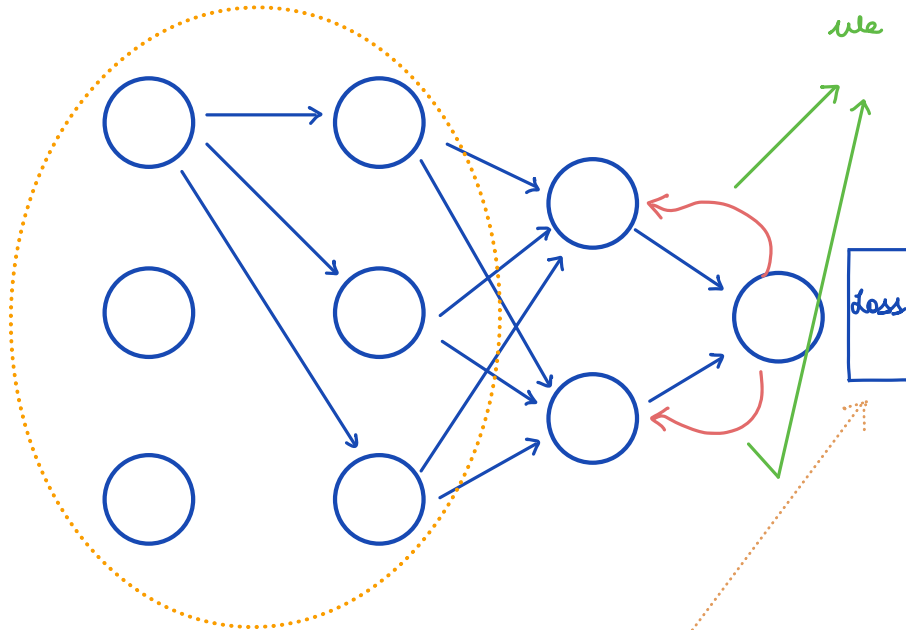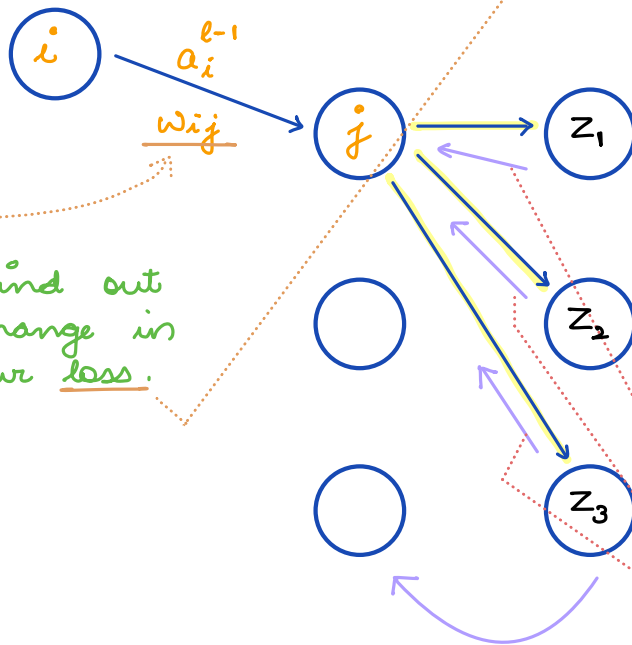# Case 2    Hidden Layer



we want to figure out how weights associated with this will get updated.

we have already seen how these weights will get updated.

$\ell-1$     $\ell$     $\ell+1$
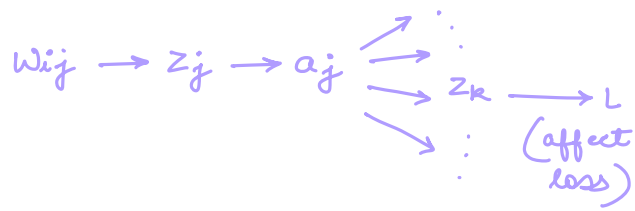
$a_i^{\ell-1}$

$W_{ij}$

we want to find out how small change in $W_{ij}$ affects our loss.

Let us assume that we know error upto this layer. we will see how to back propagate this error to previous layers.
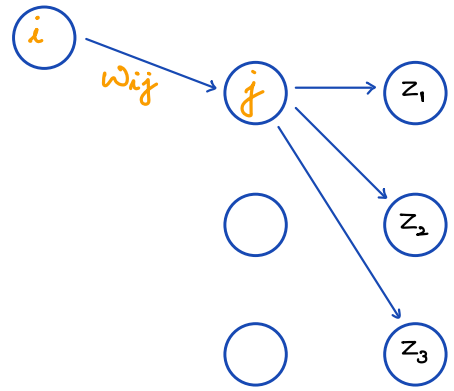
All these 3 neurons will backpropagate some loss.

$$\frac{\partial L}{\partial W_{ij}^{\ell}} = \sum_{k} \underbrace{\left(\frac{\partial L}{\partial Z_k^{\ell+1}}\right)}_{\delta_k^{\ell+1}} \cdot \frac{\partial Z_k^{\ell+1}}{\partial a_j} \cdot \underbrace{\frac{\partial a_j}{\partial z_j}}_{②} \cdot \underbrace{\frac{\partial z_j}{\partial W_{ij}}}_{③}$$

①

$$W_{ij} \rightarrow z_j \rightarrow a_j \nearrow \vdots$$
$$\rightarrow z_k \rightarrow L$$
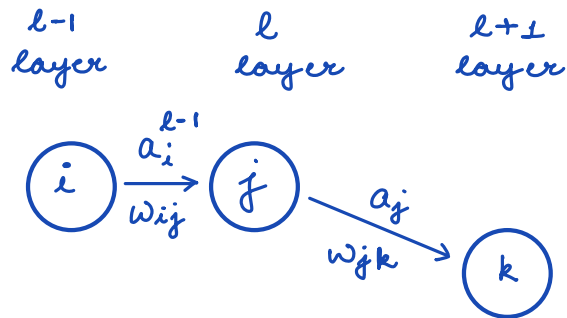$$\searrow \vdots \quad (\text{affect loss})$$

how loss changes according to $W_{ij}$ will depend on how loss changes according to $z_k$ * how $z_k$ of $l+1$ changes according to $a_j$ * how $a_j$ changes according to $z_j$ * how $z_k$ changes according to $W_{ij}$.



---

$l-1$ layer    $l$ layer    $l+1$ layer



$$Z_k^{l+1} = \sum_j W_{jk} \, a_j^l$$

① $\dfrac{\partial Z_k^{l+1}}{\partial a_j^l}$

$$Z_k^{l+1} = \sum_j W_{jk} \cdot a_j^l$$

$$\dfrac{\partial Z_k^{l+1}}{\partial a_j^l} = W_{jk}$$

② $\dfrac{\partial a_j}{\partial z_j}$

$$a_j = \sigma(z_j)$$

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

③ $\dfrac{\partial z_j}{\partial W_{ij}}$

$$z_j = \sum_i \omega_{ij} \, a_i^{\ell-1}$$

$$\frac{\partial z_j}{\partial \omega_{ij}} = a_i^{\ell-1}$$

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

$$= \boxed{\sum_k \left( \delta_k^{\ell+1} \cdot \omega_{jk} \right) \cdot \sigma'(z_j^\ell)} \cdot a_i^{\ell-1}$$

$$\downarrow$$

$$\delta_j^\ell = \sum_k \left( \delta_k \cdot \omega_{jk} \right) \cdot \sigma'(z_j^\ell)$$

$$\frac{\partial L}{\partial \omega_{ij}} = \delta_j^\ell \cdot a_i^{\ell-1}$$

Bias Update Rule

$$\frac{\partial L}{\partial b} = \boxed{\sum_k \frac{\partial L}{\partial z_k^{\ell+1}} \cdot \frac{\partial z_k^{\ell+1}}{\partial a_j} \cdot \frac{\partial a_j}{\partial z_j}} \cdot \frac{\partial z_j}{\partial b}$$

$$\delta_j^\ell$$

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

$$\delta_j^\ell = \sum_k \left( \delta_k^{\ell+1} \cdot \omega_{jk}^{\ell+1} \right) \odot \sigma'(z_j^\ell)$$

# Final Result for Hidden Neuron:

$$\frac{\partial L}{\partial w_{ij}} = \delta_j^\ell \, a_i^{\ell-1}$$
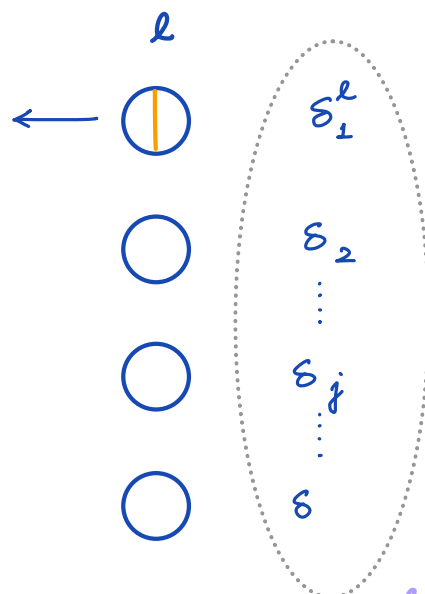
$$\frac{\partial L}{\partial b} = \delta_j^\ell$$

$$\delta_j^\ell = \sum_k \left( w_{jk}^{\ell+1} \, \delta_k^{\ell+1} \right) \odot \sigma'(z_j^\ell)$$

Results are same for output and hidden layer, only the value of $\delta_j^\ell$ will change.

## Matrix Representation:

While writing code we won't use for loop to iterate over all neuron, instead we will take matrix and vector and do dot product. Like this we can do parallel computation

$\delta_j^\ell$ = loss associated with one unit



$\ell$

$\delta_1^\ell$

$\delta_2$

$\vdots$

$\delta_j$

$\vdots$

$\delta$

These neuron have some associated loss which they will propagate back

This entire vector is called $\delta_{\ell \times 1}^L$ and has dimension of $\ell \times 1$

$L^{th}$ layer and $\ell$ units

$\ell$ = no. of hidden layer

Earlier we calculated $\frac{\partial L}{\partial w_{ij}}$ , now we want to calculate $\frac{\partial L}{\partial w^l}$ ?
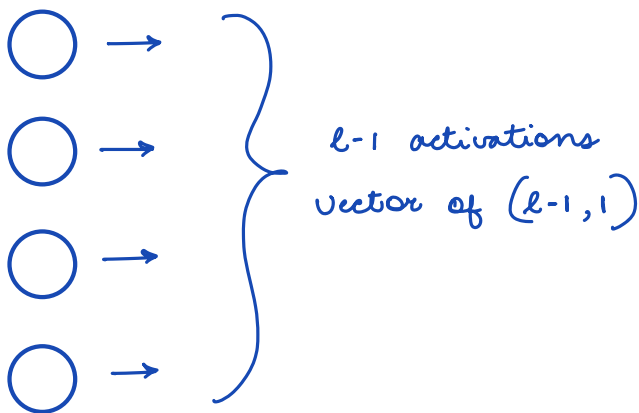
$$\frac{\partial L}{\partial w^l} = a^{l-1} \cdot (\delta^l)^T$$

- $\frac{\partial L}{\partial w^l}$ : $(l-1, l)$
- $a^{l-1}$ : $(l-1, 1)$
- $(\delta^l)^T$ : $(l, 1) \rightarrow (1, l)$
- $a^{l-1} \cdot (\delta^l)^T$ : $(l-1, l)$

$$W = \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix} \quad (l-1, l)$$

$L-1$ layer

for activation dimension:



$l-1$ activations vector of $(l-1, 1)$

why the above formula is correct ?

activations of previous layer $\rightarrow$

$$\begin{bmatrix} a_1^{l-1} \\ a_2^{l-1} \\ a_3^{l-1} \\ \vdots \end{bmatrix} \begin{bmatrix} \delta_1^l & \delta_2^l & \delta_3^l & \cdots & \delta_l^l \end{bmatrix}$$

$$\begin{bmatrix} a_1^{l-1} \delta_1^l & a_1^{l-1} \delta_2^l & a_3^{l-1} \delta_3^l \\ a_2^{l-1} \delta_1^l & a_2^{l-1} \delta_2^l & \cdots \\ \vdots & & \\ \vdots & & \end{bmatrix}$$

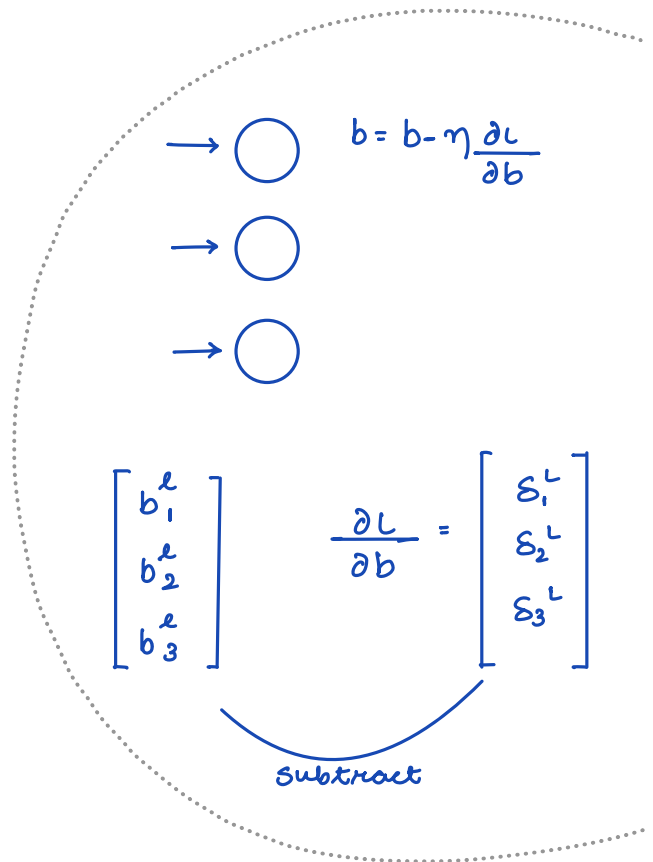Each element is $\frac{\partial L}{\partial w_{ij}}$

$$\frac{\partial L}{\partial w^l} = a^{l-1} \cdot (\delta^l)^T \quad \Bigg\} \text{ This formula is true both for input layer and output layer}$$

for biases, $\dfrac{\partial L}{\partial b_j^l} = \delta_j^l$

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

↓ <span style="color:magenta">for entire bias matrix</span>

<span style="color:blue">$b = b - \eta \dfrac{\partial L}{\partial b}$</span>

$$\begin{bmatrix} b_1^l \\ b_2^l \\ b_3^l \end{bmatrix} \quad \frac{\partial L}{\partial b} = \begin{bmatrix} \delta_1^L \\ \delta_2^L \\ \delta_3^L \end{bmatrix}$$

<span style="color:blue">subtract</span>

We have seen how to compute $\dfrac{\partial L}{\partial w^l}$ and $\dfrac{\partial L}{\partial b^l}$

<span style="color:magenta">→ vector/matrix notation</span>

<span style="color:orange">vector</span>

$$\frac{\partial L}{\partial w^l} = \underbrace{a^{l-1} (\delta^l)^T}_{\text{matrix}} \quad \Bigg\} \text{ true for both output and hidden layer.}$$

$$\frac{\partial L}{\partial b^l} = \delta^L$$
↓ <span style="color:orange">vector</span>

<span style="color:orange">Element wise Product (it scales up $a^L - y^L$ vector)</span>

$$\delta^L = \underbrace{(a^L - y^L)}_{\text{vector}} \odot \underbrace{\sigma'(z^l)}_{\text{vector}}$$

You find out the difference between prediction and actual value and you just scale that difference by multiplying with derivative of $\underset{\rightarrow z^\ell}{\underline{\text{activation } f x^n}}$ for a given input.

$$\delta^\ell = \left( \omega^{\ell+1} \, \delta^{(\ell+1)} \right) \odot \sigma'(z^\ell) \quad \Big\} \text{ for hidden layer}$$