points :

if weight matrix $W$ is initialized too large, the output of the matrix multiply too could probably have a very large range, which in turn will make all the outputs in the vector $z$ almost binary either '1' or '0' .

1. sigmoid, $\sigma(z) = \dfrac{1}{1+e^{-z}}$

2. local gradient, is a derivative that measures how much the output of the sigmoid function changes with respect to its input $z$.

$$\sigma'(z) = \dfrac{d}{dz}\left(\dfrac{1}{1+e^{-z}}\right)$$

$f(g(z))$; where $f(u) = \dfrac{1}{u}$ , $g(z) = 1+e^{-z}$

$f'(u) = -\dfrac{1}{u^2}$    $g'(z) = \dfrac{d}{dz}(1+e^{-z}) = -e^{-z}$

$\dfrac{d\sigma(z)}{dz} = f'(g(z)) \cdot g'(z) = -\dfrac{1}{(1+e^{-z})^2} \cdot (-e^{-z})$

$\dfrac{d\sigma(z)}{dz} = \dfrac{e^{-z}}{(1+e^{-z})^2}$

we know, $\sigma(z) = \dfrac{1}{1 + e^z}$

re written as, Subtract the sigmoid from 1.

$$1 - \sigma(z) = 1 - \dfrac{1}{1 + e^{-z}}$$

express 1 with common denominator, $1 = \dfrac{1 + e^{-z}}{1 + e^{-z}}$

$$1 - \sigma(z) = \dfrac{1 + e^{-z}}{1 + e^{-z}} - \dfrac{1}{1 + e^{-z}}$$

$$1 - \sigma(z) = \dfrac{(1 + e^{-z}) - 1}{(1 + e^{-z})} = \dfrac{e^{-z}}{(1 + e^{-z})(1 + e^{-z})} \quad \times (1 - \sigma(z))$$

$$\boxed{\dfrac{d\sigma(z)}{dz} = \sigma(z) \times (1 - \sigma(z))}$$

if $\sigma(z)$ is close to 1 or 0, the term $\sigma(z) \times (1 - \sigma(z))$ becomes very small approaching to zero.

In;

→ Back propagation, gradients are propagated backward through the network by multiplying them together (as per the chain rule)

→ if local gradient is small, when you multiply it with other gradients in the backward pass, the overall gradient becomes effectively zero.

→ This causes the entire backward pass to have nearly zero gradients, preventing the network from learning.

# ReLU (Non-linear) activation Function.

$$ReLU(z) = \max(0, z) \quad \begin{cases} z & \text{if } z > 0 \\ 0 & \text{if } z \leq 0 \end{cases}$$

z = np.maximum (0, np.dot (W, x))   #forward pass
      <sub>min</sub>   <sub>max</sub>

dW = np.outer (z > 0, x)   # representing backward pass:
            local gradient for w.

(<span style="color:red">↓ how much each,<br>weight should be adjusted )</span>  <span style="color:red">↑ which elements in<br>'z' are greater than zero that pass through ReLU</span>

Calculates dot product between weight matrix 'w' and input vector 'x'.

→ if any negative value it return '0'.

→ np.outer (z > 0, x) : Outer product (z > 0) (a vector of 'True' or 'False') and the input vector x.

The outer product is used to determine how much each weight 'w' should be adjusted based on the input 'x' and the result 'z'.

## Chain Rule:
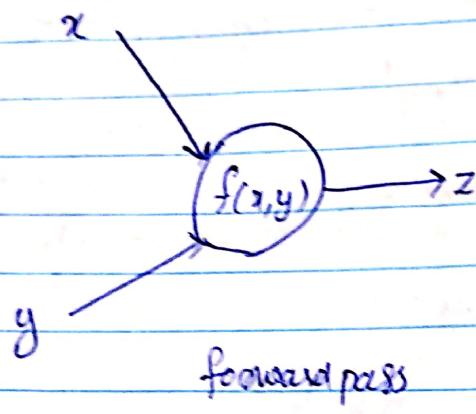
$$F(x) = f(g(x)) \text{ for all } x,$$

$$F'(x) = f'(g(x)) \cdot g'(x)$$

Leibniz notation, $\dfrac{dz}{dx} = \dfrac{dz}{dy} \cdot \dfrac{dy}{dx}$   z dependent on y<br>                   y   "   on x

Sigmoid function :- $z = \dfrac{1}{(1 + np.exp(-np.dot(W, x)))}$   #forward pass

Computes the gradient db ; dx = np.dot (W.T, z*(1-z)))   # backward<br>the z w.r.t input x                   local<br>            ↓    ↓  pass gradient x<br>         (transpose db (derivative db<br>         weight matrix)   sigmoid function)

output z w.r.t to<br>gradients , w,  dW = np.outer (z*(1-z), x)   #backward<br>                ↓      pass,<br>             Sigmoid.    local gradient w.

$$\frac{dL}{dx} = \frac{dL}{dz} \frac{dz}{dx}$$

$$\frac{dL}{dz}$$

$$\frac{dL}{dy} = \frac{dL}{dz} \frac{dz}{dy}$$

forward pass

Backward pass