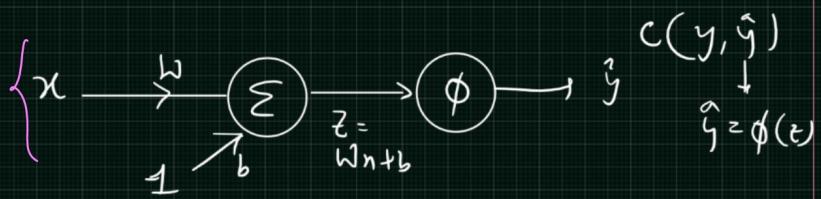


# Activation Functions



forward pass ✓

$$\boxed{\hat{y} = \underline{\phi}(z)}$$

where  $z = \underbrace{w_n}_\uparrow x + \underbrace{b}_\uparrow$

domain  $\rightarrow (-\infty, \infty)$

range  $\rightarrow (0, 1)$  sigmoid

Backward pass

$$\frac{\partial c}{\partial w} = \frac{\partial c}{\partial a} \cdot \frac{\partial a}{\partial z}, \frac{\partial z}{\partial w}$$

$$a' = \underline{\phi'(z)} = \frac{\partial \phi}{\partial z}$$

domain  $\rightarrow (-\infty, \infty) \leftarrow$

range  $\rightarrow (0, 0.5) \leftarrow$

sigmoid

$$\sigma'(z) = \sigma(z) \{1 - \sigma(z)\}$$

$$\text{at } z=0 \quad \sigma'(z=0) = 0.5$$

$$z = \underset{-\infty}{+\infty} \quad \sigma'(z) = 0$$

for forward pass :-

use  $\phi(z)$   $\phi \rightarrow \text{antif}$

consider its  
domain & range

for backward pass

use  $\underline{\phi'(z)}$

Consider its domain &  
range.

# 1) Sigmoid activation $f^n$

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

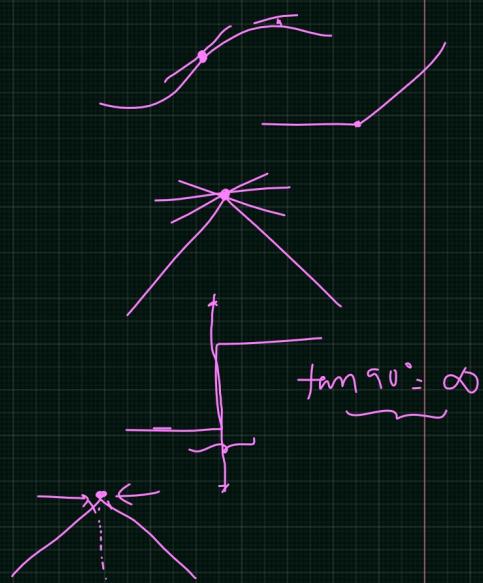
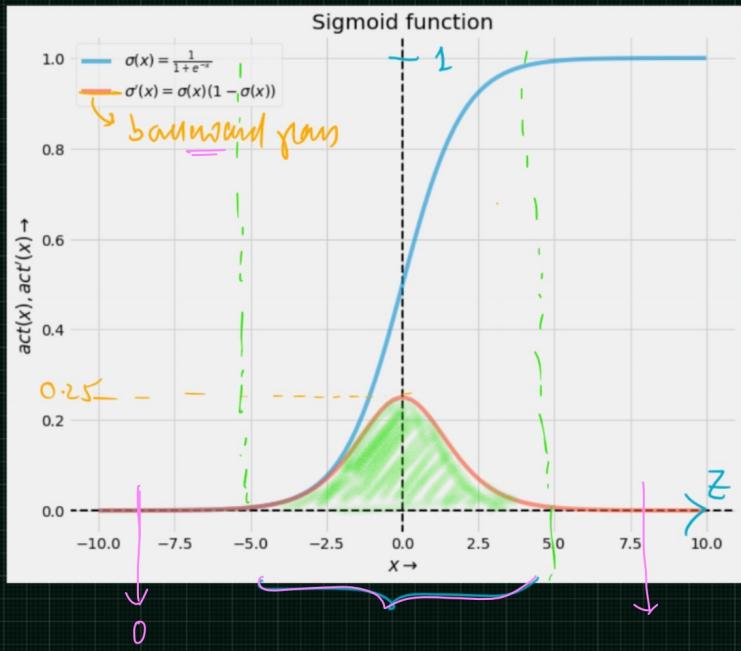
range  $\sigma(z) \in (0, 1)$

domain  $z \in (-\infty, \infty)$

$$\sigma'(z) = \underbrace{\sigma(z)}_{\text{range}} \left\{ 1 - \underbrace{\sigma(z)}_{\text{range}} \right\}$$

$\sigma'(z) \in (0, 0.25)$

domain  $z \in (-\infty, \infty)$



## Advantages:-

- 1) It has smooth gradient  $\Rightarrow$  prevents "jumps" in output.  
continuous  $f^n \Rightarrow$  derivable everywhere.
- 2) Output value or Range  $\sigma(z)$  is between 0 & 1.  
 $\Rightarrow$  normalizing the o/p of each neuron  
 $\Rightarrow$  helps us to reduce the solution space.

## Disadvantages:-

- $\nexists$  1) Prone to gradient vanishing

$$\sigma'(z) \in (0, 0.25] \quad 0.25 < 1$$

- 2) Power operation  $\Rightarrow$  exponential terms  
makes calculation time consuming

latency ↑

$$\left\{ \frac{1}{1+e^{-z}} \right\}$$

## 2. Hyperbolic Tangent activation function

$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

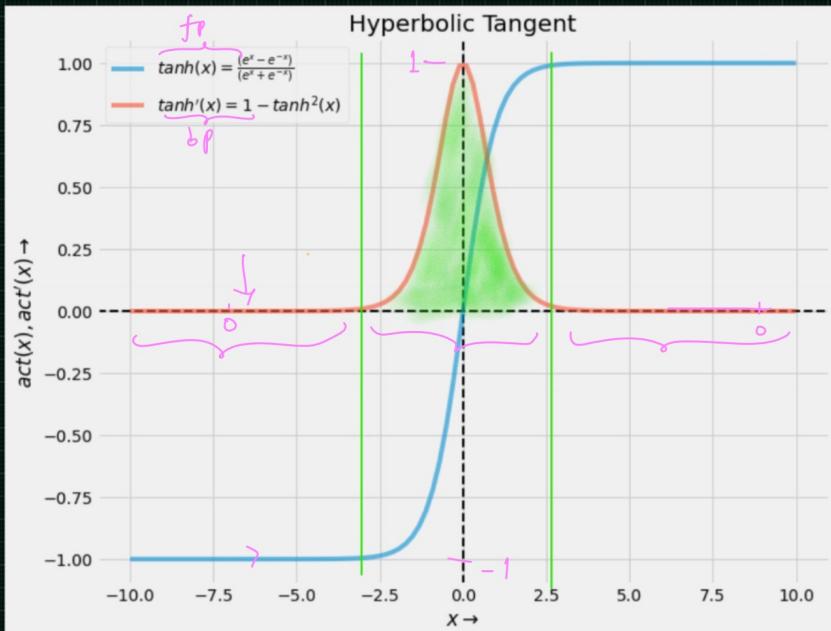
range  $\rightarrow \tanh(z) \in (-1, 1)$

domain  $\rightarrow z \in (-\infty, \infty)$

$$\tanh'(z) = 1 - \tanh^2(z)$$

range  $\rightarrow \tanh'(z) \in (0, 1)$

domain  $\rightarrow z \in (-\infty, \infty)$



Advantages:-

1) Smooth gradients

2) Symmetric to origin  $\Rightarrow$  zero mean

3) Solves vanishing gradient issue to some extent provided there is a proper wt. initialization and also because its range for  $\tanh'(x)$  is between  $(0, 1)$

4) Suitable function for hidden layers

Disadvantage.

1) It also has saturation region

2) It introduces latency because of exponential terms.

### 3) Rectified Linear unit (ReLU) activation fn.

$$\text{ReLU}(z) = \max(z, 0)$$

— or —

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

domain  $z \in (-\infty, \infty)$

range  $\text{ReLU}(z) \in [0, \infty)$

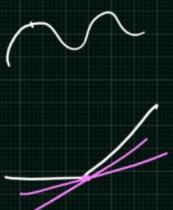
$$\text{ReLU}'(z) = \begin{cases} 1 & z > 0 \\ 0 & z \leq 0 \end{cases}$$

{ derivative is not defined  
at  $z = 0$

domain  $z \in (-\infty, \infty)$   
except  $z = 0$

range

$\text{ReLU}'(z) \in 0 \text{ and } 1$ .



Advantages:-

- 1) If input is +ve  $\Rightarrow$  There is no gradient division problem
- 2) Calculation is much faster {forward or backward} as compared to tanh & Sigmoid.

Disadvantage:-

- 1) For -ve inputs ReLU is completely inactive

↓  
Dying ReLU problem.

- 2) Not a zero centric function.

- 3) At zero its derivative is not defined  $\Rightarrow$  jump

## 4) Leaky ReLU function

$$\frac{d\alpha z}{dz} = \alpha \underset{-\alpha}{\cancel{1}}$$

$$\text{Leaky ReLU}(z) = \max(\tilde{z}, \alpha z)$$

— or —

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

domain  $z \in (-\infty, \infty)$

Range of LeakyReLU(z)  $\in \underline{(-\infty, \infty)}$

$$\alpha = \underline{0.01} \quad 0.01z$$

$\downarrow$

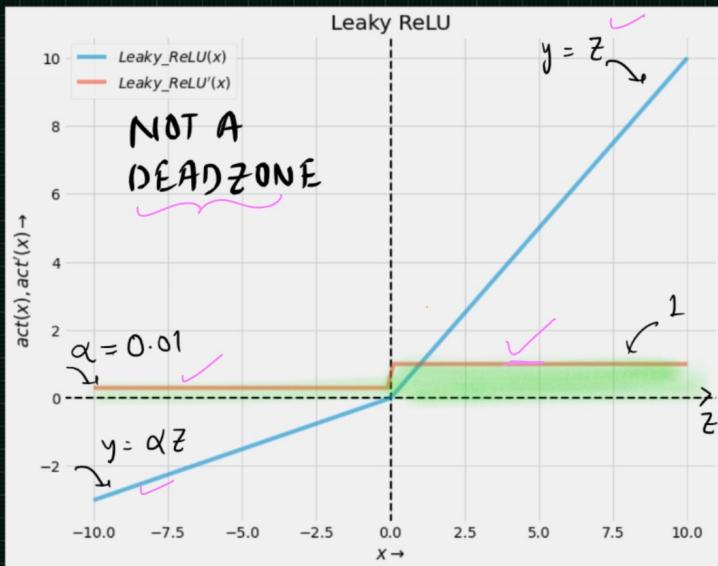
$$\text{LeakyReLU}'(z) = \begin{cases} 1 & z > 0 \\ \underline{\alpha} & z < 0 \end{cases}$$

derivative is not defined  
at  $z = 0$

domain  $z \in (-\infty, \infty)$   
except  $z = 0$

range  $0.01 \text{ or } 1$

$$\text{LeakyReLU}'(z) \in \underline{\alpha} \text{ or } \underline{1}.$$



Advantages:-

for all input

- 1) If input is +ve  $\Rightarrow$  There is no gradient admission problem
- 2) Calculation is much faster {forward or backward} as compared to tanh & sigmoid.

Disadvantage:-

1) for -ve inputs ReLU is completely inactive

Dying ReLU problem

- 2) Not a zero centric function.
- 3) At zero its derivative is not defined  $\Rightarrow$  jump

## 5) PReLU (Parametric ReLU)

$$\text{PReLU}(z) = \begin{cases} z & z \geq 0 \\ \alpha z & z < 0 \end{cases}$$

$\rightarrow \alpha \rightarrow \text{learnable parameter}$

or  
learnable

$$\omega = \omega - \eta \frac{\partial e}{\partial \omega}$$

$$\rightarrow \alpha = \alpha - \eta \frac{\partial e}{\partial \alpha}$$

if  $\alpha = 0 \Rightarrow \text{ReLU}$

if  $\alpha > 0 \Rightarrow \text{Leaky ReLU}$

if  $\alpha$  is learnable parameter  $\Rightarrow \text{PReLU}$

**Advantage:**

1)  $\alpha$  is now intelligently updated }  $\checkmark$   
learnt during the training

↓  
No manual fine-tuning required }  $\checkmark$   
to set  $\alpha$

↓  
 $\alpha$  will be adaptable as per }  $\checkmark$   
the data

**Disadvantage:**

1) You have entire parameter  $\underline{\alpha}$

↓  
to train

↓  
learning time ↑