

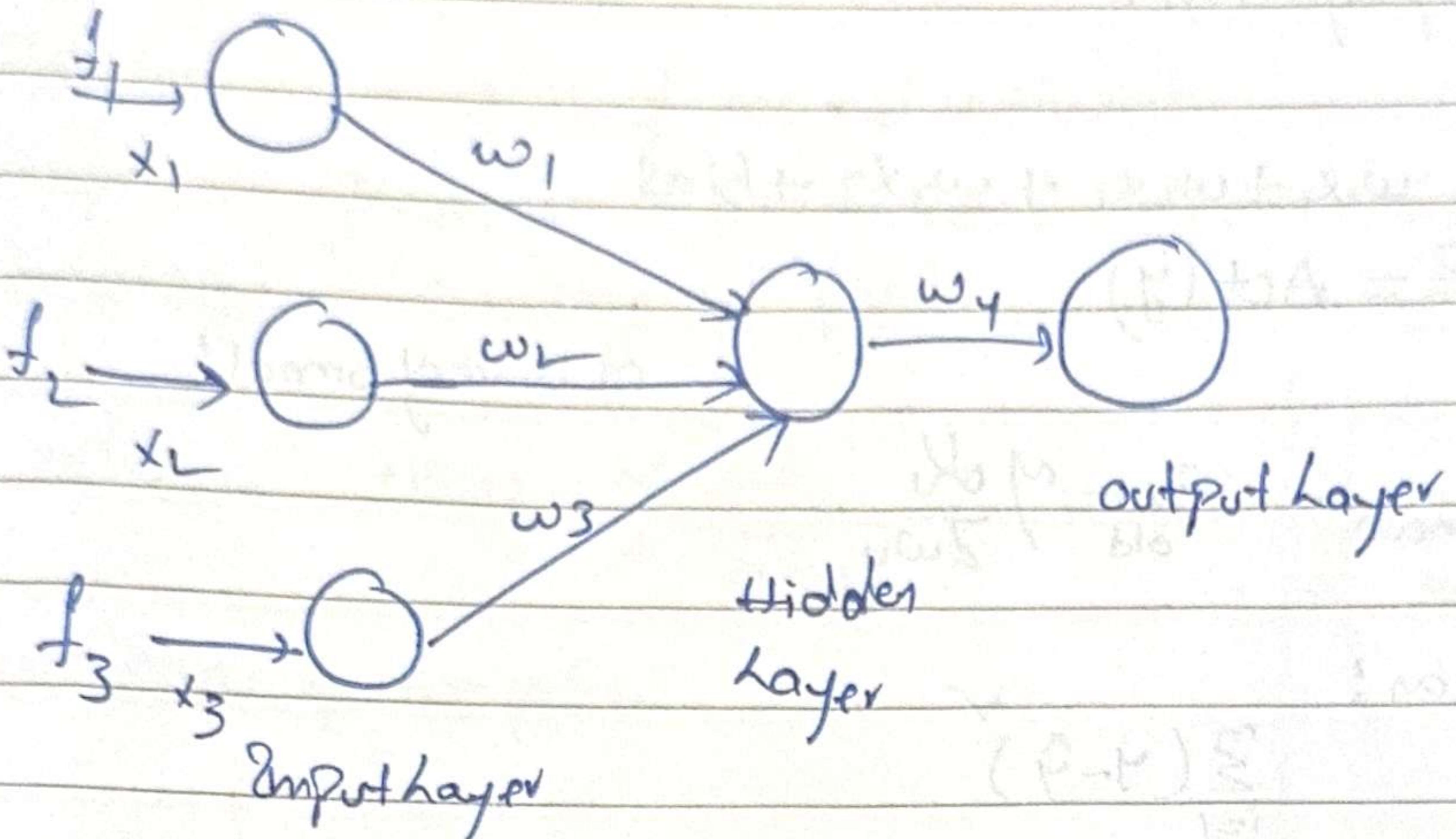
## Deep Learning.

ANN  
CNN  
RNN

classmate

Date \_\_\_\_\_  
Page \_\_\_\_\_

- \* the first Neural network created is Perceptron

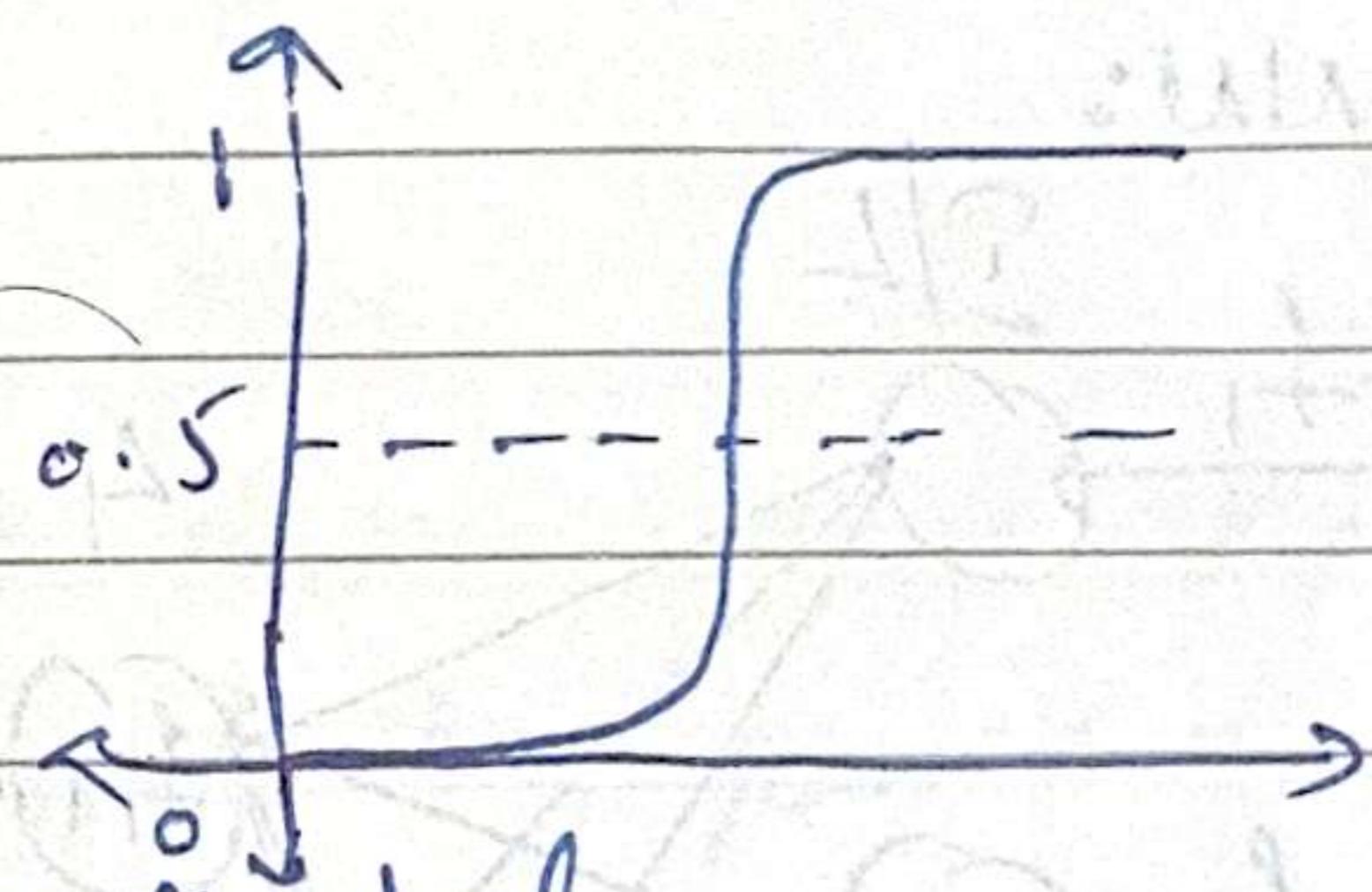


$$y = w_1x_1 + w_2x_2 + w_3x_3 + \text{bias} \quad | \quad z = \text{Act}(y)$$

Activation functions:

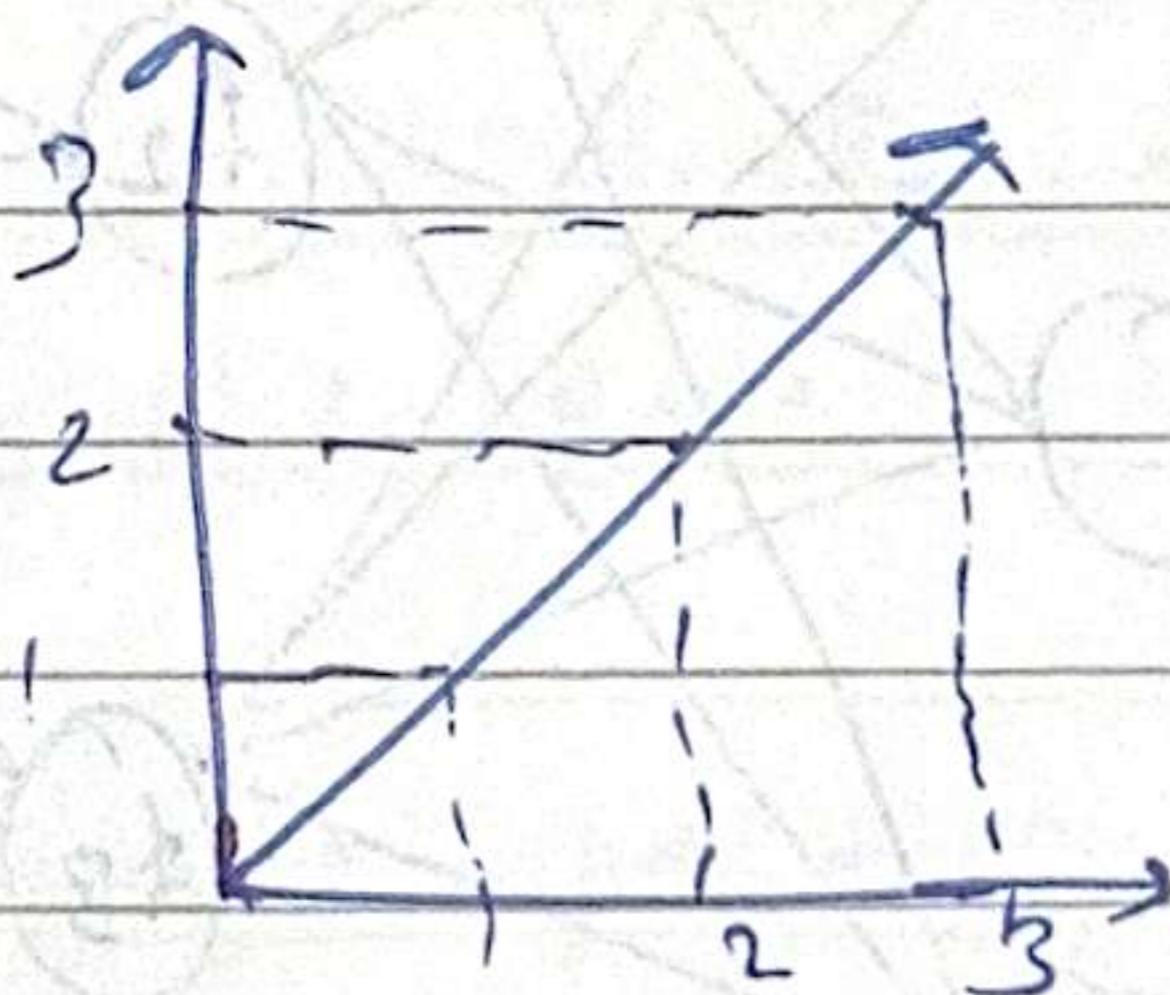
1. Sigmoid Af:  $= \frac{1}{1+e^{-y}}$

If the output is '1' then the Neuron is activated else not activated.



2. ReLU Af:  $\max(y, 0)$

If  $y = -ve$  then its always 0, else



In classification problem it always uses the Sigmoid fn at out layer if binary.

Multi-class, classification we use Softmax

Loss function:

Back propagation.

$$y = w_0x_0 + w_1x_1 + w_2x_2 + \text{bias}$$

$$\hat{y} = \text{Act}(y)$$

$\gamma$  is very small

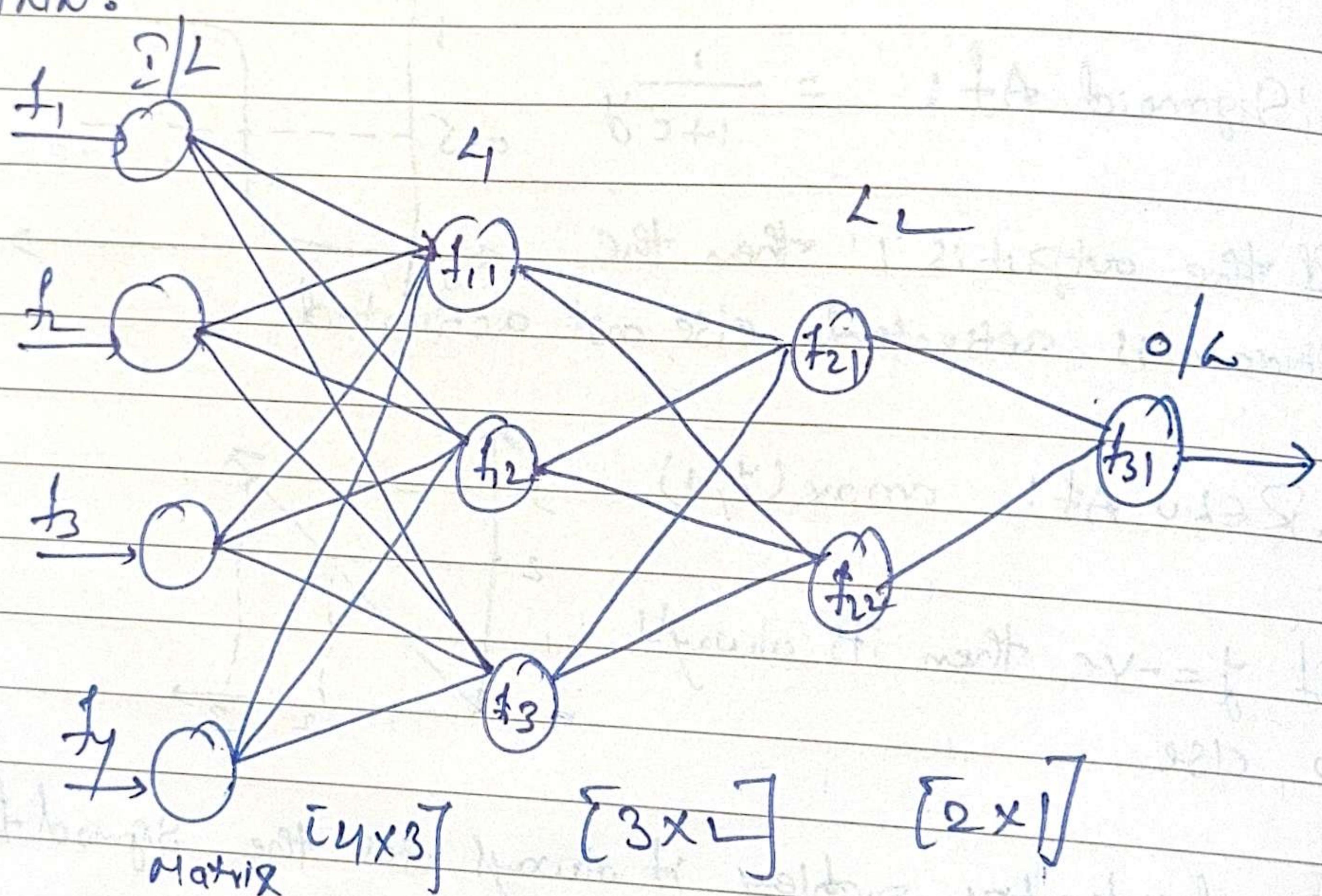
$$w_{\text{new}} = w_{\text{old}} - \gamma \frac{\partial L}{\partial w_i}$$

Cost function:

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2$$

In back propagation we find how much the weight is responsible for the loss by the derivation of the weight.

MNN:



We can select different Act function at different layers every connection have different weights and every layer have different bias.

## Gradient Descent : (optimizers)

Hill - loss function

e.g. A man on hill

Position - weights of neural network.

Slope — gradient (derivative of loss wrt weights)

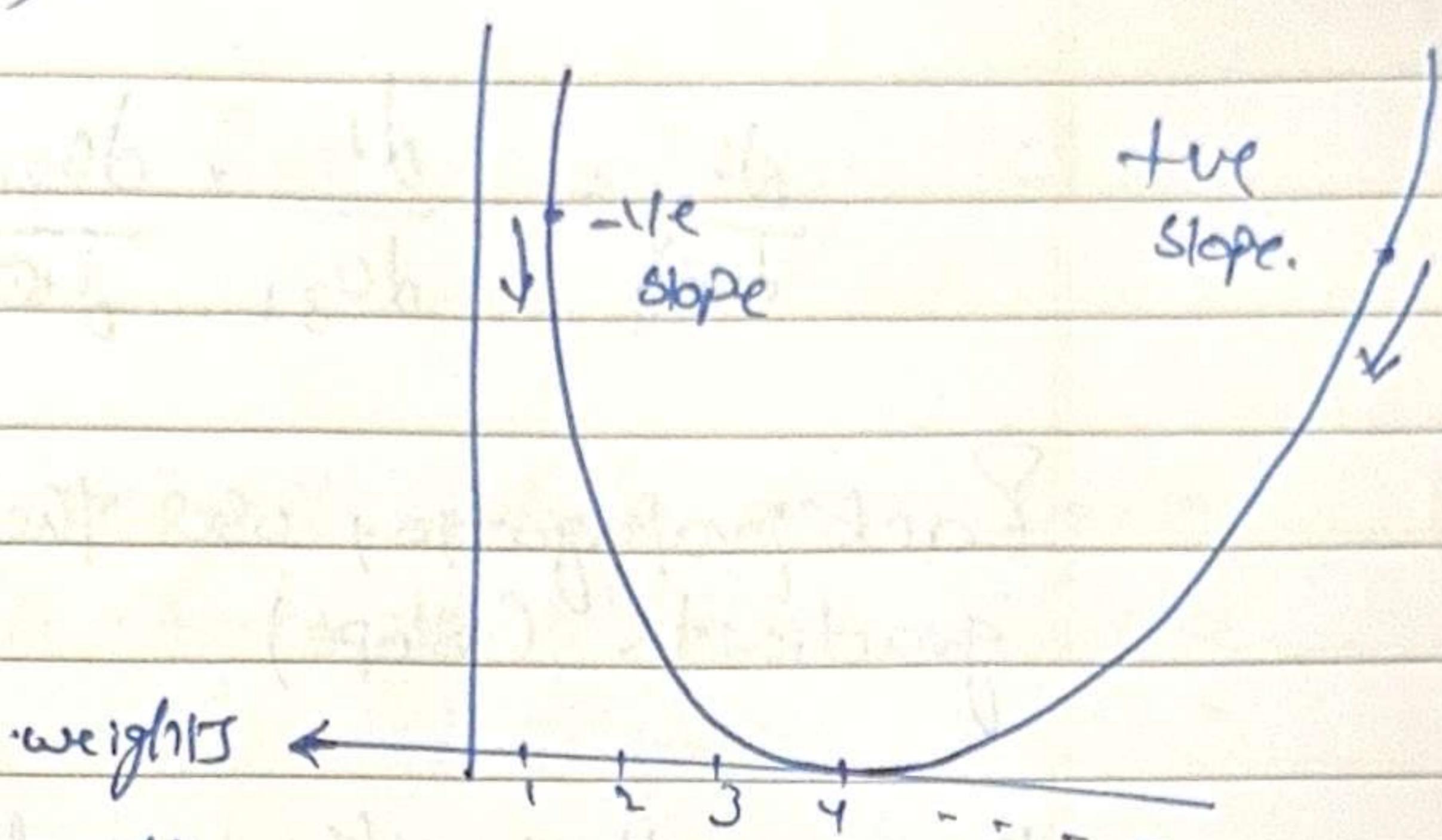
Stepsize - learning rate ( $\gamma$ )

$$w_{\text{new}} = w_{\text{old}} - \gamma \frac{dw}{dw_{\text{old}}}$$

-ve slope  $\uparrow$  weight

+ve slope  $\downarrow$  weight.

How fast  $\uparrow$  or  $\downarrow$  the weights depends on the  $\gamma$  learning rate.



If we take bigger value of  $\gamma$  then we may never reach the global min and if we take less  $\gamma$  then it may take lot of time and computational power to compute best weights.

**Batch gradient Descent :** uses the whole dataset but very slow for complex data (or) large data sets.

**Stochastic Gradient Descent:** speeds up the process by updating weights after one random training sample till all the batches are done.

**Mini batch gradient descent:** used small batches like 32 (or) 64 samples: good speed and the stability.

## Chain Rule:

$$\text{Loss function} = (y - g)^2$$

2. Hidden Layer Example

$$w_{ii}^{\text{new}} = w_{ii}^{\text{old}} - \gamma \frac{dw}{dw_{ii}}$$

$$\frac{dw}{dw_{ii}} = \frac{dL}{d\omega_{31}} \times \frac{d\omega_{31}}{d\omega_{22}} \times \frac{d\omega_{22}}{dw_{ii}}$$

Backpropagation uses the chain rule to compute the gradients (slope)

\* we multiply the gradient from the next layer by the derivative of the activation function of the current layer.

\* Each layer passes its error signal to the previous layer.

\* Backpropagation computes gradients of the loss with respect to every weight in the network.

\* Gradient tells us how much small change in that weight changed the loss. So it measures the proportion of the error caused by the weight.

→ we are multiplying the gradients every time so if slope  $> 1$  then grows exponentially else shrinks exponentially so gradient value so the first few layers can learn to change according to the loss function.

## Vanishing Gradient Problem:

\* In early days in deep learning they always used the Sigmoid function as the Act.

⇒ the Sigmoid function converges the values to  $0 \rightarrow 1$

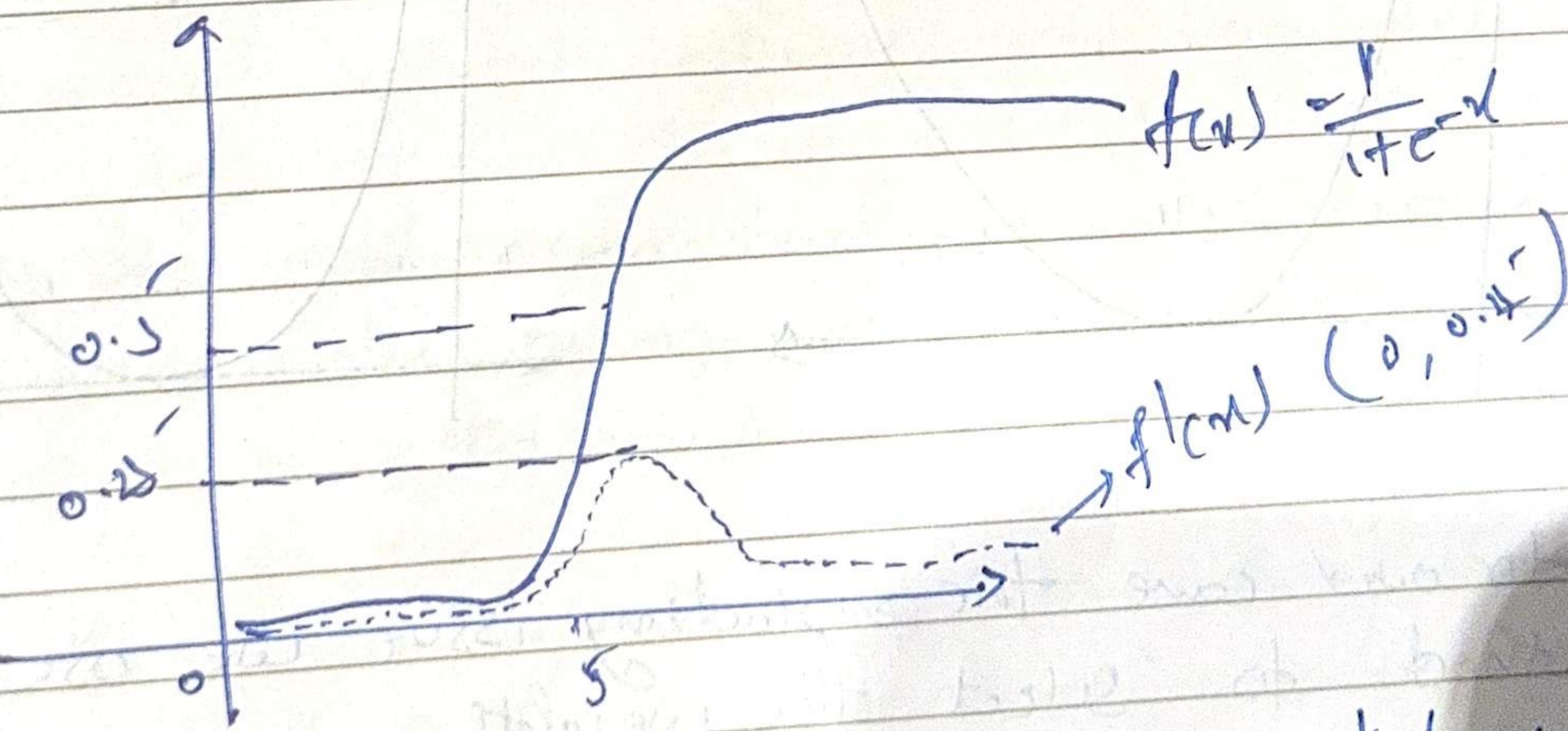
⇒ derivative of the Sigmoid function always lies between 0 to 0.25

\* when we have more hidden layers and by chain rule multiply the gradients it tends to 0 there is no change in new weight and old weight.

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

$$f(x) = \tanh(x) \quad f'(x) = (1-\tanh^2(x))$$

Sigmoid (0, 1)    tanh (-1, 1)



every time we differentiate the value reduced to 1/4 of its original which eventually tends to 0

Page

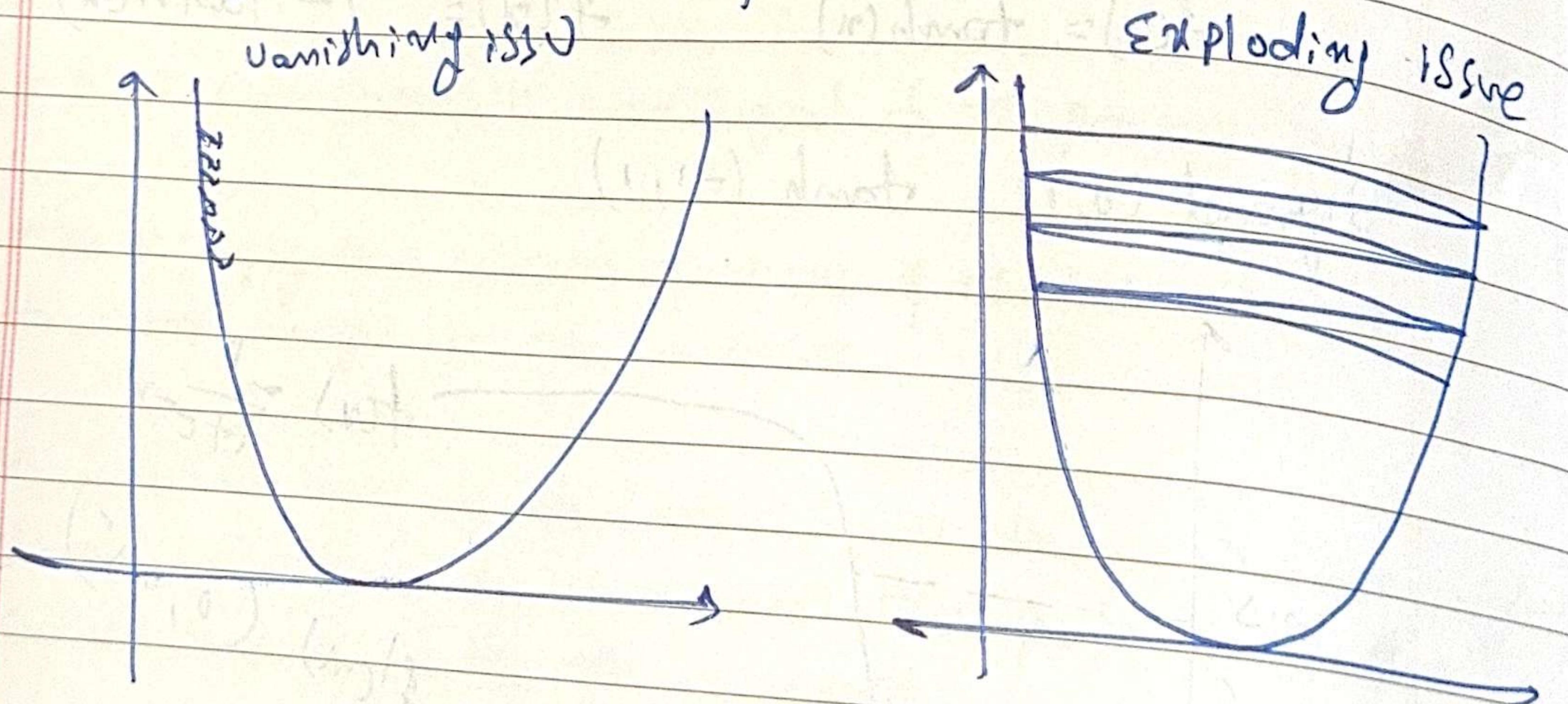
Having deep layered network and using the sigmoid (or) tanh as Act in every layer causes the vanishing problem.

⇒ the layers which are closer input layer the weights won't be changing (or) never converge to global min or it may take infinite time to reach the global min.

### Gradient exploding Problem!

⇒ In vanishing issue the layers closer to input layer learn anything because  $w_{old} = w_{new}$

⇒ In exploding due the larger +ve (or) -ve weight it will never go to the global min if will just zig zag at the top of the hill.



⇒ To overcome the exploding issue we use different method to select the weight.

Learn and far out to calculate the approximated weight which we can assign to the model.

## Drop out & Regularization:

Drop out ratio  $0 \leq p \leq 1$

⇒ we select the p-value from  $(0, 1)$

⇒ Dropout is the regularization technique to prevent overfitting.

⇒ During the training we randomly drop a fraction of neurons.

Dropped Neurons will not participate in forward or back word propagation. for the training step.

⇒ this forces the network not to rely on the few neurons, this will make the system more robust.

In the input layer the p-value  $(0.1, 0.2)$  we don't want to lose more input info, in hidden layer  $(0.3-0.5)$

Each mini batch sees a slightly different network therefore Learned network is the average all the networks.

As the no of neurons drops as p-value the total act are the average info also drops to maintain we scale up the value by multiplying with  $1/(1-p)$

⇒ This scaling is applied in every single training step not just in the final averaging for all mini batches.

to maintain inference.

Act functions:

Sigmoid functions:

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

This function is not '0' zero centered only positive

\* works good with the binary classification problem and having the very less hidden layers

\*  $\alpha'(x) = \alpha(x)(1-\alpha(x))$  in every preceding step the gradient value is decreased by  $1/4$

\* This function mostly used at output layer

threshold function:

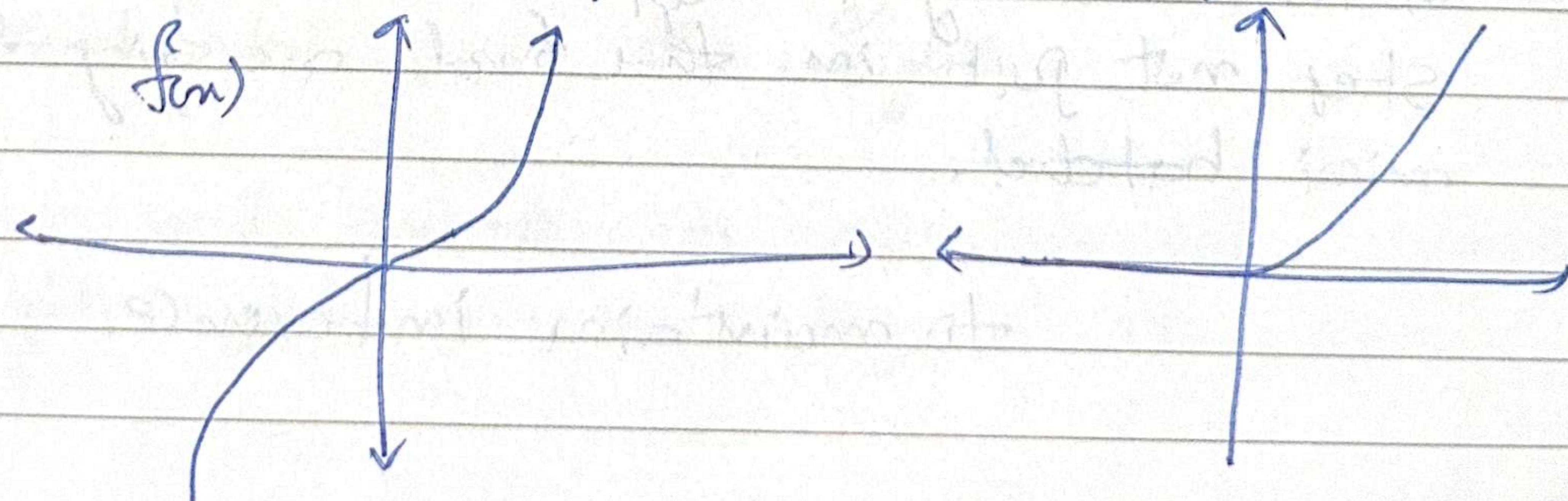
$$\alpha(x) = \tanh x \quad \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$f(x) = 1 - \tanh x$$

this is 0 zero centered function (-1, 1)

If the act function is not zero centered it may converge to global min the weights zigzag at top of the curve.

Not differentiable at  $x=0$



ReLU:

$$f(x) = \begin{cases} x & ; x > 0 \\ 0 & ; x \leq 0 \end{cases}$$

$$f(x) = \max(0, x)$$

$$f'(x) = \begin{cases} 1 & ; x > 0 \\ 0 & ; x \leq 0 \end{cases}$$

the ReLU is piecewise linear function

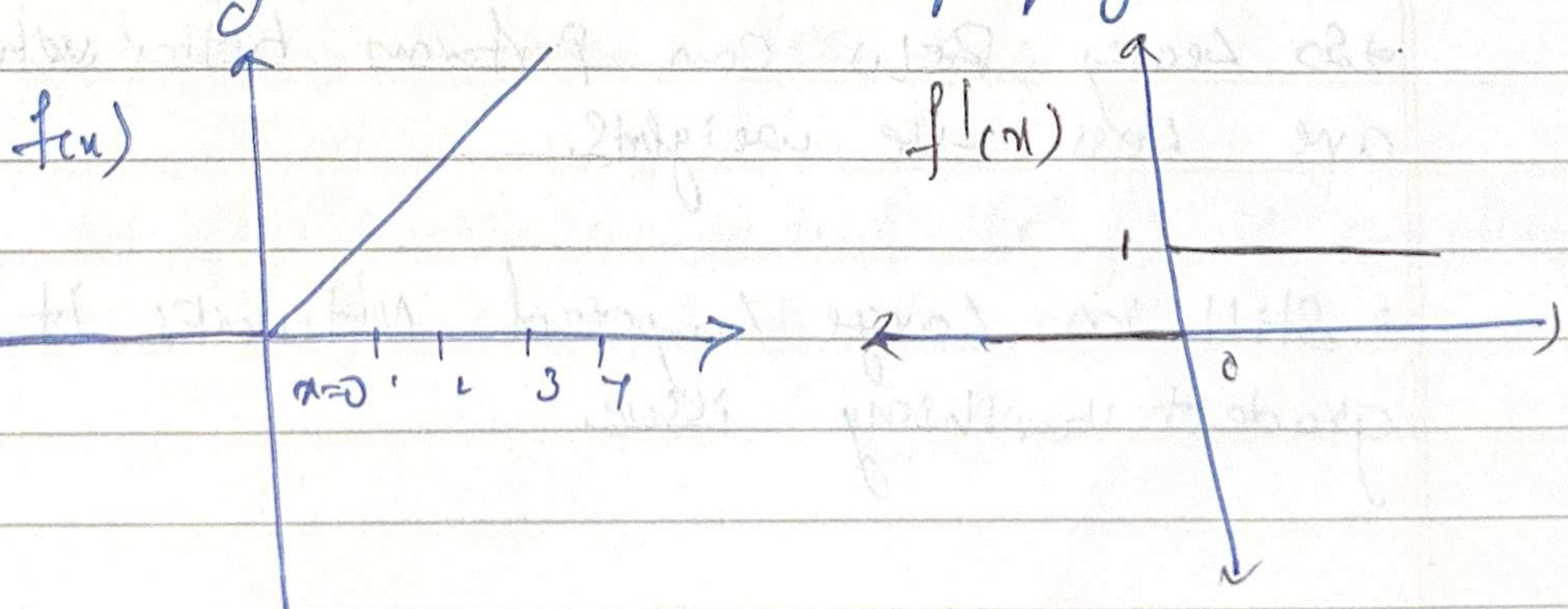
$\Rightarrow$  for the positive value's the ReLU is Linear

$\Rightarrow$  At  $x=0$  there is the hinge which makes the ReLU Non-Linear

$\Rightarrow$  The non-linearity rises as the no. of hidden layers increase

$\Rightarrow$  The ReLU solve the dying gradient issue by not Reducing the Act Value while differentiating.

$\rightarrow$  But for the -ve weight the df is always 0 so the large -ve weight can never converge to global min which creates the dead neurons which learn nothing in the back propagation.



## Leaky Relu Act:

⇒ The Leaky Relu Sloves the dying Relu issue by adding the 0.01( $x$ ) to the negative portions then the derivative will be 0.01 not the zero as in the Relu.

$$f(x) = \max(0.01x, x)$$

$$f(x) = \begin{cases} x & ; x > 0 \\ 0.01 & ; x \leq 0 \end{cases}$$

⇒ But still we are multiplying the gradient with 0.01 this leads to vanishing gradient issue in the -ve weights

⇒ it converges for the +ve value towards the global min will be far faster than -ve weights and in some cases -ve weights may not reach the global min.

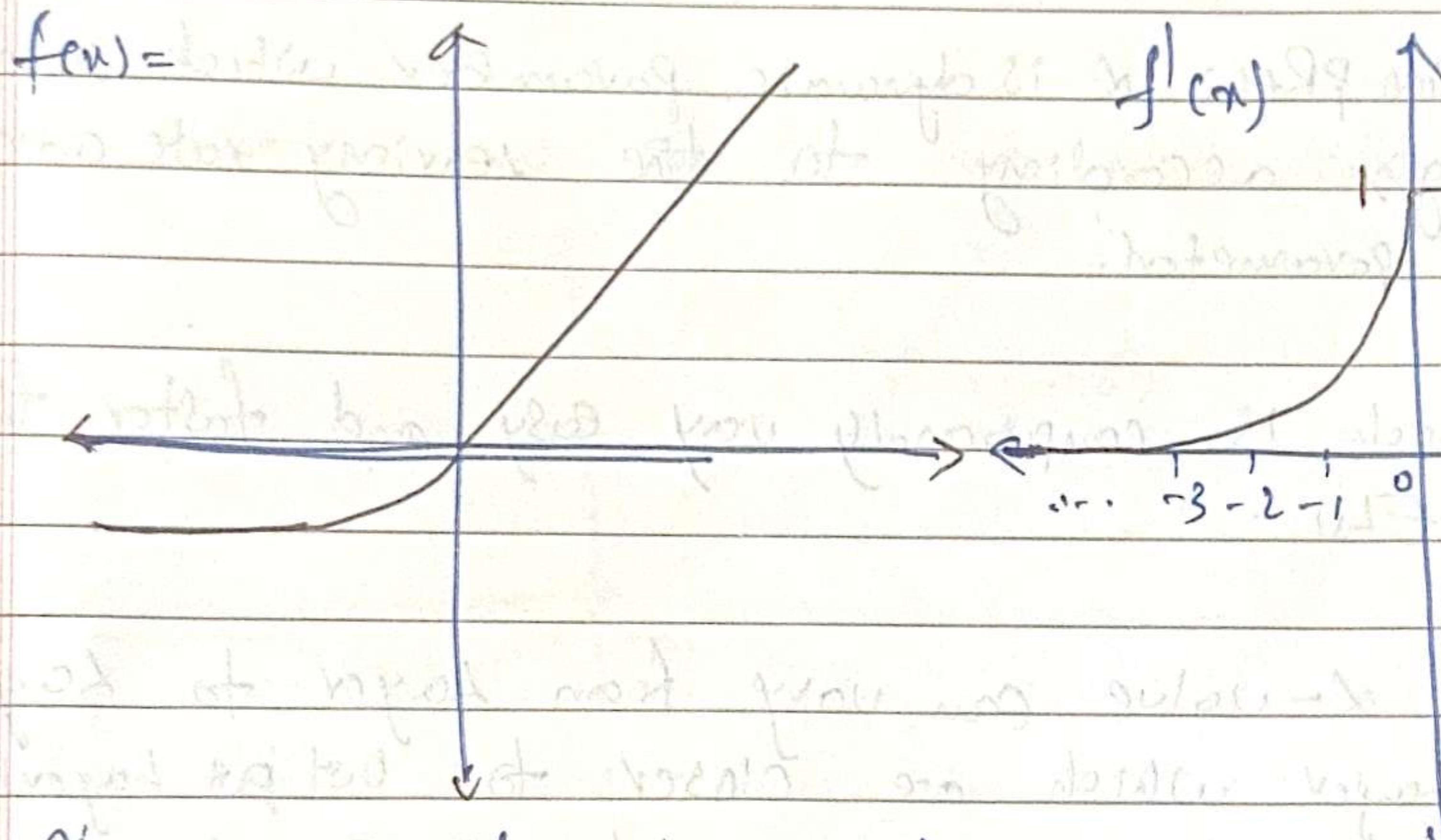
⇒ Leaky Relu is better than the Relu But more sequence of the weights are -ve than every time we are multiplying the gradient by  $10^{-2}$  times.

⇒ So Leaky Relu can perform better when there are less -ve weights.

⇒ Still in Large Layered networks it leads to gradient vanishing issue.

Elu (Exponential Linear Unit):

$$f(x) = \begin{cases} x & : x > 0 \\ \alpha(e^x - 1) & ; x \leq 0 \end{cases}$$



$\alpha$  - alpha is hyper parameter.

→ In ReLU and Leaky ReLU the derivative of  $x$  at  $x=0$  is undefined because the slope = 1 when  $x \rightarrow 0^+$  and slope = 0 when  $x \rightarrow 0^-$ . So the

$\frac{dx}{dt} \neq \frac{dx}{dt}$  and derivative of  $x=0$  not defined

but in some cases it  $f'(x=0)$  is considered as 0 in the training data

→ we can overcome the issue in the ELU

→ the derivative of the exponential  $(e^{x-1})$  is complex if requires more computational power.

## PReLU: (Parametric ReLU)

if  $\alpha=0$  ReLU

$$f(x) = \begin{cases} \alpha x & x > 0 \\ \alpha x + \beta & x \leq 0 \end{cases}$$

if  $\alpha=0.01$  Leaky ReLU

But in PReLU  $\alpha$  is dynamic parameter which changes according to the learning rate and other parameters.

→ which is computationally very easy and faster than the ELU.

→ the  $\alpha$ -value can vary from layer to layer the layers which are closer to output layers may have the bigger  $\alpha$  and as we go deeper towards the input layer the  $\alpha$ 's are smaller so we can overcome the gradient vanishing issue and the very easy to compute.

Swish:

$$f(x) = x * \text{sigmoid}(\alpha x) \quad \text{Self gating.}$$

$$f(x) = x \cdot \alpha(x)$$

$$f'(x) = \alpha(x) + x \cdot \alpha'(x)(1 - \alpha(x))$$

so gradient has two parts:

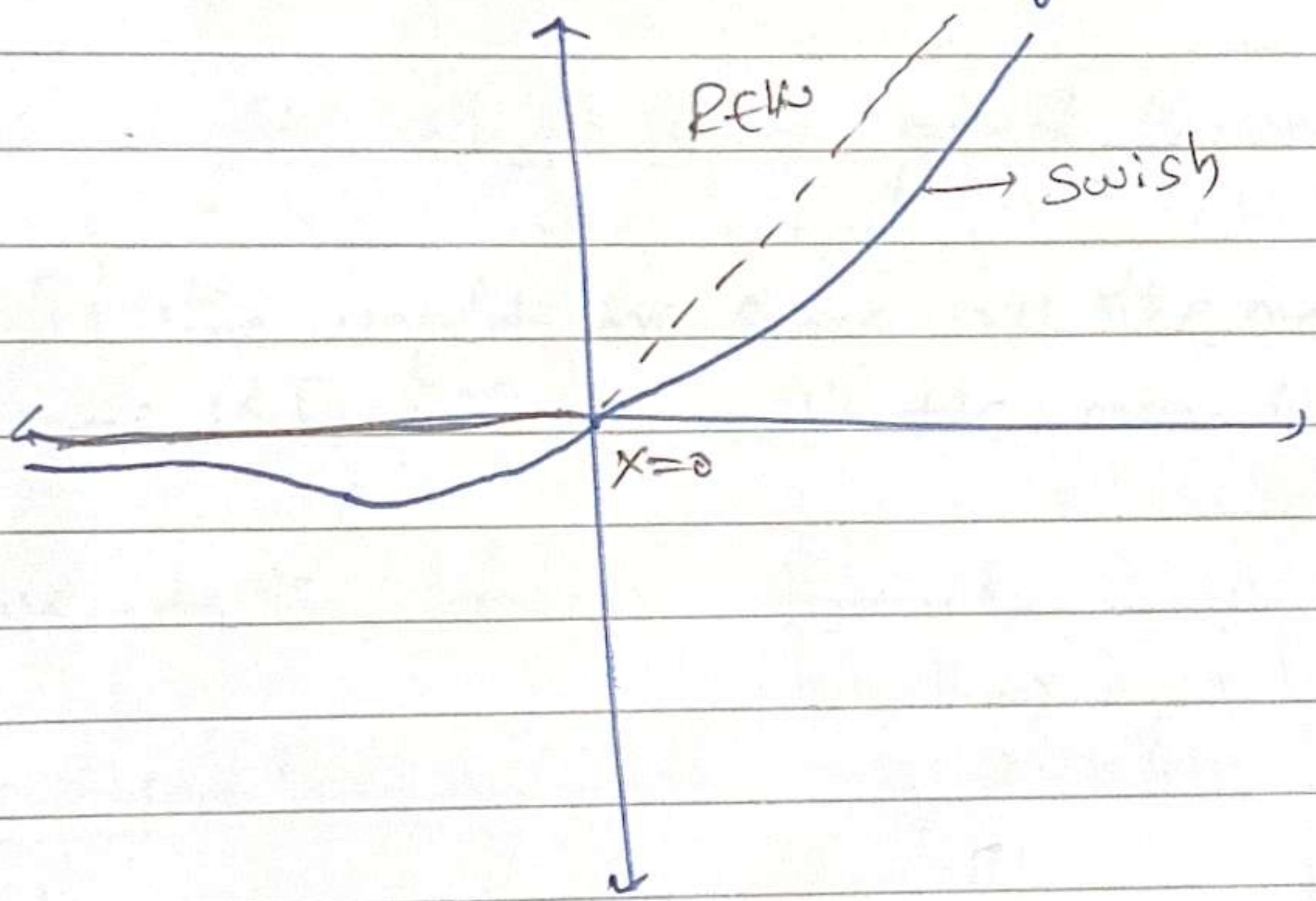
$\alpha(x)$  ensures some gradient always flows (no dead neurons)

$x \cdot \alpha'(x)(1 - \alpha(x))$  introduced a self modulation effect.

Sigmoid acts as a Soft gate (between 0, 1)

→ instead of just cutting off negative like ReLU Swish allows some negative values to leak through, but scaled by Sigmoid.

⇒ 2x better work for the layer  $\geq 40$  not less than  $< 10$



Softplus:

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

The derivative of the Softplus is the Sigmoid function (0, 1)

## weight initialization.

- ⇒ when curr we are initializing the weight should not be large we will have the exploding issue , if too small use it will take lot of time to compute.
- ⇒ weight<sup>+</sup> should not be same
- ⇒ weights should have good variance.
- ⇒ If the weights are same all the neurons receive the same info which kill the non-linearity.

we use the fan in fanout

uniform Distribution:  $w_{ij} \sim \text{uniform}$

$$\left[ \frac{-1}{\text{fanin}} \quad \frac{1}{\text{fanin}} \right]$$

weights are selected between the range of  $(a, b)$  uniformly.

### b. Xavier/Glorot Distribution:

Xavier Normal

$$w_{ij} \sim N(0, \sigma^2)$$

$$\sigma = \sqrt{\frac{2}{(\text{fanin} + \text{fanout})}}$$

3. He init;

He uniform

$$w_{ij} \sim U \left[ -\frac{\sqrt{6}}{\text{fanin}} , \frac{\sqrt{6}}{\text{fanin}} \right]$$

Xavier Uniform.

$$w_{ij} \sim U$$

$$\left[ \frac{-\sqrt{6}}{\sqrt{\text{fanin} + \text{fanout}}} , \frac{\sqrt{6}}{\sqrt{\text{fanin} + \text{fanout}}} \right]$$

the Normal

$$\sigma = \sqrt{\frac{2}{\text{fanin}}}$$