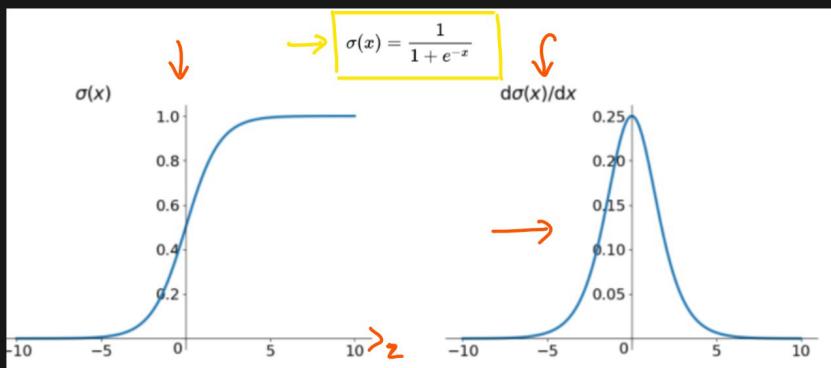


## Activation Functions

Sigmoid, Softmax are more prone to vanishing gradient as  $\sigma_z \leq 0.25$

① Sigmoid Activation function  $[0 \text{ to } 1]$

$$z = \sum_{i=1}^n w_i^T x + b$$

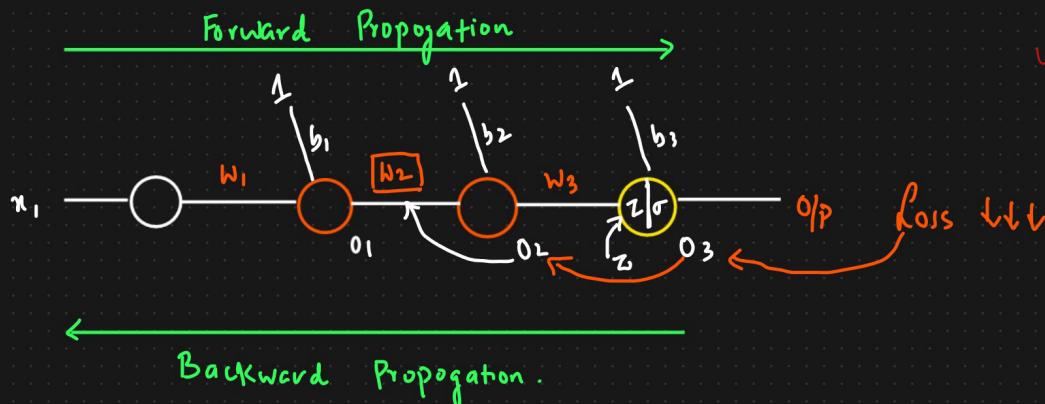


$$\sigma(z) \Rightarrow 0 \text{ to } 1$$

$$\phi(z) \Rightarrow$$

$$\frac{\partial \sigma(z)}{\partial z} = 0 \text{ to } 0.25$$

problem for  
vanishing  
Gradient.



$$w_{2\text{new}} = w_{2\text{old}} - \eta \left[ \frac{\partial L}{\partial w_{2\text{old}}} \right] \quad \Rightarrow \quad w_{2\text{new}} \approx w_{2\text{old}}$$

$$\frac{\partial L}{\partial w_{2\text{old}}} = \frac{\partial L}{\partial o_3} * \boxed{\frac{\partial o_3}{\partial o_2}} * \frac{\partial o_2}{\partial w_2}$$

$0.20 \downarrow \downarrow \downarrow \downarrow \quad 0.01 \quad \times \quad \det z = (o_2 * w_3) + b_3$

$$\frac{\partial o_3}{\partial o_2} = \frac{\partial (\sigma(z))}{\partial z} * \frac{\partial z}{\partial o_2} \quad [0 \text{ to } 1]$$

$$= [0 - 0.25] * \frac{\partial [(o_2 * w_3) + b_3]}{\partial o_2}$$

$$= [0 - 0.25] * w_3 \Rightarrow \text{Small value} \Rightarrow$$

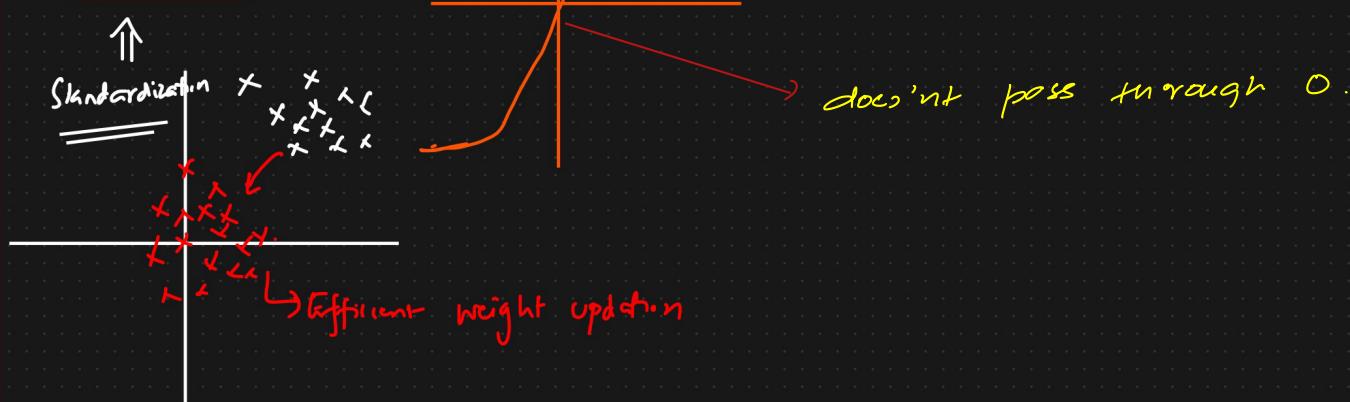
## Advantages

- ① Binary classification Suitable.
- ② Clear prediction i.e. very close 1 or 0

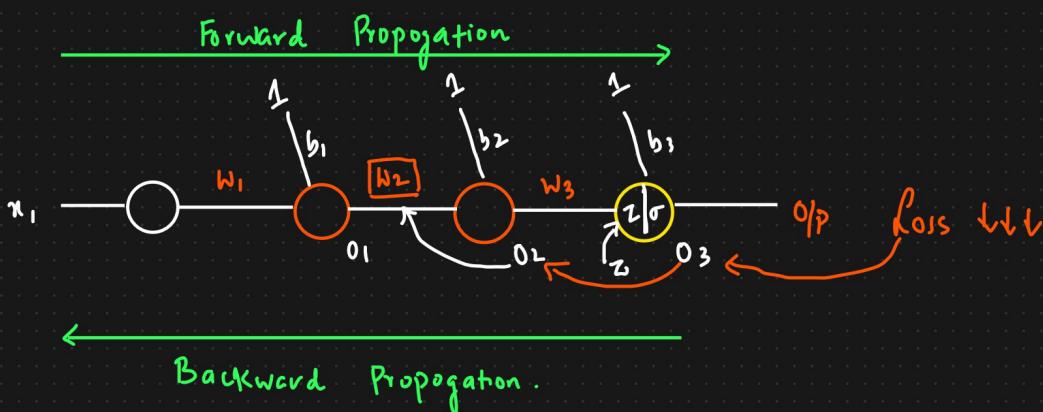
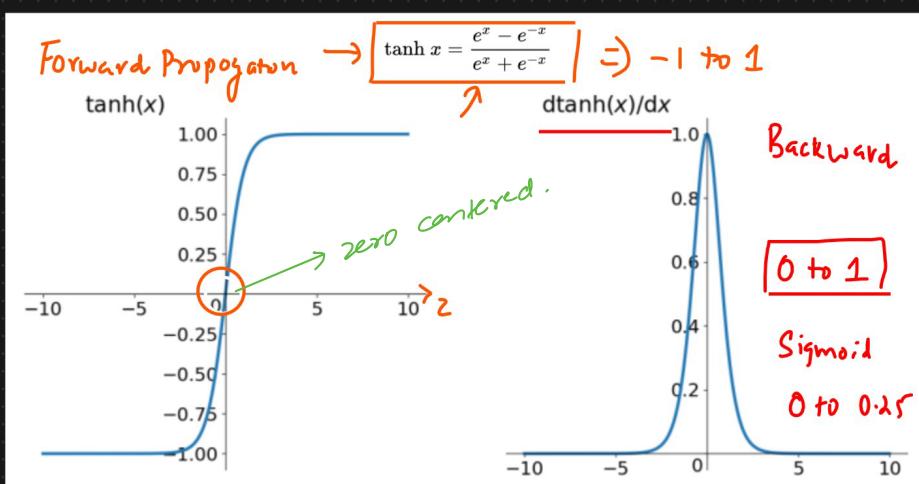
## Disadvantages

- ① Prone to vanishing Gradient Problem.
- ② Function output is not Zero centered  $\Rightarrow$  Efficient weight update
- ③ Mathematical operation are relatively time consuming

## Zero centered



## ② Tanh Activation Function



$$\frac{\partial h}{\partial w_{2,0,1,0}} = \frac{\partial h}{\partial o_3} \neq \boxed{\frac{\partial o_3}{\partial o_2}} \neq \frac{\partial o_2}{\partial w_2}$$

$\downarrow \downarrow \downarrow$

$$0 \cdot 20_{11} * 0 \cdot 01 \times \text{det } z = (o_2 * w_3) + b_3$$

$$\frac{\partial o_3}{\partial o_2} = \boxed{\frac{\partial (\tanh(z))}{\partial z} * \frac{\partial z}{\partial o_2}} [0 \rightarrow 1]$$

$$= [0 - 1] * \frac{\partial [(o_2 * w_3) + b_3]}{\partial o_2}$$

$$= [0 - 1] * w_3 \Rightarrow \text{Small value} \Rightarrow$$

### Advantages

① Zero Centric  $\Rightarrow$  Weight Updation is Efficient

① Prone to Vanishing Gradient Problem

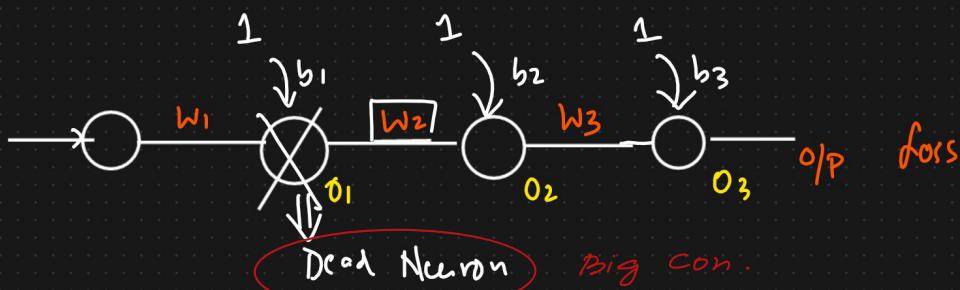
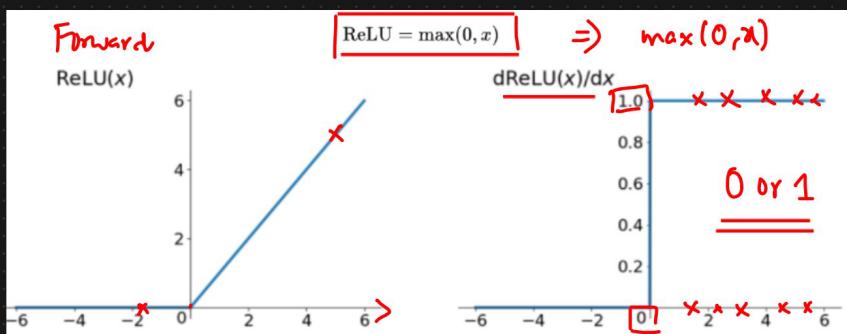
② Time Complexity

[Rectified Linear Unit].

### ③ ReLU Activation Function

Tanh  $\Rightarrow 0 \rightarrow 1$

Sigmoid  $\Rightarrow 0 \rightarrow 0.25$



$$\frac{\partial h}{\partial w_{201d}} = \frac{\partial h}{\partial o_3} * \boxed{\frac{\partial o_3}{\partial o_2}} * \frac{\partial o_2}{\partial w_2}$$

↓ ↓ ↓ ↓

$$0.20_{11} * 0.01 * \text{det } z = (o_2 * w_3) + b_3$$

$$\frac{\partial o_3}{\partial o_2} = \boxed{\frac{\partial (\text{relu}(z))}{\partial z}} * \frac{\partial z}{\partial o_2}$$

[0 to 1]

$$= [0 \text{ or } 1] * \frac{\partial [(o_2 * w_3) + b_3]}{\partial o_2}$$

$$= \begin{bmatrix} -ve & +ve \\ 0 \text{ or } 1 \end{bmatrix} * w_3 \Rightarrow \text{Small value} \Rightarrow$$

If Derivative of ReLU output is  $\boxed{1}$   $\Rightarrow$  Weight updation will happen

If ReLU output is  $\boxed{0}$   $\Rightarrow$  Dead neuron

$$w_{2\text{new}} = w_{201d} - \eta \boxed{\frac{\partial h}{\partial w_{201d}}} \Rightarrow 0$$

If Derivative of  $\text{ReLU}(z)$  is 0

$$\boxed{w_{2\text{new}} \approx w_{201d}} \Rightarrow \text{Dead neuron}$$

If  $\underline{z} = +ve$        $\frac{\partial \text{ReLU}(z)}{\partial z} = 1$

If  $\underline{z} = -ve$        $\frac{\partial \text{ReLU}(z)}{\partial z} = 0$

## Advantages

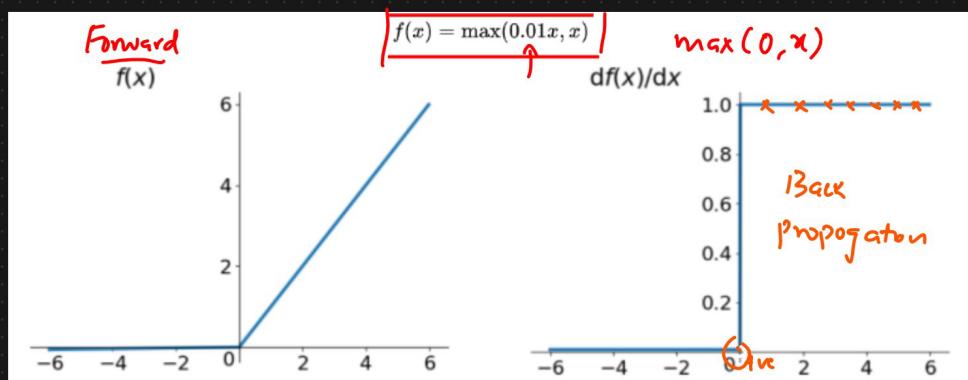
- ✓ ① Solving Vanishing Gradient Problem
- ②  $\text{Max}(0, x) \rightarrow$  Calculation is Superfast. The ReLU function has a linear relationship.
- ③ It is much faster than Sigmoid or Tanh.

## Disadvantages

- ① Dead Neuron
- ② ReLU function is  $0/P$   
 $(0, x) \Rightarrow 0$  or zero number  
 $\downarrow$   
 It is not zero centric

## ④ Leaky ReLU And Parametric ReLU

$\max(\lambda x, x)$   $\rightarrow$  hyperparameter  $\lambda = \alpha = 0.01, 0.02, \dots, 0.03$



ReLU  $\rightarrow$  Dead Neuron  $\rightarrow$  Dead ReLU Problem

## Advantages

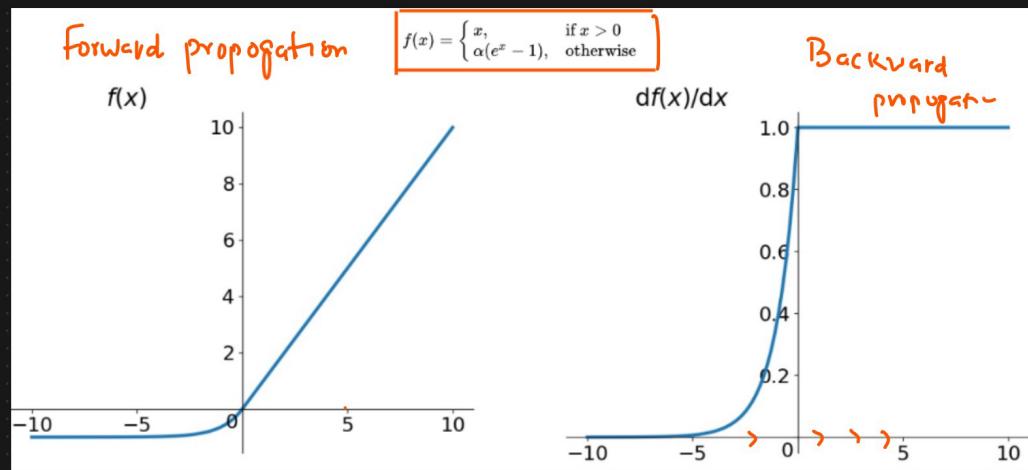
- ① Leaky ReLU has all the advantages of ReLU.
- ② It removes the Dead ReLU Problem.

## Disadvantage

- ① It is not zero centric

## ⑤ ELU (Exponential Linear Units)

$$\alpha(e^x - 1).$$



Advantages

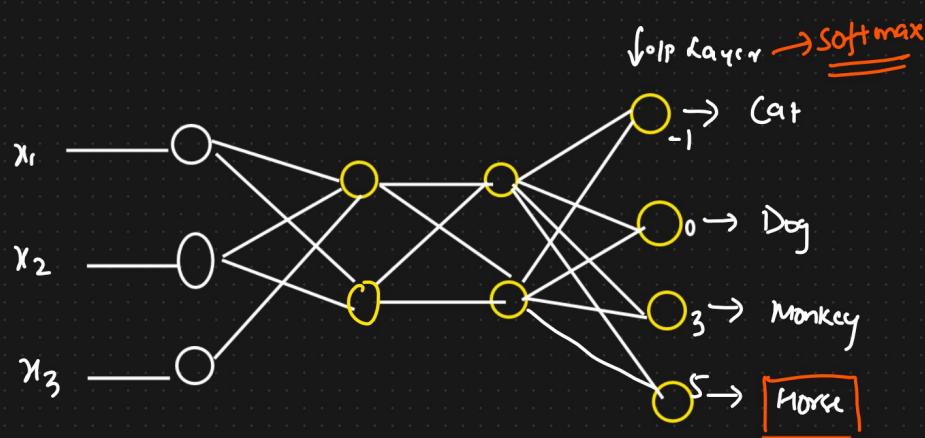
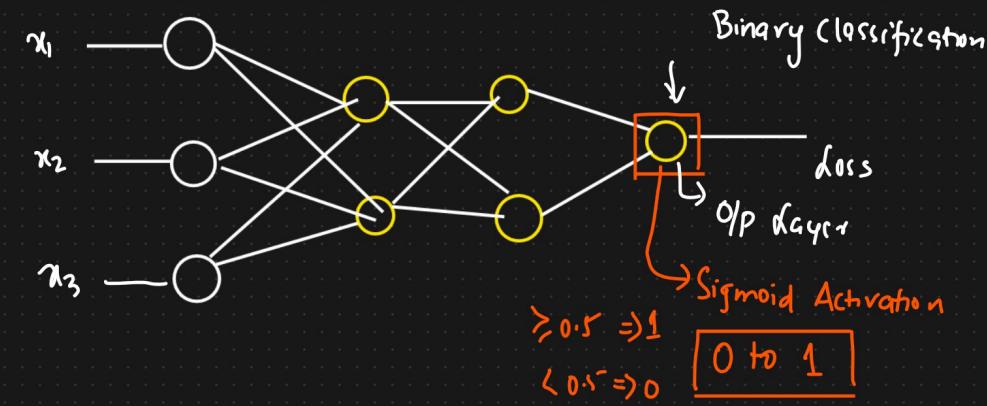
[It is used to solve  
ReLU problems]

Disadvantage

- ① No Dead ReLU Issues
- ② Zero centered

i) Slightly more computationally intensive.

## ⑥ Softmax Activation function [Multiclass classification problem]



$$\text{Softmax} = \frac{e^{y_i}}{\sum_{k=0}^n e^{y_k}}$$

$$y_i = \theta * w + b$$

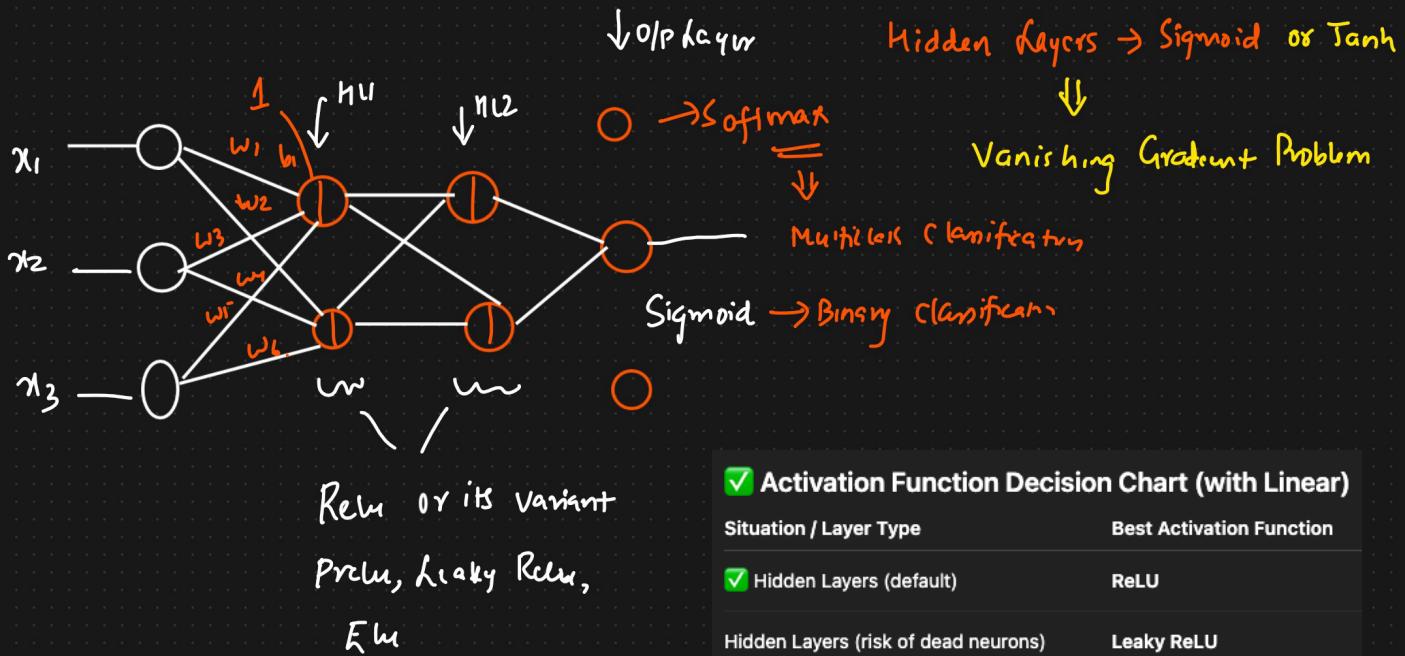
$$\text{Softmax} \Rightarrow \text{Cat} = \frac{e^{-1}}{e^{-1+0+3+5}} = 0.00033 \quad \Pr(\text{Horse}) = \frac{0.1353}{0.00033 + 0.0024 + 0.0183 + 0.1353}$$

$$\text{Dog} = \frac{e^0}{e^{-1+0+3+5}} = 0.0024 \quad \approx 86\%$$

$$\text{Monkey} = \frac{e^3}{e^{-1+0+3+5}} = 0.0183$$

$$\text{Horse} = \frac{e^5}{e^{-1+0+3+5}} = 0.1353$$

## 7 Which Activation Function To Use When?



In short;

Hidden Layer — ReLU

Output Layer (Classification) — Sigmoid  
Softmax

Output Layer (Regression) — Regression Activator

### Activation Function Decision Chart (with Linear)

Situation / Layer Type	Best Activation Function
Hidden Layers (default)	ReLU
Hidden Layers (risk of dead neurons)	Leaky ReLU
Hidden Layers (RNNs / time series)	tanh
Output – Binary Classification	Sigmoid
Output – Multi-class Classification	Softmax
Output – Regression (real number prediction)	Linear
Advanced models (deep / experimental)	Swish / ELU