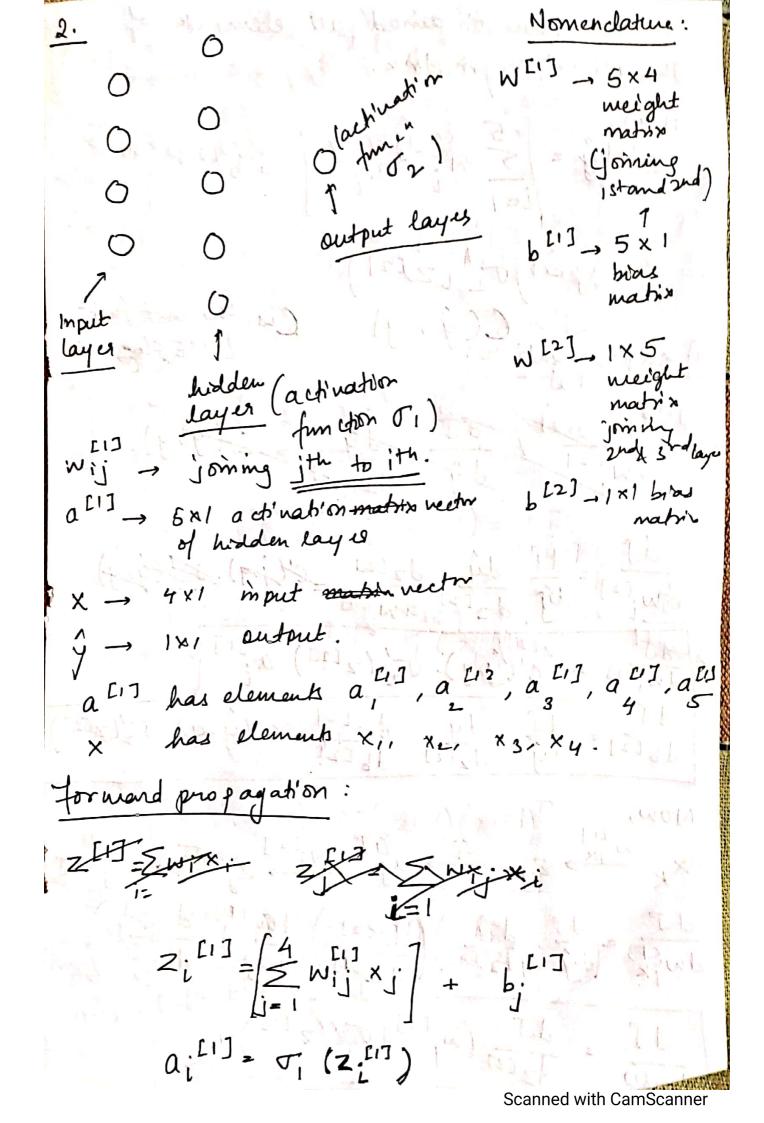1. Forward Propagation →

Forward propagation refers to the process of calculating the values of all the activation neurons and subsequently, the output of a neural network, for a given set of weights and bias. The term earns its name from the fact that we start with the input values, which are feed-forwarded to the next layer and so on, until we reach the last layer.

Backward Propagation:

Backward propagation refers to the process of calculating the gradient of the cost function with respect to a given set of weights and bias. These gradients are then used to optimize the value of weights and bias. The algorithm makes use of chain rule to compute the gradient. The algorithm earns its name from the fact that we compute the gradients layer by layer, starting from the last one and then iterating "backwards".

**2.**

$W^{[1]} \rightarrow 5 \times 4$
weight matrix
(joining 1st and 2nd)

$b^{[1]} \rightarrow 5 \times 1$
bias matrix

Input layer

output layer

O (activation fun." $\sigma_2$)

hidden layer (activation function $\sigma_1$)

$W^{[2]} \rightarrow 1 \times 5$
weight matrix
joining 2nd & 3rd layer

$b^{[2]} \rightarrow 1 \times 1$ bias matrix

$W_{ij}^{[1]} \rightarrow$ joining $j$th to $i$th.

$a^{[1]} \rightarrow 5 \times 1$ activation matrix vector of hidden layer

$x \rightarrow 4 \times 1$ input matrix vector

$\hat{y} \rightarrow 1 \times 1$ output.

$a^{[1]}$ has elements $a_1^{[1]}, a_2^{[1]}, a_3^{[1]}, a_4^{[1]}, a_5^{[1]}$

$x$ has elements $x_1, x_2, x_3, x_4$.

**Forward propagation :**

$z^{[1]} = \sum w_j^T x_i$    $z_i^{[1]} = \sum_{i=1} w_{ij} x_i$

$$z_i^{[1]} = \left[ \sum_{j=1}^{4} w_{ij}^{[1]} x_j \right] + b_j^{[1]}$$

$$a_i^{[1]} = \sigma_1 (z_i^{[1]})$$

So, we have obtained all elements of the activation matrix.

Now,

$$z^{[2]} \, \hat{y} = \left[ \sum_{i=1}^{5} w_i^{[2]} a_i^{[1]} \right] + b^{[2]}$$

$$\hat{y} = \sigma_2 (z^{[2]})$$

$$J = \mathcal{C}(\hat{y}, y) \qquad \mathcal{C} \text{ is the cost func}^n$$
$$(MSE \text{ log}, etc.)$$

Back propagation:

$$\frac{dJ}{dw_i^{[2]}} = \frac{dJ}{d\hat{y}} \quad \frac{d\hat{y}}{dw_i^{[2]}} = \mathcal{C}'(\hat{y}, y).$$

$$\frac{dJ}{dw_i^{[2]}} = \frac{dJ}{d\hat{y}} \frac{d\hat{y}}{dz^{[2]}} \frac{dz^{[2]}}{dw_i^{[2]}} = \mathcal{C}'(\hat{y}, y) . \sigma_2'(\hat{y}, y)$$

$$\Rightarrow \boxed{\frac{dJ}{dw_i^{[2]}} = \mathcal{C}'(\hat{y}, y) . \sigma_2'(z^{[2]}) \, a_i^{[1]}}$$

$$\boxed{\frac{dJ}{db^{[2]}} = \frac{dJ}{d\hat{y}} . \frac{d\hat{y}}{dz^{[2]}} . \frac{dz^{[2]}}{db^{[2]}} = \mathcal{C}'(\hat{y}, y) . \sigma_2'(z^{[2]})}$$

Now,

$$x_k \xrightarrow{w_{jk}^{[1]}} z_j \xrightarrow{\sigma_1} a_j^{[1]} \xrightarrow{w_j^{[2]}} z_j^{[2]} \xrightarrow{\sigma_2} a\hat{y} \rightarrow J$$

$$\frac{dJ}{dw_{jk}^{[1]}} = \left( \frac{dJ}{d\hat{y}} . \frac{d\hat{y}}{dz^{[2]}} \right) \left( \frac{dz^{[2]}}{da_j^{[1]}} \right) \left( \frac{da_j^{[1]}}{dz_j^{[1]}} \frac{dz_j^{[1]}}{dw_{jk}^{[1]}} \right)$$

$$\boxed{\frac{dJ}{dw_{jk}^{[1]}} = \frac{dJ}{dz^{[2]}} \left( w_j^{[2]} \right) \sigma_1'(z_j^{[1]}) x_k}$$

$$\frac{dJ}{dz^{[2]}} = C'(\hat{y}, y) \cdot \sigma_2'(z^{[2]}) = \frac{dJ}{db^{[2]}}$$

$$\frac{dJ}{db_j^{[1]}} = \frac{dJ}{dz^{[2]}} \cdot W_{0j}^{[2]} \cdot \sigma_1'(z_j^{[1]})$$

Now, vectorizing them:

$$\frac{dJ}{dW^{[2]}} = C'(\hat{y}, y) \cdot \sigma_2'(z^{[2]})(a^{[1]})^T$$

$$\frac{dJ}{db^{[2]}} = C'(\hat{y}, y) \cdot \sigma_2'(z^{[2]}) \cdot \frac{dJ}{dz^{[2]}}$$

$$\frac{dJ}{dW^{[1]}} = (W^{[2]T}) \cdot \frac{dJ}{dz^{[2]}}$$

$$\frac{dJ}{dw^{[1]}} = \left(\frac{dJ}{dz^{[2]}}\right)\left((w^{[2]})^T * \sigma_1'(z^{[1]})\right) \times F$$

element wise
multiplication

$$\frac{dJ}{db^{[1]}} = \left(\frac{dJ}{dz^{[2]}}\right)\left((w^{[2]})^T * \sigma_1'(z^{[1]})\right)$$

**3. (a) Sigmoid:**

$$f(x) = \frac{1}{1 + e^{-x}}$$

$$f'(x) = \frac{e^{-x}}{(1 + e^{-x})^2} = f(x) \times (1 - f(x))$$

**(b) ReLU:**

$$f(x) = \begin{cases} x & x \geq 0 \\ 0 & x < 0 \end{cases} \quad \text{or } f(x) = max(0, x)$$

$$\Rightarrow f'(x) = \begin{cases} 1 & x \geq 0 \\ 0 & x < 0 \end{cases}$$

**(c) Leaky ReLU :**

$$f(x) = \begin{cases} x & x \geq 0 \\ 0.01\, x & x < 0 \end{cases}$$

$$f'(x) = \begin{cases} 1 & x \geq 0 \\ 0.01 & x < 0 \end{cases}$$

**(d) Tanh :**

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

$$\Rightarrow f'(x) = \frac{(e^x + e^{-x})^2 - (e^x - e^{-x})^2}{(e^x + e^{-x})^2}$$

$$\ast\, f'(x) = 1 - (f(x))^2$$

(e) Softmax:

$$\vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ | \\ x_n \end{bmatrix}$$

$\sigma(\vec{x})$ is also a vector of dimensions $n \times 1$.

$$\boxed{[\sigma(\vec{x})]_i = \frac{e^{x_i}}{\sum\limits_{j=1}^{n} e^{x_j}}}$$

1st case: derivative wrt $x_k$, $K \neq i$.

$$\therefore \frac{d[\sigma(\vec{x})]_i}{dx_k} = -\frac{e^{x_i} \cdot e^{x_k}}{\left(\sum\limits_{j=1}^{n} e^{x_j}\right)^2} = -[\sigma(\vec{x})]_i \, [\sigma(\vec{x})]_k$$

2nd case: derivative wrt $x_k$, $K = i$

$$\therefore \frac{d[\sigma(\vec{x})]_i}{dx_{ki}} = \frac{\left[\left(\sum\limits_{j=1}^{n} e^{x_j}\right) - (e^{x_i})\right](e^{x_i})}{d\left(\sum\limits_{j=1}^{n} e^{x_j}\right)^2}$$

$$= [\sigma(\vec{x})]_i \left(1 - [\sigma(\vec{x})]_i\right)$$

$$\boxed{\frac{d[\sigma(\vec{x})]_i}{dx_k} = \begin{cases} [\sigma(\vec{x})]_i \, (1 - \sigma(\vec{x}))_k \, , & i = k \\ -[\sigma(\vec{x})]_i \, [\sigma(\vec{x})]_k \, , & i \neq k \end{cases}}$$