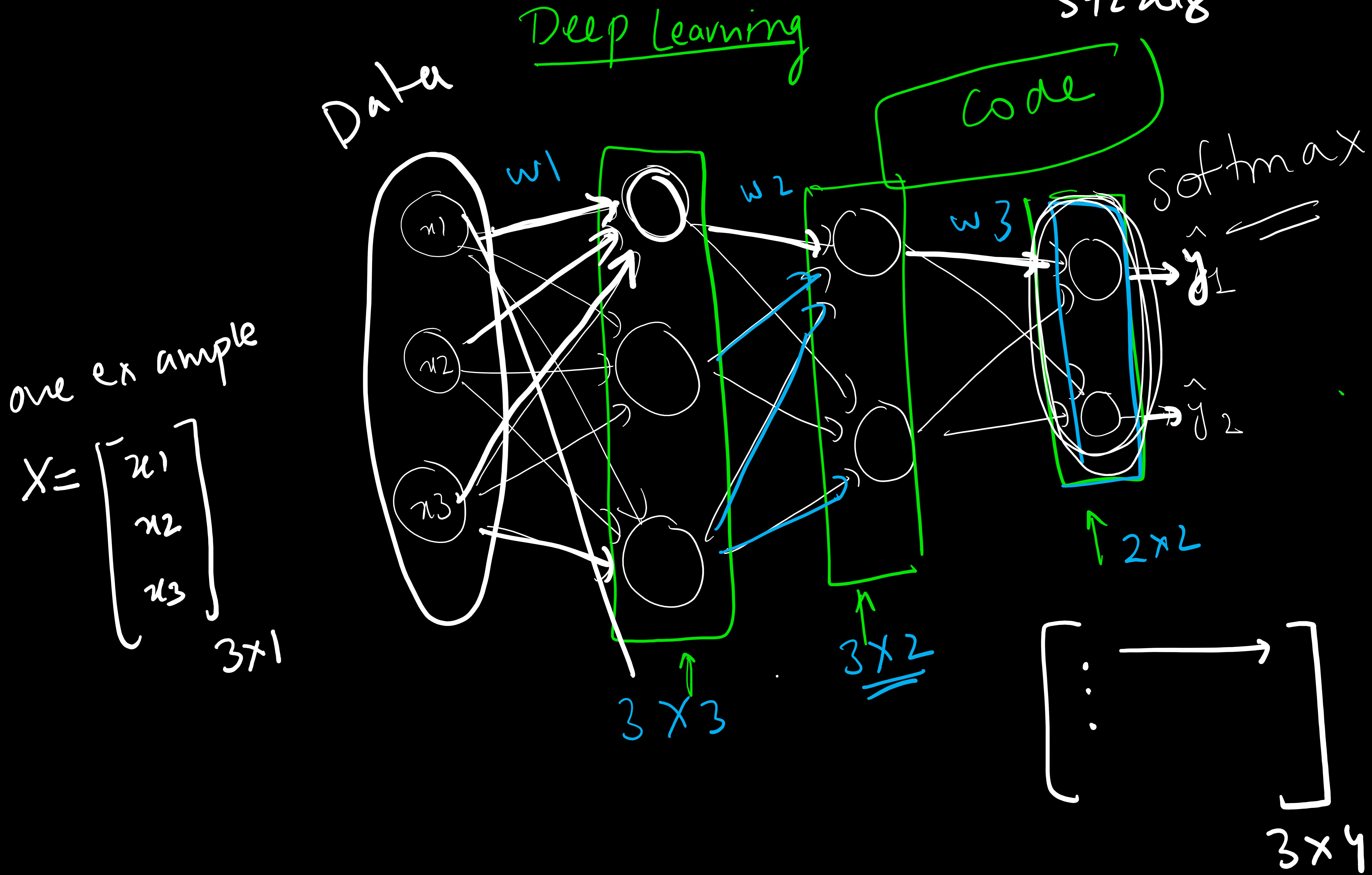


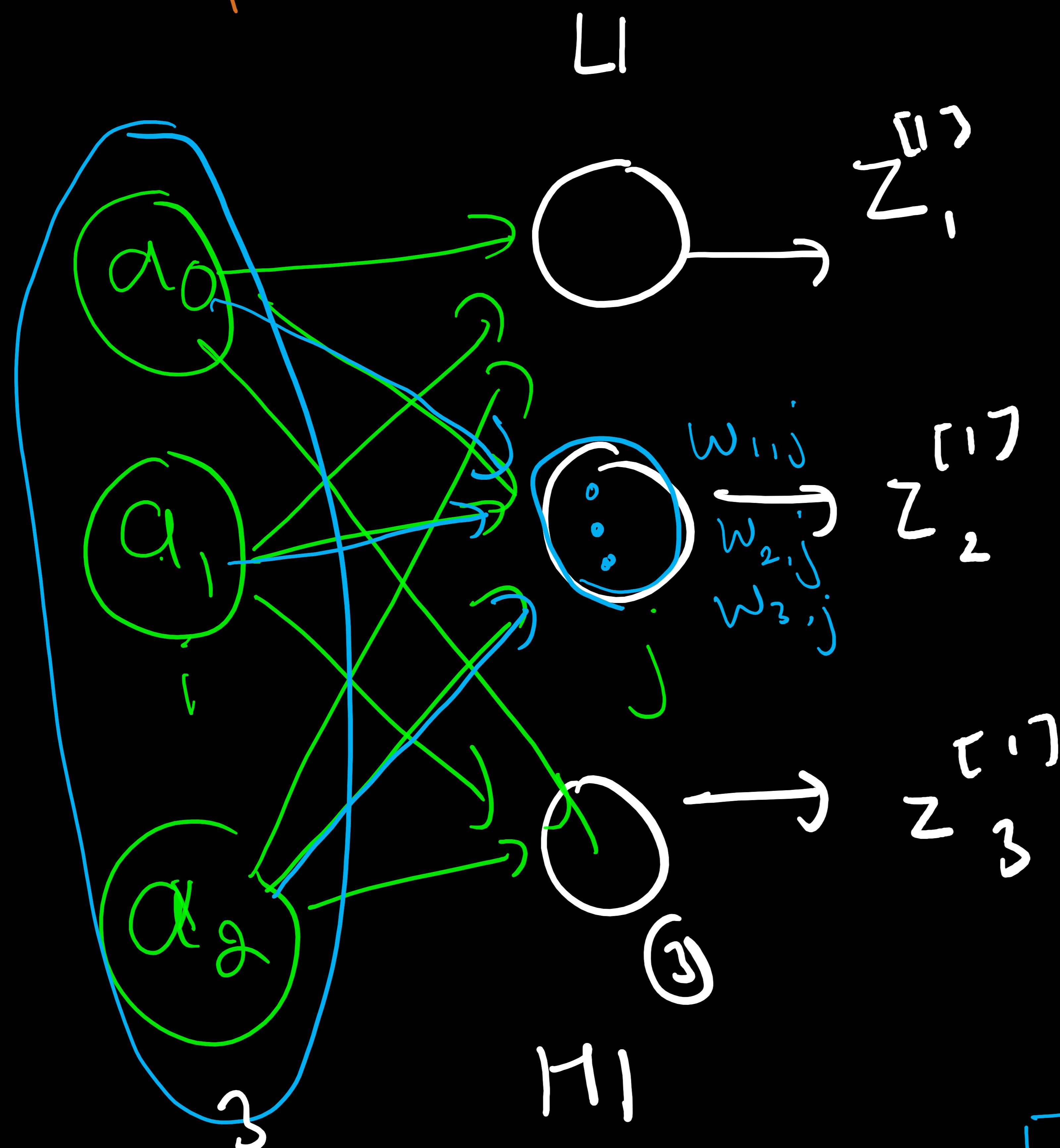
Forward & Back Propagation Derivation

Piadeek Navang
5-12-2018



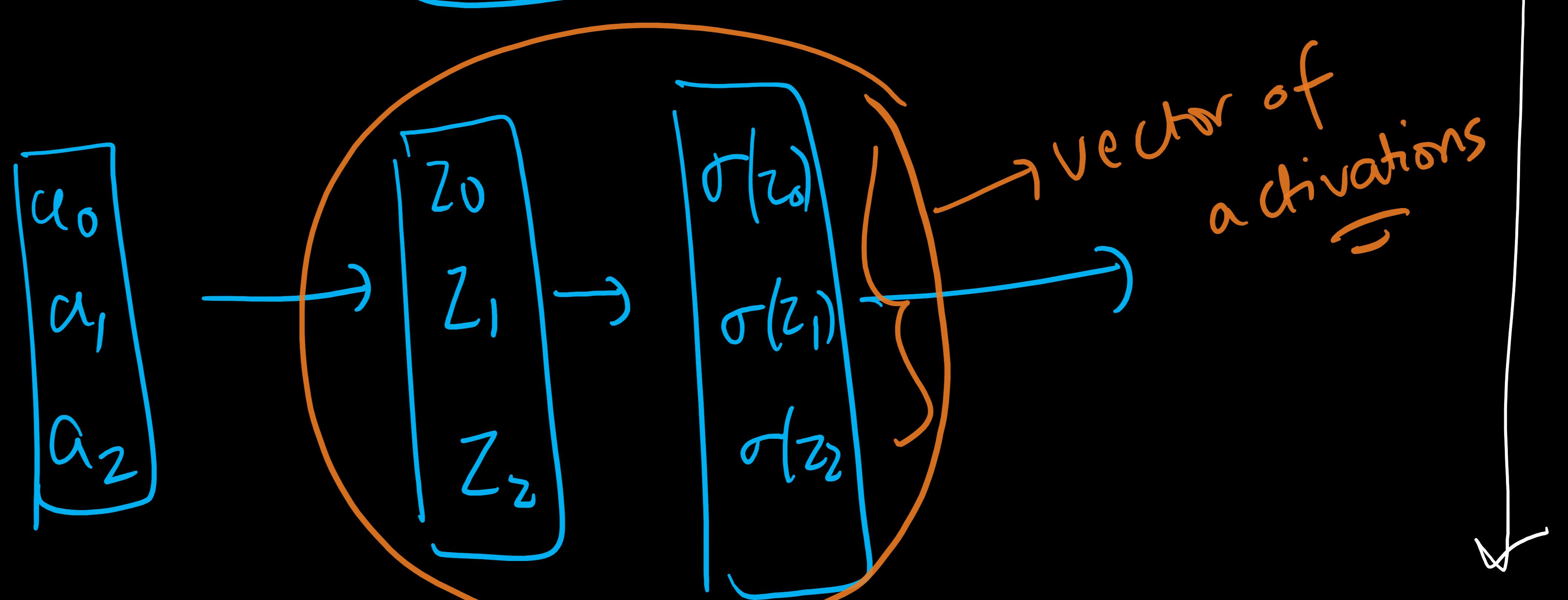
$$z = \sum w_i x_i + b$$

for one input example



$$w^l = \begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} \\ w_{2,1} & w_{2,2} & w_{2,3} \\ w_{3,1} & w_{3,2} & w_{3,3} \end{bmatrix}_{(3 \times 3)}$$

$$z_j^{[l]} = \sum_i w_{i,j}^{[l]} a_i^{[l-1]} + b$$



Vector
Notation?

$$z^{[e]} = \underline{w^e}^\top \underline{a}^{l-1} + \underline{b}^e$$

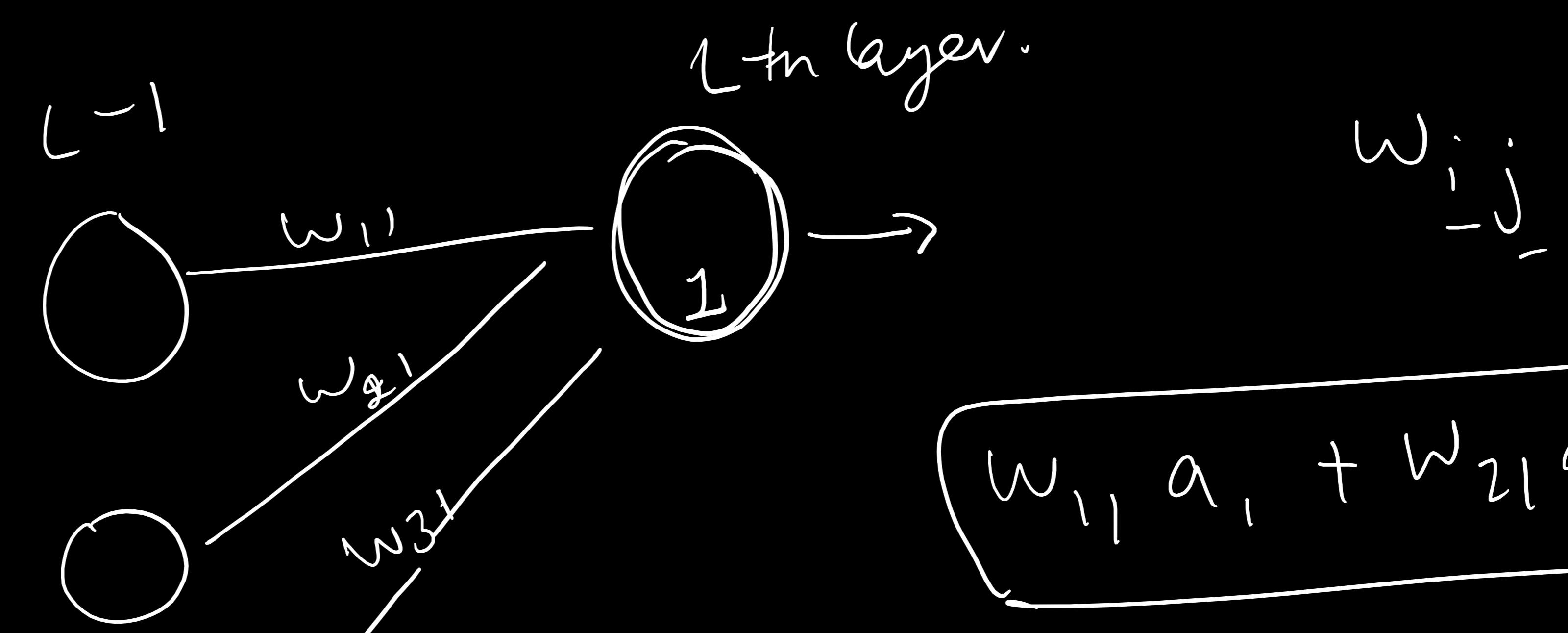
Dimensions:
 $(l, l-1) \times (l+1)$
 $(l, 1)$
 $(w^e)^\top \rightarrow (l, l-1)$

↓
 $z^{[e]}$

= \underline{a}^{l-1}

$\vdots l \text{ outputs}$

$$(w^L)^\top = \begin{bmatrix} w_{11}^L & w_{21}^L & w_{31}^L \\ w_{12} & w_{22} & w_{32} \\ w_{13} & w_{23} & w_{33} \end{bmatrix} \begin{bmatrix} a_1^{l-1} \\ a_2^{l-1} \\ a_3^{l-1} \end{bmatrix} = \begin{bmatrix} z_1^{[e]} \\ z_2^{[e]} \\ z_3^{[e]} \end{bmatrix} = z^{[e]}$$



Forward
Propagation
over 1 example

$$\begin{aligned} a^{[1]} &= g(z^{[0]}) = \sigma(z^{[0]}) = \left[\begin{array}{c} \sigma(z_1^{[0]}) \\ \sigma(z_2^{[0]}) \\ \vdots \\ \sigma(z_n^{[0]}) \end{array} \right] \quad \text{Output after hidden layer} \\ z_1 &= (w^{[1]})^T a^{[0]} + b^{[1]} \\ a_1 &= \sigma(z_1) \\ z_2 &= (w^{[2]})^T a^{[1]} + b^{[2]} \\ a_2 &= \sigma(z_2) \\ z_3 &= (w^3)^T a^{[2]} + b^{[3]} \\ \hat{y} &= \text{Softmax}(z_3) \end{aligned}$$

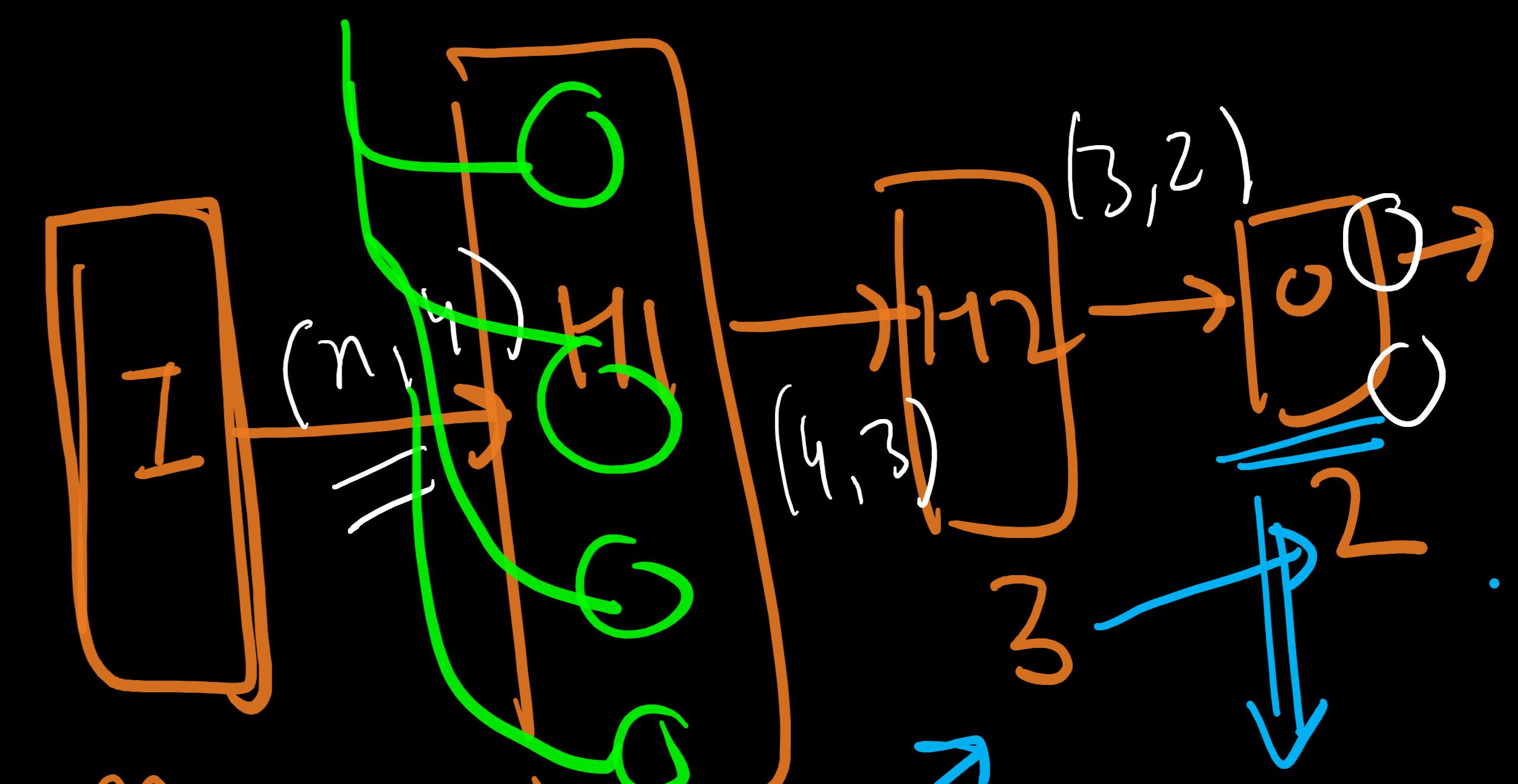
$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}_{(n,1)} \quad X = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \dots & x_n^{(1)} \\ \vdots & \vdots & & \vdots \\ x_1^{(i)} & x_2^{(i)} & \dots & x_n^{(i)} \\ \vdots & \vdots & & \vdots \\ x_1^{(m)} & x_2^{(m)} & \dots & x_n^{(m)} \end{bmatrix}_{m \times n}$$

$$A^{[0]} = X$$

$$z_i = A^{[0]} \underbrace{w^{[1]}}_{(m,n) \times (n,1)} + b^{[1]}_{(1,1)}$$

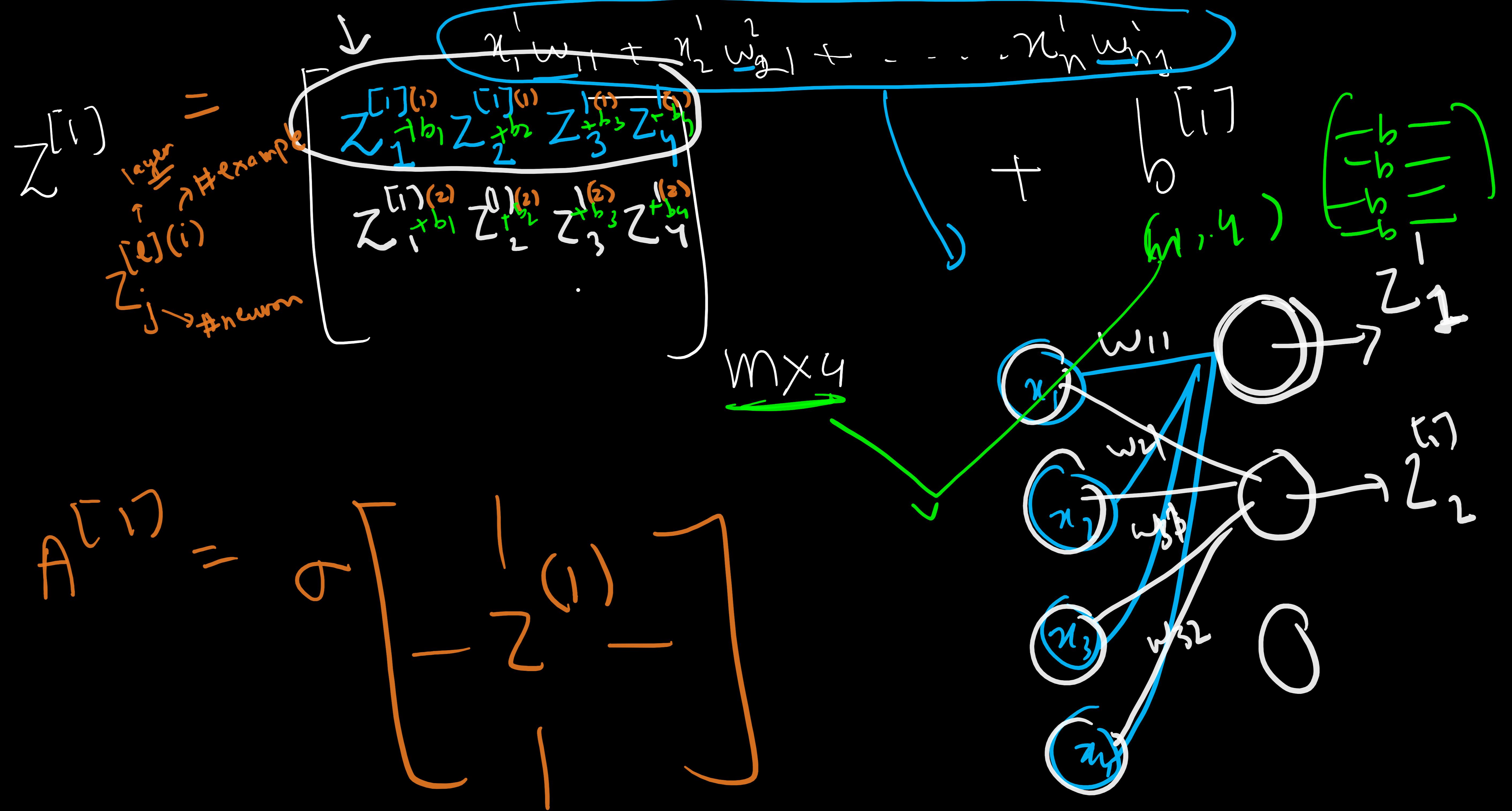
$$X = A^{[0]} = \begin{bmatrix} x_1^1 & x_2^1 & x_3^1 & \dots & x_n^1 \\ x_1^2 & x_2^2 & x_3^2 & \dots & x_n^2 \end{bmatrix}_{m \times n}$$

$w^{[1]} = \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1n} \\ w_{21} & w_{22} & \dots & w_{2n} \end{bmatrix}_{n \times 4}$

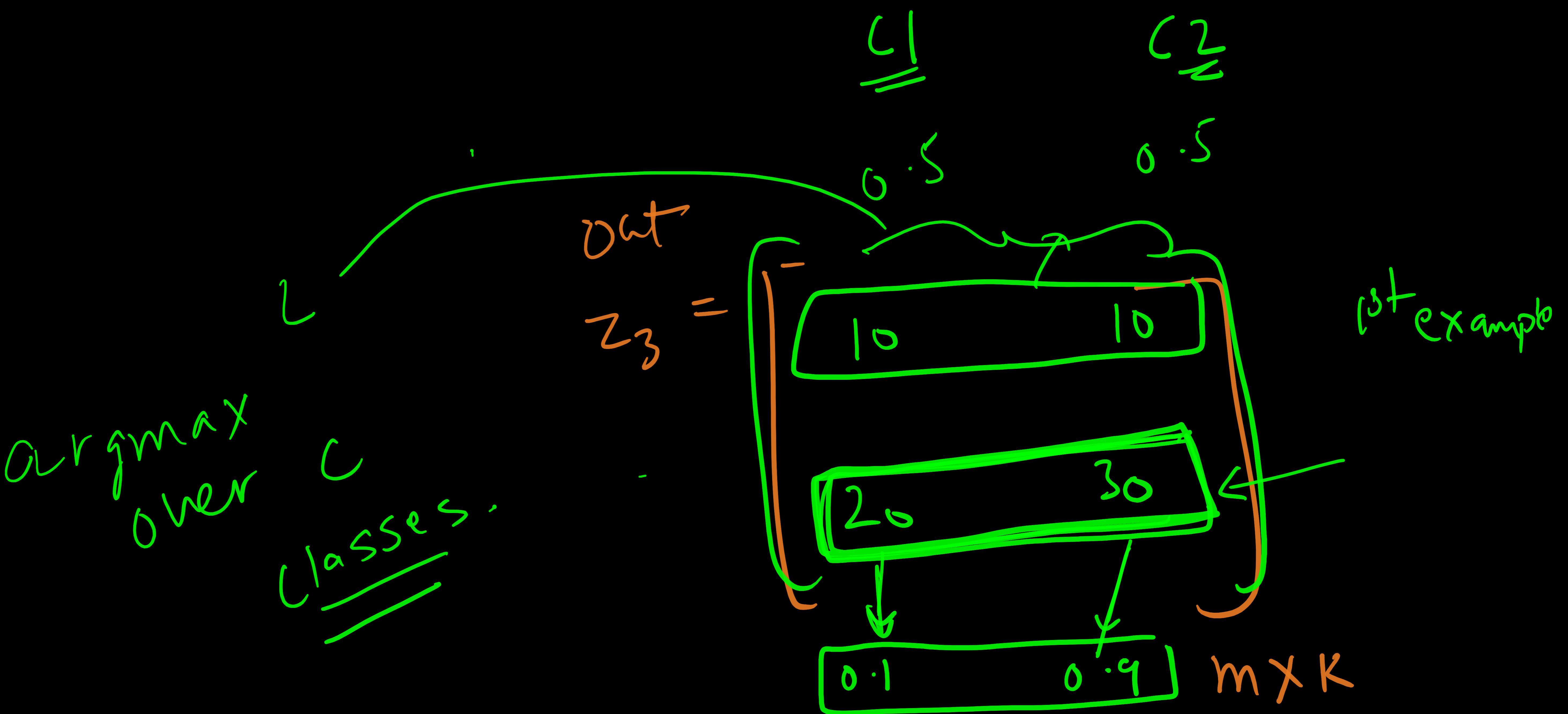


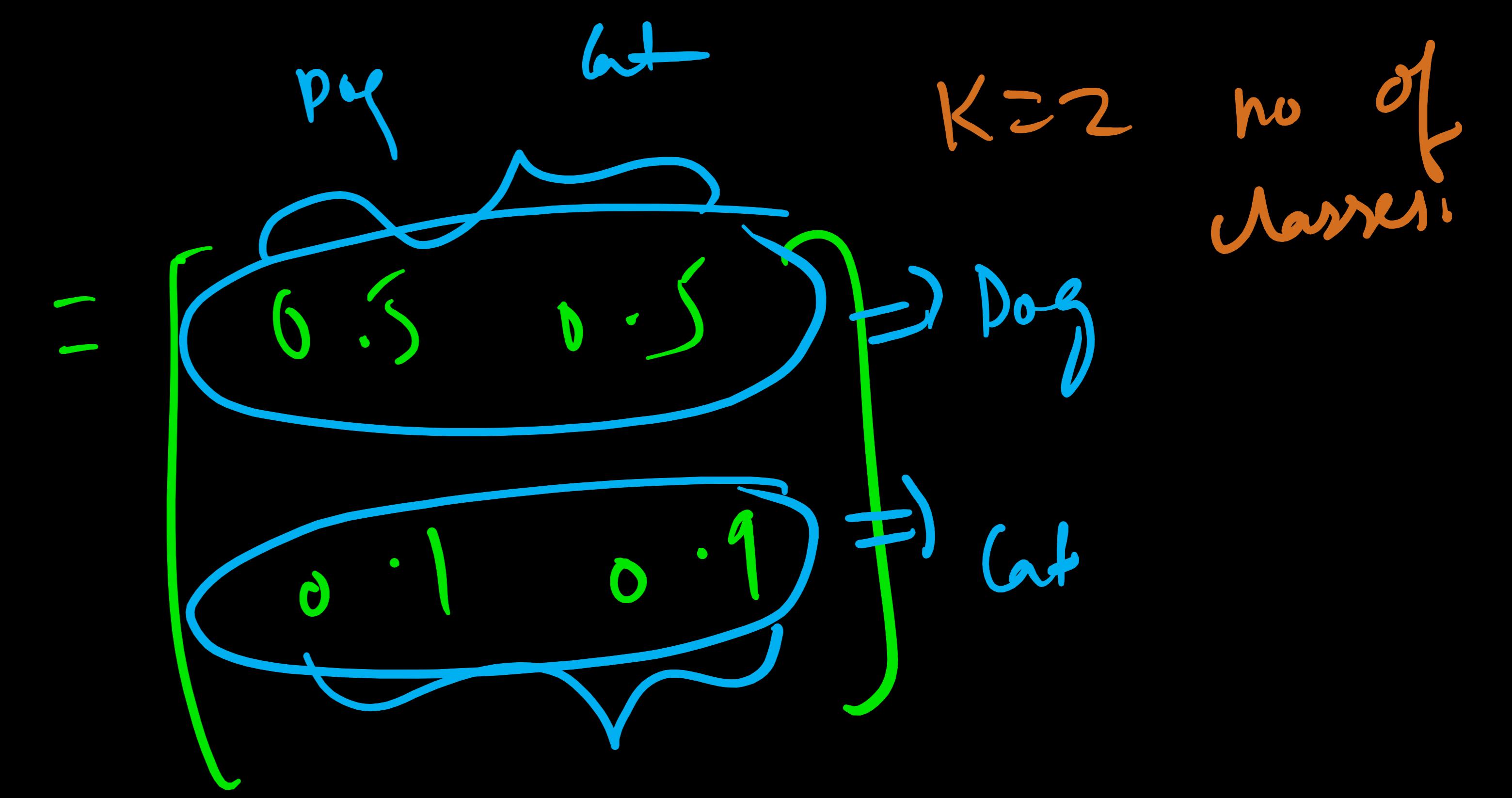
$$w_{ij} = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 2 & 0 & 3 & 0 \end{bmatrix}_{m \times 2}$$

$x_j^i = j^{\text{th}} \text{ feature}$
for i^{th} example



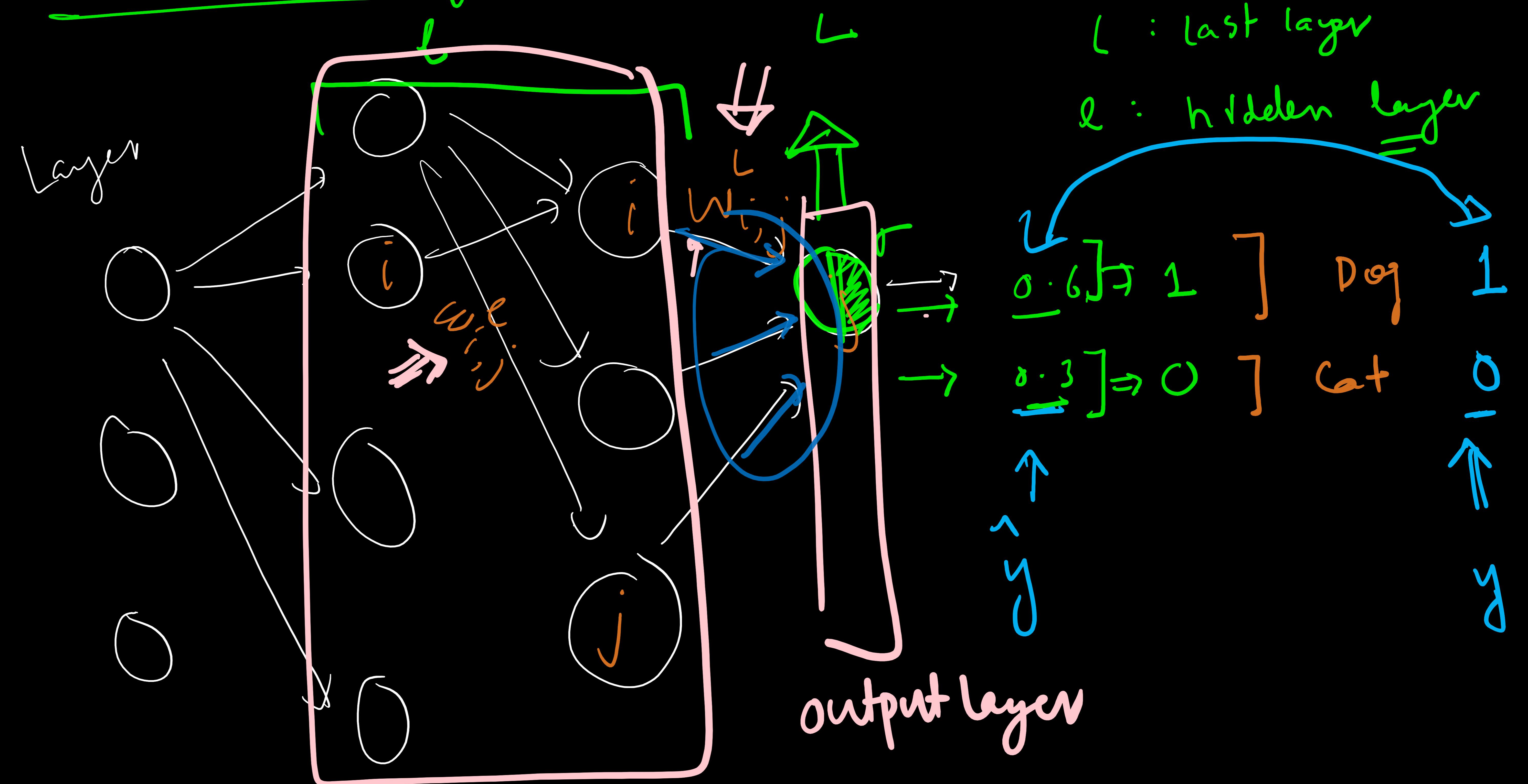
$$A^{(1)} = \sigma \begin{bmatrix} 1 \\ -z^{(1)} \\ 1 \end{bmatrix}$$





PART-3

Backpropagation



hidden neurons

For simplification -

$$MSE = L(\hat{y}, y) =$$

$$\frac{1}{2} \sum_{i=1}^m (\hat{y}^{(i)} - \hat{y}^{(i)})^2$$

Case-1 output layer:

$$\frac{\partial L}{\partial W_{i,j}^L} = ?$$

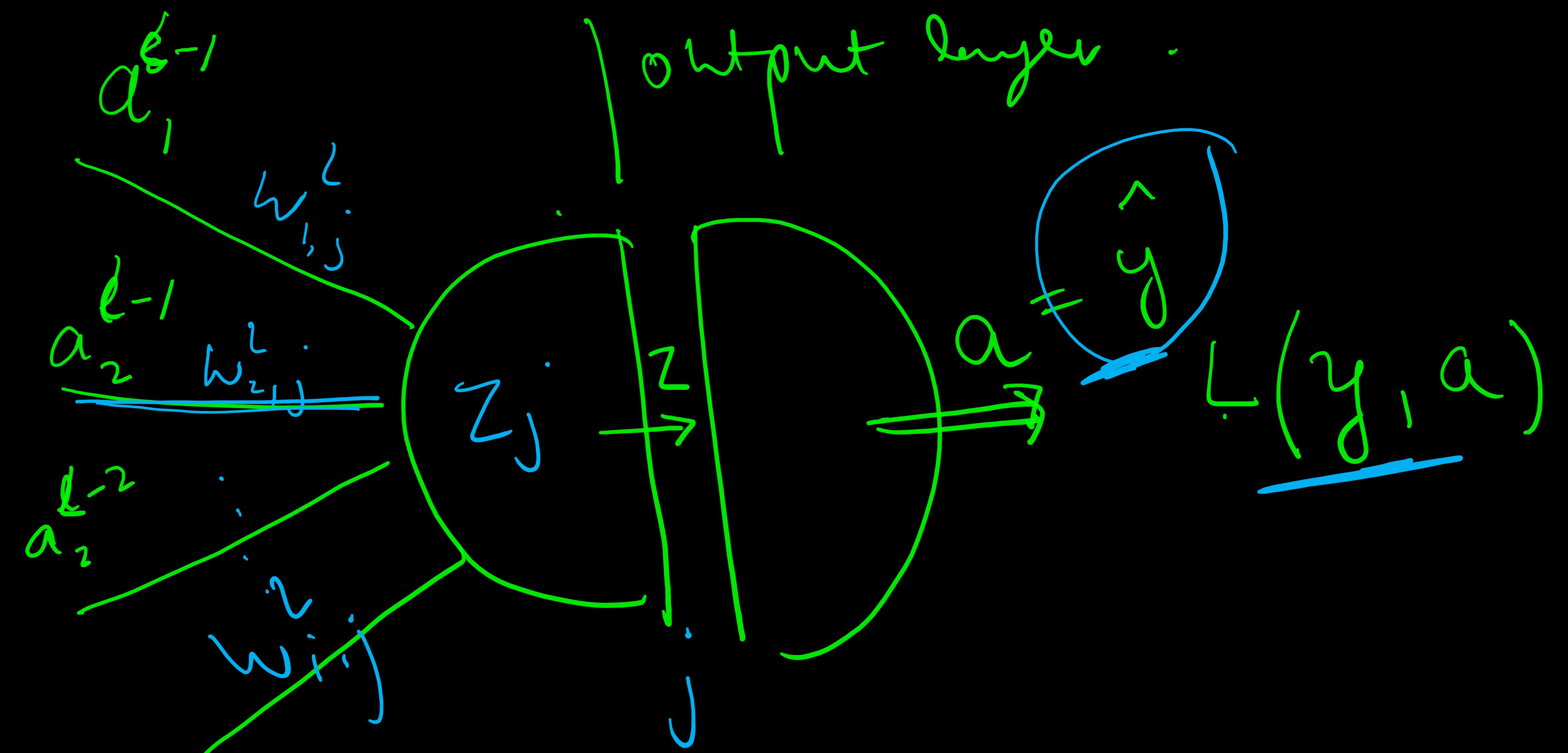
$$\frac{\partial L}{\partial b^L} = ?$$

Case-2 hidden layer:

$$\frac{\partial L}{\partial W_{i,j}^l} = ?$$

$$\frac{\partial L}{\partial b^l} = ?$$

Case-1



$$w_{ij} \rightarrow z_j \rightarrow a_i \rightarrow L(a, y)$$

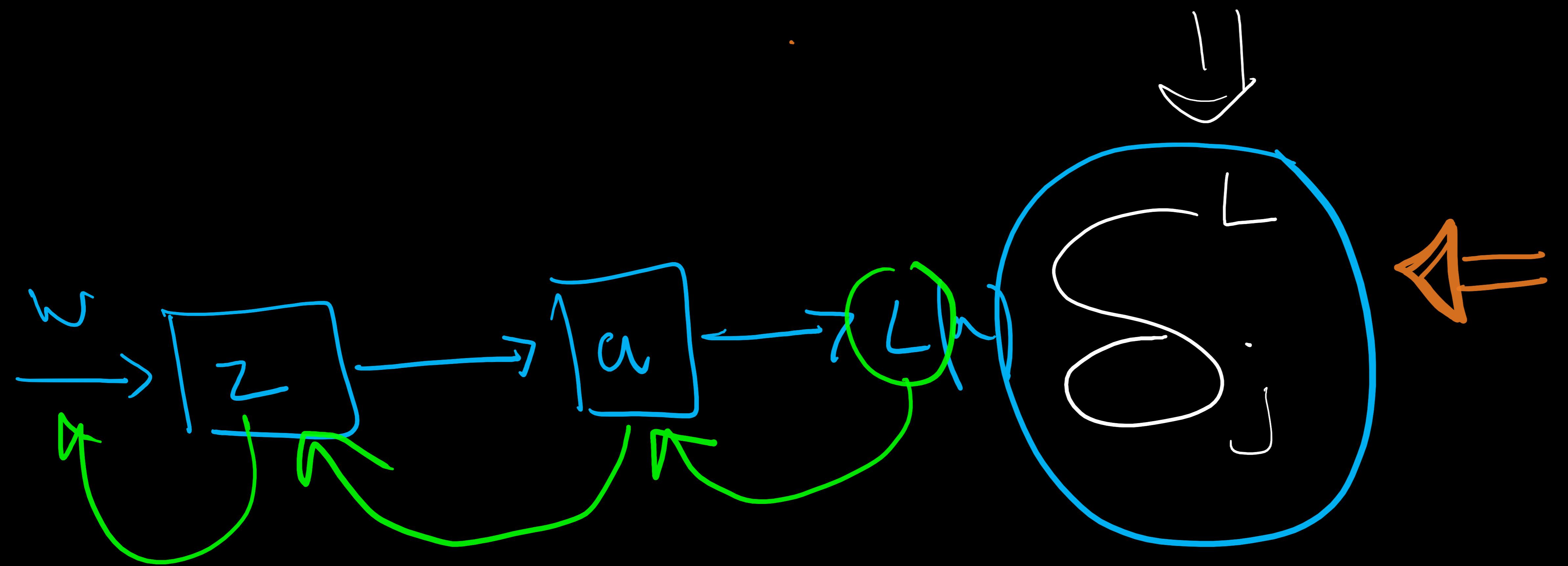
$$z_j =$$

$$\sum_i^l w_{ij} a_i^{l-1} + b_j^l$$

$$\frac{\partial L}{\partial w_{ij}}$$

$$= \boxed{\frac{\partial L}{\partial a_i} \cdot \frac{\partial a_i}{\partial z_j} \cdot \frac{\partial z_j}{\partial w_{ij}}}$$

Chain
Rule



$$\textcircled{1} \quad \left\{ \begin{array}{l} \frac{\partial L}{\partial a_j} = (y_j - a_j) \\ \frac{\partial a_j}{\partial z_j} = \sigma'(z_j) \end{array} \right.$$

$$\textcircled{2} \quad \left\{ \begin{array}{l} \frac{\partial a_j}{\partial z_j} = \sigma'(z_j) \\ \frac{\partial z_j}{\partial w_{i,j}} = a_i^{l-1} \end{array} \right.$$

$$\textcircled{3} \quad \left\{ \begin{array}{l} \frac{\partial z_j}{\partial w_{i,j}} = a_i^{l-1} \\ \frac{\partial z_j}{\partial b_j} = 1 \end{array} \right.$$

$$\frac{\partial L}{\partial z_j}$$

L	a_L	a_i^l	y^l
0	0.3	0.5	1
i	0.5	0.2	0
0	0.2	0	0

$$\textcircled{1} \quad L = \frac{1}{2} \sum_i (y_i^l - a_i^l)^2$$

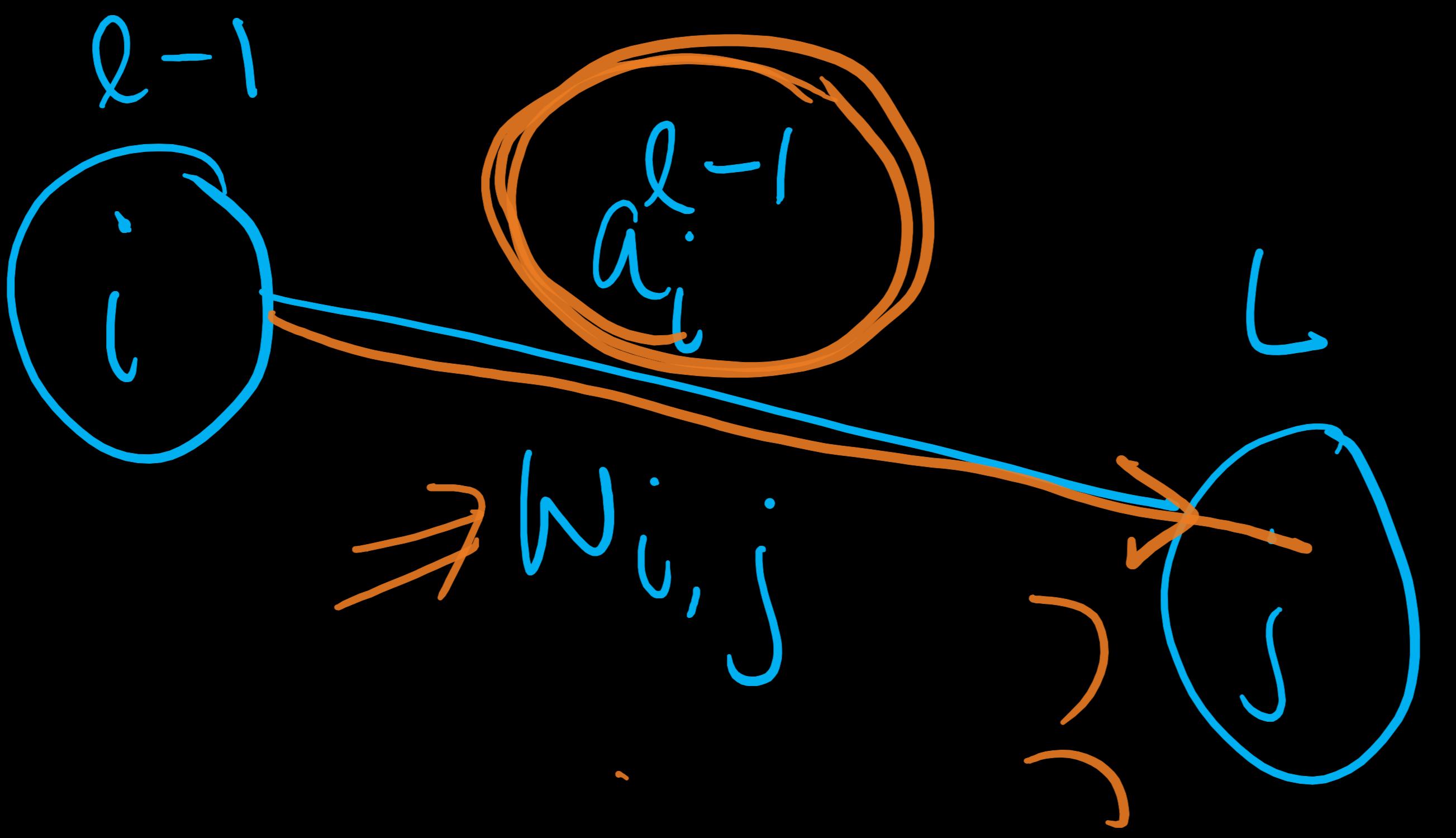
$$\frac{\partial L}{\partial a_j} = \frac{1}{2} \times (y_j - a_j)$$

$$\textcircled{2} \quad a_j = \sigma(z_j)$$

$$\frac{\partial a_j}{\partial z_j} = \underline{\sigma'(z_j)}$$

$$\frac{\partial z_j}{\partial w_{i,j}} = (1 - \sigma(z_j)) \sigma(z_j)$$

$$\textcircled{3} \quad z_j = \sum_i w_{i,j} a_i^{l-1} + b_j^l$$



$$\frac{\partial z_j}{\partial w_{i,j}} = a_i^{l-1}$$

$$\frac{\partial z_j}{\partial b_j} = 1$$

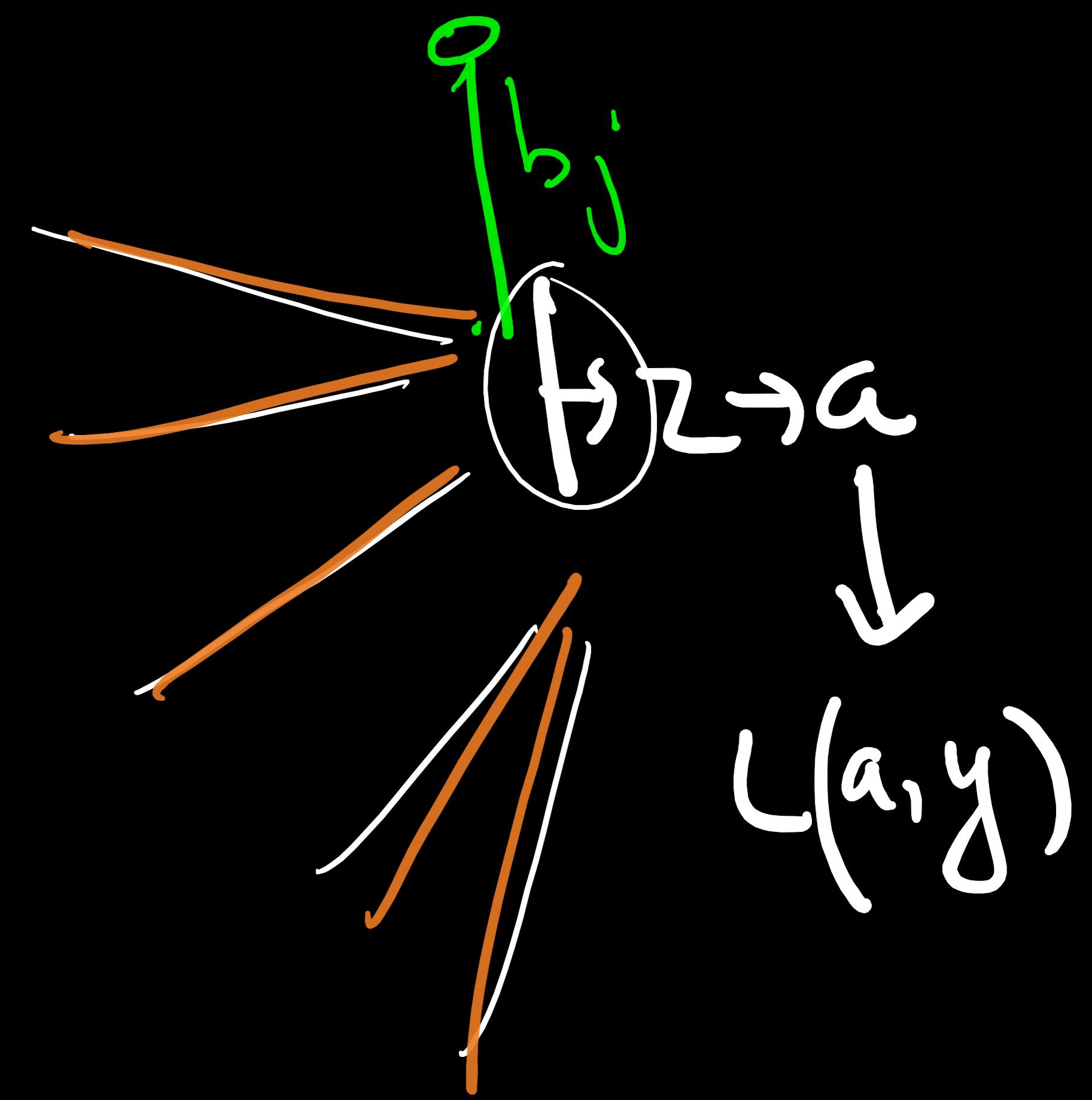
$$\frac{\partial z_j}{\partial a_i} = l-1$$

$$\frac{\partial L}{\partial w_{i,j}} = \boxed{(y_j + a_j) \sigma'(z_j)} a_i$$

$\downarrow s_{i,j}$

$$\Rightarrow \frac{\partial L}{\partial w_{i,j}} = s_j^l a_i^{l-1}$$

$$\Rightarrow \frac{\partial L}{\partial b_j} = \left(\frac{\partial L}{\partial a_j} \cdot \frac{\partial a_j}{\partial z_j} \right) \cdot \frac{\partial z_j}{\partial b_j}$$



$$\frac{\partial L}{\partial b_j} = \delta_j^L \cdot 1$$

$z_j \rightarrow a \rightarrow L$

For output layer

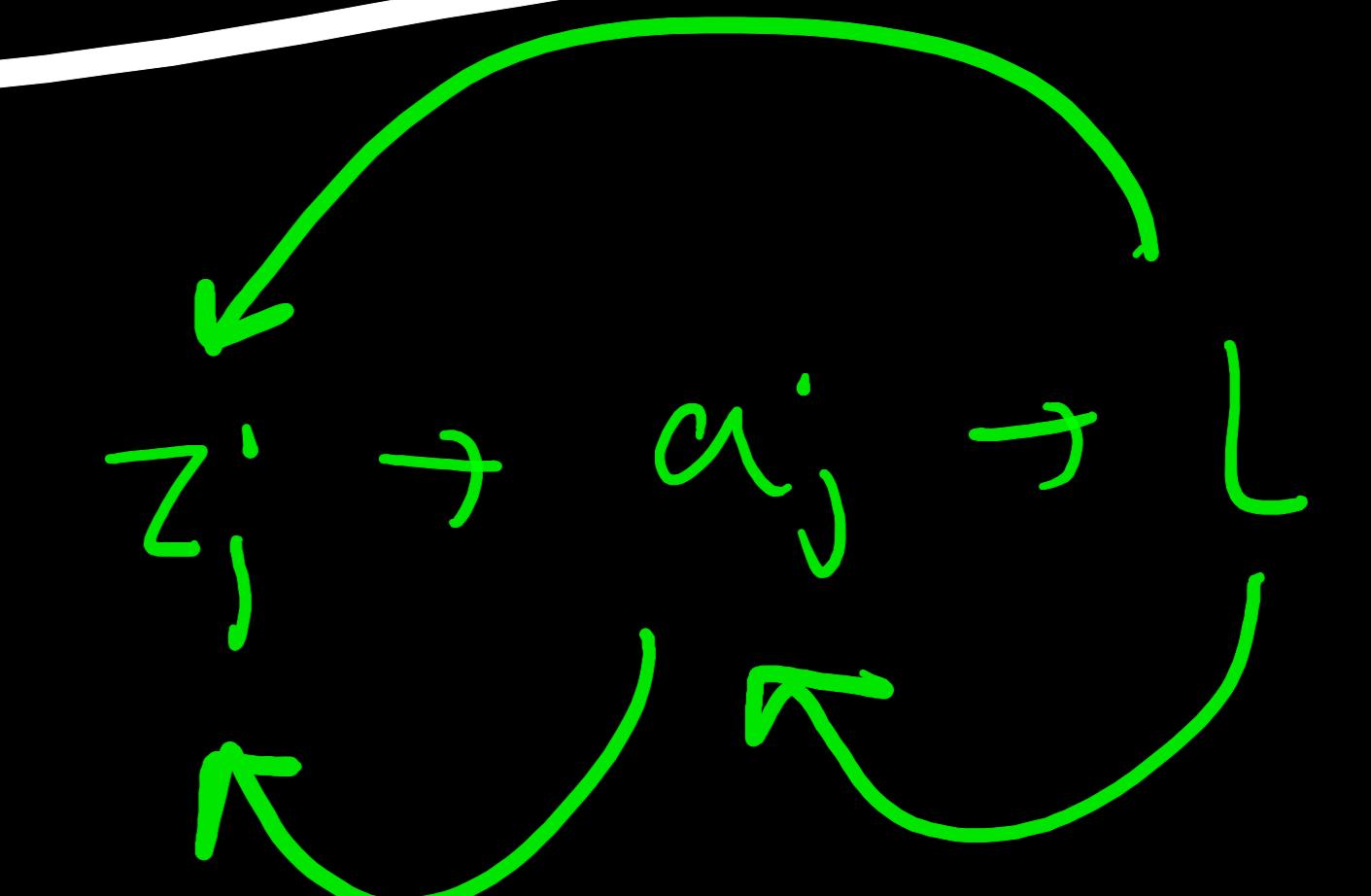
$$\therefore \frac{\partial L}{\partial z_j} = \frac{\partial L}{\partial a_j} \cdot \frac{\partial a_j}{\partial z_j}$$

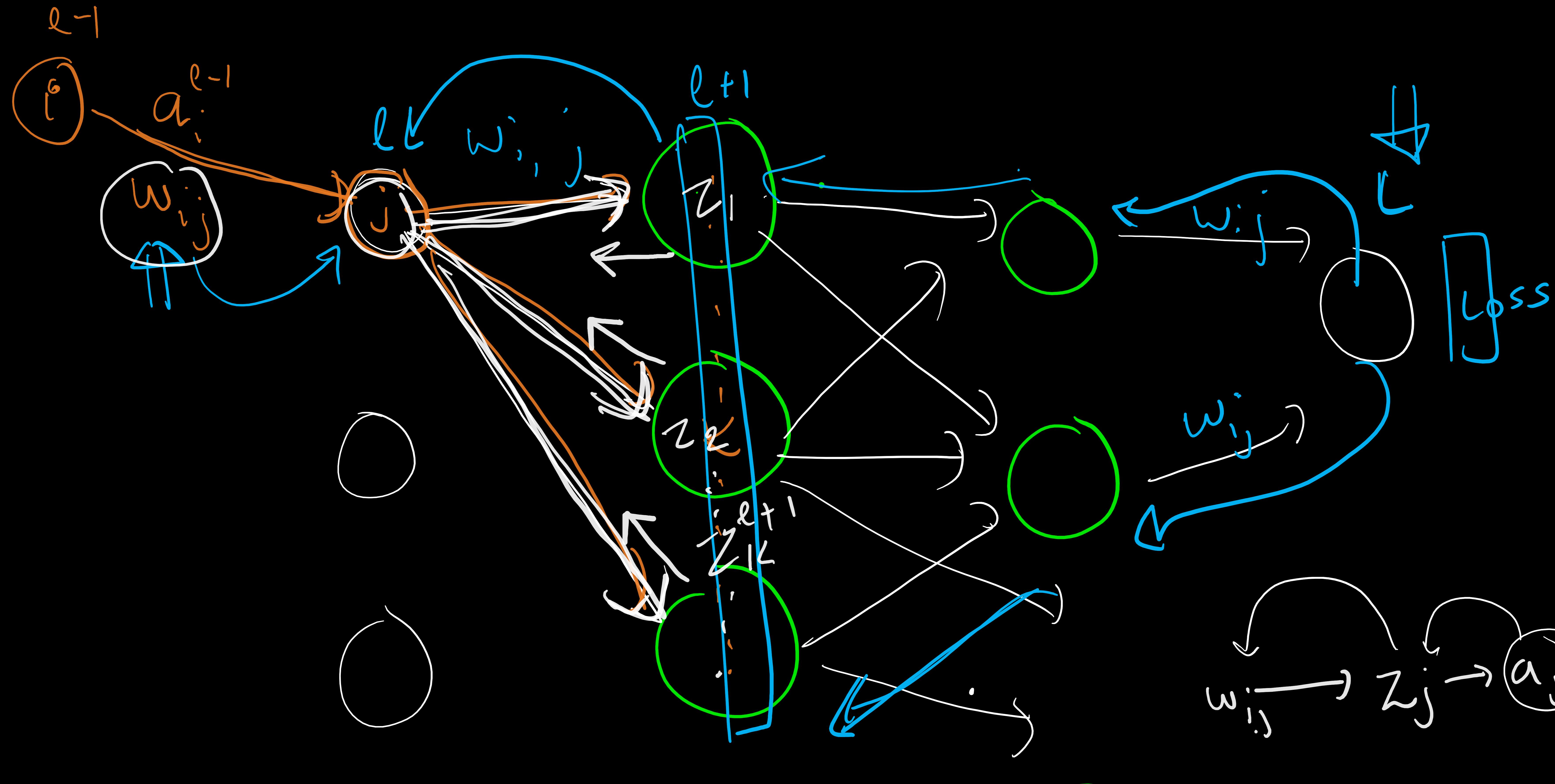
where -

$$\delta_j^L = (a_j - y_j) \circ \sigma'(z_j)$$

$$1 \quad \frac{\partial L}{\partial w_{ij}} = \delta_j^L \cdot a_i$$

$$2 \quad \frac{\partial L}{\partial b_j} = \delta_j^L \cdot 1$$





$$\frac{\partial L}{\partial w_{ij}} = \sum_k \frac{\partial L}{\partial z_k^{l+1}} \cdot \frac{\partial z_k^{l+1}}{\partial a_j} \cdot \frac{\partial a_j}{\partial z_j} \cdot \frac{\partial z_j}{\partial w_{ij}}$$

Diagram illustrating a neural network layer with nodes \$i\$, \$j\$, and \$K\$ across three layers (\$l-1\$, \$l\$, \$l+1\$).
 - Layer \$l-1\$ has node \$i\$ with value \$a_i^{l-1}\$ and weight \$w_{ij}^l\$ to node \$j\$.
 - Layer \$l\$ has node \$j\$ with value \$a_j^l = \sigma(z_j)\$ and weight \$w_{jk}^l\$ to node \$K\$.
 - Layer \$l+1\$ has node \$K\$ with value \$a_K^{l+1} = \sum_j w_{jk}^l \cdot a_j^l\$.
 - The forward pass equation is \$z_j^l = \sum_i w_{ij}^l \cdot a_i^{l-1}\$.
 - The backpropagation equation for node \$j\$ is \$\frac{\partial Z_{K+1}}{\partial a_j^l} = w_{jk}^l\$.
 - The error term for node \$j\$ is \$\delta_j^l = \sum_k (S_k \cdot w_{jk}^l) \sigma'(z_j^l)\$.

$$\begin{aligned}
 z_j^l &= \sum_i w_{ij}^l \cdot a_i^{l-1} \\
 \frac{\partial Z_{K+1}}{\partial a_j^l} &= w_{jk}^l \\
 \delta_j^l &= \sum_k (S_k \cdot w_{jk}^l) \sigma'(z_j^l)
 \end{aligned}$$

Hidden Neurons

(1)

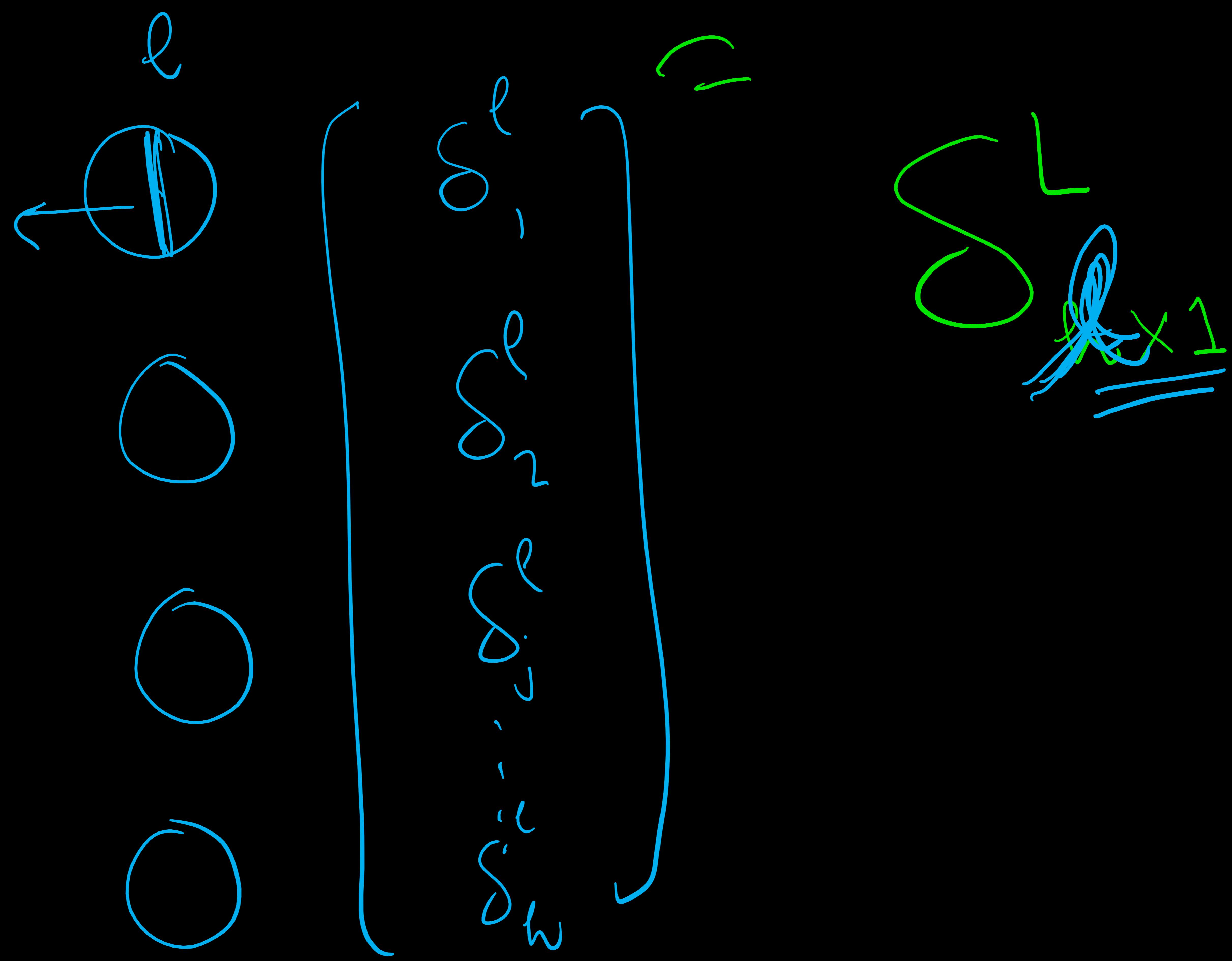
$$\frac{\partial L}{\partial w_{ij}} = \delta_j^l \cdot a_i^{l-1}$$

$$\delta_j^l = \sum_k (w_{jk} \cdot \delta_k^{l+1}) \sigma'(z_j^l)$$

(2)

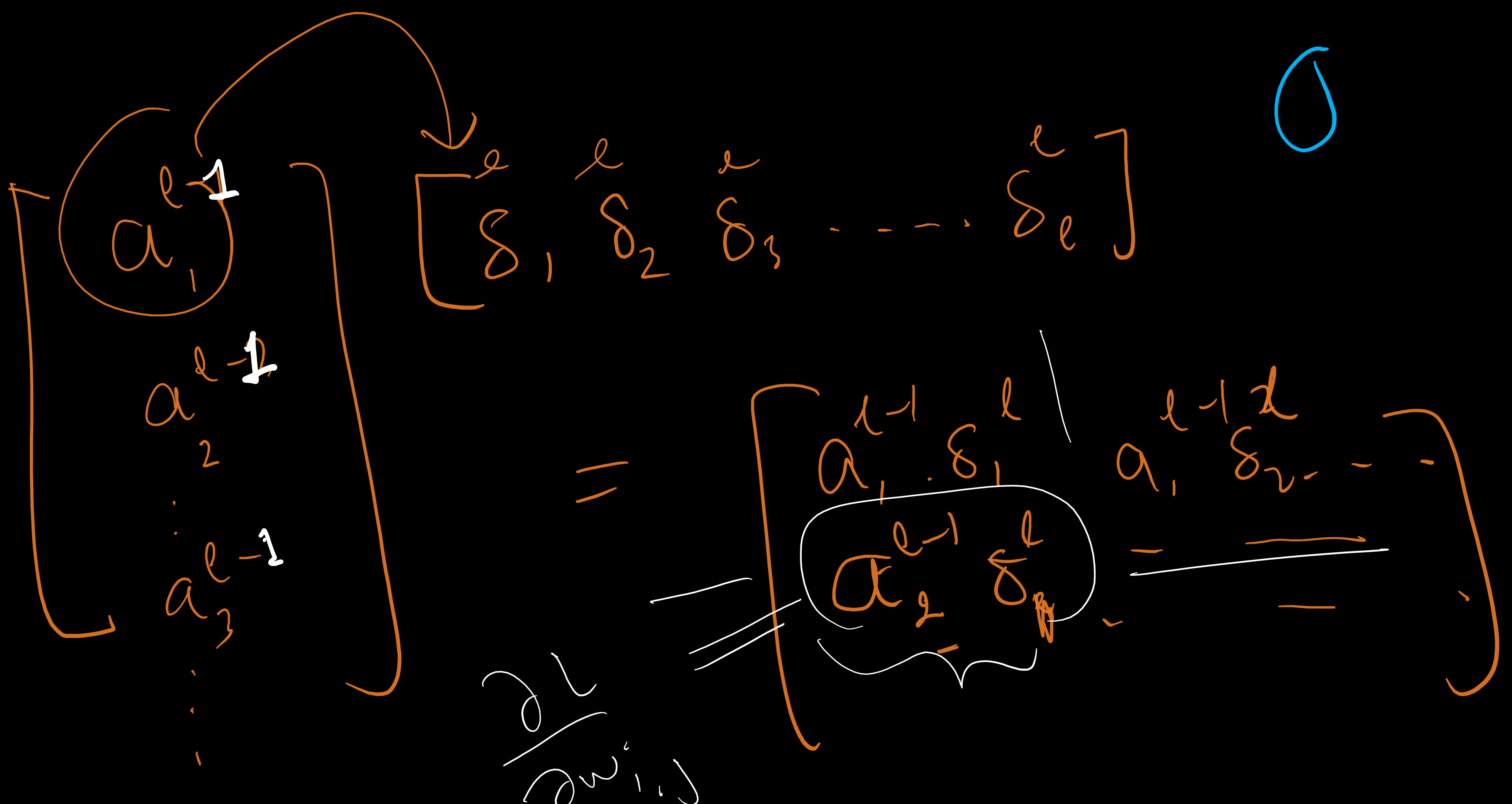
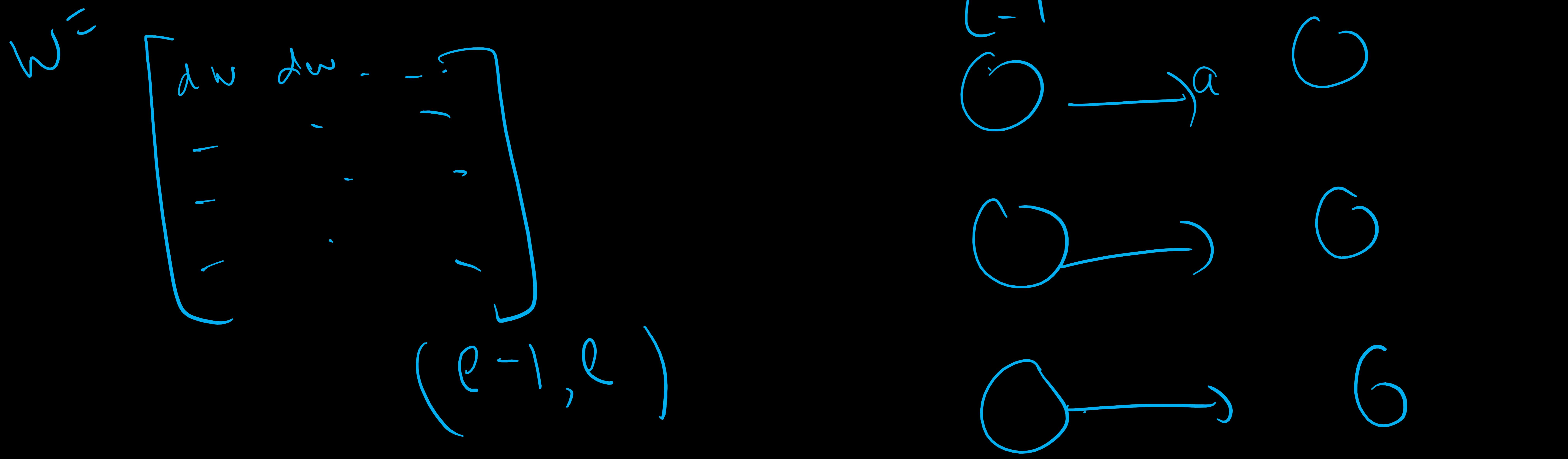
$$\frac{\partial L}{\partial b_j} = \delta_j^l$$

Matrix / Vector

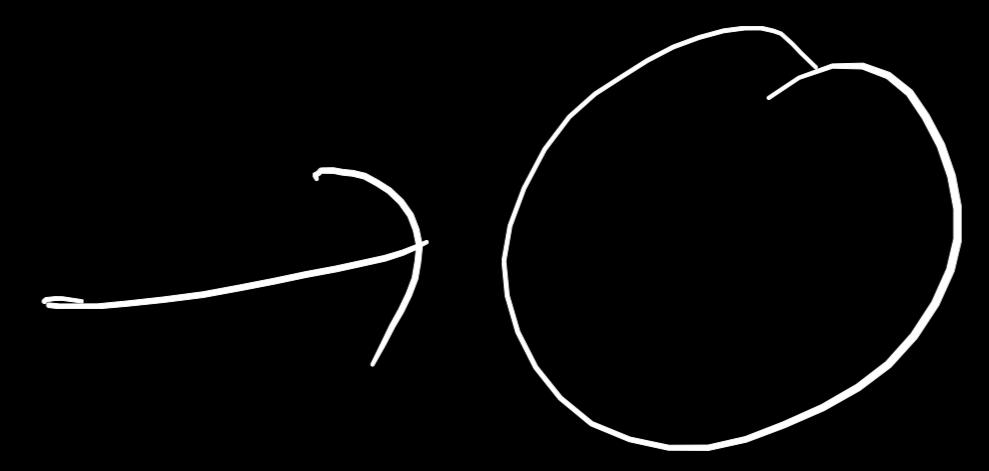


$$\frac{\partial L}{\partial w_{i,j}}$$

$$\frac{\partial L}{\partial w_i} = \underbrace{a^{l-1} \cdot (\delta^l)^T}_{(l-1, 1) \times (1 \times l)} \quad \text{Correct } \delta^{(l+1)}$$
$$(l-1, l) \iff (l-1, h)$$

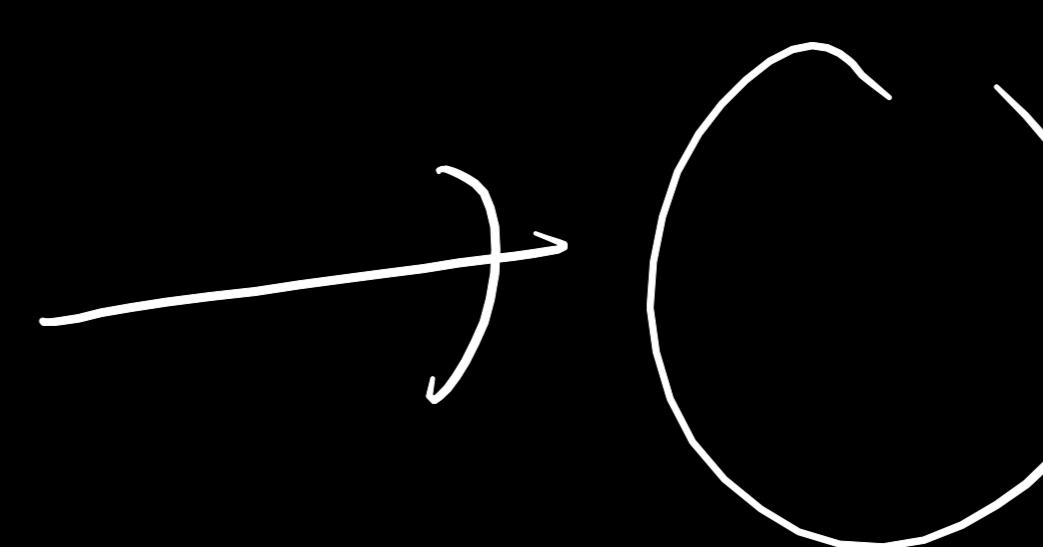
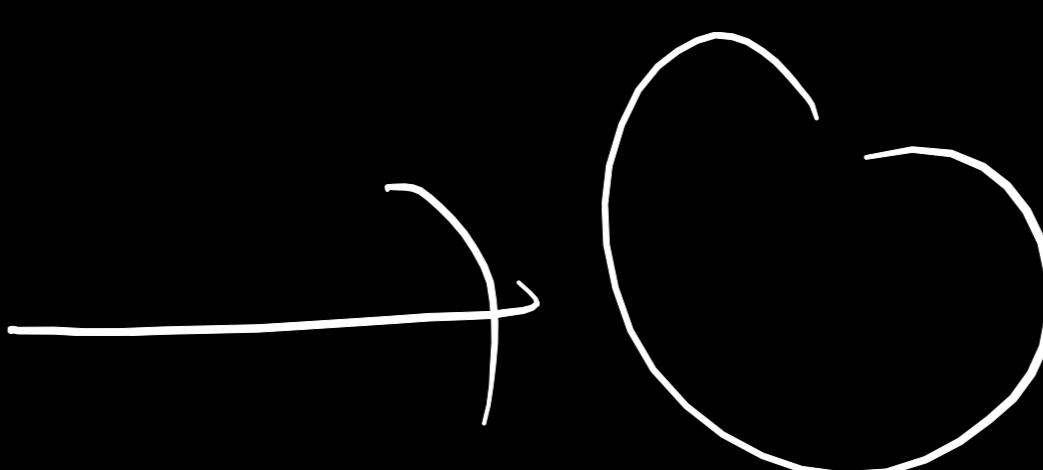
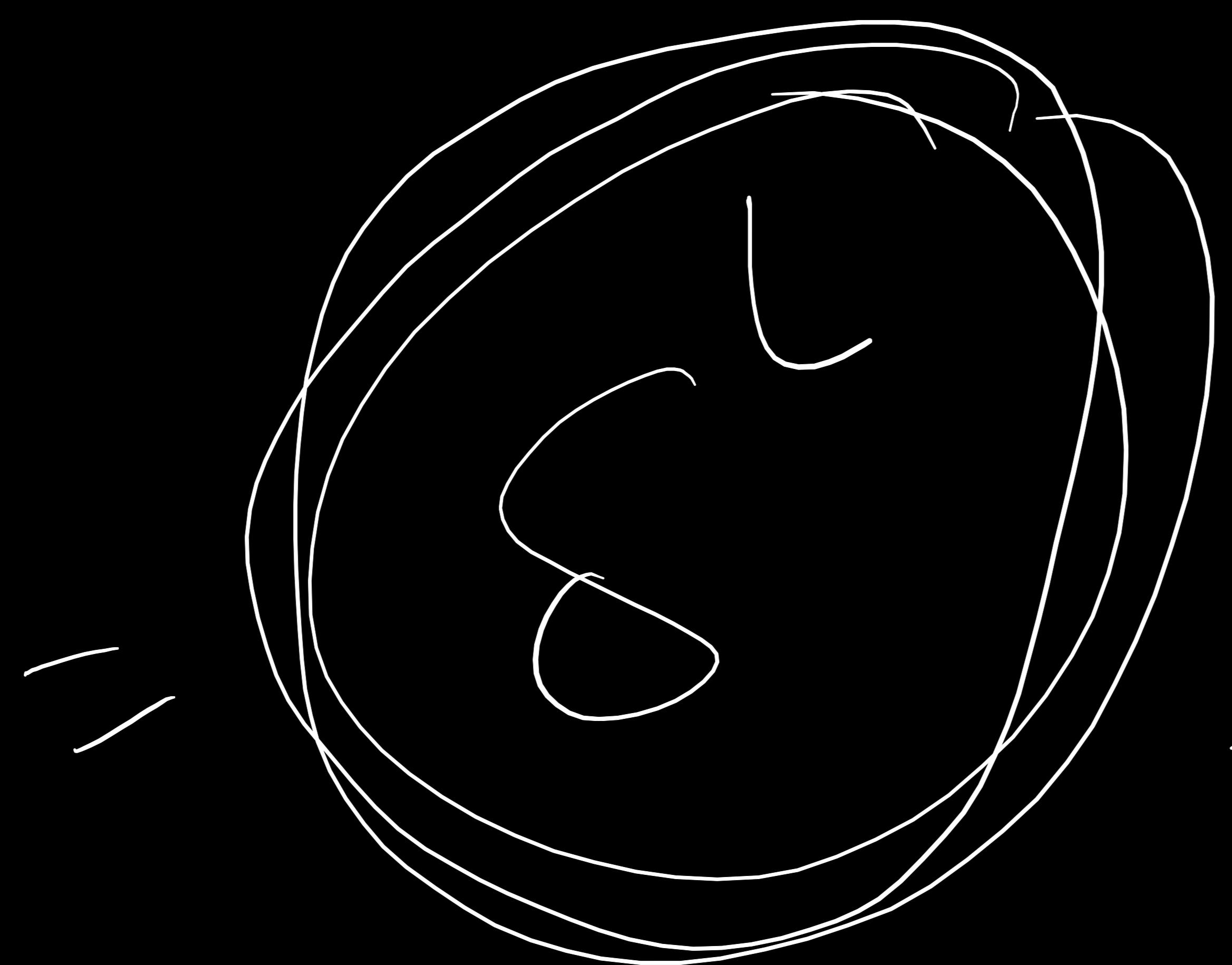


$$\frac{\partial L}{\partial b_j} =$$

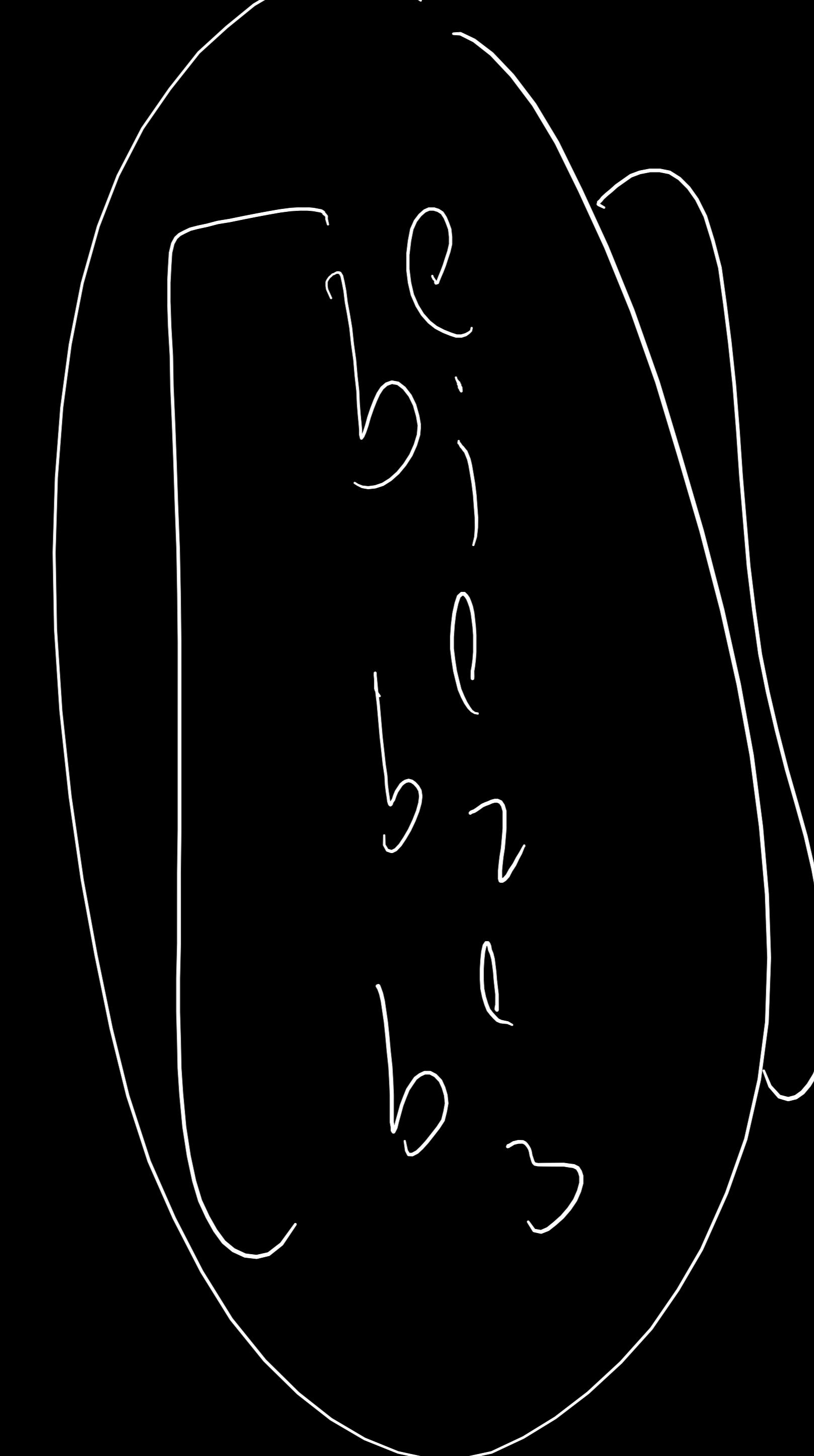


$$b = b - n_q \frac{\partial L}{\partial b}$$

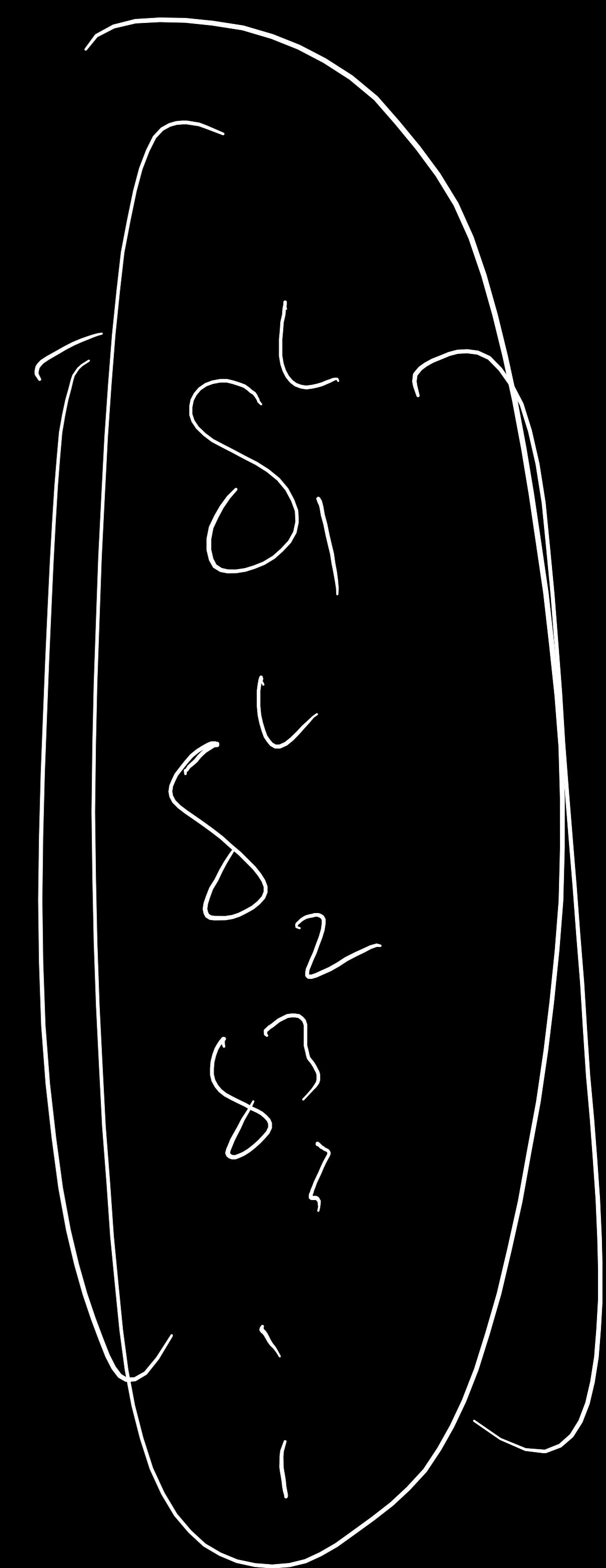
$$\frac{\partial L}{\partial b} =$$



$$\frac{\partial L}{\partial b} =$$



$$\frac{\partial L}{\partial b} =$$



$$\frac{\partial L}{\partial w_l} = \underbrace{s_{l-1}}_{\text{vector}} \cdot \underbrace{(s^T)}_{\text{matrix}} \underbrace{a^{l-1}}_{\text{vector}} \underbrace{\leftarrow \text{vector/matrix}}_{\text{vector/matrix work}} \quad \text{for both output hidden layers}$$
$$\frac{\partial L}{\partial b^e} = s^e \underbrace{\leftarrow \text{vector}}_{\text{vector/matrix}}$$

$$s^L = (a^L - y^L) \circ \sigma'(z^e)$$

$$s^e = (\omega^{et+1} \cdot s^{(et+1)}) \circ \sigma'(z^e)$$

