

ANN  $\rightarrow$  DL

Date:

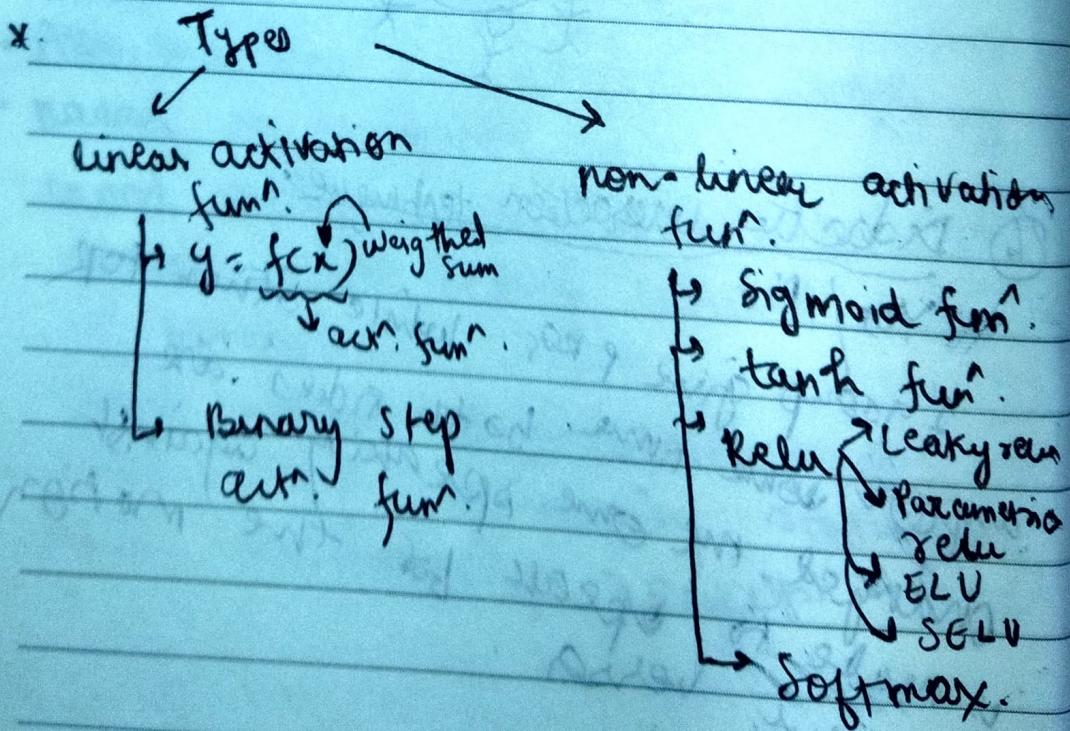
Lecture 6

## ACTIVATION FUNCTION

23<sup>rd</sup>

July 2023

- \* why do we need activation fun".
  - To achieve non-linearity in the data. (make classification or regression in case of non linear dataset).
  - To activate the neuron (Acts as a gate with triggers when we pass the threshold).

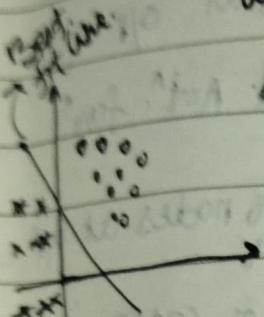


**Note:** weighted sum represent the line in 2D space (for linear separable dataset).

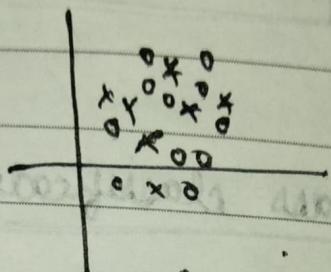
$$(x_1 w'_1 + x_2 w'_2 + b_1) = 0_{12}$$

$ax + by + c$

Point line:  $w^T x + b$ .



linearly  
separable  
dataset.



non linearly  
separable  
dataset.

**Note:** For non-linearly separable dataset, weighted sum represent curve in 2D space.

↳ Ques.

\* So, To conclude if dataset is linearly separable then weighted sum alone can be applied on top of regression or classification.

↳ Ans.

\* But if dataset is non linearly separable then we need to pass the weighted sum through some activation fun? (non linear act. fun.) so that we are able to implement respective regression/ classification user case.

- \* Note: Activation fun<sup>n</sup>. at o/p node:  
Regression  $\rightarrow$  1 node/neuron at o/p layer  
 and  
 $\hookrightarrow$  we use linear act<sup>n</sup>. fun<sup>n</sup>.

Binary classification  $\rightarrow$  1 nodes at o/p layer  
 $\hookrightarrow$  Sigmoid: Act<sup>n</sup>. fun<sup>n</sup>.

Multi class classification  $\rightarrow$  no. of nodes at o/p layer  
 $=$  no. of categories in o/p.  
 $\hookrightarrow$  Softmax: Act<sup>n</sup>. fun.

- \* Note: Act<sup>n</sup>. fun<sup>n</sup>. for nodes/neuron at hidden layer:

In old days: sigmoid & tanh

In current time: ReLU and its diff. variants since it's so powerful that it replaced the earlier variants.

- \* (Sigmoid & tanh in hidden layer)  $\rightarrow$  linear separable dataset

\* in hidden layer

Date:

1. (cell & its variants)  $\rightarrow$  non linear differentiable dataset (Regression & Classification)

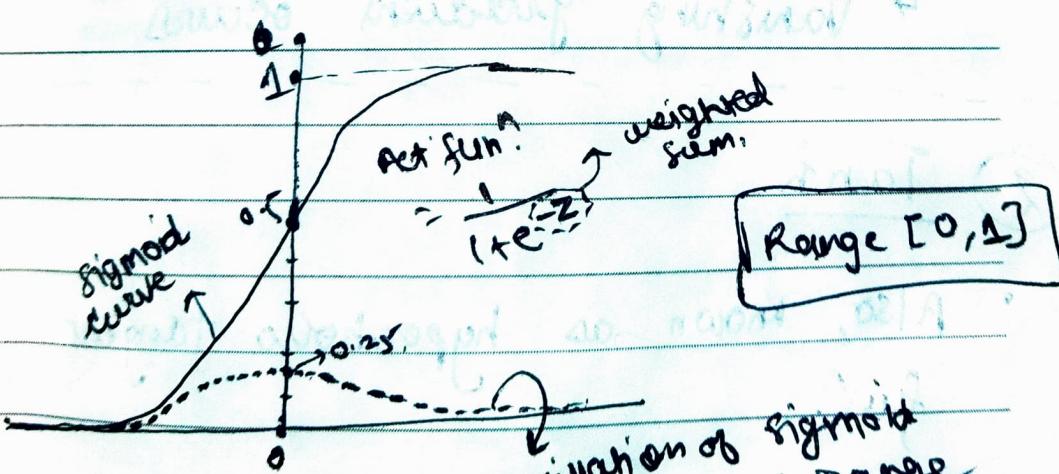
Sigmoid fun<sup>n</sup> } vanishing Gradient  
tanh fun<sup>n</sup>.

tanh & its variants } dying neuron.  
~~overfitted~~

## \* Properties of Act<sup>n</sup> fun<sup>n</sup>:

- ① Non-linearity
- ② Differentiable. (since it will be required to tweak weights in back propagation)
- ③ Computationally inexpensive
- ④ Zero-centered (since then training will be faster due to normalization)
- ⑤ Non-Saturating / (To counter Vanishing gradient)  
non-squeezed

## ① SIGMOID FUN<sup>n</sup>:



derivation of sigmoid  
curve having range  
[-0.25, 0.25].

squeezed [0 - 0.25]

This is the squeezing fun<sup>n</sup> when we are squeezing values to 0/1. This violates the non-squeezed property.

TCS  
CONSULTANCY  
SERVICES

of act. fun. due to this vanishing gradient. This is the reason why we do not use sigmoid in hidden layer and instead use ReLU & its variants.

Date:

- Advantage:

- ① Non linear ( $f(x) = e^{-x}$  and not  $y = f(x) = x$ )
- ② Differentiable
- ③ O/P layer  $\rightarrow$  Binary Classification

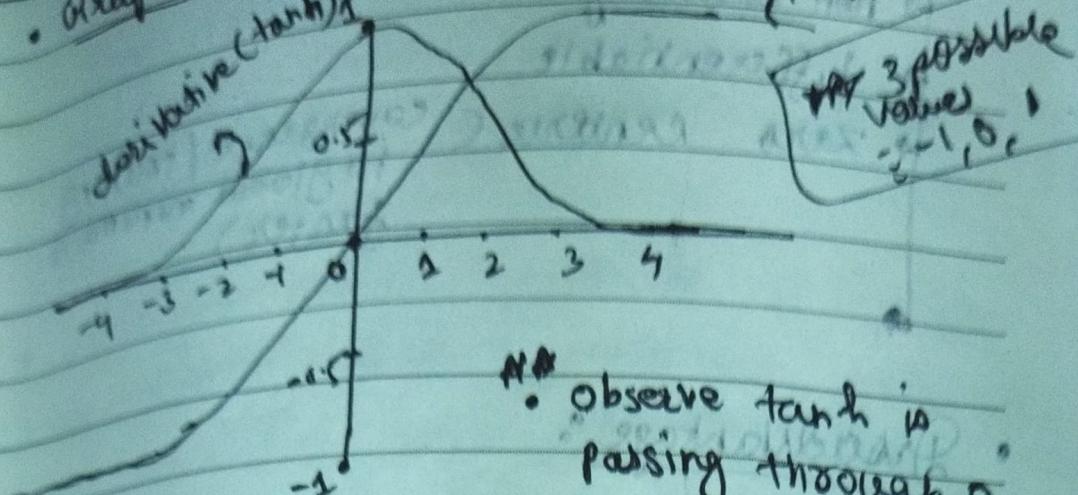
- Dissadvantage:

- Saturated / Squeezed fun<sup>n</sup>.
- computationally expensive
- non zero centric hence training is slow.
- Vanishing gradient occurs

## 2 Tanh

- Also, known as hyperbolic tangent fun<sup>n</sup>.

## Graphical Representation :



• Observe tank is passing through 0 hence it is zero centric Act" fun"

- This is saturating in nature since here values are squeezed in the range  $(-1, +1)$ . In sigmoid it was  $(0, 1)$ .

## • formula :

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

↓ sigmoid

- ## • Range of Derivatives

$\tanh = (0, 1)$  and  $\text{sigmoid} = \log(0.25)$

## Advantage :

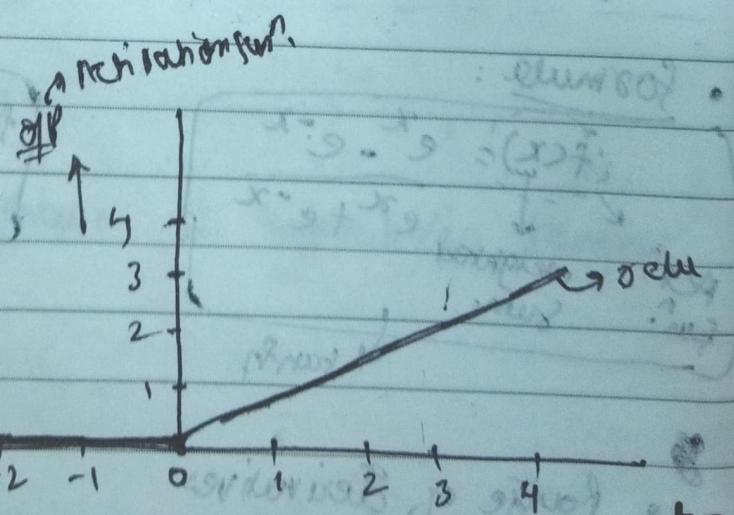
- non linearity
- differentiable.
- zero centered ( convergence / tracking to global minima is faster)

## Disadvantage :

- Saturated / Squeezing nature.
- Vanishing gradient.
- Computationally expensive
- Do not use tanh as activation fun" at o/p layer for binary classification.

## ③ RELU:

Graph:



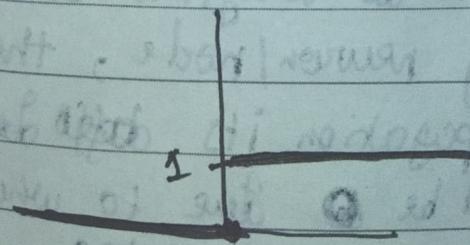
Possible Values,  
weighted sum

- By observing the graph we can say that RELU produces o/p as 0 for -ve values, whereas it gives same value of +ve values.
- we might think that it is linear in nature. But in actual when we are combining multiple RELU then it becomes non-linear in nature. This is the reason why we use it in hidden layers whereas others in o/p layer & also since RELU solves vanishing gradient problem.

Formula :

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

- This is non differentiable.
- Derivative will be either 0 or 1 because of which a problem occurs called dying neuron and then we use its variant leaky relu.



### Advantage:

- Non linear (on adding multiple relu).
- non-saturated saturated in the region.
- computational inexpensive.
- Handles vanishing gradient problem.
- Convergence to global minima.
- fastest ~~reaches zero gradient~~.

ReLU > tanh > Sigmoid

Order of convergence →

(Hence sigmoid is slowest.)

### Disadvantage:

#### Dying neuron Problem:

\* So differentiation of ReLU art? fun.

It will be  $\begin{cases} x \leq 0 \Rightarrow 0 \\ x > 0 \Rightarrow 1 \end{cases}$  weight sum off of art. fun.

\* So, when there is any -ve I/P

or in other words weighted sum is -ve for any neuron/node, then in back propagation its ~~difficult~~ derivative will be 0 due to which weights will not get updated updated for that neuron and

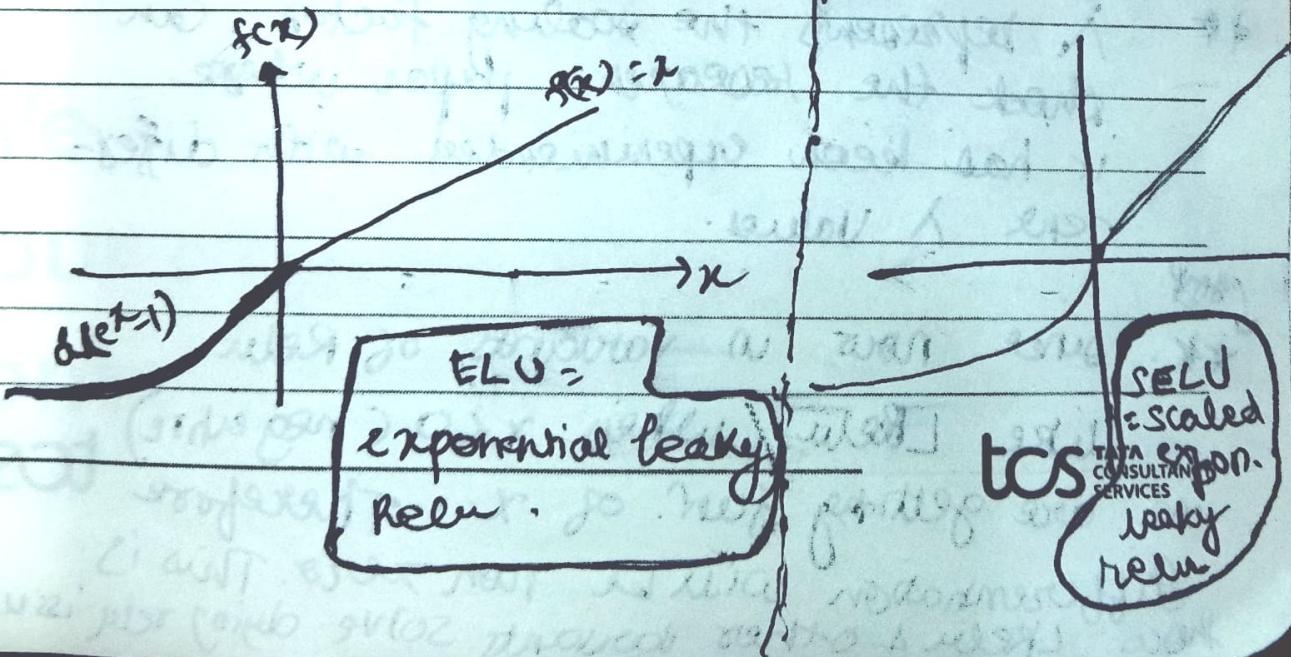
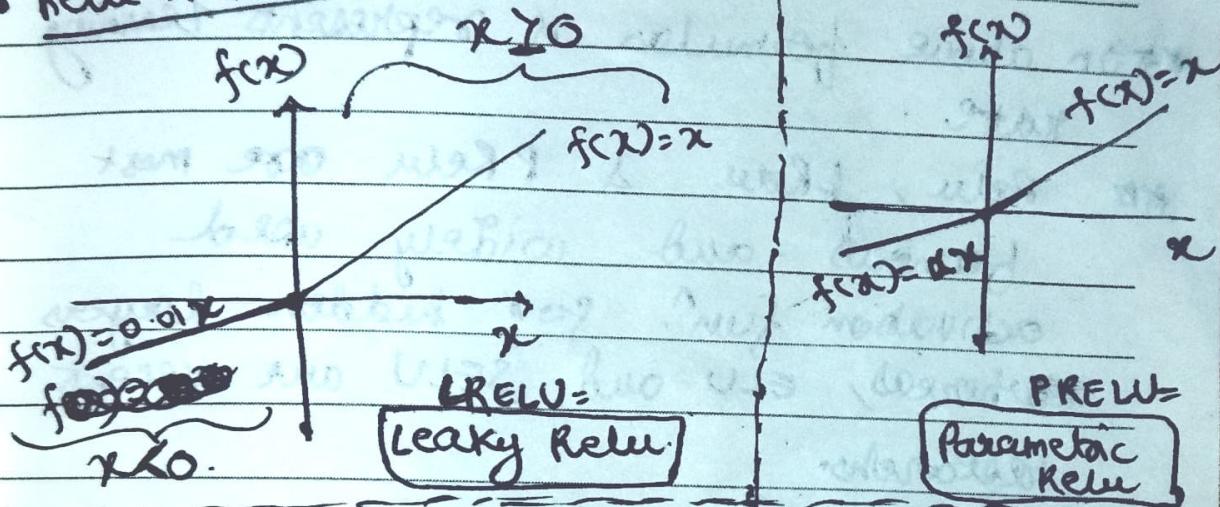
This is called dying neuron problem.

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w} \xrightarrow{\text{Chain rule}} \text{Dense}$$

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

Another disadvantage is that it is non-zero centric. ReLU passed from zero but it is not symmetric. That's why non-zero centric. But still it converges faster since we use Batch normalization.

### ReLU Variants:



## ReLU Variants formula comparison :

$$\text{ReLU} \Rightarrow \max(0, x)$$

$$\text{LReLU} \Rightarrow \max(0.01x, x)$$

$$\text{PReLU} \Rightarrow \max(\alpha x, x)$$

$$\text{ReLU} \Rightarrow \begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$

$$\text{SELU} \Rightarrow \begin{cases} \lambda x & x \geq 0 \\ \lambda \alpha(e^x - 1) & x < 0 \end{cases}$$

\* In above formulas  $\alpha$  represents learning rate.

# ReLU, LReLU & PReLU are most famous and widely used activation fun. for hidden layers whereas, ELU and SELU are recent researches.

#  $\lambda$  represents the scaling factor. Can check the research paper where it has been experimented with different  $\lambda$  values.

\*\* Since now in variants of ReLU like LReLU when  $x < 0$  (negative) we are getting 'fun' of  $x$ . Therefore differentiation will be non zero. This is how LReLU & others variants solve dying ReLU issue.

\*: Acc. to research value optimal  $\alpha$  and  $\lambda$  values are:

$$\lambda \approx 1.050700987355\ldots$$

$$\alpha \approx 1.67326324235\ldots$$

\*: We use ReLU and its variants in the hidden layer since they address the vanishing gradient and dying neuron issue.

\*: Can check the derivative graph of different variants of ReLU in the ipynb file shared by Sunny.