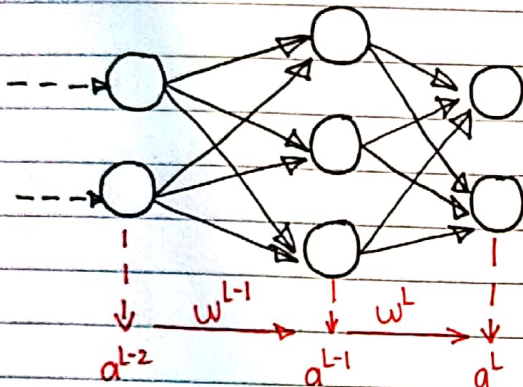# Back propagation in neural network explained:

↳ Let say we have the following neural network:



→ Let say we have the activated outputs at each layer are $a^L$, $a^{l-1}$ and $a^{l-2}$

* Activated output = output after activation.

→ Example: To calculate activated output at layer L, we have:

$$\circledast \begin{cases} z^L = a^{l-1} \cdot w^L \quad ① \\ a^L = \sigma(z^L) \quad ② \end{cases}$$

·) Where $z^L$ in ① calculates unactivated output, $a^{l-1}$ is activated output in previous layer, $w^L$ is the set of weights at layer L. In ②, $\sigma(\cdot)$ represents the activation function.

## ① Gradients:

↳ To perform gradient descent, we need the gradient of a loss function L with respect to all sets of weights $w^L$, $w^{l-1}$, ...

↳ Let say our loss function is a simple squared difference loss, let's calculate the gradients for $w^L$.

·) $\dfrac{\partial L}{\partial w^L} = \dfrac{\partial}{\partial w^L}[(y-a^L)^2]$  (where y is the labels).

·) $\dfrac{\partial L}{\partial w^L} = \dfrac{\partial L}{\partial a^L} \cdot \dfrac{\partial a^L}{\partial z^L} \cdot \dfrac{\partial z^L}{\partial w^L}$  (chain rule in differentiation).

↳ Let's examine each component:

① $\dfrac{\partial L}{\partial a^L} = \dfrac{\partial}{\partial a^L}[(y-a^L)^2] = -2(y-a^L) = 2(a^L-y)$.

② $\dfrac{\partial a^L}{\partial z^L} = \dfrac{\partial}{\partial z^L}[\sigma(z^L)] = \sigma'(z^L)$.

③ $\dfrac{\partial z^L}{\partial w^L} = \dfrac{\partial}{\partial w^L}[(a^{l-1} \cdot w^L)] = a^{l-1}$

⟹ From the above:  $\dfrac{\partial L}{\partial w^L} = 2(a^L-y)\sigma'(z^L)a^{l-1}$  #

## ② Back-propagation.

⌐ Okay now we know $\frac{\partial L}{\partial w^L}$. However, we need to know $\frac{\partial L}{\partial w^k}$ with $k$ being any layer in our neural network. Let's state the formula again:

$$\cdot) \quad \frac{\partial L}{\partial w^k} = \frac{\partial L}{\partial a^k} \cdot \frac{\partial a^k}{\partial z^k} \cdot \frac{\partial z^k}{\partial w^k} \quad \text{④}$$

⌐ We know $\frac{\partial a^k}{\partial z^k} = \sigma'(z^k)$ is just the derivative of the activation function

and $\frac{\partial z^k}{\partial w^k} = a^{k-1}$ will just always be the output from previous layer.

⌐ So to calculate ④ we need $\frac{\partial L}{\partial a^k}$, let's find it using induction:

1. Base case: $\frac{\partial L}{\partial a^L} = 2(a^L - y)$ ← *We already know this

2. $\frac{\partial L}{\partial a^{L-1}} = \frac{\partial L}{\partial a^L} \cdot \frac{\partial a^L}{\partial z^L} \cdot \frac{\partial z^L}{\partial a^{L-1}} = \boxed{\frac{\partial L}{\partial a^L}} \cdot \sigma'(z^L) \cdot w^L$

   *We know this →

3. $\frac{\partial L}{\partial a^{L-2}} = \frac{\partial L}{\partial a^{L-1}} \cdot \frac{\partial a^{L-1}}{\partial z^{L-1}} \cdot \frac{\partial z^{L-1}}{\partial a^{L-2}} = \boxed{\frac{\partial L}{\partial a^{L-1}}} \cdot \sigma'(z^{L-1}) \cdot w^{L-1}$

$$\Rightarrow \boxed{\frac{\partial L}{\partial a^k} = \frac{\partial L}{\partial a^{k+1}} \cdot \sigma'(a^{k+1}) \cdot w^{k+1}} \qquad (\text{for } k < L)$$

⟹ Final note: Since as long as we know $\frac{\partial L}{\partial a^{k+1}}$, we can always iteratively propagate our gradients backwards to calculate the gradients for previous layers, it is called "Back-propagation".