

→ "Activation functions" :-

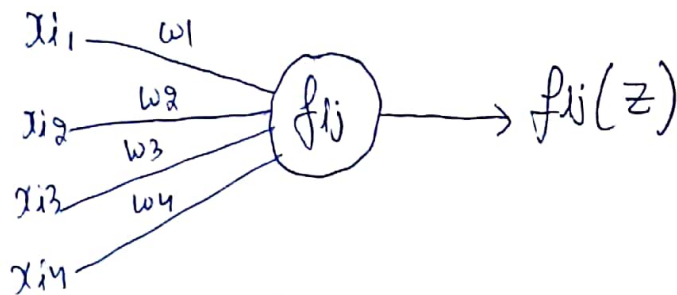
①

from 1980's till early 2000's there are two activation function which are most popular & these are :-

↳ Sigmoid function (Sigmoid activation unit)

↳ Tanh function (Tanh activation unit)

let's assume we have an activation unit, let's call it " f_{ij} ." let's assume we have some inputs & some weights associated with these inputs.



Now when we multiply w_i & x_i & summate them, that, is equivalent to $w^T x$

$$\sum_i w_i x_i = w^T x$$

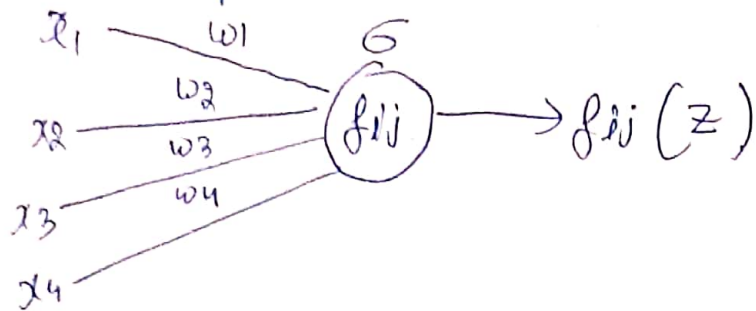
$$\text{where } w = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \end{bmatrix} \text{ and } x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}$$

So, the input that we get at " f_{ij} " is the weighted sum of w_i & x_i .

let's call that input as, $z = \sum_i w_i x_i = w^T x$.

If " f_{ij} " is the Sigmoid Unit

lets represent it with " σ " sign



$$z = \sum_i w_i x_i = w^T x$$

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

or

$$\frac{1}{1+e^{-z}}$$

→ This is what we used in "logistic regression" to squashes inputs.

"Requirements of the activation function" →

1) It needs to be differentiable.

2) It should be easy & fast to differentiate.

We have already seen the usage of sigmoid function in "logistic regression".

lets now look at the advantages of sigmoid from "differentiation perspective".

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

$$\text{+ } \boxed{z = w^T x}$$

Now if we compute derivative of this " σ " w.r.t. z . we get.

P.T.O

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

(3)

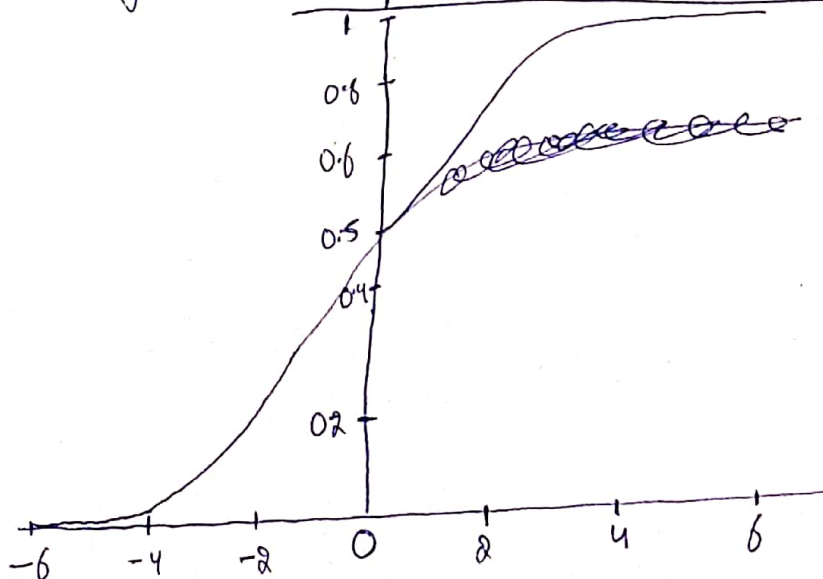
The important part here is, the derivative of the sigmoid function is represented in terms of sigmoid function itself.

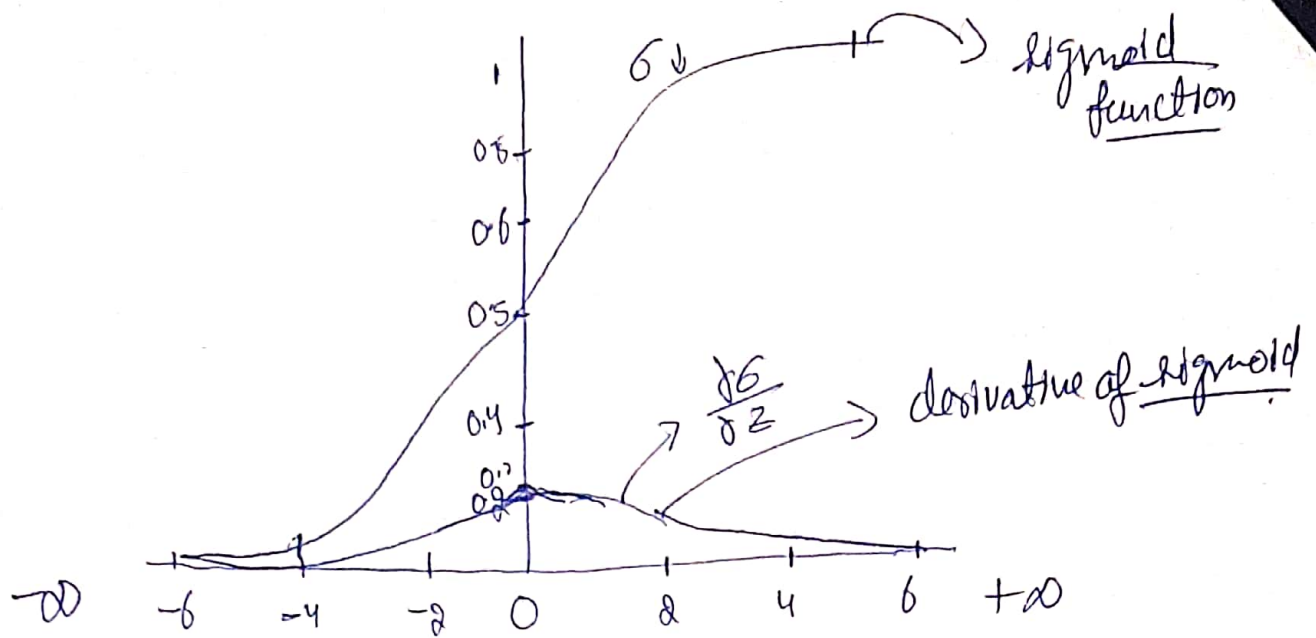
Let's say you write some code to implement the sigmoid function(z).

Sigmoid(z) \rightarrow [This function can be used both during forward propagation and backward propagation.]

\hookrightarrow We can use the same function during 'forward propagation'. as well as computing derivatives during ['backward propagation'] for updating the weights.

The sigmoid function looks like an 'S' shaped curve.





Sigmoid function is an S-shaped curve, it goes from $-\infty$ to $+\infty$, minimum value is '0' & maximum value is '1'.

If we notice the derivative of sigmoid the max. value is ~~at~~ 0.25 & the derivative is only significant b/w -4 & $+4$ because after -4 & $+4$ the derivative is extremely small.

Derivative basically says how fast is the sigmoid function changing.

$$0 \leq \frac{d\sigma}{dz} < 1$$

(precisely less than 0.3)

There is another popular activation unit called tanh function.

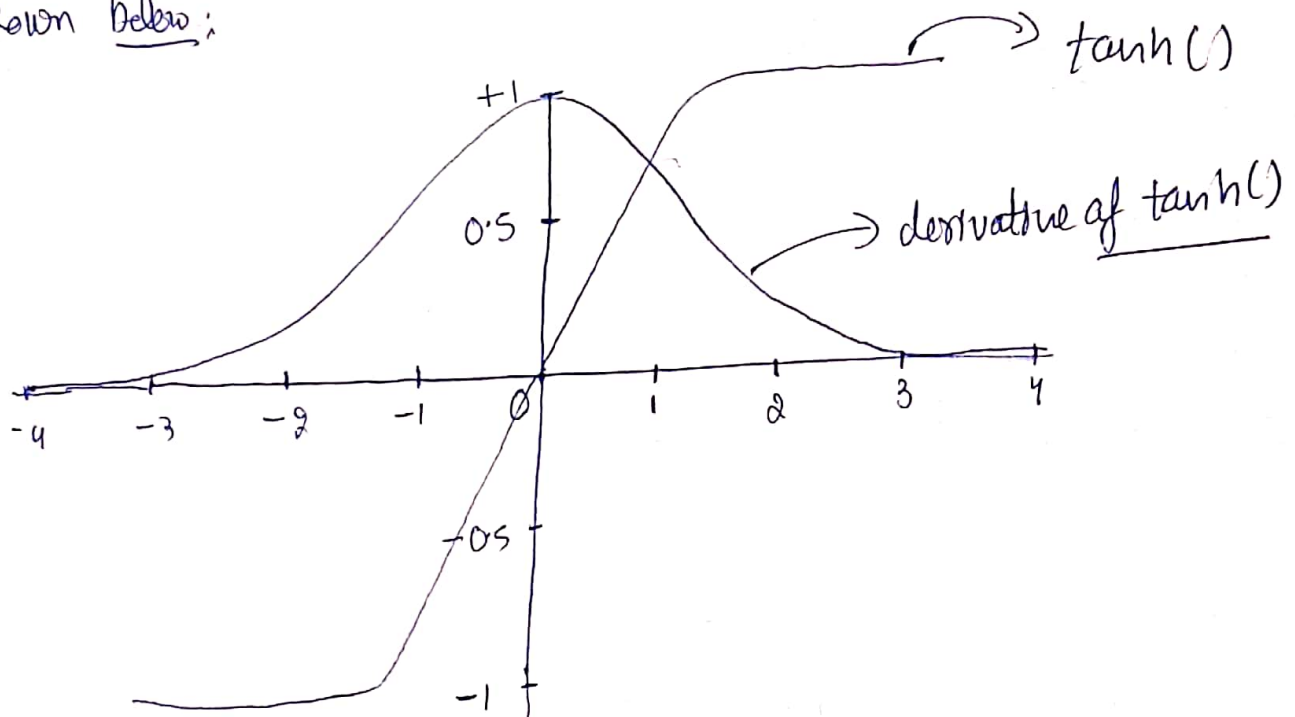
$$a_z = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$\frac{da}{dz} = 1 - a^2 = 1 - (\tanh^2(z))^2$$

$\tanh(z)$ function is very very similar to sigmoid function (5)

We know that derivative of sigmoid function can be represented in terms of sigmoid function itself, similarly derivative of $\tanh()$ function can also be represented in terms of $\tanh()$ function itself.

The plot of $\tanh()$ function & its derivative looks like as shown below:



$\tanh()$ function is also an S-shaped curve, but the 'minimum value' is "-1" & 'maximum value' is "1".

Derivative of $\tanh(z)$ is more peaked

It lies between "0" & "1" the maximum value of derivative of $\tanh()$ is at "1" & minimum value of its derivative is at "0".

& Between "-3" & "3", it is having some reasonable value.

After '+3' & '-3', the derivative of "tanh()" becomes very very small.

So traditionally, "tanh()" & "sigmoid()" are the two mostly used activation functions, bcoz they are differentiable & they are easily differentiable, bcoz of that we can represent them using the same function. (the function as well as its derivative).

Note:-> Derivative of sigmoid function can be represented in terms of sigmoid function itself.

& Derivative of tanh() function can also be represented in terms of tanh() function itself.

Even though these two are very popular, but the most popular activation function, we use now a days is the

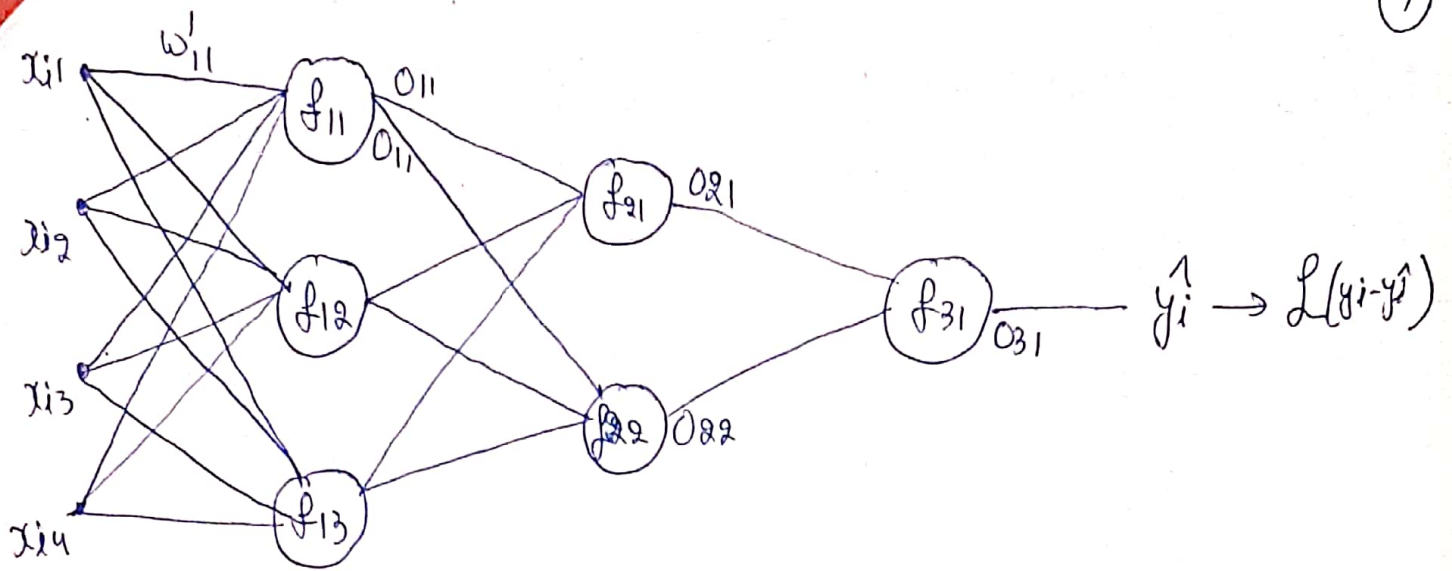
ReLU

↳ "Rectified Linear Unit".

→ Vanishing Gradients (Problem) :-

One of the biggest problem with neural networks in 80's till early 2000s was this problem of vanishing gradients.

This is literally one of the most important reason, why people lost interest in neural networks in late 90's & early 2000's. (This is not the only reason but one of the important reasons)



All of f_{ij} 's in earlier days, people used sigmoid or $\tanh()$ then typically as activation units.

There is one major problem with ["sigmoid()"] which is the same problem with ["tanh()"] also.

Let's understand the problem.

Suppose we want to update the weight " w'_{11} ".

We know that, to update this weight, we need to compute the derivative of " L " w.r.t. " w'_{11} ".

As our update equation is:

$$\{ (w'_{11})_{\text{new}} = (w'_{11})_{\text{old}} - \eta \left(\frac{\partial L}{\partial w'_{11}} \right)_{\text{old}} \}$$

P.T.O

$$\frac{\partial L}{\partial w'_{11}} = \frac{\partial L}{\partial o_{31}} \left[\left(\frac{\partial o_{31}}{\partial o_{21}} \times \frac{\partial o_{21}}{\partial o_{11}} \times \frac{\partial o_{11}}{\partial w'_{11}} \right) + \left(\frac{\partial o_{31}}{\partial o_{22}} \times \frac{\partial o_{22}}{\partial o_{11}} \times \frac{\partial o_{11}}{\partial w'_{11}} \right) \right]$$

Now, we know that at all the activation units, we have used 'sigmoid' or 'tanh function' & in order to update the weights we need to compute the derivatives

& we know that derivative of sigmoid is $0 \leq \frac{\partial \sigma}{\partial z} \leq 1$

strictly speaking $\frac{\partial \sigma}{\partial z}$ is never reaching even 0.3 value

& in case of $\tanh()$ the derivative is strictly

$$\text{b/w } \left[0 \leq \frac{\partial \tanh(z)}{\partial z} \leq 1 \right]$$

Since each of these activation units are sigmoid & tanh, its derivatives of these will be b/w 0 & 1. Typically speaking about sigmoid, its derivative will always be less than 0.3.

Here since we are multiplying these derivatives & since these derivatives are small so the multiplication term will be smaller even more.

Let's only focus on the first part here.

In the first multiplication term we have three numbers all less than 0.3 (as it is sigmoid)

(9)

$$\left[\frac{\delta O_3}{\delta O_2} \times \frac{\delta O_2}{\delta O_1} \times \frac{\delta O_1}{\delta w'_{11}} \right]$$

let $\frac{\delta O_1}{\delta w'_{11}} \uparrow$

$$0.2 \times 0.1 \times 0.05 = 0.0010 = 10^{-3}$$

$$\downarrow \quad \downarrow \quad \downarrow$$

Very very less value

It is a miniscule value

& Remember, we are using this miniscule value in the update function.

$$\left[(w'_{11})_{\text{new}} = (w'_{11})_{\text{old}} - \eta \left(\frac{dL}{\delta w'_{11}} \right) \right]$$

let's say $(w'_{11})_{\text{old}}$ is 2.5

$$\& \eta = 1 \quad \& \text{let } \frac{\delta L}{\delta w'_{11}} = 0.001$$

$$\text{Now } (w'_{11})_{\text{new}} = 2.5 - 1 * 0.001$$

$$= \underline{2.499}$$

↓

which is very similar to the old value

∴ New weight value is very very close to the old value (old weight value)..

This is happening becoz derivative is small, & derivatives is small becoz individual derivatives are small & individual derivatives are small because we are using sigmoid() or tanh() as our activation units.

So, updating the weights become very very hard & training becomes very very slow.

Here we have used just a three layered network.

Imagine if we have a 50 layered or 10 layered network then updating weights for such a deep network would be even more tough.

Note ⇒ # of multiplication of derivative =
of hidden layers + 1

if we have 2 hidden layers then no. of multiplication of derivative is = 3

Similarly if we have 10 hidden layers, then no. of multiplication of derivative is = 11

P.T.O

If just with 3 multiplications, we are getting such a small value of derivative (11)

Then, imagine how minuscule the value would be if 11 such multiplication operations are performed.

In this example with 3 multiplications, we get derivative value as 10^{-3}

If we have eleven multiplications, then the value becomes even more smaller = 10^{-12} .

vvv small

So, if the derivative is very small, then newer weights & older weights become almost the same.

Imagine the derivative was 10^{-10}

then our new weight would have been

$$\begin{array}{rcl} 2.5 - 0.0000 \dots 1 & & \\ \text{old.} \downarrow & = & 2.499999 \dots 9 \\ & & \downarrow \text{new.} \end{array}$$

which is exactly same.

When this problem is called 'Vanishing Gradient Problem', because the final gradient $\left(\frac{\partial L}{\partial w_{11}}\right)$ we are getting here is vanishing or it is becoming very very small. This is because of chain rule of multiplication of the sigmoid function we are using.

→ ReLU (Rectified Linear Units)

One of the big problems in classical neural networks era (pre-2000) have been the 'vanishing gradient problem'

which is very often seen when we are using a sigmoid() activation or a tanh() activation function.

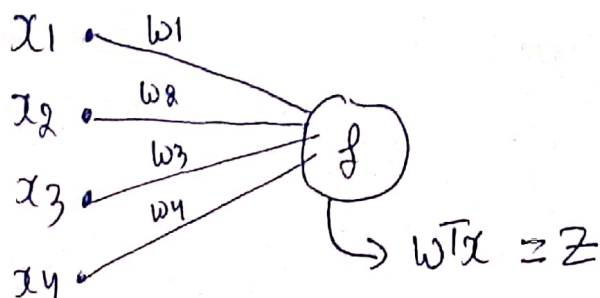
When we have this problem of 'vanishing gradient' our convergence also slows down & training ~~becomes~~ ^{becomes} longer & longer.

The reason is, if we have vanishing gradient problem, the update is very very small.

'ReLU' were introduced in 2011 & they are one of the most important ideas in 'Deep Learning'

Note → Today, all the activation problems are by default 'ReLU'

Let's now see how ReLU looks like:-



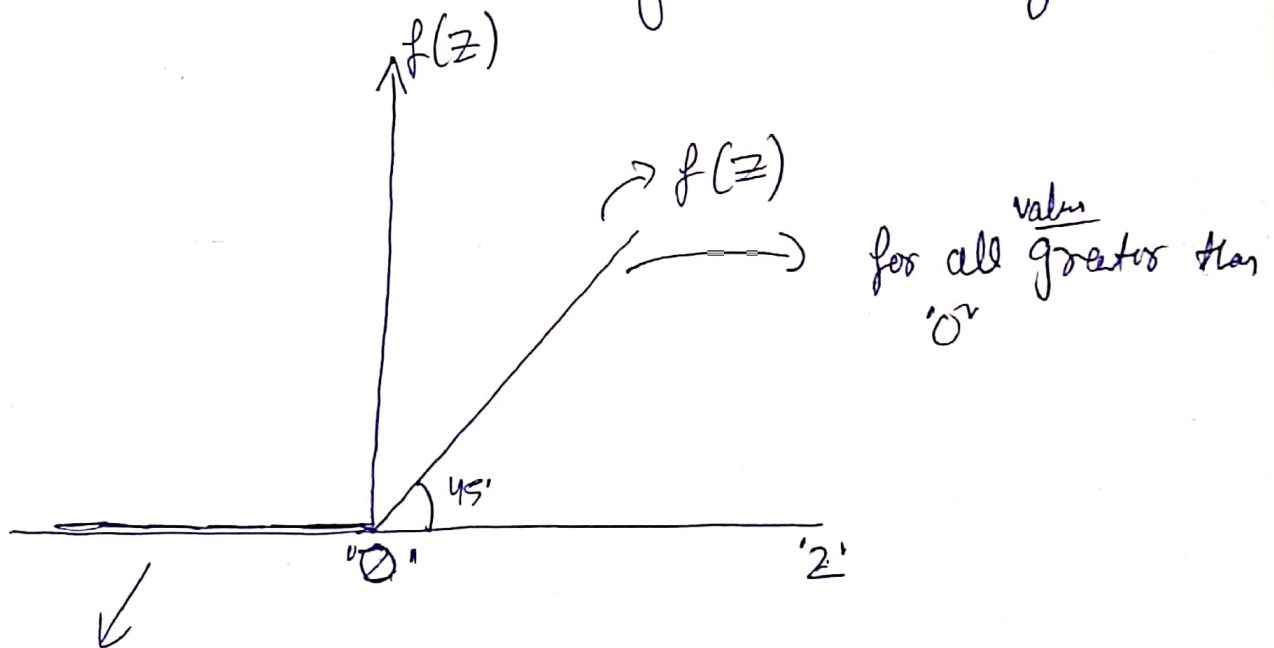
Now function $f(z)$ for ReLU it often written as (13)

$$f(z) = z^+ = \max(0, z)$$

ReLU definition (max value of $0, z$)

let's see how does it work \Rightarrow

$$f(z) = \max(0, z) = \begin{cases} 0, & \text{if } z \leq 0 \\ z, & \text{otherwise} \end{cases}$$



for all the values which are negative (less than 0)
the value of the function is '0'

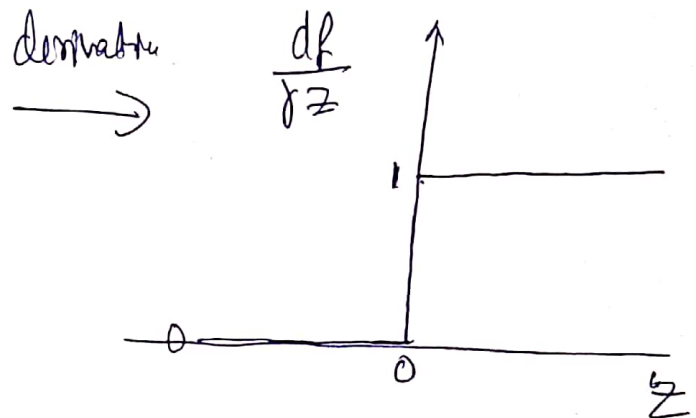
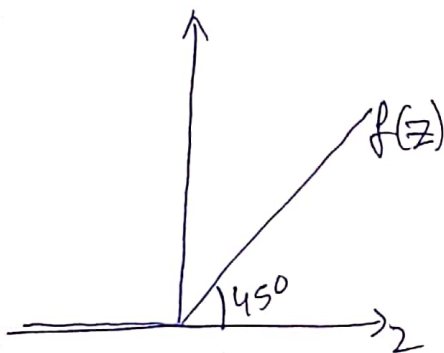
& whenever 'z' is greater than '0', its value simply
will be 'z'

This is basically a "45°" line

Computing derivative of the 'ReLU function' is also very very simple.

The condition for a function to be used as an activation unit is that it should be differentiable & it should be differentiated easily.

So if this 'ReLU' function is $f(z)$ then its derivative looks like



for a 45° line, the slope is exactly equal to '1'

Derivative is not defined at '0', and everywhere else it is defined:

The derivative of the activation (ReLU) function is either '1' or '0'.

$$\frac{df_{\text{ReLU}}}{dz} \in \{0, 1\}$$

So here we will not be getting values like 0.2, 0.3, 0.1 & we are not multiplying these by the

problem of vanishing gradient will not occur.
———— 0 ————

(15)