

Lecture4 142p – Backpropagation

박주찬

DICE Lab
School of Computer Science and Engineering
KOREATECH
green261535@gmail.com

Introduction

CS231n Lecture4 142pg에 나오는 computational graph에서 W_1, W_2 에 대하여 backpropagation을 증명한다.

Computation graph는 다음 Fig. 1. 과 같다.

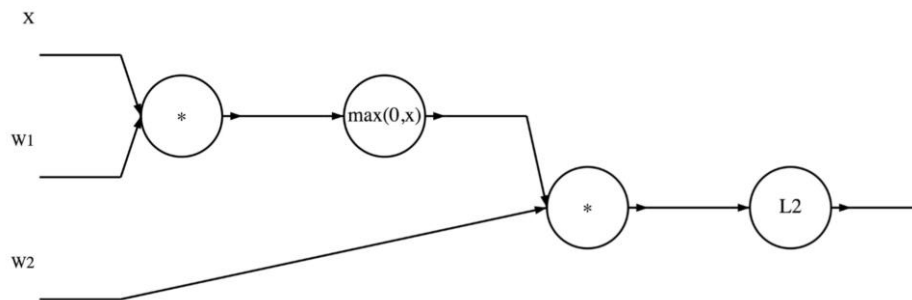


Fig. 1.

Fig. 2. Computational graph의 수식은 다음과 같다.

$$z = X * W_1$$

$$h_1 = \text{ReLU}(z)$$

$$y = h_1 * W_2$$

$$L = \|y\|^2$$

$X, z, W_1, W_2, h_1, y, L$ 을 다음과 같이 정의 한다.

$$X = \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{bmatrix}, W_1 = \begin{bmatrix} w_{11} \\ w_{12} \end{bmatrix}, z = \begin{bmatrix} x_{11} * w_{11} + x_{12} * w_{12} \\ x_{21} * w_{11} + x_{22} * w_{12} \end{bmatrix} = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$$

$$h = \begin{bmatrix} \max(0, z_1) \\ \max(0, z_2) \end{bmatrix} = \begin{bmatrix} h_1 \\ h_2 \end{bmatrix}, W_2 = \begin{bmatrix} w_{21} & w_{22} \end{bmatrix}$$

$$y = \begin{bmatrix} \max(0, z_1) * w_{21} & \max(0, z_1) * w_{22} \\ \max(0, z_2) * w_{21} & \max(0, z_2) * w_{22} \end{bmatrix} = \begin{bmatrix} h_1 * w_{21} & h_1 * w_{22} \\ h_2 * w_{21} & h_2 * w_{22} \end{bmatrix} = \begin{bmatrix} y_{11} & y_{12} \\ y_{21} & y_{22} \end{bmatrix}$$

$$L = \|y\|^2$$

Derivation for $\frac{\partial L}{\partial w_2}$

$\frac{\partial L}{\partial w_2}$ 를 구하려면 먼저 곱셈 연산에 대한 upstream과 local을 구해야 한다.

upstream은 다음과 같다.

$$\begin{aligned}\frac{\partial L}{\partial y} &= 2 \times y \\ &= \begin{bmatrix} 2 \times y_{11} & 2 \times y_{12} \\ 2 \times y_{21} & 2 \times y_{22} \end{bmatrix}\end{aligned}\tag{1}$$

w_2 의 각 원소에 대한 local은 다음과 같다.

$$\frac{\partial y}{\partial w_{21}} = \begin{bmatrix} h_1 \\ h_2 \end{bmatrix}, \quad \frac{\partial y}{\partial w_{22}} = \begin{bmatrix} h_1 \\ h_2 \end{bmatrix}\tag{2}$$

upstream(1)과 local(2)을 곱해서 downstream을 구한다.

$$\begin{aligned}\frac{\partial L}{\partial w_{21}} &= \sum_{i=1}^2 \sum_{j=1}^2 \frac{\partial L}{\partial y_{ij}} * \frac{\partial y_{ij}}{\partial w_{21}} \\ &= 2 * y_{11} * h_1 + 2 * y_{21} * h_2\end{aligned}\tag{3}$$

$$\begin{aligned}\frac{\partial L}{\partial w_{22}} &= \sum_{i=1}^2 \sum_{j=1}^2 \frac{\partial L}{\partial y_{ij}} * \frac{\partial y_{ij}}{\partial w_{22}} \\ &= 2 * y_{12} * h_1 + 2 * y_{22} * h_2\end{aligned}\tag{4}$$

$$\begin{aligned}\frac{\partial L}{\partial w_2} &= \frac{\partial L}{\partial L} * \frac{\partial L}{\partial y} * \frac{\partial y}{\partial w_2} \\ &= \begin{bmatrix} \frac{\partial L}{\partial w_{21}} & \frac{\partial L}{\partial w_{22}} \end{bmatrix} \\ &= [2 * y_{11} * h_1 + 2 * y_{21} * h_2 \quad 2 * y_{12} * h_1 + 2 * y_{22} * h_2] \\ &= 2 * \begin{bmatrix} h_1 & h_2 \end{bmatrix} * \begin{bmatrix} y_{11} & y_{12} \\ y_{21} & y_{22} \end{bmatrix} \\ &= 2 * h^T * y\end{aligned}\tag{5}$$

$$\frac{\partial L}{\partial w_2} = 2 * h^T * y \quad (6)$$

Derivation for $\frac{\partial L}{\partial w_1}$

$\frac{\partial L}{\partial w_1}$ 를 구해주기 위해 먼저 $\frac{\partial L}{\partial h}$ 를 구해준다. 방식은 upstream을 구하고 local을 구해서 서로 곱한 후 downstream을 구하는 방식으로 한다.

upstream은 (1)과 같다.

h_1, h_2 에 대한 local은 다음과 같다.

$$\begin{aligned} \frac{\partial y_{11}}{\partial h_1} &= w_{21}, \quad \frac{\partial y_{12}}{\partial h_1} = w_{22}, \quad \frac{\partial y_{21}}{\partial h_1} = 0, \quad \frac{\partial y_{22}}{\partial h_1} = 0, \quad \frac{\partial y_{11}}{\partial h_2} = 0, \quad \frac{\partial y_{12}}{\partial h_2} = 0, \\ \frac{\partial y_{21}}{\partial h_2} &= w_{21}, \quad \frac{\partial y_{22}}{\partial h_2} = w_{22} \end{aligned} \quad (7)$$

upstream(1)과 local(7)을 곱해서 downstream을 구한다.

$$\begin{aligned} \frac{\partial L}{\partial h_1} &= \sum_{i=1}^2 \sum_{j=1}^2 \frac{\partial L}{\partial y_{ij}} * \frac{\partial y_{ij}}{\partial h_1} \\ &= 2 * y_{11} * w_{21} + 2 * y_{12} * w_{22} \end{aligned} \quad (8)$$

$$\begin{aligned} \frac{\partial L}{\partial h_2} &= \sum_{i=1}^2 \sum_{j=1}^2 \frac{\partial L}{\partial y_{ij}} * \frac{\partial y_{ij}}{\partial h_2} \\ &= 2 * y_{21} * w_{21} + 2 * y_{22} * w_{22} \end{aligned} \quad (9)$$

$$\begin{aligned} \frac{\partial L}{\partial h} &= \frac{\partial L}{\partial h_1} * \frac{\partial h_1}{\partial h} + \frac{\partial L}{\partial h_2} * \frac{\partial h_2}{\partial h} \\ &= \begin{bmatrix} \frac{\partial L}{\partial h_1} \\ \frac{\partial L}{\partial h_2} \end{bmatrix} \\ &= \begin{bmatrix} 2 * y_{11} * w_{21} + 2 * y_{12} * w_{22} \\ 2 * y_{21} * w_{21} + 2 * y_{22} * w_{22} \end{bmatrix} \end{aligned} \quad (10)$$

$$\begin{aligned}
\frac{\partial L}{\partial \mathbf{h}} &= 2 * \begin{bmatrix} y_{11} & y_{12} \\ y_{21} & y_{22} \end{bmatrix} * \begin{bmatrix} w_{21} \\ w_{22} \end{bmatrix} \\
&= 2 * \mathbf{y} * \mathbf{w}_2^T
\end{aligned} \tag{11}$$

$\frac{\partial L}{\partial w_1}$ 를 구해주기 위해 $\frac{\partial h}{\partial z}$ 를 구해준다. 방식은 upstream을 구하고 local을 구해서 서로 곱한 후 downstream을 구하는 방식으로 한다

upstream은 (11)과 같다.

local은 다음과 같다.

ReLU(z)를 z에 대하여 미분을 하면 다음과 같이 두가지 경우가 나온다.

$$h_1(z) = \text{ReLU}(z) = \max(0, z) = \begin{cases} z & z > 0 \\ 0 & z \leq 0 \end{cases}$$

if) $z_1 > 0, z_2 > 0$ 일 때, $\frac{\partial h_1}{\partial z}, \frac{\partial h_2}{\partial z}$ 는 다음과 같다.

$$\frac{\partial h_1}{\partial z} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad \frac{\partial h_2}{\partial z} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \tag{12}$$

if) $z_1 > 0, z_2 \leq 0$ 일 때, $\frac{\partial h_1}{\partial z}, \frac{\partial h_2}{\partial z}$ 는 다음과 같다.

$$\frac{\partial h_1}{\partial z} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad \frac{\partial h_2}{\partial z} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \tag{13}$$

if) $z_1 \leq 0, z_2 > 0$ 일 때, $\frac{\partial h_1}{\partial z}, \frac{\partial h_2}{\partial z}$ 는 다음과 같다.

$$\frac{\partial h_1}{\partial z} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \quad \frac{\partial h_2}{\partial z} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \tag{14}$$

if) $z_1 \leq 0, z_2 \leq 0$ 일 때, $\frac{\partial h_1}{\partial z}, \frac{\partial h_2}{\partial z}$ 는 다음과 같다

$$\frac{\partial h_1}{\partial z} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \quad \frac{\partial h_2}{\partial z} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \tag{15}$$

위 조건들(12), (13), (14), (15)을 합해서 간단히 보면 다음과 같다

$$\begin{aligned} \frac{\partial h_1}{\partial z} = h'_1(z) &= \begin{bmatrix} \frac{\partial h_1}{\partial z_1} \\ \frac{\partial h_1}{\partial z_2} \end{bmatrix} = \begin{cases} \begin{bmatrix} 1 \\ 0 \end{bmatrix} & z_1 > 0, z_2 > 0 \\ \begin{bmatrix} 1 \\ 0 \end{bmatrix} & z_1 > 0, z_2 \leq 0 \\ \begin{bmatrix} 0 \\ 0 \end{bmatrix} & z_1 \leq 0, z_2 > 0 \\ \begin{bmatrix} 0 \\ 0 \end{bmatrix} & z_1 \leq 0, z_2 \leq 0 \end{cases} \\ \frac{\partial h_2}{\partial z} = h'_2(z) &= \begin{bmatrix} \frac{\partial h_2}{\partial z_1} \\ \frac{\partial h_2}{\partial z_2} \end{bmatrix} = \begin{cases} \begin{bmatrix} 0 \\ 1 \end{bmatrix} & z_1 > 0, z_2 > 0 \\ \begin{bmatrix} 0 \\ 0 \end{bmatrix} & z_1 > 0, z_2 \leq 0 \\ \begin{bmatrix} 0 \\ 1 \end{bmatrix} & z_1 \leq 0, z_2 > 0 \\ \begin{bmatrix} 0 \\ 0 \end{bmatrix} & z_1 \leq 0, z_2 \leq 0 \end{cases} \end{aligned} \quad (16)$$

$$\frac{\partial h_1}{\partial z} = h'_1(z), \quad \frac{\partial h_2}{\partial z} = h'_2(z) \quad (17)$$

upstream(10)과 local(17)을 곱해서 downstream을 구한다.

$$\begin{aligned} \frac{\partial L}{\partial z} &= \frac{\partial L}{\partial h_1} * \frac{\partial h_1}{\partial z} + \frac{\partial L}{\partial h_2} * \frac{\partial h_2}{\partial z} \\ &= (2 * y_{11} * w_{21} + 2 * y_{12} * w_{22}) * h'_1(z), \\ &\quad + (2 * y_{21} * w_{21} + 2 * y_{22} * w_{22}) * h'_2(z) \end{aligned} \quad (18)$$

위 downstream식을 z_1, z_2 두가지에 대하여 모든 경우에 수를 확인해보면 아래와 같다.

if) $z_1 > 0, z_2 > 0$ 일 때, $\frac{\partial L}{\partial z}$ 는 다음과 같다.

$$\begin{aligned} \frac{\partial L}{\partial z} &= \begin{bmatrix} (2 * y_{11} * w_{21} + 2 * y_{12} * w_{22}) \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ (2 * y_{21} * w_{21} + 2 * y_{22} * w_{22}) \end{bmatrix} \\ &= \begin{bmatrix} (2 * y_{11} * w_{21} + 2 * y_{12} * w_{22}) \\ (2 * y_{21} * w_{21} + 2 * y_{22} * w_{22}) \end{bmatrix} \end{aligned} \quad (19)$$

if) $z_1 > 0, z_2 \leq 0$ 일 때, $\frac{\partial L}{\partial z}$ 는 다음과 같다.

$$\begin{aligned} \frac{\partial L}{\partial z} &= \begin{bmatrix} (2 * y_{11} * w_{21} + 2 * y_{12} * w_{22}) \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \end{bmatrix} \\ &= \begin{bmatrix} (2 * y_{11} * w_{21} + 2 * y_{12} * w_{22}) \\ 0 \end{bmatrix} \end{aligned} \quad (20)$$

if) $z_1 \leq 0, z_2 > 0$ 일 때, $\frac{\partial L}{\partial z}$ 는 다음과 같다.

$$\begin{aligned}\frac{\partial L}{\partial z} &= \begin{bmatrix} 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ (2 * y_{21} * w_{21} + 2 * y_{22} * w_{22}) \end{bmatrix} \\ &= \begin{bmatrix} 0 \\ (2 * y_{21} * w_{21} + 2 * y_{22} * w_{22}) \end{bmatrix}\end{aligned}\quad (21)$$

if) $z_1 \leq 0, z_2 \leq 0$ 일 때, $\frac{\partial L}{\partial z}$ 는 다음과 같다.

$$\begin{aligned}\frac{\partial L}{\partial z} &= \begin{bmatrix} 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \end{bmatrix} \\ &= \begin{bmatrix} 0 \\ 0 \end{bmatrix}\end{aligned}\quad (22)$$

$\frac{\partial L}{\partial w_1}$ 를 구해주기 위해 $\frac{\partial L}{\partial z}$ 를 구해줬다. 이제 마지막으로 $z = X * W_1$ 연산에 대한 upstream을 구하고 local을 구해서 서로 곱한 후 downstream을 구한다.

upstream은 (19), (20), (21), (22)와 같다.

local은 다음과 같다.

$$\frac{\partial z_1}{\partial w_{11}} = x_{11}, \quad \frac{\partial z_2}{\partial w_{11}} = x_{21}, \quad \frac{\partial z_1}{\partial w_{12}} = x_{12}, \quad \frac{\partial z_2}{\partial w_{12}} = x_{22} \quad (23)$$

upstream(19), (20), (21), (22)과 local(23)을 곱해서 downstream을 구한다.

먼저 W_1 의 각 원소 별로 $\frac{\partial L}{\partial w_{11}}, \frac{\partial L}{\partial w_{12}}$ 를 구해준다.

if) $z_1 > 0, z_2 > 0$ 일 때, $\frac{\partial L}{\partial w_{11}}$ 는 다음과 같다.

$$\begin{aligned}\frac{\partial L}{\partial w_{11}} &= \sum_{i=1}^2 \frac{\partial L}{\partial z_i} * \frac{\partial z_i}{\partial w_{11}} \\ &= (2 * y_{11} * w_{21} + 2 * y_{12} * w_{22}) * x_{11} \\ &\quad + (2 * y_{21} * w_{21} + 2 * y_{22} * w_{22}) * x_{21}\end{aligned}\quad (24)$$

if) $z_1 > 0, z_2 > 0$ 일 때, $\frac{\partial L}{\partial w_{12}}$ 는 다음과 같다.

$$\begin{aligned}
\frac{\partial L}{\partial w_{12}} &= \sum_{i=1}^2 \frac{\partial L}{\partial z_i} * \frac{\partial z_i}{\partial w_{12}} \\
&= (2 * y_{11} * w_{21} + 2 * y_{12} * w_{22}) * x_{12} \\
&\quad + (2 * y_{21} * w_{21} + 2 * y_{22} * w_{22}) * x_{22}
\end{aligned} \tag{25}$$

if) $z_1 > 0, z_2 \leq 0$ 일 때, $\frac{\partial L}{\partial w_{11}}$ 는 다음과 같다.

$$\begin{aligned}
\frac{\partial L}{\partial w_{11}} &= \sum_{i=1}^2 \frac{\partial L}{\partial z_i} * \frac{\partial z_i}{\partial w_{11}} \\
&= (2 * y_{11} * w_{21} + 2 * y_{12} * w_{22}) * x_{11} + 0 * x_{21} \\
&= (2 * y_{11} * w_{21} + 2 * y_{12} * w_{22}) * x_{11}
\end{aligned} \tag{26}$$

if) $z_1 > 0, z_2 \leq 0$ 일 때, $\frac{\partial L}{\partial w_{12}}$ 는 다음과 같다.

$$\begin{aligned}
\frac{\partial L}{\partial w_{12}} &= \sum_{i=1}^2 \frac{\partial L}{\partial z_i} * \frac{\partial z_i}{\partial w_{12}} \\
&= (2 * y_{11} * w_{21} + 2 * y_{12} * w_{22}) * x_{12} + 0 * x_{22} \\
&= (2 * y_{11} * w_{21} + 2 * y_{12} * w_{22}) * x_{12}
\end{aligned} \tag{27}$$

if) $z_1 \leq 0, z_2 > 0$ 일 때, $\frac{\partial L}{\partial w_{11}}$ 는 다음과 같다.

$$\begin{aligned}
\frac{\partial L}{\partial w_{11}} &= \sum_{i=1}^2 \frac{\partial L}{\partial z_i} * \frac{\partial z_i}{\partial w_{11}} \\
&= 0 * x_{11} + (2 * y_{21} * w_{21} + 2 * y_{22} * w_{22}) * x_{21} \\
&= (2 * y_{21} * w_{21} + 2 * y_{22} * w_{22}) * x_{21}
\end{aligned} \tag{28}$$

if) $z_1 \leq 0, z_2 > 0$ 일 때, $\frac{\partial L}{\partial w_{12}}$ 는 다음과 같다.

$$\begin{aligned}
\frac{\partial L}{\partial w_{12}} &= \sum_{i=1}^2 \frac{\partial L}{\partial z_i} * \frac{\partial z_i}{\partial w_{12}} \\
&= 0 * x_{12} + (2 * y_{21} * w_{21} + 2 * y_{22} * w_{22}) * x_{22} \\
&= (2 * y_{21} * w_{21} + 2 * y_{22} * w_{22}) * x_{22}
\end{aligned} \tag{29}$$

if) $z_1 \leq 0, z_2 \leq 0$ 일 때, $\frac{\partial L}{\partial w_{11}}$ 는 다음과 같다.

$$\begin{aligned}\frac{\partial L}{\partial w_{11}} &= \sum_{i=1}^2 \frac{\partial L}{\partial z_i} * \frac{\partial z_i}{\partial w_{11}} \\ &= 0 * x_{11} + 0 * x_{21} \\ &= 0\end{aligned}\tag{30}$$

if) $z_1 \leq 0, z_2 \leq 0$ 일 때, $\frac{\partial L}{\partial w_{12}}$ 는 다음과 같다.

$$\begin{aligned}\frac{\partial L}{\partial w_{12}} &= \sum_{i=1}^2 \frac{\partial L}{\partial z_i} * \frac{\partial z_i}{\partial w_{12}} \\ &= 0 * x_{12} + 0 * x_{22} \\ &= 0\end{aligned}\tag{31}$$

if) $z_1 > 0, z_2 > 0$ 일 때, $\frac{\partial L}{\partial w_1}$ 는 다음과 같다. (24), (25)사용

$$\begin{aligned}\frac{\partial L}{\partial w_1} &= \frac{\partial L}{\partial L} * \frac{\partial L}{\partial y} * \frac{\partial y}{\partial h} * \frac{\partial h}{\partial z} * \frac{\partial z}{\partial w_1} \\ &= \begin{bmatrix} \frac{\partial L}{\partial w_{11}} \\ \frac{\partial L}{\partial w_{12}} \end{bmatrix} \\ &= \begin{bmatrix} (2 * y_{11} * w_{21} + 2 * y_{12} * w_{22}) * x_{11} + (2 * y_{21} * w_{21} + 2 * y_{22} * w_{22}) \\ (2 * y_{11} * w_{21} + 2 * y_{12} * w_{22}) * x_{12} + (2 * y_{21} * w_{21} + 2 * y_{22} * w_{22}) \end{bmatrix} \\ &= 2 * \begin{bmatrix} x_{11} & x_{21} \\ x_{12} & x_{22} \end{bmatrix} * \begin{bmatrix} y_{11} & y_{12} \\ y_{21} & y_{22} \end{bmatrix} * \begin{bmatrix} w_{21} \\ w_{22} \end{bmatrix} \\ &= 2 * X^T * y * W_2^T\end{aligned}\tag{32}$$

if) $z_1 > 0, z_2 \leq 0$ 일 때, $\frac{\partial L}{\partial w_1}$ 는 다음과 같다. (26), (27)사용

$$\begin{aligned}
\frac{\partial L}{\partial w_1} &= \frac{\partial L}{\partial L} * \frac{\partial L}{\partial y} * \frac{\partial y}{\partial h} * \frac{\partial h}{\partial z} * \frac{\partial z}{\partial w_1} \\
&= \begin{bmatrix} \frac{\partial L}{\partial w_{11}} \\ \frac{\partial L}{\partial w_{12}} \end{bmatrix} \\
&= \begin{bmatrix} (2 * y_{11} * w_{21} + 2 * y_{12} * w_{22}) * x_{11} \\ (2 * y_{11} * w_{21} + 2 * y_{12} * w_{22}) * x_{12} \end{bmatrix} \\
&= 2 * \begin{bmatrix} x_{11} \\ x_{12} \end{bmatrix} * \begin{bmatrix} y_{11} & y_{12} \end{bmatrix} * \begin{bmatrix} w_{21} \\ w_{22} \end{bmatrix} \\
&= 2 * (X^T \text{의 첫번째 column}) * (y \text{의 첫번째 row}) * W_2^T
\end{aligned} \tag{33}$$

if) $z_1 \leq 0, z_2 > 0$ 일 때, $\frac{\partial L}{\partial w_1}$ 는 다음과 같다. (28), (29)사용

$$\begin{aligned}
\frac{\partial L}{\partial w_1} &= \frac{\partial L}{\partial L} * \frac{\partial L}{\partial y} * \frac{\partial y}{\partial h} * \frac{\partial h}{\partial z} * \frac{\partial z}{\partial w_1} \\
&= \begin{bmatrix} \frac{\partial L}{\partial w_{11}} \\ \frac{\partial L}{\partial w_{12}} \end{bmatrix} \\
&= \begin{bmatrix} (2 * y_{21} * w_{21} + 2 * y_{22} * w_{22}) * x_{21} \\ (2 * y_{21} * w_{21} + 2 * y_{22} * w_{22}) * x_{22} \end{bmatrix} \\
&= 2 * \begin{bmatrix} x_{21} \\ x_{22} \end{bmatrix} * \begin{bmatrix} y_{21} & y_{22} \end{bmatrix} * \begin{bmatrix} w_{21} \\ w_{22} \end{bmatrix} \\
&= 2 * (X^T \text{의 두번째 column}) * (y \text{의 두번째 row}) \\
&\quad * W_2^T
\end{aligned} \tag{34}$$

if) $z_1 \leq 0, z_2 \leq 0$ 일 때, $\frac{\partial L}{\partial w_1}$ 는 다음과 같다. (30), (31)사용

$$\begin{aligned}
\frac{\partial L}{\partial w_1} &= \frac{\partial L}{\partial L} * \frac{\partial L}{\partial y} * \frac{\partial y}{\partial h} * \frac{\partial h}{\partial z} * \frac{\partial z}{\partial w_1} \\
&= \begin{bmatrix} \frac{\partial L}{\partial w_{11}} \\ \frac{\partial L}{\partial w_{12}} \end{bmatrix} \\
&= \begin{bmatrix} 0 \\ 0 \end{bmatrix}
\end{aligned} \tag{35}$$

References

1. <http://cs231n.stanford.edu/> lecture 4