

Deep Learning

AI: Application which can do its own task without any human Intervention.

Ex: Netflix APP, Self Driving Cars.

Subset

ML:

Statics tool to analyze, visualize the data, Predictive, forecasting, clustering.

Subset

DL:

To Mimic the human brain behaviour.
And learn from it.

DS: A Part of AI, ML, DL

Why Deep learning?

Ans: 2005 → Facebook, insta, whatsapp generated a massive data which need to analyse and used for developing business. To mimic the human behavior and memorise the data and learning from data, deep learning is required.

Hardware Advancement: GPU's used to train the DL Models.

Ex: RTX Titan

RTX 3090

CPU: has several cores, Serial processing, large memory capacity. Small batch operation. Ex: Intel
= AMD

GPU: (Graphics Processing Unit) it's works as an performance accelerators for CPU. has thousands of cores. High throughput, parallel processing. executing thousands of operations at once.

Ex: NVIDIA

AMD

Google

BROADCOM

TPU: Tensor Processing unit. Application specific integrated circuit. fast at performing dense vector & matrix computations to accelerate neural network.

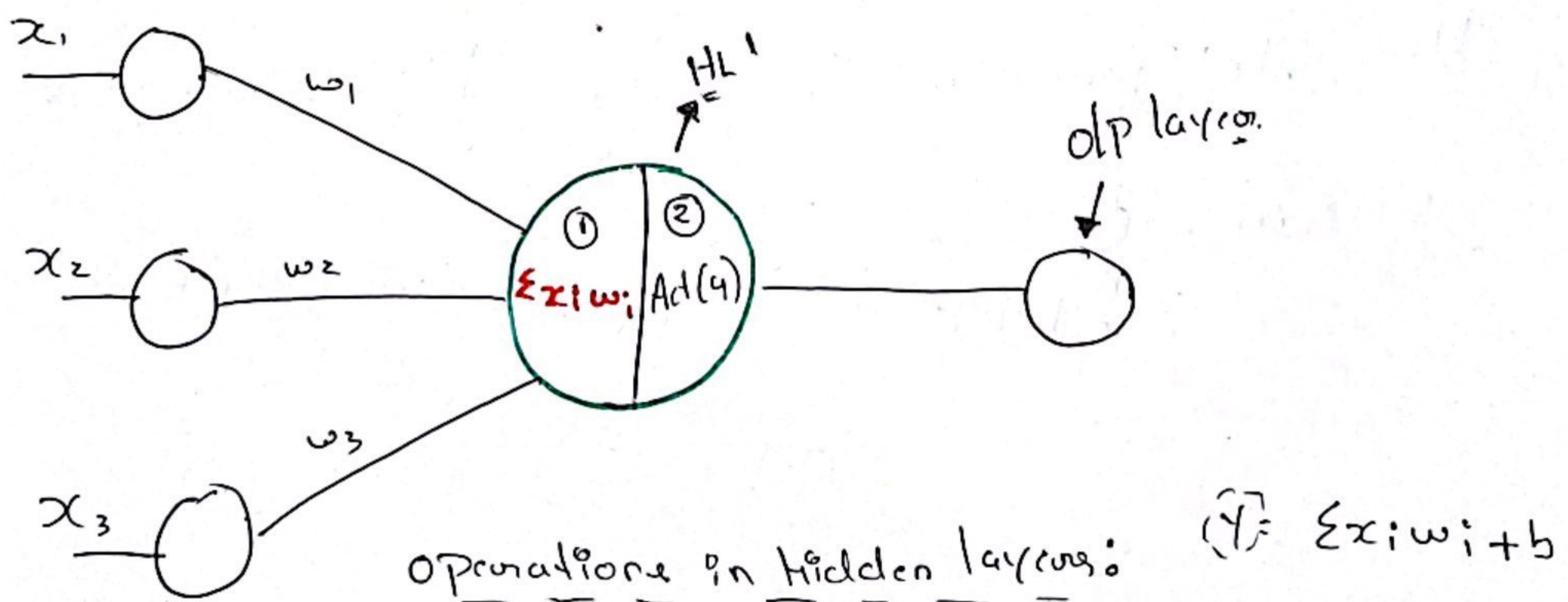
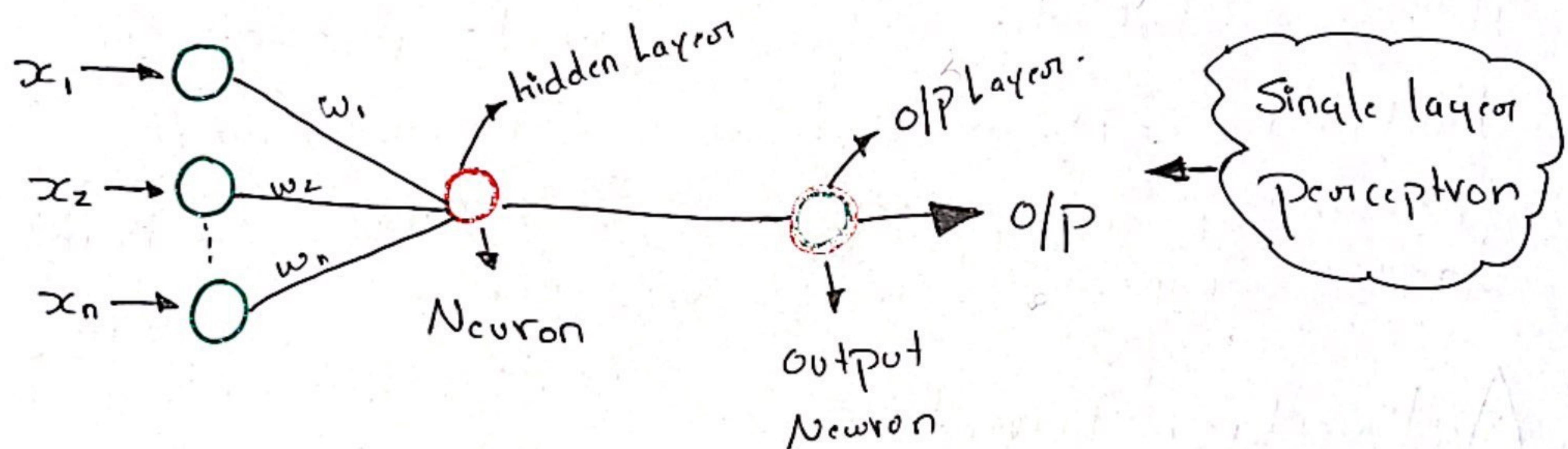
- very high throughput
- extreme parallelism.
- highly optimized for large batch

Ex: Google
= Coral
• HALO

Perceptron

Perceptron is also understood as an "Artificial Neuron" or "neural networks" unit that helps to detect certain input data computation in bimolecular intelligence.

I/P Layer



$$\textcircled{1} \rightarrow \sum x_i w_i = x_1 w_1 + x_2 w_2 + x_3 w_3 \quad (\text{Summation}) \\ = \vec{w}^T \vec{x}$$

(2) Activation function.

Scanned with CamScanner

(3)

Weight importance?

Ans: which weight should have what value

So that the neuron should be activated till that level. (Make the weight to the important process of important to neuron).

Bias?

Each hidden layer has a bias

It is added to weights for to give certain best value. If weight are 0 to overcome the problem of $x_i w_i$. Bias is added. (It is a constant) added for product of feature and weight) to help to shift.

Activation Function

AF to positive or negative side.

It classifying the output to 0 or 1 on the the Act function (Sigmoid). ($y = \frac{1}{1 + e^{-(\sum x_i w_i + b)}}$)

$$\text{Sigmoid} \approx \frac{1}{1 + e^{-y}}$$

$$\frac{1}{1 + e^{-(\sum x_i w_i + b)}}$$

$$\geq 0.5 = 1$$

$$< 0.5 = 0$$

Forward Propagation

- ① Features are the inputs. x_i
- ② Initialize the weights to each x_i component w_i
- ③ summation
Multiple the weights with x_i , add bias.
$$y = (\sum x_i w_i) + b$$
- ④ Pass it to activation function
- ⑤ Output layer receives a value.

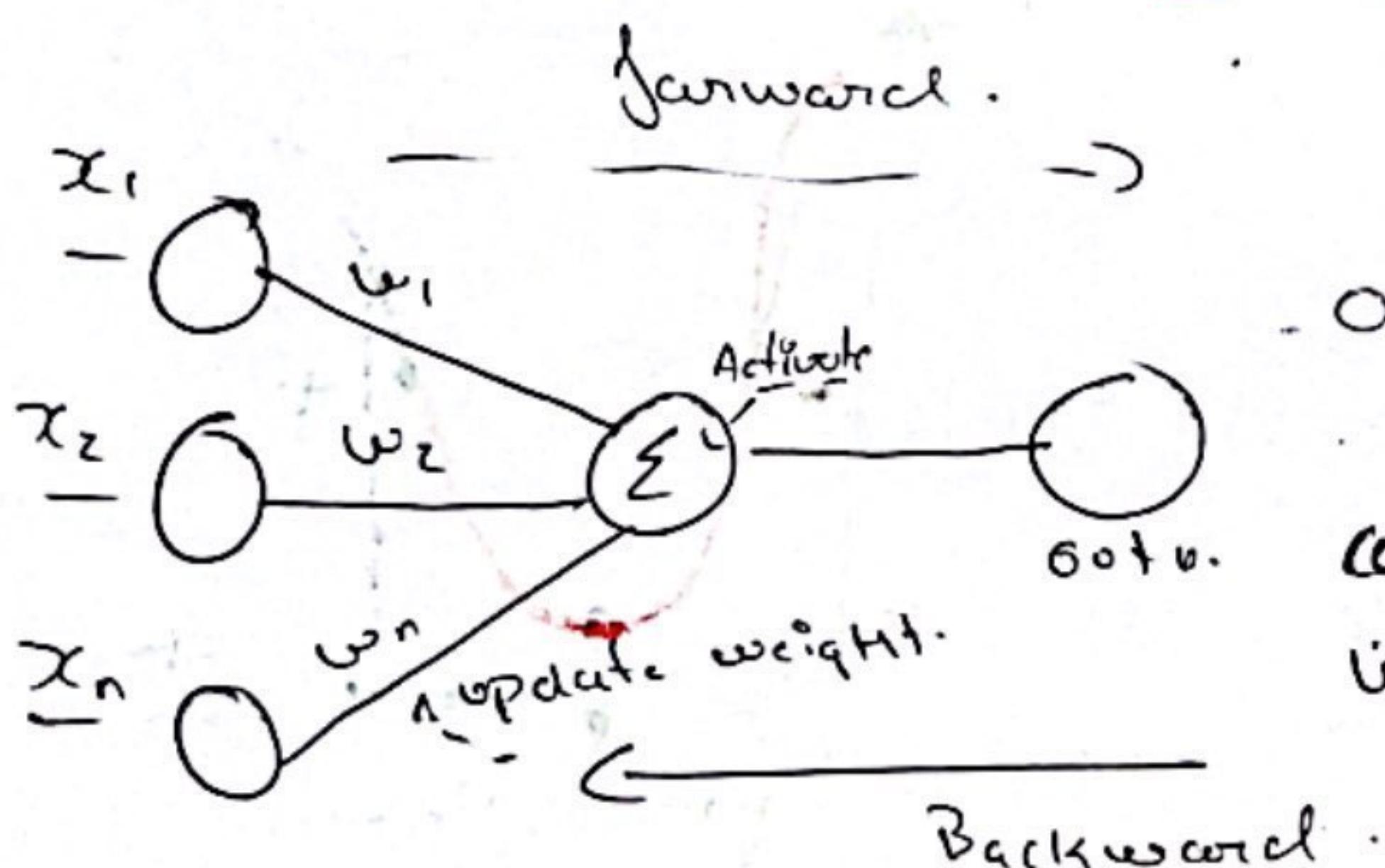
Back Propagation

① After the output is generated if the loss function is high. we use Back propagation.

② Basically it will update the weights.

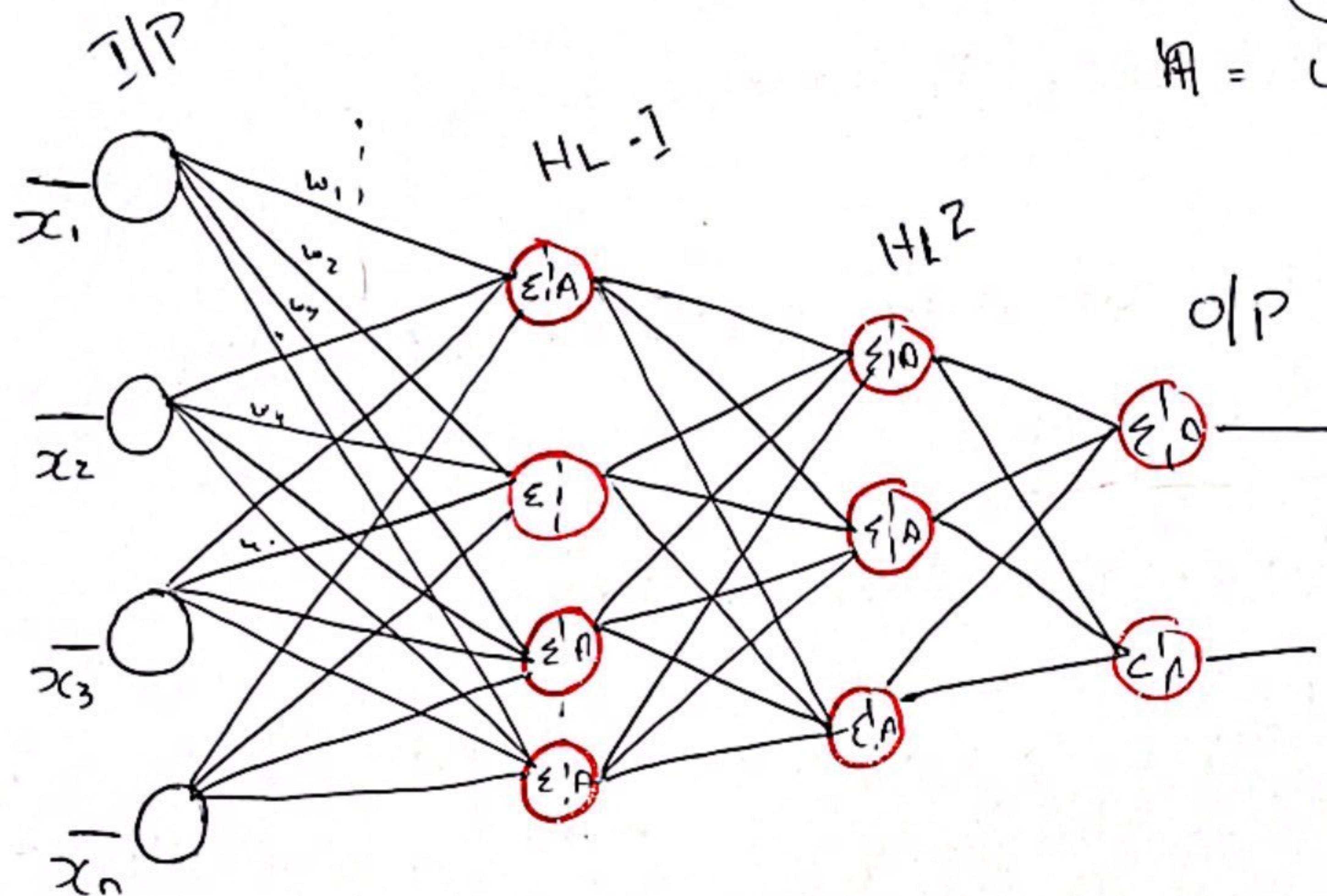
③ Optimizer is used reduce the loss function.

Optimizer: (Gradient descent).



calculate a loss. if high
use optimizer to reduce loss.

Multi Layer Neural network.



$$\Sigma = (\sum w_i x_i) + b$$

$$y_A = q = \frac{1}{1 + e^{-\Sigma}}$$

$$q = \frac{1}{1 + e^{-(\sum w_i x_i + b)}}$$

Back Propagation or weight updation

formula: GD

$$W_{new} = W_{old} - h \frac{\partial H}{\partial w_{old}}$$

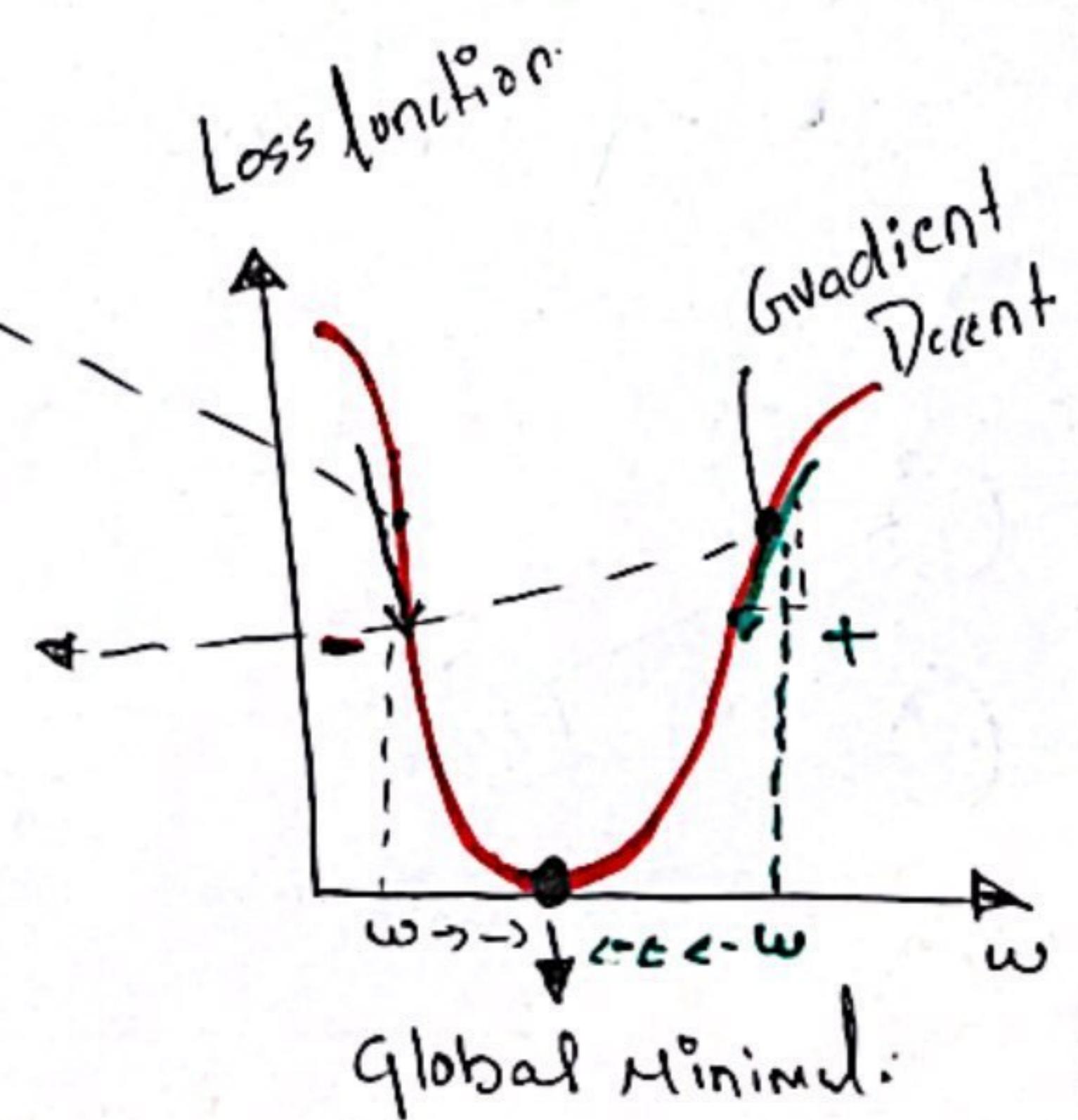
→ learning rate.

$$W_{new} = W_{old} - h(-v_e)$$

$$\Rightarrow W_{new} = W_{old} + h(+v_e)$$

$$W_{new} = W_{old} - n(+v_e)$$

$$W_{new} < W_{old}$$



(6)

② Chain Rule of Differentiation

$$W_{\text{new}} = W_{\text{old}} - \eta \frac{\partial L}{\partial W_{\text{old}}} \rightarrow \text{Learning rate} \cdot \frac{\partial L}{\partial W_{\text{old}}} \rightarrow \text{derivative of loss}$$

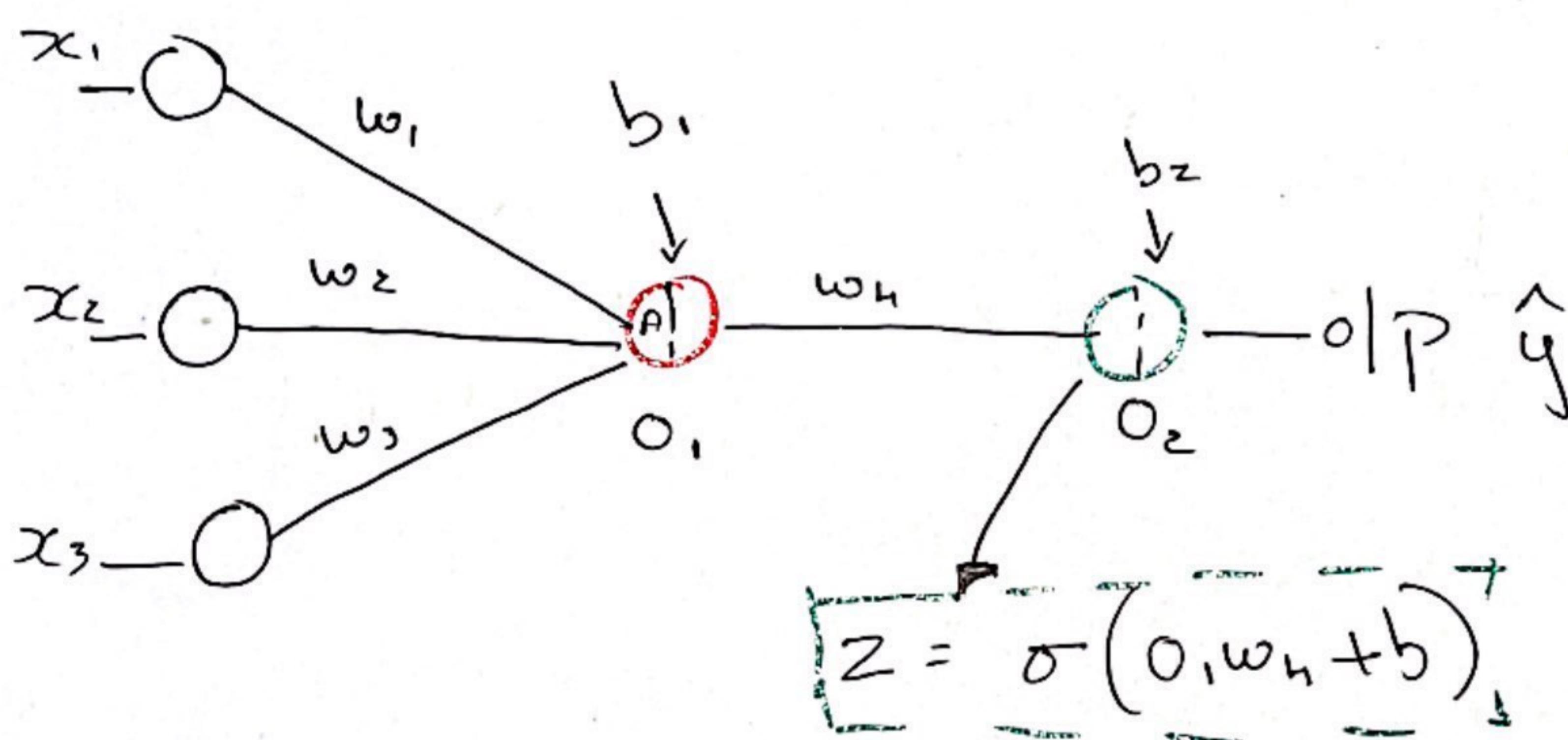
$\frac{\partial L}{\partial W_{\text{old}}} \rightarrow \text{derivative of } W_{\text{old}}$

$$(\sum w_i x_i) + b$$

Learning rate $\boxed{n = 0.01}$

or

$\boxed{n = 0.001}$ best



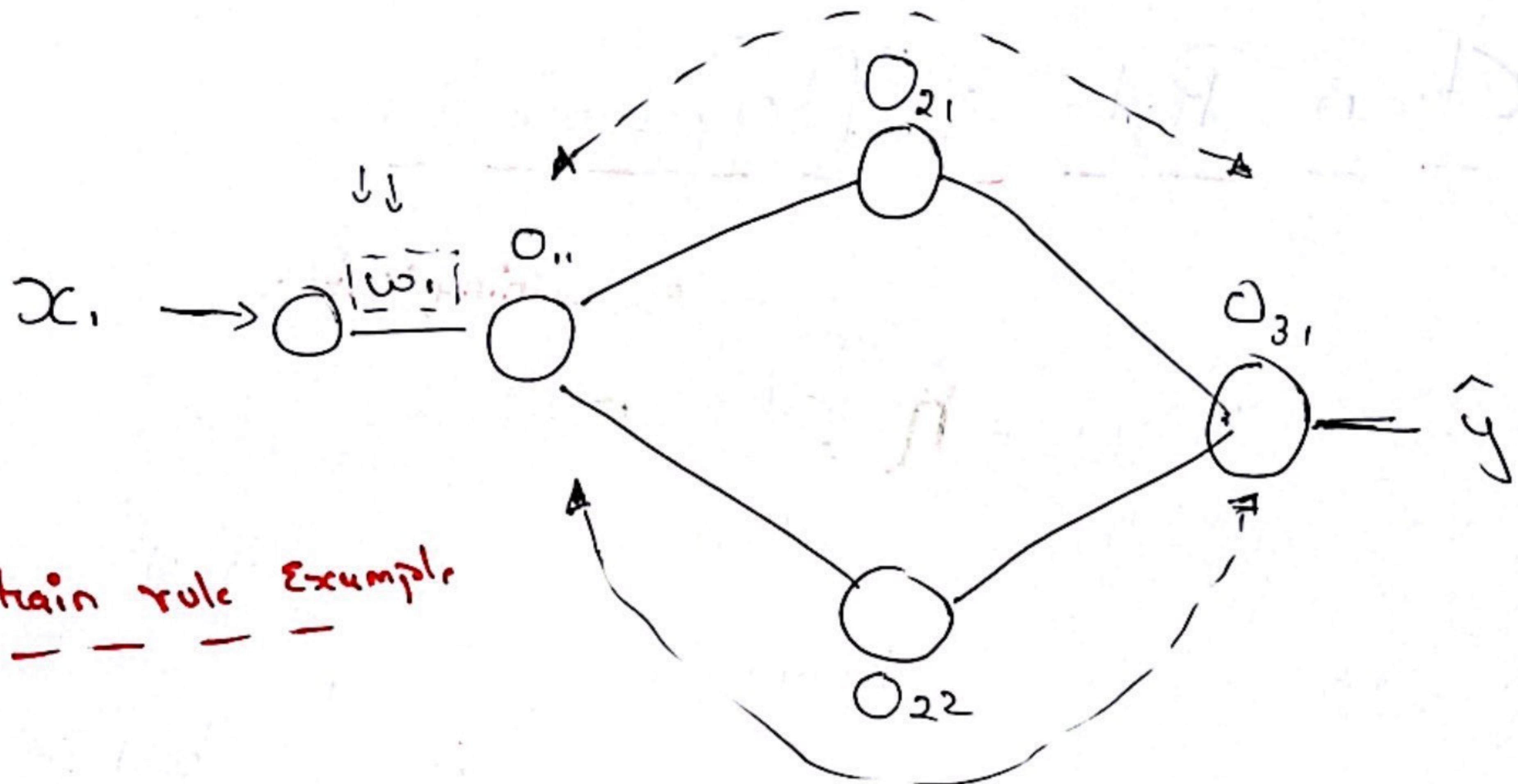
$n = 0.001$

$$W_{4,\text{new}} = W_{4,\text{old}} - n \frac{\frac{\partial L}{\partial Z}}{\frac{\partial Z}{\partial w_{4,\text{old}}}} \quad O1P_2 = \sigma((w_4 * O_1) + b_2)$$

$$\frac{\frac{\partial L}{\partial H}}{\frac{\partial H}{\partial w_{4,\text{old}}}} = \frac{\frac{\partial L}{\partial Z}}{\frac{\partial Z}{\partial O_2}} * \frac{\frac{\partial O_2}{\partial w_{4,\text{old}}}}{\frac{\partial w_{4,\text{old}}}{\partial w_{4,\text{old}}}} \rightarrow \text{chain rule.}$$

Bias update

$$b_{2,\text{new}} = b_{2,\text{old}} - n \frac{\frac{\partial L}{\partial b_2}}{\frac{\partial b_2}{\partial b_{2,\text{old}}}}$$



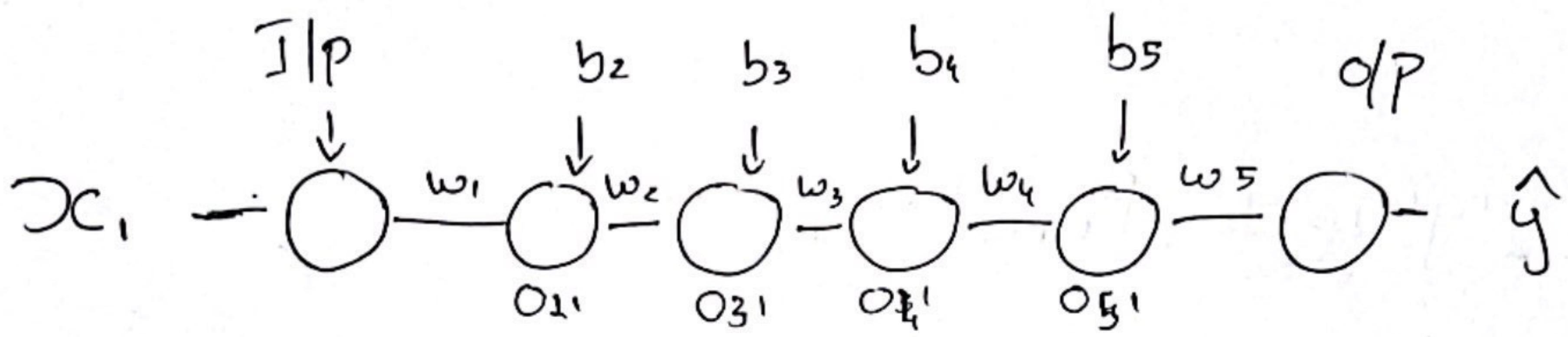
Chain rule Example

$$w_{1, \text{new}} = w_{1, \text{old}} - n \left(\frac{\partial L}{\partial w_{1, \text{old}}} \right)$$

$$\frac{\partial L}{\partial w_{1, \text{old}}} = \left[\frac{\partial L}{\partial O_{31}} \times \frac{\partial O_{31}}{\partial O_{21}} \times \frac{\partial O_{21}}{\partial O_{11}} \times \frac{\partial O_{11}}{\partial w_1} \right] +$$

$$\left[\frac{\partial L}{\partial O_{31}} \times \frac{\partial O_{31}}{\partial O_{22}} \times \frac{\partial O_{22}}{\partial O_{11}} \times \frac{\partial O_{11}}{\partial w_1} \right].$$

Vanishing Gradient Problem



$$\text{MSE} = \boxed{\text{Loss} = \frac{1}{2} (y - \hat{y})^2}$$

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

$$\frac{\partial L}{\partial w_{\text{old}}} = \frac{\partial L}{\partial o_{51}} \times \frac{\partial o_{51}}{\partial o_{41}} \times \frac{\partial o_{41}}{\partial o_{31}} \times \frac{\partial o_{31}}{\partial o_{21}} \times \frac{\partial o_{21}}{\partial w_{\text{old}}}.$$

$$o_{51} = \sigma [(o_{41} \times w_4) + b]$$

↳ Sigmoid funct

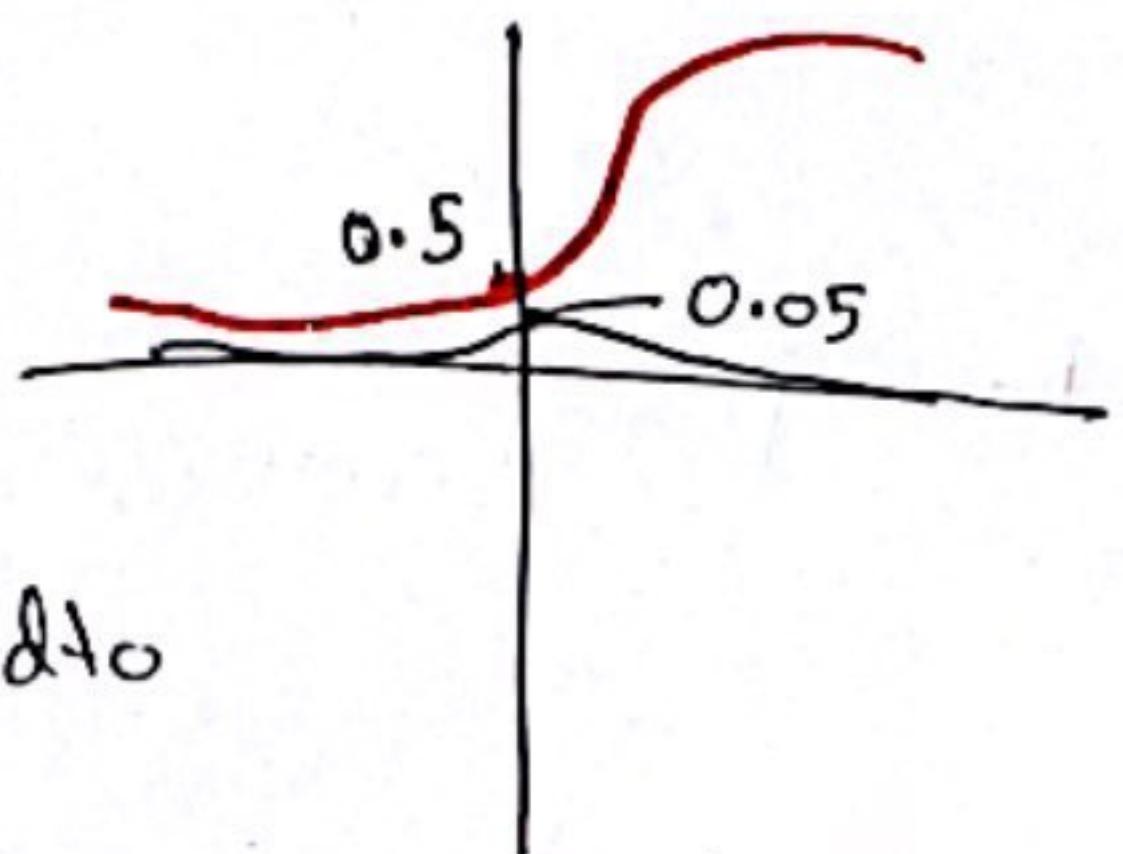
Derivation role is value will be ranging in 0 to 0.25

$$\frac{\partial L}{\partial w_{\text{old}}} = 0.25 \times 0.15 \times 0.10 \times 0.05.$$

When values are hardly updating.

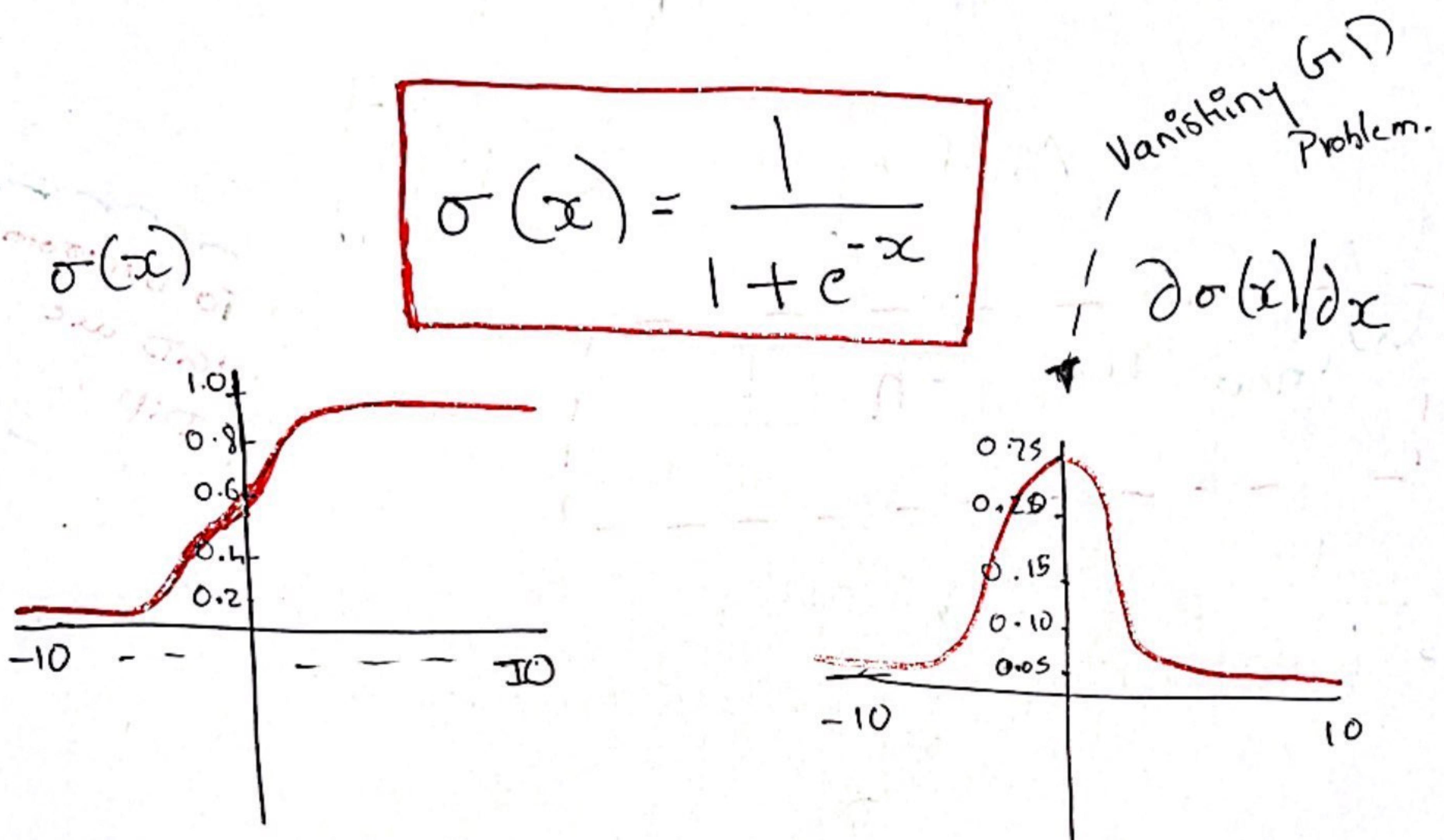
When $w_{\text{old}} \approx w_{\text{new}}$ it is said to

be ~~VGD~~ VGD when the values are too close or not updating.



Activation Functions

Sigmoid Function

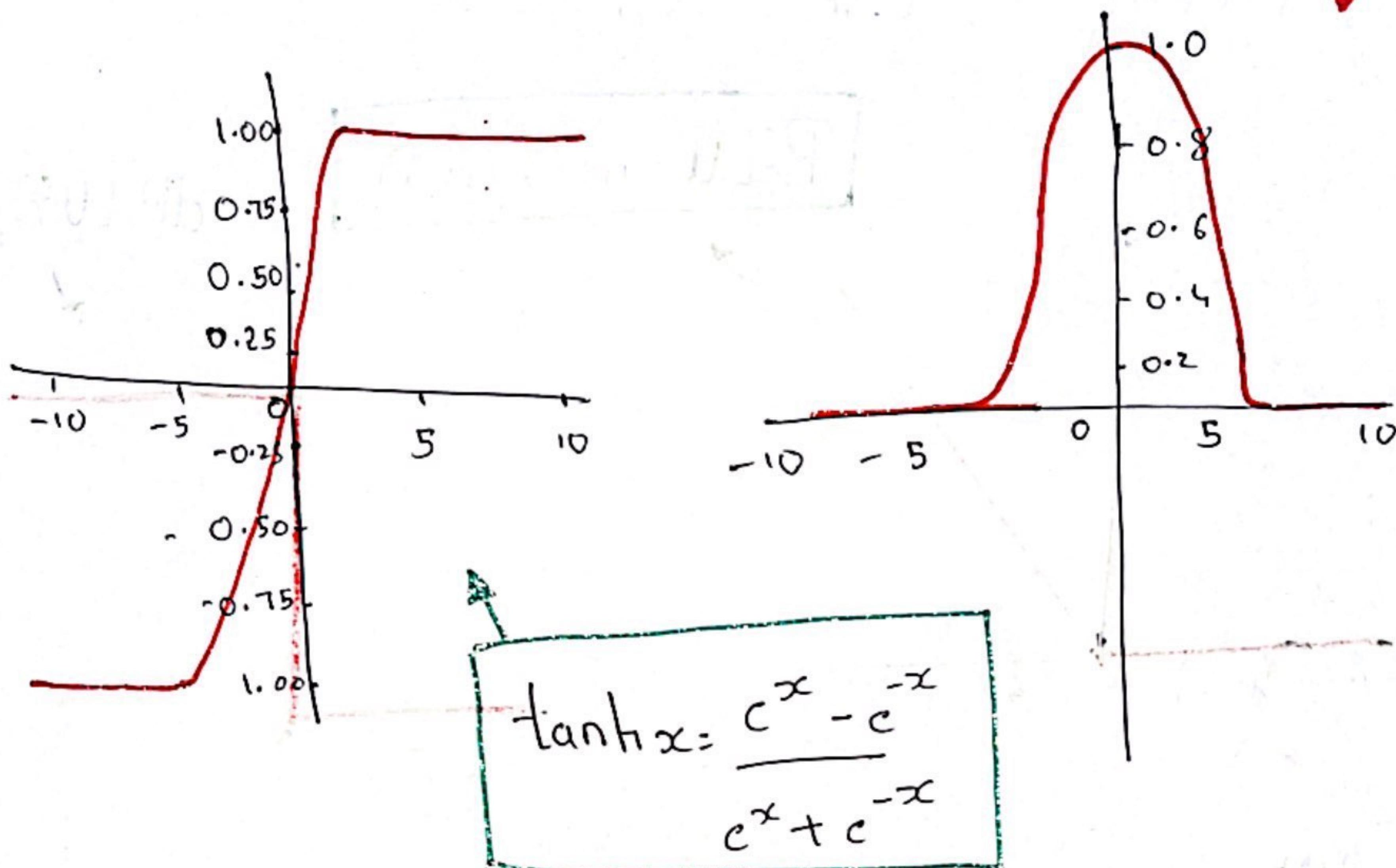


- It is used when we have to predict the probability as an output.
- To add non-linearity in ML Model
- Mainly used for binary classification.
- Advantage: Smooth gradient
Prevent jumps in output value.
- Disadvantage:
 - Prone to gradient vanishing.
 - Power operators are time consuming

②

tanh Function

$$\frac{d \tanh(x)}{dx}$$



- It is a hyperbolic tanh curve.
- Value (-1 to 1) derivative will be in range (0 to 1)

It overcomes the centric problem. (It is non centric).

Note important

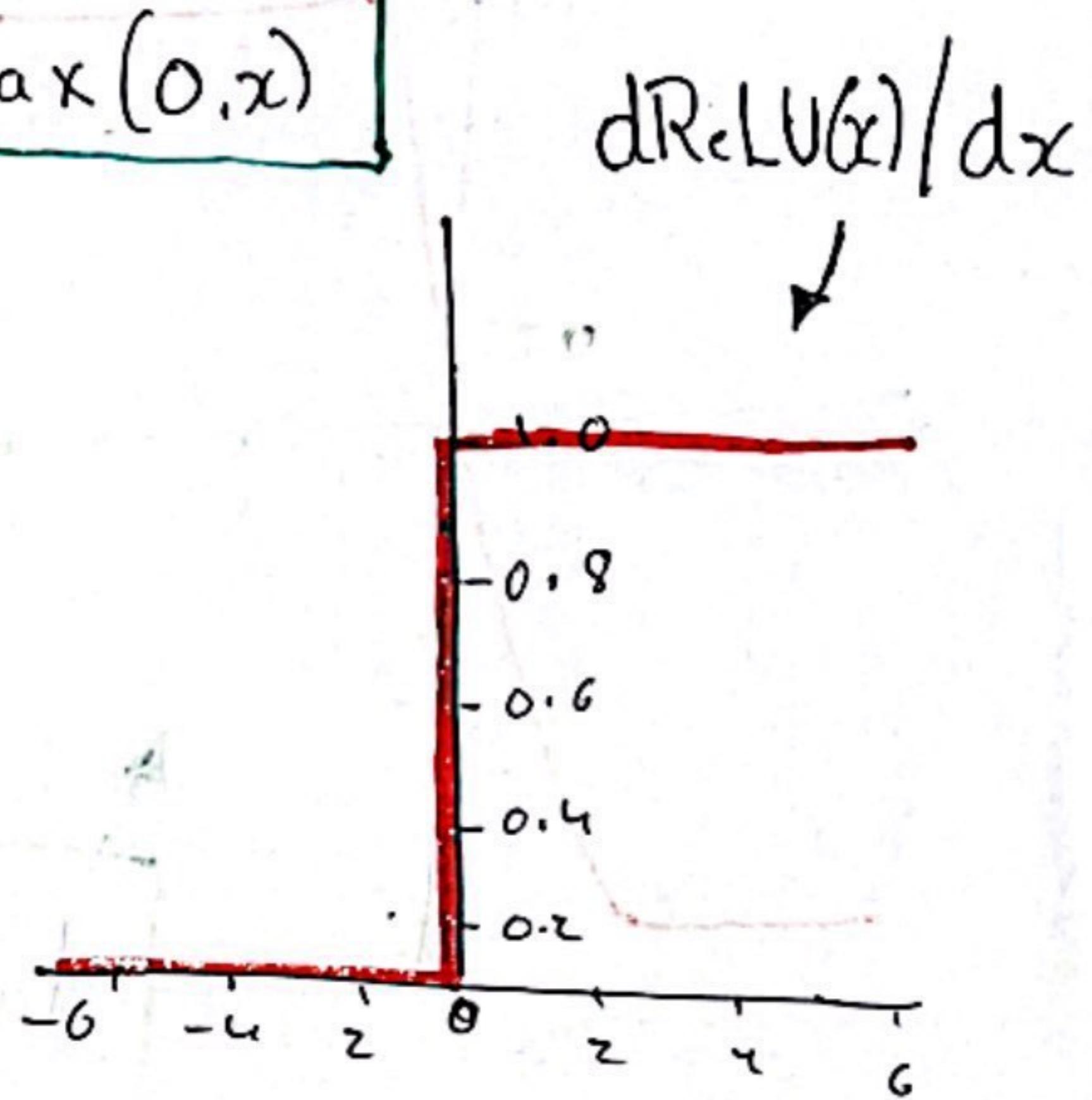
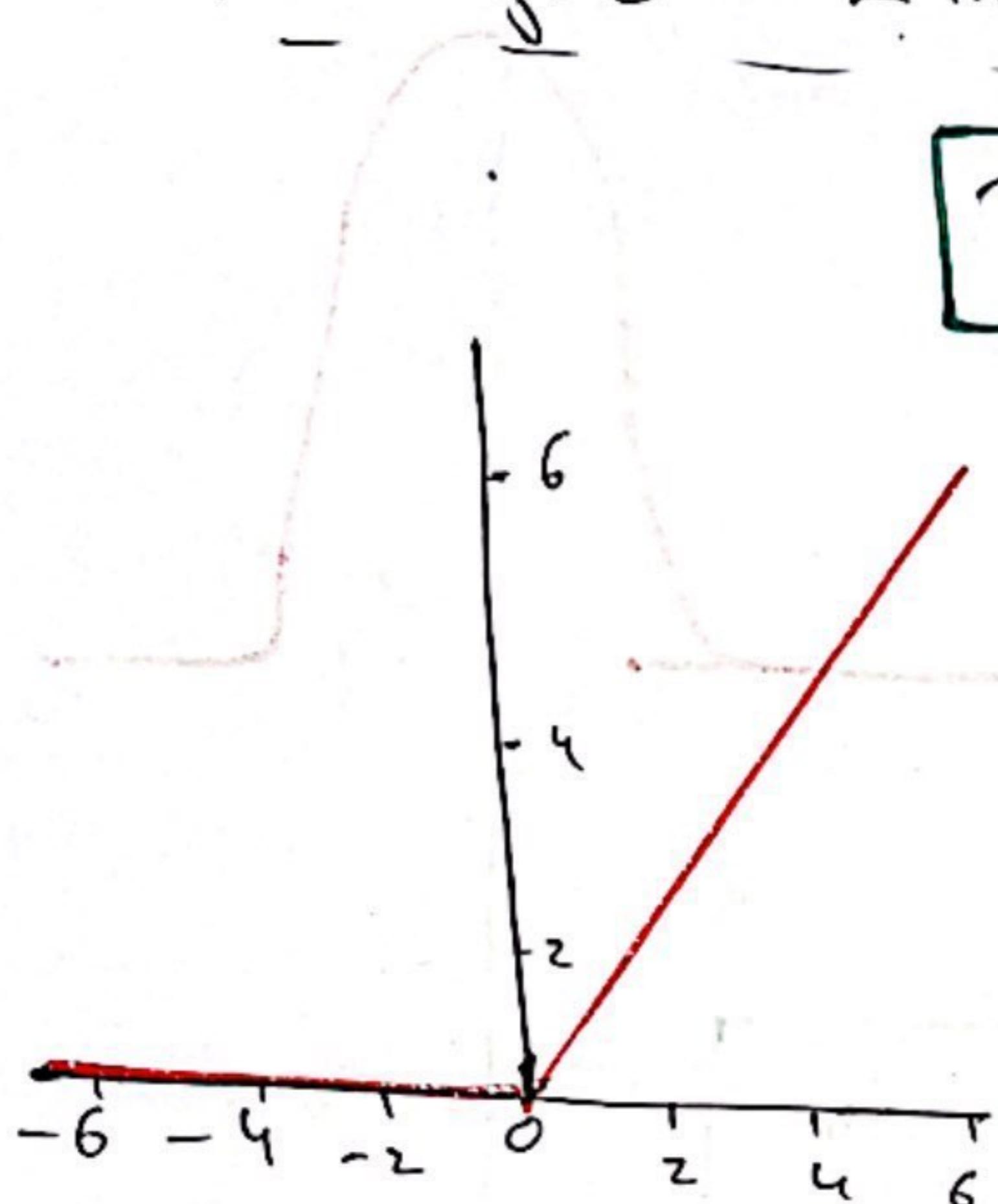
: - in Binary classification problem
the "tanh function is used for the
"Hidden layer" & "Sigmoid function is used for
the output."

- There will be a some small amt of VGD problem

③ ReLU Function

"Rectified Linear Unit"

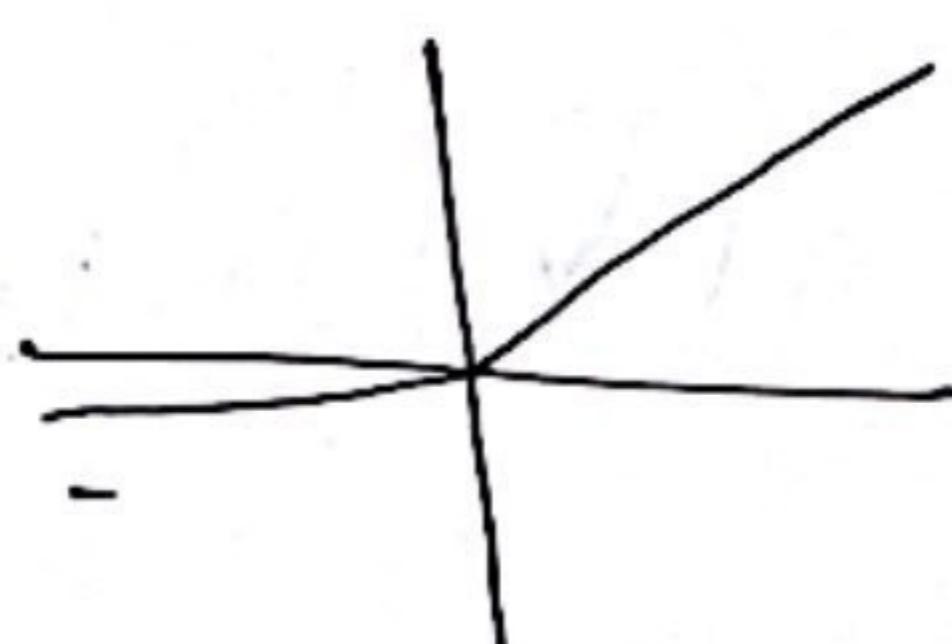
$$\text{ReLU} = \max(0, x)$$



- When x value is -ve ReLU will convert to 0.
- ④ When x value is +ve x value will be same.
- derivative will be 0 or 1
- It solves the problem of vanishing gradient.
- When input is -ve ReLU will be inactive. Neuron will die.

④ Leaky ReLU Function

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

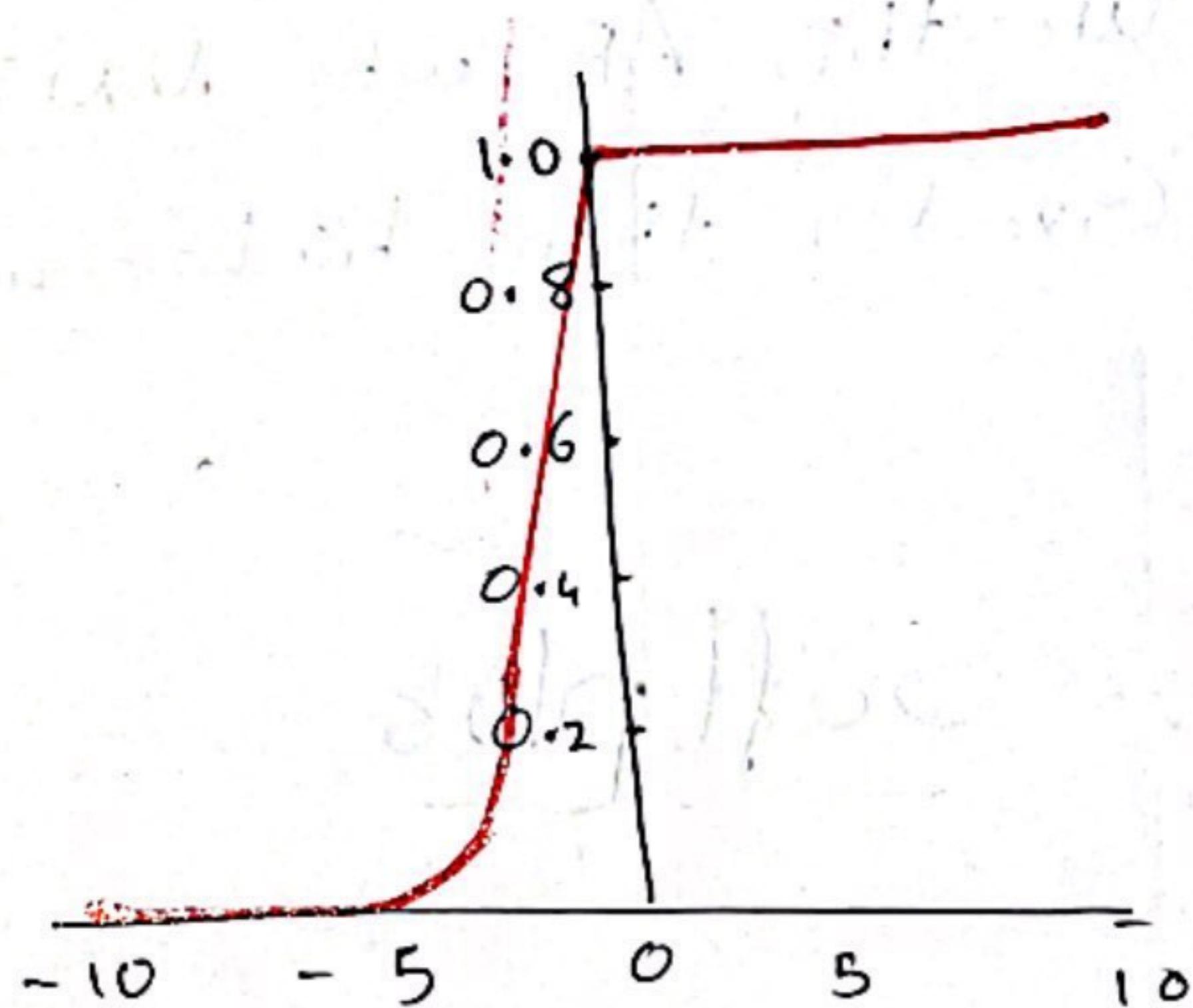
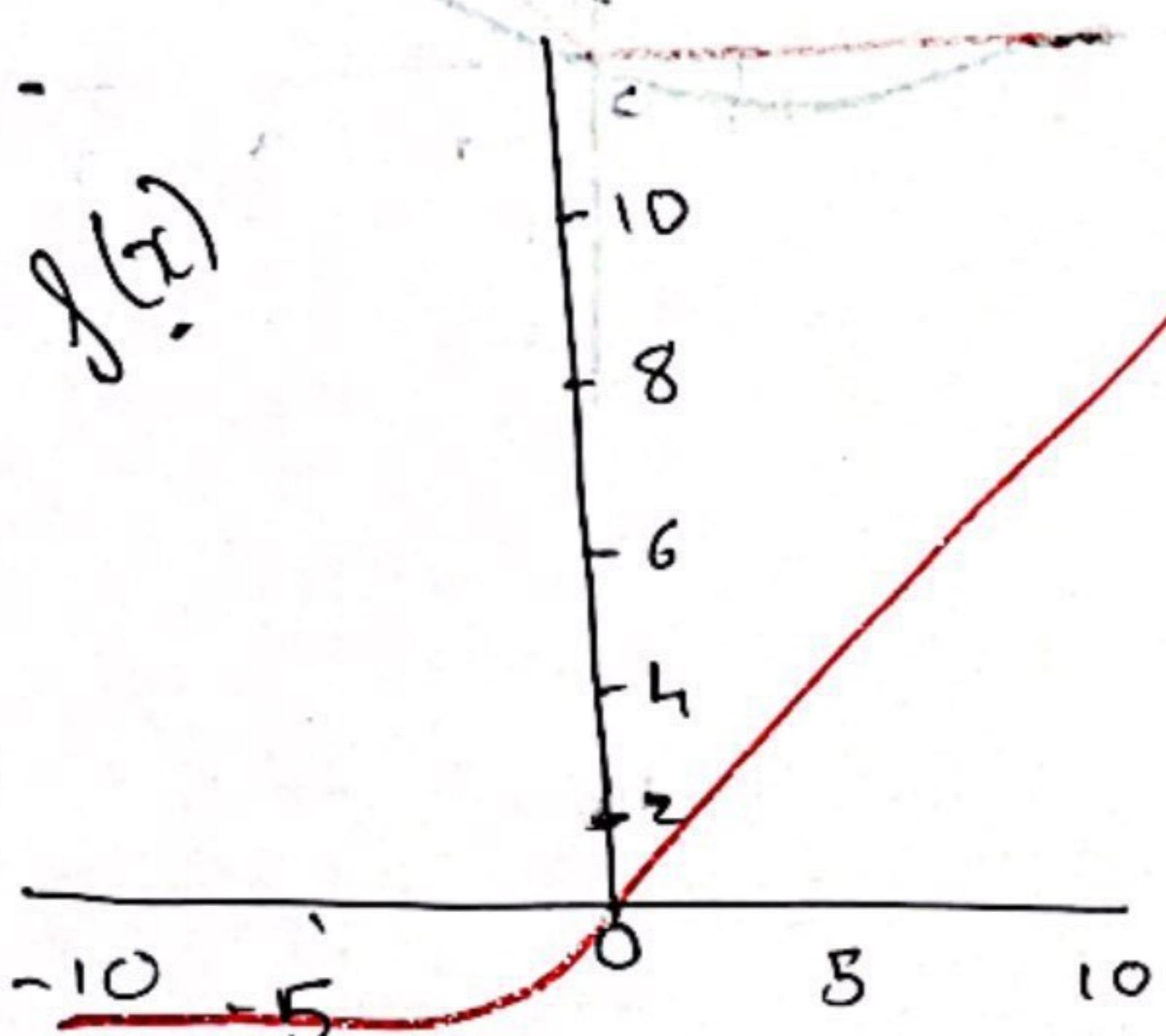


- -ve value will be there in small amt.
- It solves the neuron problem.

⑤ ELU (Exponential Linear Unit) Function

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha(e^{-x} - 1) & \text{otherwise} \end{cases}$$

$$\partial f(x) / \partial x$$

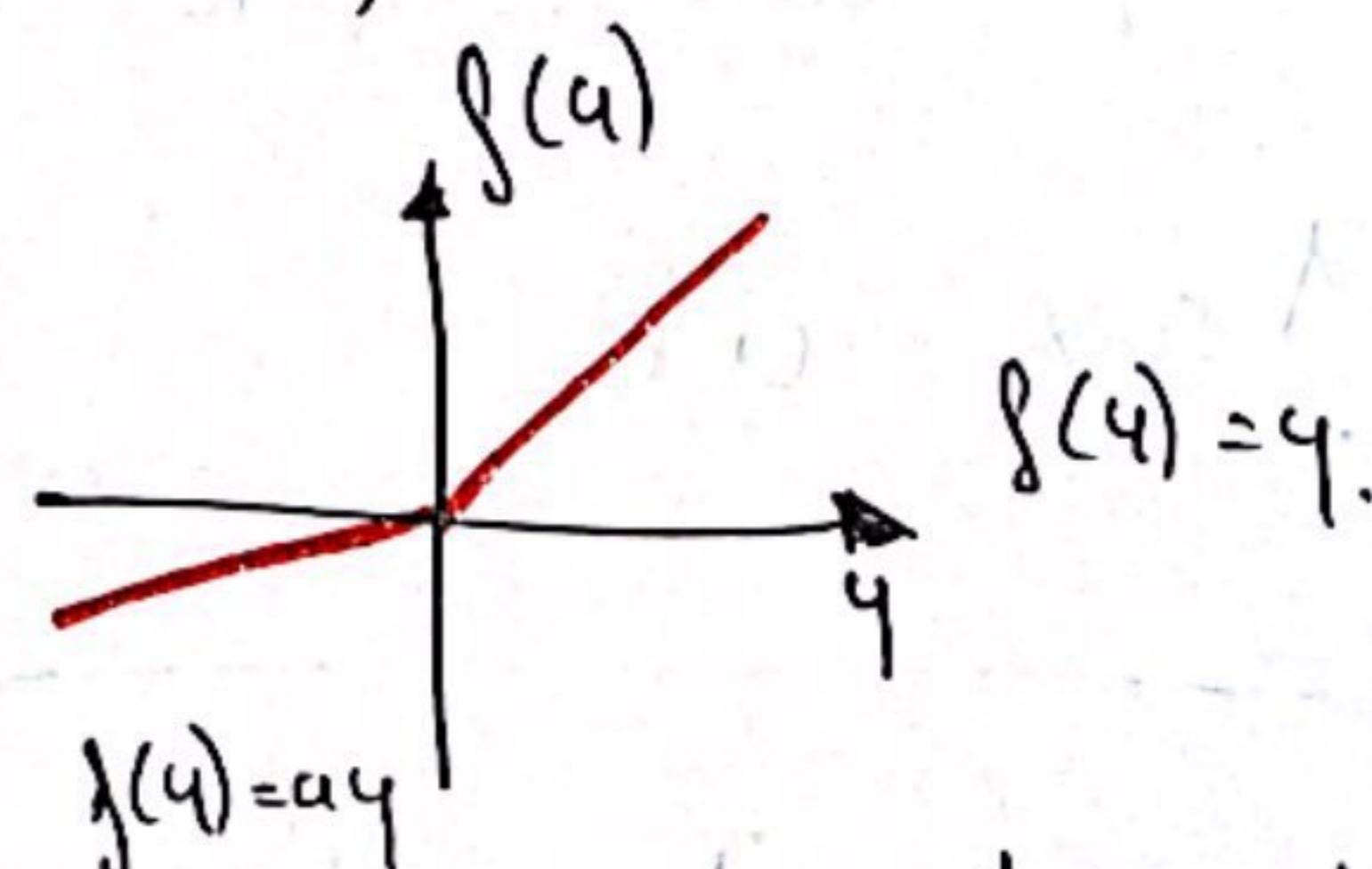
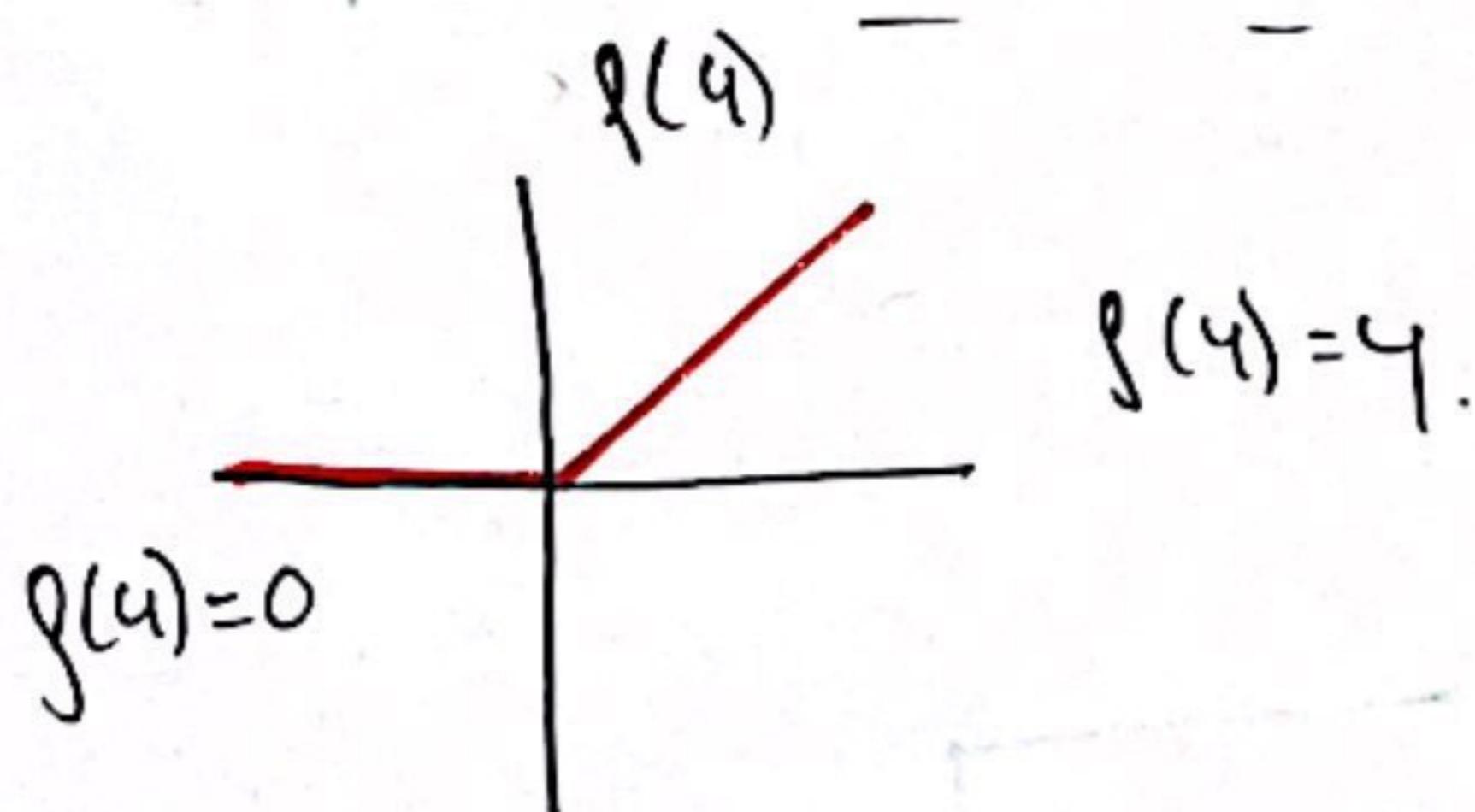


→ ELU is also proposed to solve ReLU problem.

- No Dead ReLU issue.

→ Disadvantage: more computational intensive.

⑥ PReLU (Parametric ReLU)



$$f(y_i) = \begin{cases} y_i & \text{if } y_i > 0 \\ a_i y_i & \text{if } y_i \leq 0 \end{cases}$$

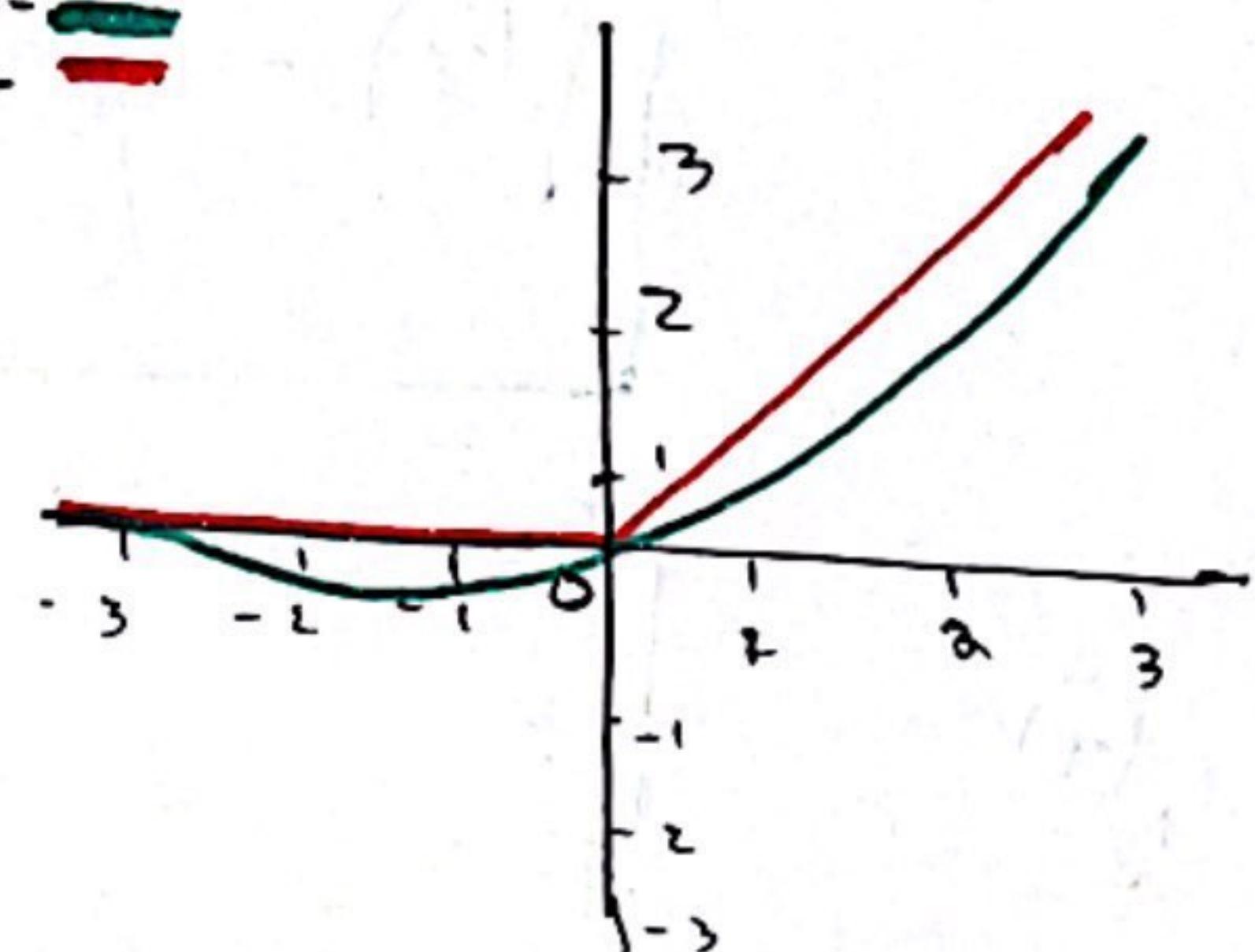
- if $a_i = 0$, it's ReLU
- if $a_i > 1$, it's leaky ReLU
- if a_i is a PReLU

- Value will be in Range $(0 \text{ or } 1)$

⑦ Swish (A Self Gated) Function. ReLU

$$Y = X * \text{Sigmoid}(X)$$

Swish
ReLU



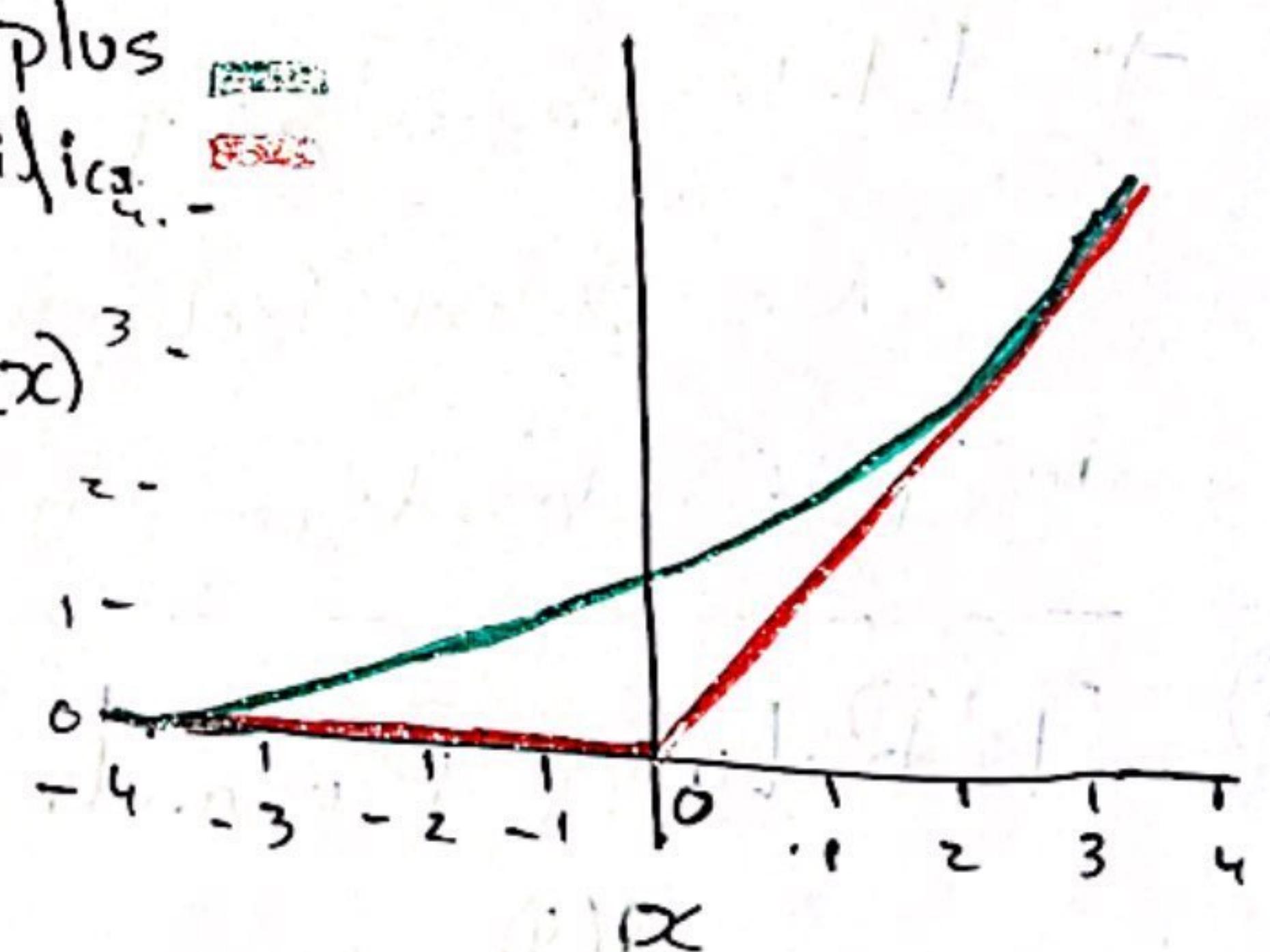
We this AF whe NN is
Greater than 1 layer.

⑧ Softplus

$$f(x) = \ln(1 + e^x)$$

Softplus
Rectifica.

Nonlinearity.



① Range $(0, +\infty)$

① Similar & relatively Smooth.

① It is unilaterial Suspension.

⑨ Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

① Parameters K.

① ReLU & Leaky ReLU is Generalised Here.

① Output will be largest activation value.

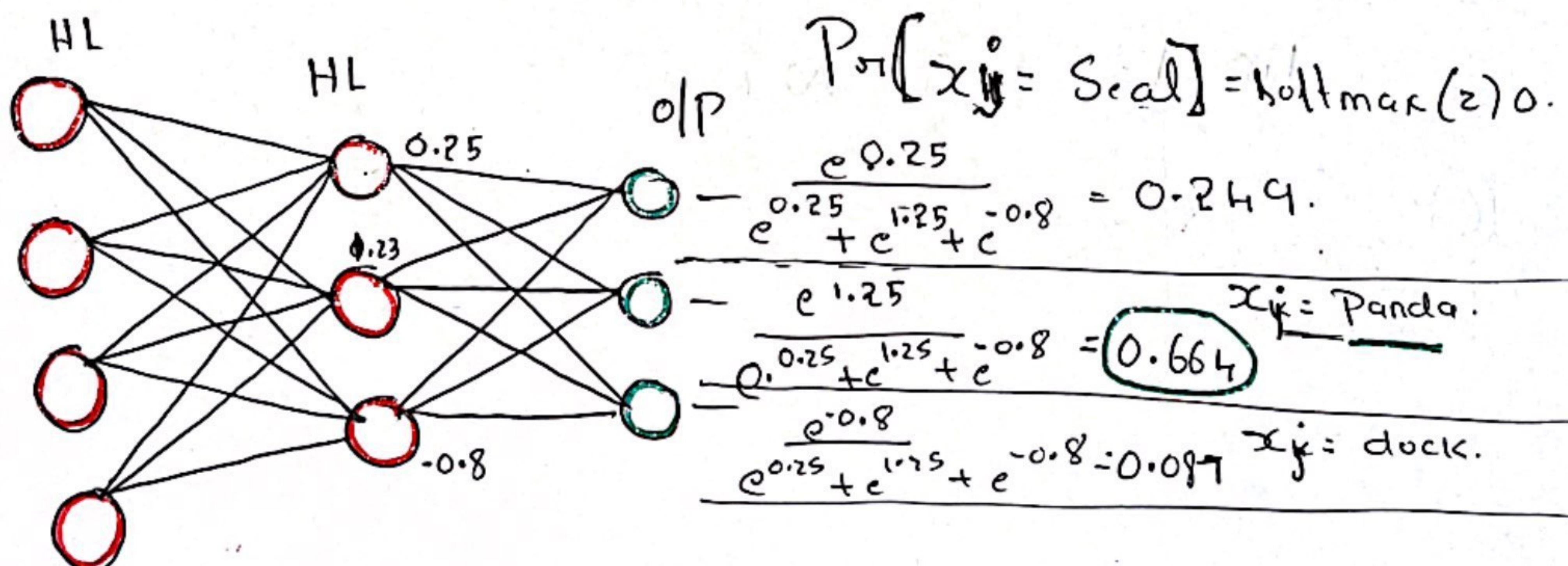
10

Softmax

$$S(x_j) = \frac{e^{x_j}}{\sum_{k=1}^K e^{x_k}}, \quad k = 1, 2, \dots, K.$$

- Softmax is used when output > 2 layers.
- Output will be in the form of probability. Highest Probability will be considered.
- It is always applied at "output layer".

Example: identify.



Note : Sigmoid and Softmax are always used for output layers if ≤ 2 out layer Sigmoid
if > 2 out layer Softmax.

15

Notes

- ① ReLU is Most Preferred for hidden layers.
- ②
 - Binary classification output layer use Sigmoid.
 - Multi classification output layer use Softmax.
- ③ for regression like Problem. to achieve continuous value use linear activation function in Output layer.
- ④ Tanh is used in Hidden layer in Binary classification.
- ⑤ Swish is used only when there is more than 40 layers in N.N.
- ⑥

Evaluation Methods.

① Regression.

a) Mean Square Error (MSE)

$$\text{Loss function: } \frac{1}{2} (y - \hat{y})^2$$

Single Record.

$$\text{cost function: } \frac{1}{2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Batch of Records.

$$\text{Loss function: } \frac{1}{2n} (y - \hat{y})^2$$

$$\text{cost function: } \frac{1}{2n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Advantage: ① Differentiable.

- ② It has only one local or Global Minima.
- ③ It converges faster.

disadvantage: Robust to outliers.

b) Mean absolute error (MAE)

$$\text{Loss function: } \frac{1}{n} |y - \hat{y}|$$

$$\text{cost function: } \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Advantage: overcomes the Robust to outliers.

- No Squar. is done.

disadvantage: Time consume.

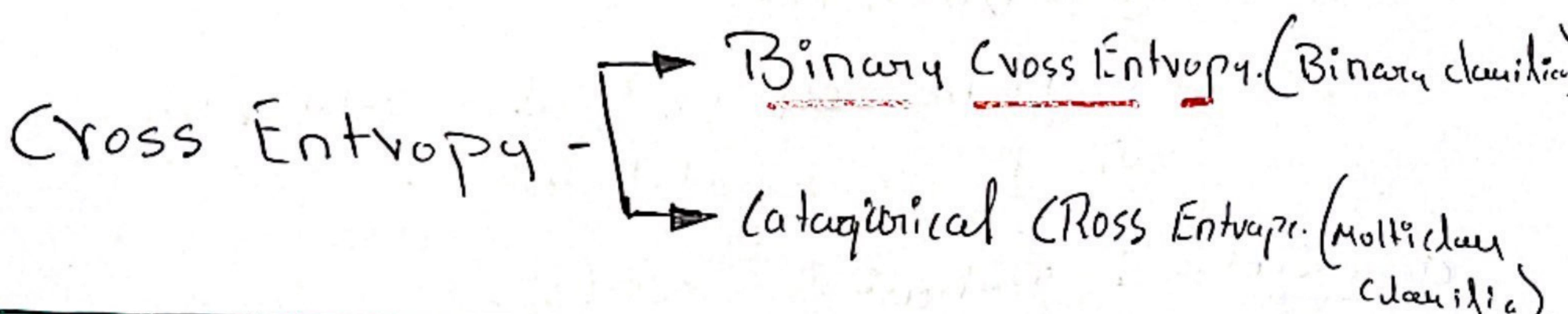
⑥ Huber loss

combination of both. MSE & MAE. Hyperparameters

$$\text{Loss} = \begin{cases} \frac{1}{2}(\hat{y} - y)^2 & \text{if } |\hat{y} - y| \leq \delta \\ |\hat{y} - y| - \frac{1}{2}\delta^2 & \text{otherwise.} \end{cases}$$

No outliers.

② Classification



a) Binary Cross Entropy.

$$\text{Loss} = -y * \log(\hat{y}) - (1-y) * \log(1-\hat{y}). \approx \text{LR}$$

$$\text{loss} = \begin{cases} -\log(1-\hat{y}) & \text{if } y=0 \\ -\log(\hat{y}) & \text{if } y=1 \end{cases}$$

$$\hat{y} = \text{Sigmoid function: } \frac{1}{1+e^{-x}}$$

b

Categorical Cross Entropy.

$$\text{Loss}(x_i, y_i) = - \sum_{j=1}^c y_{ij} * \ln(\hat{y}_{ij})$$

$$y_i = [y_{i1}, y_{i2}, y_{i3}, \dots, y_{ic}]$$

$$y_{ij} = \begin{cases} 1 & \text{if the element is in class.} \\ 0 & \text{otherwise.} \end{cases}$$

$$\hat{y}_{ij} = \text{Soft Max Activation function}$$

Deep learning Optimizers terms

- ① **Epoch** - The number of times algorithm runs on the whole training dataset.
- ② **Sample** - A single row of dataset
- ③ **Batch** - the number of samples to be taken for updating model parameters.
- ④ **Learning Rate** - It is a parameter that provides the model a scale of how much model parameters weight should be updated.

Optimization

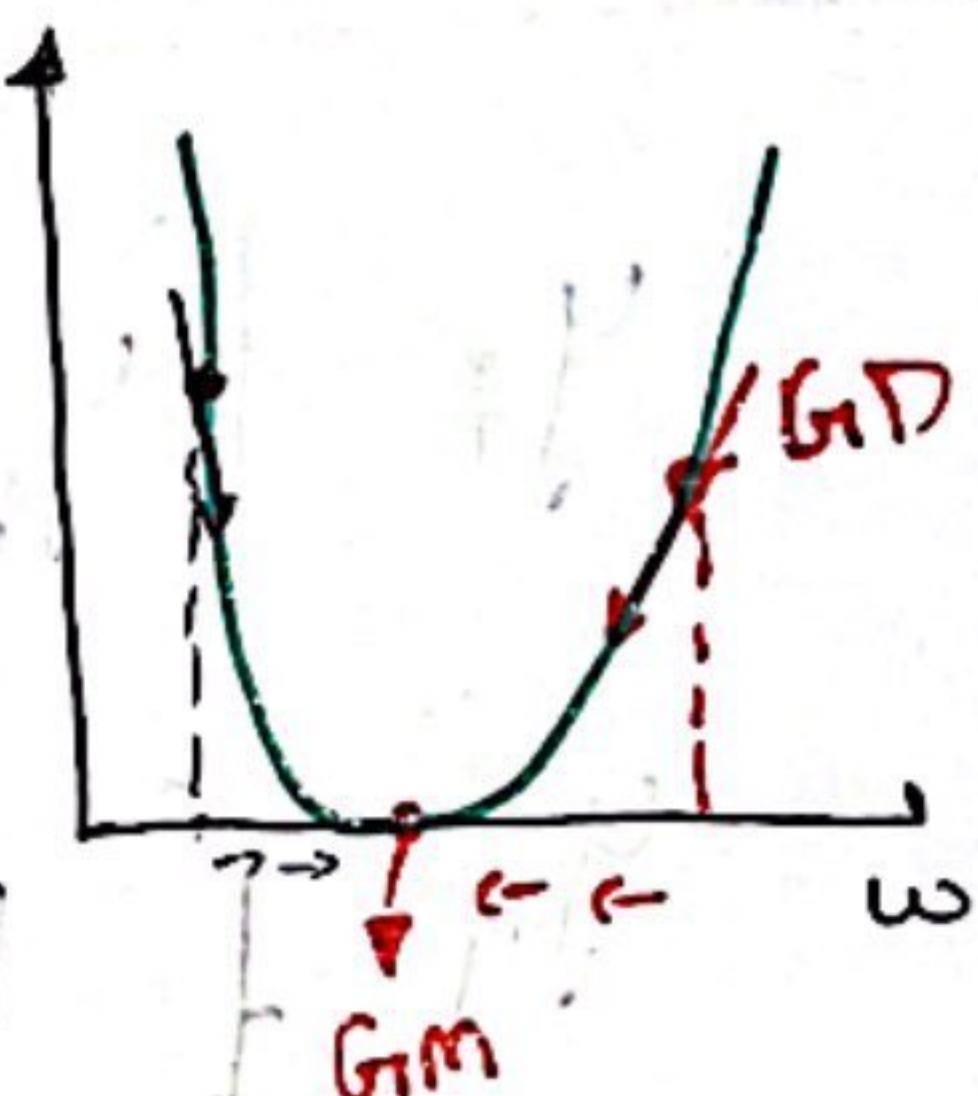
① Gradient Descent

weight updation formula.

$$W_{\text{new}} = W_{\text{old}} - h \frac{\partial h}{\partial w_{\text{old}}}$$

learning rate chain rule derivative.

loss/cost



$$\underline{(\sum w_i x_i) + b}$$

• Disadvantages

* Adotilue this if dataset is huge.

→ Resources Expensive.

② Stochastic Gradient Descent

- To overcome GD computation cost Δ GD was derived.
- Stochastic means randomness. we randomize Select the Samples & make that as Batch.
"Batch size = No of Random Samples."

$$W_{\text{new}} = W_{\text{old}} - h \frac{\partial h}{\partial w_{\text{old}}}$$

- SGD takes more iterations to reach local minima.

- It is in increase training time but decreases computation cost of fed Single record to algorithm

Disadvantage

- Increase in training time.
- Increase in noise cuz of randomness.

③ Mini-Batch SGD

- Batch Size is a Parameter to make group of Random Sample.

- How to choose Perfect batch size?

Ans:-

$$\text{Batch Size} = 10,000 \quad \text{Dataset} = 1,00,000$$

Iteration = 10. Randomly Selected.

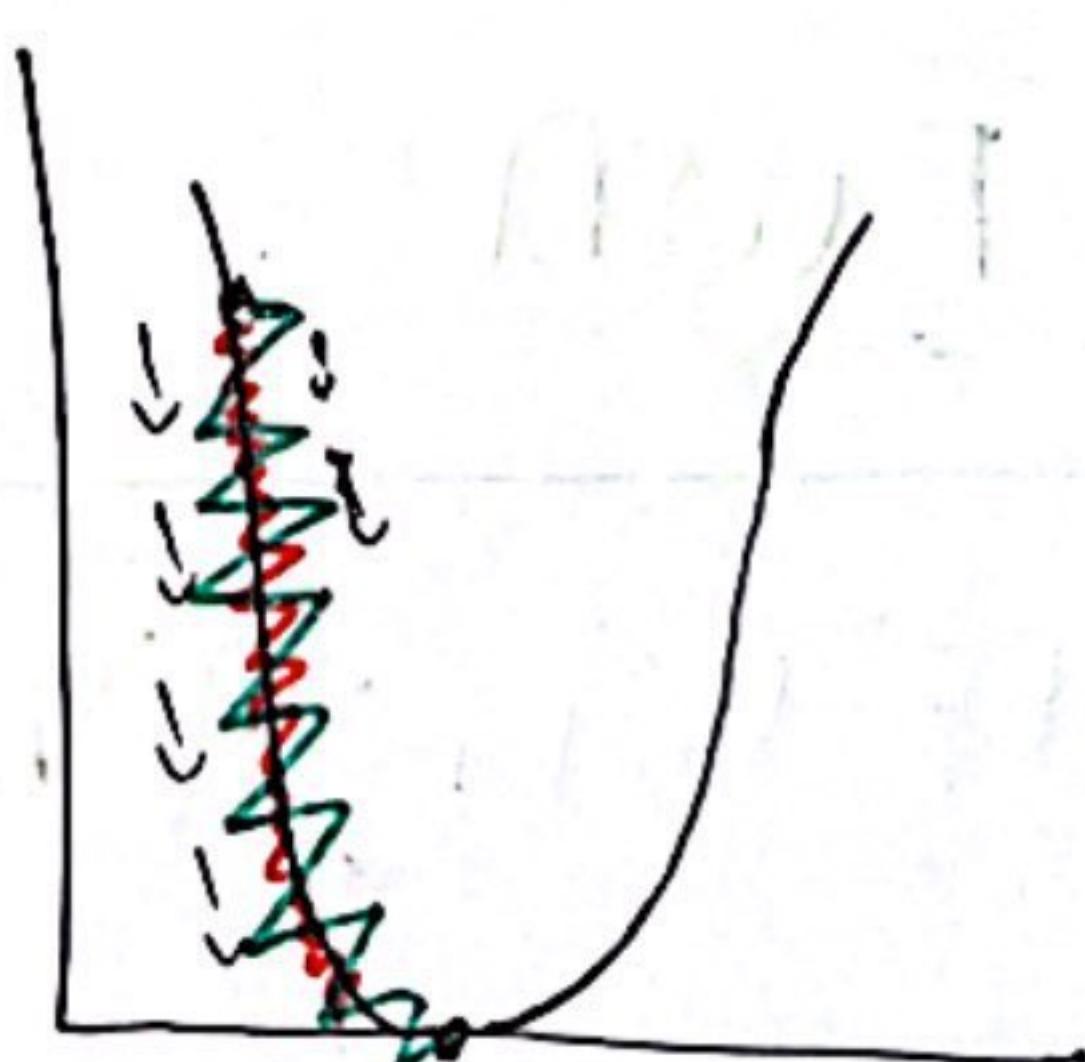
Once cycle is 1 Epoch.

- Convergence is better.
- time will improve.
- comparatively less resources expensive.
- Reduces the noise

Disadvantage

- increasing in train time.
- Noise in data.

MB SGD - —
SGD. - —



④ SGD with Momentum

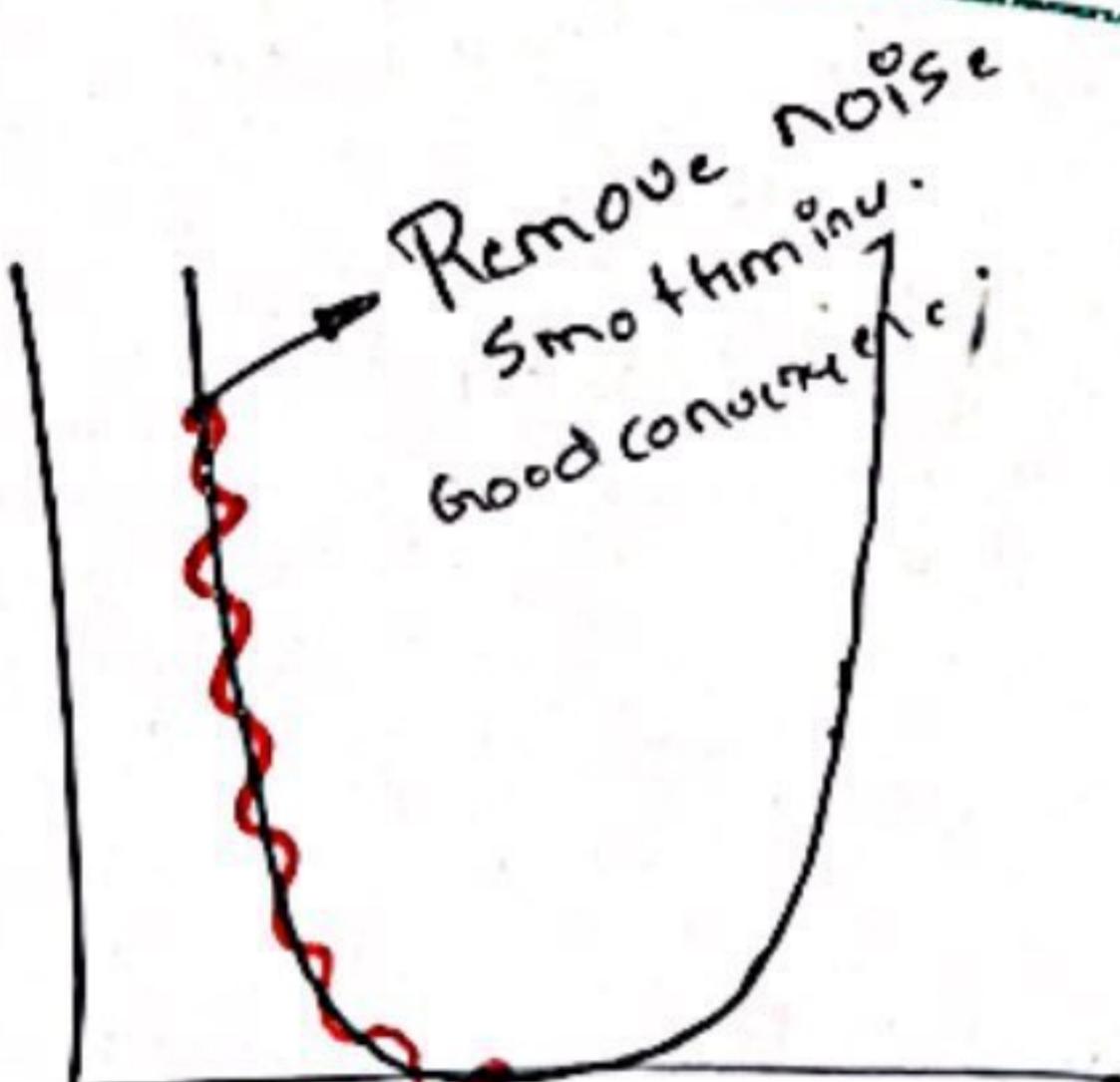
To overcome the Problem of Noise we use SGD with Momentum.

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

$$b_{\text{new}} = b_{\text{old}} - \eta \left[\frac{\partial L}{\partial b_{\text{old}}} \right]$$

We use "Exponential weighted average" mixed in time series.

$$w_t = w_{t-1} - \eta \frac{\partial L}{\partial w_{t-1}}$$



$$t_1, t_2, t_3, t_4, \dots, t_n$$

$$a_1, a_2, a_3, a_4, \dots, a_n$$

$$\beta = 0.10 \\ = 0.95$$

$$V_{t_1} = a_1$$

$$V_{t_2} = \beta * V_{t_1} + (1-\beta) * a_2 \\ = (0.95) * V_{t_1} + (0.05) * a_2$$

$$V_{t_3} = \beta * V_{t_2} + (1-\beta) * a_3$$

EWA

$$w_t = w_{t-1} - \eta \bar{V}_{dw}$$

$$V_{dw} = \beta * V_{dw_{t-1}} + (1-\beta) * \frac{\partial L}{\partial w_{t-1}}$$

weight updating formula

⑤ AdaGrad (Adaptive Gradient Descent)

It mainly concentrate on learning rate η

- Learning rate is constant. in this AdaGrad.
h will be adaptive.

$$w_t = w_{t-1} - \eta \frac{\partial L}{\partial w_{t-1}}$$

$$\eta' = \frac{\eta}{\sqrt{\sum_i (\frac{\partial L}{\partial w_i})^2 + \epsilon}}$$

Small num
not to make
zero.

$$h_t = \sum_{i=1}^n \left(\frac{\partial L}{\partial w_i} \right)^2$$

- Note: η' Should be

decreased when it reaches near
to Global minima.

Disadvantage

• Decreasing fast so hard to find GM

→ In deep neural networks. when η will decrease.
it might effect

Advantage

- It reaches convergence at a higher speed.

⑥ Adadelta and RMSProp

Exponential weight Average.

$$S\omega = 0$$

$$S\omega_t = \beta S\omega_{t-1} + (1-\beta) \left(\frac{\partial L}{\partial \omega_{t-1}} \right)^2$$

① Main aim is to control increase in $\left(\frac{\partial L}{\partial \omega_t}\right)^2$ ↑↑ R
 So h will be controlled.

② Adam Optimizer ★★★

Adaptive moment estimation.

It is the combination of Momentum + RMSProp.

Both Adaptive learning rate and Smoothing is also done together.

① Smoothing

② Learning is Adaptive.

Procs

$$w_t = w_{t-1} - \eta' V_{\partial w} \rightarrow 1a$$

$$\eta' = \frac{\eta}{\sqrt{S_{\partial w}} + \epsilon} \rightarrow 2$$

$$b_t = b_{t-1} - \eta' V_{\partial b} \rightarrow 1b$$

$$V_{\partial w_t} = \beta * V_{\partial w_{t-1}} + (1-\beta) \frac{\partial L}{\partial w_{t-1}} \rightarrow 3$$

$$V_{\partial b_t} = \beta * V_{\partial b_{t-1}} + (1-\beta) \frac{\partial L}{\partial b_{t-1}} \rightarrow 4$$

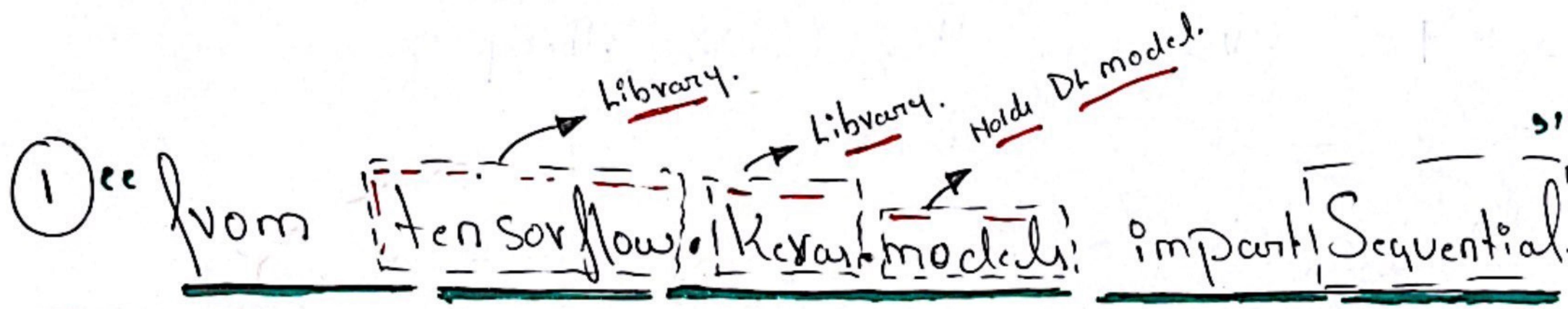
where $S_{\partial w}$ is

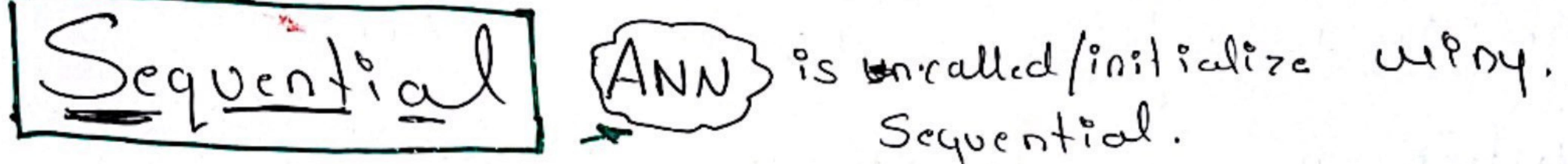
$$S_{\partial w_t} = \beta S_{\partial w_{t-1}} + (1-\beta) \left(\frac{\partial L}{\partial w_{t-1}} \right)^2 \rightarrow 5$$

③ in ②

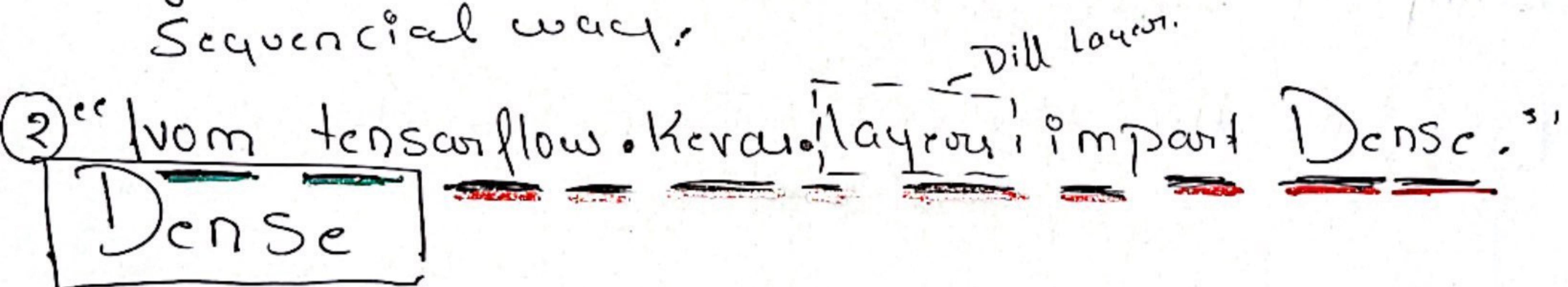
②, ③, ④ in 1a & 1b

Deep learning Libraries

① "from tensorflow.keras.models import Sequential" 

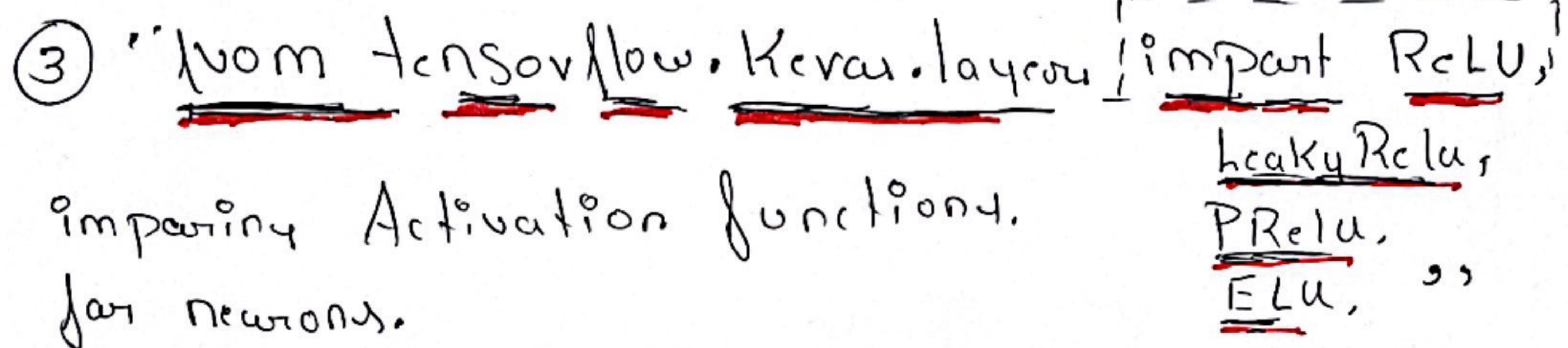
Sequential  is unrolled/initialize w.r.t. Sequential.

whole set of input will be taken in the form of blocks. It can perform all the forward and Backward Propagation in Sequential way.

② "from tensorflow.keras.layers import Dense." 

Creating and managing the neurons. And for modifying the hidden layer neurons.

- Create hidden layers
- Input and Output layers creation.

③ "from tensorflow.keras.layers import ReLU," 

improving Activation function for neurons.

ReLU,
LeakyReLU,
PReLU,
ELU,

④ "From Tensorflow. Keras. layers import Dropout."

- To overcome the Overfitting Problem in Dataset we use Dropout.
- used to solve complex Problem.
- Similar to regularization in ML.

When we specify the Parameter of dropout.
Ex: if 0.3 is specified. ^{Assigned} 30% of the neurons gets deactivated.

⑤ Import Tensorflow.

opt = tensorflow.keras.optimizers.Adam (learning_rate=0.01)
→ Specify in `compile()`