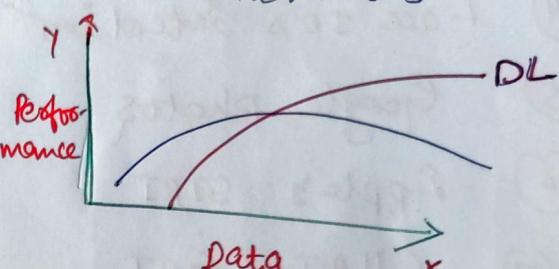


# \* \* Deep Learning \* \*

## ML V/S DL

- (1)  $ML \rightarrow FE$  (domain  
↓  
Knowledge)  
FS
- (2) Data increases but  
Performance not increases  
at that level.
- (3) ML models unable to  
identify Complex (nonlinear)  
Pattern.

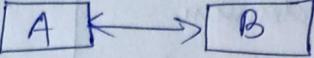
- (1) Automated FE
- (2) Data and Performance  
both increases.  

- (3) DL adept at identifying  
Complex Patterns.  
⇒ Complex data

DL → It is a method in AI that teaches  
computers to process data in a way that  
is inspired by the human brain.

## Need for Deep learning ⇒

- (1) To solve very complex problem.
- (2) Data management becomes easier in terms of hardware.
- (3) Handling Large data sets.
- (4)

## \* Uses cases in DL →

- (1) YouTube Recommendation → By YouTube history  
→ for Experience Top Notch
- (2) Google Translate → 
- (3) Netflix Recommendation
- (4) Face ID's detection
- (5) Google Photos
- (6) Apple's Siri
- (7) CHAT - GPT

\* Deep Learning ← → Neural Network (NN)

→ Any deep neural network called as DL.

## \* Characterise DL →

### (1) Artificial NN

- (1) Tabular Data
- (2) Fraud Detection
- (3) Price optimisation

### (3) Recurrent NN (RNN)

(1) Sequential data

- Text (Running text)
- Video
- Speeches
- Trading data

### (2) Convolution NN →

- (1) Images
- (2) Object Detection
- (3) Stable Diffusion
- (4) Fraud Recognition
- (5) Traffic Violation

(2) Speech Recognition

(3) Translation

(4) Forecast Prices

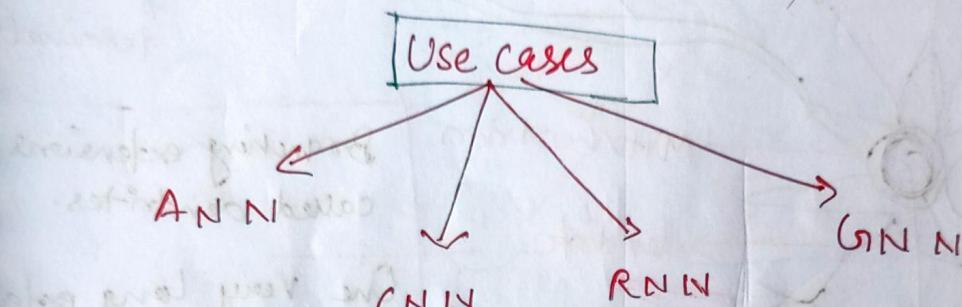
(5) Language Modeling

(6) Transformer → SEQ.

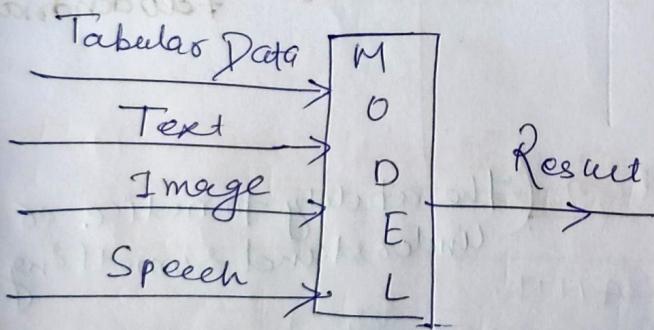
# Transformer — Advanced RNN

## ④ Graph NN

- Social Media Data
- Molecular Structure



## \* Multi-Model Nature →

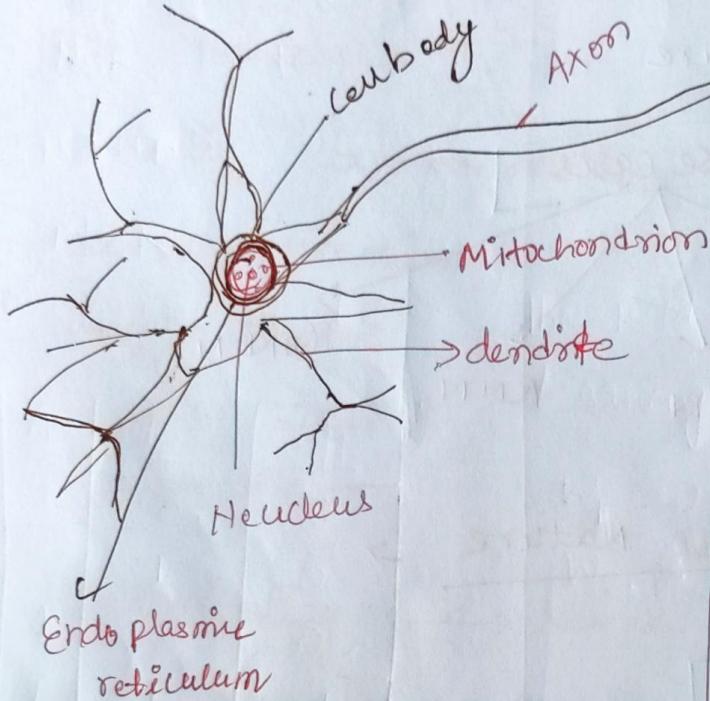


- here not only tabular data but also all types of data are used.

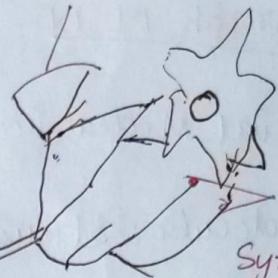
## \* Biological Neuron <→ Artificial Neuron

- Artificial neuron is based on Biological neuron, structure.
- Biological neuron is also called as spiking neuron.

## Biological Neurons →



Telodendria



Synaptic terminals

Branching extensions called dendrites.

One very long extension called Axon.

Axon splits off into many branches called telodendria.

(Fran Rosenblatt 1957) Perception (the ability to notice or understand something)

→ It is building block of ANN.

→ It works as linear classifier to make linear boundary.

(Stepping stone to solve complex Problems)

$$\frac{O(X_1, X_2, X_3, X_4, \dots)}{I/P - \text{Perception}} = (w_1x_1 + w_2x_2 + \dots + w_n)$$

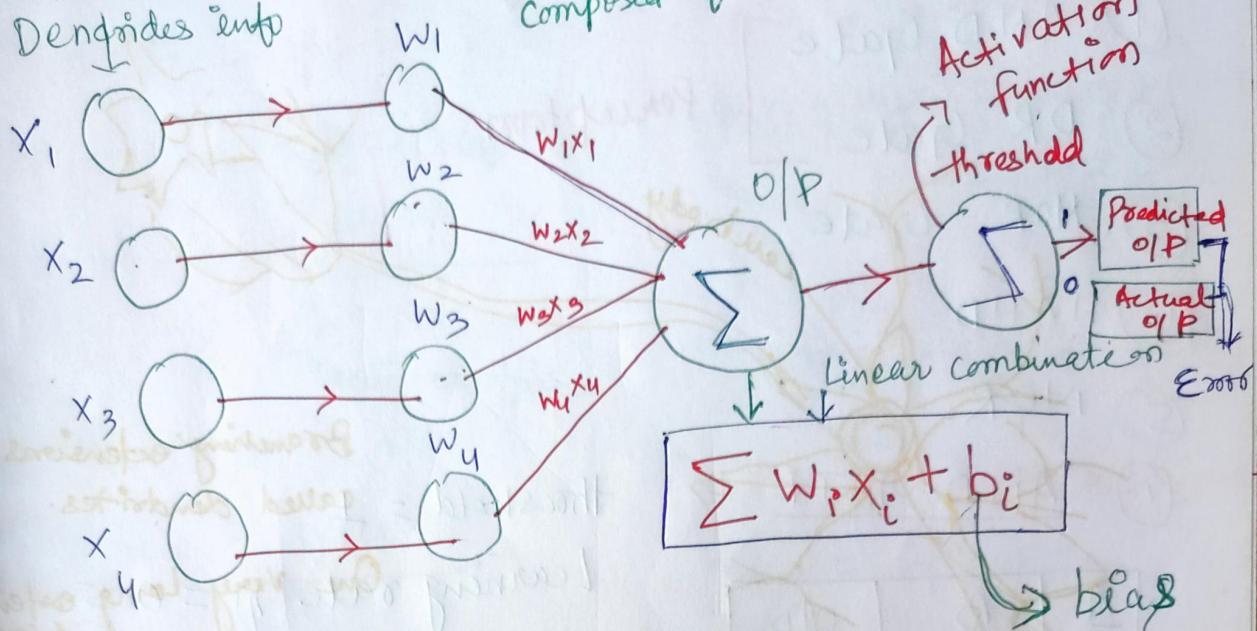
O/P → linear combination (Lc)

Fully connected layer also called dense layer.

Compare with threshold  
threshold > 1 ↓

$$\begin{cases} Lc > 1 \rightarrow O/P \rightarrow 1 \\ Lc < 0 \rightarrow O/P \rightarrow 0 \end{cases}$$

Architecture of Perceptron  $\Rightarrow$  A perceptron is simply composed of a single layer of TLUs.



$$\Rightarrow \text{Difference between Predicted O/P \& Actual O/P} = E_{\text{error}}$$

$\Rightarrow$  If  $LC < \text{threshold}$  or 0  $\Rightarrow$  Weight is updated by ~~error~~ to error.

\* Weight Update Formulae  $\Rightarrow$  [Training]

$$w_{\text{new}} = w_{\text{old}} + \eta ( \text{Target O/P} - \text{Pred. O/P} ) * x_i$$

$$E_{\text{error}} = \text{Target O/P} - \text{Pred. O/P}$$

$\eta \rightarrow$  Learning rate [0 to -1]

$x_i \rightarrow$  Input

\* Perceptron also can be multilayer.

\* It can be also called threshold logic unit.

$\Rightarrow$  Logical Gates solved by Perceptron  $\Rightarrow$

- (1) AND Gate ]  
 (2) OR Gate ] Perception  
 (3) NOR Gate  
 (4) NAND  
 (5) ~~NOR~~

(1) AND GATE  $\rightarrow$  threshold = 1

Learning rate ( $\eta$ ) = 0.5

	A	B	O/p
case 1	0	0	0
" 2	0	1	0
" 3	1	0	0
" 4	1	1	1

(1) case I :-  $A = B = 0 \mid p = 0$   $\xrightarrow{\text{Actual}}$

$$w_1x_1 + w_2x_2 > 1 \rightarrow 1$$

$x_1 = A = 0$

$x_2 = B = 0$

$$w_1x_1 + w_2x_2 = 0 < 1 \Rightarrow 0 \mid p = 0 \xrightarrow{\text{Predicted}}$$

Actual O/p = 0

$$\text{Error} = 0 - 0 = 0$$

$$w_{\text{new}} = w_{\text{old}} - \eta(0) * x_i$$

$\downarrow$

$$w_{\text{new}} = w_{\text{old}}$$

(2) Case - II

$$A = 0, B = 1, O/p = 0$$

$$w_1x_1 + w_2x_2 = 0 \times 0.2 + 1 \times 0.6 = 0.6 < 1$$

Predicted O/p = 0  
Error = 0

$$\therefore w_{\text{new}} = w_{\text{old}}$$

③ Case 4 :  $A = B = 1 = O/p = 1$

④ Case 3 :

$$W_1x_1 + W_2x_2 = 0.2 \times 1 + 0.6 \times 1 \\ = 0.8 > 1$$

$$A = 1, B = 0, O/p = 0$$

$$\text{Error} = 1 - 1 = 0$$

$$\therefore W_{\text{new}} = W_{\text{old}}$$

$$\text{Predicted } O/p = 1$$

$$W_1x_1 + W_2x_2 = 1.2 \times 1 + 0.6 \times 0 \\ = 1.2 > 1$$

$$O/p = 1$$

$$\text{Error} = 0.1 \rightarrow -1$$

Update weight

$$W_{\text{new}} = W_i + \eta (1) * 1$$

$$W_1_{\text{new}} = 1.2 + 0.5 \times 1 \\ = 1.2 - 0.5$$

$$W_1_{\text{new}} = 0.7$$

$$W_2_{\text{new}} = 0.6 + 0.5 \times 1 \\ = 0.6 - 0.5$$

$$W_2_{\text{new}} = 0.1$$

here in case ③

$$\frac{1}{1} \rightarrow 1.2 \rightarrow 0.7$$

### Advantages

- Single Layer Perceptrons can learn only linearly separable patterns.
- Easy to implement
- Low computation Power
- Perceptrons can implement Logic Gates like AND, OR

### Disadvantages

- Perception networks have several limitations.
- the output of a Perception can take on only one of two values (0 or 1) due to the hard-limit activation function.
- Here we just use the step function as an activation
- Perceptrons can only classify linearly separable set of vectors.

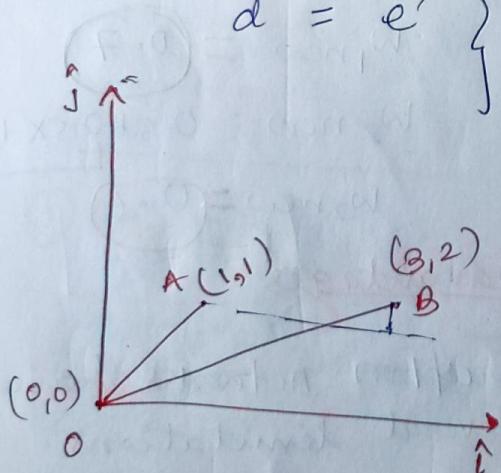
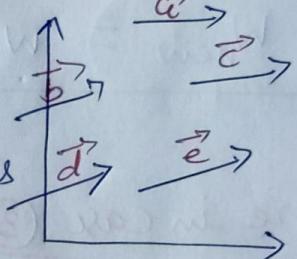
## Mathematical Concepts :-

- (1) Vectors
- (2) Differentiation
- (3) Partial Differentiation
- (4) Gradient
- (5) Maxima and Minima

\* Vectors :→ Any object which have two characteristic  
one of magnitude and 2nd is direction  
Called Vectors.

→ if magnitude & direction of any vectors  
are same then those are equal.

$$\vec{d} = \vec{e} \quad \left. \begin{array}{l} \text{magnitude & dirn same.} \\ \text{ } \end{array} \right\}$$



$$\overrightarrow{OA} = 1\hat{i} + 1\hat{j}$$

$$\overrightarrow{OB} = 3\hat{i} + 2\hat{j}$$

$$\overrightarrow{AB} = \overrightarrow{OB} - \overrightarrow{OA}$$

$$\boxed{\overrightarrow{AB} = (3\hat{i} + 2\hat{j}) - (1\hat{i} + 1\hat{j})}$$

$$|\overrightarrow{AB}| \text{ Magnitude} = \sqrt{g^2 + f^2} = \sqrt{5}$$

### Scales

- (1) Only magnitude exists

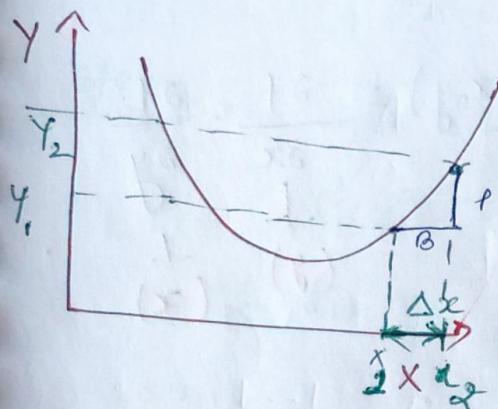
Ex :- Speed

### Vectors

- (1) Magnitude & direction both exists.

Ex:- Velocity, Gravity  
Displacement etc.

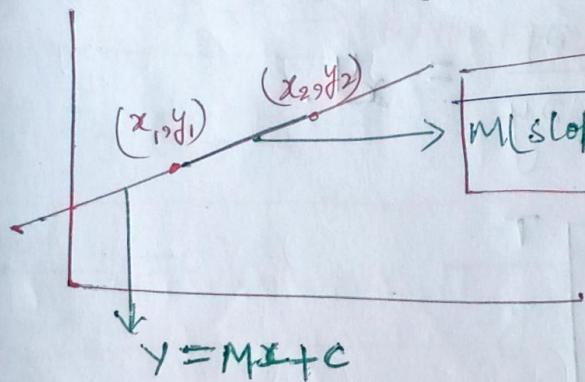
\* Differentiation  $\Rightarrow$  Rate of change of variable in any function.



Perpendicular Base  $= \tan \theta$

$$\textcircled{1} \quad \frac{\Delta y}{\Delta x} = \frac{dy}{dx}$$

$$\textcircled{2} \quad \tan \theta = \frac{\Delta y}{\Delta x} = \frac{y_2 - y_1}{x_2 - x_1}$$



$$\textcircled{3} \quad \lim_{x \rightarrow 0} \frac{f(x_1 + \Delta x) - f(x_1)}{\Delta x}$$

$$\therefore \boxed{\text{Differentiation} = \lim_{x \rightarrow 0} \frac{f(x_1 + \Delta x) - f(x_1)}{\Delta x}} \quad [x_2 = x_1 + \Delta x]$$

$$y = x^5 \quad \textcircled{1} \quad \text{Use Power Rule}$$

$$\frac{dy}{dx} = ?$$

$$\boxed{y = x^n, \frac{dy}{dx} = nx^{n-1}}$$

$$\text{Ex: } \frac{dy}{dx} = 5x^4 \quad \textcircled{2} \quad \text{Product Rule}$$

$$\boxed{f(x)g(x) = f'(x)g(x) + f(x)g'(x)}$$

$$\text{Ex: } f(x) = (x+2) \cdot \cos x \quad \frac{dy}{dx} = 4x^3 \cos x + (x+2) \sin x \\ = 4x^3 \cos x - x \sin x - 2 \sin x$$

### ③ Partial Differentiation $\Rightarrow$

$$f(x, y, z) = \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z} . f(x, y) = \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}$$

↓                    ↓                    ↓                    ↓                    ↓                    ↓  
 (y, z) const      (x, z)      (x, y)      (y)      (x)

Ex:  $f(x, y) = x^4 y$

$$\textcircled{1} \quad \frac{\partial f}{\partial x} = 4x^3 \cdot y . \textcircled{2} \quad \frac{\partial f}{\partial y} = x^4 .$$

### ④ Gradient $\Rightarrow$

$$\nabla f = \left[ \frac{\partial f}{\partial x} \parallel \frac{\partial f}{\partial y} \right]^T$$

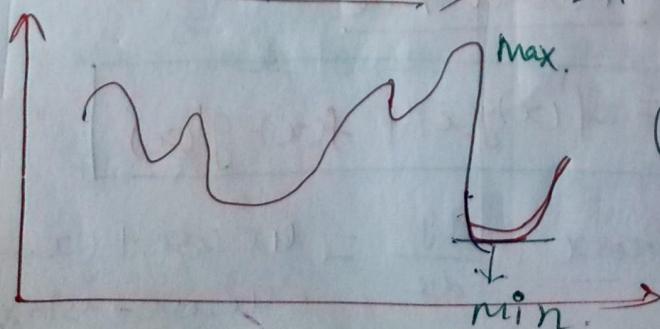
\* Partial derivatives are written in vector format.

$$\nabla f = \left[ 4x^3 y \parallel x^4 \right]^T = \begin{bmatrix} 4x^3 y \\ x^4 \end{bmatrix} .$$

$$f(x, y) = 2x^2 + 4y$$

$$\nabla f = \left[ \begin{array}{c} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{array} \right] = \begin{bmatrix} 4x^3 \\ 4 \end{bmatrix}$$

### ⑤ Maxima & Minima $\Rightarrow$ used by $\rightarrow$ error $\rightarrow$ min profit $\rightarrow$ max.



$$\textcircled{1} \quad y = f(x) \rightarrow \text{Univariate}$$

Step 1:  $f'(x) = 0$

Step 2: find critical points (a), (b)

Step 3:  $f(x) = x^5$

$$f'(x) = 5x^4$$

$$f''(x) = 20x^3$$

\*  $f(x) \rightarrow$  Bivariate

①  $y = f(x, y)$

$$\left[ \frac{\partial f}{\partial x} = 0, \quad \frac{\partial f}{\partial y} = 0 \right] \rightarrow \text{find critical points}$$

$\downarrow$

$f(x)$                        $f(y)$

②  $f_{xx}$

$$f_{yy}$$

$$\left[ \frac{\partial^2 f}{\partial x^2} = f_{xx}, \quad \frac{\partial^2 f}{\partial y^2} = f_{yy} \right] \rightarrow \left( \frac{\partial f}{\partial x} \right) \rightarrow (x, y) \rightarrow \frac{\partial^2 f}{\partial y^2}$$

③ Critical Point  $(a, b) \rightarrow$  Check Minima or Maxima.

Case-1  $\rightarrow \Delta = \sigma^2 > 0 \rightarrow$  Minima/Maxima

$\sigma > 0 \rightarrow$  Local min.  $\sigma < 0 \rightarrow$  Local max.

Case-2  $\rightarrow \Delta = \sigma^2 \leq 0 \rightarrow$  Point of ~~repetition~~

Case-3  $\rightarrow \Delta = \sigma^2 < 0 \rightarrow$  neither max nor min saddle pt.

$\gamma > 0 \rightarrow$  Local Minima }  
 $\gamma < 0 \rightarrow$  Local Maxima }

### \* Activation Functions $\Rightarrow$

The purpose of the activation function is to introduce non-linearity into the output of a neuron. The activation function does the non-linear transformation to the input making it capable to learn and perform more complex tasks.

Ex :- XOR operation

\* Activation functions are : There are most used activation functions.

(1) Sigmoid

(2) ReLU (Rectified Linear Unit)

(3) Leaky ReLU

(4) tanh (Hyperbolic Tangent)

$$w_i(x) + b_i \text{ali}$$

(5) Soft max

(6) Parametric ReLU

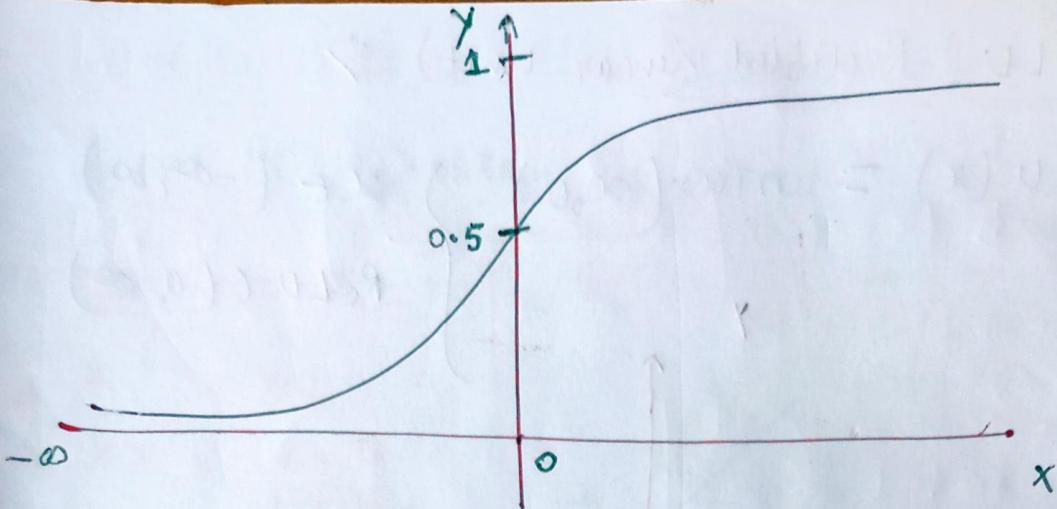
(7) Sigmoid function  $\Rightarrow$  It is used in binary classification

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

$$; x \in (-\infty, \infty)$$

$$\sigma(x) \in [0, 1]$$

In this function, output belongs to  $[0, 1]$



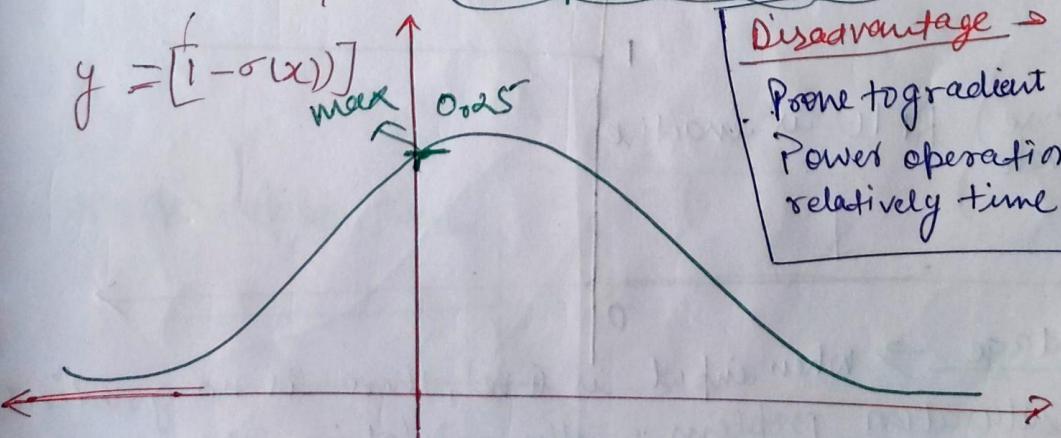
\* We would need to understand the derivative of every activation function, Because derivative helps our NN models to optimise Weights for our NN.

Derivative of sigmoid functions →

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad [\text{Power product}]$$

$$\sigma'(x) = \frac{\sigma(x)}{1 - \sigma(x)} \quad (\frac{u}{v})' = \frac{u'v - uv'}{v^2}$$

Derivative  $\{0, 1\}$  Graphical Structure



Disadvantage →

Prone to gradient vanishing  
Power operations are relatively time consuming

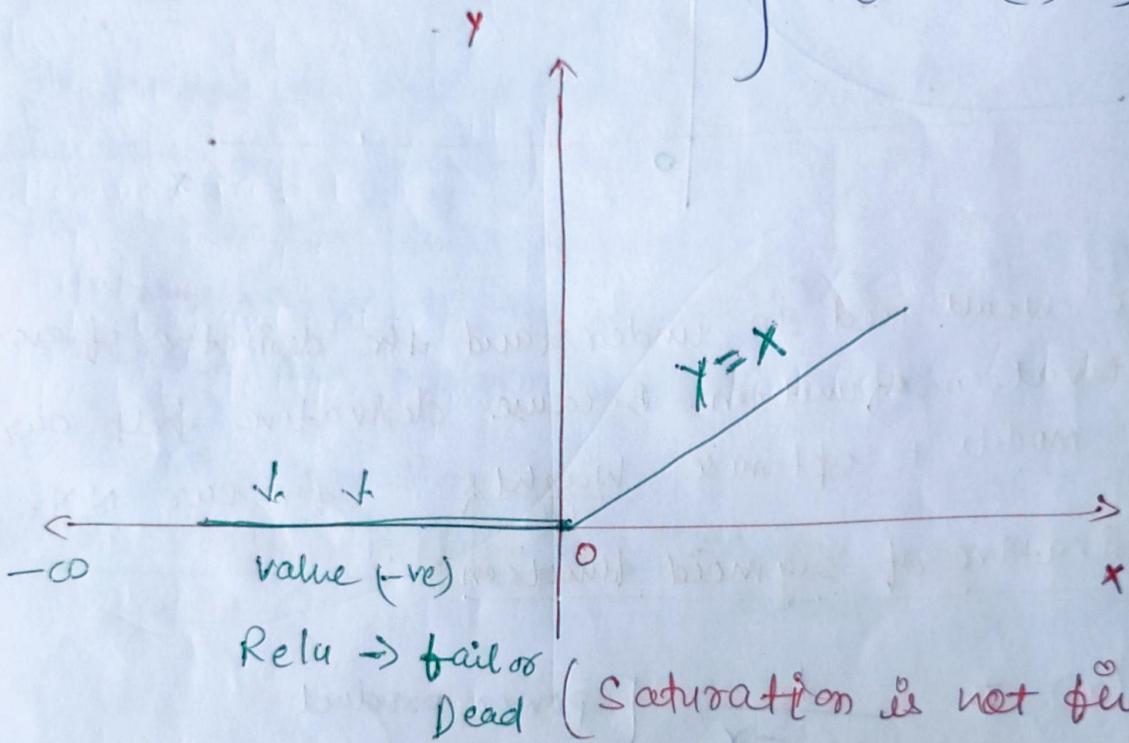
Advantage →

derivative of sigmoid

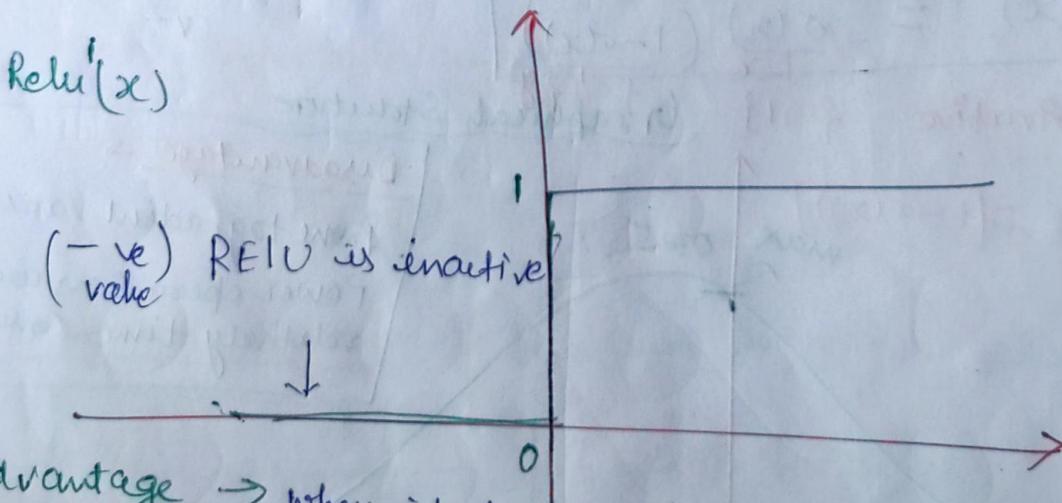
→ Smooth Gradient, preventing 'jumps' in output values.

## ② RELU (Rectified Linear Unit) $\Rightarrow$

$$\text{ReLU}(x) = \max(x, 0) \quad \left. \begin{array}{l} x \in (-\infty, \infty) \\ \text{ReLU} \in [0, \infty) \end{array} \right\}$$



Derivative of ReLU  $\rightarrow$



Advantage  $\rightarrow$  When input is (+ve), there is no gradient saturation problem. The calculation speed is much faster.

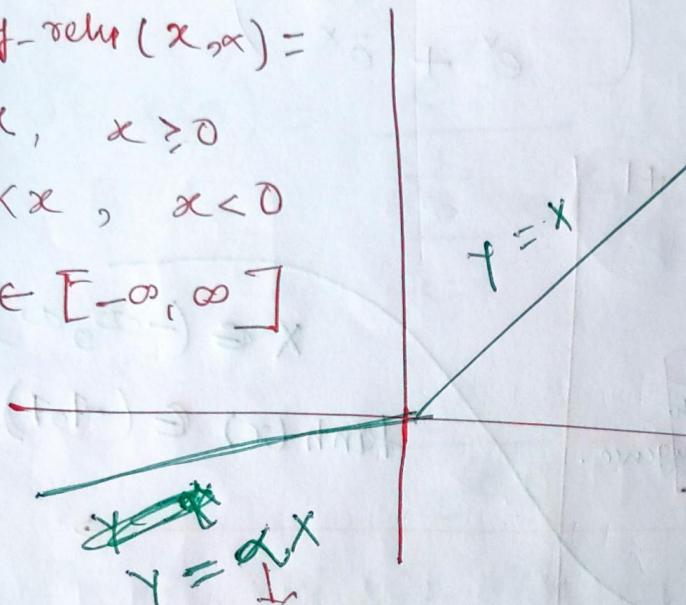
Disadvantage  $\rightarrow$  When the input is -ve, ReLU is completely inactive (than sigmoid & tanh)

(3) LEAKY RELU  $\rightarrow$  It is special case of Parametric RELU.  
 $\rightarrow$  Dead RELU problem is solved.

Leaky-relu( $x, \alpha$ ) =

$$\begin{cases} x, & x \geq 0 \\ \alpha x, & x < 0 \end{cases}$$

$$x \in [-\infty, \infty]$$

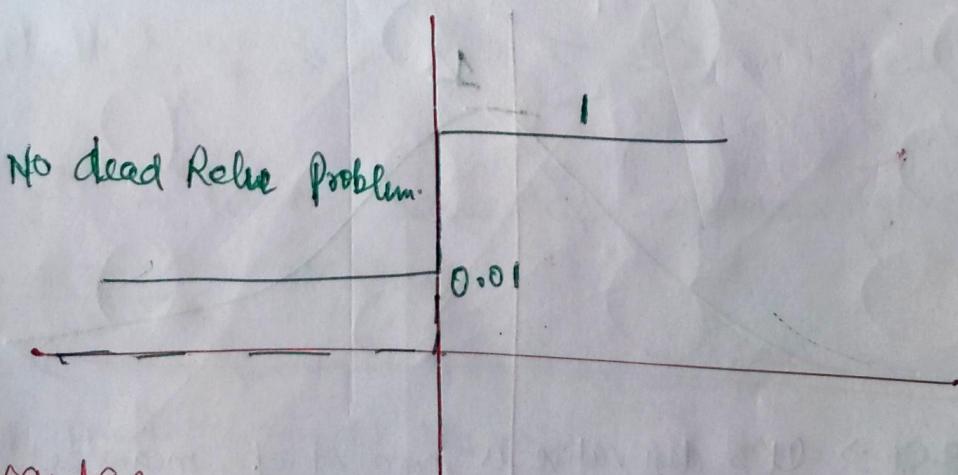


If  $\alpha = 0.01 \rightarrow$  Leaky Relu

if  $\alpha \in (\text{Other than } 0.01) \rightarrow$  Randomized Relu

if  $\alpha = \text{Learnable} \rightarrow$  Parametric Relu.

Derivative of Leaky Relu  $\rightarrow$



Advantage  $\rightarrow$  Gradient vanishing Problem is not present.

#### (4) Hyperbolic tangent activation function (tanh)

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

This activation fun is

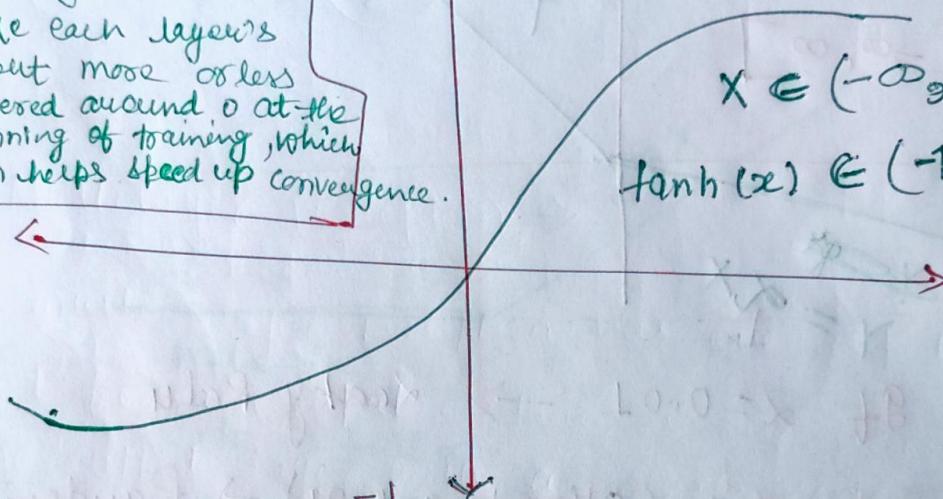
S-shaped, continuous, and differentiable.

The range tends to make each layer's output more orders centered around 0 at the beginning of training, which often helps speed up convergence.

+1

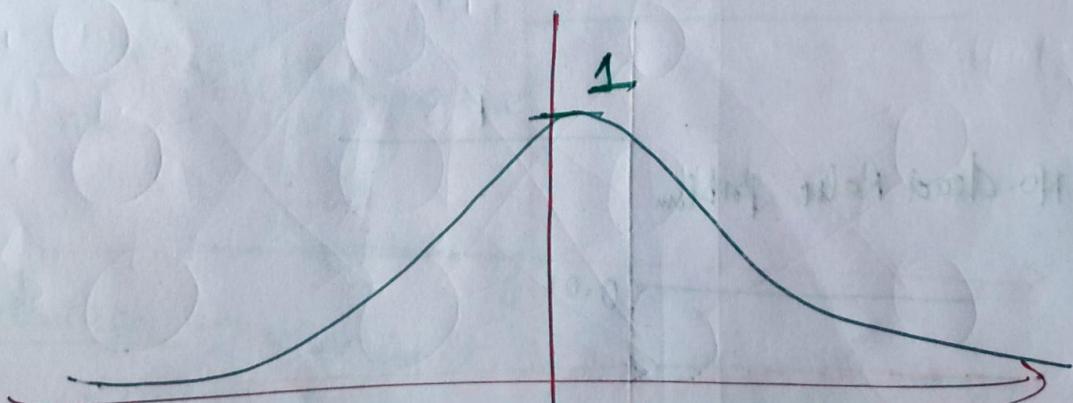
$$x \in (-\infty, \infty)$$

$$\tanh(x) \in (-1, 1)$$



→ Now step curve is found than Sigmoid fun.

$$\text{tanh}'(x) = 1 - \tanh^2(x)$$



Advantage → Its derivative is more steep, means it gets

→ It has a wider range for faster learning more value.

Disadvantage → More computation expensive.

## (5) Soft Max activation function : $\rightarrow$

\*  $\rightarrow$  It is used mostly in multiclass problem.

$$S(x_j \phi) = \frac{e^{x_j}}{\sum_{k=1}^K e^{x_k}}$$

If  $K = 1$ ,

$$\frac{e^x}{e^x} = 1$$

$$K = 2, P(X = 1/x) = \frac{e^{x_1}}{e^{x_1} + e^{x_2}}$$

If multiple class  $\rightarrow$

Case 1  $\rightarrow$

$$\frac{e^{x_1}}{e^{x_1} + e^{x_2} + e^{x_3}}$$

Case 2

$$\frac{e^{x_2}}{\sum e^{x_i}}$$

Case 3

$$\frac{e^{x_2}}{\sum e^{x_3}}$$

$\rightarrow$  It is also same as sigmoid function.

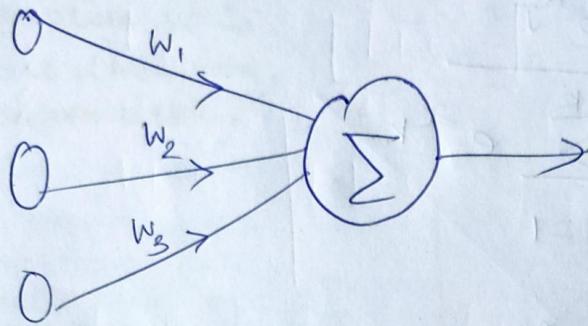
Advantage  $\rightarrow$  ~~gradient~~

The main advantage of using softmax is the output Probabilities range. The range will 0 to 1, and the sum of all the probabilities will be equal to one.

Disadvantage  $\rightarrow$  Vanishing gradient problem.

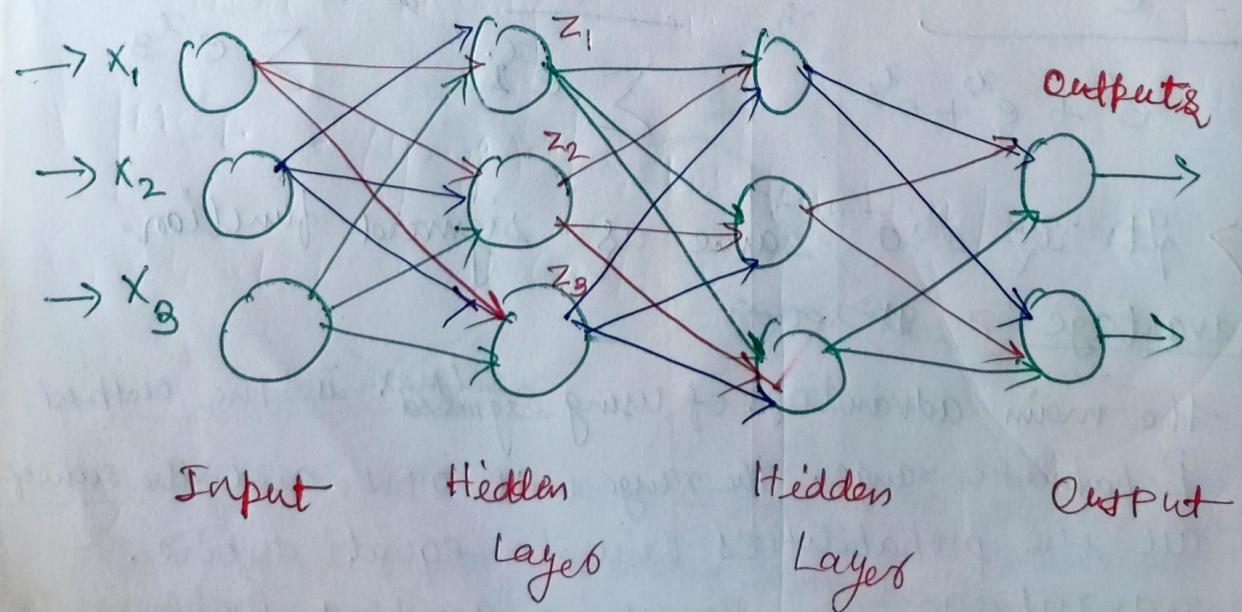
## \* Forward Propagation →

Information is fed through a network, in a forward direction, to generate an output.



## \* Multilayer Perceptron →

A multi-layer perceptron has one input layer and for each input, there is one neuron. It has one output layer with a single node for each output and it can have any no. of hidden layers and each hidden layer can have any no. of nodes.



$$x_1 \rightarrow \underbrace{(w_{ax_1} + w_{bx_2} + w_{cx_3})}_z$$

$$z_1 \rightarrow w_{ag_1} z_1 + w_{bg_2} z_2 + w_{cg_3} z_3$$

It is able to compute the gradient of the network's error with regard to every single model parameter.

Backward Propagation  $\Rightarrow$  Back propagation requires

a known, desired output for each input value in order to calculate the loss function gradients - that how a prediction differs from actual results - (as a type of supervised ML).

Chain Rule  $\Rightarrow$

$$f(x) = \frac{1}{1+e^x} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{Let } P = -x \quad p(x)$$

$$q = 1 + e^{-x} = q(P)$$

$$\boxed{f(q) = \frac{1}{q}} \quad \leftarrow \boxed{q = 1 + e^P}$$

$f(q(p(x))) \rightarrow$  Sigmoid function

$$\frac{dt}{dq} = \bar{q} = (-1) \times \bar{q}^{-1} = \frac{-1}{q^2}$$

$$\frac{dq}{dp} = 0 + e^p = e^p$$

$$\frac{dp}{dx} = -1$$

$$\frac{df}{dq} \times \frac{dq}{dp} \times \frac{dp}{dx} = \frac{df}{dx}$$

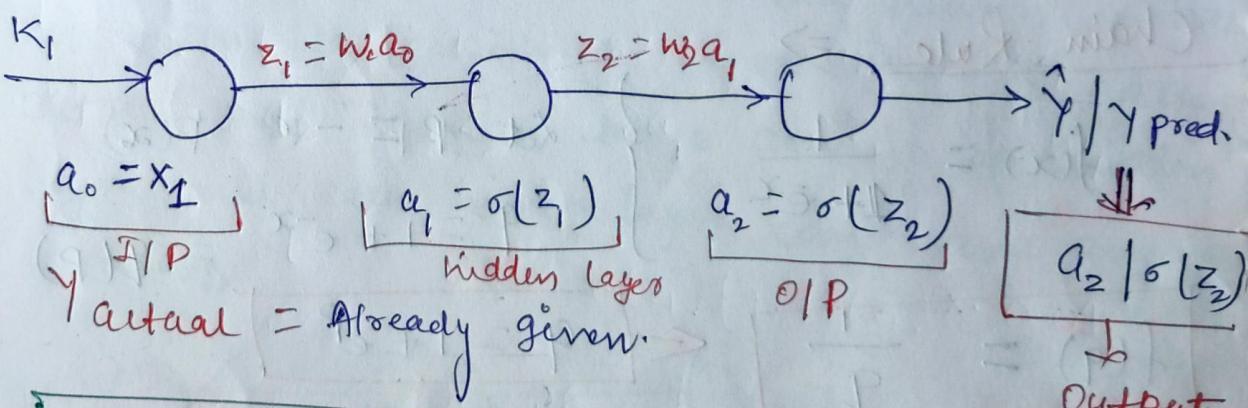
$$= \frac{e^P}{q^2} = \frac{e^{-x}}{[1+e^x]}$$

$$\Rightarrow \boxed{1+e^x = \frac{1}{1+e^x} \left[ 1 - \frac{1}{1+e^x} \right]}$$

$$\sigma|x| = \sigma|x| = (1 - \sigma|x|)$$

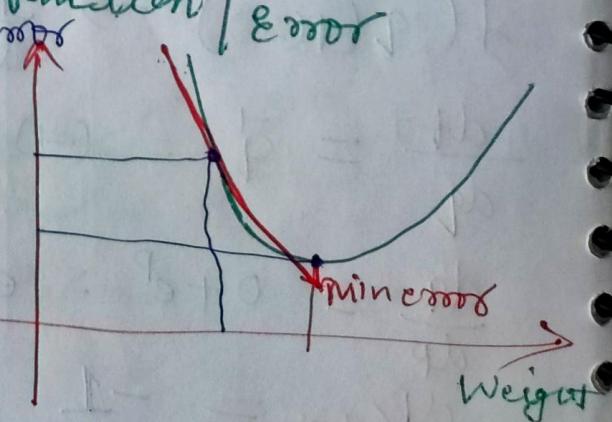
\* Back Propagation  $\Rightarrow$  using Chain Rule

In Simple NN  $\Rightarrow$



\*  $MSE = (y - \hat{y})^2 \Rightarrow$  Loss function / error

Gradient Descent Graph  $\Rightarrow$



$$W_{\text{new}} = W_{\text{old}} - \eta \frac{\partial E}{\partial W}$$

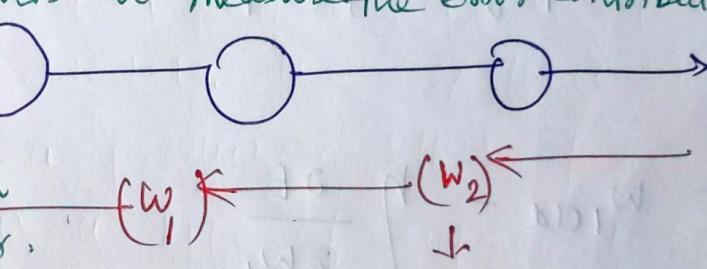
Gradient Descent weight formulae.

$\eta \rightarrow \text{learning rate}$   
 $\in (0, 1)$

Now

$\rightarrow$  We will go backward and find weight

$\Rightarrow$  Back propagation algo first makes a prediction (forward pass) and measures the algo, then goes through each layer in reverse to measure the error contribution from each connection and finally tweaks the connection weights to reduce the error.



$w_2$  update :

$$w_{2 \text{ new}} = w_{2 \text{ old}} - \eta \frac{\partial E}{\partial w_2}$$

$$\frac{\partial E}{\partial w_2} = \frac{\partial E}{\partial a_2} \times \frac{\partial a_2}{\partial z_2} \times \frac{\partial z_2}{\partial w_2}$$

$$\textcircled{1} E = (y - \hat{y})^2$$

$$= -2(y - a_2)(\sigma(z_2)(1 - \sigma(z_2)) \times a_1)$$

$$w_{2 \text{ new}} = w_{2 \text{ old}} - \eta \frac{\partial E}{\partial w_2} \rightarrow \text{find}$$

$$w_{2 \text{ new}} = w_{2 \text{ old}} - \eta \cancel{a_1}$$

$$\begin{aligned} \textcircled{1} \quad & \text{err} = (y - \hat{y})^2 \\ \textcircled{2} \quad & \hat{y} = \sigma(z_2) \\ & \sigma(z_2) = \frac{1}{1 + e^{-z_2}} \\ & \hat{y} = \sigma(w_2 a_1) \end{aligned}$$

$$\text{err} = (y - a_2)^2$$

$$\text{err} = f(a_2)$$

$$\text{err} = f(z_2)$$

$$\textcircled{3} \quad z_2 = w_2 a_1 \rightarrow f(w_2)$$

$$a_1$$

$w_1$  update  $\rightarrow$

① Error =  $(y - \hat{y})^2 = (y - a_2)^2 \rightarrow f(a_2)$

②  ~~$z_2$~~  =  $w_2 a_1 \Rightarrow f(a_1)$

③  $a_1 = \sigma(z_1) = f(z_1)$

④  $a_2 = \sigma(z_2) \Rightarrow f(z_2)$

⑤  $z_1 = w_1 a_0 = f(w_1)$

⑥

$$w_{1\text{new}} = w_{1\text{old}} - \eta \frac{\partial E}{\partial w_1}$$

$$\frac{\partial E}{\partial w_1} = \frac{\partial E}{\partial a_2} \times \frac{\partial a_2}{\partial z_2} \times \frac{\partial z_2}{\partial a_1} \times \frac{\partial a_1}{\partial z_1} \times \frac{\partial z_1}{\partial w_1}$$

$$= \frac{\partial E}{\partial w_2} \times w_2 \times \sigma(z_1)(1 - \sigma(z_1)) * 0.01$$

$$w_{1\text{new}} = w_{1\text{old}} - \eta \frac{\partial E}{\partial w_1}$$