Jin Hyun Park
677001829

1.

$$E_{in}(\omega) = \frac{1}{N}\sum_{n=1}^{N}(\tanh(\omega^T x_n) - y_n)^2$$

$$\nabla E_{in}(\omega) = \frac{1}{N}\times 2 \cdot \sum_{n=1}^{N}(\tanh(\omega^T x_n) - y_n)\cdot\nabla\tanh(\omega^T x_n) \quad \cdots ①$$

Let $\tanh(x) = g(x)$. We will find the derivative of $g$ and use it to solve eq.①

$$\frac{d}{dx}\cdot g(x) = \frac{d}{dx}\cdot\frac{e^x - e^{-x}}{e^x + e^{-x}}$$

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

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

$$\therefore \frac{d}{dx}g(x) = \nabla g(x) = 1 - \left(\frac{e^x - e^{-x}}{e^x + e^{-x}}\right)^2 = 1 - (\tanh x)^2 \quad \cdots ②$$

Plugging the result ② to ①, we get

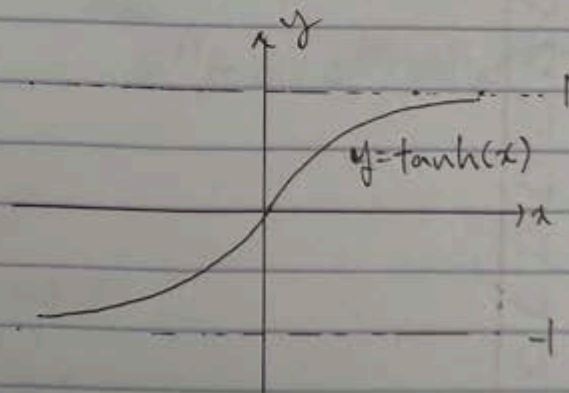$$\nabla E_{in}(\omega) = \frac{2}{N}\sum_{n=1}^{N}(\tanh(\omega^T x_n) - y_n)(1 - \tanh^2(\omega^T x_n))\cdot\frac{d}{d\omega}(\omega^T x_n)$$

$$\therefore \nabla E_{in}(\omega) = \frac{2}{N}\cdot\sum_{n=1}^{N}(\tanh(\omega^T x_n) - y_n)(1 - \tanh^2(\omega^T x_n))\cdot x_n$$

$$= (\text{gradient of in-sample error})$$

If $\omega \to \infty$, $\tanh^2(\omega^T x_n) \to 1$

Then,

$$\left(1 - \tanh^2(\omega^T x_n)\right) \to 0$$

This implies that $\nabla E_{in}(\omega)$ becomes 0 and this results in vanishing gradient issue. The weights won't be updated properly.


$y = \tanh(x)$

**Q2.**

$$x^0 \xrightarrow{w^1} S^1 \xrightarrow{\theta} x^1 \xrightarrow{w^2} \cdots \rightarrow x^{(L)} = h(x)$$

Weight matrices are:

$$S^1 = w^{1^T} x^0 \qquad \qquad S^{(L)} \xrightarrow{\theta} x^{(L)}$$

$$W^{(1)} = \begin{bmatrix} 0.1 & 0.2 \\ 0.7 & 0.4 \end{bmatrix} \qquad W^{(2)} = \begin{bmatrix} 0.2 \\ 1 \\ -3 \end{bmatrix} \qquad W^{(3)} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

and $\lambda = 2$, $y = 1$.

$$x^{(0)} = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \qquad S^{(1)} = \underset{W^{(1)^T}}{\begin{bmatrix} 0.1 & 0.7 \\ 0.2 & 0.4 \end{bmatrix}} \underset{x^{(0)}}{\begin{bmatrix} 1 \\ 2 \end{bmatrix}} = \begin{bmatrix} 0.7 \\ 1 \end{bmatrix}$$

$$x^{(1)} = \begin{bmatrix} 1 \\ 0.6 \\ 0.76 \end{bmatrix} \qquad S^{(2)} = \underset{W^{(2)^T}}{[0.2 \quad 1 \quad -3]} \underset{x^{(1)}}{\begin{bmatrix} 1 \\ 0.6 \\ 0.76 \end{bmatrix}} = [-1.48]$$

$$x^{(2)} = \begin{bmatrix} 1 \\ -0.9 \end{bmatrix} \qquad S^{(3)} = \underset{W^{(3)^T}}{[1 \quad 2]} \underset{x^{(2)}}{\begin{bmatrix} 1 \\ -0.9 \end{bmatrix}} = [-0.8], \quad x^{(3)} = [-0.8]$$

output transformation is identity.

Output transformation is identity. This means that $\theta'(s^{(3)}) = 1$

We use the following equation to calculate $\delta$. also,

things to remember.

$$\bigstar \begin{bmatrix} \delta^{(\ell)} = \theta'(s^{(\ell)}) \otimes \left[ W^{(\ell+1)} \delta^{(\ell+1)} \right]_1^{d^{(\ell)}} & \ell = [-1 \text{ to } 1] \\ \delta^{(L)} = 2(x^{(L)} - y) \theta'(s^{(L)}) \end{bmatrix}$$

$$\delta^{(3)} = 2(x^{(3)} - 1) \cdot 1 = 2(-1.8) = [-3.6]$$

$$\delta^{(2)} = \theta'(s^{(2)}) \otimes [w^{(3)} \delta^{(3)}] = (1 - \tanh^2(-1.48)) \otimes 2 \cdot (-3.6) = [-1.368]$$

$$\delta^{(1)} = \theta'(s^{(1)}) \otimes [w^{(2)} \cdot \delta^{(2)}] = \begin{bmatrix} (1 - \tanh^2(0.7)) \cdot 1 \\ (1 - \tanh^2(1)) \cdot (-3) \end{bmatrix} \cdot (-1.368) = \begin{bmatrix} -0.876 \\ 1.774 \end{bmatrix}$$

We ignore the 0th row of W

✦

$$\frac{\partial e}{\partial w^{(l)}} = x^{(l-1)} \cdot (\delta^{(l)})^T \quad : \text{thing to remember.}$$

$$\frac{\partial e}{\partial w^{(1)}} = x^{(0)} \cdot (\delta^{(1)})^T = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \cdot \begin{bmatrix} -0.876 & 1.734 \end{bmatrix} = \begin{bmatrix} 0.876 & 1.734 \\ -1.752 & 3.468 \end{bmatrix}$$

$$\underset{2\times 1}{} \qquad \underset{1\times 2}{}$$

$$\frac{\partial e}{\partial w^{(2)}} = x^{(1)} \cdot (\delta^{(2)})^T = \begin{bmatrix} 1 \\ 0.6 \\ 0.76 \end{bmatrix} \cdot \begin{bmatrix} -1.368 \end{bmatrix} = \begin{bmatrix} -1.368 \\ -0.821 \\ -1.040 \end{bmatrix}$$

$$\frac{\partial e}{\partial w^{(3)}} = x^{(2)} \cdot (\delta^{(3)})^T = \begin{bmatrix} 1 \\ -0.9 \end{bmatrix} \cdot \begin{bmatrix} -3.6 \end{bmatrix} = \begin{bmatrix} -3.6 \\ 3.24 \end{bmatrix}$$

3

① Standard residual block

input dim. = output dim

$$= 128 \times 16 \times 16 \times 32$$

<u>Batch</u> <u>feature</u> <u>channels</u>
size     map
         size

<span style="background:lightgreen">case1: Num. of param. with bias.= 18496</span>  since,

1st
layer  $\left( (\underline{32 \times 3 \times 3}) + 1 \right) \cdot \underline{32} = 9248$
       <u>channel</u>  <u>filter size</u>   32 filters
       depth      (kernel size)

2nd
layer  $\left( (32 \times 3 \times 3) + 1 \right) \cdot 32 = 9248$

<span style="background:lightgreen">case2: Num of param. without bias= 18432</span>  since,

1st   $(32 \times 3 \times 3) \times 32 = 9216$
layer

2nd   $(32 \times 3 \times 3) \times 32 = 9216$
layer

② Bottleneck block
input dim= output dim

$$= 128 \times 16 \times 16 \times 128$$
<u>Batchsize</u>              <u>channels</u>

<span style="background:lightgreen">case1: Num of param. with bias = 17600</span>  since

1st layer: $\left( (128 \times 1 \times 1) + 1 \right) \times 32 = 4128$

2nd layer: $\left( (32 \times 3 \times 3) + 1 \right) \times 32 = 9248$

3rd layer $\left( (32 \times 1 \times 1) + 1 \right) \times 128 = 4224$

case 2. Num of param without bias = 17408 since,

1st layer : $(128 \times 1 \times 1) \times 32 = 4096$

2nd layer : $(32 \times 3 \times 3) \times 32 = 9216$

3rd layer : $(32 \times 1 \times 1) \times 128 = 4096$

We can summarize the num of param in the following table

| | standard residual block | bottle-neck block |
|---|---|---|
| bias O | 18496 | 17600 |
| bias X | 18432 | 17408 |

Advantage of bottleneck over standard residual :
  · As we use less parameters, its computational cost is lower than the standard residual block
  · It can be used to obtain a representation with reduced dimensionality
  · This is similar to Autoencoder where the latent space(= vector or layer) contains(= encodes) the important features of the image!

Disadvantage of using bottleneck
  · We might loose some important features of the image because it
    $\underset{\text{(information)}}{}$ uses identity convolution.

(because we force dimensionality to be reduced which may lead to loose some important features)

**4.**

(a) shape of mean and variance is $1 \times C$

(b) As we normalise all the activation in a batch (batch size $= N$) and this normalisation cover all pixels (or element) in $H \times W$, the shape of mean and variance is $1 \times 1 \times 1 \times C$

**5.(a)**

$$X = \begin{bmatrix} \lambda_{11} & \lambda_{12} & \lambda_{13} & \lambda_{14} \\ \lambda_{21} & \lambda_{22} & \lambda_{23} & \lambda_{24} \end{bmatrix} \in R^{2 \times 4}$$

$$W^{ij} = [\, w_1^{ij}, \; w_2^{ij}, \; w_3^{ij} \,], \quad i = 1,2 \; / \; j = 1,2$$

$$Y = \begin{bmatrix} y_{11} & y_{12} \\ y_{21} & y_{22} \end{bmatrix} \in R^{2 \times 2}$$

$$y_{11} = w_1^{11}\lambda_{11} + w_2^{11}\cdot\lambda_{12} + w_3^{11}\lambda_{13} + w_1^{21}\lambda_{21} + w_2^{21}\cdot\lambda_{22} + w_3^{21}\lambda_{23}$$

$$y_{12} = w_1^{11}\lambda_{12} + w_2^{11}\cdot\lambda_{13} + w_3^{11}\cdot\lambda_{14} + w_1^{21}\lambda_{22} + w_2^{21}\cdot\lambda_{23} + w_3^{21}\cdot\lambda_{24}$$

$$y_{21} = w_1^{12}\lambda_{11} + w_2^{12}\lambda_{12} + w_3^{12}\lambda_{13} + w_1^{22}\lambda_{21} + w_2^{22}\lambda_{22} + w_3^{22}\lambda_{23}$$

$$y_{22} = w_1^{12}\lambda_{12} + w_2^{12}\lambda_{13} + w_3^{12}\lambda_{14} + w_1^{22}\lambda_{22} + w_2^{22}\lambda_{23} + w_3^{22}\lambda_{24}$$

$$\tilde{Y} = A\tilde{X} \quad, \quad A \in R^{4 \times 8}$$

$$\underset{4 \times 1}{} \quad \underset{8 \times 1}{}$$

$$\underset{4 \times 8}{} \qquad \underset{4 \times 8}{}$$

$$\underset{4\times1}{\begin{bmatrix} y_{11} \\ y_{12} \\ y_{21} \\ y_{22} \end{bmatrix}} = \begin{bmatrix} w_1^{11} & w_2^{11} & w_3^{11} & 0 & w_1^{21} & w_2^{21} & w_3^{21} & 0 \\ 0 & w_1^{11} & w_2^{11} & w_3^{11} & 0 & w_1^{21} & w_2^{21} & w_3^{21} \\ w_1^{12} & w_2^{12} & w_3^{12} & 0 & w_1^{22} & w_2^{22} & w_3^{22} & 0 \\ 0 & w_1^{12} & w_2^{12} & w_3^{12} & 0 & w_1^{22} & w_2^{22} & w_3^{22} \end{bmatrix} \underset{8\times1}{\begin{bmatrix} \lambda_{11} \\ \lambda_{12} \\ \lambda_{13} \\ \lambda_{14} \\ \lambda_{21} \\ \lambda_{22} \\ \lambda_{23} \\ \lambda_{24} \end{bmatrix}}$$

5(b)

$$
\begin{bmatrix} \dfrac{\partial L}{\partial \lambda_{11}} \\[6pt] \dfrac{\partial L}{\partial \lambda_{12}} \\[6pt] \dfrac{\partial L}{\partial \lambda_{13}} \\[6pt] \dfrac{\partial L}{\partial \lambda_{14}} \\[6pt] \dfrac{\partial L}{\partial \lambda_{21}} \\[6pt] \dfrac{\partial L}{\partial \lambda_{22}} \\[6pt] \dfrac{\partial L}{\partial \lambda_{23}} \\[6pt] \dfrac{\partial L}{\partial \lambda_{24}} \end{bmatrix}
=
\begin{bmatrix}
w_1^{11} & 0 & w_1^{12} & 0 \\[4pt]
w_2^{11} & w_1^{11} & w_2^{12} & w_1^{12} \\[4pt]
w_3^{11} & w_2^{11} & w_3^{12} & w_3^{12} \\[4pt]
0 & w_3^{11} & 0 & w_3^{12} \\[4pt]
w_1^{21} & 0 & w_1^{22} & 0 \\[4pt]
w_2^{21} & w_1^{21} & w_2^{22} & w_1^{22} \\[4pt]
w_3^{21} & w_2^{21} & w_3^{22} & w_3^{22} \\[4pt]
0 & w_3^{21} & 0 & w_3^{22}
\end{bmatrix}
\begin{bmatrix} \dfrac{\partial L}{\partial y_{11}} \\[6pt] \dfrac{\partial L}{\partial y_{12}} \\[6pt] \dfrac{\partial L}{\partial y_{21}} \\[6pt] \dfrac{\partial L}{\partial y_{22}} \end{bmatrix}
$$

$8 \times 1$      $8 \times 4$      $4 \times 1$

$\dfrac{\partial L}{\partial \tilde{x}}$      $B$      $\dfrac{\partial L}{\partial \tilde{y}}$

5(b)

$$\frac{\partial L}{\partial \tilde{Y}} = \left[ \frac{\partial L}{\partial y_{11}} \quad \frac{\partial L}{\partial y_{12}} \quad \frac{\partial L}{\partial y_{21}} \quad \frac{\partial L}{\partial y_{22}} \right]^T$$

$$\frac{\partial L}{\partial \tilde{X}} = B \cdot \frac{\partial L}{\partial \tilde{Y}} \Rightarrow \underbrace{\frac{\partial L}{\partial \tilde{X}}}_{8\times1} = \underbrace{\frac{\partial \tilde{Y}}{\partial \tilde{X}}}_{8\times4} \cdot \underbrace{\frac{\partial L}{\partial \tilde{Y}}}_{4\times1}$$

$$\frac{\partial L}{\partial x_{11}} = \frac{\partial y_{11}}{\partial x_{11}} \cdot \frac{\partial L}{\partial y_{11}} + \cancel{\frac{\partial y_{12}}{\partial x_{11}} \cdot \frac{\partial L}{\partial y_{12}}}^{0} + \frac{\partial y_{21}}{\partial x_{11}} \cdot \frac{\partial L}{\partial y_{21}} + \cancel{\frac{\partial y_{22}}{\partial x_{11}} \cdot \frac{\partial L}{\partial y_{22}}}^{0} \quad \cdots \ (1)$$

$$\frac{\partial L}{\partial x_{12}} = \frac{\partial y_{11}}{\partial x_{12}} \cdot \frac{\partial L}{\partial y_{11}} + \frac{\partial y_{12}}{\partial x_{12}} \cdot \frac{\partial L}{\partial y_{12}} + \frac{\partial y_{21}}{\partial x_{12}} \cdot \frac{\partial L}{\partial y_{21}} + \frac{\partial y_{22}}{\partial x_{12}} \cdot \frac{\partial L}{\partial y_{22}} \quad \cdots \ (2)$$

$$\frac{\partial L}{\partial x_{13}} = \frac{\partial y_{11}}{\partial x_{13}} \cdot \frac{\partial L}{\partial y_{11}} + \frac{\partial y_{12}}{\partial x_{13}} \cdot \frac{\partial L}{\partial y_{12}} + \frac{\partial y_{21}}{\partial x_{13}} \cdot \frac{\partial L}{\partial y_{21}} + \frac{\partial y_{22}}{\partial x_{13}} \cdot \frac{\partial L}{\partial y_{22}} \quad \cdots \ (3)$$

$$\frac{\partial L}{\partial x_{14}} = \cancel{\frac{\partial y_{11}}{\partial x_{14}} \cdot \frac{\partial L}{\partial y_{11}}}^{0} + \frac{\partial y_{12}}{\partial x_{14}} \cdot \frac{\partial L}{\partial y_{12}} + \cancel{\frac{\partial y_{21}}{\partial x_{14}} \cdot \frac{\partial L}{\partial y_{21}}}^{0} + \frac{\partial y_{22}}{\partial x_{14}} \cdot \frac{\partial L}{\partial y_{22}} \quad \cdots \ (4)$$

By simplifying eq (1), we get $w_1^{11} \frac{\partial L}{\partial y_{11}} + w_1^{12} \frac{\partial L}{\partial y_{21}}$

" (2), " $w_2^{11} \frac{\partial L}{\partial y_{11}} + w_1^{11} \frac{\partial L}{\partial y_{12}} + w_2^{12} \frac{\partial L}{\partial y_{21}} + w_1^{12} \frac{\partial L}{\partial y_{22}}$

(3), " $w_3^{11} \frac{\partial L}{\partial y_{11}} + w_2^{11} \frac{\partial L}{\partial y_{12}} + w_3^{12} \frac{\partial L}{\partial y_{21}} + w_2^{12} \frac{\partial L}{\partial y_{22}}$

(4), " $w_3^{11} \frac{\partial L}{\partial y_{12}} + w_3^{12} \cdot \frac{\partial L}{\partial y_{22}}$

And for $\frac{\partial L}{\partial x_{2k}}$ where $k = \{1, 2, 3, 4\}$, perform the same step. See the next page.

$$\frac{\partial L}{\partial \lambda_{21}} = \frac{\partial y_{11}}{\partial \lambda_{21}} \cdot \frac{\partial L}{\partial y_{11}} + \frac{\partial y_{12}}{\partial \lambda_{21}} \cancel{\frac{\partial L}{\partial y_{12}}}^{0} + \frac{\partial y_{21}}{\partial \lambda_{21}} \frac{\partial L}{\partial y_{21}} + \frac{\partial y_{22}}{\partial \lambda_{21}} \cancel{\frac{\partial L}{\partial y_{22}}}^{0} \quad \cdots (5)$$

$$\frac{\partial L}{\partial \lambda_{22}} = \frac{\partial y_{11}}{\partial \lambda_{22}} \cdot \frac{\partial L}{\partial y_{11}} + \frac{\partial y_{12}}{\partial \lambda_{22}} \cdot \frac{\partial L}{\partial y_{12}} + \frac{\partial y_{21}}{\partial \lambda_{22}} \cdot \frac{\partial L}{\partial y_{21}} + \frac{\partial y_{22}}{\partial \lambda_{22}} \cdot \frac{\partial L}{\partial y_{22}} \quad \cdots (6)$$

$$\frac{\partial L}{\partial \lambda_{23}} = \frac{\partial y_{11}}{\partial \lambda_{23}} \cdot \frac{\partial L}{\partial y_{11}} + \frac{\partial y_{12}}{\partial \lambda_{23}} \cdot \frac{\partial L}{\partial y_{12}} + \frac{\partial y_{21}}{\partial \lambda_{23}} \cdot \frac{\partial L}{\partial y_{21}} + \frac{\partial y_{22}}{\partial \lambda_{23}} \cdot \frac{\partial L}{\partial y_{22}} \quad \cdots (7)$$

$$\frac{\partial L}{\partial \lambda_{24}} = \frac{\partial y_{11}}{\partial \lambda_{24}} \cancel{\frac{\partial L}{\partial y_{11}}}^{0} + \frac{\partial y_{12}}{\partial \lambda_{24}} \cdot \frac{\partial L}{\partial y_{12}} + \frac{\partial y_{21}}{\partial \lambda_{24}} \cancel{\frac{\partial L}{\partial y_{21}}}^{0} + \frac{\partial y_{22}}{\partial \lambda_{24}} \cdot \frac{\partial L}{\partial y_{22}} \quad \cdots (8)$$

By simplifying eq. (5), we get $\quad w_1^{21} \cdot \frac{\partial L}{\partial y_{11}} + w_1^{22} \cdot \frac{\partial L}{\partial y_{21}}$

$\qquad \qquad '' \qquad (6) \qquad '' \qquad w_2^{21} \cdot \frac{\partial L}{\partial y_{11}} + w_1^{21} \cdot \frac{\partial L}{\partial y_{12}} + w_2^{22} \frac{\partial L}{\partial y_{21}} + w_1^{22} \cdot \frac{\partial L}{\partial y_{22}}$

$\qquad \qquad '' \qquad (7) \qquad '' \qquad w_3^{21} \cdot \frac{\partial L}{\partial y_{11}} + w_2^{21} \frac{\partial L}{\partial y_{12}} + w_3^{22} \cdot \frac{\partial L}{\partial y_{21}} + w_2^{22} \cdot \frac{\partial L}{\partial y_{22}}$

$\qquad \qquad '' \qquad (8) \qquad '' \qquad w_3^{21} \cdot \frac{\partial L}{\partial y_{12}} + w_3^{22} \frac{\partial L}{\partial y_{22}}$

$B \in R^{8 \times 4}$

$$B = \begin{bmatrix} w_1^{11} & 0 & w_1^{12} & 0 \\ w_2^{11} & w_1^{11} & w_2^{12} & w_1^{12} \\ w_3^{11} & w_2^{11} & w_3^{12} & w_2^{12} \\ w_1^{21} & 0 & w_1^{22} & 0 \\ w_2^{21} & w_1^{21} & w_2^{22} & w_1^{22} \\ w_3^{21} & w_2^{21} & w_3^{22} & w_2^{22} \\ 0 & w_3^{21} & 0 & w_3^{22} \end{bmatrix}$$

$\xrightarrow{\qquad} \quad {}'' \; 0 \; w_3^{11} \; 0 \; w_3^{12} \; ''$

Relationship between A and B
$$A = B^T$$

**5(C)**

$$\text{Conv}_{out} = \frac{\text{conv}_{in} - K}{S} + 1 \qquad \cdots \text{(A)}$$

(spatial) input size of conv layer (includes padding.)
kernel size
stride

output size of a conv layer (spatial)

$$\frac{\partial L}{\partial \tilde{X}} = B \cdot \frac{\partial L}{\partial \tilde{Y}}$$

$R^{2\times4}$  $R^{2\times2}$

✱ We don't have kernels for convolution if there is no padding. Consider (A), then

$$4 = \frac{2-K}{S} + 1 \quad , \quad 3 = \frac{2-K}{S} \quad , \quad 3S = 2-K . \text{ We cannot find } S \text{ and } K \text{ that}$$
(B)    such

satisfies equation (B).

$\boxed{\text{Consider (A) again ✱ with padding on } \frac{\partial L}{\partial y}}$

$$(2+2\cdot2) \quad 4 = \frac{6-3}{1} + 1 . \quad K=3. \; S=1 \text{ then this condition satisfies 5(c). Therefore,}$$

padding left & right

original conv$_{in}$

$$\frac{\partial L}{\partial y} = \begin{bmatrix} 0 & 0 & \frac{\partial L}{\partial y_{11}} & \frac{\partial L}{\partial y_{12}} & 0 & 0 \\ 0 & 0 & \frac{\partial L}{\partial y_{21}} & \frac{\partial L}{\partial y_{22}} & 0 & 0 \end{bmatrix}$$

withpadding

$$\frac{\partial L}{\partial X} = \begin{bmatrix} \frac{\partial L}{\partial x_{11}} & \frac{\partial L}{\partial x_{12}} & \frac{\partial L}{\partial x_{13}} & \frac{\partial L}{\partial x_{14}} \\ \frac{\partial L}{\partial x_{21}} & \frac{\partial L}{\partial x_{22}} & \frac{\partial L}{\partial x_{23}} & \frac{\partial L}{\partial x_{24}} \end{bmatrix}$$

$$\text{kernel} \Rightarrow \quad W^{\bar{i}\bar{j}} = \begin{bmatrix} W_3^{\bar{j}\bar{i}} , & W_2^{\bar{j}\bar{i}} , & W_1^{\bar{j}\bar{i}} \end{bmatrix} \quad \bar{i}=1,2. \; \bar{j}=1,2.$$

also consider that $A = B^T$! And equation (3) on qt5 in HW2 is $W^{\bar{i}\bar{j}} = \begin{bmatrix} W_1^{\bar{i}\bar{j}} , & W_2^{\bar{i}\bar{j}} , & W_3^{\bar{i}\bar{j}} \end{bmatrix} \begin{matrix} \bar{i}=1,2 \\ \bar{j}=1,2 \end{matrix}$