(2)} w_{21}^{(3)} + a_2^{(2)} w_{22}^{(3)} + b_2^{(3)})}{\partial b_2^{(3)}} = 1$
--

출력 층에서의 오차역전파

출력 층에서의 오차역전파 - 정리

은닉층 출력 값 벡터	$A2 = (a_1^{(2)} \quad a_2^{(2)})$
출력층 가상 손실 벡터	$loss_3 = ((a_1^{(3)} - t_1^{(3)})a_1^{(3)}(1 - a_1^{(3)}) \quad (a_2^{(3)} - t_2^{(3)})a_2^{(3)}(1 - a_2^{(3)}))$

$$\bullet \frac{\partial E}{\partial W^{(3)}} = \begin{bmatrix} \frac{\partial E}{\partial w_{11}^{(3)}} & \frac{\partial E}{\partial w_{21}^{(3)}} \\ \frac{\partial E}{\partial w_{12}^{(3)}} & \frac{\partial E}{\partial w_{22}^{(2)}} \end{bmatrix} = \begin{bmatrix} (a_1^{(3)} - t_1^{(3)})a_1^{(3)}(1 - a_1^{(3)})a_1^{(2)} & (a_2^{(3)} - t_2^{(3)})a_2^{(3)}(1 - a_2^{(3)})a_1^{(2)} \\ (a_1^{(3)} - t_1^{(3)})a_1^{(3)}(1 - a_1^{(3)})a_2^{(2)} & (a_2^{(3)} - t_2^{(3)})a_2^{(3)}(1 - a_2^{(3)})a_2^{(2)} \end{bmatrix}$$

$$= \begin{bmatrix} a_1^{(2)} \\ a_2^{(2)} \end{bmatrix} [(a_1^{(3)} - t_1^{(3)})a_1^{(3)}(1 - a_1^{(3)}) \quad (a_2^{(3)} - t_2^{(3)})a_2^{(3)}(1 - a_2^{(3)})]$$

$$= A2^T \cdot loss_3$$

$$\bullet \frac{\partial E}{\partial b^{(3)}} = \begin{bmatrix} \frac{\partial E}{\partial b_1^{(3)}} & \frac{\partial E}{\partial b_2^{(3)}} \end{bmatrix} = [(a_1^{(3)} - t_1^{(3)})a_1^{(3)}(1 - a_1^{(3)}) \quad (a_2^{(3)} - t_2^{(3)})a_2^{(3)}(1 - a_2^{(3)})]$$

$$= loss_3$$

출력 층에서의 오차역전파

☞ 출력 층에서의 오차역전파 - 정리

은닉층 출력 값 벡터	$A2 = (a_1^{(2)} \quad a_2^{(2)})$
출력층 가상 손실 벡터	$loss_3 = ((a_1^{(3)} - t_1^{(3)})a_1^{(3)}(1 - a_1^{(3)}) \quad (a_2^{(3)} - t_2^{(3)})a_2^{(3)}(1 - a_2^{(3)}))$

- $\frac{\partial E}{\partial W^{(3)}} = A2^T \times loss_3$
- $\frac{\partial E}{\partial b^{(3)}} = loss_3$

➤ 출력층 가중치 계산 :

$$W^{(3)} := W^{(3)} - \alpha \frac{\partial E}{\partial W^{(3)}} = W^{(3)} - \alpha (A2^T \cdot loss_3)$$

➤ 출력층 편향 계산 :

$$b^{(3)} := b^{(3)} - \alpha \frac{\partial E}{\partial b^{(3)}} = b^{(3)} - \alpha (loss_3)$$

CONTENTS

- 1 오차역전파 개념
- 2 출력층에서의 오차역전파
- 3 은닉층에서의 오차역전파
- 4 오차역전파 이용 MNIST 검증

은닉층에서의 오차역전파

$w_{11}^{(2)}$ 에 대한 E 의 변화율 - $\frac{\partial E}{\partial w_{11}^{(2)}}$

$$W^{(2)} := W^{(2)} - \alpha \frac{\partial E}{\partial W^{(2)}}$$

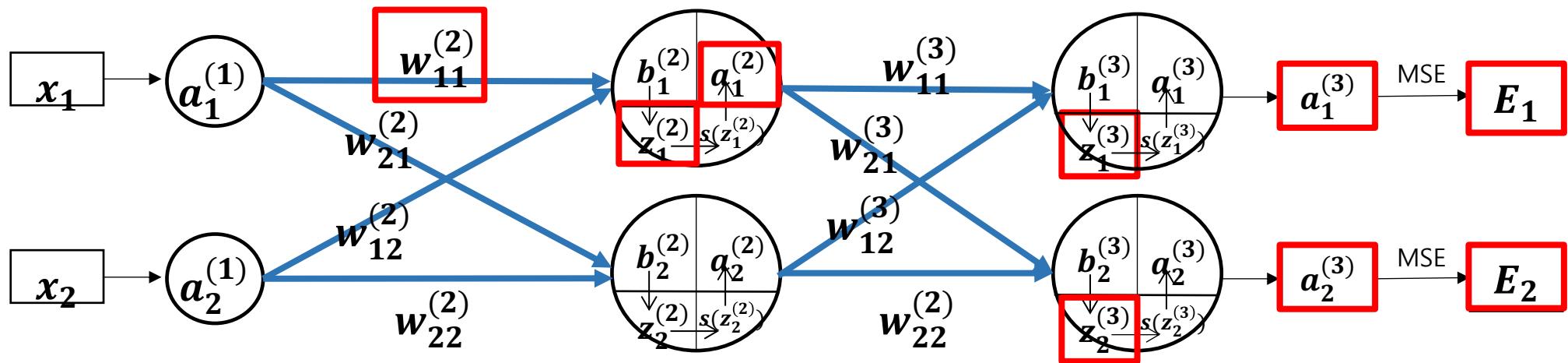
$$b^{(2)} := b^{(2)} - \alpha \frac{\partial E}{\partial b^{(2)}}$$

$$W^{(2)} = \begin{bmatrix} w_{11}^{(2)} & w_{21}^{(2)} \\ w_{12}^{(2)} & w_{22}^{(2)} \end{bmatrix} \quad b^{(2)} = [b_1^{(2)} \quad b_2^{(2)}]$$

$\frac{\partial E}{\partial W^{(2)}} \rightarrow \frac{\partial E}{\partial b^{(3)}}$

$$\frac{\partial E}{\partial W^{(2)}} = \begin{bmatrix} \frac{\partial E}{\partial w_{11}^{(2)}} & \frac{\partial E}{\partial w_{21}^{(2)}} \\ \frac{\partial E}{\partial w_{12}^{(2)}} & \frac{\partial E}{\partial w_{22}^{(2)}} \end{bmatrix} \quad \frac{\partial E}{\partial b^{(2)}} = \begin{bmatrix} \frac{\partial E}{\partial b_1^{(2)}} & \frac{\partial E}{\partial b_2^{(2)}} \end{bmatrix}$$

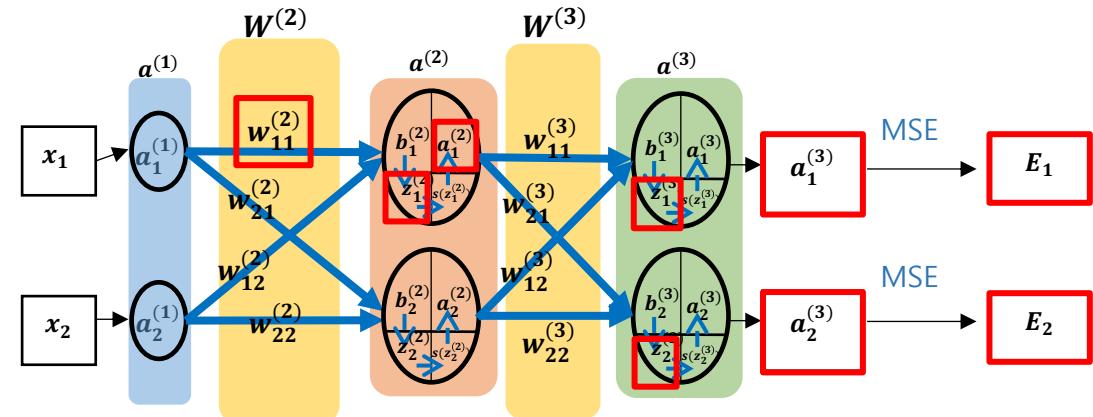
$$\frac{\partial E}{\partial w_{11}^{(2)}} = ?$$



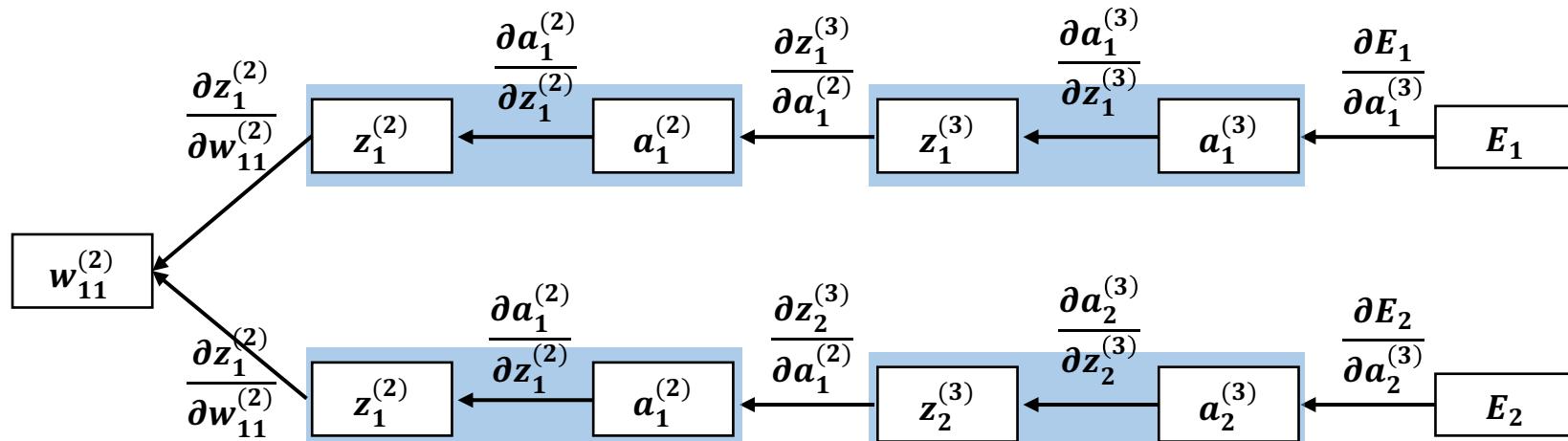
은닉층에서의 오차역전파

$w_{11}^{(2)}$ 에 대한 E 의 변화율 - $\frac{\partial E}{\partial w_{11}^{(2)}}$

$$\frac{\partial E}{\partial W^{(2)}} = \begin{bmatrix} \frac{\partial E}{\partial w_{11}^{(2)}} & \frac{\partial E}{\partial w_{21}^{(2)}} \\ \frac{\partial E}{\partial w_{12}^{(2)}} & \frac{\partial E}{\partial w_{22}^{(2)}} \end{bmatrix} \quad \frac{\partial E}{\partial b^{(2)}} = \begin{bmatrix} \frac{\partial E}{\partial b_1^{(2)}} & \frac{\partial E}{\partial b_2^{(2)}} \end{bmatrix}$$



Back Propagation

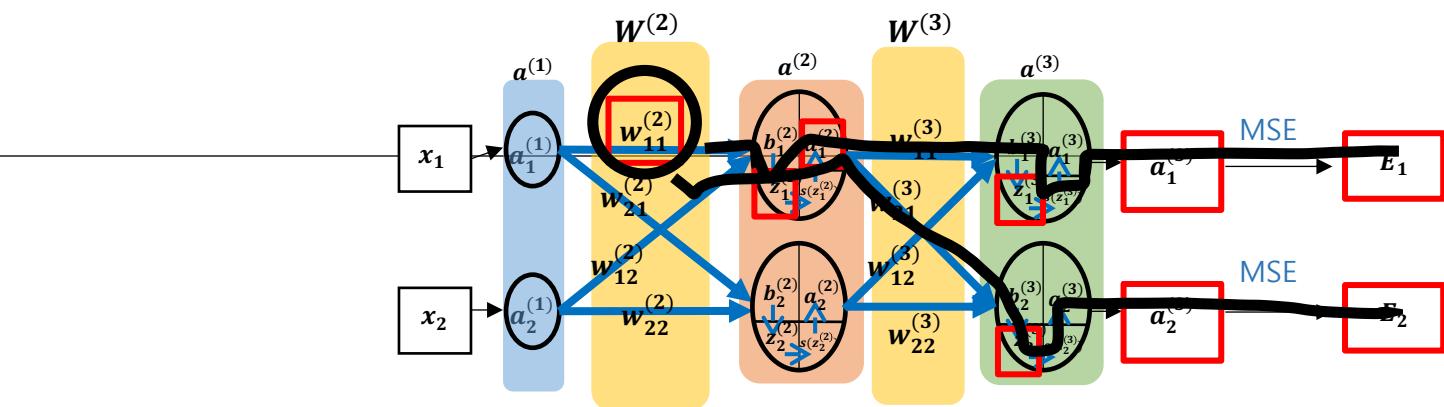


$$\Rightarrow \frac{\partial E}{\partial w_{11}^{(2)}} = \frac{\partial E_1}{\partial w_{11}^{(2)}} + \frac{\partial E_2}{\partial w_{11}^{(2)}} = \frac{\partial E_1}{\partial a_1^{(3)}} \cdot \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} \cdot \frac{\partial z_1^{(3)}}{\partial a_1^{(2)}} \cdot \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \cdot \frac{\partial z_1^{(2)}}{\partial w_{11}^{(2)}} + \frac{\partial E_2}{\partial a_2^{(3)}} \cdot \frac{\partial a_2^{(3)}}{\partial z_2^{(3)}} \cdot \frac{\partial z_2^{(3)}}{\partial a_2^{(2)}} \cdot \frac{\partial a_2^{(2)}}{\partial z_1^{(2)}} \cdot \frac{\partial z_1^{(2)}}{\partial w_{11}^{(2)}}$$

은닉층에서의 오차역전파

$w_{11}^{(2)}$ 에 대한 E 의 변화율 - $\frac{\partial E}{\partial w_{11}^{(2)}}$

$$\begin{aligned} & \frac{\partial E}{\partial w_{11}^{(2)}} = \underline{\frac{\partial E_1}{\partial w_{11}^{(2)}}} + \underline{\frac{\partial E_2}{\partial w_{11}^{(2)}}} = \frac{\partial E_1}{\partial a_1^{(3)}} \cdot \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} \cdot \frac{\partial z_1^{(3)}}{\partial a_1^{(2)}} \cdot \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \cdot \frac{\partial z_1^{(2)}}{\partial w_{11}^{(2)}} + \frac{\partial E_2}{\partial a_2^{(3)}} \cdot \frac{\partial a_2^{(3)}}{\partial z_2^{(3)}} \cdot \frac{\partial z_2^{(3)}}{\partial a_2^{(2)}} \cdot \frac{\partial a_2^{(2)}}{\partial z_1^{(2)}} \cdot \frac{\partial z_1^{(2)}}{\partial w_{11}^{(2)}} \\ &= \underbrace{(a_1^{(3)} - t_1^{(3)}) \cdot S(z_1^{(3)})}_{\text{Jacobian}} \left(1 - S(z_1^{(3)})\right) \cdot w_{11}^{(3)} \cdot S(z_1^{(2)}) \left(1 - S(z_1^{(2)})\right) \cdot a_1^{(1)} \\ &\quad + \underbrace{(a_2^{(3)} - t_2^{(3)}) \cdot S(z_2^{(3)})}_{\text{Jacobian}} \left(1 - S(z_2^{(3)})\right) \cdot w_{21}^{(3)} \cdot S(z_1^{(2)}) \left(1 - S(z_1^{(2)})\right) \cdot a_1^{(1)} \end{aligned}$$

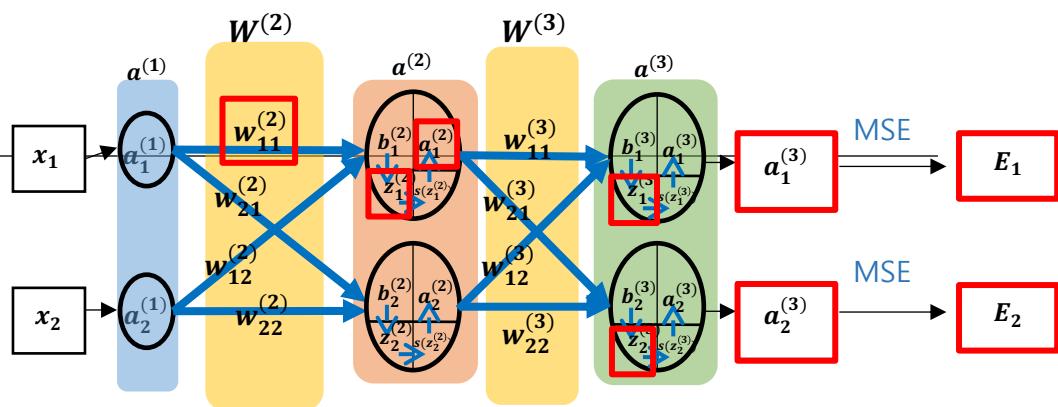


- $\frac{\partial E_1}{\partial a_1^{(3)}} = \frac{\partial \left\{ \frac{1}{2}(t_1^{(3)} - a_1^{(3)})^2 \right\}}{\partial a_1^{(3)}} = (a_1^{(3)} - t_1^{(3)})$
- $\frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} = \frac{\partial S(z_1^{(3)})}{\partial z_1^{(3)}} = S(z_1^{(3)}) \cdot (1 - S(z_1^{(3)}))$
- $\frac{\partial z_1^{(3)}}{\partial a_1^{(2)}} = \frac{\partial (a_1^{(2)}w_{11}^{(3)} + a_2^{(2)}w_{12}^{(3)} + b_1^{(3)})}{\partial a_1^{(2)}} = w_{11}^{(3)}$
- $\frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} = \frac{\partial S(z_1^{(2)})}{\partial z_1^{(2)}} = S(z_1^{(2)}) \cdot (1 - S(z_1^{(2)}))$
- $\frac{\partial z_1^{(2)}}{\partial w_{11}^{(2)}} = \frac{\partial (a_1^{(1)}w_{11}^{(2)} + a_2^{(1)}w_{12}^{(2)})}{\partial w_{11}^{(2)}} = a_1^{(1)}$

- $\frac{\partial E_2}{\partial a_2^{(3)}} = \frac{\partial \left\{ \frac{1}{2}(t_2^{(3)} - a_2^{(3)})^2 \right\}}{\partial a_2^{(3)}} = (a_2^{(3)} - t_2^{(3)})$
- $\frac{\partial a_2^{(3)}}{\partial z_2^{(3)}} = \frac{\partial S(z_2^{(3)})}{\partial z_2^{(3)}} = S(z_2^{(3)}) \cdot (1 - S(z_2^{(3)}))$
- $\frac{\partial z_2^{(3)}}{\partial a_2^{(2)}} = \frac{\partial (a_1^{(2)}w_{21}^{(3)} + a_2^{(2)}w_{22}^{(3)} + b_2^{(3)})}{\partial a_2^{(2)}} = w_{21}^{(3)}$
- $\frac{\partial a_2^{(2)}}{\partial z_1^{(2)}} = \frac{\partial S(z_1^{(2)})}{\partial z_1^{(2)}} = S(z_1^{(2)}) \cdot (1 - S(z_1^{(2)}))$
- $\frac{\partial z_1^{(2)}}{\partial w_{11}^{(2)}} = \frac{\partial (a_1^{(1)}w_{11}^{(2)} + a_2^{(1)}w_{12}^{(2)})}{\partial w_{11}^{(2)}} = a_1^{(1)}$

은닉층에서의 오차역전파

$w_{11}^{(2)}$ 에 대한 E 의 변화율 - $\frac{\partial E}{\partial w_{11}^{(2)}}$



$$\checkmark \Rightarrow \frac{\partial E}{\partial w_{11}^{(2)}} = \frac{\partial E_1}{\partial w_{11}^{(2)}} + \frac{\partial E_2}{\partial w_{11}^{(2)}} = \frac{\partial E_1}{\partial a_1^{(3)}} \cdot \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} \cdot \frac{\partial z_1^{(3)}}{\partial a_1^{(2)}} \cdot \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \cdot \frac{\partial z_1^{(2)}}{\partial w_{11}^{(2)}} + \frac{\partial E_2}{\partial a_2^{(3)}} \cdot \frac{\partial a_2^{(3)}}{\partial z_2^{(3)}} \cdot \frac{\partial z_2^{(3)}}{\partial a_2^{(2)}} \cdot \frac{\partial a_2^{(2)}}{\partial z_1^{(2)}} \cdot \frac{\partial z_1^{(2)}}{\partial w_{11}^{(2)}}$$

$$\begin{aligned} \checkmark &= (a_1^{(3)} - t_1^{(3)}) \cdot S(z_1^{(3)}) (1 - S(z_1^{(3)})) \cdot w_{11}^{(3)} \cdot S(z_1^{(2)}) (1 - S(z_1^{(2)})) \cdot a_1^{(1)} \\ &\quad + (a_2^{(3)} - t_2^{(3)}) \cdot S(z_2^{(3)}) (1 - S(z_2^{(3)})) \cdot w_{21}^{(3)} \cdot S(z_1^{(2)}) (1 - S(z_1^{(2)})) \cdot a_1^{(1)} \end{aligned}$$

$$\begin{aligned} \checkmark &= (a_1^{(3)} - t_1^{(3)}) \cdot a_1^{(3)} (1 - a_1^{(3)}) \cdot w_{11}^{(3)} \cdot a_1^{(2)} (1 - a_1^{(2)}) \cdot a_1^{(1)} \\ &\quad + (a_2^{(3)} - t_2^{(3)}) \cdot a_2^{(3)} (1 - a_2^{(3)}) \cdot w_{21}^{(3)} \cdot a_1^{(2)} (1 - a_1^{(2)}) \cdot a_1^{(1)} \end{aligned}$$

$$\checkmark \therefore \frac{\partial E}{\partial w_{11}^{(2)}} = (a_1^{(3)} - t_1^{(3)}) \cdot a_1^{(3)} (1 - a_1^{(3)}) \cdot w_{11}^{(3)} \cdot a_1^{(2)} (1 - a_1^{(2)}) \cdot a_1^{(1)} + (a_2^{(3)} - t_2^{(3)}) \cdot a_2^{(3)} (1 - a_2^{(3)}) \cdot w_{21}^{(3)} \cdot a_1^{(2)} (1 - a_1^{(2)}) \cdot a_1^{(1)}$$

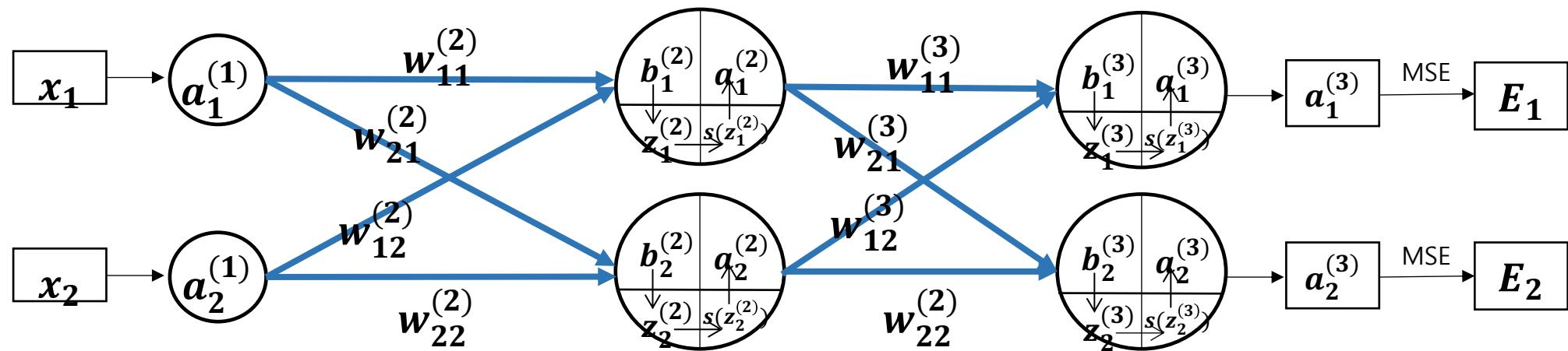
은닉층에서의 오차역전파



Do It Yourself

$$\frac{\partial E}{\partial w_{21}^{(2)}}, \frac{\partial E}{\partial w_{12}^{(2)}}, \frac{\partial E}{\partial w_{22}^{(2)}}, \frac{\partial E}{\partial b_1^{(2)}}, \frac{\partial E}{\partial b_2^{(2)}}$$

$$\frac{\partial E}{\partial W^{(2)}} = \begin{bmatrix} \frac{\partial E}{\partial w_{11}^{(2)}} \\ \frac{\partial E}{\partial w_{21}^{(2)}} \\ \frac{\partial E}{\partial w_{12}^{(2)}} \\ \frac{\partial E}{\partial w_{22}^{(2)}} \end{bmatrix}$$
$$\frac{\partial E}{\partial b^{(2)}} = \begin{bmatrix} \frac{\partial E}{\partial b_1^{(2)}} \\ \frac{\partial E}{\partial b_2^{(2)}} \end{bmatrix}$$



은닉층에서의 오차역전파

은닉층에서의 오차역전파 - 정리

입력층 출력 값 벡터	$A1 = (a_1^{(1)} \quad a_2^{(1)})$
은닉층 출력 값 벡터	$A2 = (a_1^{(2)} \quad a_2^{(2)})$
출력층 가상 손실 벡터	$loss_3 = ((a_1^{(3)} - t_1^{(3)})a_1^{(3)}(1 - a_1^{(3)}) \quad (a_2^{(3)} - t_2^{(3)})a_2^{(3)}(1 - a_2^{(3)}))$
출력층 가중치	$W3 = \begin{pmatrix} w_{11}^{(3)} & w_{21}^{(3)} \\ w_{12}^{(3)} & w_{22}^{(3)} \end{pmatrix}$
은닉층 가상 손실 벡터	$loss_2 = (loss_3 \cdot W3^T) \times A2(1 - A2)$

- $$\frac{\partial E}{\partial W^{(2)}} = \begin{bmatrix} \frac{\partial E}{\partial w_{11}^{(2)}} & \frac{\partial E}{\partial w_{21}^{(2)}} \\ \frac{\partial E}{\partial w_{12}^{(2)}} & \frac{\partial E}{\partial w_{22}^{(2)}} \end{bmatrix} = A1^T \cdot ((loss_3 \cdot W3^T) \times (A2 \times (1 - A2))) = A1^T \cdot loss_2$$

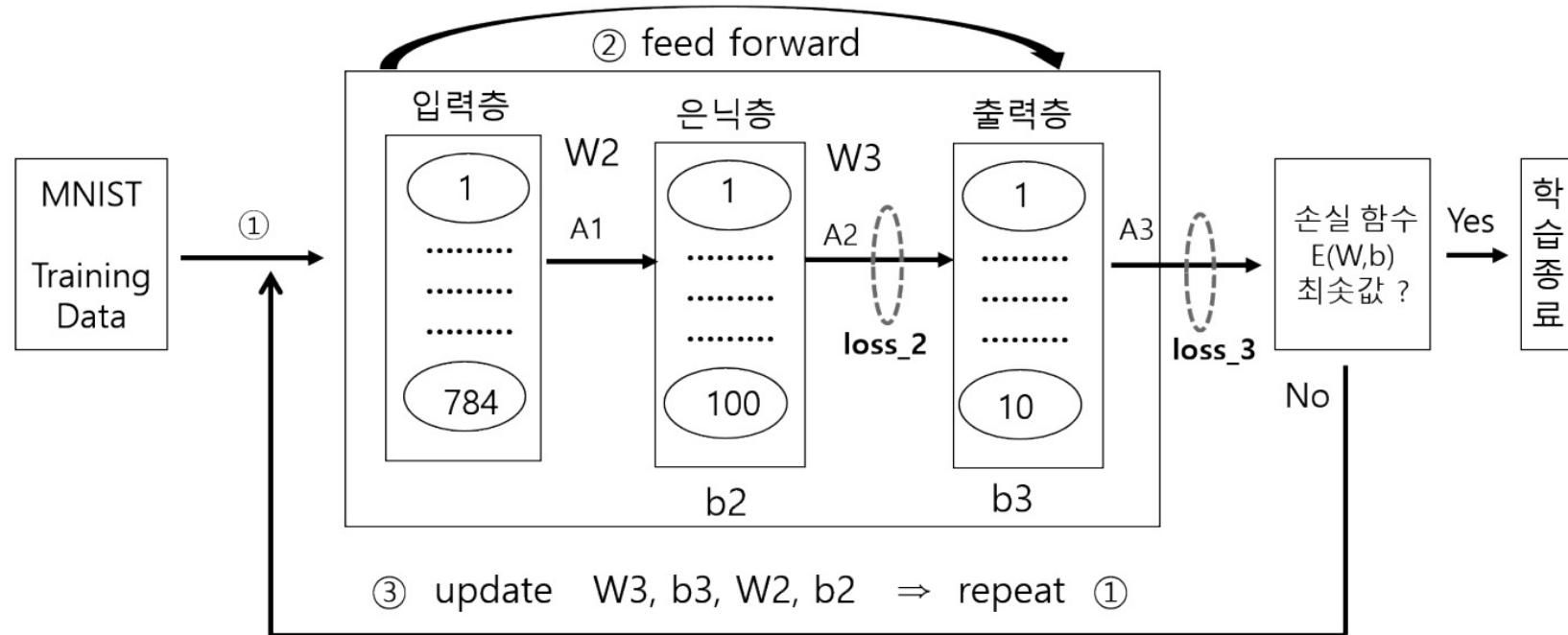
- $$\frac{\partial E}{\partial b^{(2)}} = \begin{bmatrix} \frac{\partial E}{\partial b_1^{(2)}} & \frac{\partial E}{\partial b_2^{(2)}} \end{bmatrix} = ((loss_3 \cdot W3^T) \times (A2 \times (1 - A2))) = loss_2$$

CONTENTS

- 1 오차역전파 개념
- 2 출력층에서의 오차역전파
- 3 은닉층에서의 오차역전파
- 4 오차역전파 이용 MNIST 검증

오차역전파를 이용한 MNIST 분류

① 오차역전파를 이용한 딥러닝 architecture



- $W^{(3)} := W^{(3)} - \alpha (A^{(2)T} \cdot loss_3)$
- $W^{(2)} := W^{(2)} - \alpha (A^{(2)T} \cdot loss_2)$
- $b^{(3)} := b^{(3)} - \alpha \times loss_3$
- $b^{(2)} := b^{(2)} - \alpha \times loss_2$
- $A^{(1)} = [a_1^{(1)} \quad a_2^{(1)}]$
- $A^{(2)} = [a_1^{(2)} \quad a_2^{(2)}]$
- $loss_3 = [(a_1^{(3)} - t_1) a_1^{(3)} (1 - a_1^{(3)}) \quad (a_2^{(3)} - t_2) a_2^{(3)} (1 - a_2^{(3)})]$
- $loss_2 = (loss_3 \cdot W^{(3)T}) \times (A^{(2)} (1 - A^{(2)}))$

오차역전파를 이용한 MNIST 분류

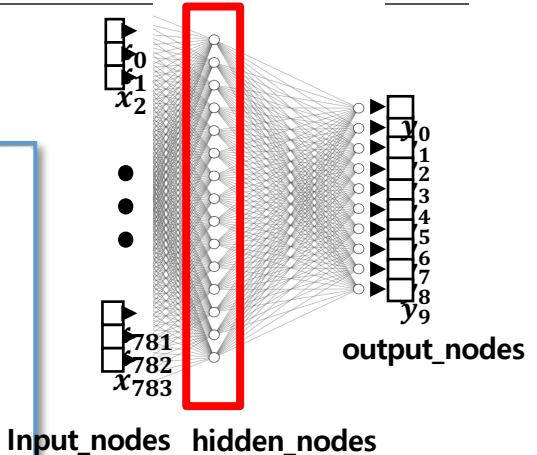
CODE – NeuralNetwork class

```
● ● ●  
1 class NeuralNetwork:  
2     def __init__(self, input_nodes, hidden_nodes, output_nodes, learning_rate):  
3         def feed_forward(self):                      #Feed Forward 수행  
4             def loss_val(self):                      #손실함수 값 계산  
5             def train(self):                         #가중치, 편향 업데이트  
6             def predict(self, input_data):           #미래 값 예측  
7             def accuracy(self, input_data, target_data):#정확도 측정
```

오차역전파를 이용한 MNIST 분류

CODE – NeuralNetwork class – 생성자 `__init__()`

```
● ● ●  
1 class NeuralNetwork:  
2  
3     def __init__(self, input_nodes, hidden_nodes, output_nodes, learning_rate):  
4  
5         self.input_nodes = input_nodes  
6         self.hidden_nodes = hidden_nodes  
7         self.output_nodes = output_nodes  
8  
9         # 은닉층 가중치 W2 = (784 X 100) Xavier 방법으로 self.W2 가중치 초기화  
10        self.W2 = np.random.randn(self.input_nodes, self.hidden_nodes) / np.sqrt(self.input_nodes)  
11        self.b2 = np.random.rand(self.hidden_nodes)  
12  
13        # 출력층 가중치는 W3 = (100X10) Xavier 방법으로 self.W3 가중치 초기화  
14        self.W3 = np.random.randn(self.hidden_nodes, self.output_nodes) / np.sqrt(self.hidden_nodes)  
15        self.b3 = np.random.rand(self.output_nodes)  
16  
17        # 출력층 선형회귀 값 Z3, 출력값 A3 정의 (모두 행렬로 표시)  
18        self.Z3 = np.zeros([1,output_nodes])  
19        self.A3 = np.zeros([1,output_nodes])  
20  
21        # 은닉층 선형회귀 값 Z2, 출력값 A2 정의 (모두 행렬로 표시)  
22        self.Z2 = np.zeros([1,hidden_nodes])  
23        self.A2 = np.zeros([1,hidden_nodes])  
24  
25        # 입력층 선형회귀 값 Z1, 출력값 A1 정의 (모두 행렬로 표시)  
26        self.Z1 = np.zeros([1,input_nodes])  
27        self.A1 = np.zeros([1,input_nodes])  
28  
29        # 학습률 learning_rate 초기화  
30        self.learning_rate = learning_rate
```



Input_nodes hidden_nodes

오차역전파를 이용한 MNIST 분류

CODE – NeuralNetwork class – feed_forward

```
● ● ●  
1 def feed_forward(self):  
2  
3     delta = 1e-7    # log 무한대 발산 방지  
4  
5     # 입력층 선형회귀 값 Z1, 출력값 A1 계산  
6     self.Z1 = self.input_data  
7     self.A1 = self.input_data  
8  
9     # 은닉층 선형회귀 값 Z2, 출력값 A2 계산  
10    self.Z2 = np.dot(self.A1, self.W2) + self.b2  
11    self.A2 = sigmoid(self.Z2)  
12  
13    # 출력층 선형회귀 값 Z3, 출력값 A3 계산  
14    self.Z3 = np.dot(self.A2, self.W3) + self.b3  
15    self.A3 = sigmoid(self.Z3)  
16  
17    return -np.sum(self.target_data*np.log(self.A3 + delta) + (1-self.target_data)*np.log((1 - self.A3)+delta))
```

CrossEntropy

오차역전파를 이용한 MNIST 분류

CODE – NeuralNetwork class – loss_val

```
● ● ●  
1 def loss_val(self):  
2  
3     delta = 1e-7    # log 무한대 발산 방지  
4  
5     # 입력층 선형회귀 값 Z1, 출력값 A1 계산  
6     self.Z1 = self.input_data  
7     self.A1 = self.input_data  
8  
9     # 은닉층 선형회귀 값 Z2, 출력값 A2 계산  
10    self.Z2 = np.dot(self.A1, self.W2) + self.b2  
11    self.A2 = sigmoid(self.Z2)  
12  
13    # 출력층 선형회귀 값 Z3, 출력값 A3 계산  
14    self.Z3 = np.dot(self.A2, self.W3) + self.b3  
15    self.A3 = sigmoid(self.Z3)  
16  
17    return -np.sum(self.target_data*np.log(self.A3 + delta) + (1-self.target_data)*np.log((1-self.A3)+delta))
```

오차역전파를 이용한 MNIST 분류

CODE – NeuralNetwork class – train

```
● ● ●  
1 def train(self, input_data, target_data): # input_data : 784 개, target_data : 10개  
2  
3     self.target_data = target_data  
4     self.input_data = input_data  
5  
6     # 먼저 feed forward 를 통해서 최종 출력값과 이를 바탕으로 현재의 에러 값 계산  
7     loss_val = self.feed_forward()  
8  
9     # 출력층 loss 인 loss_3 구함  
10    loss_3 = (self.A3-self.target_data) * self.A3 * (1-self.A3) •  $loss_3 = [(a_1^{(3)} - t_1) a_1^{(3)} (1 - a_1^{(3)}) \quad (a_2^{(3)} - t_2) a_2^{(3)} (1 - a_2^{(3)})]$   
11  
12    # 출력층 가중치 W3, 출력층 바이어스 b3 업데이트  
13    self.W3 = self.W3 - self.learning_rate * np.dot(self.A2.T, loss_3) •  $W^{(3)} := W^{(3)} - \alpha (A^{(2)T} \cdot loss_3)$   
14  
15    self.b3 = self.b3 - self.learning_rate * loss_3 •  $b^{(3)} := b^{(3)} - \alpha \times loss_3$   
16  
17    # 은닉층 loss 인 loss_2 구함  
18    loss_2 = np.dot(loss_3, self.W3.T) * self.A2 * (1-self.A2) •  $loss_2 = (loss_3 \cdot W^{(3)T}) \times (A^{(2)} (1 - A^{(2)}))$   
19  
20    # 은닉층 가중치 W2, 은닉층 바이어스 b2 업데이트  
21    self.W2 = self.W2 - self.learning_rate * np.dot(self.A1.T, loss_2) •  $W^{(2)} := W^{(2)} - \alpha (A^{(2)T} \cdot loss_2)$   
22  
23    self.b2 = self.b2 - self.learning_rate * loss_2 •  $b^{(2)} := b^{(2)} - \alpha \times loss_2$ 
```

오차역전파를 이용한 MNIST 분류

CODE – NeuralNetwork class – predict



```
1  def predict(self, input_data):          # input_data 는 행렬로 입력됨 즉, (1, 784) shape 을 가짐
2
3      Z2 = np.dot(input_data, self.W2) + self.b2
4      A2 = sigmoid(Z2)
5
6      Z3 = np.dot(A2, self.W3) + self.b3
7      A3 = sigmoid(Z3)
8
9      predicted_num = np.argmax(A3)
10
11     return predicted_num
```

오차역전파를 이용한 MNIST 분류

CODE – NeuralNetwork class – accuracy

```
● ● ●  
1  def accuracy(self, test_data):  
2      matched_list = []  
3      not_matched_list = []  
4  
5      for index in range(len(test_data)):  
6          label = int(test_data[index, 0])  
7          # one-hot encoding을 위한 데이터 정규화 (data normalize)  
8          data = (test_data[index, 1:] / 255.0 * 0.99) + 0.01  
9          # predict 를 위해서 vector 을 matrix 로 변환하여 인수로 넘겨줌  
10         predicted_num = self.predict(np.array(data, ndmin=2))  
11  
12         if label == predicted_num:  
13             matched_list.append(index)  
14         else:  
15             not_matched_list.append(index)  
16  
17         print("Current Accuracy = ", 100*(len(matched_list)/(len(test_data)))), "%")  
18  
19     return matched_list, not_matched_list
```

오차역전파를 이용한 MNIST 분류



CODE – NeuralNetwork class – 학습

```
1 # 0~9 숫자 이미지가 784개의 숫자 (28X28)로 구성되어 있는 training data 읽어옴
2 training_data = np.loadtxt('./mnist_train.csv', delimiter=',', dtype=np.float32)
3
4 # 0~9 숫자 이미지가 784개의 숫자 (28X28)로 구성되어 있는 test data 읽어옴
5 test_data = np.loadtxt('./mnist_test.csv', delimiter=',', dtype=np.float32)
6
7 input_nodes = 784          #input nodes 개수
8 hidden_nodes = 100         #hidden nodes 개수
9 output_nodes = 10          #output nodes 개수
10 learning_rate = 0.3        #learning rate
11 epochs = 1                #반복 횟수
12
13 nn = NeuralNetwork(input_nodes, hidden_nodes, output_nodes, learning_rate)
14
15 start_time = datetime.now()
16
17 for i in range(epochs):
18
19     for step in range(len(training_data)): # train
20
21         # input_data, target_data normalize
22         target_data = np.zeros(output_nodes) + 0.01
23         target_data[int(training_data[step, 0])] = 0.99
24
25         input_data = ((training_data[step, 1:] / 255.0) * 0.99) + 0.01
26
27         nn.train( np.array(input_data, ndmin=2), np.array(target_data, ndmin=2) )
28
29         if step % 400 == 0:
30             print("step = ", step, " loss_val = ", nn.loss_val())
31
32 end_time = datetime.now()
33 print("\nelapsed time = ", end_time - start_time)
```

오차역전파를 이용한 MNIST 분류

결과

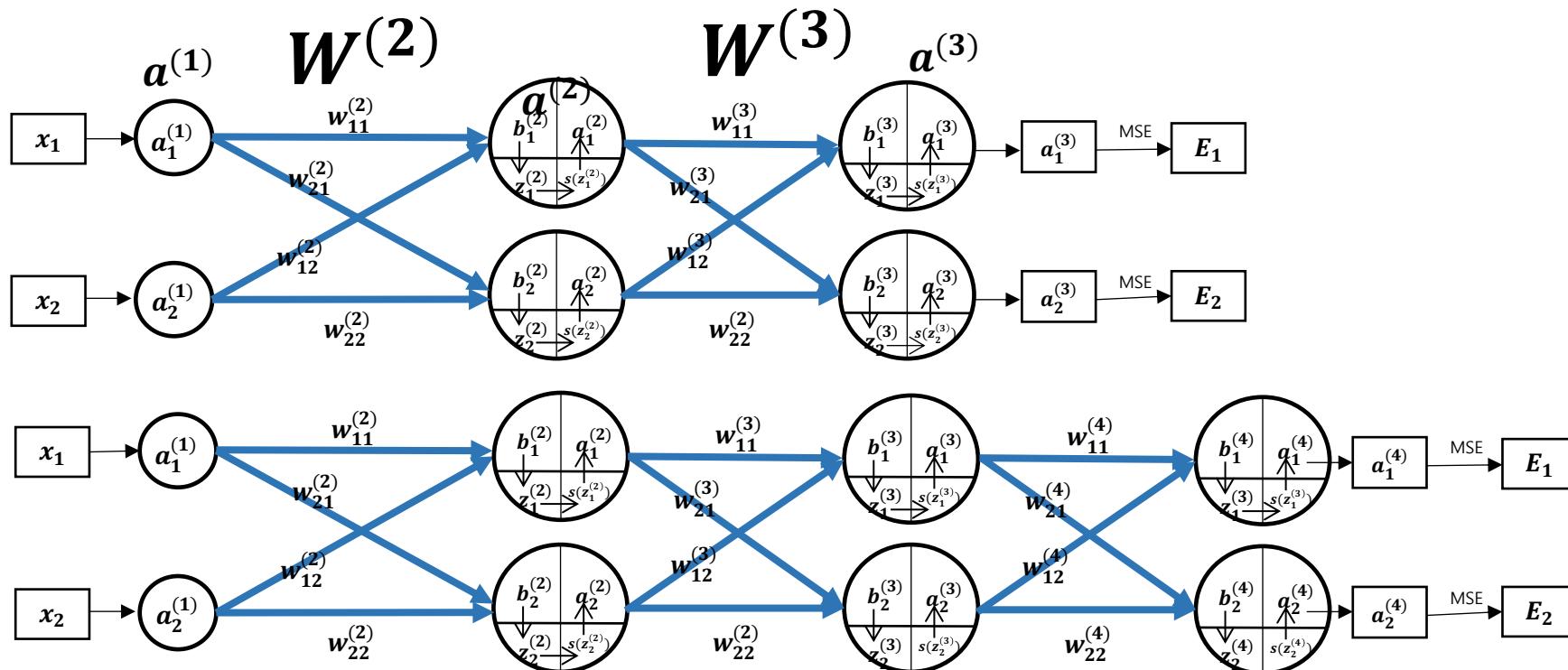
학습 결과	epochs = 0 , step = 0 , current loss_val = 3.596431774902841 epochs = 0 , step = 1000 , current loss_val = 0.8424336448981657 epochs = 0 , step = 2000 , current loss_val = 1.13223305094993 ----- epochs = 0 , step = 57000 , current loss_val = 1.1467274015653924 epochs = 0 , step = 58000 , current loss_val = 0.9788280232573221 epochs = 0 , step = 59000 , current loss_val = 0.935433991132101
검증 코드	test_data = np.loadtxt('./mnist_test.csv', delimiter=',', dtype=np.float32) test_input_data = test_data[:, 1:] test_target_data = test_data[:, 0] (true_list, false_list) = obj.accuracy(test_input_data, test_target_data) ⑤
검증 결과	Current Accuracy = 0.9432

CONTENTS

- 1 오차역전파 개념
- 2 출력층에서의 오차역전파
- 3 은닉층에서의 오차역전파
- 4 오차역전파 이용 MNIST 검증
- 5 Appendix: 은닉층이 여러개일경우?

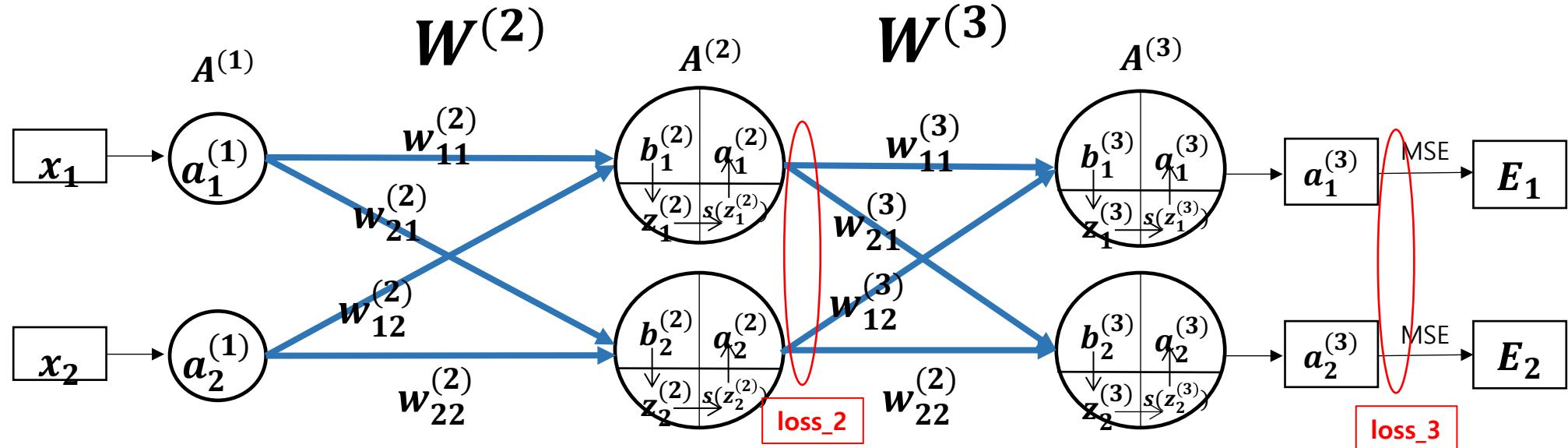
은닉층이 여러 개일 경우의 오차 역전파

은닉층이 한 개 vs 두 개인 FFNN구조 비교



은닉층이 여러 개일 경우의 오차 역전파

은닉층이 한 개인 FFNN의 일반화: 이전층 | 현재층 | 다음층 의 개념



- $A^{(1)} = [a_1^{(1)} \ a_2^{(1)}]$

- $A^{(2)} = [a_1^{(2)} \ a_2^{(2)}]$

- $A^{(3)} = [a_1^{(3)} \ a_2^{(3)}]$

- $\bullet \ loss_2 = (loss_3 \cdot W^{(3)^T}) \times A^{(2)}(1 - A^{(2)})$

- $\bullet \ loss_3 = [(a_1^{(3)} - t_1) a_1^{(3)} (1 - a_1^{(3)}) \quad (a_2^{(3)} - t_2) a_2^{(3)} (1 - a_2^{(3)})]$

- $\bullet \ W^{(2)} := W^{(2)} - \alpha (A^{(2)^T} \cdot loss_2)$

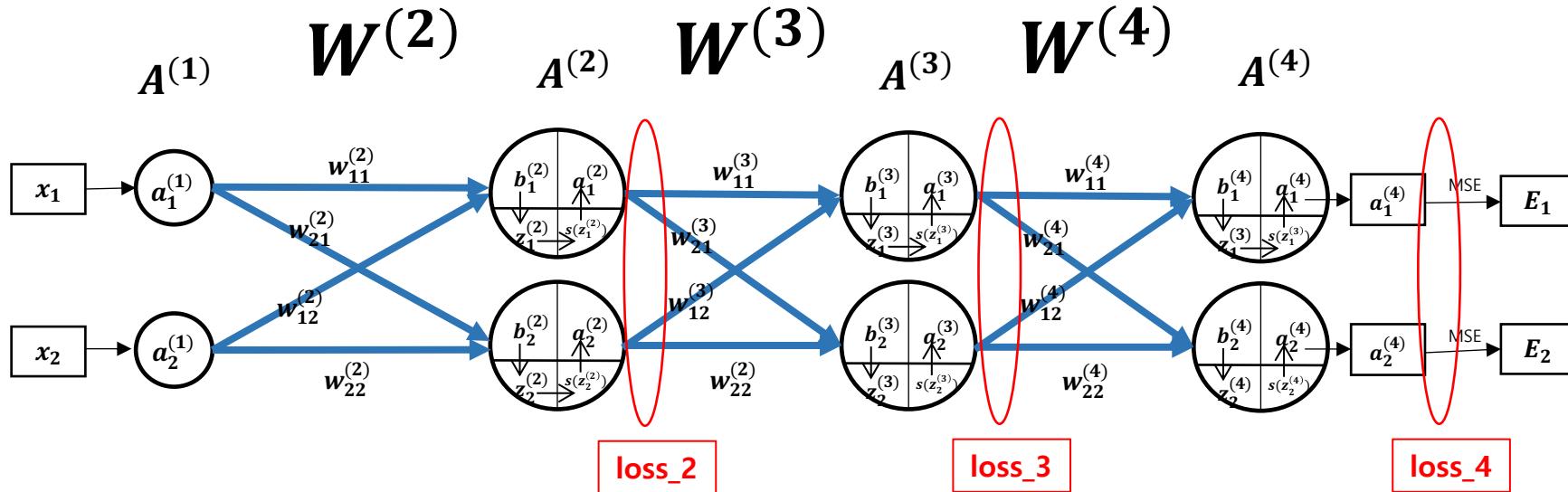
- $\bullet \ b^{(2)} := b^{(2)} - \alpha \times loss_2$

- $\bullet \ W^{(3)} := W^{(3)} - \alpha (A^{(2)^T} \cdot loss_3)$

- $\bullet \ b^{(3)} := b^{(3)} - \alpha \times loss_3$

은닉층이 여러 개일 경우의 오차 역전파

은닉층이 두 개인 FFNN의 일반화: 이전층 | 현재층 | 다음층 의 개념

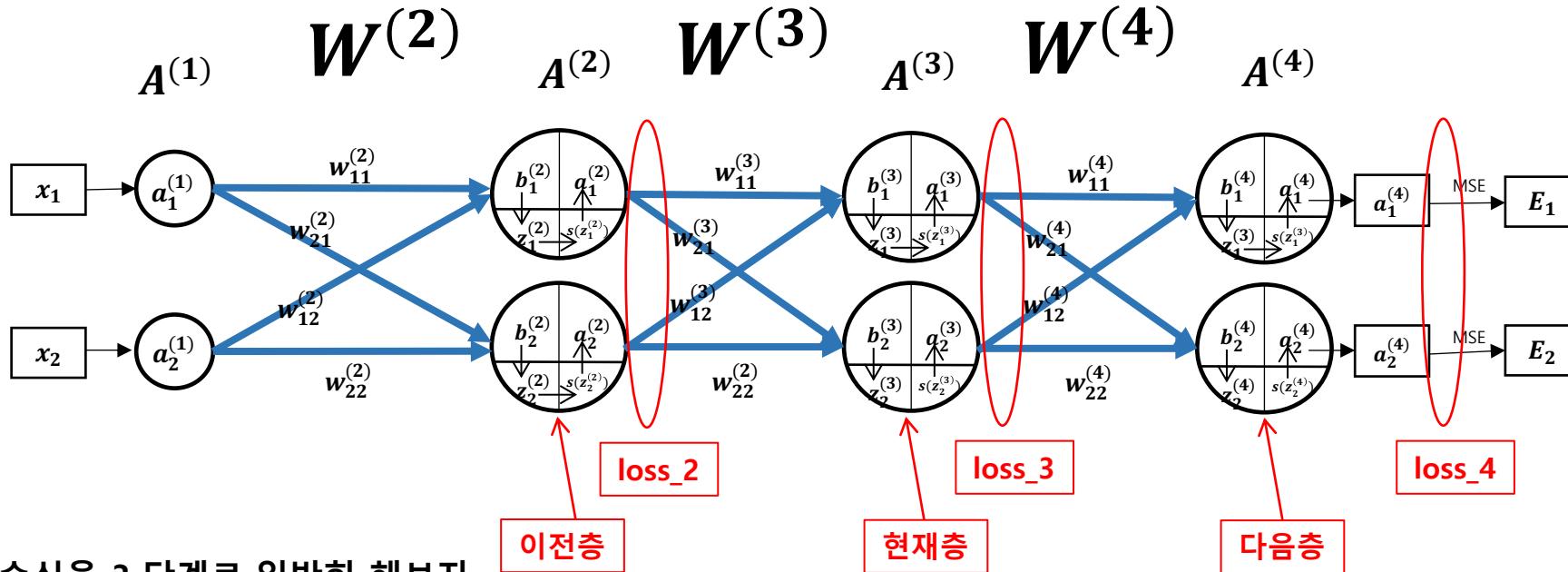


오차 역전파의 수식을 3 단계로 일반화 해보자

- 출력층의 손실계산 일반화: 출력층 손실 = (출력층 출력 - 정답) \times 출력층 출력(1-출력층 출력)
- 은닉층의 손실계산 일반화: 은닉층의 현재손실 = (다음층손실 \cdot 다음층적용 가중치 W^T) * 현재층 출력(1-현재층 출력)
- 현재층의 바이어스 변화율 $\frac{\partial E}{\partial b(\text{현재층})}$ = 현재층 손실, $\frac{\partial E}{\partial W(\text{현재층})}$ 현재층에 적용되는 가중치 변화율 = (이전층 출력) $^t \cdot$ 현재층 손실

은닉층이 여러 개일 경우의 오차 역전파

은닉층이 두 개인 FFNN의 일반화: 이전층 | 현재층 | 다음층 의 개념

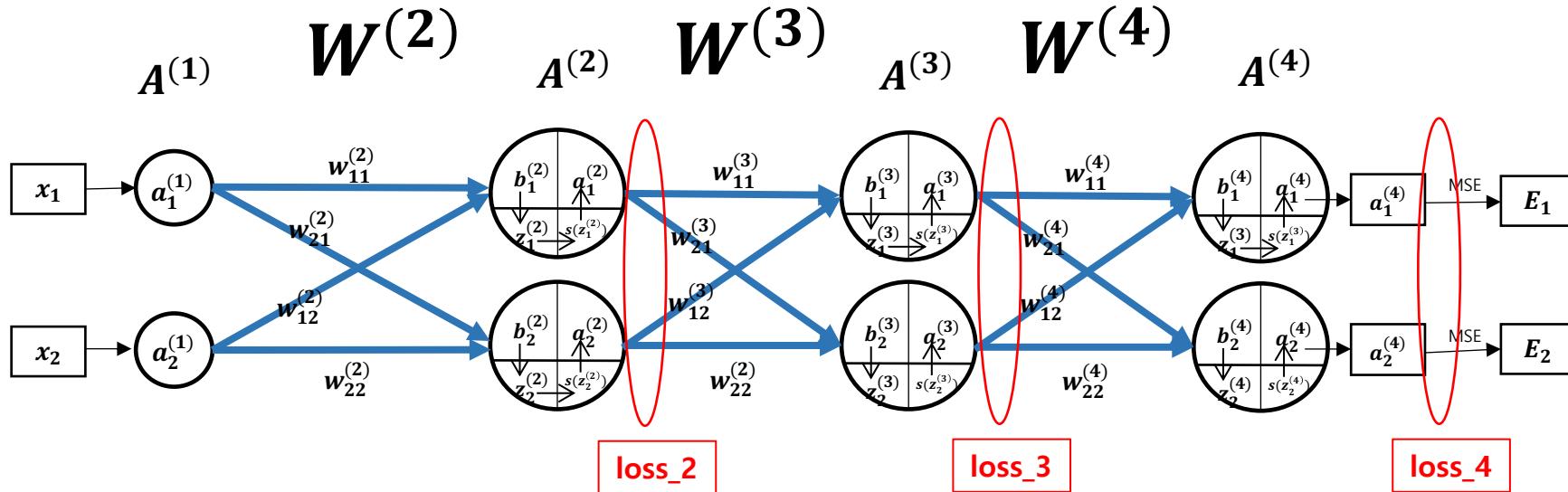


오차 역전파의 수식을 3 단계로 일반화 해보자

- 출력층의 손실계산 일반화: 출력층 손실 = (출력층 출력 - 정답) × 출력층 출력(1 - 출력층 출력) • $loss_4 = (A^{(4)} - Target) \times A^{(4)}(1 - A^{(4)})$
- 은닉층의 손실계산 일반화: 은닉층의 현재손실 = (다음층손실 · 다음층적용 가중치 W^T) * 현재층 출력(1 - 현재층 출력)
 - $loss_2 = [(loss_3 \cdot W^{(3)^T}) \times A^{(2)}(1 - A^{(2)})]$ • $loss_3 = [(loss_4 \cdot W^{(4)^T}) \times A^{(3)}(1 - A^{(3)})]$
- 현재층의 바이어스 변화율 $\frac{\partial E}{\partial b(\text{현재층})} = \text{현재층 손실}, \frac{\partial E}{\partial W(\text{현재층})}$ 현재층에 적용되는 가중치 변화율 = (이전층 출력) t · 현재층 손실
 - $b^{(2)} := b^{(2)} - \alpha \times loss_2$ • $W^{(2)} := W^{(2)} - \alpha (A^{(1)^T} \cdot loss_2)$

은닉층이 여러 개일 경우의 오차 역전파

은닉층이 두 개인 FFNN의 일반화: 이전층 | 현재층 | 다음층 의 개념

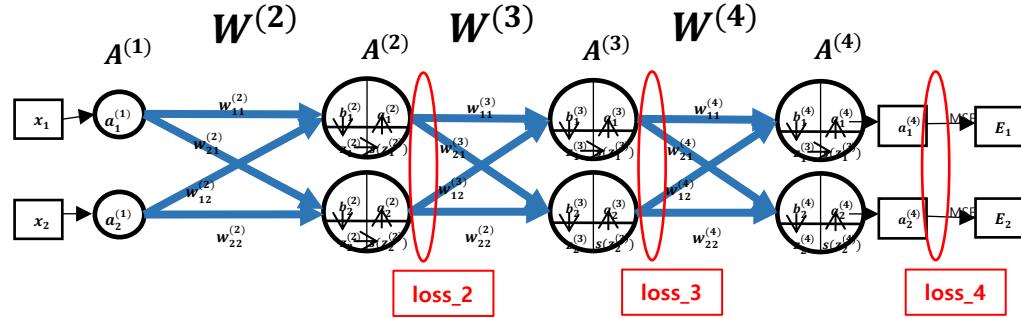


정리

- $W^{(4)} := W^{(4)} - \alpha (A^{(3)^T} \cdot loss_4)$
- $W^{(3)} := W^{(3)} - \alpha (A^{(2)^T} \cdot loss_3)$
- $W^{(2)} := W^{(2)} - \alpha (A^{(1)^T} \cdot loss_2)$
- $b^{(4)} := b^{(4)} - \alpha \times loss_4$
- $b^{(3)} := b^{(3)} - \alpha \times loss_3$
- $b^{(2)} := b^{(2)} - \alpha \times loss_2$
- $loss_4 = (A^{(4)} - Target) \times A^{(4)}(1 - A^{(4)})$
- $loss_3 = [(loss_4 \cdot W^{(4)^T}) \times A^{(3)}(1 - A^{(3)})]$
- $loss_2 = [(loss_3 \cdot W^{(3)^T}) \times A^{(2)}(1 - A^{(2)})]$

은닉층이 여러 개일 경우의 오차 역전파

은닉층이 두 개인 FFNN의 일반화 구현



정리

- $W^{(4)} := W^{(4)} - \alpha (A^{(3)^T} \cdot loss_4)$
- $b^{(4)} := b^{(4)} - \alpha \times loss_4$
- $loss_4 = (A^{(4)} - Target) \times A^{(4)}(1 - A^{(4)})$
- $W^{(3)} := W^{(3)} - \alpha (A^{(2)^T} \cdot loss_3)$
- $b^{(3)} := b^{(3)} - \alpha \times loss_3$
- $loss_3 = [(loss_4 \cdot W^{(4)^T}) \times A^{(3)}(1 - A^{(3)})]$
- $W^{(2)} := W^{(2)} - \alpha (A^{(1)^T} \cdot loss_2)$
- $b^{(2)} := b^{(2)} - \alpha \times loss_2$
- $loss_2 = [(loss_3 \cdot W^{(3)^T}) \times A^{(2)}(1 - A^{(2)})]$

어떻게 구현할까 고민해 봅시다!



감사합니다.