

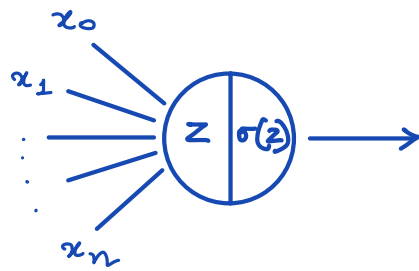
## Vanishing Gradient in Backpropagation:

We will learn about Vanishing Gradient problem, how it happens and how you can overcome this problem.

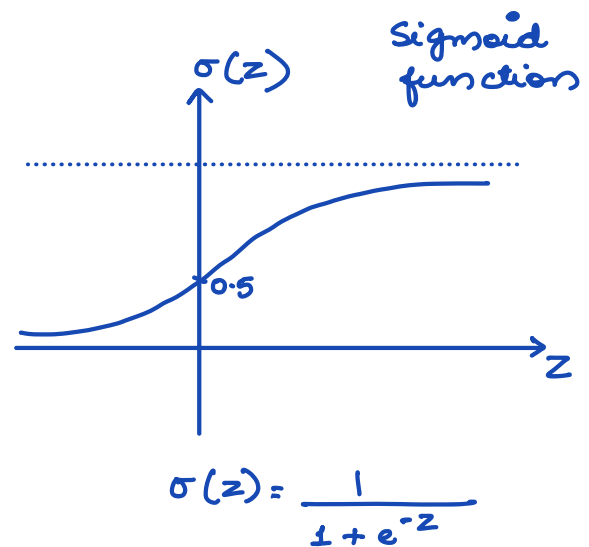
This problem arises particularly due to usage of activation function like sigmoid.

### PROBLEM:

Let's say you have a neural network and there are some inputs coming in ( $z$ ) and you apply sigmoid ( $\sigma$ ) function and produce some output.



$$Z = \sum_{i=0}^n \omega_i x_i$$



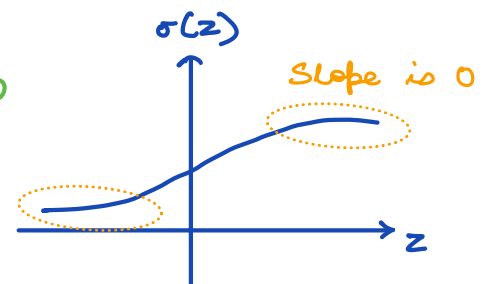
If  $z = +\infty$  function approaches 1

If  $z = -\infty$  function approaches 0

If  $z = 0$  function has value 0.5

Problem is due to following reasons:

- If  $z$  is very large then slope is 0
- If  $z$  is very small then gradient of this function becomes 0



$$\sigma'(z) = \sigma(z) [1 - \sigma(z)]$$

- If  $z$  is large then  $\sigma(z) = 1$      $z \rightarrow \infty$      $\sigma(z) = 1$

$$\sigma'(z) = \sigma(z) \underbrace{[1 - \sigma(z)]}_{\text{becomes } 0}$$

$$\sigma'(z) = 0$$

- If  $z$  is small then  $\sigma(z) = 0$      $z \rightarrow -\infty$      $\sigma(z) = 0$

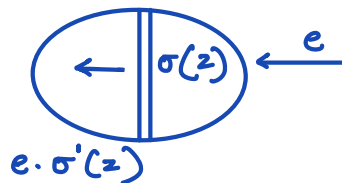
$$\sigma'(z) = \underbrace{\sigma(z)}_{\text{becomes } 0} [1 - \sigma(z)]$$

$$\sigma'(z) = 0$$

Why do we need slope?

In a multilayer perceptron when you propagate errors back, error is multiplied with derivative of activation function.

$$\delta^L = (\delta^{L+1}, \omega^{L+1^T}) \odot \sigma'(z^L)$$



$$e \cdot \underbrace{\sigma(z)(1 - \sigma(z))}_0$$

if  $z$  is large  $\sigma(z) = 1$

$$\underbrace{\hspace{10em}}_0$$

$$e \cdot \underbrace{\sigma(z)(1 - \sigma(z))}_0$$

if  $z$  is small  $\sigma(z) = 0$

$$\underbrace{\hspace{10em}}_0$$

Eventually the gradient that you are propagating back becomes zero.

This is why we don't use sigmoid functions. To avoid this one thing that we generally do is when you are initialising weights you make sure all the weights are small and they should be in range of  $[-0.5 \text{ to } 0.5]$

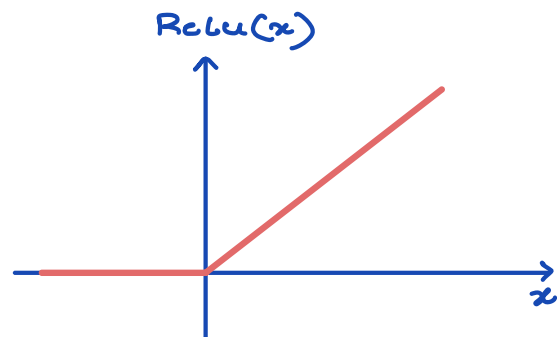
$$\begin{array}{c} \sum w_i x_i \\ \downarrow \\ \text{small} \\ [-0.5, 0.5] \end{array}$$

You can't have weights in range  $[2, 5]$  because:

- all of them are positive.  
Some weights should be positive and some negative.
- they will make value of  $z$  very large.

To avoid this, people use another activation function called as ReLU.

$$\text{ReLU}(x) = \begin{cases} x & x \geq 0 \\ 0 & x < 0 \end{cases}$$



ReLU helps because gradient becomes 1 if  $x$  is +ve and 0 if  $x$  is -ve.

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

Relu will activate some neurons and deactivate some neurons while back propagating.

Instead of using Sigmoid as activation function you should use Relu.

Sigmoid doesnot work well when you have a deep network, during back propagation it will create trouble.