

Activation Function

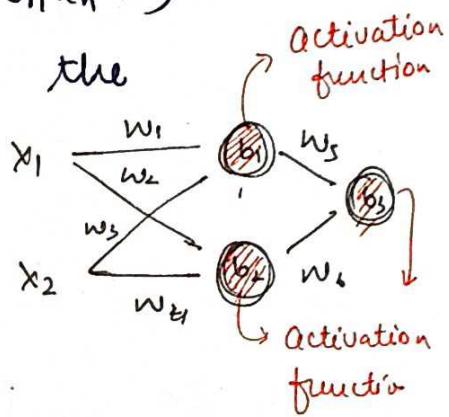
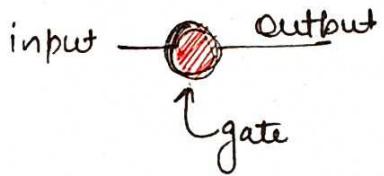
What are Activation functions?

In artificial neural Networks, each neuron forms a weighted sum of its inputs and passes the resulting scalar value through a function referred to as an activation then the output or activation of a neuron is

$$a = g(w_1x_1 + w_2x_2 + w_3x_3 + \dots + w_nx_n + b)$$

The function g is referred to as the activation function.

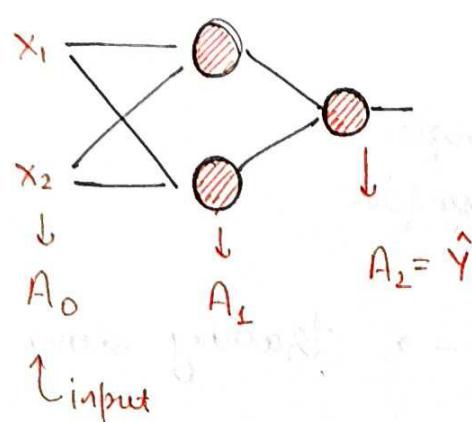
basically, Activation function is a mathematically gate



If the function g is taken as the linear regression function $g(z) = z$ then the neuron perform linear regression or classification. In general g is taken to be a non-linear function if the function g to do non-linear regression and solve classification problems that are not linearly separable.

Why activation functions are needed?

Forward propagation



$$z_1 = w_1 A_0 + b_1$$

$$A_1 = g(z_1) = z_1$$

We are not using any activation function. So, our output is same (z_1).

$$A_2 = g(w_2 A_1 + b_2)$$

$$= w_2 A_1 + b_2 = w_2 (w_1 A_0 + b_1) + b_2$$

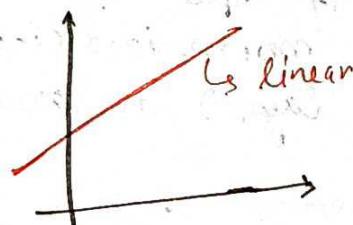
$$= \underbrace{w_2 w_1}_{w'} A_0 + \underbrace{w_2 b_1 + b_2}_{b'}$$

Output

Input

$$A_2 = \hat{y} = w' A_0 + b'$$

one degree polynomial



Ideal Activation function

1.) Non-linear. \rightarrow Good activation functions are non-linear.

$y = f(\alpha) \rightarrow$ relation betn y and α should be non-linear

for eg:- Sigmoid $\rightarrow \sigma(z) = \frac{1}{1+e^{-z}}$

* Universal Approx? \rightarrow With the help of enough nodes in a layer we can capture any type of non-linear pattern.

2) Differentiable \rightarrow Activation function should be differentiable

3) Computationally inexpensive \rightarrow derivative \rightarrow simple
easy fast

If computation expensive \rightarrow training slow

4) Zero centered \rightarrow output of activation function
should be

Zero center

Normalize

mean = 0

example \rightarrow tanh

Neural Network
fasterly converge
if input data is
normalized (input
layer) or later layer

5) Non-Saturating

Saturating means to squeeze the input in particular range. For eg:- Sigmoid $\rightarrow [0, 1]$, tanh $\rightarrow [-1, 1]$

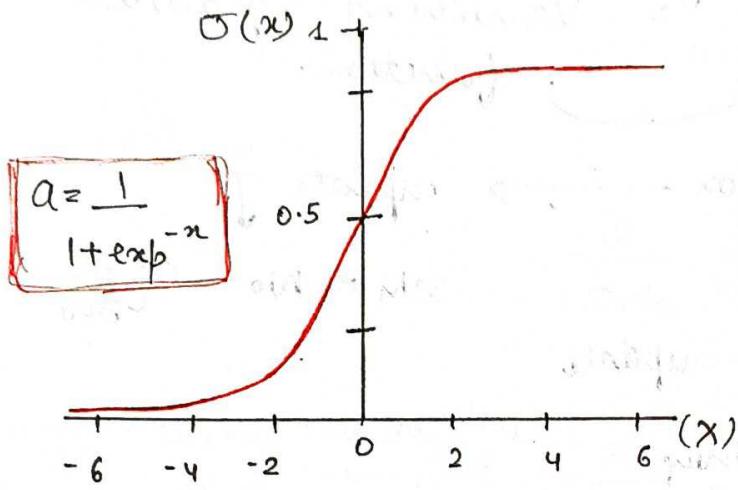
Saturating function fail \rightarrow Vanishing Gradient problem

Non-Saturating \rightarrow Relu $\rightarrow (0, \infty)$

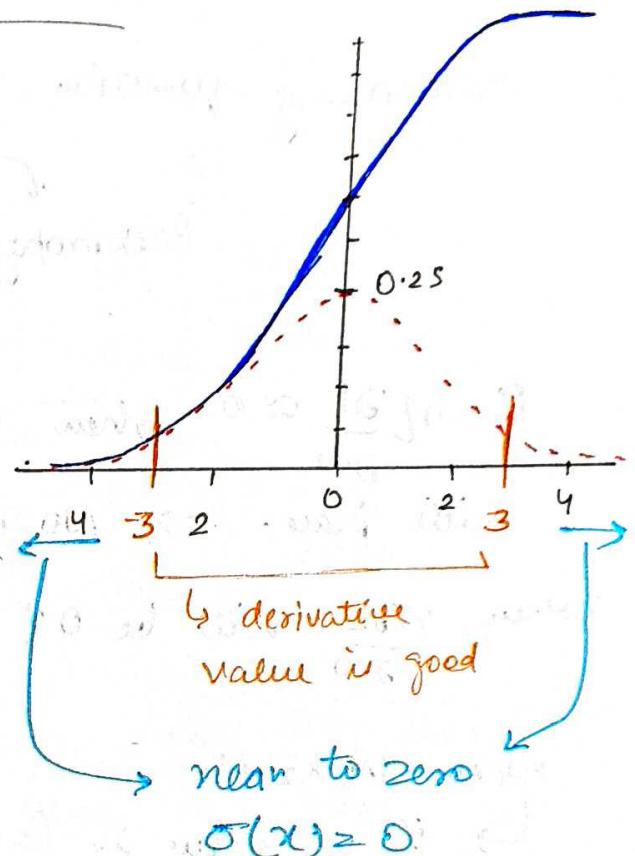
No range

Sigmoid Activation Function

(50)



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



Advantages

- 1) Output of sigmoid is betⁿ $(0, 1)$

So, we can treat as a probability and use at output layer.

use in Binary classification

- 2) Non-linear \rightarrow If your data is non-linear so, sigmoid function can capture non-linearity. (good option).

- 3) Differentiable \rightarrow Backpropagation

$$\frac{\partial L}{\partial w}$$

Disadvantage

1) Saturating function $\xrightarrow{\text{facing}}$ Vanishing gradient function.

Backpropagation \longrightarrow update \downarrow

$$w_n = w_0 - \eta \frac{\partial L}{\partial w}$$

If $\eta \frac{\partial L}{\partial w} \approx 0$ then No update

will place. \rightarrow NO Training

when $\frac{\partial L}{\partial w}$ will be 0? $\rightarrow \alpha \uparrow \uparrow$ (large)

$$w_1x_1 + w_2x_2 + b$$

\hookrightarrow if w value is large then n value is also large
 when we go to other layers sigmoid output will be 0. (Layer ke piche jate jayenge output chhoti hoti jaygi. fir ek time pe 0 ho jayega.)

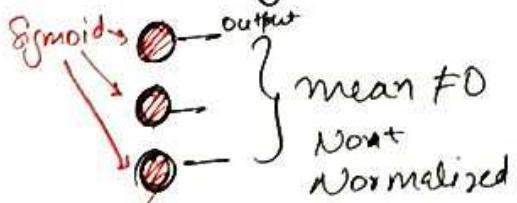
New weight == Old weight

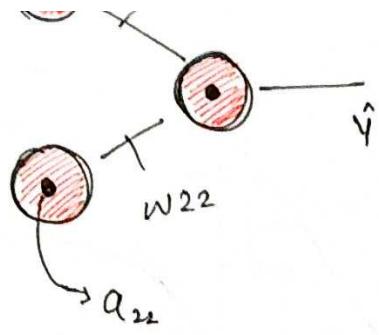
That's why sigmoid not use in hidden layer only use in output layer.

2) Non-zero centered

After applying activation function (Sigmoid)^{all values} is not normalized and mean $\neq 0$.

Training is slow bcz of not normalized and mean $\neq 0$





$$\sigma(z) \leftarrow \begin{cases} w_{21} \\ a_{21} \end{cases}$$

$$\frac{\partial L}{\partial w_{21}} = \boxed{\frac{\partial L}{\partial \hat{y}} \quad \frac{\partial \hat{y}}{\partial z_{31}}} \quad \frac{\partial z_{31}}{\partial w_{21}}$$

Same ↕

$$\frac{\partial z_{31}}{\partial w_{21}} = a_{21} w_{21} + a_{22} w_{22} + b_{31}$$

↓ ↓ ↓

1 0 1

$$\frac{\partial L}{\partial w_{22}} = \boxed{\frac{\partial L}{\partial \hat{y}} \quad \frac{\partial \hat{y}}{\partial z_{31}}} \quad \frac{\partial z_{31}}{\partial w_{22}}$$

$$\frac{\partial z_{31}}{\partial w_{21}} = a_{21}$$

$$\frac{\partial z_{31}}{\partial w_{22}} = a_{22}$$

Output of
Sigmoid
(0 → 1) → +ve

$$\frac{\partial L}{\partial w_{21}} = \boxed{\frac{\partial L}{\partial \hat{y}} \quad \frac{\partial \hat{y}}{\partial z_{31}}} \quad (a_{21}) \rightarrow +ve$$

$$\frac{\partial L}{\partial w_{22}} = \boxed{\frac{\partial L}{\partial \hat{y}} \quad \frac{\partial \hat{y}}{\partial z_{31}}} \quad (a_{22}) \rightarrow +ve$$

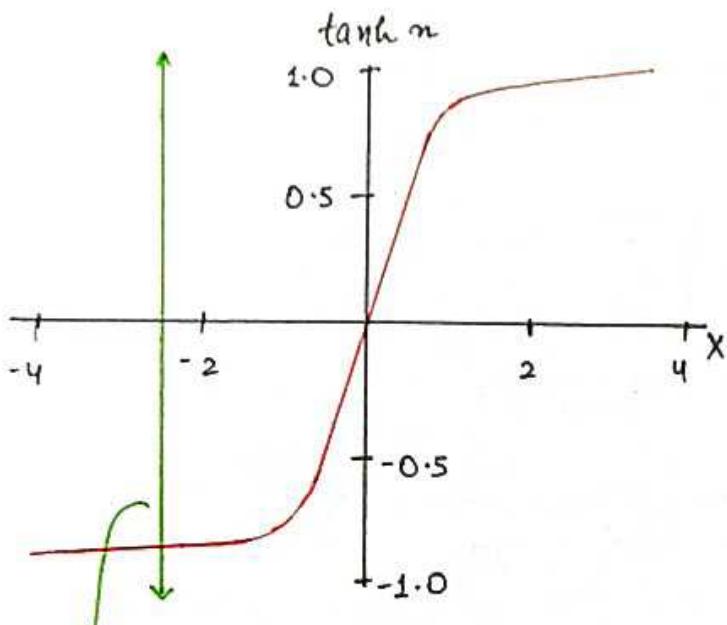
- Weight's derivation same (+ve) charge
- \hat{y}_a same (-ve) charge.
- (1 positive or -ve) net charge.

That's mean
this is restricted

if ① is (-ve) then ② also (-ve)
if ① is positive then ② also (+ve)
(2 both are same.)

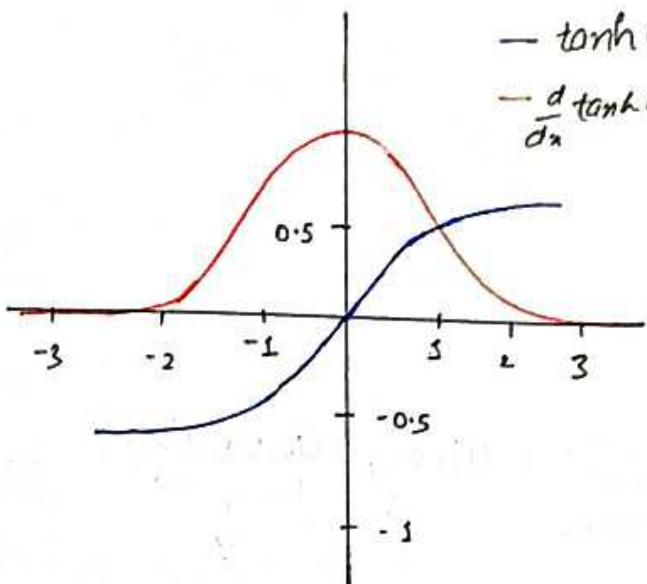
3) Sigmoid formula $\rightarrow \frac{1}{1+e^{-x}} \rightarrow$ computationally expensive

Tanh Activation Function



[$-1, 1$] formula

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



↳ after derivative

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

Advantage

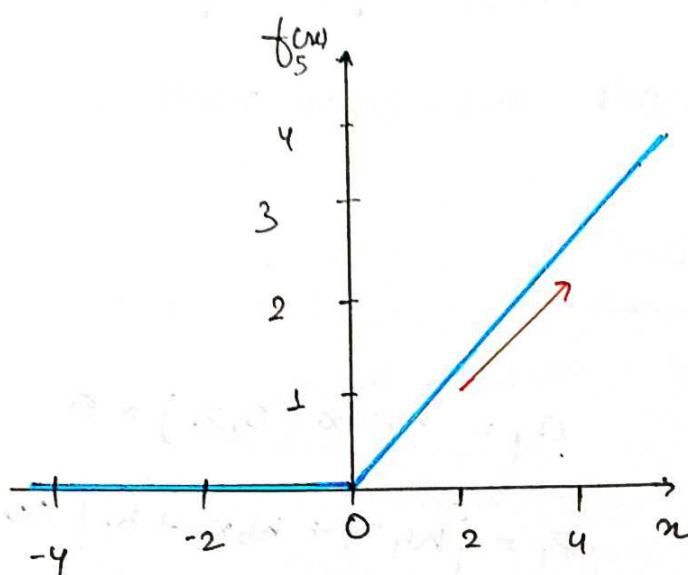
- 1) Non-linear
- 2) Differentiable
- 3) Zero-centered \rightarrow gradient is (ve) and (ve) too
↳ training faster than Sigmoid

Disadvantage

1. Saturating function \rightarrow vanishing gradient problem
2. Computational expensive

ReLU Activation function

(52)



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

Advantage:

- 1) Non-linear
- 2) Not saturated in the positive region.
Positive region
kitne bhi age ja sakte hai. No limit
- 3) Computationally inexpensive
- 4) Convergence \rightarrow faster than Sigmoid and tanh.

Disadvantage:

- 1) Differentiate \rightarrow not zero.

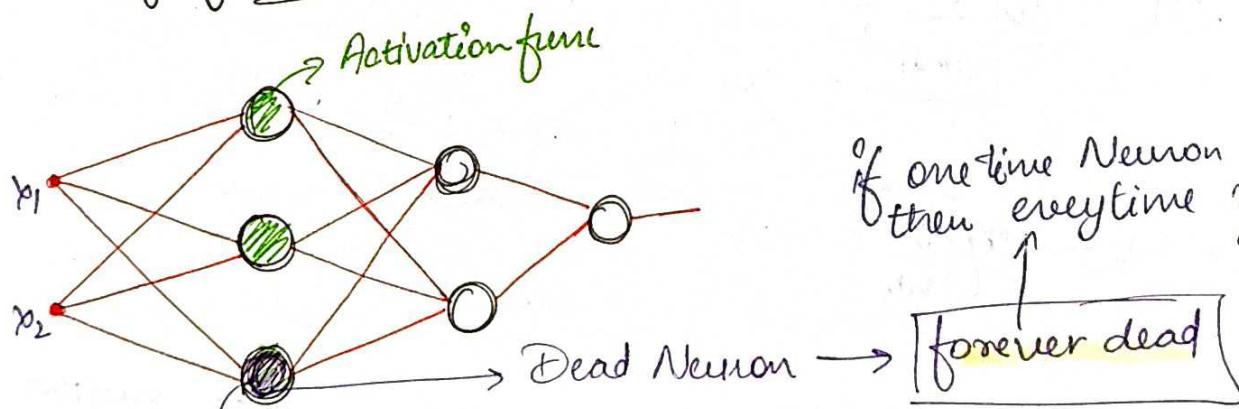
So assumption for this

$$x < 0 \rightarrow 0$$

$$x > 0 \rightarrow 1$$

- 2) Not zero centered

Dying ReLU Problem



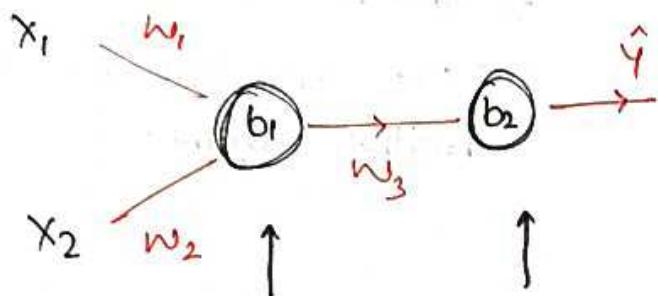
if one time Neuron is dead
then everytime neuron is dead
↑
[forever dead]

And Not depend
on input bcz you can give any value still output is 0.

If 50% or more than 50% of neurons^{or node} are dead then model or architecture cannot capture pattern.

Worst case → 100% of nodes are dead neurons

Why dying ReLU Problem happen?



$$a_1 = \max(0, z_1) \quad \text{if } z_1 > 0$$

$$z_1 = [w_1 x_1 + w_2 x_2 + b_1] \quad \text{if } z_1 \leq 0$$

Let assume

$$z_1 = [w_1 x_1 + w_2 x_2 + b_1] \leq 0$$

then

$$a_1 = \max(0, z_1) = 0$$

→ If a_1 is 0. Then w_1 and w_2 are not update.

$$w_1 = w_1 - \eta \frac{\partial L}{\partial w_1}$$

if a_1 is 0 then
derivative of a_1
w.r.t z_1 is 0.

$$\frac{\partial a_1}{\partial z_1} = 0$$

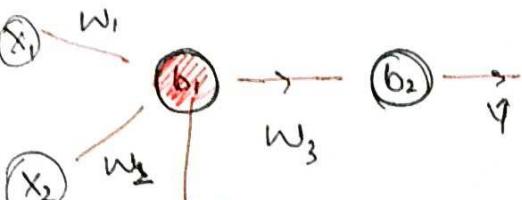
Use this term to
solve this both
derivation.

$$w_2 = w_2 - \eta \frac{\partial L}{\partial w_2}$$

$$w_1 = w_1 - \eta (0) \longrightarrow$$

$w_1 = w_1 \quad \} \quad$ both weight
are same.
 $w_2 = w_2 \quad \} \quad$ No change
in weight

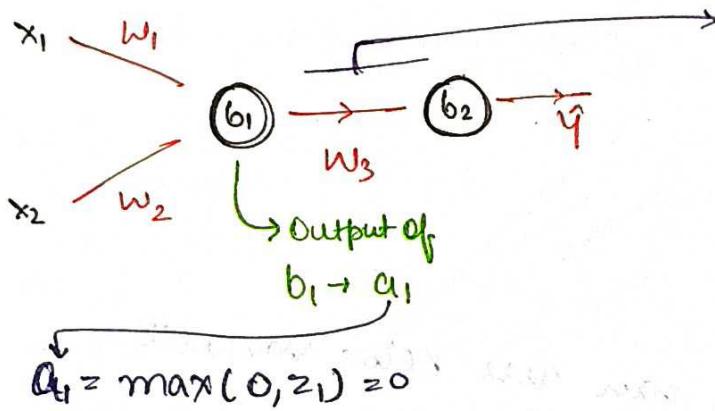
$$w_2 = w_2 - \eta (0) \longrightarrow$$



dead node
(w_1 and w_2 always same)

→ let assume neural network is regression
for prove that $\frac{\partial a_1}{\partial z_1}$ term is use to

Solve $\frac{\partial L}{\partial w_1}$ and $\frac{\partial L}{\partial w_2}$



$$z_1 = w_1 x_1 + w_2 x_2 + b_1$$

$$\text{new } w_1 = \text{old } w_1 - \eta \left(\frac{\partial L}{\partial w_1} \right)$$

$$\text{new } w_2 = \text{old } w_2 - \eta \left(\frac{\partial L}{\partial w_2} \right)$$

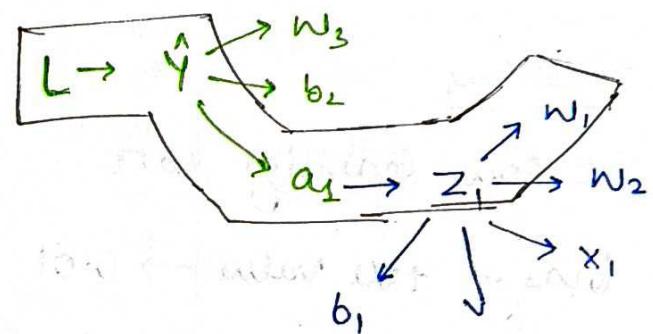
$w_{\text{new}_1} = w_{\text{old}_1} \rightarrow \text{Same}$

$w_{\text{new}_2} = w_{\text{old}_2} \rightarrow \text{Same}$

When is this term $z_1 = w_1 x_1 + w_2 x_2 + b_1$ negative?

Reason → 1. High learning rate

regression Neural Network



$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial a_1} \frac{\partial a_1}{\partial z_1} \frac{\partial z_1}{\partial w_1}$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial a_1} \frac{\partial a_1}{\partial z_1} \frac{\partial z_1}{\partial w_2}$$

$$w_1 = w_1 - n \frac{\partial L}{\partial w_1}$$

any number

Small Number

Large Number

Very high positive number $\approx 10^9$

$w_1 = (-ve) \rightarrow$ Subtract small number - large number.

In next cycle z_1 will be negative $\approx -10^9$ if w_1 is negative.

2) High Negative Bias

$$z_1 = w_1 x_1 + w_2 x_2 + b_1 \rightarrow \text{high -ve}$$

$z_1 = \text{also negative}$

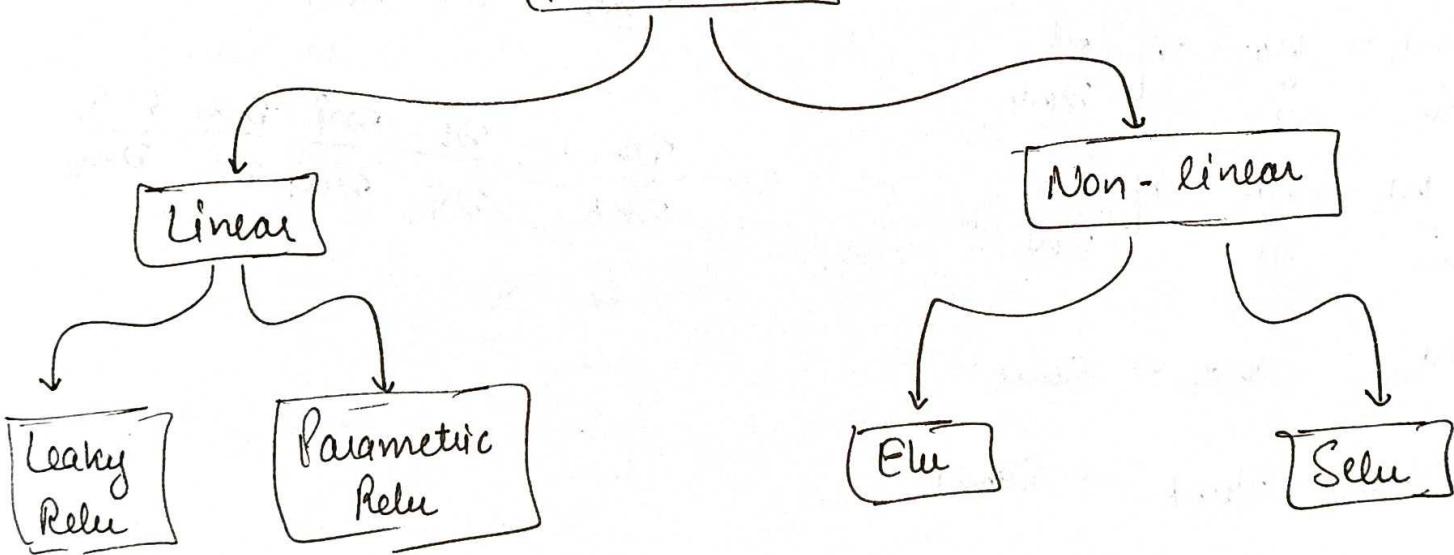
Solutions

→ Set low learning rate

→ bias → +ve value → 0.01

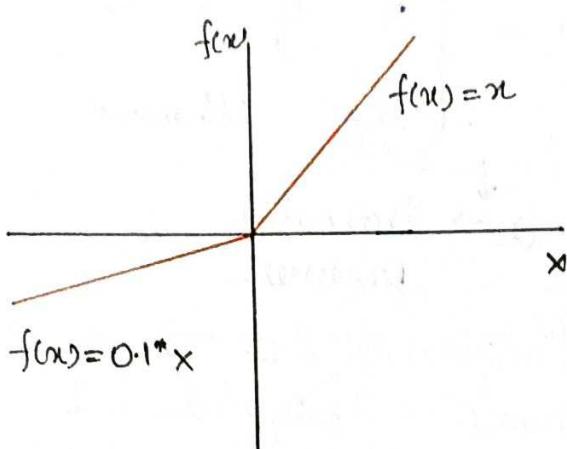
→ Don't use value → You can use ReLU variant.

ReLU Variant



Leaky Relu

(54)



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

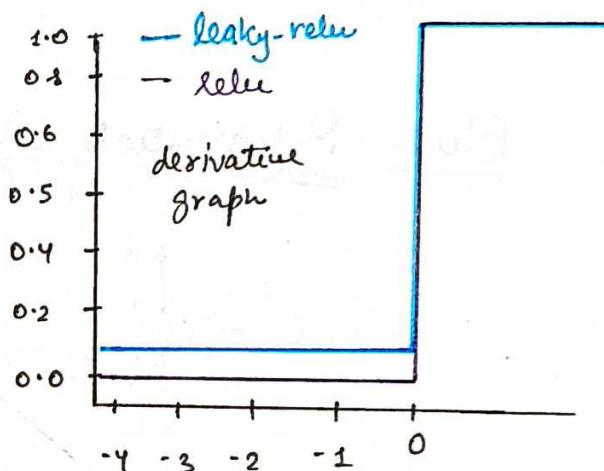
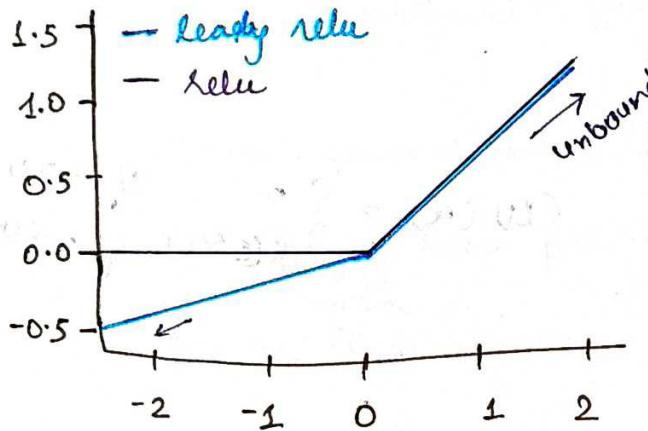
$$z \geq 0 \rightarrow z$$

$$z < 0 \rightarrow \frac{1}{100}z \text{ (fraction of } z\text{)}$$

when (z) \rightarrow -ve then derivative is not 0. it's 0.01
so, in system small gradient flow

$$\downarrow w_i = w_i - \eta \left(\frac{\partial L}{\partial w_i} \right) \rightarrow \text{not } 0 \rightarrow 0.01$$

So, not facing any dying ReLU problem.



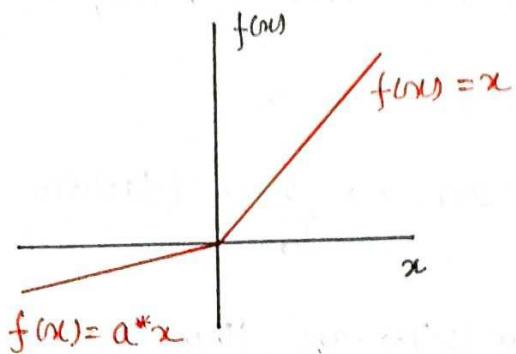
Advantages

- (i) Non-saturated (both direction \rightarrow unbounded (not restriction))
- (ii) Easily computed
- (iii) No dying ReLU problem
- (iv) Close to center \rightarrow (z) both type of values present
-ve +ve

why we take 0.01 in formula?

↳ will discuss in parametric ReLU.

Parametric Relu



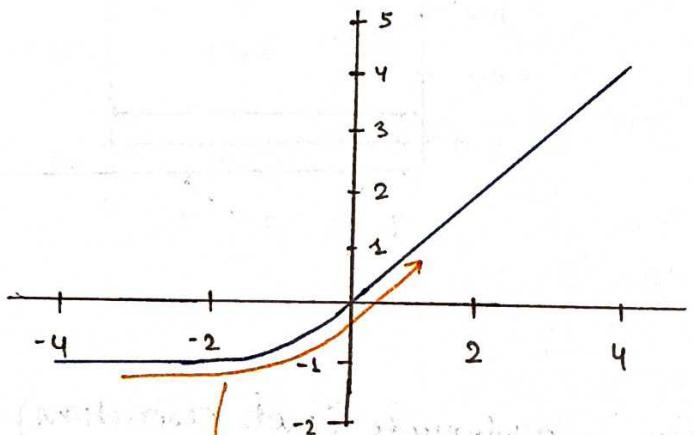
$$f(x) = \begin{cases} ax & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$

a → trainable parameter
0.01 in leaky relu

Parametric ReLU is exactly same → Leaky relu. Only difference is leaky relu and parametric relu

in negative side $f(x) = 0.01 * x$
in negative side $f(x) = a * x$

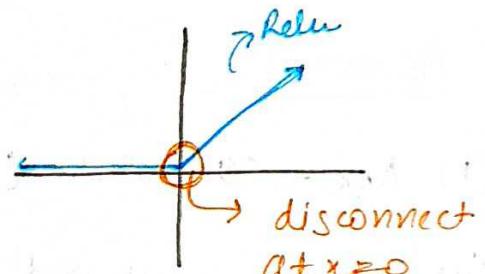
Elu - Exponential Linear Unit



Continuous (Not disconnect at $x=0$) like Relu

$$\text{ELU}(x) = \begin{cases} x & \text{if } x > 0 \\ a(e^{x-1}) & \text{if } x \leq 0 \end{cases}$$

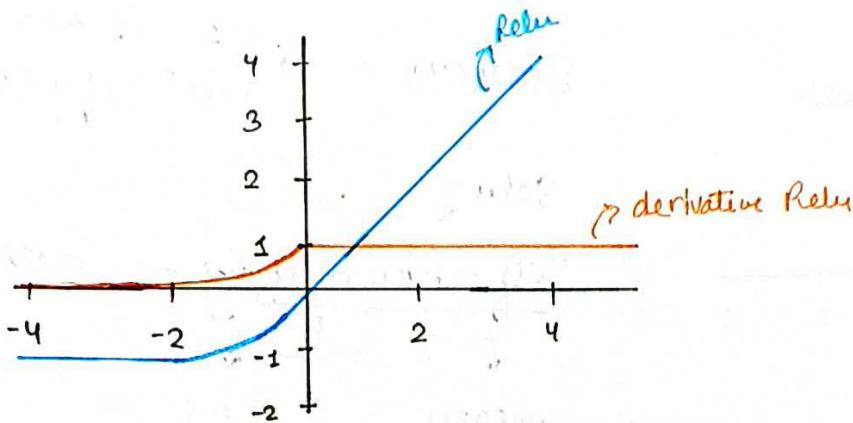
Constant (0.1 to 0.5)



→ Elu →

- ↳ always continuous
- ↳ always differentiable

Differentiable formulas of ELU \rightarrow



$$\text{ELU}'(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ \text{ELU}(x) + \alpha & \text{if } x \leq 0 \end{cases}$$

for positive
for negative

(55)

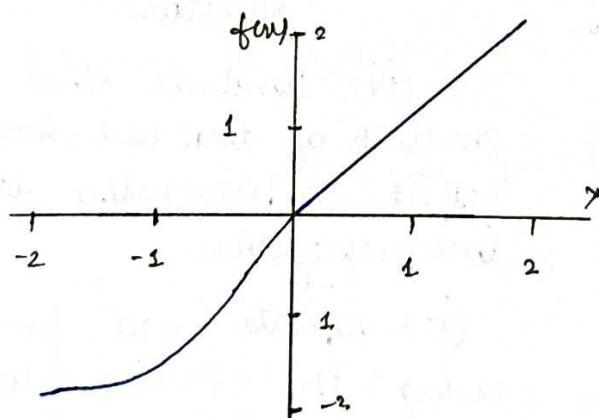
Advantage

- Close to zero centered
↳ convergence faster
- Better generalized (Good in training)
- No Dying ReLU problem
- Always continuous and always differentiable

Disadvantage

- Computation expensive

SELU - Scaled Exponential Linear Unit

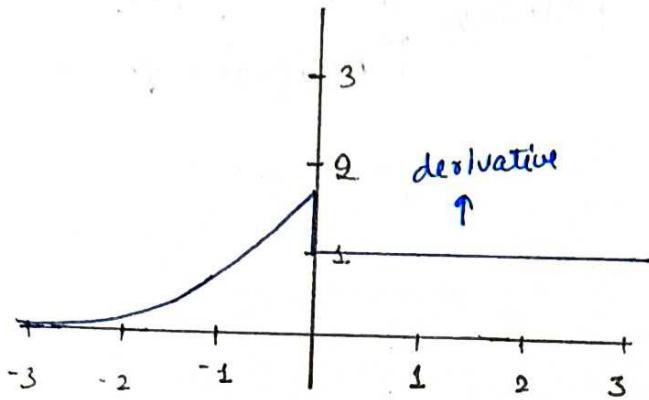


$$\text{SELU}(x) = \lambda \begin{cases} x & \text{if } x \geq 0 \\ \lambda e^x - \lambda & \text{if } x \leq 0 \end{cases}$$

experimental fix parameter

$$\lambda = 1.673263242354377$$

$$\alpha = 1.0507009873554$$



$$SFL'(x) = \lambda \begin{cases} 1 & \text{if } x > 0 \\ \lambda e^x & \text{if } x \leq 0 \end{cases}$$

Self

Self-Normalizing \rightarrow activation

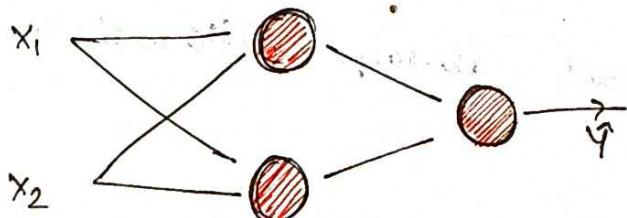
meazo

$$\sigma = 1$$

convergence faster

Weight Initialization

why weight initialization is important?



You will face problem if initial parameter is wrong.

→ vanishing gradient problem

→ exploding gradient problem

→ slow convergence

1. Initialize the parameter

2. Choose an optimization algorithm.

3. Repeat these steps:

(i) Forward propagate an input.

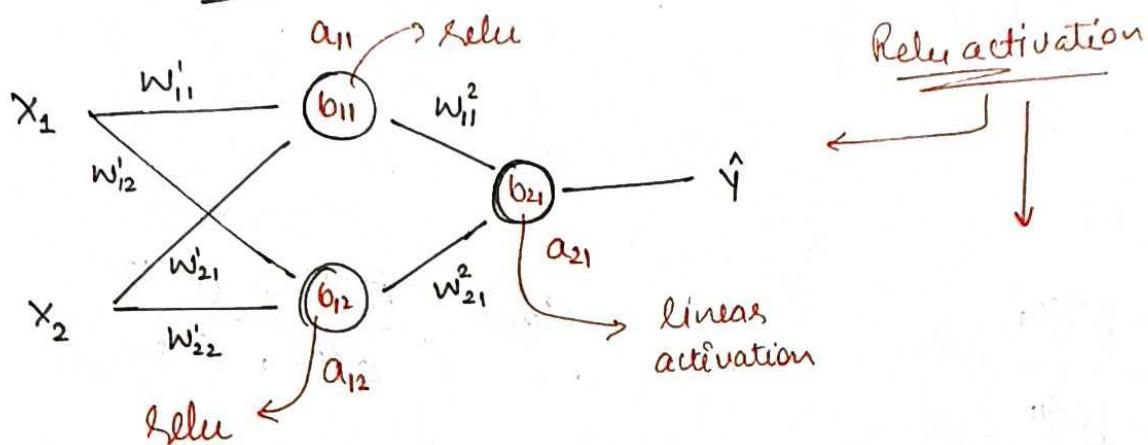
(ii) Compute the cost function.

(iii) Compute the gradients of the cost with respect to parameters using backpropagation.

(iv) Update each parameter using the gradients according to the optimization algo.

What not to do?

Cases → Zero Initialization (Regression Example)



$$a_{11} = \max(0, z_{11})$$

$$z_{11} = w_{11}'x_1 + w_{21}'x_2 + b_{11}$$

$$a_{12} = \max(0, z_{12})$$

$$z_{12} = w_{12}'x_1 + w_{22}'x_2 + b_{12}$$

We know, initial weight is 0. bz of zero initialization
bias is 0.

$$w=0 \quad b=0$$

So,

$$z_{11} = w_{11}'x_1 + w_{21}'x_2 + b_1 = 0$$

$$a_{11} = \max(0, z_{11}) \underset{z_0}{=} 0$$

$$z_{12} = w_{12}'x_1 + w_{22}'x_2 + b_{12} = 0$$

$$a_{12} = \max(0, z_{12}) \underset{z_0}{=} 0$$

$$\boxed{a_{11} = a_{12}} \rightarrow \text{equal and zero}$$

during propagation

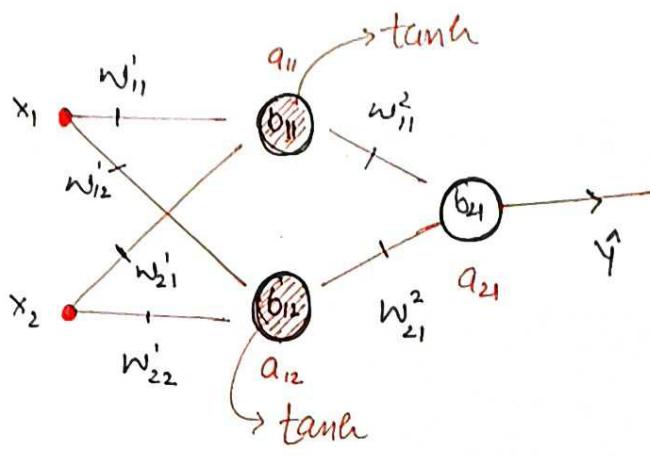
$$w_{11}' = w_{11}' - \eta \frac{\partial L}{\partial w_{11}'} \underset{z_0}{=} 0$$

for solving this, we also need a_{11} derivative but a_{11} is 0 then this term is 0.

$$w_{11}'_{\text{new}} = w_{11}'_{\text{old}} \rightarrow \text{No change or update}$$

→ No training

* In ReLU activation → Not initialize 0 in weight during starting.



tanh activation

$$a_{11} = \frac{e^{z_{11}} - e^{-z_{11}}}{e^{z_{11}} + e^{-z_{11}}}$$

$$a_{12} = \frac{e^{z_{12}} - e^{-z_{12}}}{e^{z_{12}} + e^{-z_{12}}}$$

$$z_{11} = w_{11}'x_1 + w_{21}'x_2 + b_{11}$$

$$z_{12} = w_{12}'x_1 + w_{22}'x_2 + b_{12}$$

We know that, initial weight and bias is 0. \rightarrow zero initialization

$$w=0, b=0$$

So,

$$z_{11} = w_{11}'x_1 + w_{21}'x_2 + b_{11} = 0$$

$$z_{12} = w_{12}'x_1 + w_{22}'x_2 + b_{12} = 0$$

$$a_{11} = \frac{e^{z_{11}} - e^{-z_{11}}}{e^{z_{11}} + e^{-z_{11}}} = 0$$

$$a_{12} = \frac{e^{z_{12}} - e^{-z_{12}}}{e^{z_{12}} + e^{-z_{12}}} = 0$$

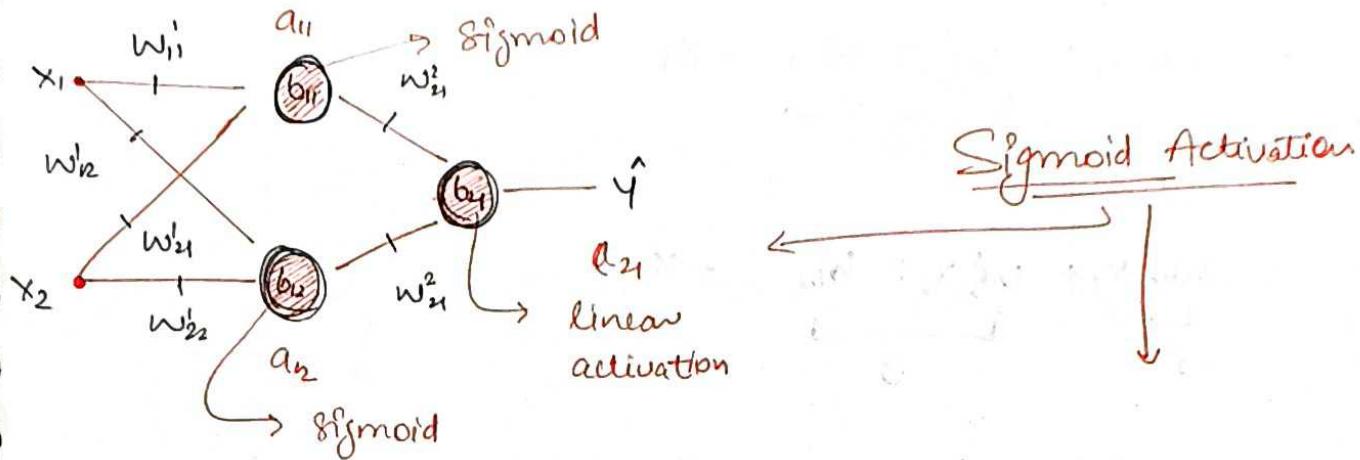
$a_{11} = a_{12} \rightarrow$ equal and zero

[during propagation]

$w_{11}' = w_{11}' - \eta \frac{\partial L}{\partial w}$ \rightarrow to solve this, we also need a_{11} derivatives but a_{11} is 0. then this term is 0.

$w_{11}' = w_{11}'$ \rightarrow No update
New \rightarrow old \rightarrow No training

\rightarrow In tanh activation \rightarrow Not initialize is 0 in weight and bias during starting.



$$a_{11} = \sigma(z_{11})$$

$$a_{12} = \sigma(z_{12})$$

$$z_{12} = w_{12}^1 x_1 + w_{22}^1 x_2 + b_{12}$$

We know that, initial weight and bias is 0.

$$w=0, b=0$$

So,

$$z_{11} = w_{11}^1 x_1 + w_{21}^1 x_2 + b_{11} = 0$$

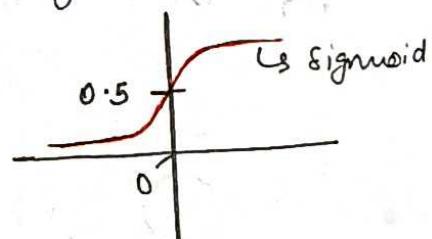
$$z_{12} = w_{12}^1 x_1 + w_{22}^1 x_2 + b_{12} = 0$$

$$a_{12} = \sigma(z_{12}) = 0.5$$

$$a_{11} = \sigma(z_{11}) = 0.5$$

a_{11} and a_{12} is 0.5 because in sigmoid at 0 is 0.5

$$\boxed{a_{11} = a_{12}}$$



$$w = w - n \boxed{\frac{\partial L}{\partial w}}$$

Solving this term for every weight

$$\frac{\partial L}{\partial w_{11}} = \boxed{\frac{\partial L}{\partial \hat{y}}} \boxed{\frac{\partial \hat{y}}{\partial a_{11}}} \boxed{\frac{\partial a_{11}}{\partial z_{11}}} \boxed{\frac{\partial z_{11}}{\partial w_{11}}}$$

$$\frac{\partial L}{\partial w_{12}} = \boxed{\frac{\partial L}{\partial \hat{y}}} \boxed{\frac{\partial \hat{y}}{\partial a_{12}}} \boxed{\frac{\partial a_{12}}{\partial z_{12}}} \boxed{\frac{\partial z_{12}}{\partial w_{12}}}$$

$$\frac{\partial L}{\partial w_{21}} = \boxed{\frac{\partial L}{\partial \hat{y}}} \boxed{\frac{\partial \hat{y}}{\partial a_{11}}} \boxed{\frac{\partial a_{11}}{\partial z_{11}}} \boxed{\frac{\partial z_{11}}{\partial w_{21}}}$$

$$\frac{\partial L}{\partial w_{22}} = \boxed{\frac{\partial L}{\partial \hat{y}}} \boxed{\frac{\partial \hat{y}}{\partial a_{12}}} \boxed{\frac{\partial a_{12}}{\partial z_{12}}} \boxed{\frac{\partial z_{12}}{\partial w_{22}}}$$

$$a_{11} = a_{12}$$

$$z_{11} = z_{12}$$

$$\frac{\partial Z_{11}}{\partial w'_{11}} = \underbrace{w'_{11}x_1}_{\perp} + \underbrace{w'_{21}x_2 + b_{11}}_{\parallel} = x_1$$

$$\frac{\partial Z_{12}}{\partial w'_{12}} = \underbrace{w'_{12}x_1}_{\perp} + \underbrace{w'_{22}x_2 + b_{12}}_{\parallel} = x_1$$

$$\frac{\partial Z_{11}}{\partial w'_{21}} = \underbrace{w'_{11}x_1}_{\perp} + \underbrace{w'_{21}x_2 + b_{11}}_{\parallel} = x_2$$

$$\frac{\partial Z_{12}}{\partial w'_{22}} = \underbrace{w'_{12}x_1}_{\perp} + \underbrace{w'_{22}x_2 + b_{12}}_{\parallel} = x_2$$

Same ↗

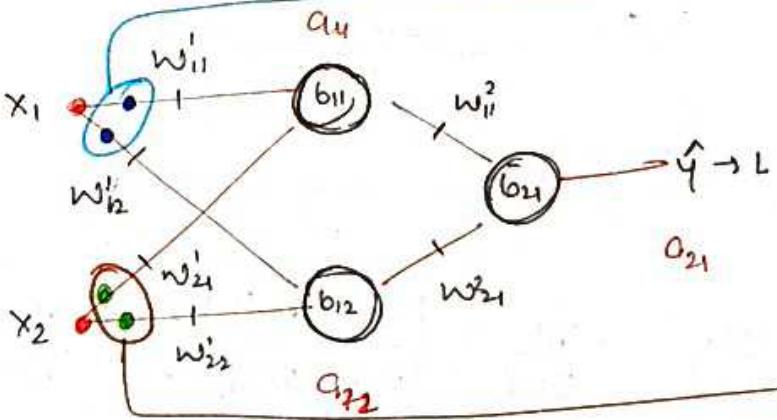
$$\frac{\partial L}{\partial w'_{11}} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial a_{11}} \frac{\partial a_{11}}{\partial Z_{11}} x_1$$

$$\frac{\partial L}{\partial w'_{12}} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial a_{12}} \frac{\partial a_{12}}{\partial Z_{12}} x_1$$

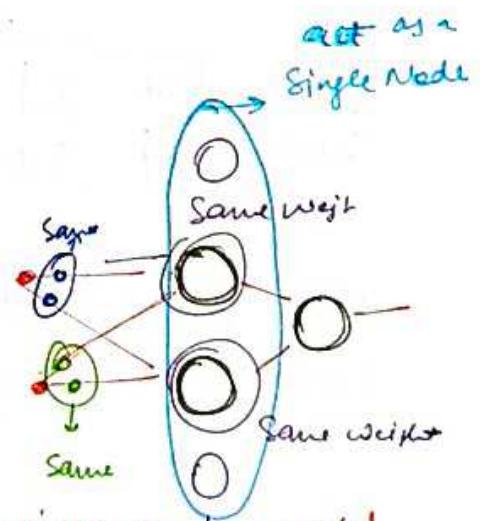
Same ↗

$$\frac{\partial L}{\partial w'_{21}} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial a_{11}} \frac{\partial a_{11}}{\partial Z_{11}} x_2$$

$$\frac{\partial L}{\partial w'_{22}} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial a_{12}} \frac{\partial a_{12}}{\partial Z_{12}} x_2$$

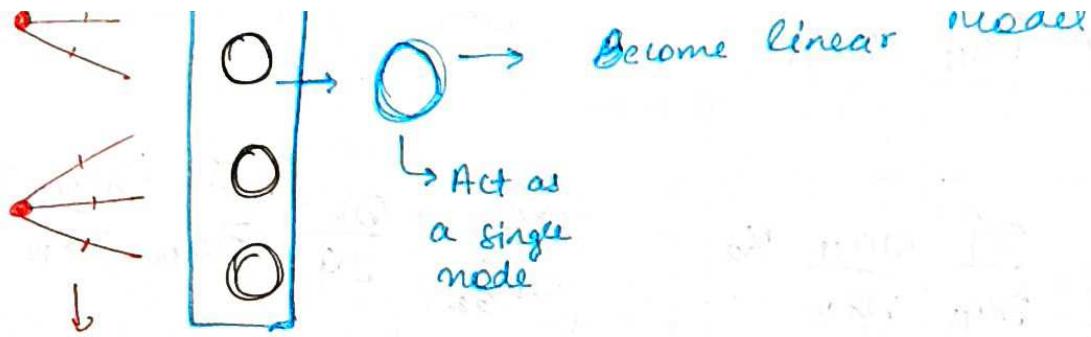


In this both weight are same so, model treat a single weight behave.

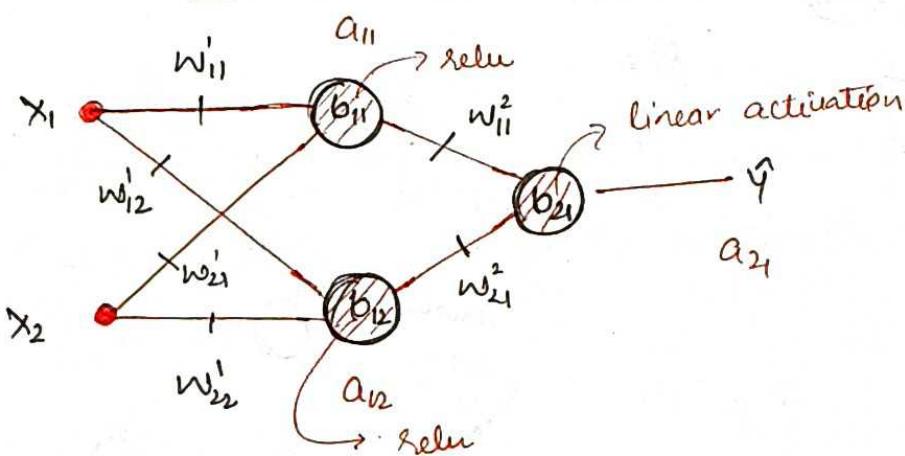


In short → If we use sigmoid and initialize (w) with 0. So, even we add 100 or 1000 node but it will act as a single Node.

Every node contain same value and this model convert into linear model



Case 2 \rightarrow Non - 0 constant value



$$\text{let } w = 0.5, b = 0.5$$

$$a_{11} = \max(0, z_{11})$$

$$z_{11} = w_{11}' x_1 + w_{21}' x_2 + b_{11}$$

$$z_{11} \neq 0 \text{ (Non-zero value)}$$

$$a_{11} = \max(0, z_{11}) \neq 0$$

$$a_{12} = \max(0, z_{12})$$

$$z_{12} = w_{12}' x_1 + w_{22}' x_2 + b_{12}$$

$$z_{12} \neq 0 \text{ (Non-zero)}$$

$$a_{12} \neq 0$$

Initial weight and bias are same in z_{11} and z_{12} . So,

$$\boxed{z_{11} = z_{12}} \rightarrow \boxed{a_{11} = a_{12}}$$

$$\frac{\partial L}{\partial w'_{11}} = \frac{\partial L}{\partial y} \frac{\partial \hat{y}}{\partial a_{11}} \frac{\partial a_{11}}{\partial z_{11}} x_1$$

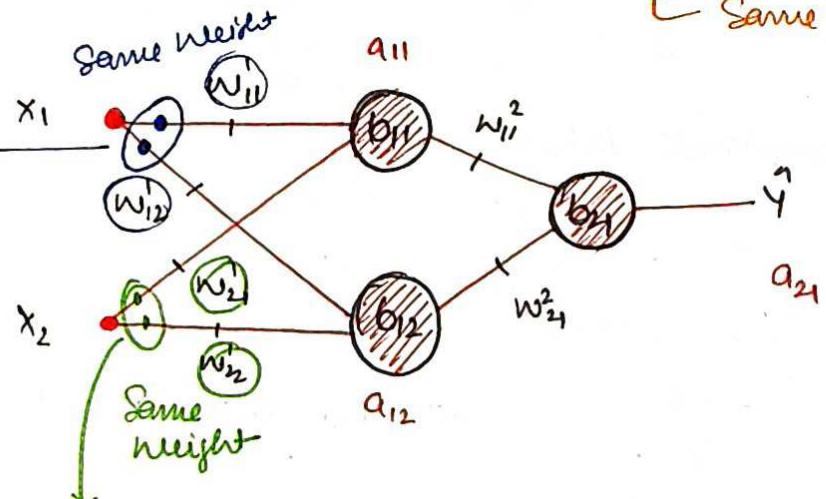
$$\frac{\partial L}{\partial w'_{12}} = \frac{\partial L}{\partial y} \frac{\partial \hat{y}}{\partial a_{12}} \frac{\partial a_{12}}{\partial z_{12}} x_1$$

↑ Same ↑

$$\frac{\partial L}{\partial w'_{21}} = \frac{\partial L}{\partial y} \frac{\partial \hat{y}}{\partial a_{21}} \frac{\partial a_{21}}{\partial z_{21}} x_2$$

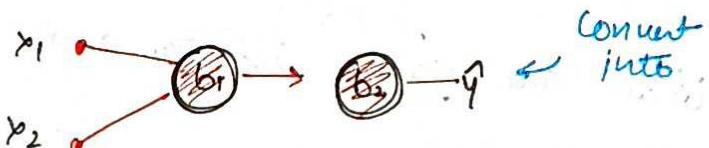
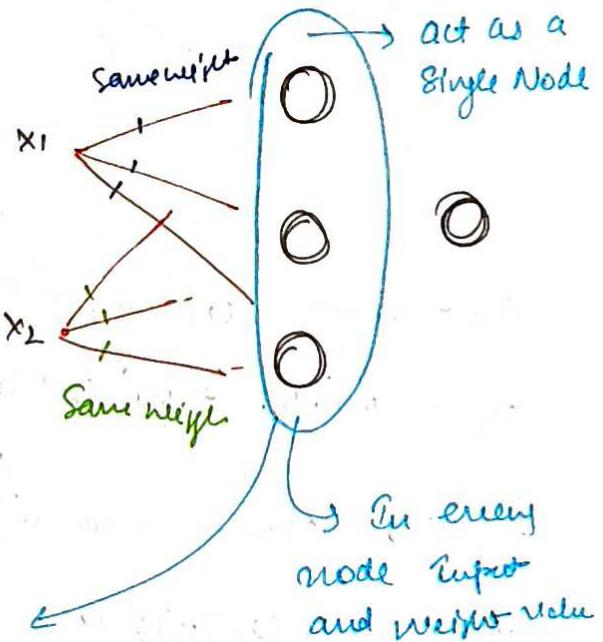
$$\frac{\partial L}{\partial w'_{22}} = \frac{\partial L}{\partial y} \frac{\partial \hat{y}}{\partial a_{22}} \frac{\partial a_{22}}{\partial z_{22}} x_2$$

↑ Same ↑



Both weight are same. So, it act as a single weight

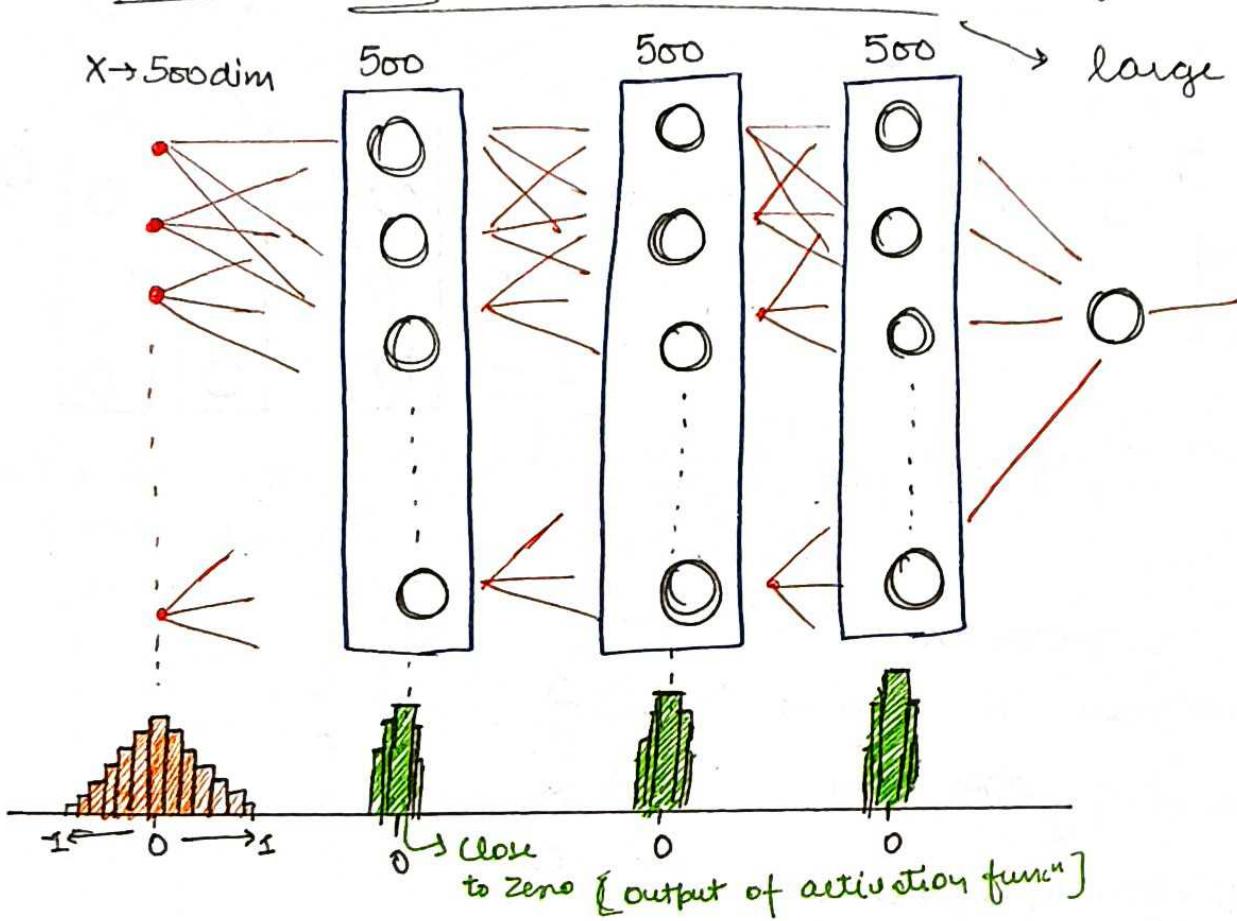
This model convert into linear model



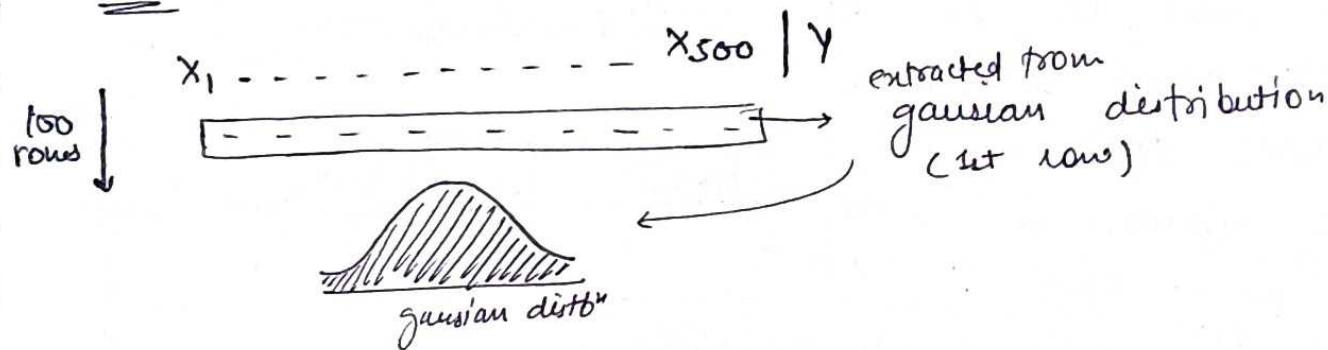
* So, don't initialize weight and bias with 0 value and any constant value. Do initialize with random.

Case 3 - Random Initialization

(59)



Data →

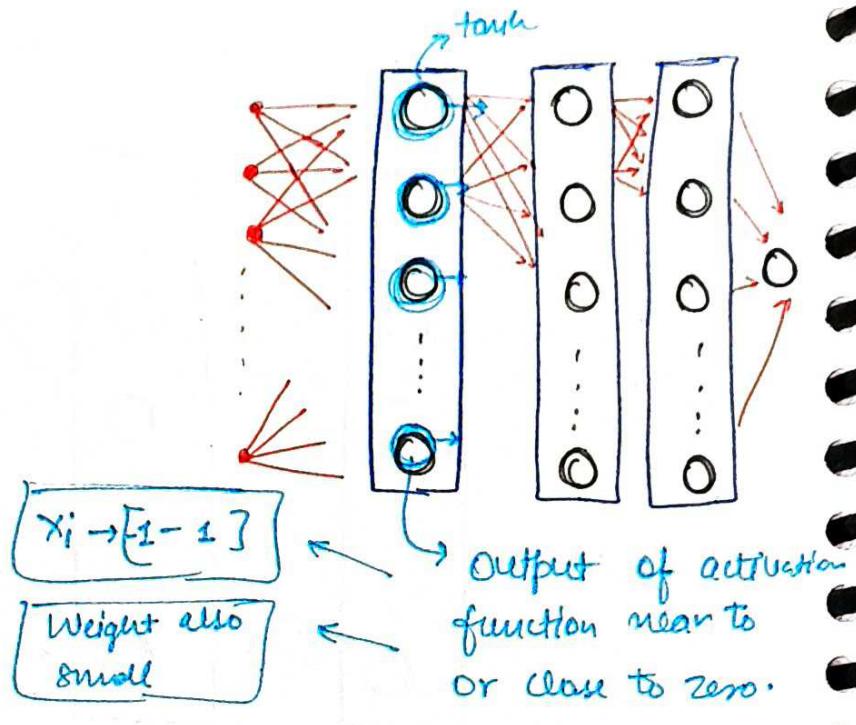
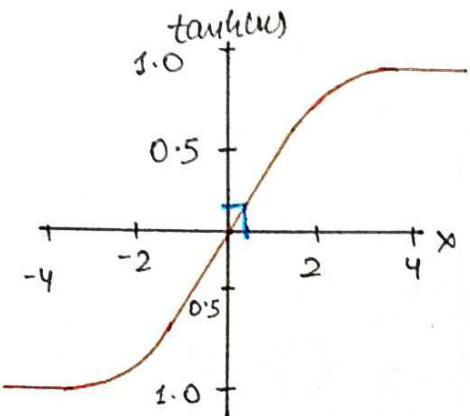


Weight Initialize = random

np.random.random $[(500, 500) * 0.01]$ after multiplying with 0.01 → weight
bias = 0 small

x_i lies betw $\rightarrow [-1 to 1]$

when we do $\sum w_i x_i$ → output is also small number.
(2) ↴ is small

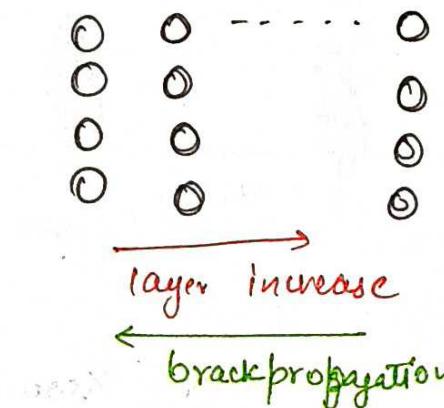


When increase the layers during backpropagation time,

Small number will multiply
eg:- $(0.1 \times 0.2 \times 0.1 \dots)$ Then gradient will be zero or close to zero.

So, when we update the weight

And no update in weight
new weight is equal to old weight



$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w} \rightarrow \text{close to zero}$$

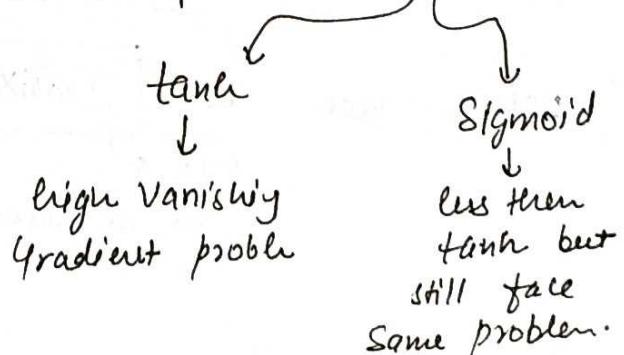
* Again Vanishing Gradient problem

Summary

1. Initializing small weights $\rightarrow 0.0007$

And layer is deep (large) or not much large

Still facing Vanishing gradient problem in



In ReLU → small value → face same vanishing gradient problem (not strong) 60
 but less than tanh and sigmoid
 coz values are not zero.
 but training will be very slow and converge ←
 also very slow. (large epochs need)

Random Initialization (Large value)



np.random.randn(500, 500)

↳ [0-1]

this range is big for deep learning

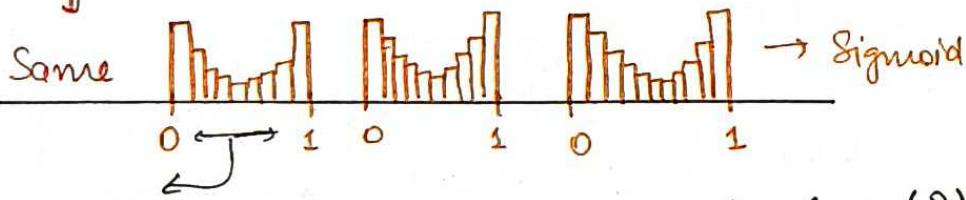
↳ tanh / sigmoid

Saturation → means

let

$$\sum w_i x_i \rightarrow 100$$

↓
Start saturating



In histogram → Maximum value is lower (0) and higher (1) in tanh and ReLU.

if value is high so take high or negative
 very low, so take negative very low value

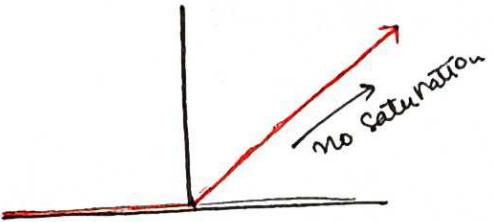
This cause → slow training

→ Vanishing gradient problem (in worst case)
not positive for every node

Saturation → If weight is large ($(\sum w_i x_i) \rightarrow 100$) send to activation function ($g(z)$). So, output is high [tanh → 1, Sigmoid → 1]

If weight is highest negative value \rightarrow after activation function value is also negative ($\tanh = -1$, Sigmoid ≈ 0)

ReLU \rightarrow ReLU is non saturating function in x .



Update \rightarrow

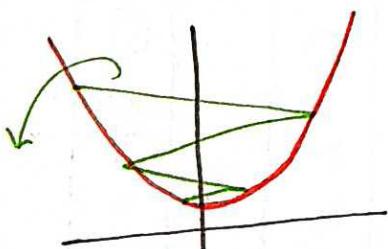
$$w = w - \eta \frac{\partial L}{\partial w}$$

gradient is high then
changes are also high

high
changes

if $\sum w_i x_i$ is 250
then ReLU output is 250

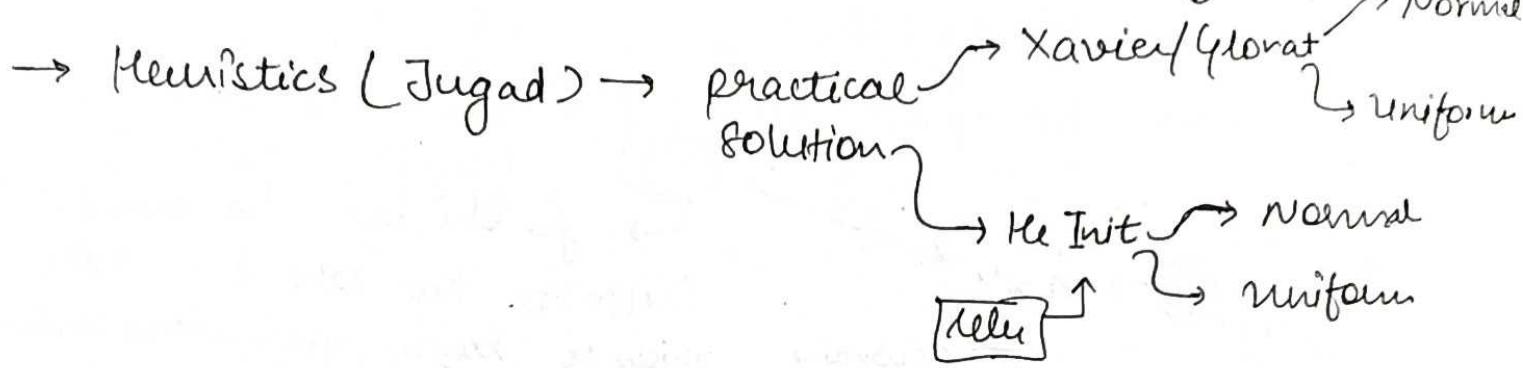
If this number is
greater than gradient
also greater



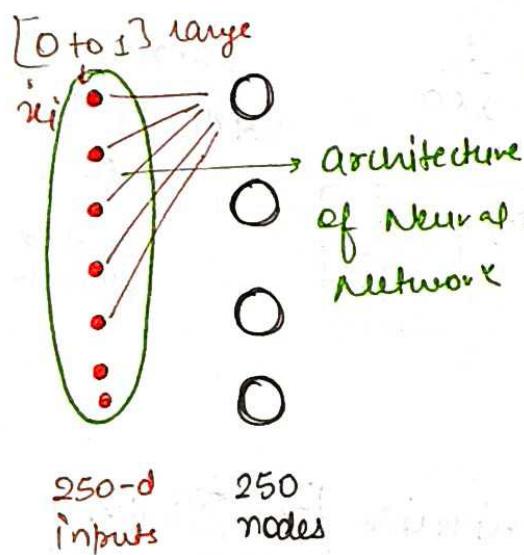
Not Initialize with

- (i) zero
- (ii) non-zero
- (iii) small random
- (iv) large random

What can be done



Intuition



Small weight (-0.01) to (0.01) constant

np.random.rand(250, 250) * 0.01

if we multiply x_i with w_i and add all $(x_i w_i)$ which is 250 * 80, the output is very small. And send into activation and activation convert into 0.

Large Weight (-3 to 3) constant

np.random.rand(250, 250) * 1

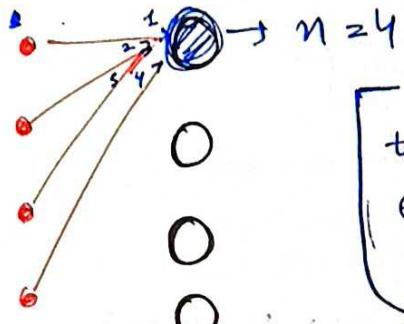
multiply weight with w_i ($x_i w_i$) and all the 250 multiply values ^{add and} lies between (-250 to 250). We use send these values into activation function. Activation function starts saturating.

if value is -250 and this value reach negative max and 250 value reach positive max.

why these problem are generally max.

because of $z = \sum x_i w_i$ formula and we can not multiply with constant. Multiple the Number, it should be dependent on architecture of NN.

So, Scientist decided that which number [imp. random. randn (250, 250)] generate in his code it should be [variance = $\frac{1}{n}$] a number's variance



according to intuition example

$$\text{Variance} = \frac{1}{n} = \frac{1}{250}$$

jis bhi layer ka weight initialize karake wo noki piche layer mai kitne nodes hain. or use particular node ko kitne input mil raha.

so, multiplying with $\sqrt{\frac{1}{250}}$

Standard deviation

$$\text{np.random. randn}(250, 250) * \sqrt{\frac{1}{250}}$$

Intuition of Variance $= \frac{1}{n}$

If $n \rightarrow \uparrow\uparrow$ then I want to minimize $Z = \sum w_i x_i$

↳ no. of inputs

So, if n is greater $\rightarrow \frac{1}{n}$ (multiply with weight) and ~~overall~~ individual weight is small.

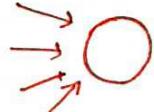
If n is smaller then individual weight is large.

let $\frac{1}{\sqrt{n}}$ \rightarrow most of the time this formula

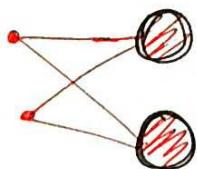
Only 2 inputs

$$\frac{1}{\sqrt{\text{fan-in}}}$$

No. of inputs coming to the node.



$$\rightarrow \text{np.random. randn}(2, 2) * \sqrt{\frac{1}{2}}$$



$$\frac{2}{\sqrt{\text{fan-in} + \text{fan-out}}} \quad \begin{matrix} \text{fan-in} + \text{fan-out} \\ \text{another formula} \end{matrix}$$

Xavier Init (Normal)

Uniform Distribution

Xavier uniform

$$[-\text{limit}, \text{limit}] \quad \text{limit} = \sqrt{\frac{6}{\text{fan-in} + \text{fan-out}}}$$

He uniform

$$[-\text{limit}, \text{limit}] \quad \text{limit} = \sqrt{\frac{6}{\text{fan-in}}} \quad \begin{bmatrix} \text{code in weight} \\ \text{initialization} \end{bmatrix}$$