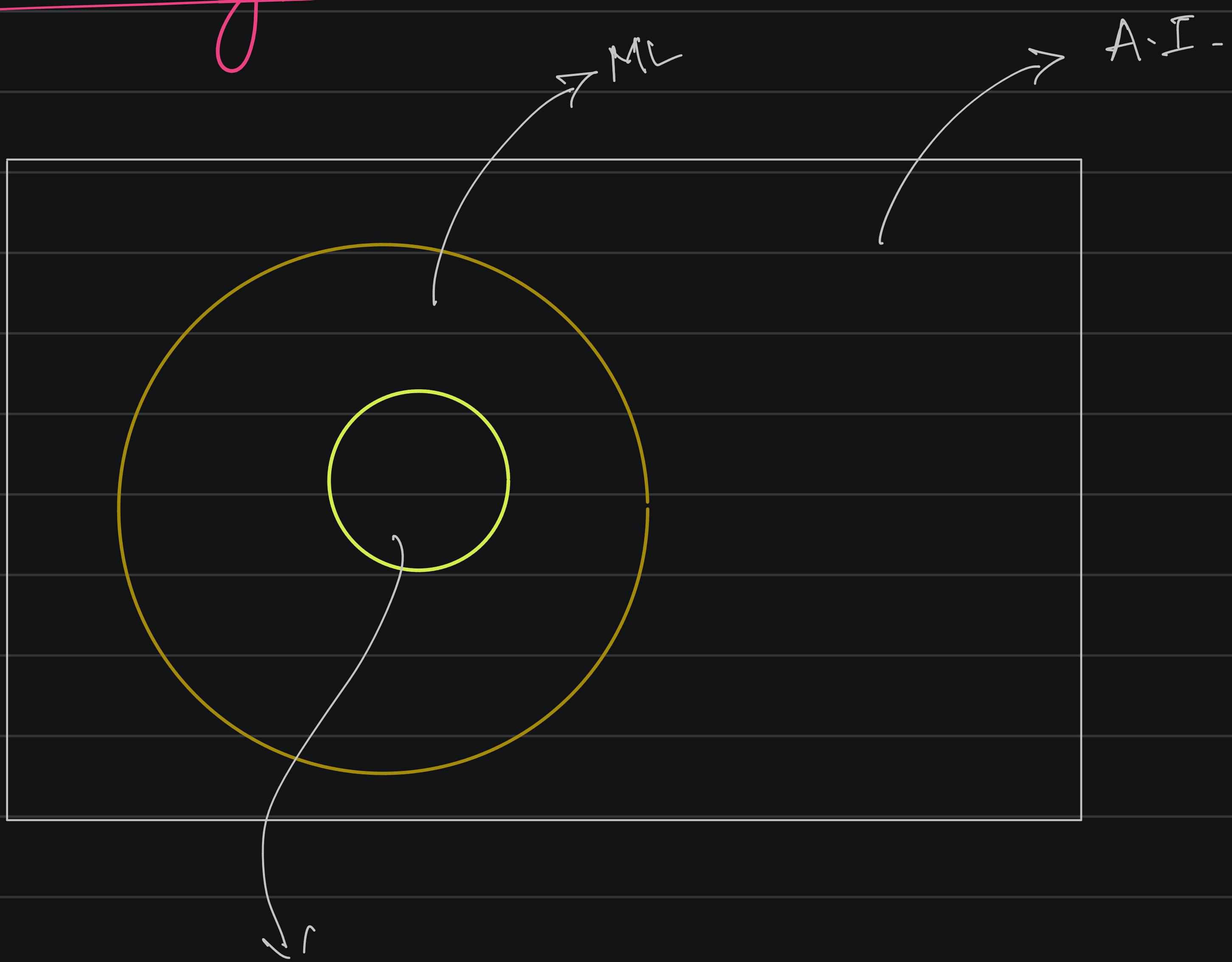


Deep Learning

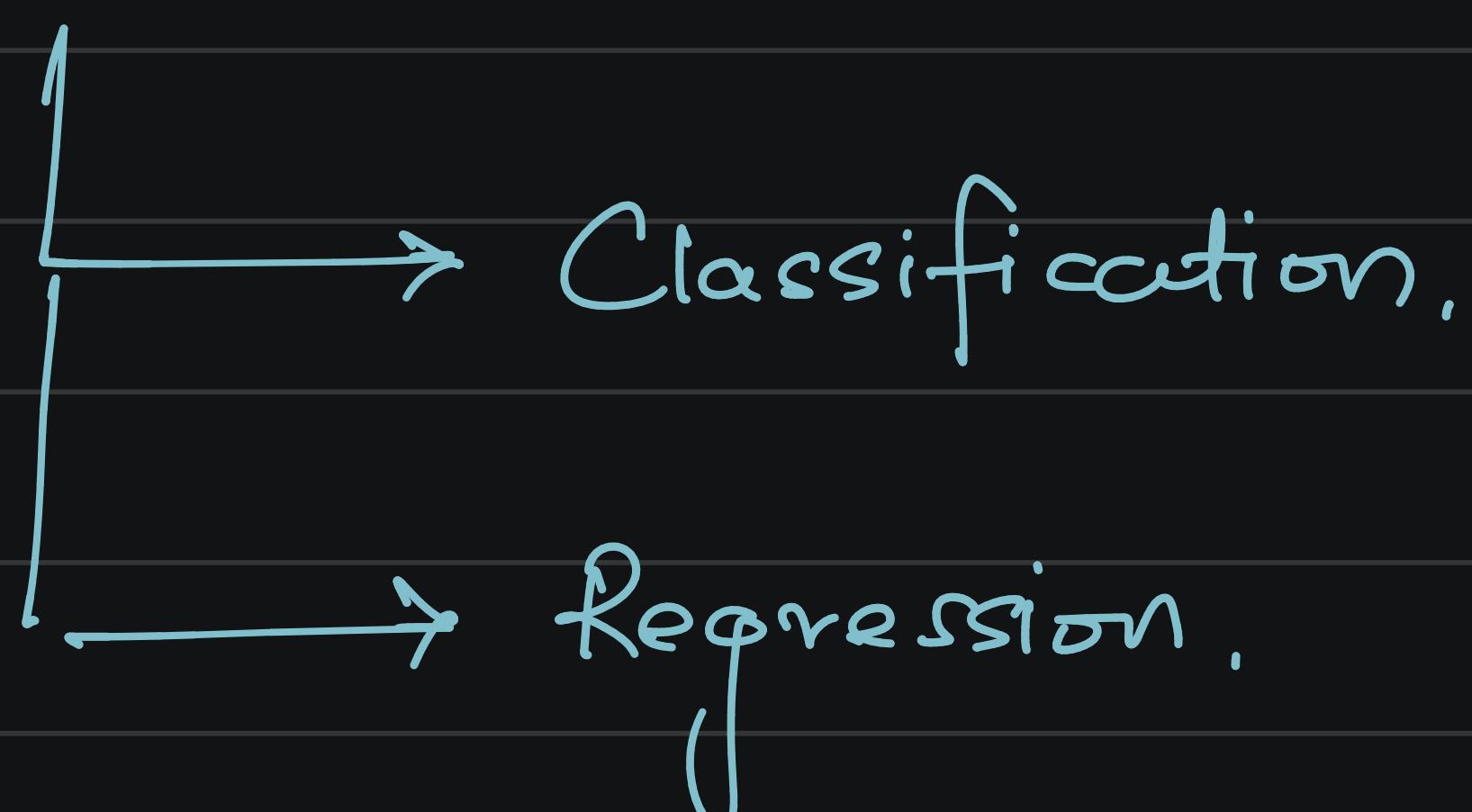


DL → Mimicing the human brain.

↓
Using multilayered Neural - Network,

Deep Learning

(1.) ANN → Artificial Neural Network.



(2.) CNN → Convolutional Neural Network.

i/p. ↗ image, video features.

Computer
vision.

e.g. RCNN, Masked RCNN, Detection,
YOLO, YOLO v8

(3.) RNN → Recurrent Neural Network.

i/p. ↳ Text data, time series data.

e.g. LSTM RNN, RNN GRU, Bidirectional

LSTM RNN, Encoder Decoder,

Transformers, BERT., Attention models

Libraries used : TensorFlow, PyTorch.

Q Why DL is becoming very popular?

(2005) → Facebook, Orkut, } Social Media
Platform.



Data generated exponentially.



(2011 → BIG DATA → Storing this data efficiently.
- 2016)

↓
Unstructured data.

↓
Structured data.

Companies. → Cloudera, Hortonworks.

2015. → What to do with the data?



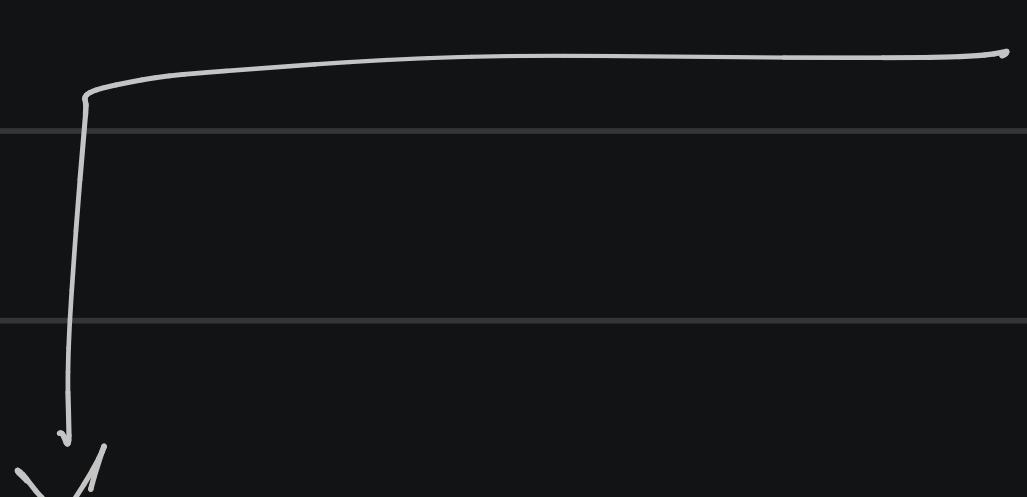
Can we use the data to
make the product better.



Data Science.



Deep Learning



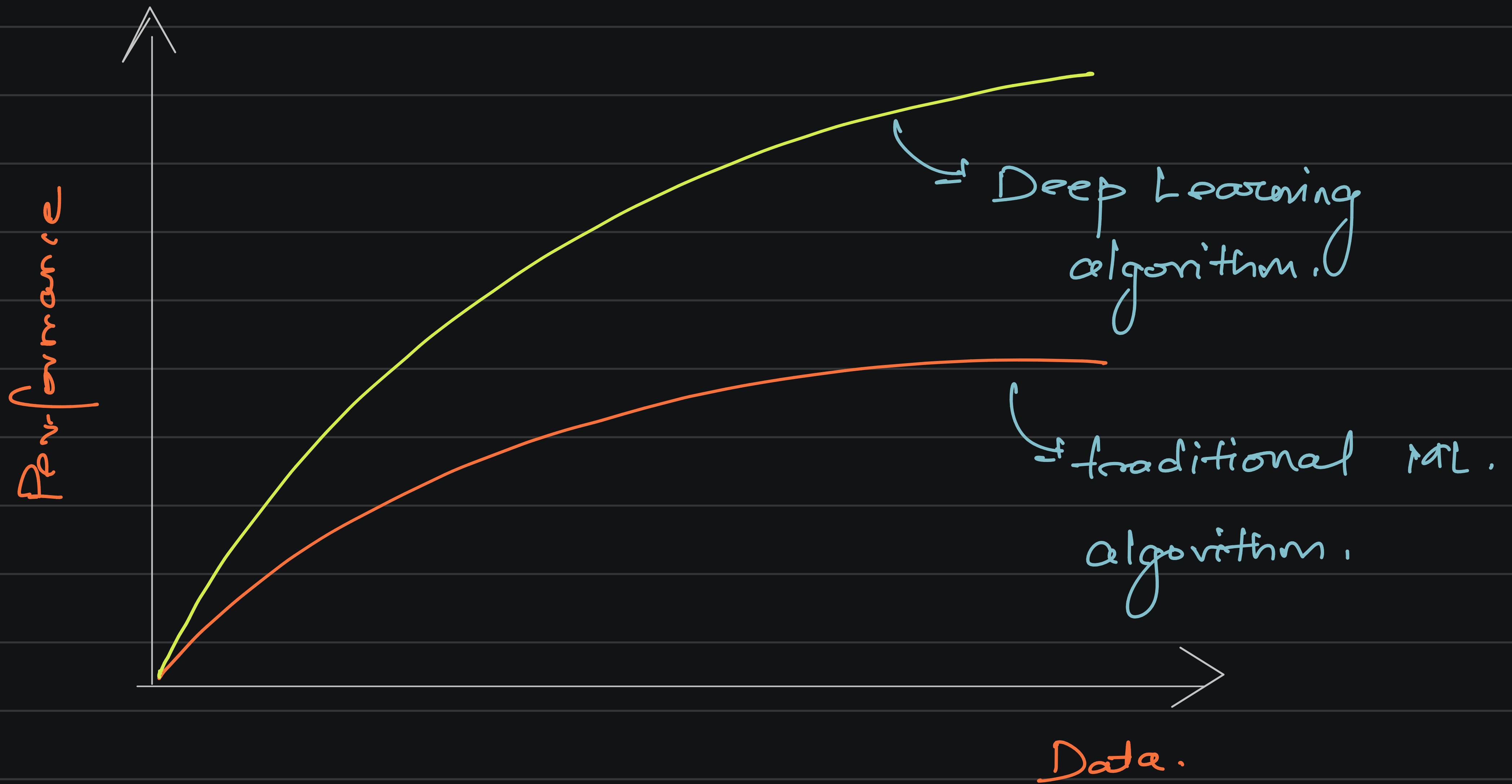
(1.) Hardware requirement \Rightarrow GPUs \Rightarrow Nvidia GPUs

* Price ↓↓↓

* Processing power ↑↑↑

(2.) Huge amount of data is generated.

→ Deep Learning model
performs well.



3. Deep learning is being used in many domains.

a. Medical → Prediction of disease
with X-ray scans.,
bone cracks, lung disease.

b. E-commerce

c. Retail.

d. Logistics.

* Perceptron. [Artificial Neuron or Neural

Network Unit]

Single layered neural network ,

Layers :

(1) I/P layer .

(2) Hidden layer .

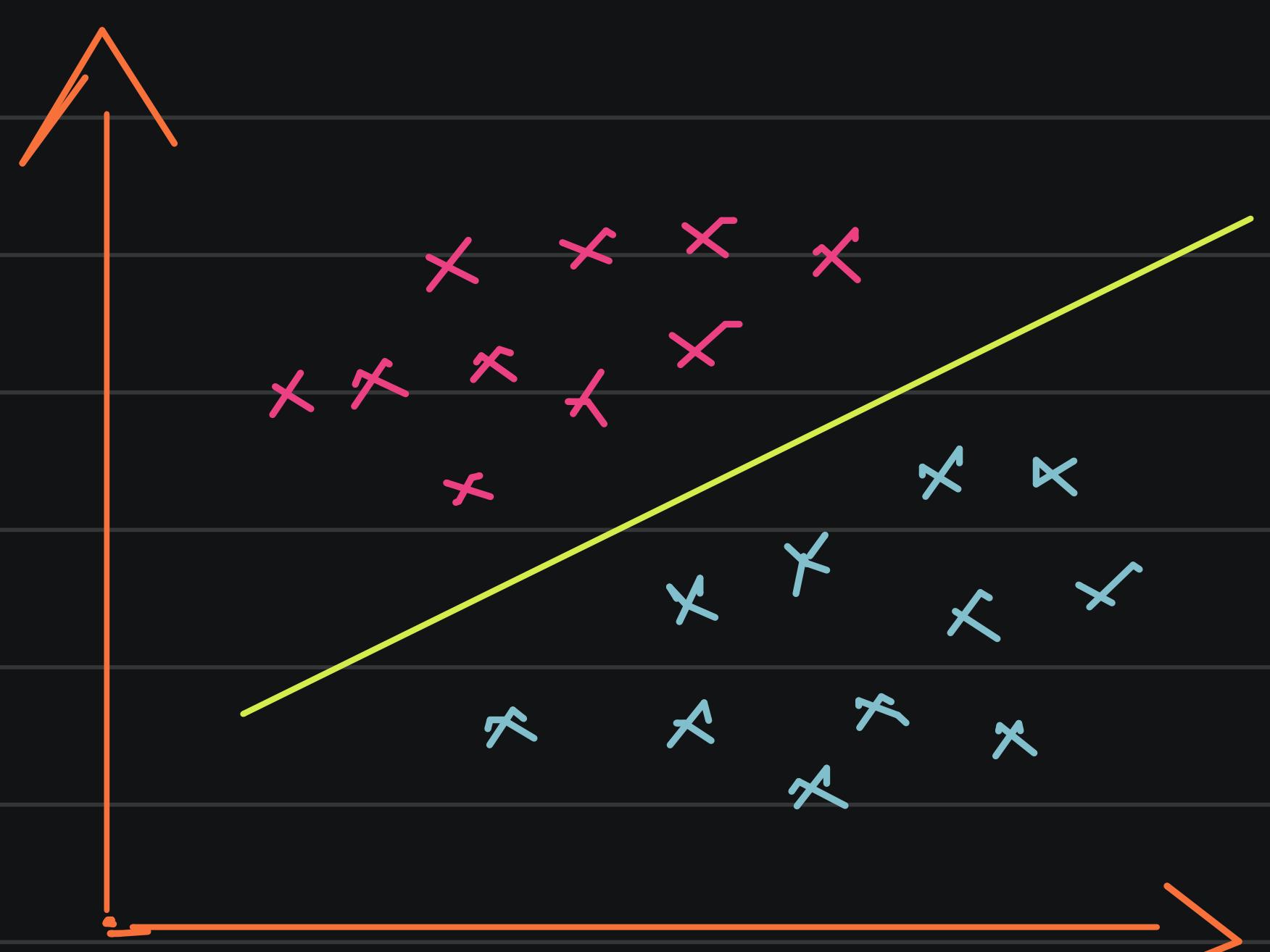
(3) Weights .

(4) Activation -function .

(5) Biases .

(6) O/P layer .

Binary classification,
(linearly separable)



DATA SET.

x_1
IQ

x_2

No. of study hrs.

y
O/P Pass/Fail.

95.

3

0

110 .

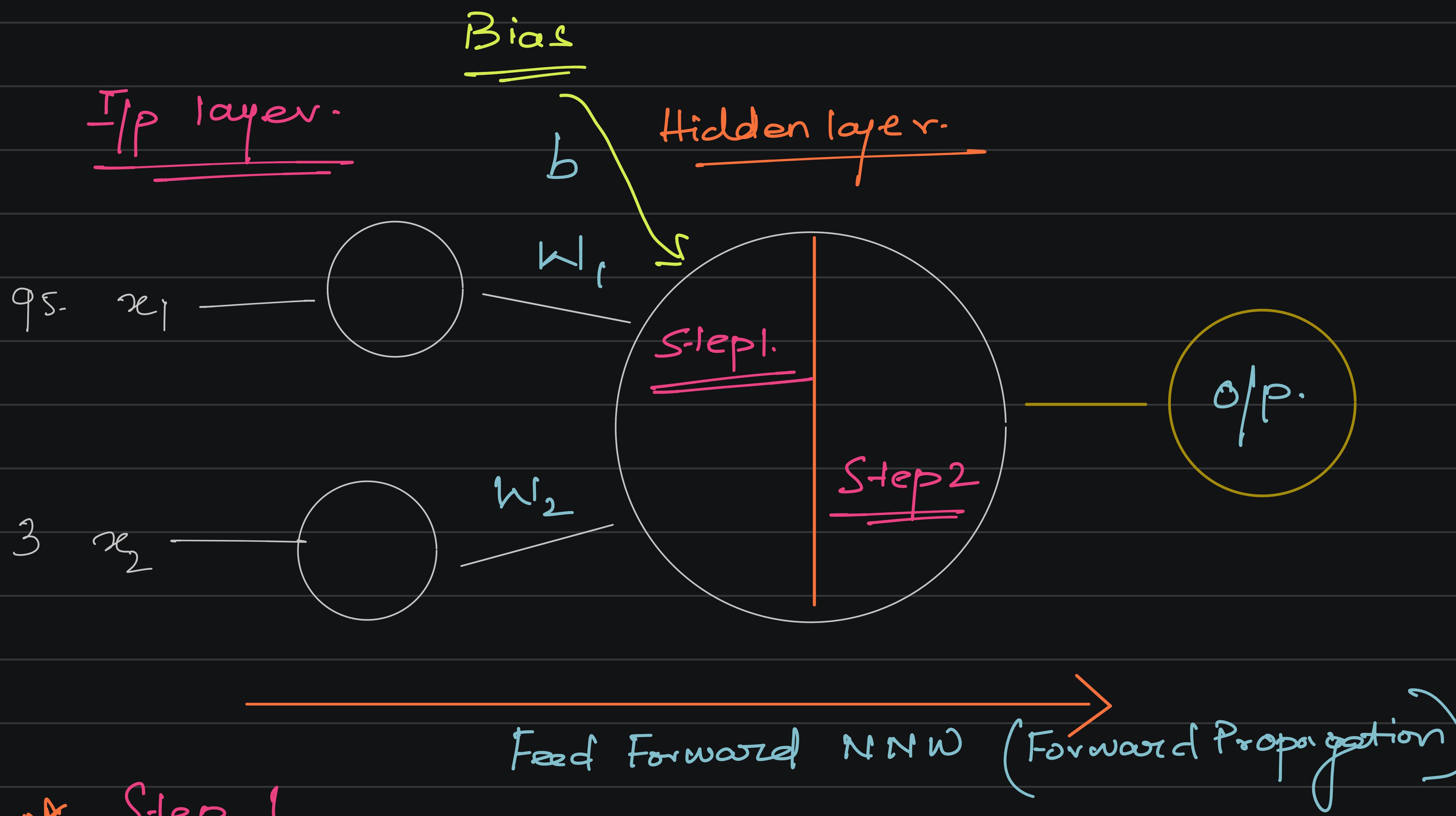
4

1 .

100

5

1 .



* Step 1.

$$Z = x_1 w_1 + x_2 w_2 + 1 \times b .$$

$$Z = \sum_{i=1}^n x_i w_i + b .$$

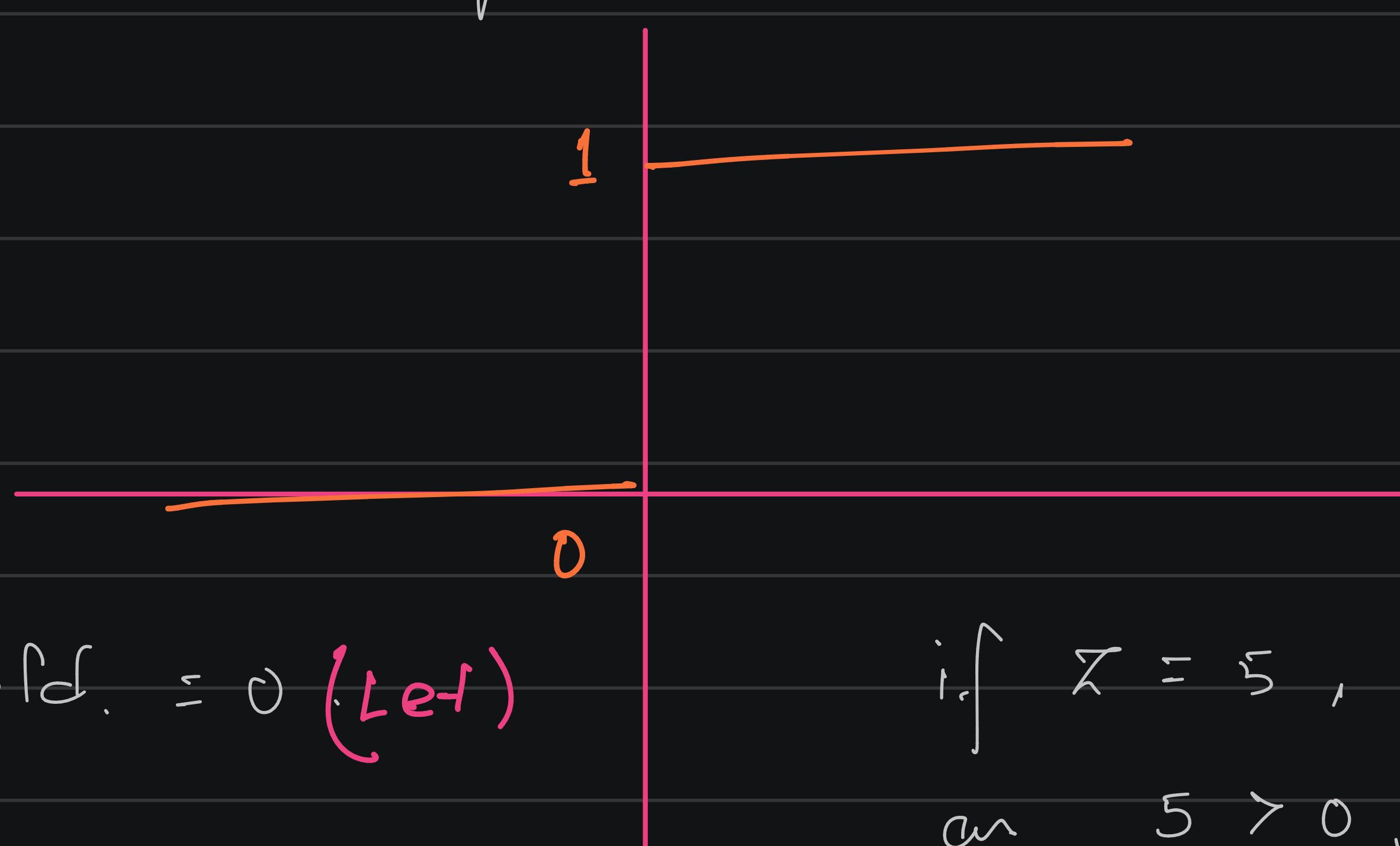
$[y = \sum mx + b]$

↓
intercept.

* Step 2

Apply an activation function, on Z-value.

eg Step function



decide a threshold. = 0 (Let)

if $Z = 5$, Step = 1.
as $5 > 0$.

Z. $\xrightarrow[\text{-function}]{\text{step}}$ -transforms into $[0 \text{ or } 1.]$
 (activation
 -function)
 ↓
 binary classifi-
 cation.

* Multi-layered Perception model. (ANN)

(1.) Forward Propagation.

(2.) Backward Propagation.

(3.) Loss function.

(4.) Activation function.

(5.) Optimizers.

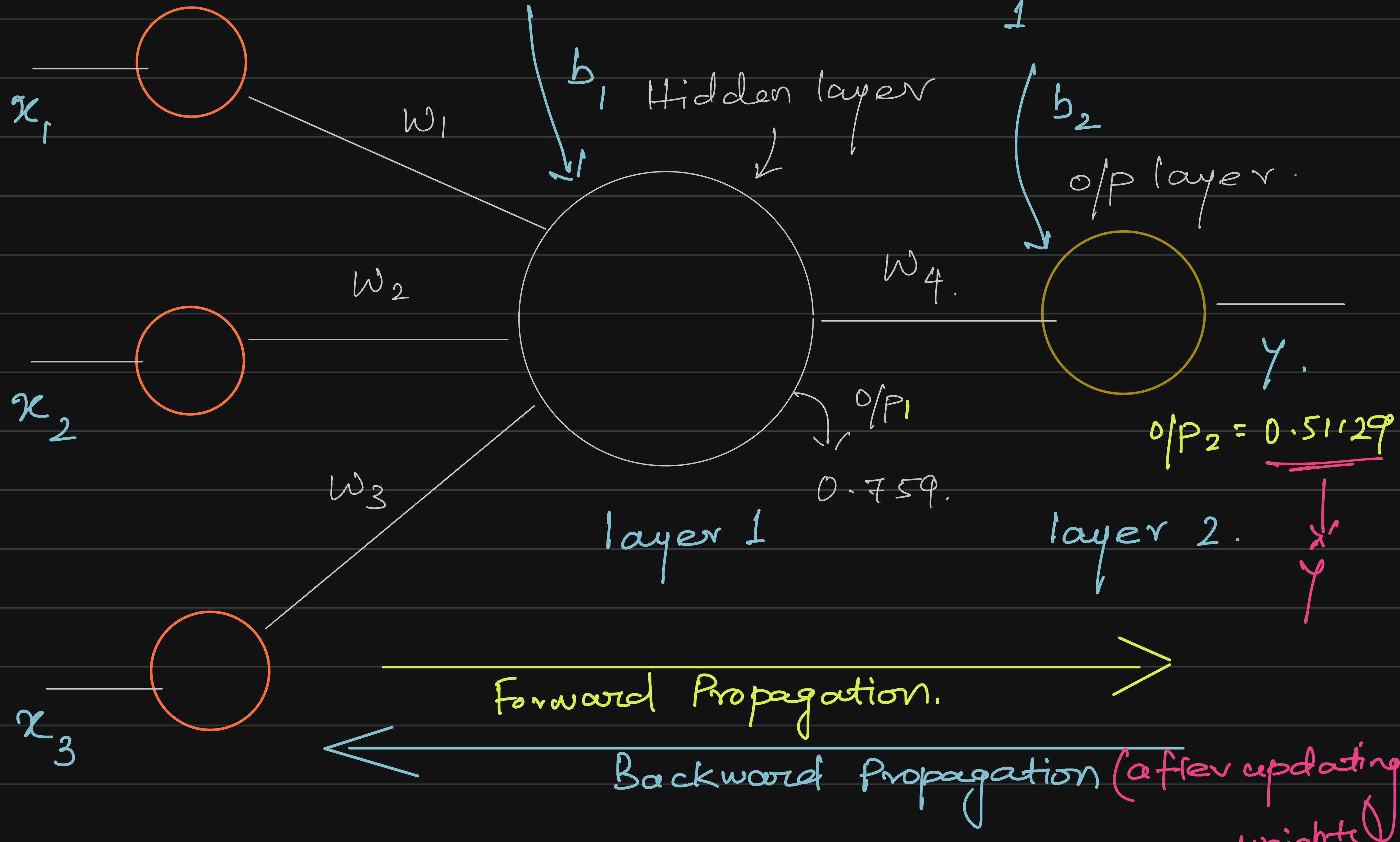
2-layered Neural N/W.

Dataset -

x_1	x_2	x_3	y	\hat{y}	Loss. ($y - \hat{y}$)
95	4	4	0/P.	1	0.51129.
100	5	2	1	1	0.49.
95	2	7	0	0	

I/P layer.

1



Forward Propagation.

Hidden layer 1.

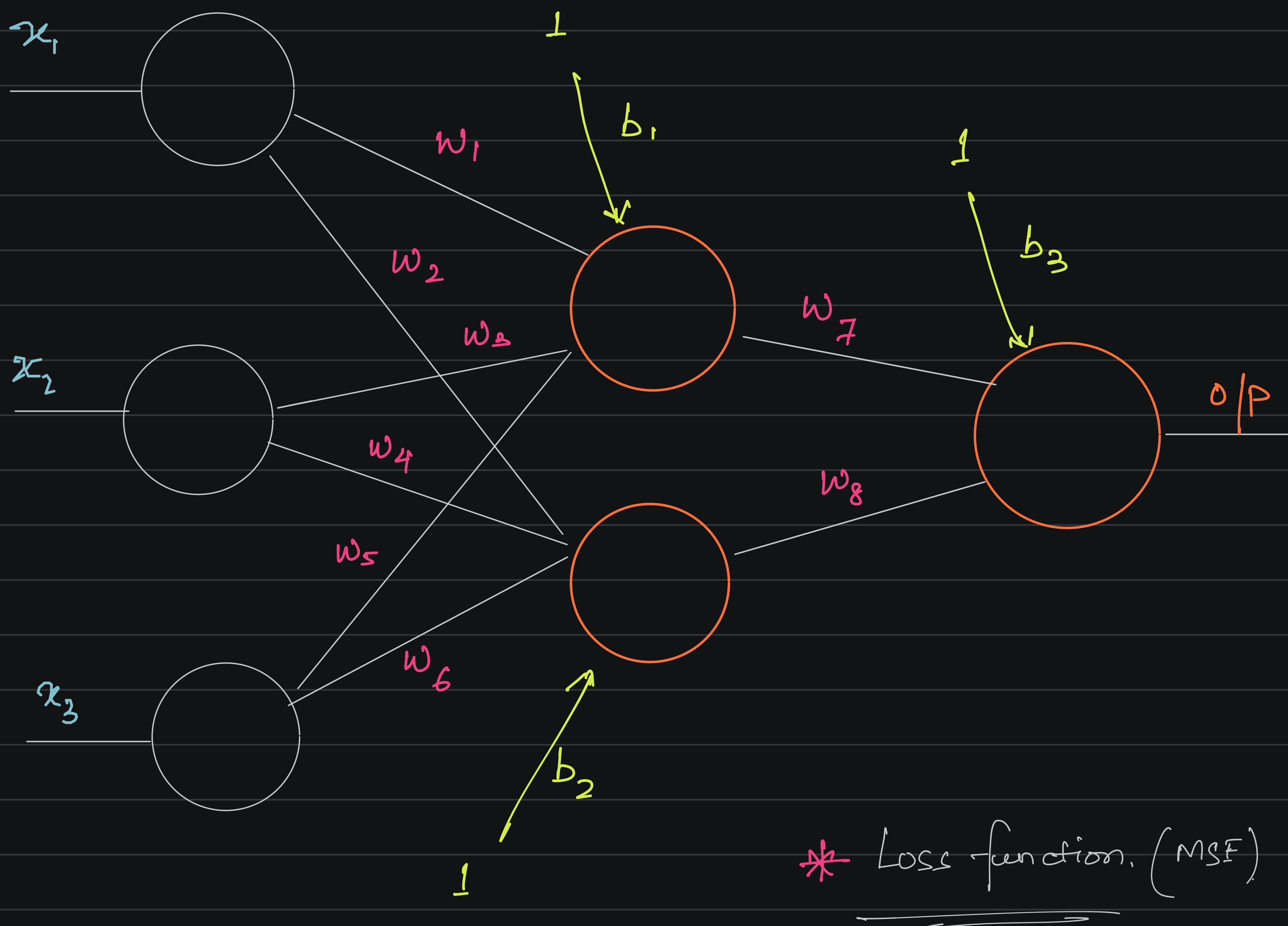
$$\begin{bmatrix} w_1 & w_2 & w_3 & b_1 \\ 0.01 & 0.02 & 0.03 & 0.01 \end{bmatrix}$$

$$\text{Step 1} = 95 * 0.01 + 4 * 0.02 + 4 * 0.03 + 0.01.$$

$$Z = 1.151.$$

Backward Propagation & weight upgradation formula

I/P layer

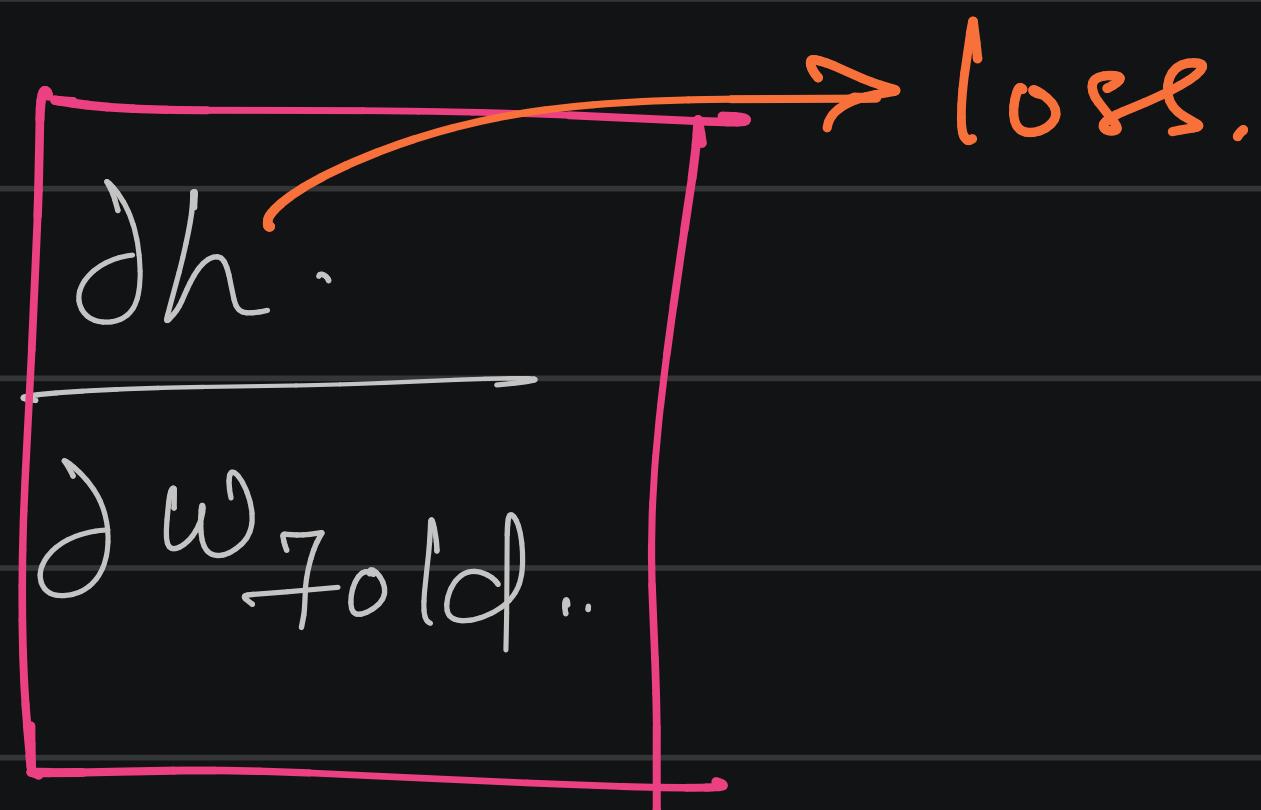


* Weight upgradation formula. :

$$(\hat{y} - y)^2 \downarrow \downarrow \downarrow$$

$$w_{7\text{ new}} = w_{7\text{ old}} - \eta$$

Learning rate



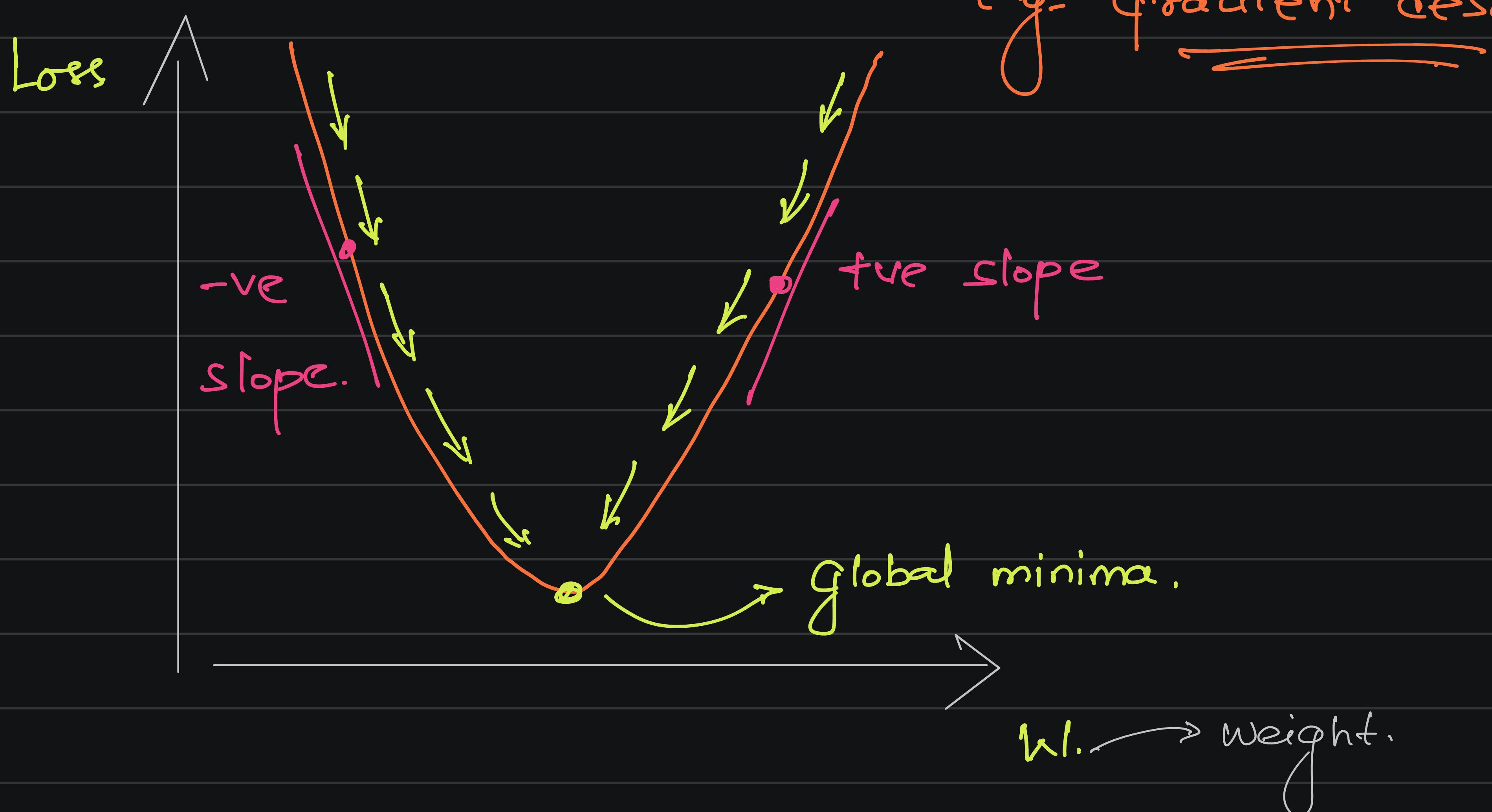
$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial h}{\partial w_{\text{old}}} \Rightarrow \text{Weight upgradation formula.}$$

Similarly,

$$b_{\text{new}} = b_{\text{old}} - \eta \frac{\partial L}{\partial b_{\text{old}}} \Rightarrow \text{Bias updation formula.}$$

* Weight updation is done using Optimizers.

e.g. Gradient descent.



* Optimizer. \rightarrow Used to reduce the loss value.

When (-)ve slope , $w_{\text{new}} = w_{\text{old}} - \eta(-)^{\text{ve}}$
 $= w_{\text{old}} + (+)^{\text{ve}}$

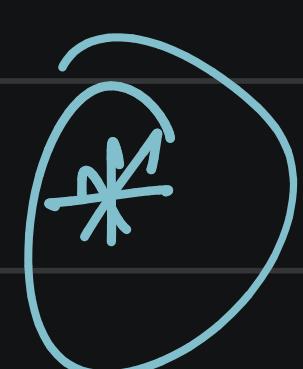
$w_{\text{new}} \gg w_{\text{old}}$, \rightarrow Moving towards global minima.

when (ϵ) re slope, $w_{\text{new}} = w_{\text{old}} - \eta (\text{slope})^{\text{rec}}$.

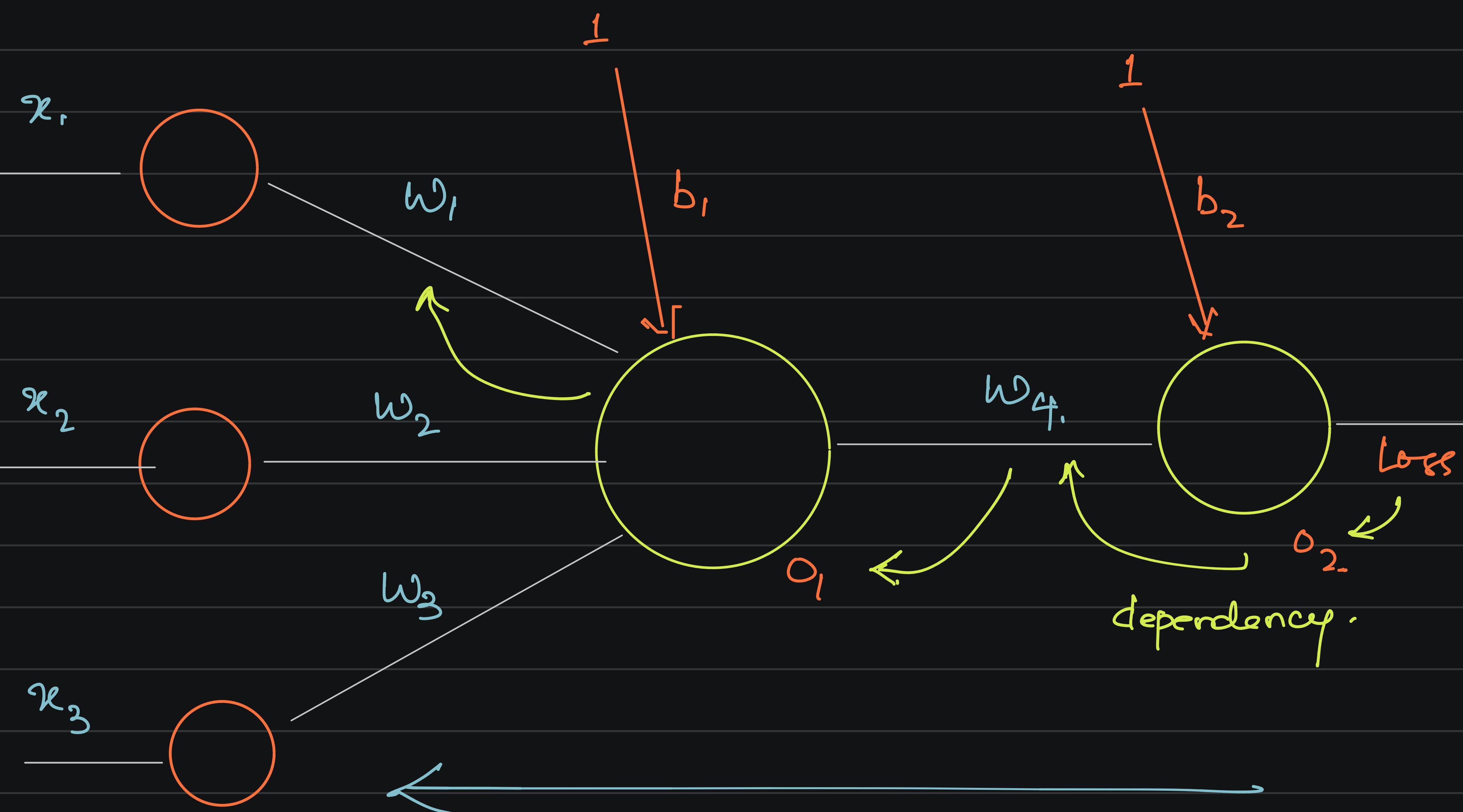
$w_{\text{new}} \ll w_{\text{old}}$.



moving toward \leftarrow global
minima.



Chain Rule of Derivative



$$w_{4 \text{ new}} = w_{4 \text{ old}} - \eta \left[\frac{\partial L}{\partial w_{\text{old}}} \right] \Rightarrow \text{slope}$$

$$\frac{\partial L}{\partial w_{4 \text{ old}}} = \frac{\partial L}{\partial o_2} * \frac{\partial o_2}{\partial w_{4 \text{ old}}} \quad \left\{ \begin{array}{l} \text{Chain} \\ \text{Rule} \\ \text{of derivative} \end{array} \right\}$$

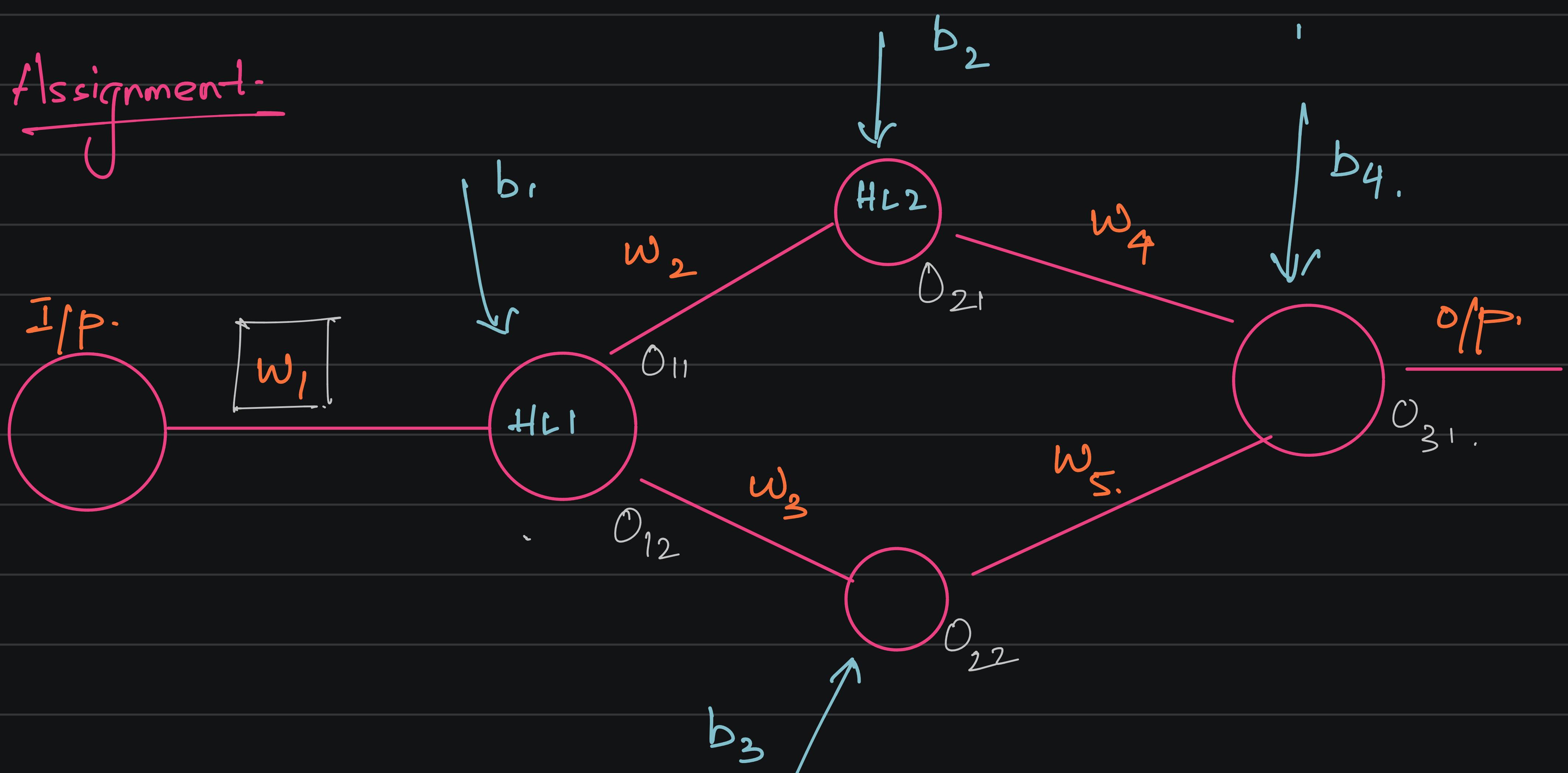
$$w_{1\text{new}} = w_{1\text{old}} - \eta \left[\frac{\partial L}{\partial w_{1\text{old}}} \right]$$

$$\frac{\partial L}{\partial w_{1\text{old}}} = \frac{\partial L}{\partial o_2} * \frac{\partial o_2}{\partial w_{4\text{old}}} * \frac{\partial w_{4\text{old}}}{\partial o_1} * \frac{\partial o_1}{\partial w_{1\text{old}}}$$

$$= \frac{\partial L}{\partial o_2} * \frac{\partial o_2}{\partial o_1} * \frac{\partial o_1}{\partial w_{1\text{old}}}$$

$w_1 \rightarrow o_1 \rightarrow w_4 \rightarrow o_2 \rightarrow h$
 dependency.

Assignment-



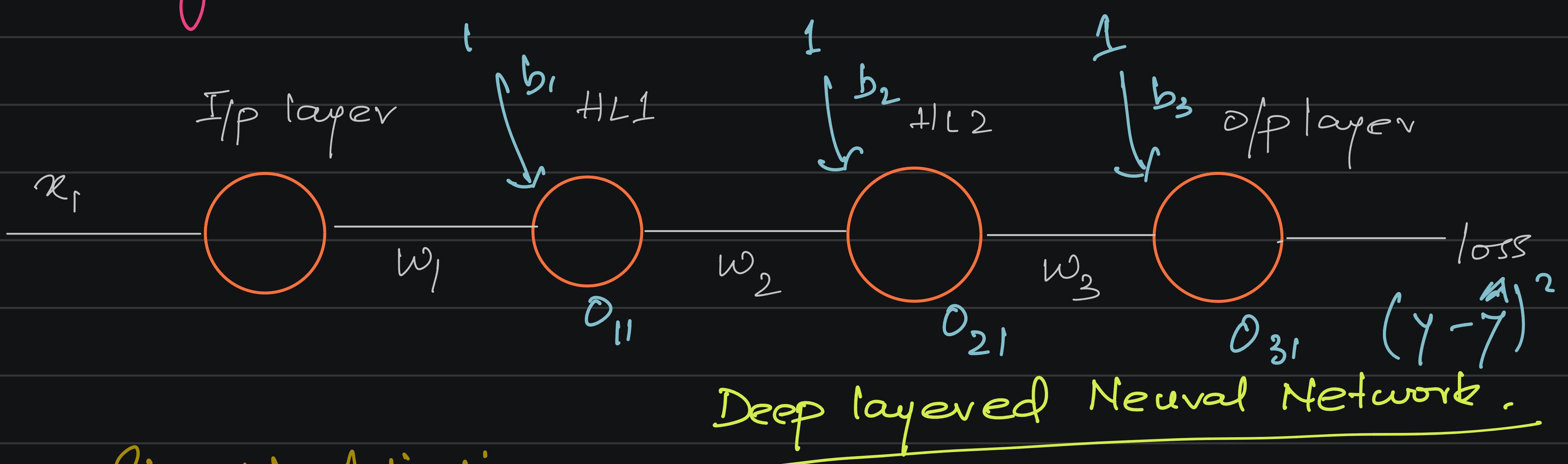
$$w_{1\text{new}} = w_{1\text{old}} - \eta \cdot \frac{\partial L}{\partial w_{1\text{old}}}$$

$$w_{\text{new}} = w_{\text{old}} - \eta \left[\frac{\partial L}{\partial o_{31}} * \frac{\partial o_{31}}{\partial o_{21}} * \frac{\partial o_{21}}{\partial o_{11}} * \frac{\partial o_{11}}{\partial w_{\text{old}}} \right] +$$

$$\left[\frac{\partial L}{\partial o_{31}} * \frac{\partial o_{31}}{\partial o_{22}} * \frac{\partial o_{22}}{\partial o_{12}} * \frac{\partial o_{12}}{\partial w_{\text{old}}} \right]$$



Vanishing Gradient Problem And Activation functions.



Sigmoid Activation.

$$\sigma(z) = \frac{1}{1 + e^{-z}} \Rightarrow [0 \text{ to } 1]$$

$$\delta_{o_1} \left[0 \leq \sigma(z) \leq 1 \right]$$

Derivative of $\sigma(z)$: $0 \leq \sigma(z) \leq 0.25$.

* In backward propagation, we take derivatives.



If we want to update weights,

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

$$\frac{\partial L}{\partial w_{i, \text{old}}} = \frac{\partial L}{\partial o_{31}} \times \begin{bmatrix} \frac{\partial o_{31}}{\partial o_{21}} \\ \frac{\partial o_{21}}{\partial o_{11}} \end{bmatrix} \times \frac{\partial o_{21}}{\partial o_{11}} \times \frac{\partial o_{11}}{\partial w_{i, \text{old}}}$$

$$o_{31} = \sigma(w_3 \times o_{21} + b_3), \quad z.$$

$$o_r, o_{31} = \sigma(z)$$

$$\text{Now, } \frac{\partial O_{31}}{\partial O_{21}} = \frac{\partial (\sigma(z))}{\partial z} \times \frac{\partial (z)}{\partial O_{21}}$$

Chain rule of derivatives was applied.

Small value.

$$\frac{\partial}{\partial} \frac{\partial (\sigma(w_3 \times O_{21} + b_3))}{\partial O_{21}} = w_3.$$

So,

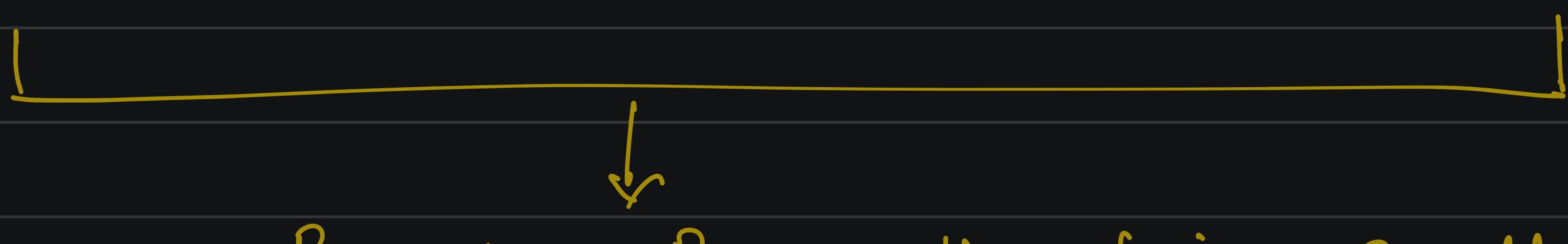
$$\frac{\partial L}{\partial w_{\text{old}}} = \frac{\partial L}{\partial O_{31}} \times \frac{\partial O_{31}}{\partial O_{21}} \times \frac{\partial O_{21}}{\partial O_{11}} \times \frac{\partial O_{11}}{\partial w_{\text{old}}}.$$

Small value

Small value

Small value

Small value



Result : Bigger the chain, smaller is the derivative of loss.
value.

Then.

$$w_{1\text{new}} = w_{1\text{old}} - \left[\eta \frac{\frac{\partial L}{\partial w_{1\text{old}}}}{\frac{\partial^2 L}{\partial w_{1\text{old}}^2}} \right]$$

↓ small.

So, $\underline{w_{\text{new}}} \approx \underline{w_{\text{old}}}$.



Resulting in gradient descent to not work properly.

So, in deep layered neural network., sigmoid activation function will not work as its derivative ranges between $0 - 0.25$.

To fix this, use other activation functions:

1. Tanh.

2. Relu.

3. Prelu.

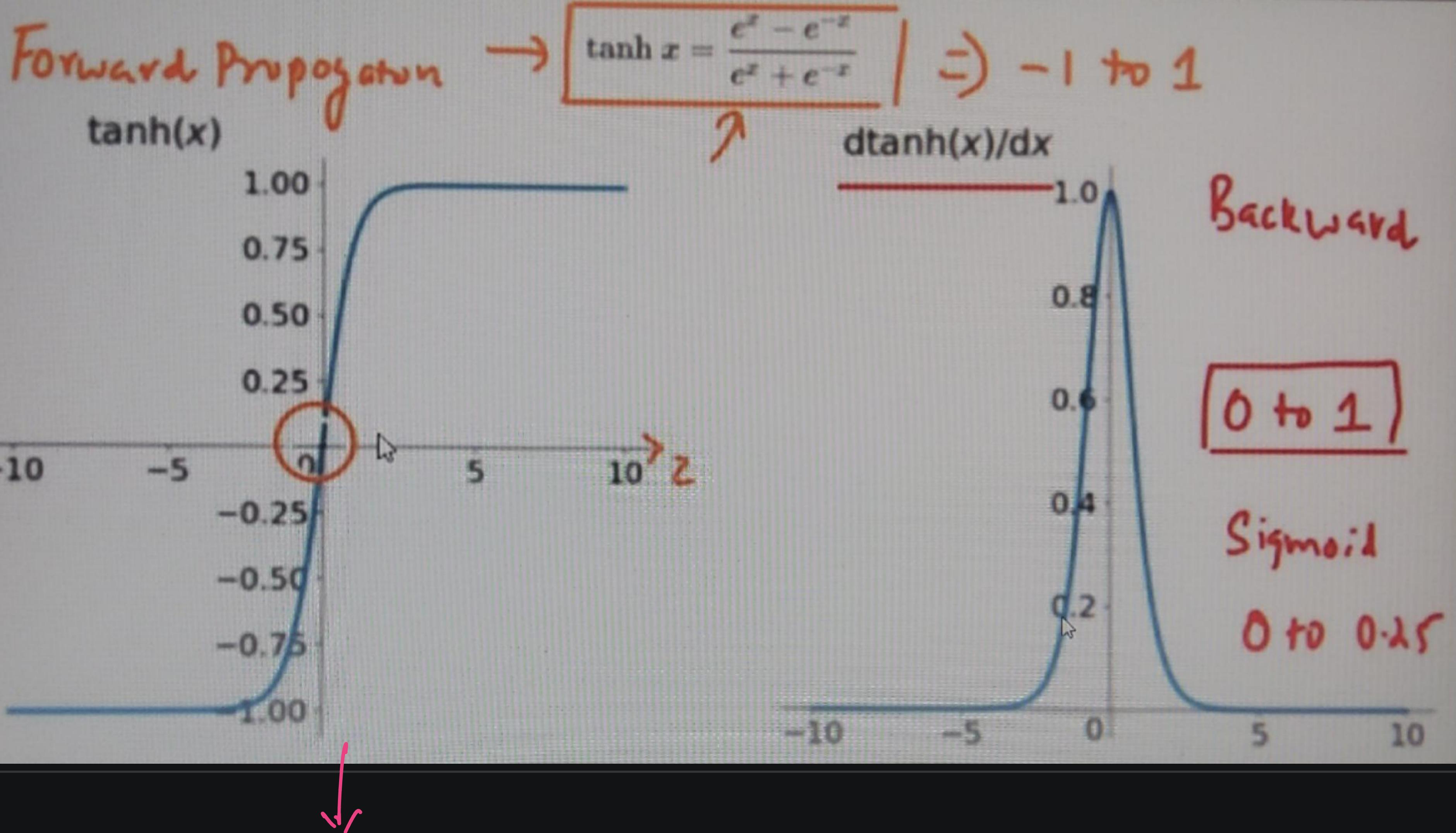
4. Elu.

5. Swiss.

* Zero centered. → Smoothes the weight updation



2. Tanh Activation Function.



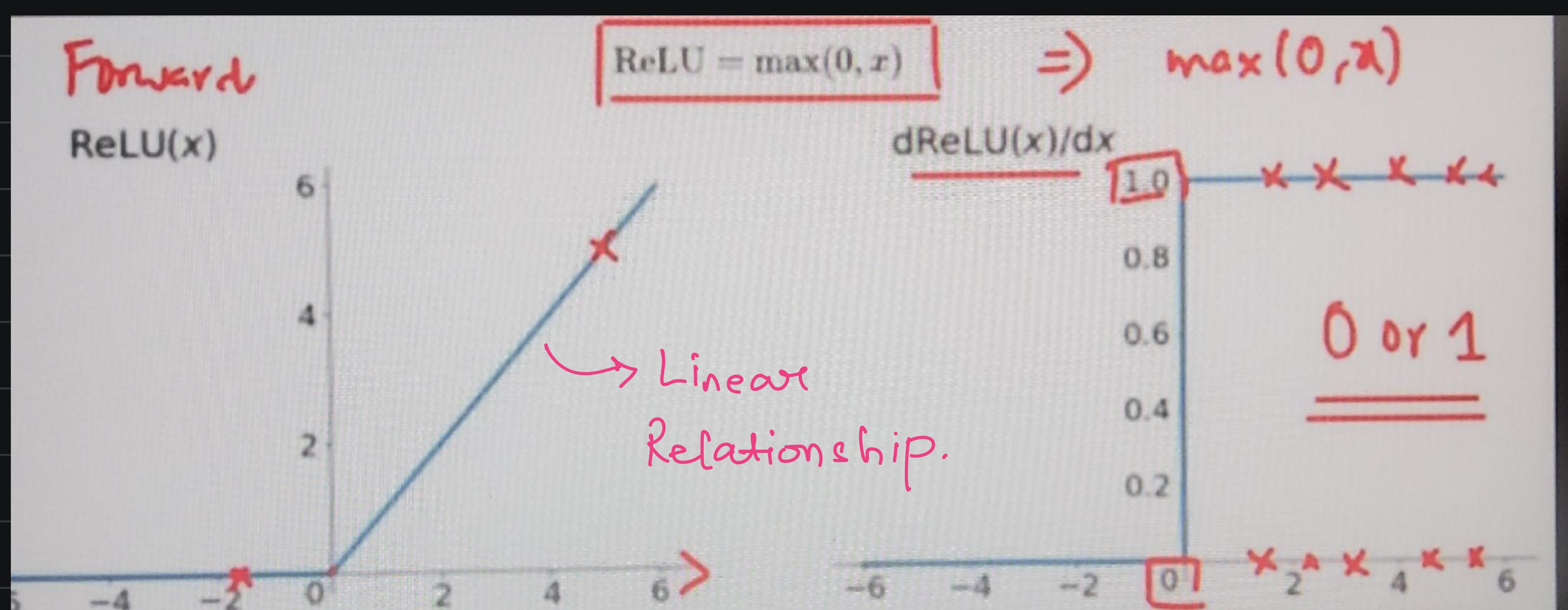
Zero centered.

Advantages:

1. Zero centered \rightarrow Efficient weight updation.

Disadvantages:

1. Time complexity is more.
2. Possibility of vanishing gradient problem for a very deep neural network.
3. ReLU Activation Function.



Advantages.

- (1.) Solves vanishing gradient problem.
- (2.) Time complexity — Fast.

Disadvantage.

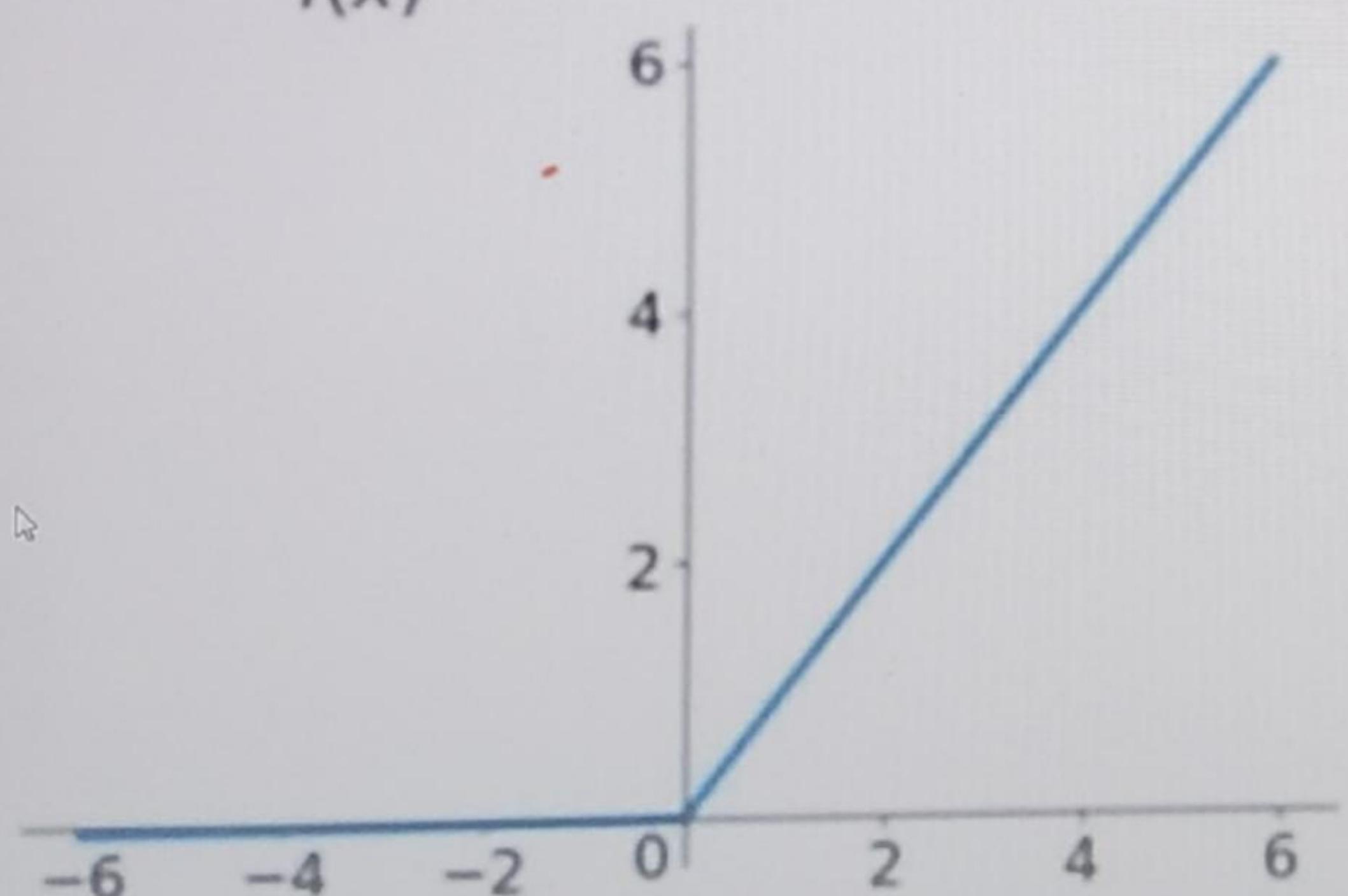
- (1.) If the derivative is 0, the neuron becomes a dead neuron.

- (2.) Not zero centric.

(4.) Leaky Relu or Parametric Relu. Activation Function.

Leaky Relu.

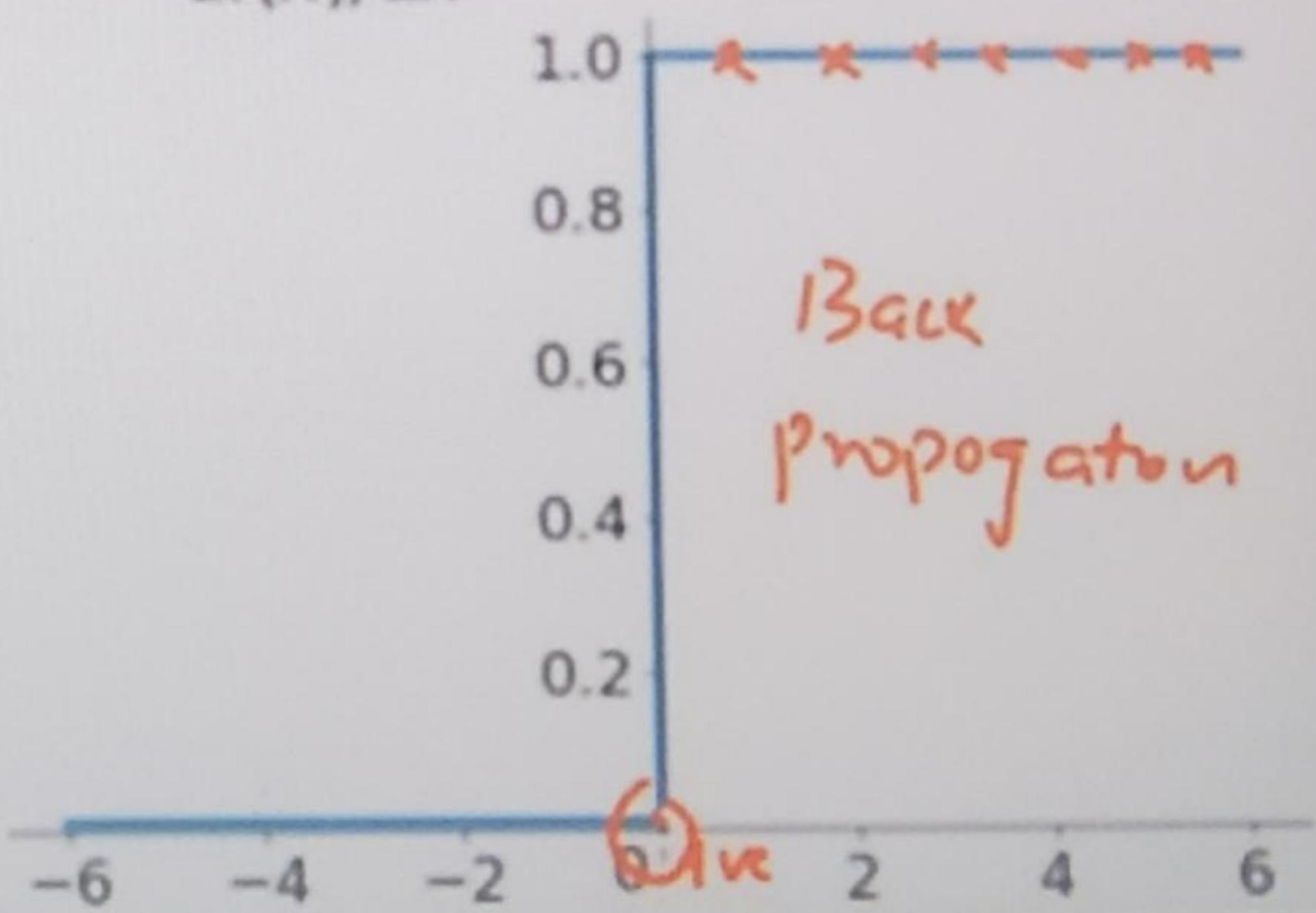
Forward
 $f(x)$



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

$$df(x)/dx$$

$$\max(0, x)$$



Parametric ReLU

hyperparameter.

$$\underline{f(x) = \max(\alpha x, x)}.$$

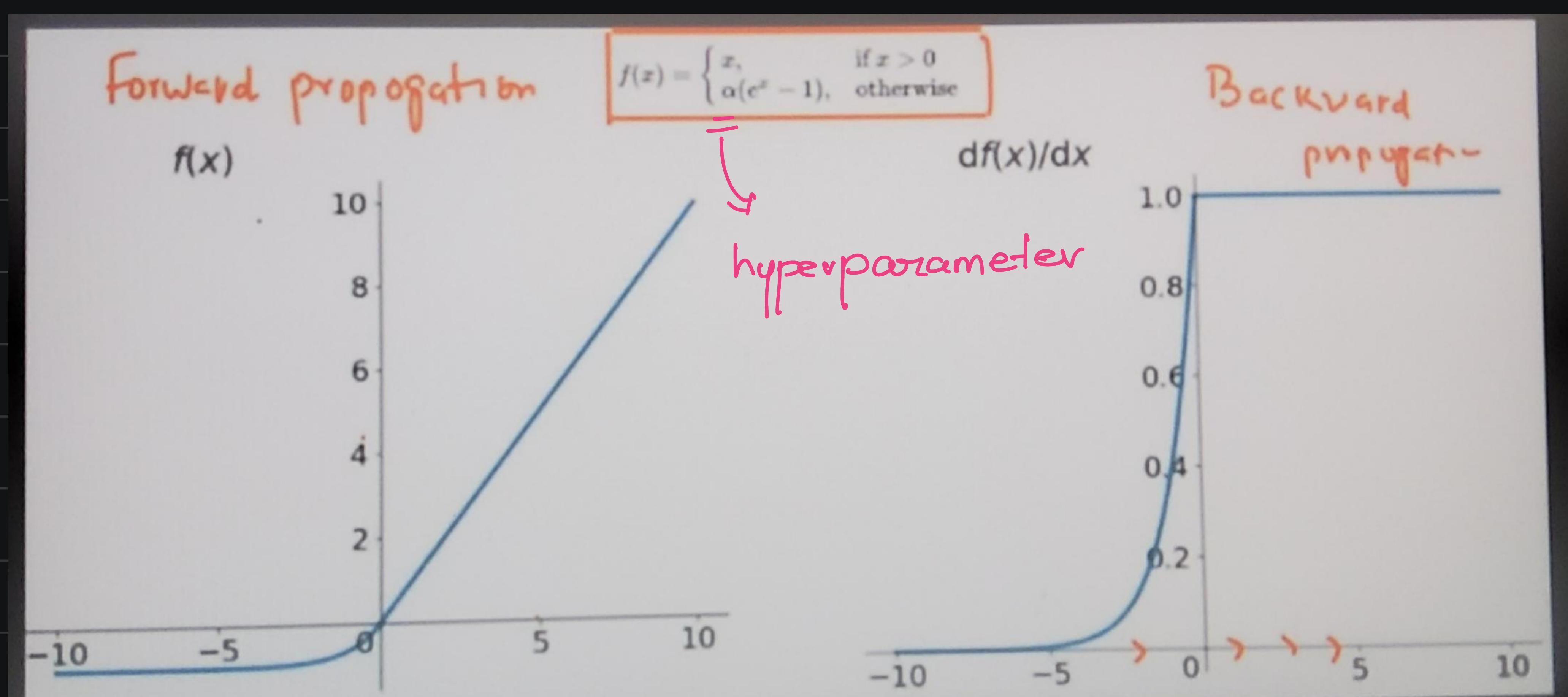
Advantage.

- (1.) Prevents the dead neuron. due to the small value. of $-f(x)$. Derivative never becomes 0.

Disadvantage.

- (2.) Not zero centric \rightarrow weight updation not efficient.

5. ELU (Exponential Linear Unit)



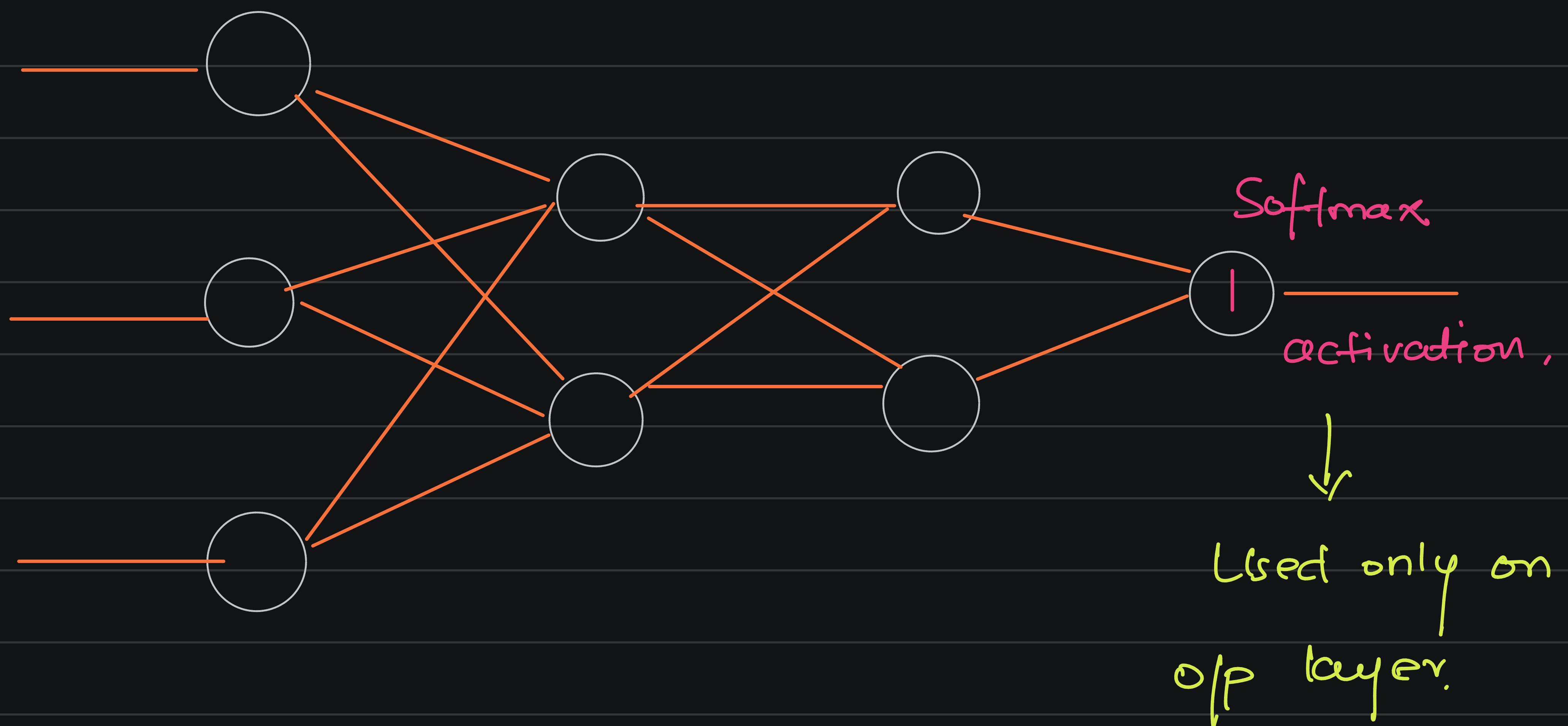
Advantage.

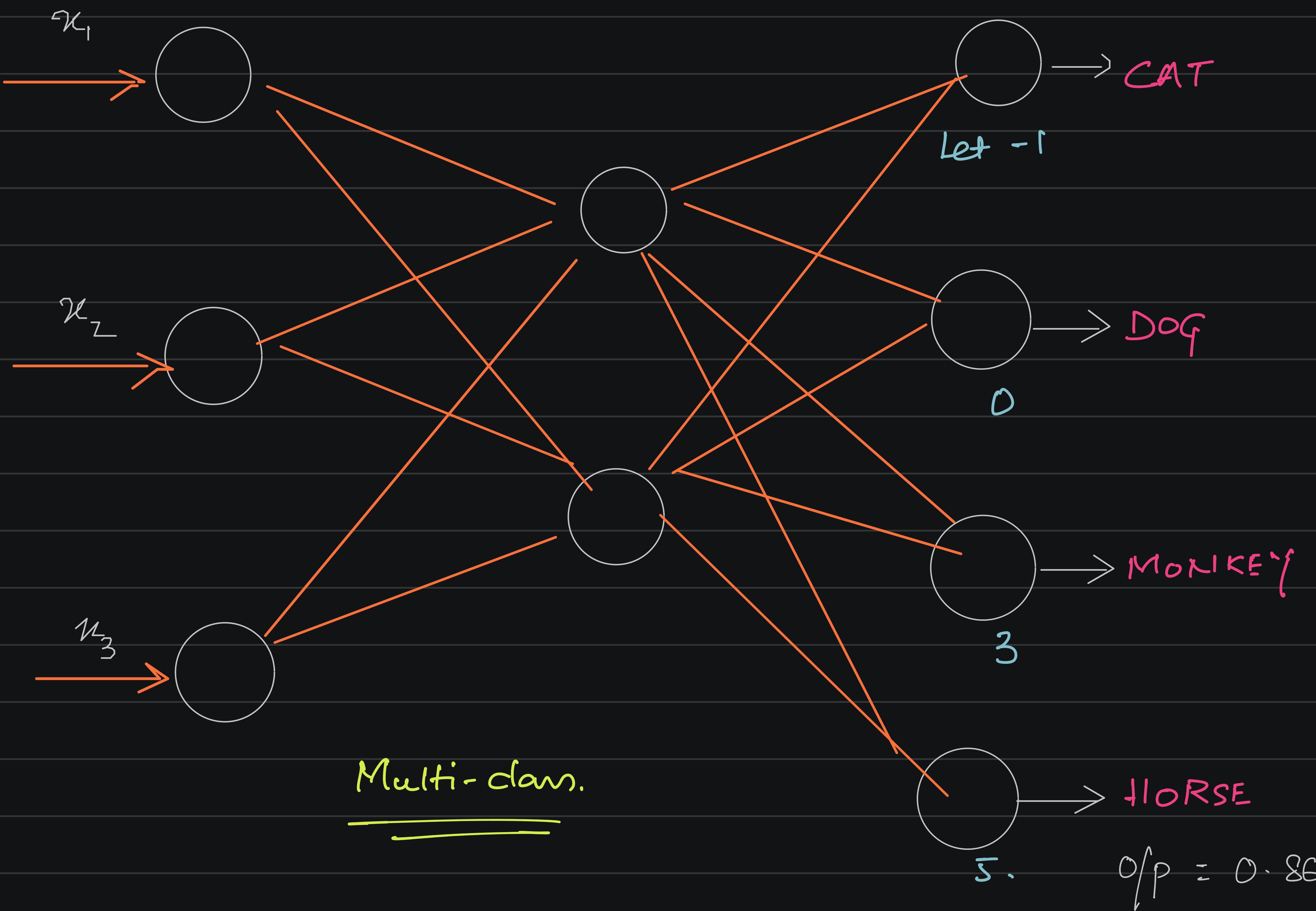
- (1.) Zero centric .
- (2.) Prevents dead neurons .
- (3.) Linear relationship .

Disadvantage

- (1.) Time complexity → high. (because of formula).

(6.) Softmax Activation Function. (Used for multi-class classification).





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

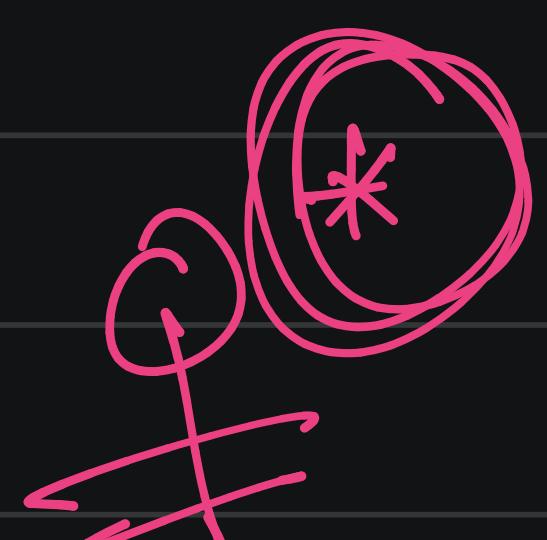
$$f_{x_1} \text{ Cat } = \frac{e^{-1}}{(e^{-1} + e^0 + e^3 + e^5)} = 0.00033$$

$$\text{Dog} = \frac{e^0}{(e^{-1} + e^0 + e^3 + e^5)} = 0.0024.$$

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

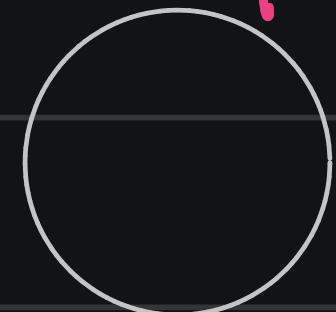
$$\text{Horse} = \frac{e^5}{(e^{-1} + e^0 + e^3 + e^5)} = 0.1353.$$

$$\Pr(\text{Horse}) = \frac{0.1353}{0.0024 + 0.0024 + 0.0183 + 0.1353} \\ = 0.86 \approx 86\%.$$

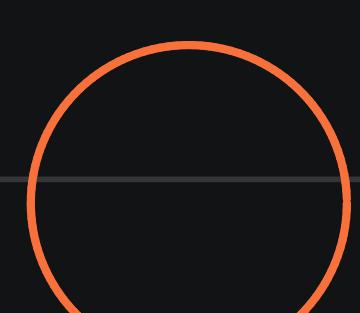
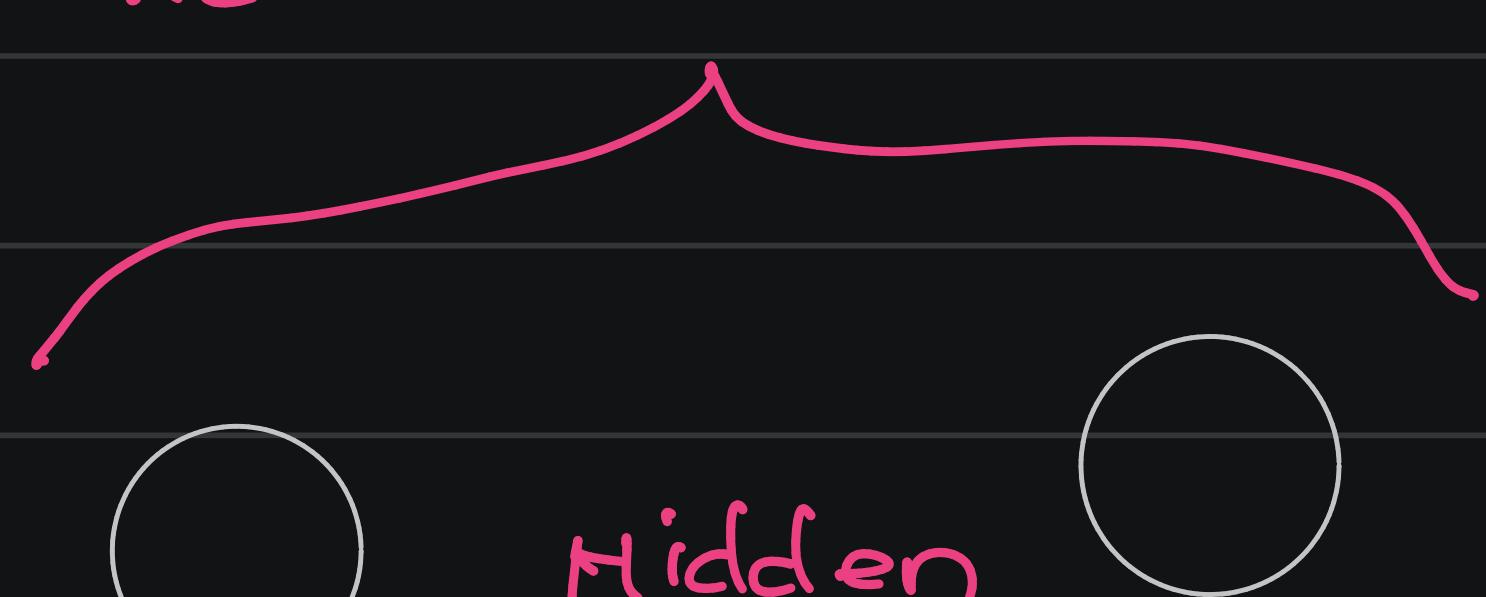


Which activation function to use when.

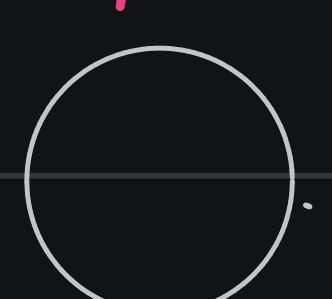
i/p layer.



ReLU and variants.

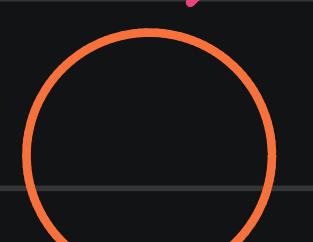


1 o/p. \rightarrow Sigmoid.



layer.

multiple
o/p.



Softmax

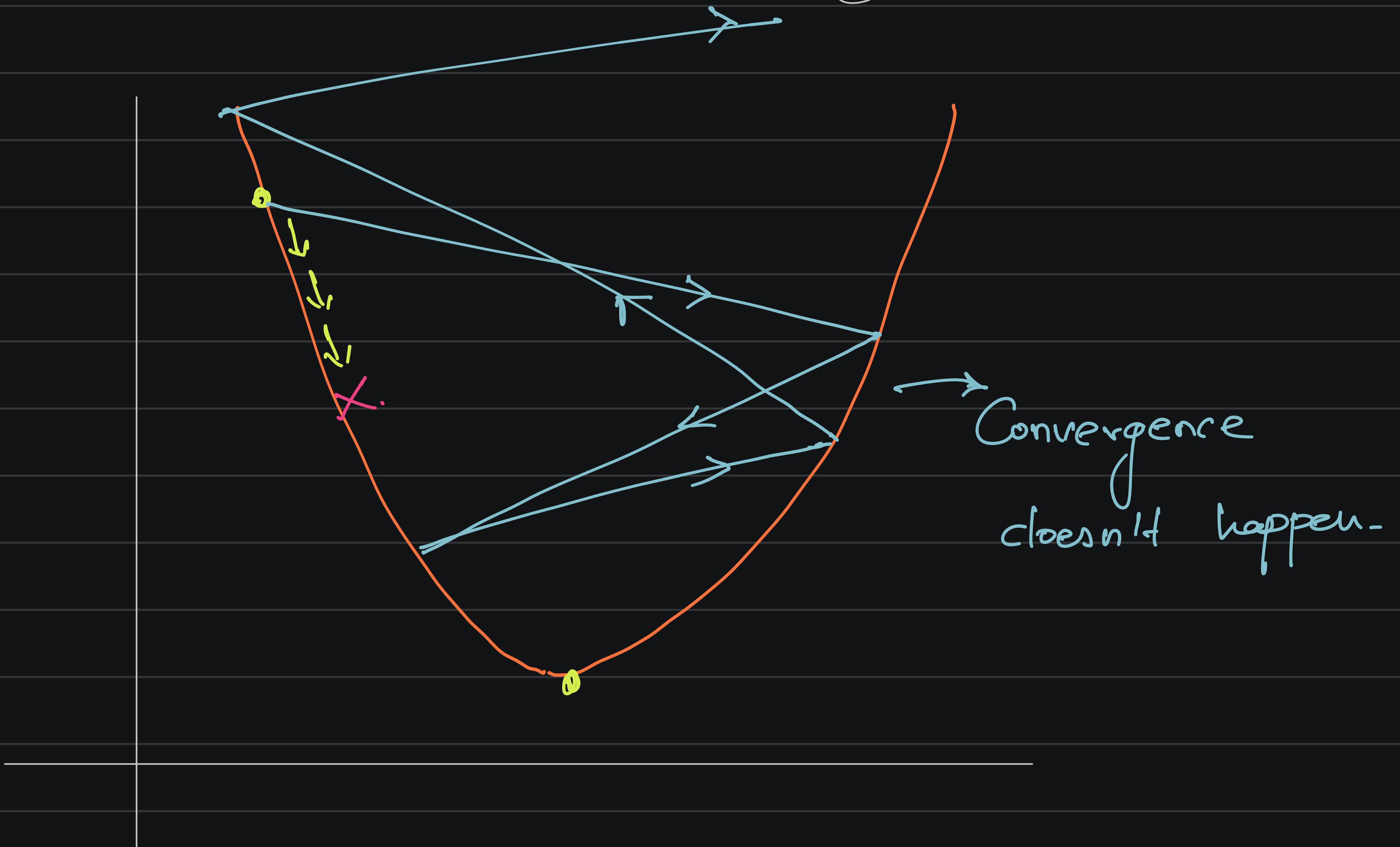
δ_0 , Binary classification in op \rightarrow Sigmoid.

Multiclass classification in op \rightarrow Softmax.

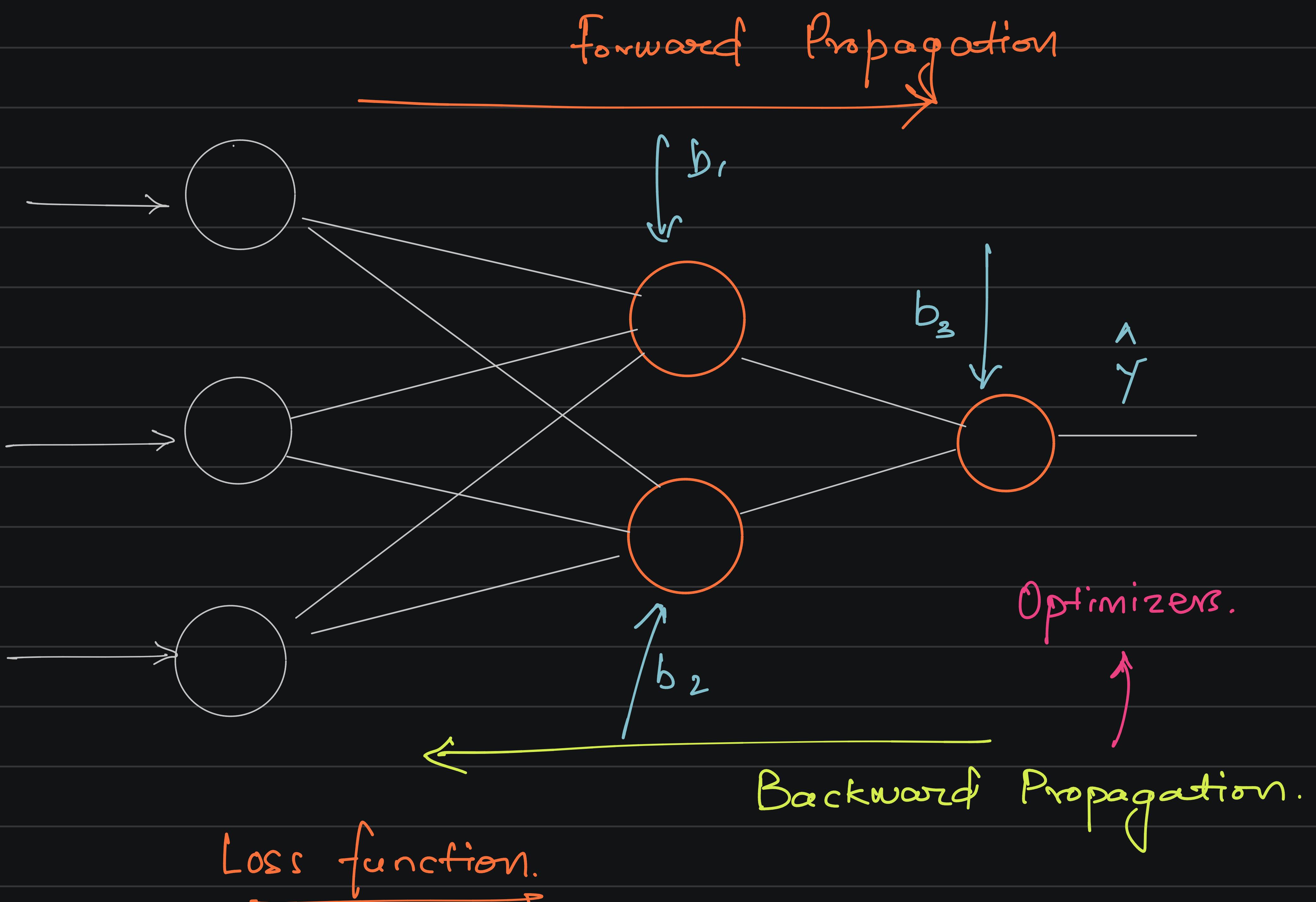
* Regression Problem \rightarrow Linear activation function.

* Exploding Gradient Problem. \rightarrow When does this happen?

e.g. weight initialization \rightarrow higher weights.



1. Loss Function & Cost Function.

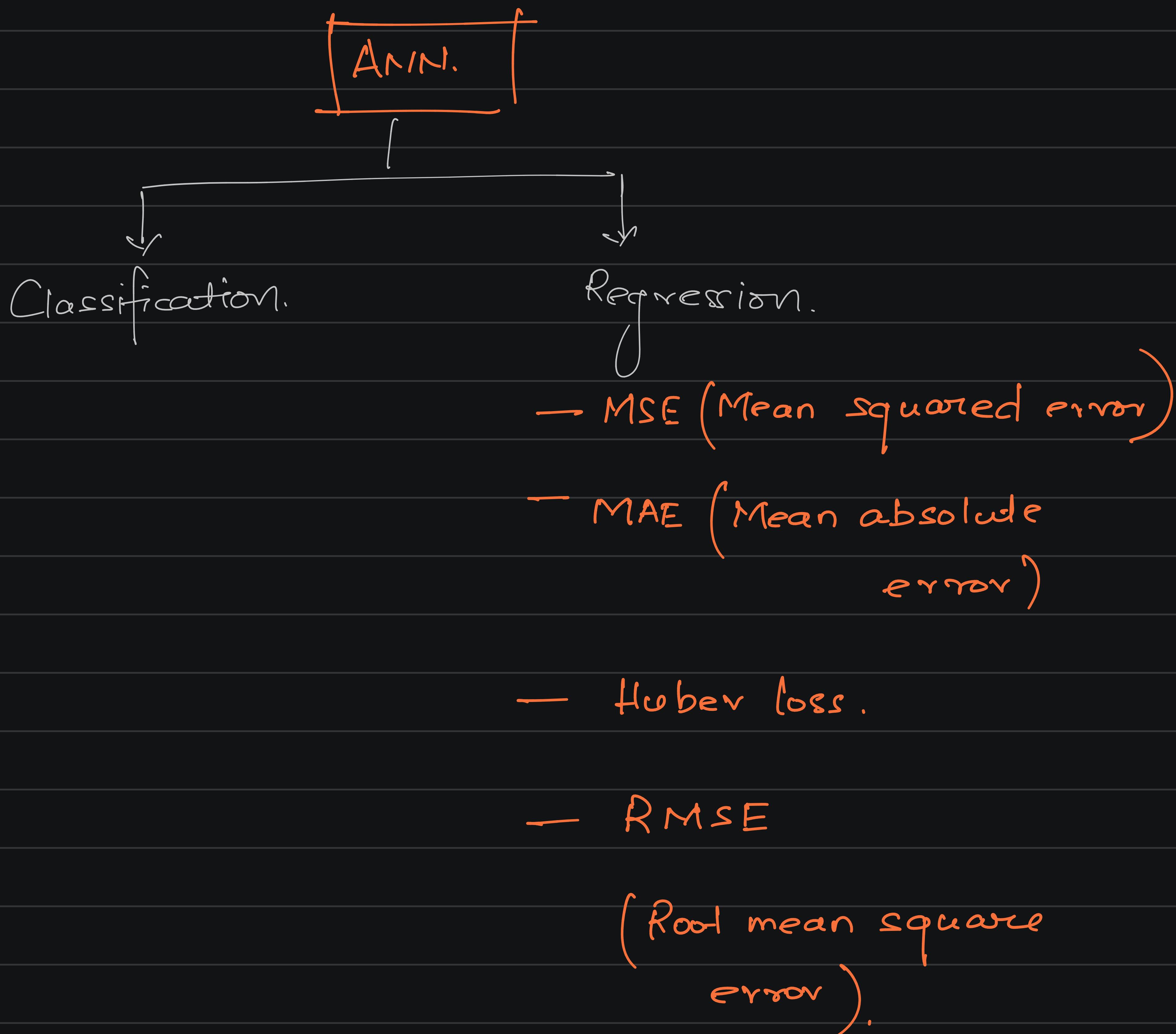


Loss function.

$$= (\gamma - \hat{\gamma})^2 \quad \xrightarrow{\text{for single datapoint.}}$$

Cost function.

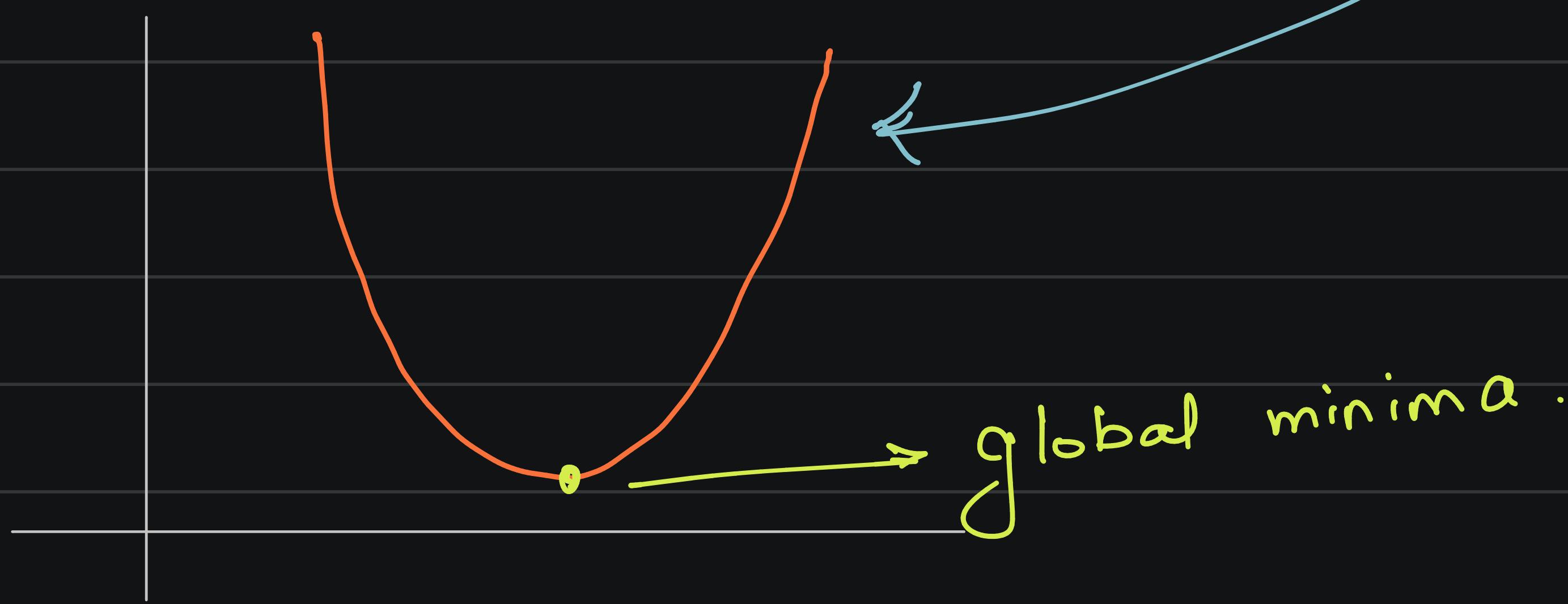
$$= \sum_{i=1}^n (\gamma_i - \hat{\gamma}_i)^2 \quad \xrightarrow{\text{for batch datapoints.}} \\ \text{(dataset).}$$



(1.) Mean squared error. (MSE).

$$\text{Loss fn} = (y - \hat{y})^2 \quad \xrightarrow{\text{Quadratic equation.}}$$

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



Advantage

Disadvantages.

(1) MSE is differentiable.

(1) Not robust to outliers.

(2) It has 1 local or 1 global minima.

(3) It converges faster.

2. Mean Absolute Error.

$$\text{Loss fn} = |y - \hat{y}|$$

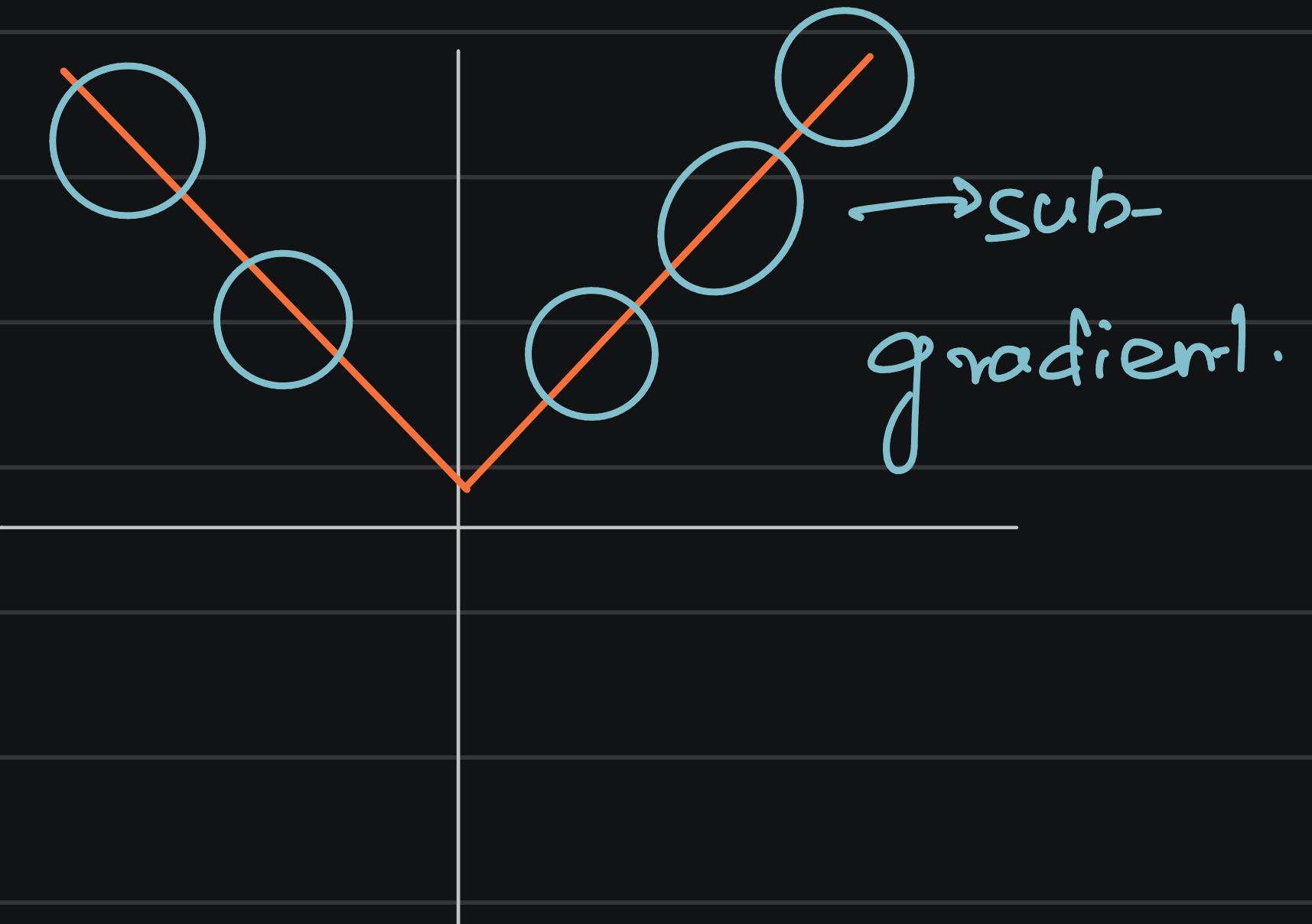
$$\text{Cost fn} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Advantage

Disadvantages.

(1) Robust to outliers.

(1) Convergence is slow.



3. Huber loss.

Combination of MSE & MAE.

hyperpara-
malev

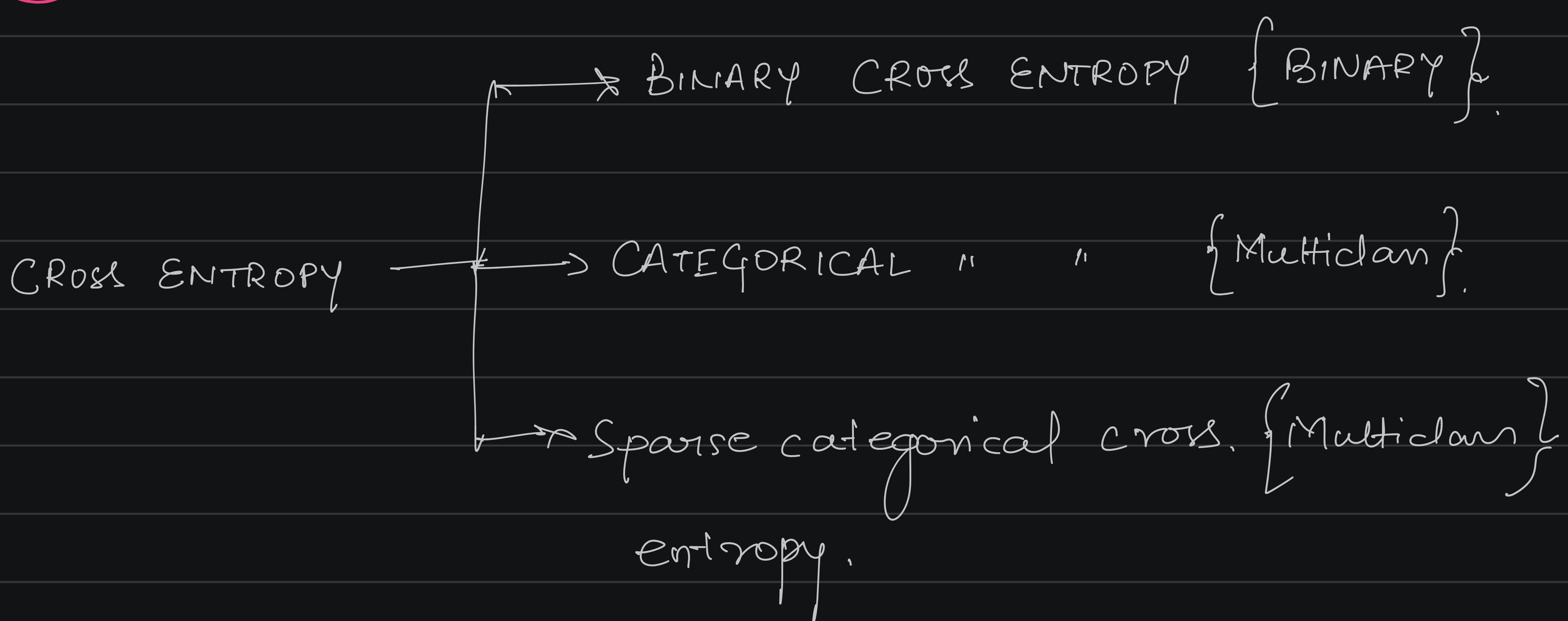
$$\text{Cost} = \begin{cases} \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - \hat{\hat{y}}_i)^2, & \text{if } |\hat{y} - \hat{\hat{y}}| \leq \delta \\ \delta |\hat{y} - \hat{\hat{y}}| - \frac{1}{2} \delta^2, & \text{otherwise.} \end{cases}$$

\Rightarrow MSE.
 \Downarrow
MAE

Advantage:

- (1) Takes care of both the issues for MSE and MAE.

* Loss & Cost function for classification problems.



1. Binary Cross Entropy.

$\Rightarrow \log$ loss function.

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

$$\text{Or Loss.} = \begin{cases} -\log(1-\hat{y}) & \text{if } \hat{y} = 0, \\ -\log(\hat{y}) & \text{if } \hat{y} = 1. \end{cases}$$

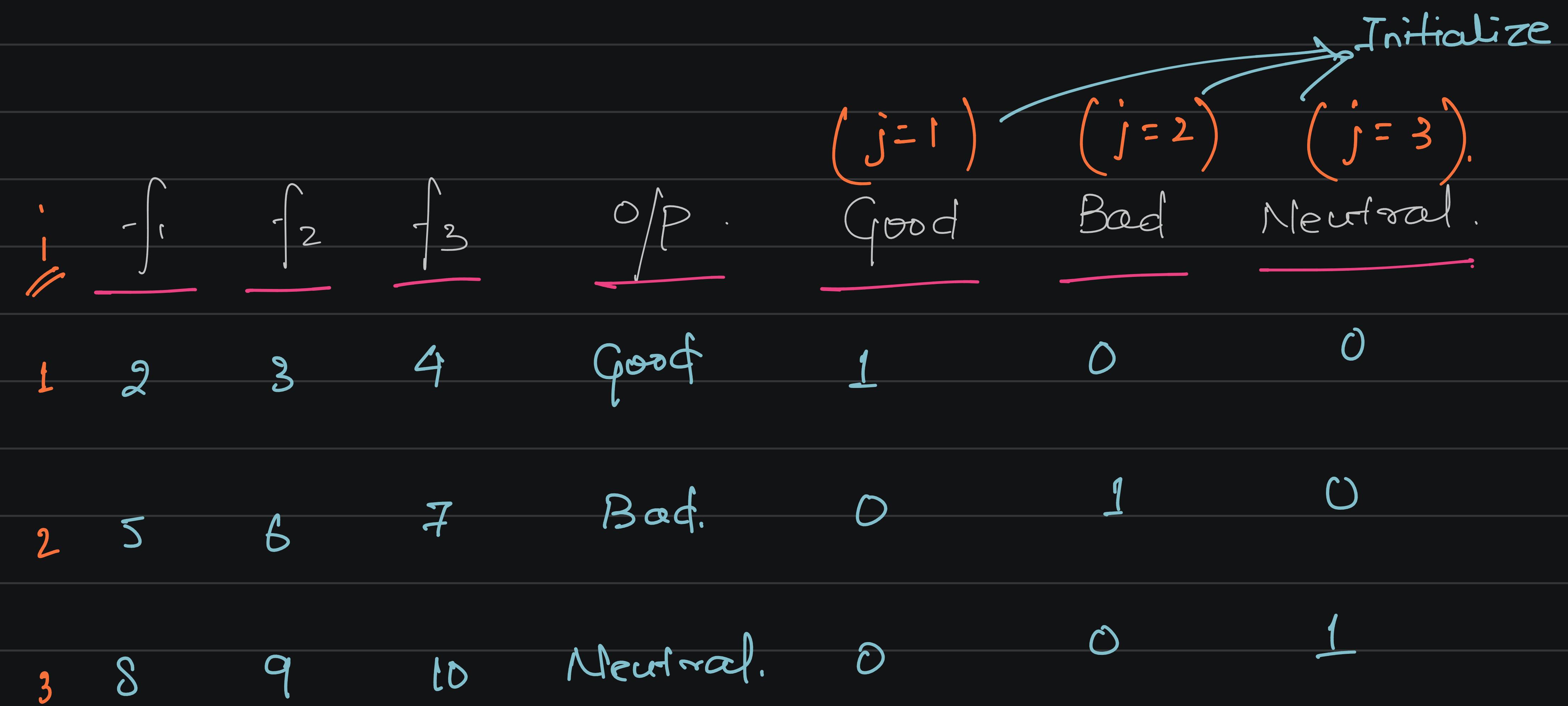
$$\hat{y} = \frac{1}{1 + e^{-z}}$$

\rightarrow Sigmoid activation
function.

(as binary classification)

2. Categorical Cross Entropy.

(Multiclass classification)



$C = \text{No. of categories.} = 3 \text{ (Initialize).}$

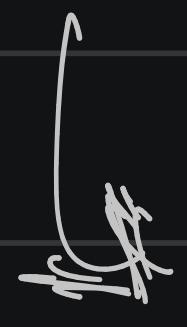
$$\text{Loss}(x_i, y_i) = -\sum_{j=1}^C y_{ij} + \ln(\hat{y}_{ij})$$

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

~~\hat{y}_{ij}~~ \Rightarrow Softmax activation = $\text{Sof}(z)$.

$$= \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}} = \frac{e^{z_i}}{e^1 + e^2 + \dots + e^k}$$

3. Sparse Categorical Cross Entropy.



Similar to categorical cross entropy.

but,

$$\text{Probabilities for } \text{Soft}(z) = [0.2, 0.3, 0.5]$$

O/P. \leftarrow Highest prob.

Ignore all
others and

just give the O/P as the
one with highest probability.

e.g. Used in image classification.

Disadvantage

- Losing all other information.

⑥ Right Combination

Activation applied



	<u>Hidden Layers</u>	<u>O/P Layer</u>	<u>Problem Statement</u>	<u>Loss function</u>
①	ReLU or its Variants	Sigmoid	Binary Classification	Binary Cross Entropy
②	ReLU or its Variants	Softmax	Multi class	Categorical or Sparse CE
③	ReLU or its Variants	linear	Regression	MSE, MAE, Huber loss, RMSE

2.

Optimizers.

↳ Used for backward propagation .

(1.) Gradient Descent .

(2.) Stochastic Gradient Descent (SGD) .

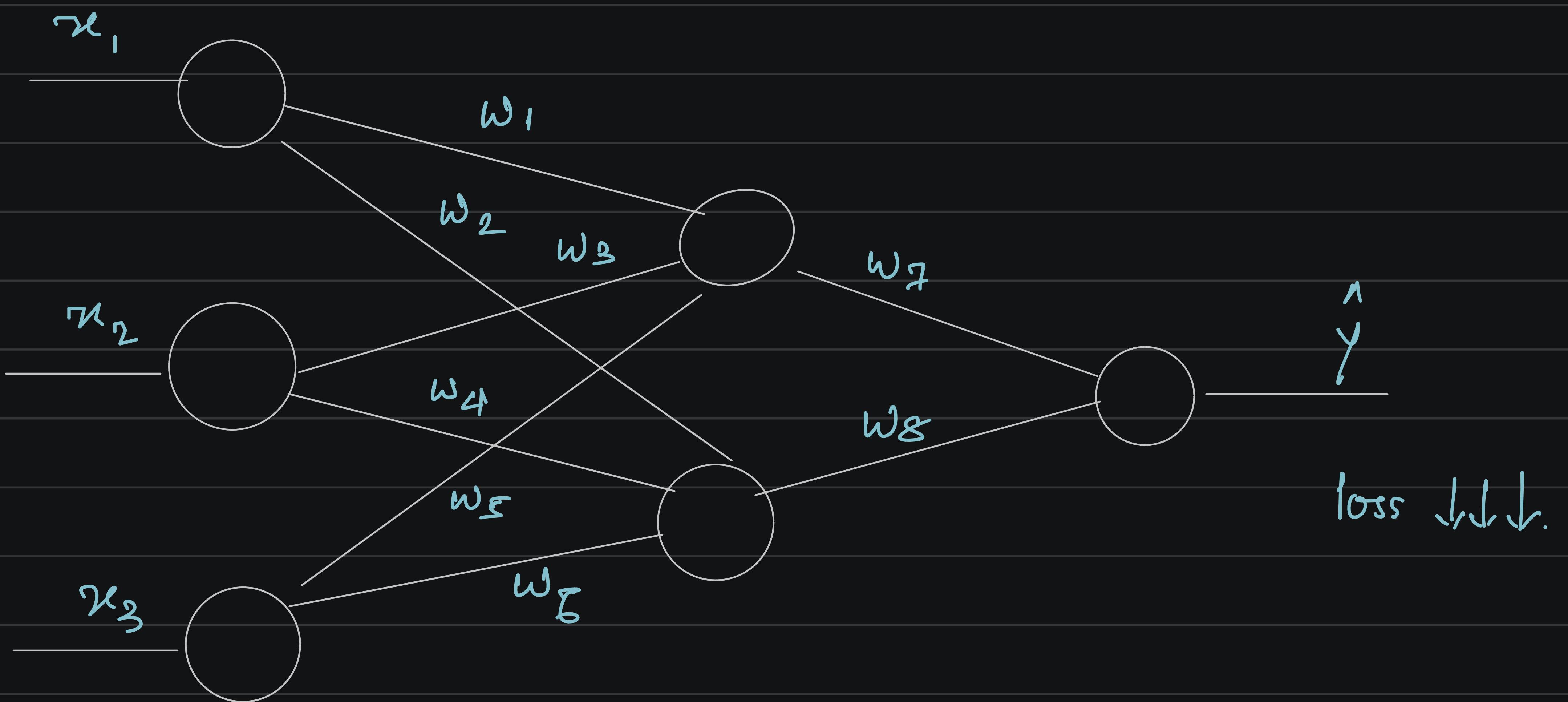
(3.) Mini batch SGD .

(4.) SGD with momentum .

(5.) Adagrad and RMS Prop .

(6.) Adam Optimizers .

(1.) Gradient Descent · Optimizer



Weight Updation formula .

$$w_{\text{new}} = w_{\text{old}} - \eta \left[\frac{\partial h}{\partial w_{\text{old}}} \right] \Rightarrow \text{slope}$$

Chain rule of derivative .



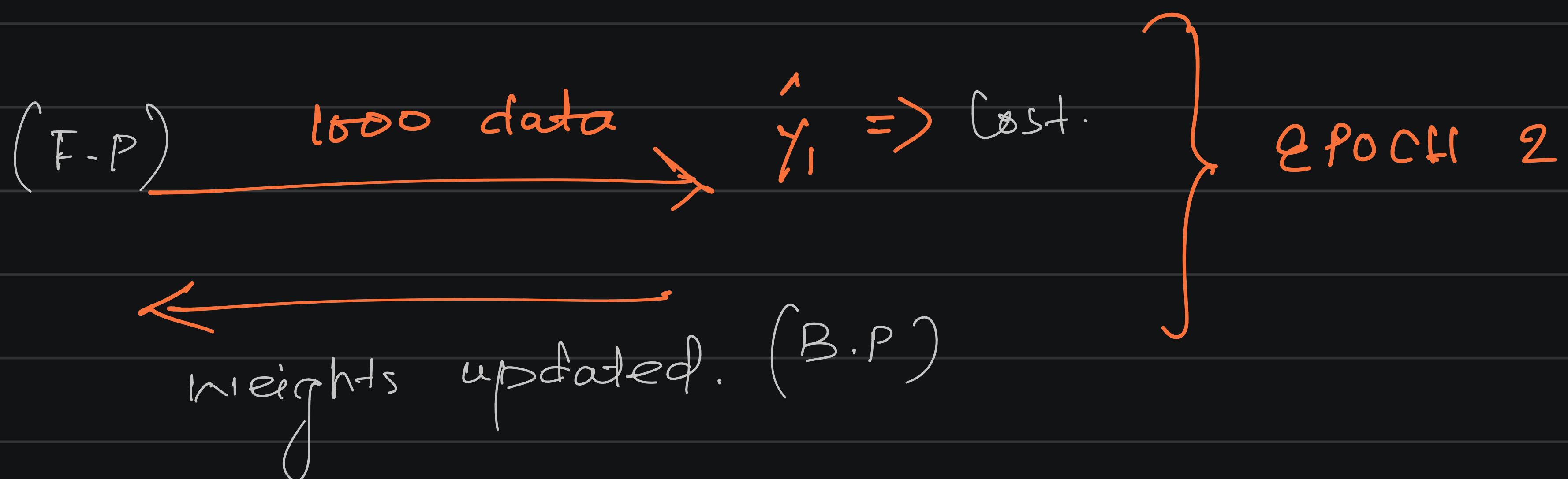
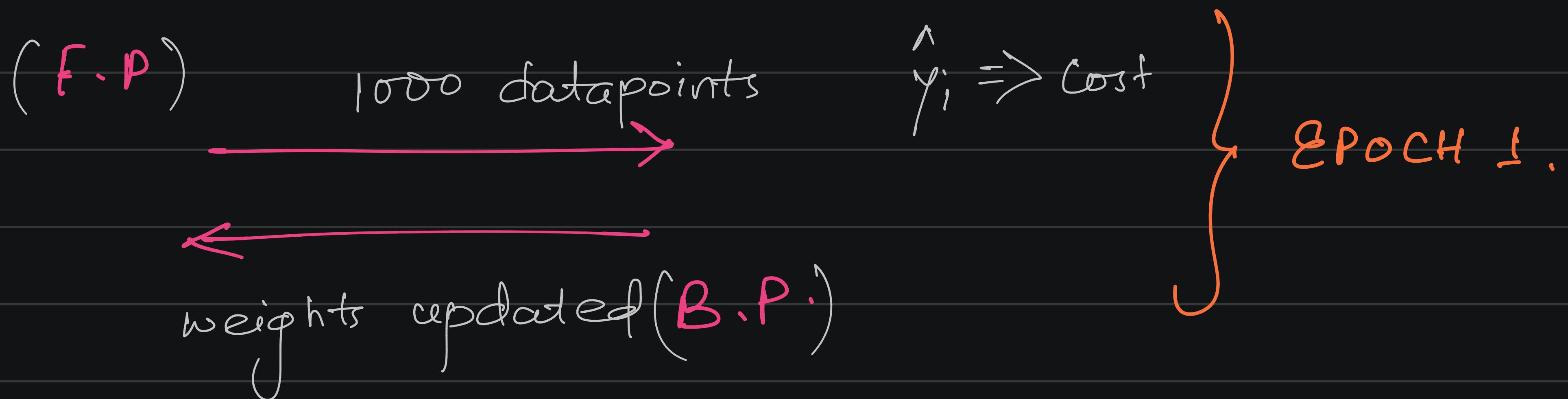
In MSE,

$$\text{Loss} = (\gamma - \hat{\gamma})^2$$

$$\text{Cost} = \frac{1}{n} \sum_{i=1}^n (\gamma_i - \hat{\gamma}_i)^2$$

Epochs vs Iteration.

[Dataset - 1000 datapoints]



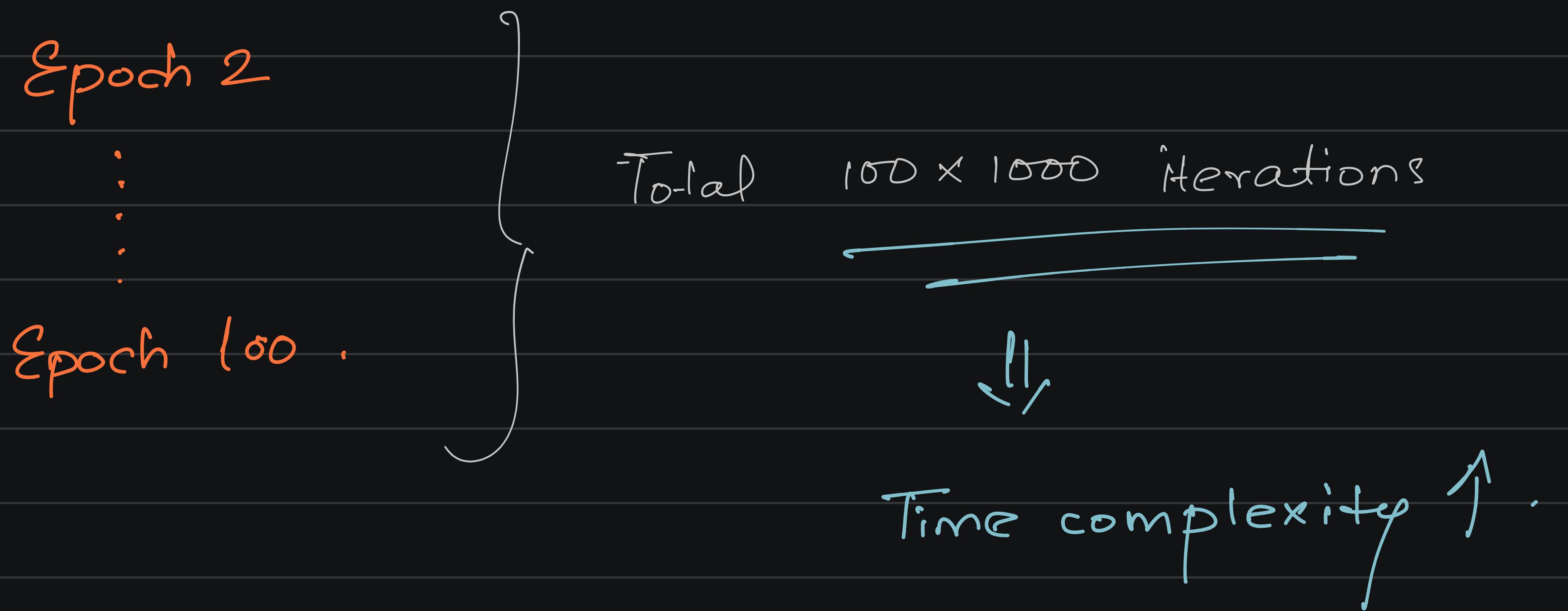
1
|
|
|
|
|
|

Until Cost ↓↓



Cost doesn't change in further

epochs. \Rightarrow Model is ready.



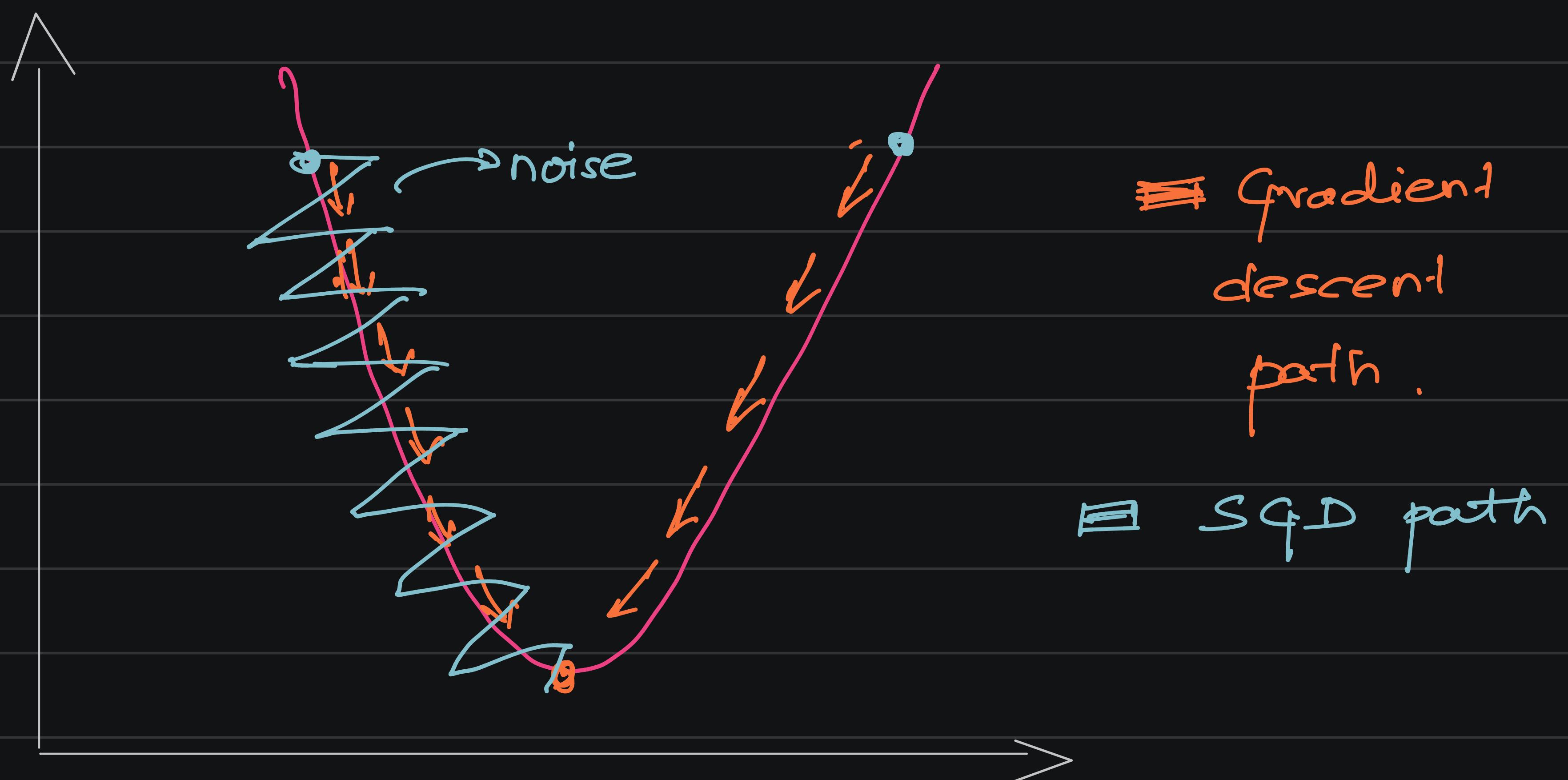
Advantage.

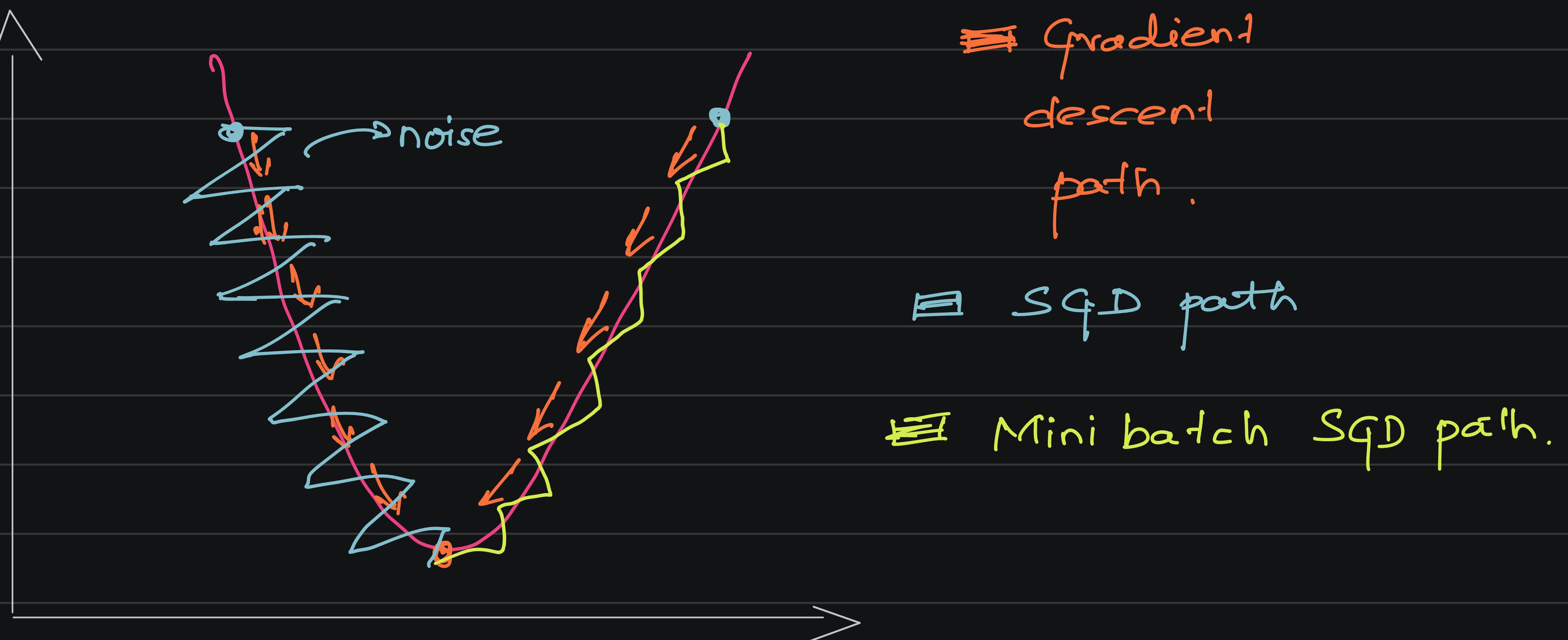
- * Solves resource issue .

Disadvantages.

- * Convergence will take time .

- * Noise will get introduced due to processing of each datapoint individually .





Disadvantage

- * Noise still exists.

4. SGD with Momentum.

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

$$b_{\text{new}} = b_{\text{old}} - \eta \frac{\partial h}{\partial b_{\text{old}}}.$$

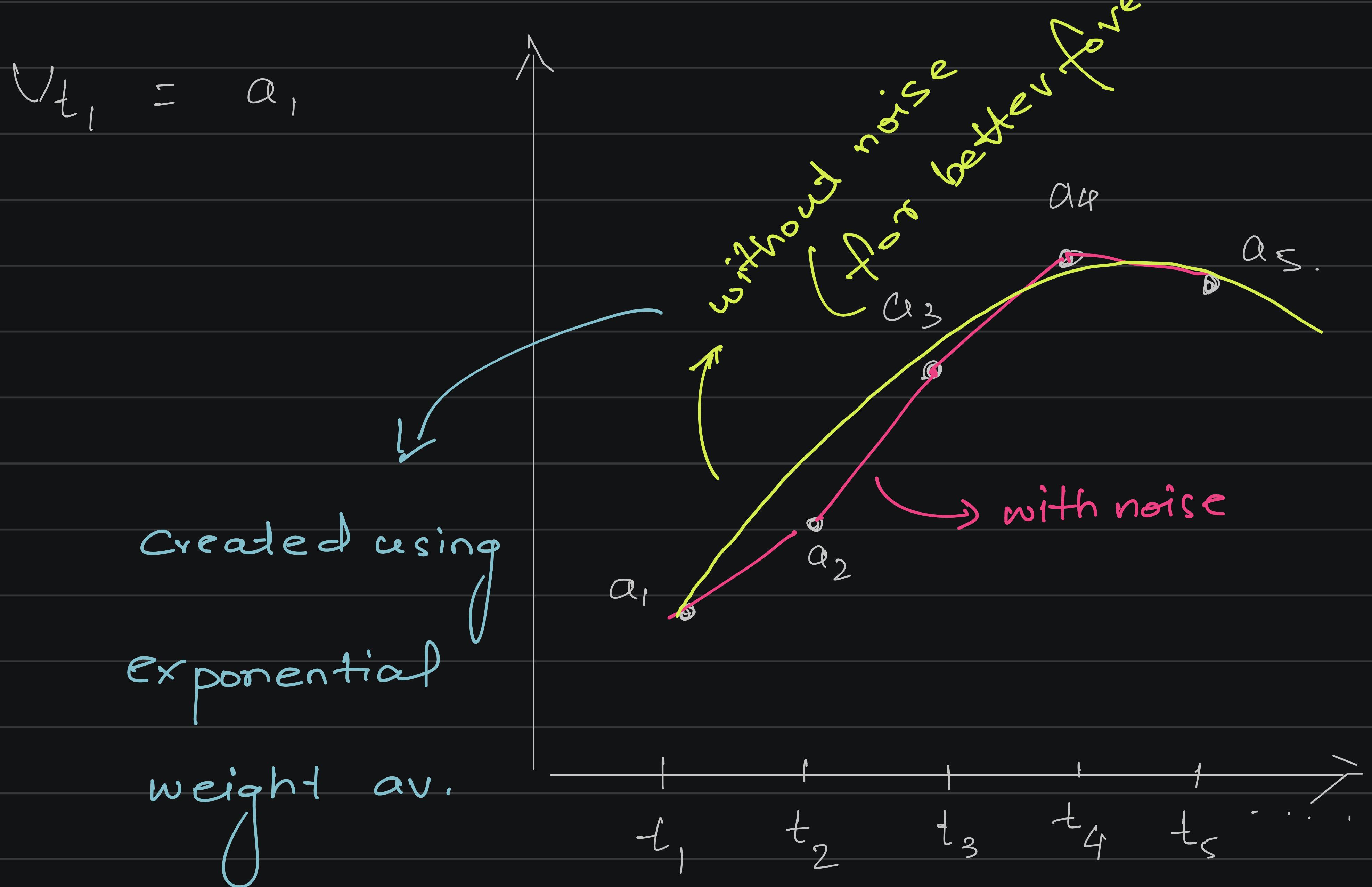
$$w_t = w_{t-1} - \eta \frac{\partial h}{\partial w_{t-1}}$$

→ Apply exponential weighted average.

Exponential Weight Average.

Time t_1 t_2 t_3 t_4 ... t_n

Value. a_1 a_2 a_3 a_4 ... a_n .



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

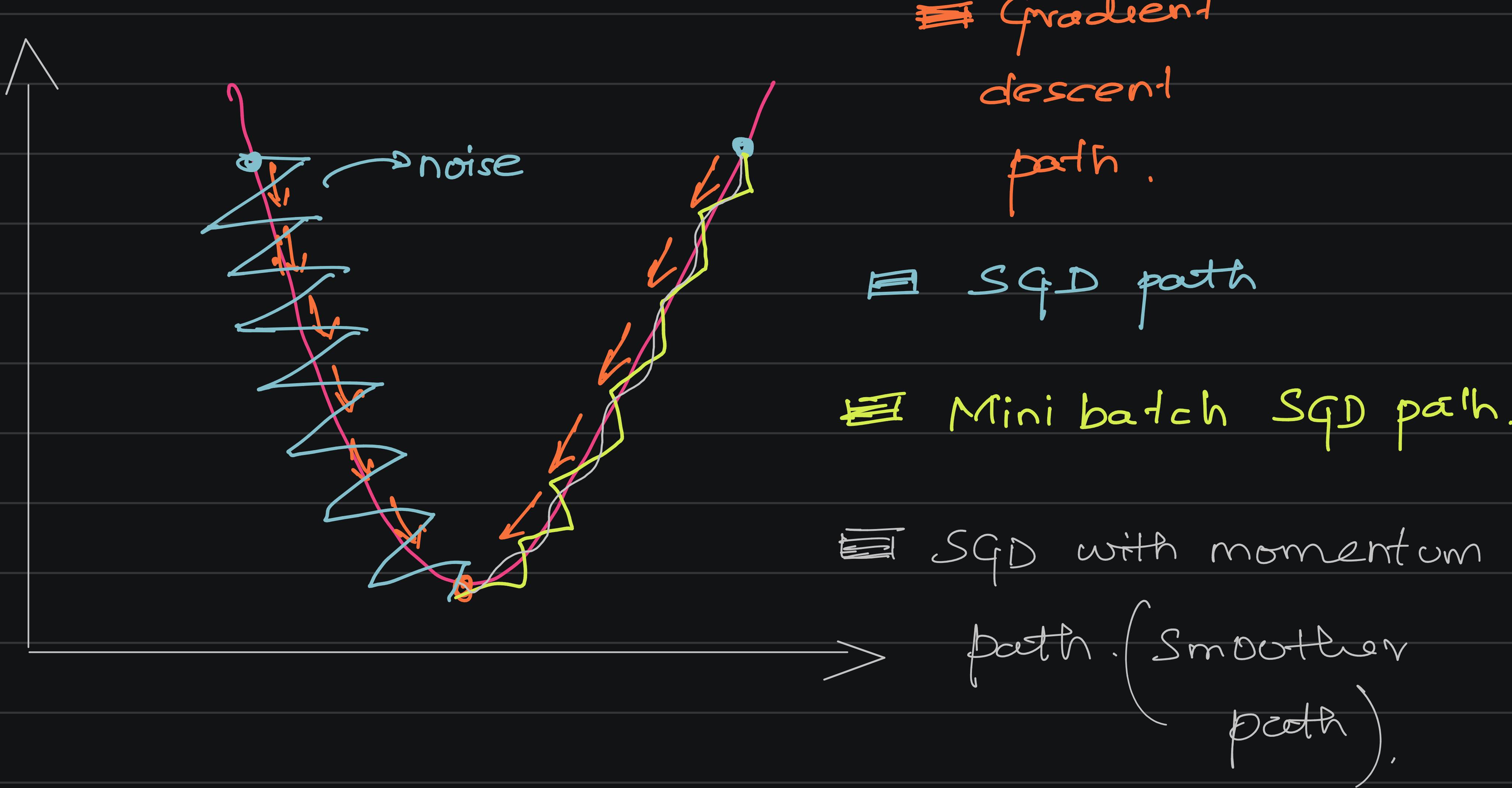
$$= 0.95 * a_1 + 0.05 * a_2 \quad [\text{let } \beta = 0.95]$$

↓
has more control.

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

$$= 0.95 * \underline{\underline{Vt_2}} + 0.05 (a_3)$$

\downarrow
more control. of the curve.



Advantage:

- * Reduces the noise
- * Smoothes the noise
- * Quick convergence.

5. Adagrad. { Adaptive Gradient Descent }

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

\downarrow

$$w_t = w_{t-1} - \eta' \frac{\partial h}{\partial w_{t-1}}$$

$$\eta' = \frac{\eta}{\sqrt{x_t + \epsilon}}$$

epsilon = small value.

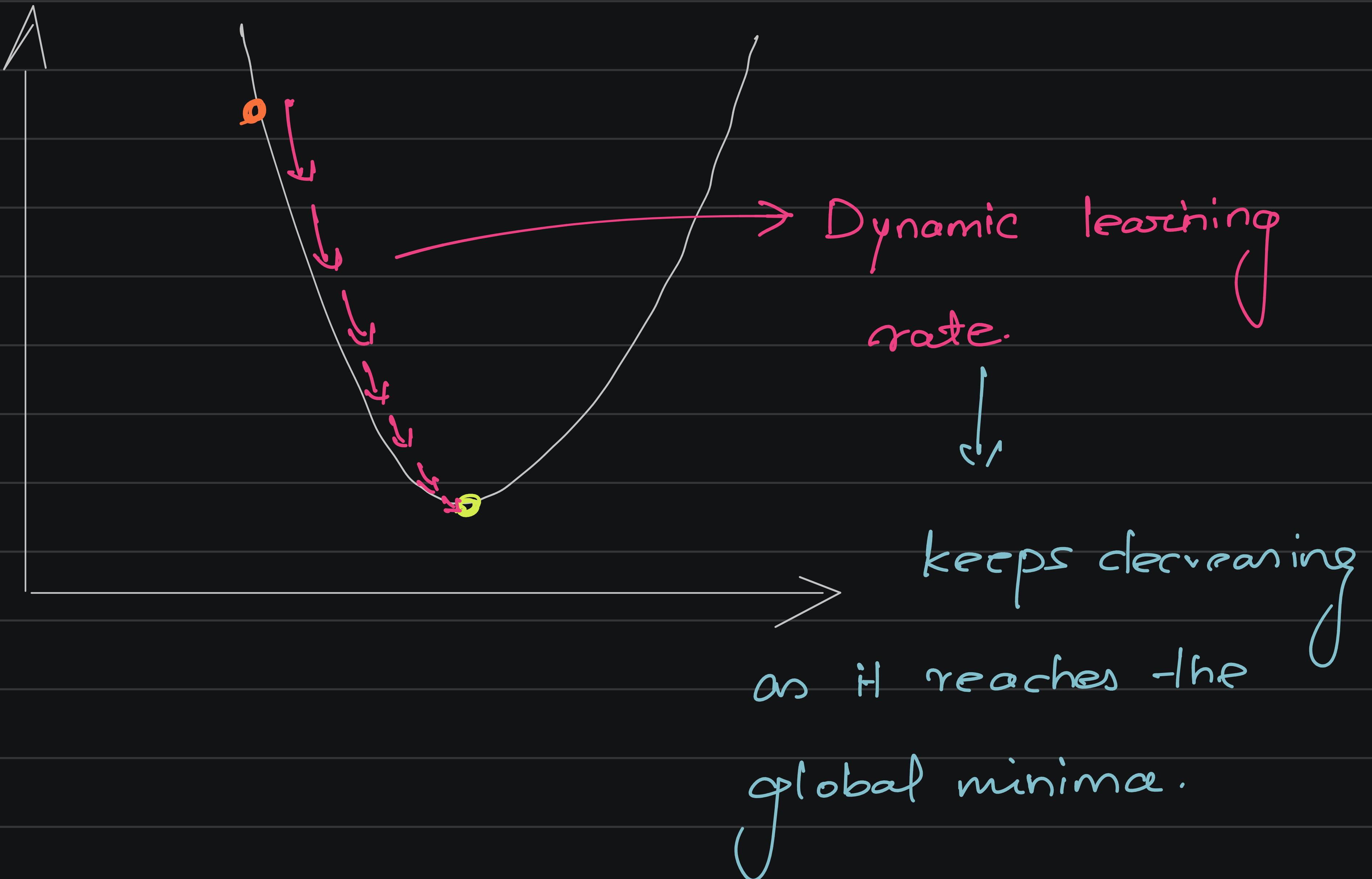
$$\alpha_t = \sum_{i=1}^t \left(\frac{\partial h}{\partial w_t} \right)^2$$

↓
to avoid dividing by 0
condition when $\alpha_t = 0$.

with each iteration.

$$\alpha_t \uparrow \cdot \equiv \eta' \downarrow \rightarrow$$

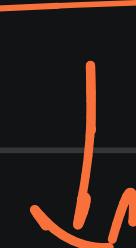
So, η' keeps getting smaller by the time of convergence.



Disadvantage

if $\alpha_t = \sum_{i=1}^t \left(\frac{\partial \mathcal{L}}{\partial w_t} \right)^2 \Rightarrow$ Large value.
 (in deep layered neural networks)

$$w_t \approx w_{t-1}$$



Vanishing gradient problem.

* So, η^t = Possibility to become a small value ≈ 0 .

6.

Adadelta and RMS Prop.

$$\eta' = \frac{\eta}{\sqrt{S_{dw} + \epsilon}}$$

↓ instead of α_t

Exponential weighted average.

(smoothening concept)
is used.

Initially, $S_{dw} = 0$

$$S_{dw} = \beta * S_{dw_{t-1}} + (1 - \beta) \left(\frac{\partial L}{\partial w_t} \right)^2$$

7.

Adam Optimizer → Best so far.



SGD with momentum + RMS prop [Dynamic Learning Rate]

$$w_t = w_{t-1} - \eta' v_{dw}$$

$$V_{dw} = \underline{V_{dw,t}} = \beta * V_{dw,t-1} - (1-\beta) \frac{\partial h}{\partial w_{t-1}}$$

$$\gamma' = \frac{\gamma}{\sqrt{s_{dw,t} + \epsilon}}$$

$$b_t = b_{t-1} - \gamma' \underline{V_{db}}.$$

$$V_{db} = \beta * V_{db,t-1} - (1-\beta) \frac{\partial h}{\partial b_{t-1}}$$

$$\gamma' = \frac{\gamma}{\sqrt{s_{db,t} + \epsilon}}$$