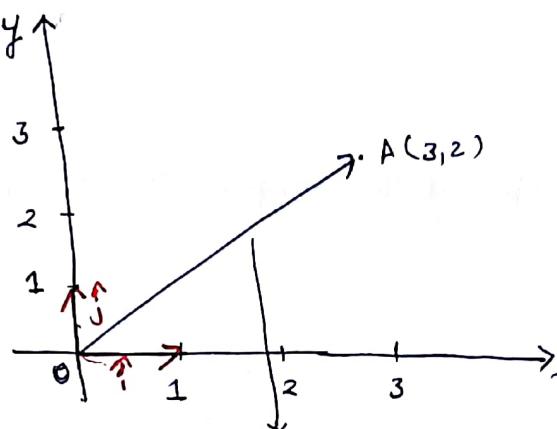


1

Vector :- Magnitude + Direction

Unit Vector  
 ✓ Magnitude + direction  
 ↓ value of magnitude is 1  
 aka directional vectors.



Unit vector in the direction of  $\vec{OA}$  is defined as  $\hat{i}$

If I want to represent a vector in  $\hat{i}$  direction  
 $\vec{OA} = 3\hat{i}$

as they provide direction to a particular magnitude.

Direction to this line

$$\vec{OA} \rightarrow 3\hat{i} + 2\hat{j} = \vec{OA}$$

$$\boxed{\frac{3}{\sqrt{13}}\hat{i} + \frac{2}{\sqrt{13}}\hat{j} = \hat{OA}}$$

Unit vector.

$$\text{magnitude of } \vec{OA} = |\vec{OA}| = \sqrt{3^2 + 2^2} = \sqrt{13}$$

Let's say the unit vector in direction of  $\vec{OA}$  is  $\hat{a}$ .

$$\vec{OA} = |\vec{OA}| \hat{a}$$

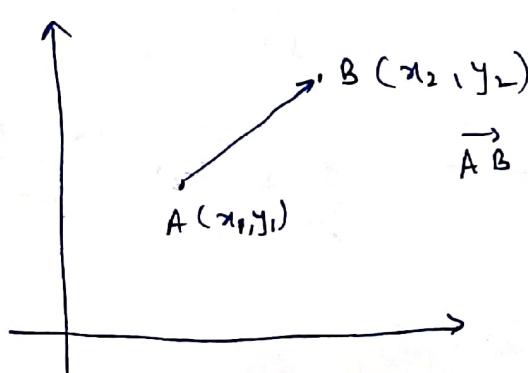
$$\frac{\vec{OA}}{|\vec{OA}|} = \hat{a}$$

→ Condition for two vectors to become equal.

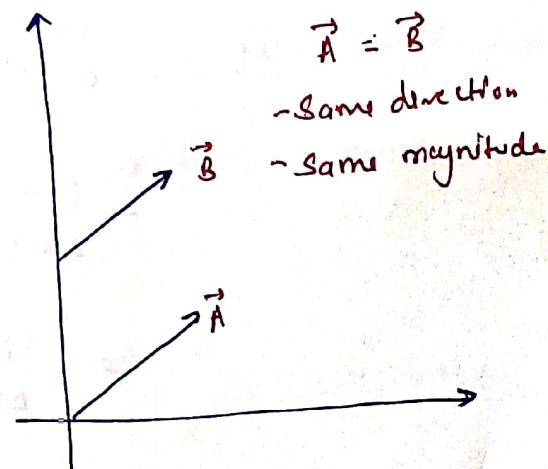
$$\vec{v} = x\hat{i} + y\hat{j} = \begin{bmatrix} x \\ y \end{bmatrix}$$

when magnitude and direction of two vectors are same then it implies those two vectors are equal.

$$\vec{v} = x\hat{i} + y\hat{j} + z\hat{k} = \begin{bmatrix} x \\ y \\ z \end{bmatrix}$$



$$\vec{AB} = (x_2 - x_1)\hat{i} + (y_2 - y_1)\hat{j}$$

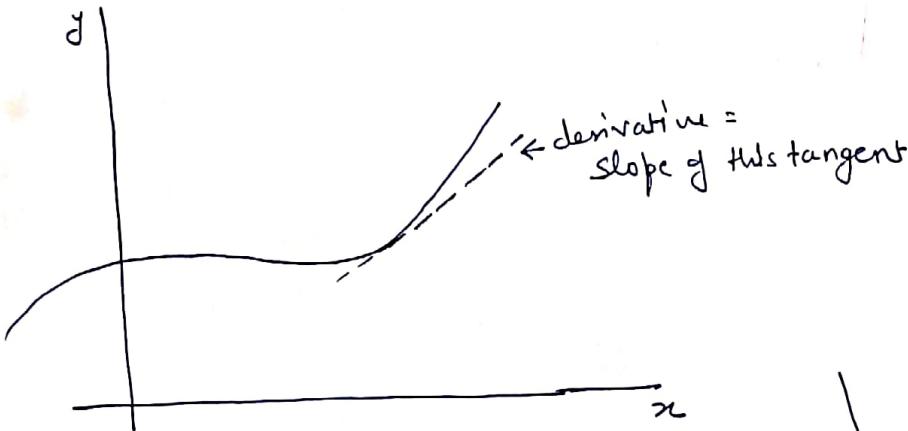


$y = f(x)$   
 ✓ Independent variable  
 dependent variable

(2)

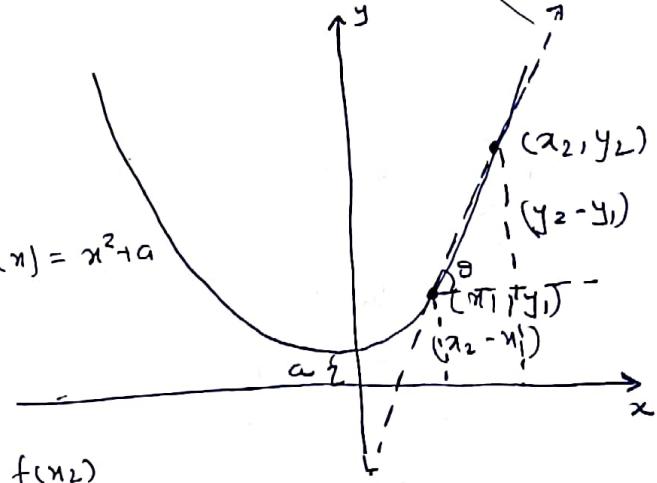
$$\frac{dy}{dx} = \frac{\text{change in } y}{\text{small change in } x}$$

Graphically :- slope of the tangent to the graph at a given point.



$$\text{slope} = \frac{y_2 - y_1}{x_2 - x_1}$$

$$y = f(x) = x^2 + a$$



$$\text{slope} = \frac{\text{change in } y}{\text{small change in } x}$$

$$\boxed{\text{slope} = \lim_{\Delta x \rightarrow 0} \frac{\Delta y}{\Delta x} = \frac{dy}{dx}}$$

↳ definition of differentiation

$$y_1 = f(x_1)$$

$$\text{slope} = \tan \theta = \frac{dy}{dx}$$

$$= \frac{x_2^2 - x_1^2}{x_2 - x_1} = x_1 + x_2$$

$$\frac{dy}{dx} = \lim_{\Delta x \rightarrow 0} \frac{y_2 - y_1}{\Delta x}$$

$$\frac{d}{dx} x^n = n x^{n-1}$$

$$= \lim_{\Delta x \rightarrow 0} \frac{f(x_2) - f(x_1)}{\Delta x}$$

$$\frac{d}{dx} e^x = e^x$$

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

$$\frac{d}{du} \frac{u}{v} = \frac{v - u'v}{v^2}$$

## Partial Differentiation

(3)

$$y = f(x)$$

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

$$\frac{\partial z}{\partial x} = 2x$$

$$\frac{\partial z}{\partial y} = 2y$$

$$\frac{\partial f}{\partial x} = \lim_{\Delta x \rightarrow 0} \frac{f(x + \Delta x, y) - f(x, y)}{\Delta x}$$

$$\frac{\partial f}{\partial y} = \lim_{\Delta y \rightarrow 0} \frac{f(x, y + \Delta y) - f(x, y)}{\Delta y}$$

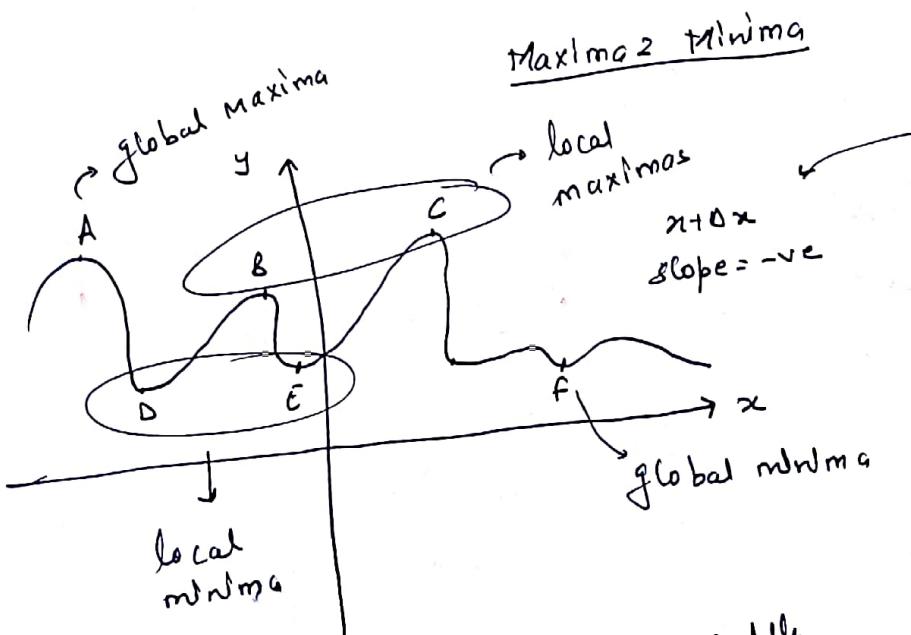
## Use of Partial Differentiation in Deep Learning

$$\text{loss } f^n = \frac{f(\text{weights, biases})}{ }$$

↳ our objective is to minimize this function

$$\text{Loss } f^n = f(w_1, w_2, \dots, w_n, b_1, b_2, \dots, b_n)$$

\* Partial derivative with one weight or bias and find out the change



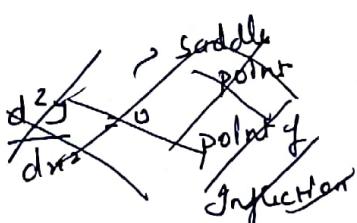
A, B, C = points of maxima.

D, E, F = points of minima  
↓  
 $x + \Delta x$  slope = +ve

①  $\frac{dy}{dx}$  at maxima & minima  $y = 0$

②  $\frac{d^2y}{dx^2} \rightarrow$  double derivative

at  $x = c_1$   
 $\frac{d^2y}{dx^2} > 0 \Rightarrow$  point of minima at  $c_1$



at  $x = c_2$   
 $\frac{d^2y}{dx^2} < 0 \Rightarrow$  point of maxima at  $c_2$

$$f(x) = x^2$$

Critical point  $x=0$

$$\frac{dy}{dx} = 2x = 0$$

$$\frac{d^2y}{dx^2} = 2 > 0 \text{ minima}$$

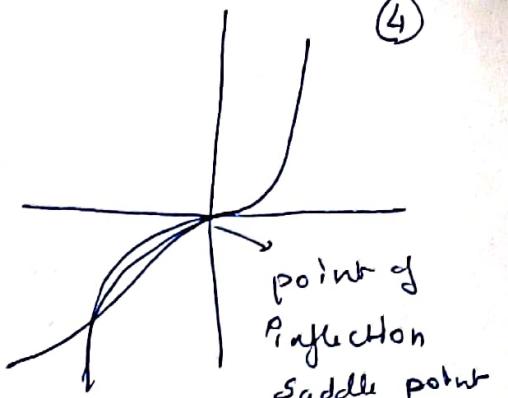
$$f(x) = x^3$$

at  $x=0$

$$\frac{dy}{dx} = 3x^2$$

at  $x=0$

$$\frac{dy}{dx} = 0$$



point of  
inflection  
saddle point

$$\frac{d^2y}{dx^2} = 6x \approx 0 \text{ (at } x=0\text{)}$$

Q Where used in Deep Learning?

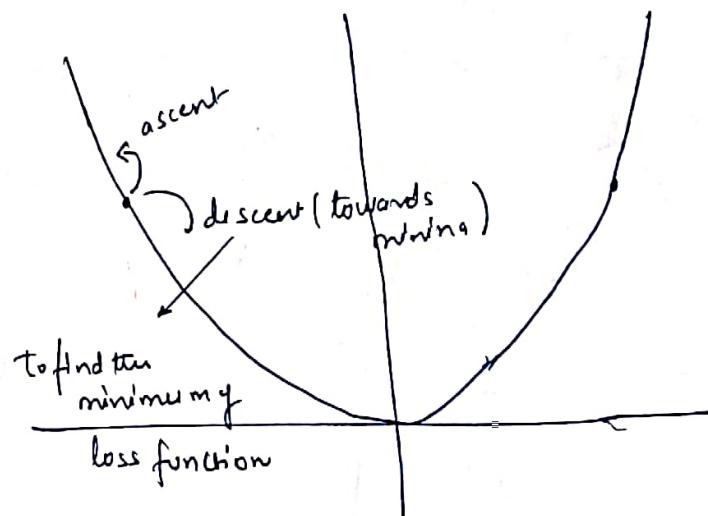
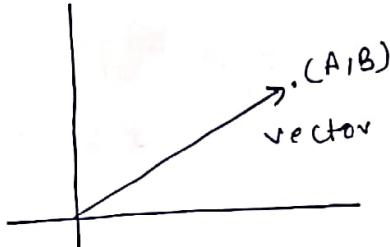
Reduction in loss function  $\rightarrow$  lowest value of the error  $\rightarrow$  global minima

### Gradient Descent

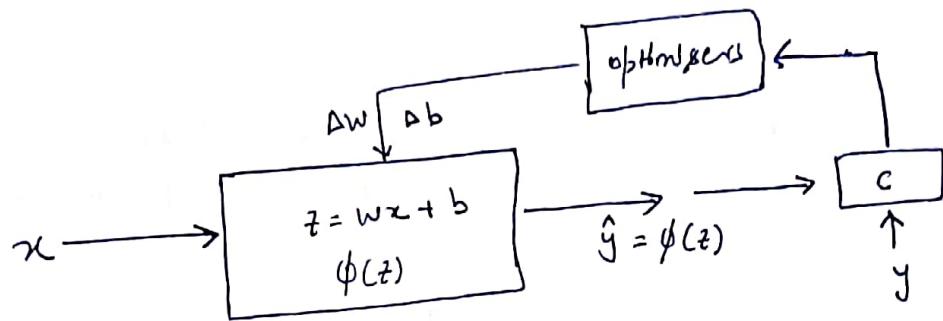
$$z = f(x, y)$$

$$\nabla f(x, y) = \left( \frac{\partial f}{\partial x} \Big|_y \hat{i} + \frac{\partial f}{\partial y} \Big|_x \hat{j} \right)$$

$$= A \hat{i} + B \hat{j}$$



(5)



$$w = w + \Delta w$$

$$b = b + \Delta b$$

$$c = (y - \phi(z))^2$$

$$= (y - \phi(wx+b))^2$$

$c \rightarrow f(w, b) \therefore c$  is nothing but function of weight and biases

Optimizers — Gradient descent

- SGD

- Adams

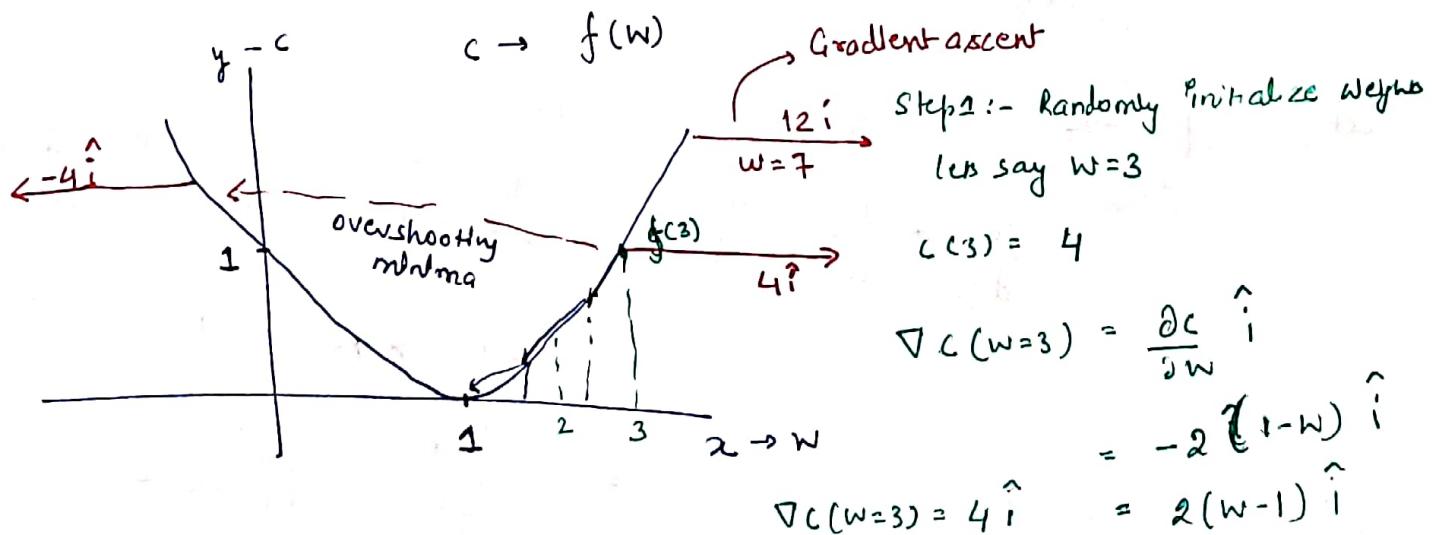
$$c = \cancel{y} (y - \phi(wx+b))^2$$

Assumptions (i) treat y as a constant  $K = 1$

(ii) bias ( $b$ )  $\approx 0$

(iii)  $\phi(wx+b) = \phi(wx) \approx \frac{1}{w}$   $\downarrow$   $\phi \rightarrow$  sigmoid  $f^n$   $0-1$   
some  $\downarrow$  RELU  $0-\infty$

$$\therefore c = (1-w)^2 \quad \text{float value}$$



(6)

## Assumption

(i) Let  $\Delta w = 4$

$$\begin{aligned} w &= w + \Delta w \\ w &= 7 \end{aligned} \quad ] \text{ leads gradient ascent}$$

$$\begin{aligned} \text{(ii)} \quad w &= w - \Delta w & w &= -1 \\ \nabla c(w=-1) &\approx -4^{\uparrow} & & \end{aligned} \quad ] \text{ overshooting minima}$$

∴ we need to use learning rate which is nothing but small number we multiply with gradients. We take the fraction of gradients.

$$\Delta w = -\eta \nabla c \quad \eta = 0.1$$

$(0 \leftrightarrow 1)$

$$w = w - \eta \nabla c = 3 - 0.1 \times 4 = 2.6$$

and every epoch we update the gradient and move towards the minima. This is called gradient descent.

Chain Rule

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

$$\frac{df}{dx}$$

$$\begin{aligned} \text{Let } p &= -x & q &= 1+e^{-x} = 1+e^p \\ \frac{df}{dq} &\times \frac{dq}{dp} \times \frac{dp}{dx} &= \frac{d}{dq} \left( \frac{1}{q} \right) &\times \frac{d}{dp} (1+e^p) \times \frac{d}{dx} e^{-x} \\ &= -\frac{1}{q^2} \times pe^p \times -1 \\ &= \frac{e^p}{q^2} = \frac{e^{-x}}{(1+e^{-x})^2} \end{aligned}$$

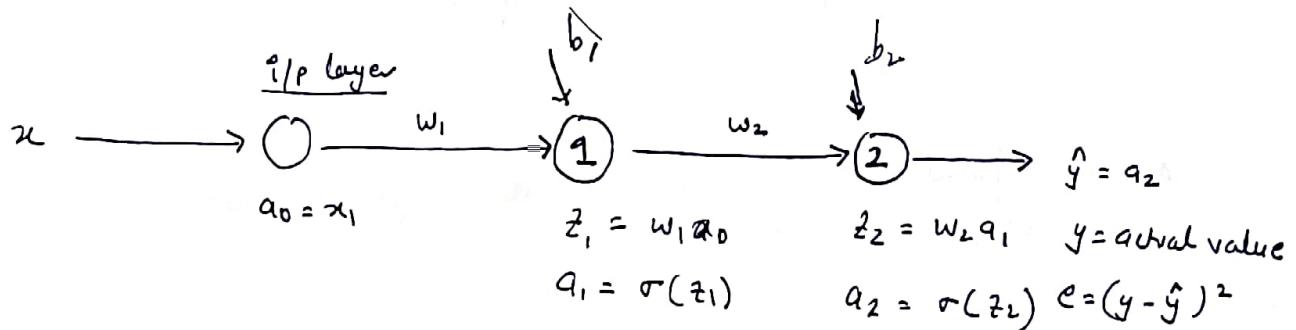
## Backpropagation

(7)

### Assumption

1) Each layer has one neuron

2) bias = 0



Weight update rule

$$w = w - \eta \nabla e = w - \eta \frac{\partial e}{\partial w}$$

Observation  $e = (y - \hat{y})^2 = (y - q_2)^2$

$q_2 = \sigma(z_2)$

$f(q_2)$   $f(z_2)$

$z_2 = f(w_2 q_1)$  ~~Let's assume~~  $\frac{\partial e}{\partial w_L} = \frac{\partial f(f(f(w_2 q_1)))}{\partial w_L}$

$$\frac{\partial e}{\partial w_2} = \frac{\partial e}{\partial q_2} \cdot \frac{\partial q_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial w_2} \Big|_{q_1}$$

$$= -2(y - q_2) \times \sigma'(z_2) [1 - \sigma(z_2)] \times q_1$$

$$w_2 = w_2 - \eta \nabla e = w_2 - \eta \frac{\partial e}{\partial w_2} = \cancel{w_2}$$

(8)

$$q_1 = \underbrace{\sigma(z_1)}_{f(z_1)}$$

$$z_1 = \underbrace{w_1}_{f(w_1, q_0)} \underbrace{a_0}_{f(w_1, q_0)}$$

$$\frac{\partial e}{\partial w_1} = \frac{\partial e}{\partial a_2} \times \frac{\partial a_2}{\partial z_2} \times \left. \frac{\partial z_2}{\partial a_1} \right|_{w_L} \times \left. \frac{\partial a_1}{\partial z_1} \right|_{q_0} \times \left. \frac{\partial z_1}{\partial w_1} \right|_{q_0}$$

?

At some point  $\frac{\partial e}{\partial w_1} \approx 0$ ; no weight update

Assumption update, bias  $\neq 0$

$$\frac{\partial e}{\partial w_L} = \frac{\partial e}{\partial a_L} \times \frac{\partial a_L}{\partial z_L} \times \boxed{\frac{\partial z_L}{\partial w_L}}$$

$$\frac{\partial e}{\partial b_2} = \frac{\partial e}{\partial a_L} \times \frac{\partial a_L}{\partial z_L} \times \boxed{\frac{\partial z_L}{\partial b_2}}$$

$$z_L = w_L a_L + b_L$$

$$\frac{\partial z_L}{\partial w_L} = a_L \quad \frac{\partial z_L}{\partial b_L} = 1$$

## \* General Problems in training Neural Network

(9)

- 1) Vanishing and Exploding Gradient
- 2) Lot of data required for training
- 3) Slow training issues
- 4) Risk of overfitting

$$w_{new} = w_{old} - \eta \nabla C$$

↳ lets say the value of this is decreasing and is very minute, weight update will become non significant

$\Delta w \approx 0$  ]  $\rightarrow$  vanishing gradient issue

If this gradient becomes very huge (in this case

$\Delta w \uparrow \uparrow$  (very high)  $\rightarrow$

We may diverge from the soln

Exploding gradient.

- \* Neural network will overfit in case of small datasets. In order to avoid this we tend to provide a huge dataset.

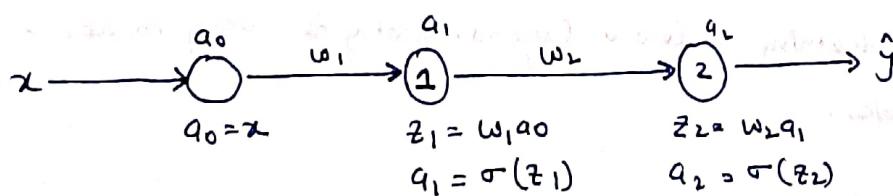
Increase the size of Neural Networks

↳ Increase in number of hidden layers

↳ Increases the size of the model.

Slow training issues

## Backpropagation Equation Brief



(10)

$$\frac{\partial e}{\partial w_2} = \frac{\partial e}{\partial a_2} \cdot \frac{\partial a_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial w_2} \quad \text{--- (1)}$$

$$\frac{\partial e}{\partial w_1} = \frac{\partial e}{\partial a_2} + \frac{\partial e}{\partial z_2} \cdot \frac{\partial z_2}{\partial a_1} + \frac{\partial a_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_1} \quad \text{--- (2)}$$

Observations: eq(1) and eq(2) is nothing but product of Ratios. These ratios can become greater than 1 or they can be between 0 < 1.

1st Case

$$0 < \text{ratio} < 1$$

eq(1)

$$\frac{\partial e}{\partial w_2} = \underbrace{\frac{\partial e}{\partial a_2}}_{0.1} \cdot \underbrace{\frac{\partial a_2}{\partial z_2}}_{0.2} \cdot \underbrace{\frac{\partial z_2}{\partial w_2}}_{0.8} \approx 0.006$$

$$w_2 = w_1 - \eta^{0.006} \quad \text{lets } \eta = 0.1$$

$$w_2 = w_1 - 0.1(0.006)$$

$$\boxed{w_2 = w_1 - 0.0006} \quad \text{almost negligible change}$$

$\Rightarrow$  There will be no weight update significantly

\* This is known as vanishing gradient problem, that means very minute or negligible weight update.

Assumption:- All the ratios in  $\frac{\partial e}{\partial w_2}$  equation(2) are b/w 0 and 1.

$$\frac{\partial e}{\partial w_2} \approx \text{negligible}$$

This also leads to slow training of lower layers, owing to very minor < negligible weight updates.

## Observation 2<sup>nd</sup>

(11)

If all the ratio are greater  $\gg 1$

Q10

$$\frac{\partial e}{\partial w_2} = \frac{\partial e}{\partial a_2} \times \frac{\partial a_2}{\partial z_2} \times \frac{\partial z_2}{\partial w_2} \approx 6000$$

$\overline{10} \quad \overline{20} \quad \overline{30}$

$$w_2 = w_2 - \eta 6000 \quad \text{very large weight update}$$

Similarly it will be further greater in case of eq(2)  
solution will diverge

This problem was first observed in 2010, by Xavier Glorot & Yoshua Bengio  
in their paper "Understanding the difficulty of training deep feed forward  
neural networks"

why vanishing and exploding gradient problem occurs?

- 1) choice of activation function  $\frac{\partial a_2}{\partial z_2}$  } derivative of activation function
- 2) Weight initialization technique  $\frac{\partial z_2}{\partial w_2}$  }

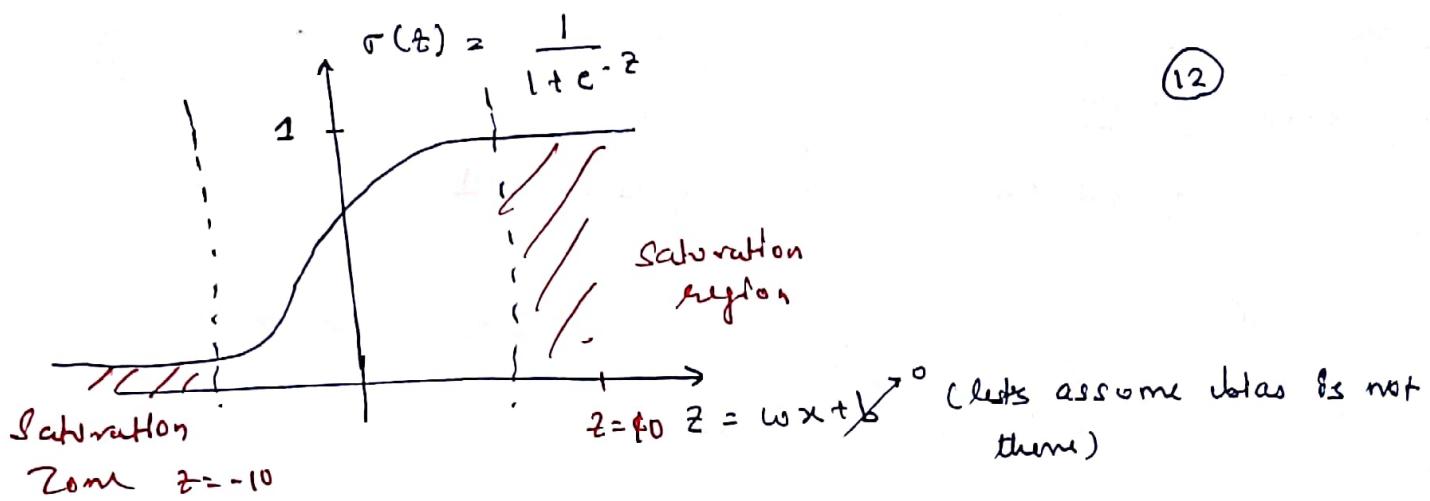
Ans

### Activation function

It always non-linearity and help us to converge to the solution, i.e.  
it helps non-linearity and help us to converge to the solution, i.e.  
help us find out the minima of loss function.

$$\frac{\partial e}{\partial w} = \frac{\partial e}{\partial a} \times \left( \frac{\partial a}{\partial z} \right) \times \frac{\partial z}{\partial w}$$

↓ this will help convergence



Case 1:-

Initialised weight  
 $w = 10$ , assume  $x = 1$

$$z = wx = 10 \cdot 1 = 10$$

$$e \approx 2.73$$

$$\sigma(z=10) = \frac{1}{1 + e^{-10}} \approx 0$$

$$= \frac{1}{1+0} = 1$$

$$\text{at point } z = 10 \quad \frac{\partial \sigma(z)}{\partial z} = \sigma(z) \{1 - \sigma(z)\}$$

$$\frac{\partial \sigma}{\partial z}(z=10) = 1 \times (1-1) = 0$$

$$\frac{\partial e}{\partial w} = \frac{\partial e}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial w} \approx 0$$

No weight update/negligible weight update

Case 2 :-

$$z = -10 \quad w = -10 \quad x = 1$$

$$\sigma(z = -10) = \frac{1}{1 + e^{10}}, y = 0$$

$$\frac{\partial \sigma(z)}{\partial z} = \sigma(z) \{1 - \sigma(z)\} = 0$$

$$\frac{\partial e}{\partial w} = \frac{\partial e}{\partial a} \cdot \cancel{\frac{\partial a}{\partial z}} \cdot \frac{\partial z}{\partial w}$$

i.e if we initialize weight very high or very low in case of sigmoid function results in vanishing gradient. (l3)

Case 3

$$z=0 \quad | \quad w=0 \quad x=1$$

$$\sigma(z=0) = \frac{1}{1+e^0} = 0.5$$

$$\frac{\partial \sigma(z)}{\partial z} = \sigma(z) \{1 - \sigma(z)\} = 0.5 \times 0.5 = 0.25$$

$$\frac{\partial e}{\partial w} = \frac{\partial e}{\partial a} \times \frac{\partial a}{\partial z} \times \frac{\partial z}{\partial w}$$

$$= 0.25 \left[ \frac{\partial e}{\partial a} \times \frac{\partial z}{\partial w} \right]$$

∴ there will be significant weight update.

⇒ Through the above example we understood that both activation function and weight initialization contribute to vanishing and exploding gradient

⇒ Previously we used to randomly initialize the weights.

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

$$\frac{\partial u}{\partial v} = \frac{u'v + v'u}{v^2}$$

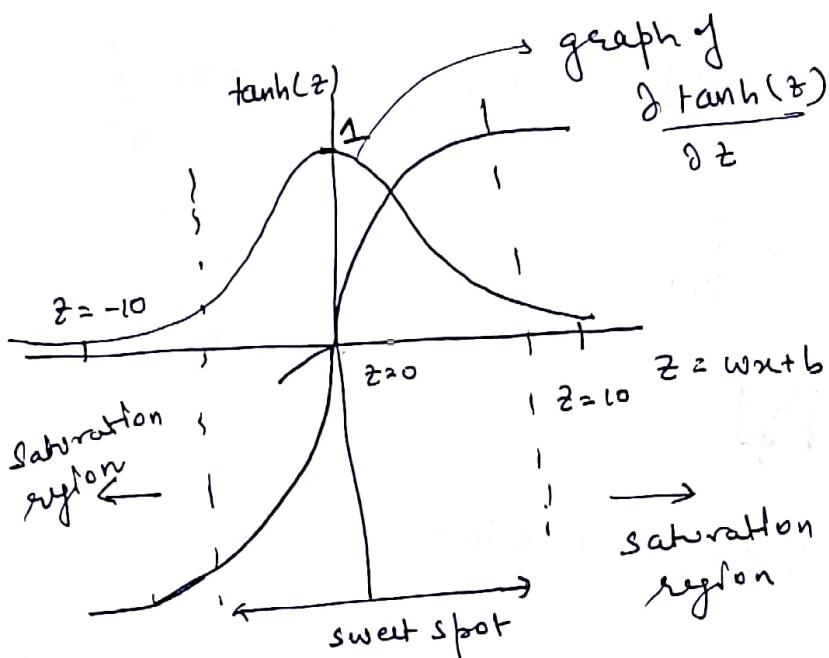
$$\frac{\partial \tanh(z)}{\partial z} = \frac{(e^z + e^{-z})(e^z + e^{-z}) - (e^z - e^{-z})(e^z - e^{-z})}{(e^z + e^{-z})^2}$$

$$= \frac{(e^z + e^{-z})^2 - (e^z - e^{-z})^2}{(e^z + e^{-z})^2} = 1 - \left( \frac{e^z - e^{-z}}{e^z + e^{-z}} \right)^2 = 1 - (\tanh(z))^2$$

(14)

$$\text{ReLU}(z) = \begin{cases} 0 & z < 0 \\ z, & z \geq 0 \end{cases}$$

$$\boxed{\frac{\partial \text{ReLU}(z)}{\partial z} = \begin{cases} 0 & z < 0 \\ 1 & z \geq 0 \end{cases}}$$



Case 1 Let  $w = 10$   $x = 1$   
 $b = 0$   
 weight initilize

$$\tanh(z=10) = \frac{e^{10} - e^{-10}}{e^{10} + e^{-10}} \approx 1$$

$$\left. \frac{\partial \tanh(z=10)}{\partial z} \right|_{z=10} = 1 - (\tanh(z=10))^2 \approx 0$$

hence issue of vanishing  
gradient

Case 2  $z = -10$

$$\tanh(z=-10) = \frac{e^{-10} - e^{10}}{e^{-10} + e^{10}} \approx -1$$

$$\left. \frac{\partial \tanh(z=-10)}{\partial z} \right|_{z=-10} = 1 - (\tanh(z=-10))^2 \approx 1$$

Case 3  $w = 0$   $b = 0$   $x = 1$

$$\tanh(z=0) = \frac{e^0 - e^0}{e^0 + e^0} \approx 0$$

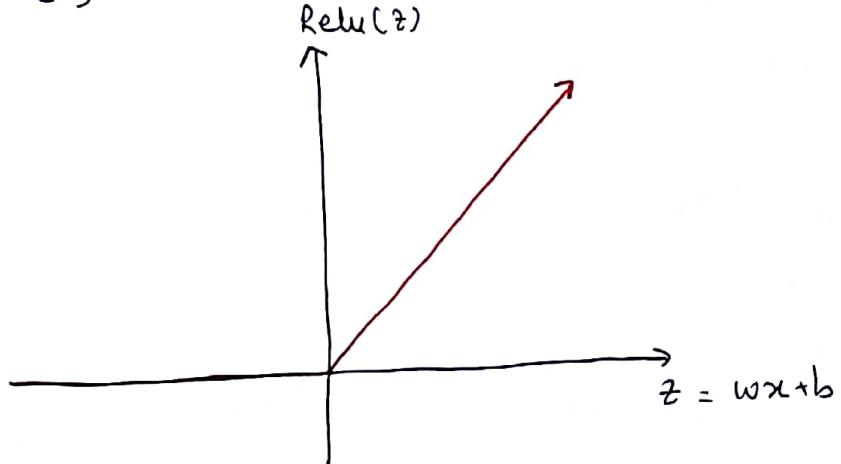
$$\left. \frac{\partial \tanh(z=0)}{\partial z} \right|_{z=0} = 1 - (0)^2 = 1$$

$$\frac{\partial e}{\partial w} = \frac{\partial e}{\partial a} \times \frac{\partial a}{\partial z} \times \frac{\partial z}{\partial w}$$

In weight update will be significant

In forward propagation information flows through  $\tanh(z)$  in backward propagation information flows through  $\frac{\partial \tanh(z)}{\partial z}$  (derivative of  $\tanh(z)$ )

Relu(z)



$$\text{Case 1} \quad w=10 \quad x=1 \quad 15$$

$$\frac{\partial \text{Relu}(z=10)}{\partial z} \Big|_{z=10} = 1$$

significant changes in weight

Case 2

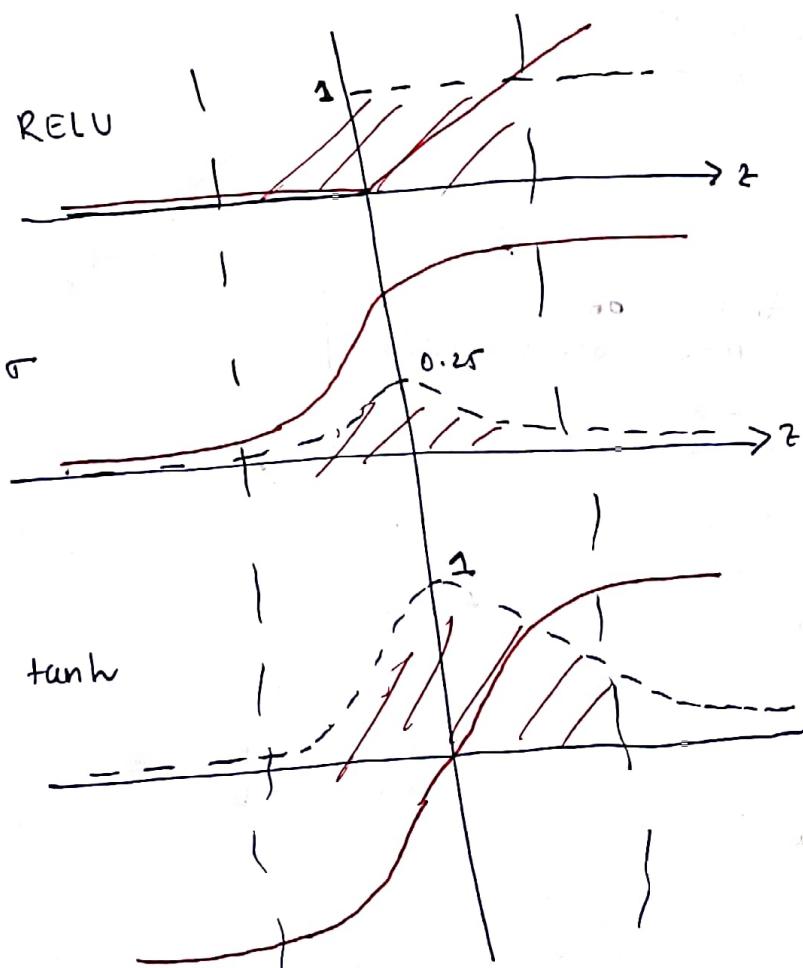
$$w < 0 \quad x = 1 \quad z < 0$$

$$\frac{\partial \text{Relu}(z)}{\partial z} \Big|_{z<0} = 0$$

- Vanishing gradient problem

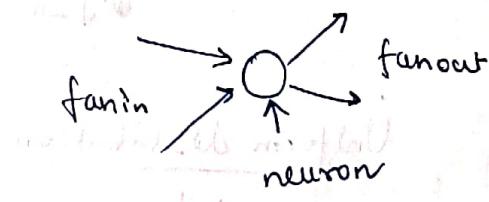
\* Observation made by the research paper by Xavier Glorot & Yoshua Bengio

-- derivative backward  
--- forward



fanin = Number of incoming edges to a layer

fanout = Number of outgoing edges from a layer



To prevent vanishing and exploding gradient:

1)  $\text{fanin} = \text{fanout}$  but not always possible to have same input & output edges.

alternative proposal:- if  $\sigma_{in}^2 = \sigma_{out}^2$  (variance of weight of input edges is equal to variance of weights of output edges)

(1b)  $\sigma_{in}^2 = \sigma_{out}^2$  : helps you maintain forward & backward information flow

Calculate

$$\text{fan}_{\text{avg}} = \frac{\text{fan}_{\text{in}} + \text{fan}_{\text{out}}}{2}$$

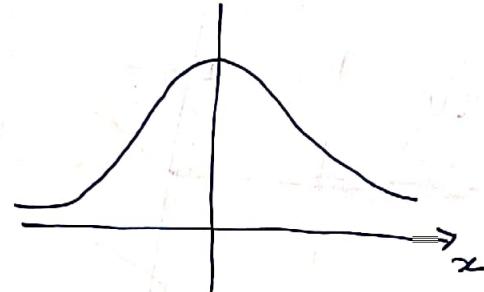
for sigmoid activation fn

- Normal distribution with  $\mu=0$  &  $\sigma = 1/\text{fan}_{\text{avg}}$
- Uniform distribution b/w  $-r \leq z \leq r$  where  $r = \sqrt{\frac{3}{\text{fan}_{\text{avg}}}}$

Glorot Initialization [or Xavier Initialization] default case of weight Initialization in Keras.

Normal distribution

$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



standard Normal distribution

$$\mu=0, \sigma=1 \text{ (std dev)}$$

$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{x^2}{2}}$$

Uniform distribution

$$a, b$$

$$f(x) = \frac{1}{b-a}$$



## Type of Initialization

Xavier Glorot

+

Youhua Bengio

Activation function

Glorot  $\rightarrow$  None, tanh, sigmoid

softmax.

$$\sigma^2(\text{normal}) = \frac{1}{\text{variance}} \text{ fanin}$$

17

Kaiming He

He  $\rightarrow$

RELU and its  $\sigma^2 = 2/\text{fanin}$   
variance  
normal dist.

Xiaonan Lecan

Lecan - SELU

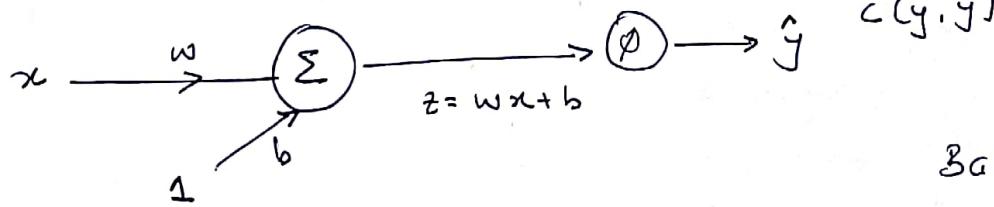
$\sigma^2 = 1/\text{fanin}$   
normal dist

If  $\text{fanin} = \text{fanout} \Rightarrow$  Glorot and Lecan Initialization are same

In tensorflow

tf.keras.layers.Dense(units, activation function, kernel\_initializer = "he-normal")

By default its glorot Initialization.



Forward Pass  $\hat{y} = \phi(z)$

when  $z = wx + b$

Sigmoid

domain  
 $(-\infty, \infty)$

range  
 $(0, 1)$

Backward Pass

$$\frac{\partial c}{\partial w} = \frac{\partial c}{\partial \phi} \times \frac{\partial \phi}{\partial z} \times \frac{\partial z}{\partial w}$$

domain range  
 $(-\infty, \infty)$   $(0, 1)$

For forward pass:

use  $\phi(z)$  and consider its domain & range

for backward pass:  
use  $\phi'(z)$  consider  
its domain and range

## 1. Sigmoid activation function

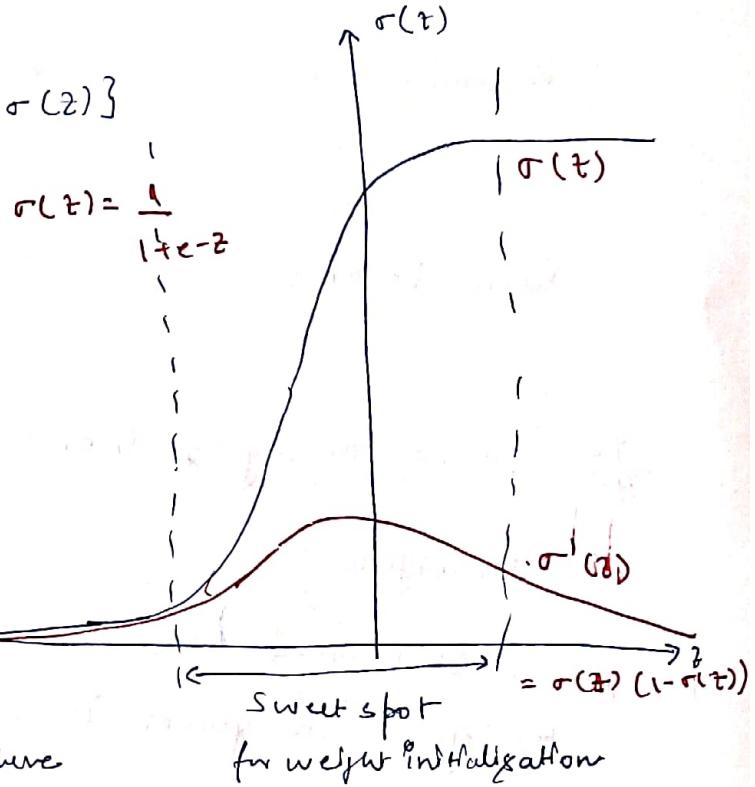
$$\sigma(z) = \frac{1}{1+e^{-z}}$$

Range  $\sigma(z) \in (0,1)$   
 domain  $z \in (-\infty, \infty)$

### Backward

$$\sigma'(z) = \sigma(z) \{ 1 - \sigma(z) \}$$

range  $\in (0,0.25)$   
 domain  $\in (-\infty, \infty)$



### Advantages

- 1) It has smooth gradient that means this prevents jumps in the output value. {continuous function} i.e. it's derivable everywhere
- 2) Output values are range of  $\sigma(z)$  b/w (0,1) this implies normalising the output of each neuron. This will help us to reduce the solution space.



### Disadvantages

- 1) It is prone to gradient vanishing, as  $\underline{\sigma'(z) \in (0,0.25)}$
- 2) ~~fraction~~ Power operation on exponential terms makes calculations time consuming i.e. it will simply increase the latency.

Mostly used in output layer

## 2) Hyperbolic Tangent activation function

(19)

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

where

$\tanh(z) \in (-1, 1)$  Range

Domain  $(-\infty, \infty)$

$$\tanh'(z) = 1 - \tanh^2(z)$$

range  $\rightarrow \tanh'(z) \in (0, 1)$

domain  $\rightarrow (-\infty, \infty)$

→ provided proper weight initialisation

\* Better than sigmoid activation function as it has lesser chance of vanishing gradient

\* Smooth gradient - derivable at every point

\* Symmetric to origin - which leads to zero mean - it will provide equal balance to negative and positive side

\* Suitable function for hidden layers.

### Disadvantages

- \* It has saturation region where gradient becomes zero.
- \* It introduces latency because of exponential terms.



## 3) Rectified Linear Unit (ReLU), most popular, highly used, most successful

### Forward Pass

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

domain  $z \in (-\infty, \infty)$  Range  $\text{ReLU}(z) \in [0, \infty)$

$$= \begin{cases} z & z \geq 0 \\ 0 & z < 0 \end{cases}$$

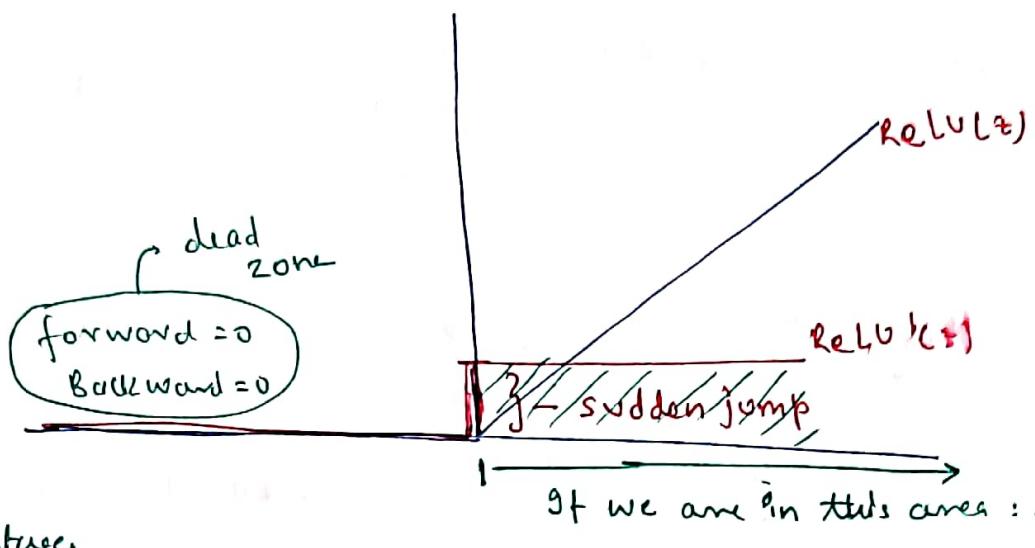
Issue with derivative as this function is not continuous at 0.

$$\text{ReLU}'(z) = \begin{cases} 1 & z \geq 0 \\ 0 & z < 0 \end{cases} \rightarrow \text{derivative is not defined at } z=0$$

Range ~~0, 1~~  $\in [0, 1]$

domain  $(-\infty, 0) \cup (0, \infty)$

From



### Advantages

- ① If input is +ve  $\Rightarrow$  There is no gradient saturation problem.
- ② Calculation is very fast compared to prev activation function both in forward and backward direction

### Disadvantages

- ① For negative inputs ReLU is completely inactive  $\Rightarrow$  Dying ReLU issue.
- ② Not a zero centric function. At zero, its derivative is not defined  $\hookrightarrow$  jump / spike

To resolve the issue of ReLU not being derivable, people have developed variants of ReLU.

### ④ Leaky ReLU functions. (Resolve dead zone issue)

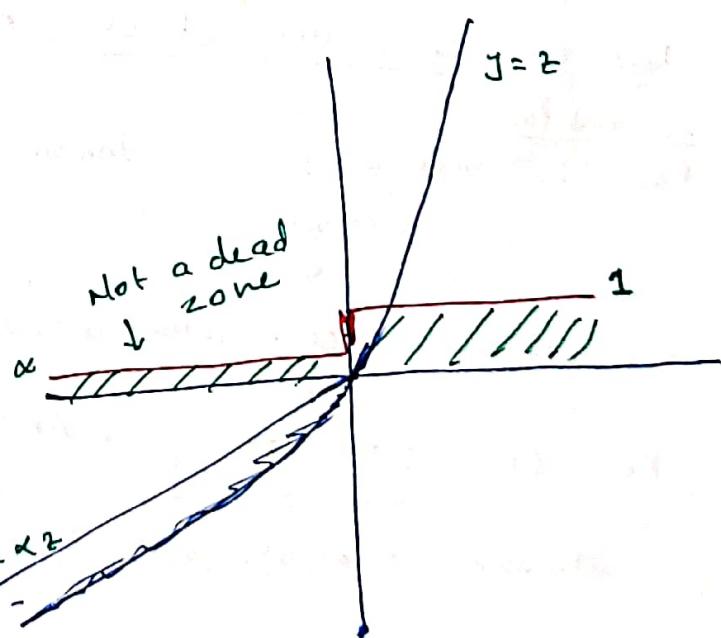
#### Forward

$$\text{Leaky ReLU}(z) = \max(z, \alpha z)$$

$$\text{Leaky ReLU}(z) = \begin{cases} z & z > 0 \\ \alpha z & z \leq 0 \end{cases}$$

domain .  $z \in (-\infty, \infty)$

range  $\in (-\infty, \infty)$   
 $\alpha \approx 0.01$



#### Backward Pass / derivative

$$\text{Leaky ReLU}'(z) = \begin{cases} 1 & z > 0 \\ \alpha & z \leq 0 \end{cases}$$

not derivable at  $z = 0$

domain  $z \in (-\infty, \infty)$

range  $\in \alpha, 1$

### (5) P.ReLU : Parametric ReLU

$$\text{PReLU} = \begin{cases} z & z \geq 0 \\ \alpha z & z < 0 \end{cases}$$

If  $\alpha = 0 \Rightarrow \text{ReLU}$

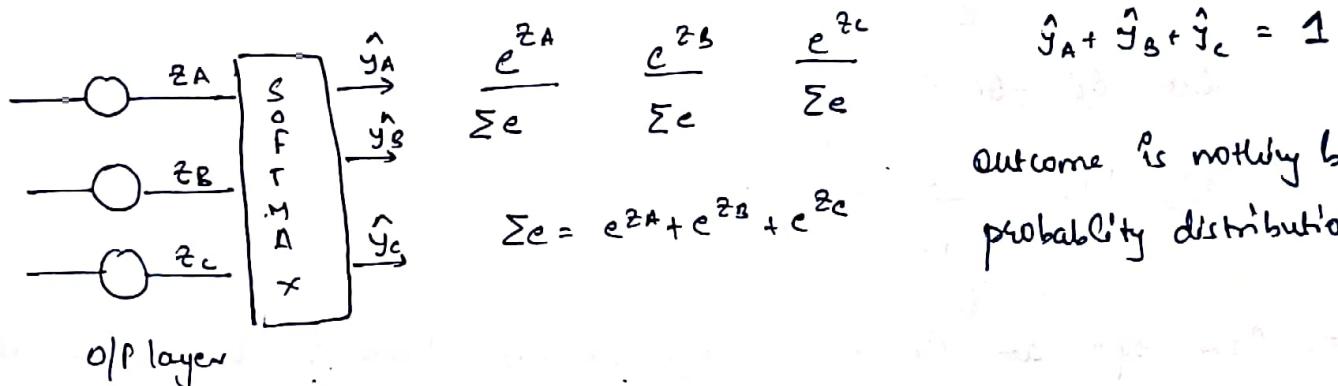
$\alpha > 0 \Rightarrow \text{Leaky ReLU}$

$\alpha$  is learnable  $\rightarrow$  P.ReLU

$\boxed{\alpha} \rightarrow$  learnable parameter  
 ↳ this will also be trained.

\* Disadv. Need to train extra parameter.

### (6) Softmax Activation Functions



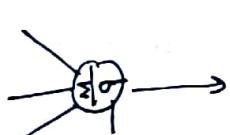
outcome is nothing but probability distribution

If number of classes  $> 2$

#### Case Study

Case 1 : binary classification  
 {sigmoid fn}

$$z = w_1 x_1 + w_2 x_2 + w_3 x_3 = \theta^T x$$



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

Sigmoid  
activation  
fn for output  
layer

$$p(y=1|x) = \frac{1}{1 + e^{-z}}$$

probability of  $y=1$

given the value of  
 $x$

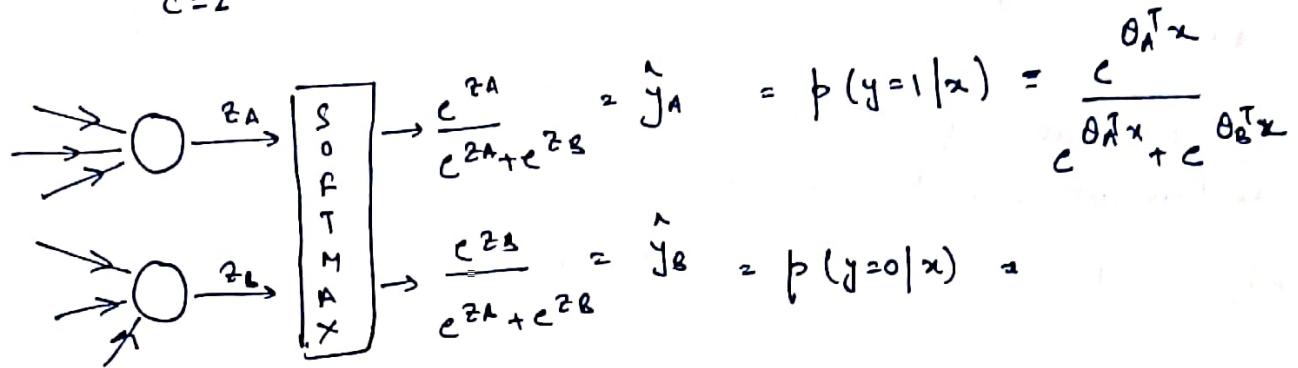
$$= \frac{1}{1 + e^{-\theta^T x}}$$

$$p(y=0|x) = 1 - p(y=1|x) = \frac{e^{-\theta^T x}}{1 + e^{-\theta^T x}}$$

Case 2 :- Binary Classification on  $f^n$ , softmax

(22)

$$c=2$$



$$p(y=1|x) = \frac{e^{\theta_A^T x}}{e^{\theta_A^T x} + e^{\theta_B^T x}} = \frac{1}{1 + e^{(\theta_B^T - \theta_A^T) \cdot x}} = \frac{1}{1 + e^{-\theta_T^T x}}$$

$$\text{Let } \theta_B^T - \theta_A^T = -\theta^T$$

$$p(y=0|x) = \frac{e^{-\theta^T x}}{1 + e^{-\theta^T x}}$$

$\Rightarrow$  These eqn are similar for sigmoid case. Hence for binary classification softmax is equivalent to sigmoid activation.

- \* Only drawback of softmax is that it contains exponential terms which simply implies that it's computationally intensive
- \* Softmax is only used in output function.

$\Rightarrow$  \* Remember unsolved issue of ReLU not derivable at  $x=0$

### ⑦ ELU (Exponential Linear Units)

Forward Pass

$$\text{elu}(z, \alpha) = \begin{cases} z & z \geq 0 \\ \alpha(e^z - 1) & z < 0 \end{cases}$$

$\alpha$  - scalar value and slope of -ve section

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

Range  $\in (-\alpha, \infty)$

backward Pass

$$\text{elu}'(z, \alpha) = \begin{cases} 1 & z \geq 0 \\ \alpha e^z & z < 0 \end{cases}$$

domain  $z \in (-\infty, \infty)$

range  $\in (0, 1]$

1) No dead zone

2) Derivable at zero.

3) Disadv. It involves exponential computation

(8) SELU

$$\text{selu}(z) = \begin{cases} z & z \geq 0 \\ \alpha e^z - \alpha & z < 0 \end{cases}$$

learnable parameter

Similar sigmoid.

domain:  $[-\infty, \infty]$

range:  $(-\infty, \infty)$

(9) SWISH

$$\text{swish}(z) = z \cdot \sigma(\beta \cdot z) = \frac{z}{1 + e^{-\beta z}}$$

(8) SELU - Scaled ELU

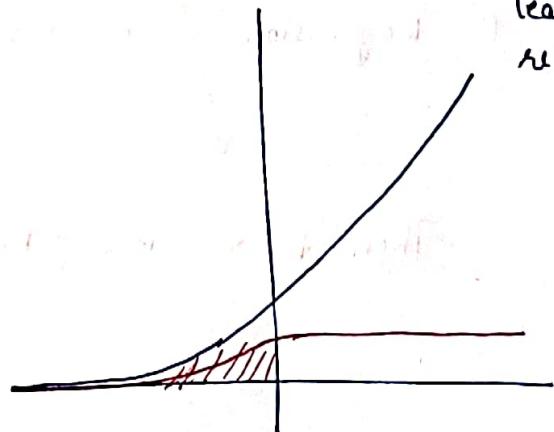
23

$$\text{selu}(z, \alpha) = \text{scale} \times \text{elu}(z, \alpha) = k \text{elu}(z, \alpha)$$

(10) SOFTPLUS

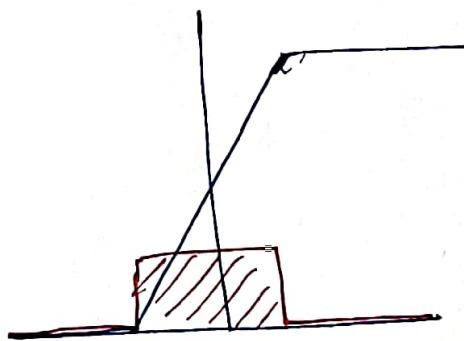
$$\text{softplus}(z) = \log(1 + e^z)$$

just like  
leaky  
relu



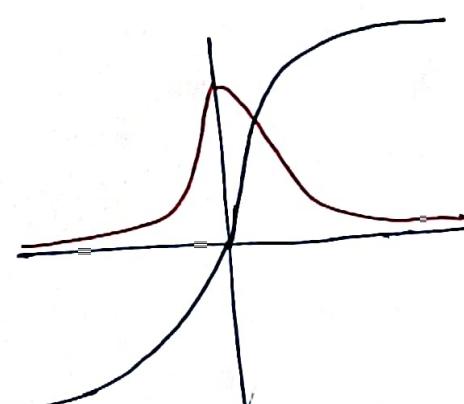
(11) Hard Sigmoid

$$\phi(z) = \begin{cases} 0 & z < -2.5 \\ 1 & z > 2.5 \\ 0.2z + 0.5 & -2.5 \leq z \leq 2.5 \end{cases}$$



(12) Soft Sign

$$\text{ss}(z) = \frac{z}{|z| + 1} = \begin{cases} \frac{z}{z+1} & z > 0 \\ \frac{z}{1-z} & z < 0 \end{cases}$$



(13) Max Out

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

drawback: 1) Compute intensive as it doubles the number of parameters

Rough idea on how to begin selection of Activation function

For any problem - start with **ReLU** for hidden layers

for classification

binary - two choices at output

Sigmoid  
softmax

- multiclass - only softmax at output

for Regression

No activation function (We are expecting continuous outcome)

Then try versions of ReLU if not getting better results.