

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w_{\text{old}}} \quad \text{--- weight update formula}$$

$$\frac{\partial L}{\partial w_{\text{old}}} = \frac{\partial L}{\partial o_{31}} \times \frac{\partial o_{31}}{\partial o_{11}} \times \frac{\partial o_{11}}{\partial w_{\text{old}}}$$

* here there are 2 roots that you can basically see.

- ① Loss is dependent on o_{31} , o_{31} is dependent on o_{21}
 $o_{21} \rightarrow o_{11}$, $o_{11} \rightarrow w_1$
- ② Loss is dependent on o_{31} , o_{31} is dependent on o_{22}
 $o_{22} \rightarrow o_{11}$ & $o_{11} \rightarrow w_1$

now

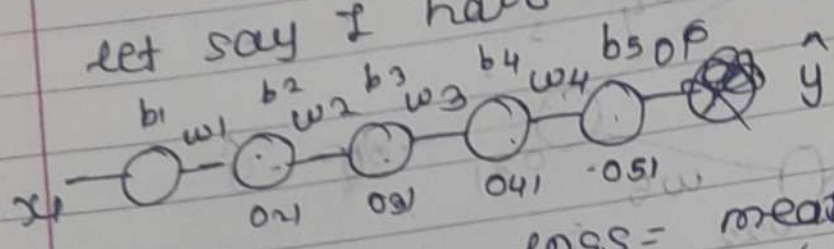
$$\frac{\partial L}{\partial w_{\text{old}}} = \left[\frac{\partial L}{\partial o_{31}} \times \frac{\partial o_{31}}{\partial o_{21}} \times \frac{\partial o_{21}}{\partial o_{11}} \times \frac{\partial o_{11}}{\partial w_{\text{old}}} \right]$$

$$\left[\frac{\partial L}{\partial o_{31}} \times \frac{\partial o_{31}}{\partial o_{22}} \times \frac{\partial o_{22}}{\partial o_{11}} \times \frac{\partial o_{11}}{\partial w_{\text{old}}} \right]$$

3)

vanishing gradient problem

let say I have a very deep neural net



loss = mean squared error

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

Initially we discussed about sigmoid activation fn
 to update w_1

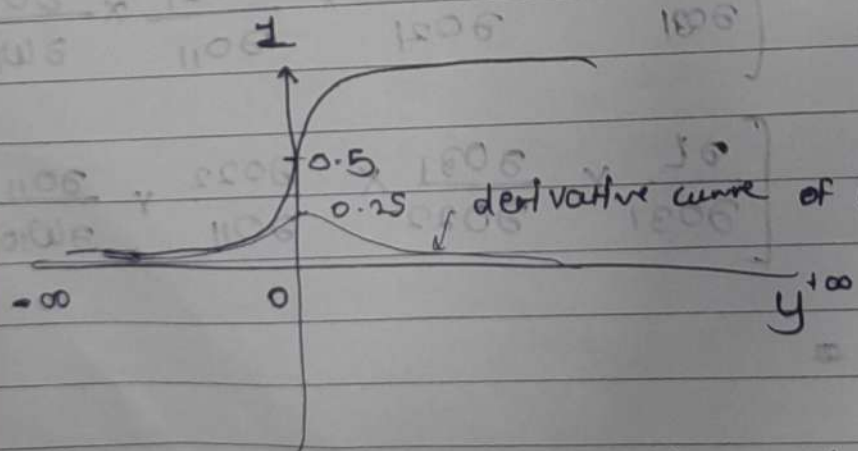
$$w_{new} = w_{old} - \frac{\partial L}{\partial w_{old}}$$

$$\frac{\partial L}{\partial w_{old}} = \frac{\partial L}{\partial 0.5} \times \frac{\partial 0.5}{\partial 0.41} \times \frac{\partial 0.41}{\partial 0.31} \times \frac{\partial 0.31}{\partial 0.21} \times \frac{\partial 0.21}{\partial w_{old}}$$

we are using here sigmoid activation fn

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

if $0.5 = y$ & $< 0.5 = 0$



derivative curve of sigmoid fn

One imp property of sigmoid fn whenever we try to find out derivative of sigmoid it will be ranging from 0 to 0.25

$$0 \leq \sigma(y) \leq 0.25$$

- bcoz of this what will happen
- In each and every value sigmoid is getting used.

e.g. $0.51 = \sigma[0.41 \times w_4 + b]$

on top of it we applying an activation fn.

- so in backpropagation when I'm finding the derivative of the value will always be ranging b/w 0 to 0.25

e.g. in eqn ② we will be getting

$$\frac{\partial L}{\partial w_{new}} = 0.25 \times 0.15 \times 0.10 \times 0.05 \times 0.02$$

- ~~if~~ the values will keep on decreasing bcoz as we are going to this chain rule, as we going to calculate derivative it is always going to reduce untill we go the end of the chain
- now what will happen bcoz of this since we are multiplying with smaller values we will be getting very small value. now bcoz of this small no.

$$w_{new} = w_{old} - \eta (\text{small no}) \rightarrow \text{no change in weights.}$$

- at one point of time $w_{new} \approx w_{old}$ it will hardly change and if it is hardly changed then weights are not getting updated & this situation where weights are not getting updated or it's just getting updated with small value. this problem is basically called as a vanishing gradient problem

then use another activation fn that is why we came up with another activation fn there are so many activation fns that we are going to learn.

- ① sigmoid
- ② Tanh
- ③ ReLU
- ④ leaky ReLU
- ⑤ pre-relu

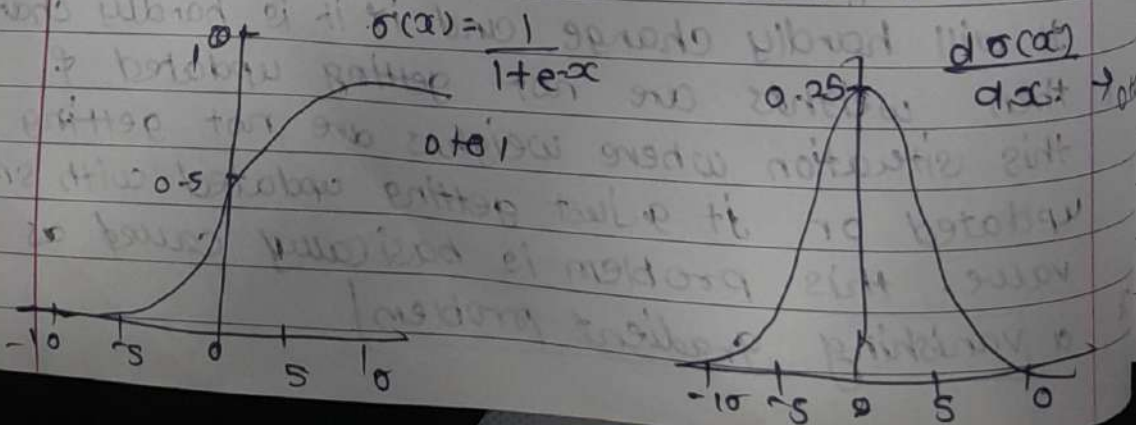
① sigmoid
→ frequently used in the beginning of DL.

advantages

- ① smooth gradient, preventing "jumps" in o/p value.
- ② o/p values bound between 0 & 1, normalizing the o/p of each neuron.
- ③ near predictions, i.e. very close to 1 or 0.

disadvantages

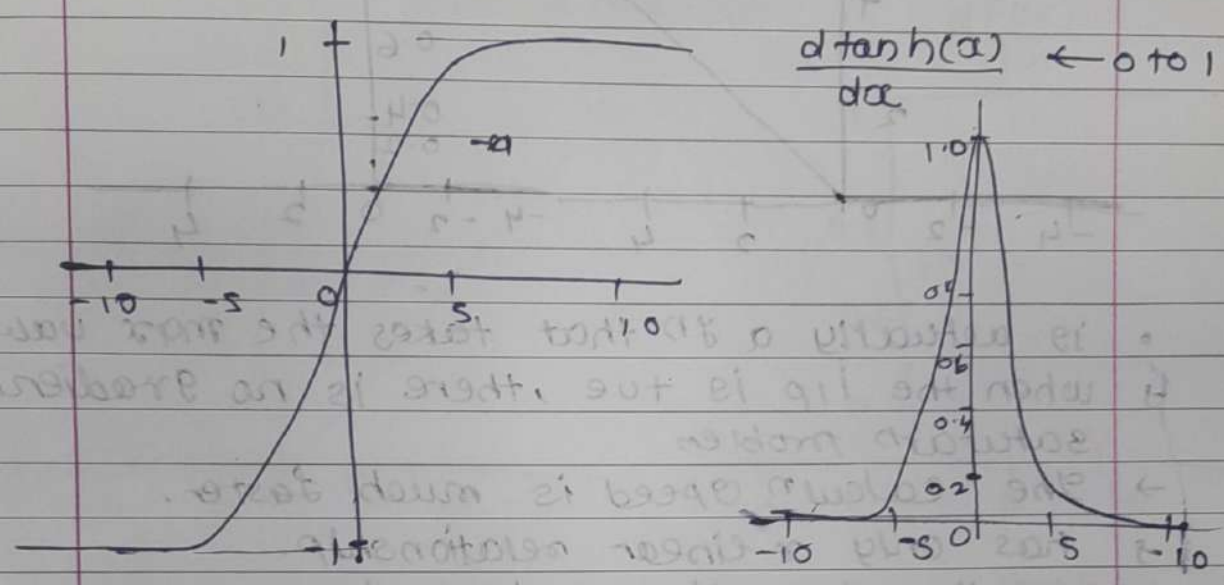
- ① prone to gradient vanishing
- ② fn o/p is not zero centered
- ③ power operations are relatively time consuming since there is an exponential fn



② tanh η (hyperbolic tangent η)

$$\tanh x = \frac{e^x - e^{-x}}{e^x + e^{-x}} \quad \text{--- 1 to +1}$$

value ranges from -1 to +1

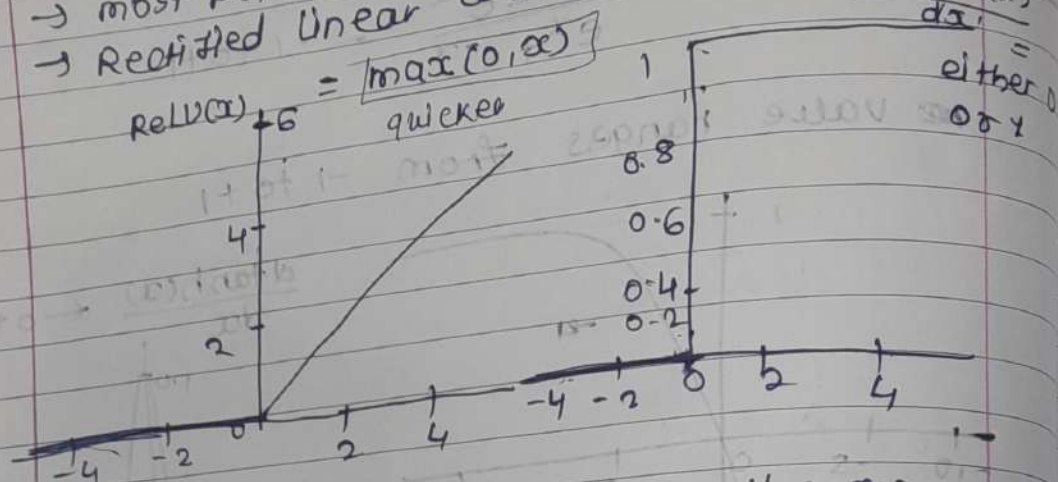


- does this prevent vanishing gradient problem?
- still there will a issue if we go on constructing a deep neural nw. or deep deep neural nw at one point of time there may be chances that vanishing G. P. still exist. bcoz of that we also not using tanh activation

③ The o/p interval of tanh is $[-1, 1]$ the whole η is 0-centric, which is better than sigmoid

In general binary classifⁿ problem, the tanh η is used for the hidden layer & sigmoid η is used for o/p layer

- ⑧ ReLU is the most popular and simple activation fn
 → Rectified Linear Unit



- is actually a fn that takes the max value
- when the i/p is +ve, there is no gradient saturation problem
- the calculation speed is much faster.
- has only a linear relationship.
- whether it is forward or backward.
- it is much faster than sigmoid & tanh. (sigmoid & tanh need to calculate the exponents, which will be slower)

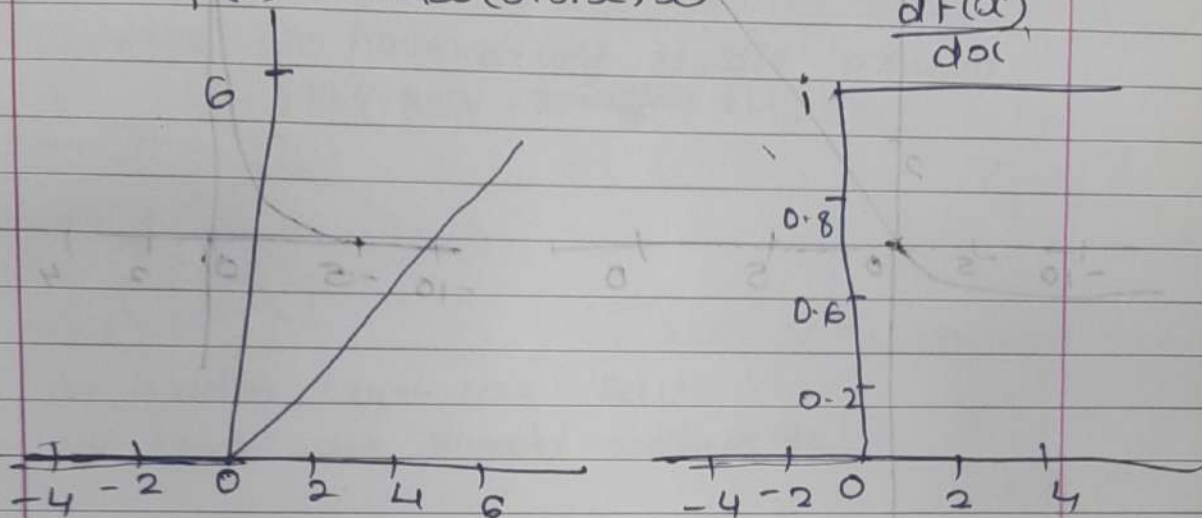
disadvantages

- when the -i/p is -ve, ReLU is completely inactive, which means that once a -ve number is entered, ReLU will die.
- In this way, in forward propagation, it is not a problem. Some areas are sensitive & some are insensitive.
- But in the B.P., if we enter a -ve number, the gradient will be completely zero, which has the same problem as the sigmoid & tanh fn.

- we find that the slope of the ReLU ϕ is either 0 or a +ve number, which means that the ReLU ϕ is not a 0 cent ϕ

(4) Leaky ReLU ϕ

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



- In order to solve the dead ReLU problem, people proposed to set the first half of ReLU $0.01x$ instead of 0.

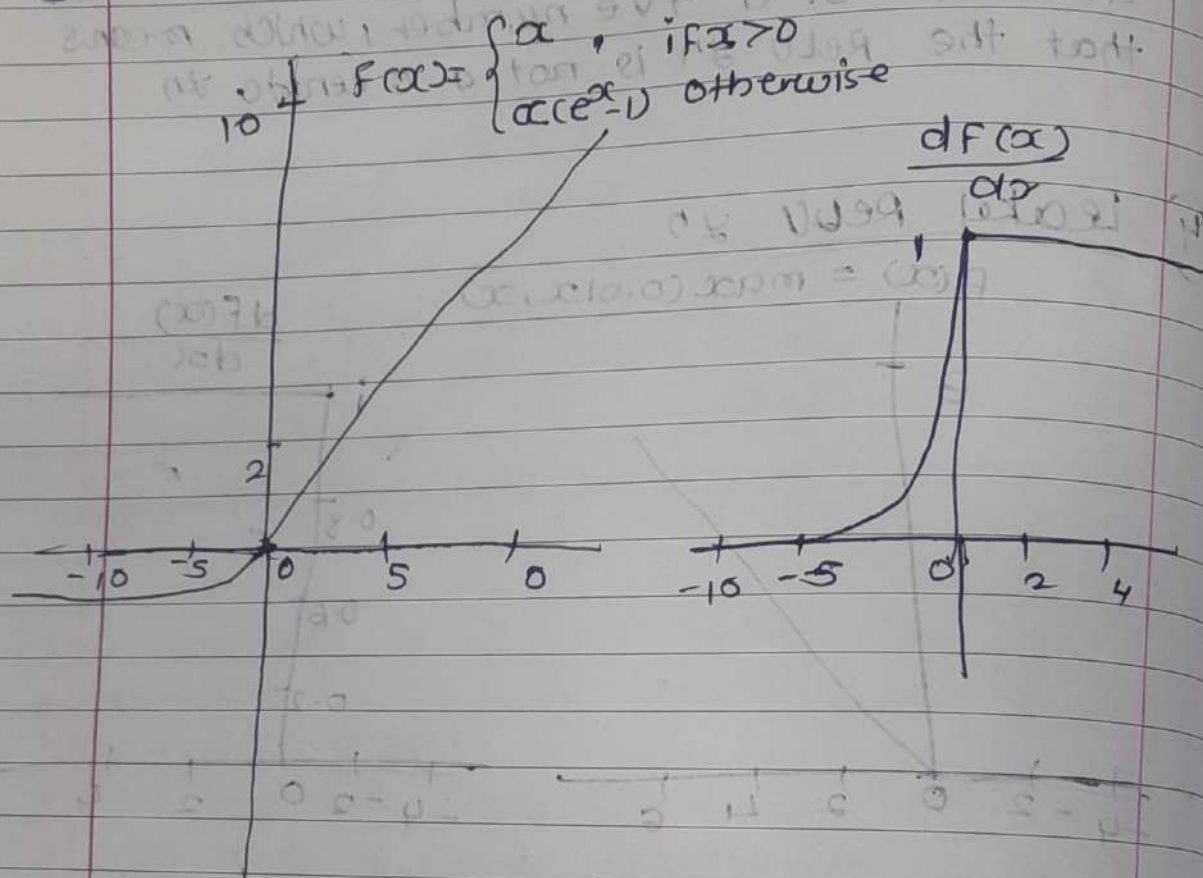
Another intuitive idea is a parameter-based method.

Parametric ReLU: $f(x) = \max(\alpha x, x)$

which α can be learned from B.P.

In theory Leaky ReLU has all the advantages of ReLU. Plus there will be no problems with dead ReLU, but in actual it has not been fully proved that Leaky ReLU is always better than ReLU.

③ ELU (Exponential Linear Units) ph.



ELU is also proposed to solve the problems of ReLU, obviously, ELU has all the advantages of ReLU.

- NO Dead ReLU issues.
- The mean of the OIP is close to zero, zero centered.
- One small problem is that it is slightly more computationally intensive, similar to leaky ReLU although theoretically better than ReLU.
- Good evidence in practice that ELU is always better than ReLU.

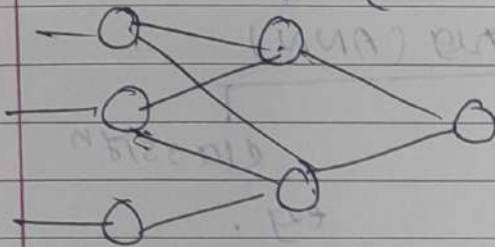
⑥ preReLU

$$f(y_i) = \begin{cases} y_i & \text{if } y_i > 0 \\ a_i y_i & \text{if } y_i \leq 0 \end{cases}$$

* Technique which Activation fn we should used?

- Suppose we have binary classifn problem

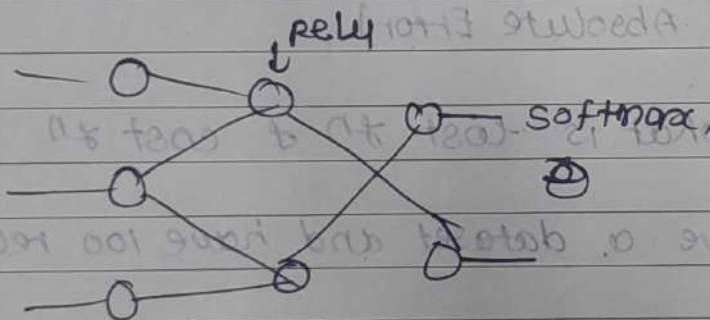
↓ HL' (ReLU, preReLU, ELU)



- In hidden layer use ReLU
o/p layer use sigmoid activation fn

multiclass classifn

use Softmax activation fn in the o/p layer.



- In the case of regression problem.

