

Deep learning

ML vs DL

OR without guidance

ML = Learning with guidance.

DL = Learning without guidance

↓ (using un-SP algo) (learn with guidance
thrives or work on big data. ↓
with SP algo)

why DL?

→ Data Explosion.

→ GPU's (Hardware Power)

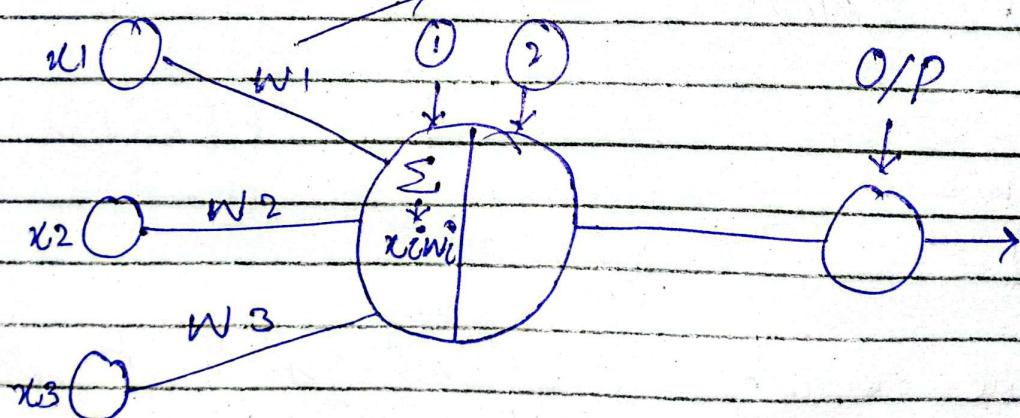
today GPU's and TPU's can
train huge data in hours
not in decade.

→ Automatic Feature learning.

→ Better performance

→ Works on every type of
data.

Perceptrons: weights



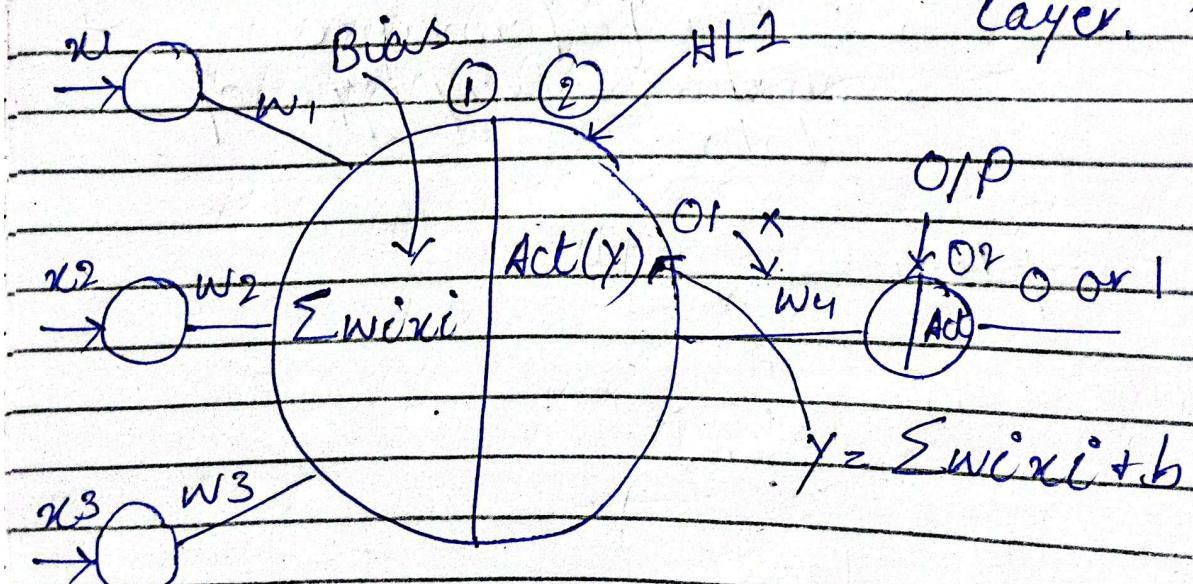
weights will help your neurons to activate at what level.

which weight should have what value so that a neuron is activated till that level or not.

In Simple:

weights says that how much neuron should activated or de activated.

→ we add Bias to overcome the situation where weight is initialized to 0, and in every situation, bias is required in every layer.



One example of activation function is sigmoid activation func, for Binary classification.

$$\text{Sigmoid} = \frac{1}{1+e^{-y}} \Rightarrow \frac{1}{1+e^{-(\sum w_i x_i + b)}}$$

this was for Binary classification
and Binary output.

For Regression we use linear
Activation function.

The entire process we have done
is called forward propagation

For instance if \hat{y} (Predicted
value for 7, 3, 7) is 0 but
should come 1) then:

$$0 \rightarrow \hat{y} \quad y \rightarrow \text{truth value} = 1$$

$$y - \hat{y} = 1 - 0 = 1 \rightarrow \textcircled{1}$$

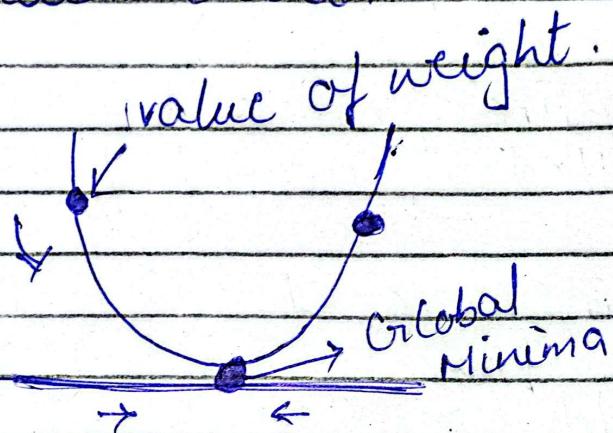
\uparrow
loss function

and our aim is always to minimize
the difference, what we'll do if
our difference is very huge? \Rightarrow then
we'll do back propagation.

In back propagation we update
the weights.

To do back propagation we use optimizers, they will ensure that each value of our weight is updated in order to get satisfactory output.

→ Gradient Descent:



will make sure that the coefficient is updated to or till to the global minima.

this G.D is type of optimizer.

→ Conclusion:

- ① I/P layer
- ② weights will get added.
- ③ Bits
- ④ Activation func.
- ⑤ output.

$y = mx + c$
Forward propagation

⑥ Loss function $\Rightarrow \{ \hat{y} - y \}$

we need to
minimize it to do
that we use:

⑦ Optimizers

↳ updates the weights.

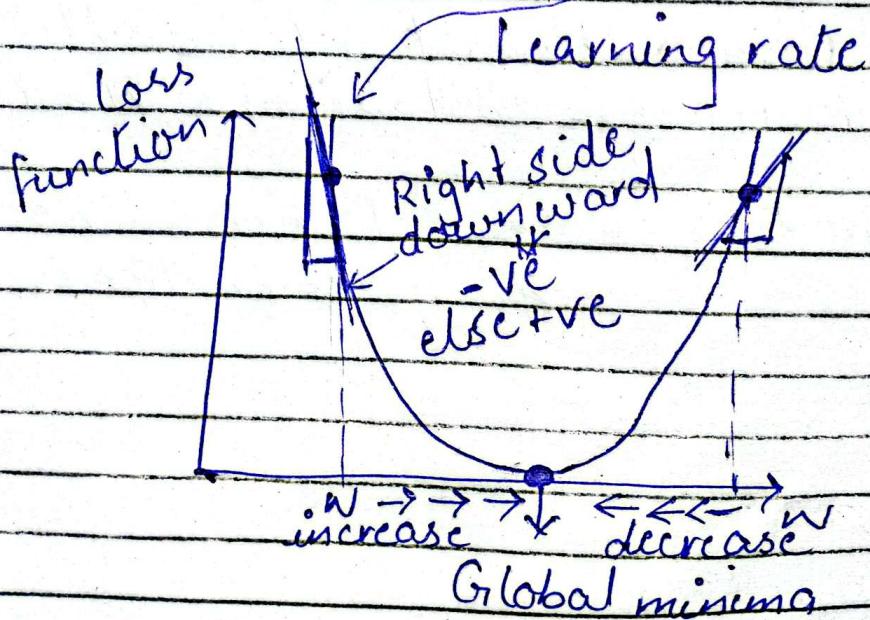
↳ By using the concept
of Gradient Descent.

Back propagation.

→ Back propagation:

Weight updation formula:

$$W_{\text{new}} = W_{\text{old}} - \eta \frac{dL}{dW_{\text{old}}} \quad \begin{matrix} \text{Calculating} \\ \text{Slope} \end{matrix}$$



for (+ve) :

$$W_{\text{new}} = W_{\text{old}} - \eta \text{ (+ve)}$$

$$W_{\text{new}} = W_{\text{old}} + \eta$$

Now my
 $W_{\text{new}} > W_{\text{old}}$

for (-ve) :

$$W_{\text{new}} = W_{\text{old}} - \eta \text{ (+ve)}$$

$$W_{\text{new}} = W_{\text{old}} - \eta \text{ (-ve)}$$

$$W_{\text{new}} < W_{\text{old}}$$

$\eta : \rightarrow \approx$ Small number
↳ it will slowly converge
to global minima (Best).

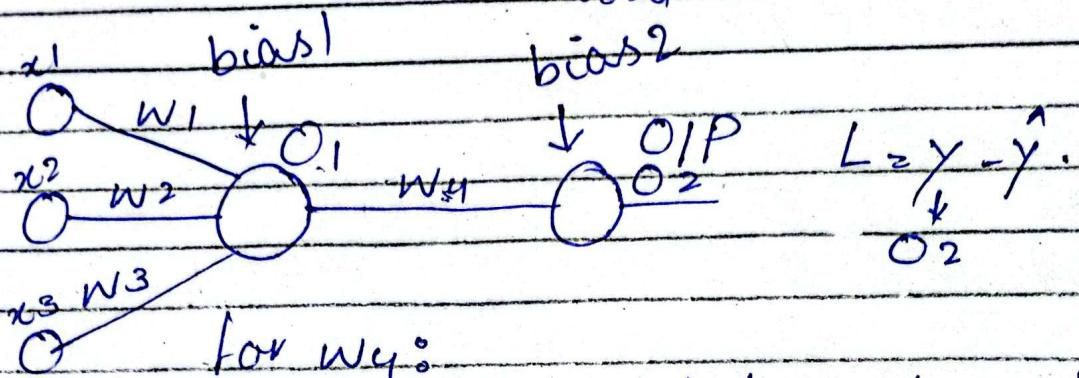
$\eta \rightarrow$ Larger number

↳ it might jump here and
there and might not be able
to reach global minima (Worst)

η should be = 0.001 or else like
this. (Best practice)

② Chain Rule of Derivative:

$$W_{\text{new}} = W_{\text{old}} - \eta \frac{dL}{dW_{\text{old}}} \rightarrow ①$$



if we want to update value of w_4 we'll take $\frac{dL}{dW_{\text{old}}}$ from eq, ①

$$\frac{dL}{dW_{\text{old}}} = \frac{dL}{dO_2} * \frac{dO_2}{dW_{\text{old}}}$$

eq ① will same formula for bias:

$$\text{bias}_{2\text{new}} = \text{bias}_{2\text{old}} - \eta \frac{\frac{dL}{d\text{bias}}}{d\text{bias}_{2\text{old}}}$$

for w_2 :

$$\frac{dL}{dW_{\text{old}}} = \frac{dL}{dO_2} * \frac{dO_2}{dO_1} * \frac{dO_1}{dW_{\text{old}}}$$

For updating w_2 :

$$w_{2,\text{new}} = w_{2,\text{old}} - \eta \frac{dL}{dW_{2,\text{old}}} \rightarrow ①$$

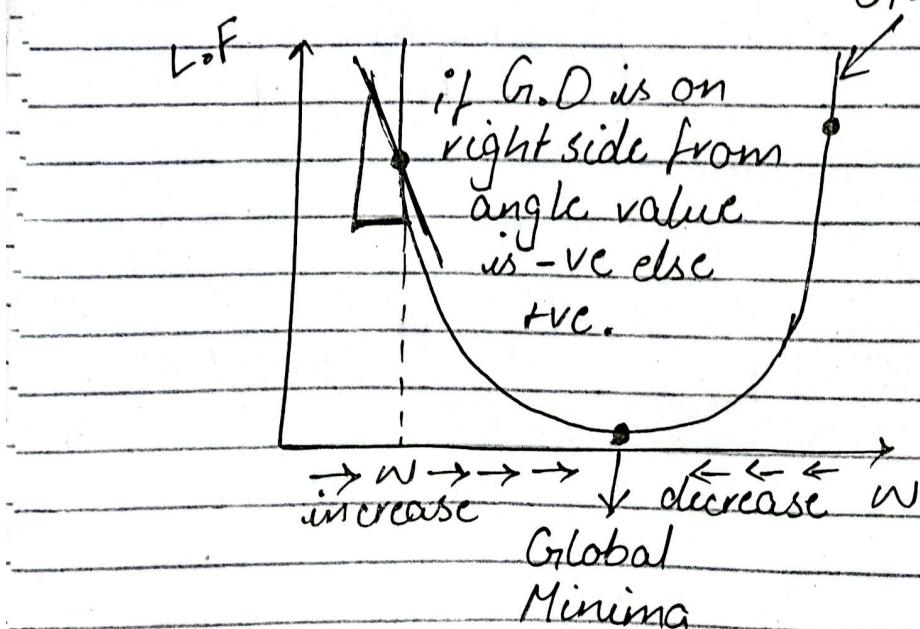
from ① take $\frac{dL}{dW_{2,\text{old}}}$.

$$\frac{dL}{dW_{2,\text{old}}} = \frac{dL}{O_2} + \frac{dO_2}{dO_1} \times \frac{dO_1}{dW_{2,\text{old}}}.$$

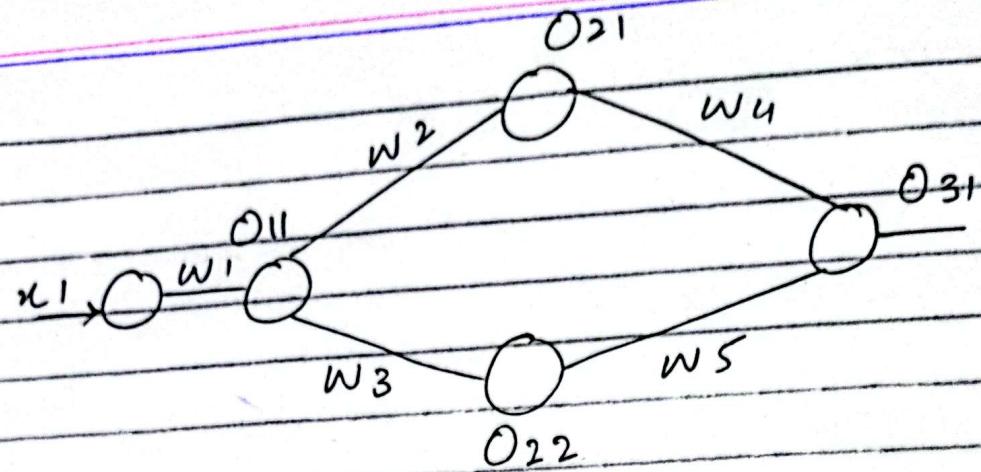
remember $\frac{dL}{dW_{\text{old}}}$ in every weight

update formula will give a slope or a tangent line in graph and from there we'll check whether the value is +ve or -ve if value is +ve then weight will decrease towards global minima else increases.

G.O.D



For 2 routes:



Now for w_1 :

$$\Delta w_{\text{new}} = w_{\text{old}} - \eta \frac{dL}{dw_{\text{old}}} \rightarrow ①$$

From ①

$$\frac{dL}{dw_{\text{old}}} = \frac{dL}{dO_{31}} * \frac{dO_{31}}{dO_{21}} * \frac{dO_{21}}{dO_{11}} * \frac{dO_{11}}{dw_{\text{old}}}$$

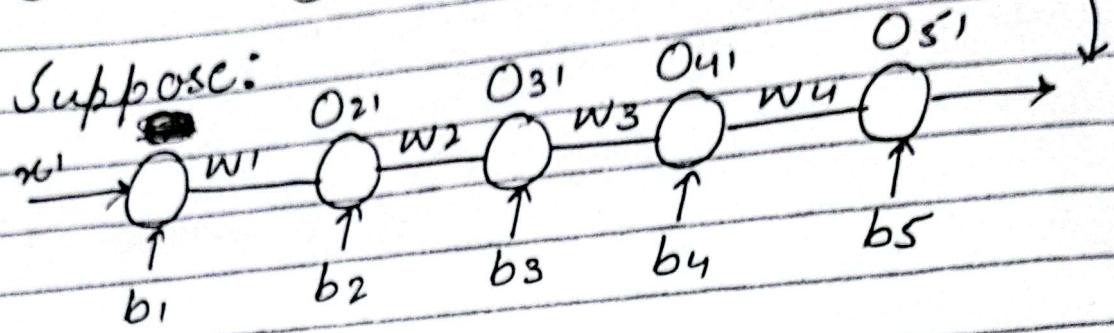
+

$$\frac{dL}{dO_{31}} * \frac{dO_{31}}{dO_{22}} * \frac{dO_{22}}{dO_{11}} * \frac{dO_{11}}{dw_{\text{old}}}$$

③ Vanishing Gradient Problem:

$$L = y - \hat{y}$$

Suppose:



$$w_{\text{new}} = w_{\text{old}} - \eta \frac{dL}{dw_{\text{old}}} \rightarrow ①$$

from ① take:

$$\frac{dL}{w_{\text{old}}} = \frac{dL}{dO_5} * \frac{dO_5}{dO_4} * \frac{dO_4}{dO_3} * \frac{dO_3}{dO_2} * \frac{dO_2}{dw_{\text{old}}}$$

Now the question is \Rightarrow Is activation function is implementing in every output?

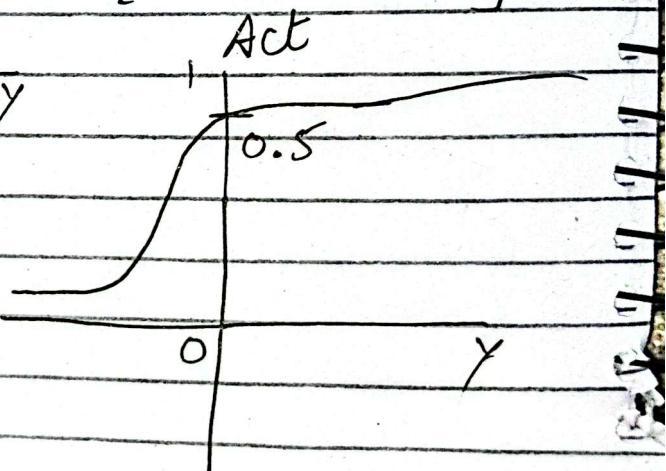
The Answer is yes.

How? \rightarrow

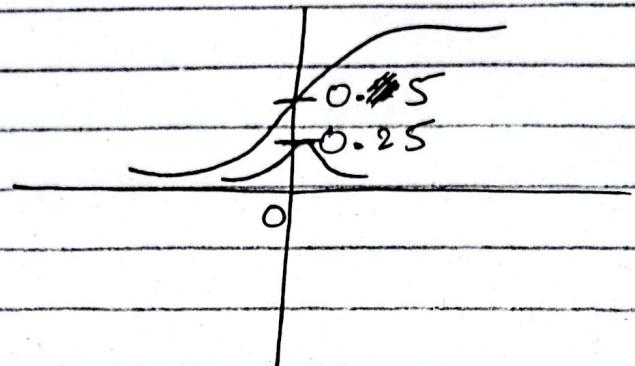
$$\text{Let take } O_5 \Rightarrow \text{Act}(O_4 * w_4) + b$$

$$\text{Act} = \frac{1}{1 + e^{-y}}$$

Act func Graph:



But the derivative of an activation function will always give value b.w 0 to 0.25.



if values of each derivation in eq, ② is like 0.24, 0.19, 0.11, 0.5, 0.001
then it's fine but what if value decreases consistently like:

$$dL = 0.24 + 0.19 + 0.11 + 0.5 + 0.001$$

w_{old}

Then:

The value of w_{new} \approx small value

if it is:

$$w_{\text{new}} = w_{\text{old}} - \eta \text{ (small value)}$$

also a small value

w_{new} = Hardly changed value.

This is called gradient descent problem.

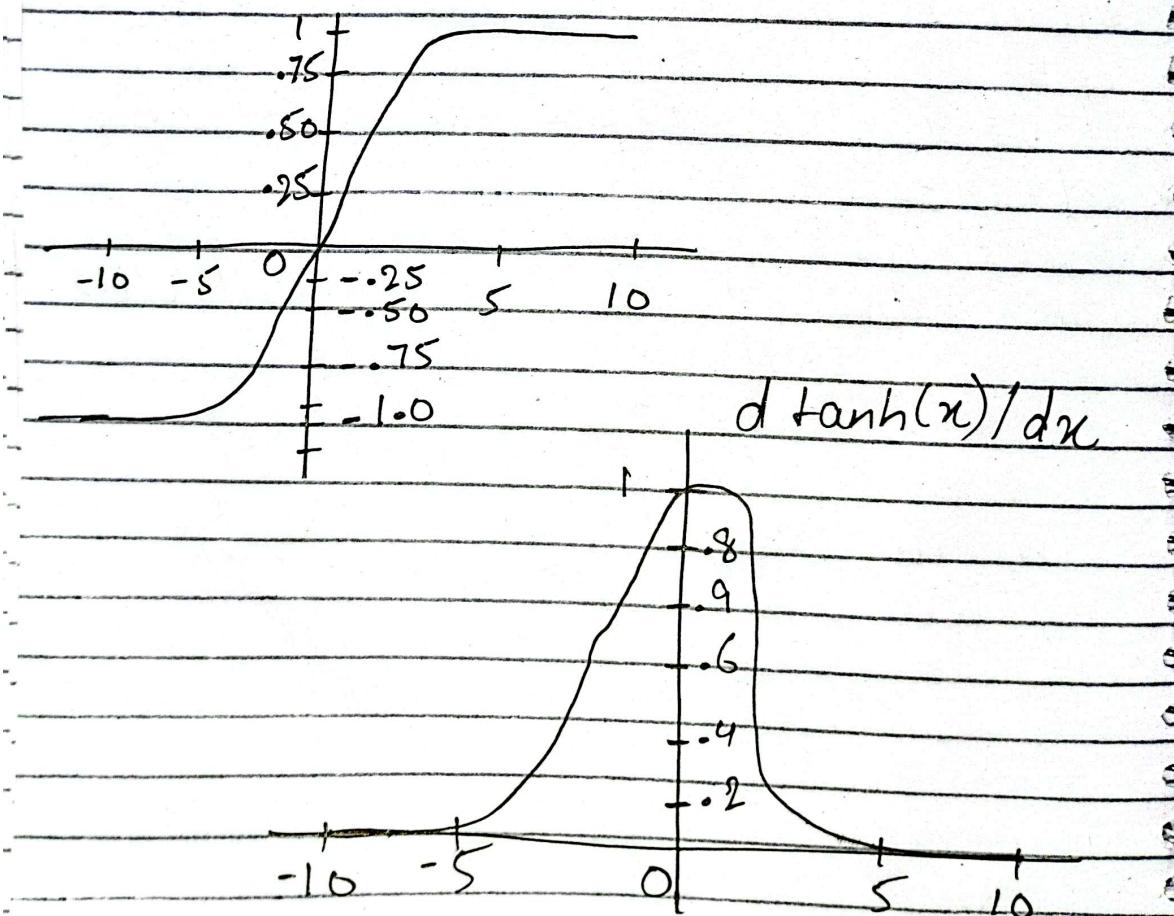
The solution to this is to try other Activation functions like:

- ① Sigmoid.
- ② Tanh.
- ③ ReLU.
- ④ Leaky ReLU.
- ⑤ PreReLU.

→ Types of Activation Functions:

- ① Tanh:

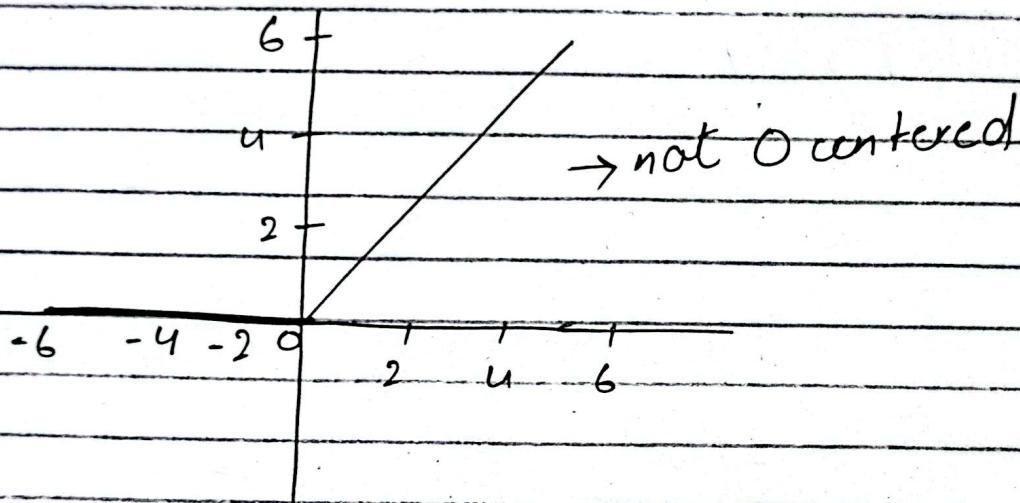
$$\tanh = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$



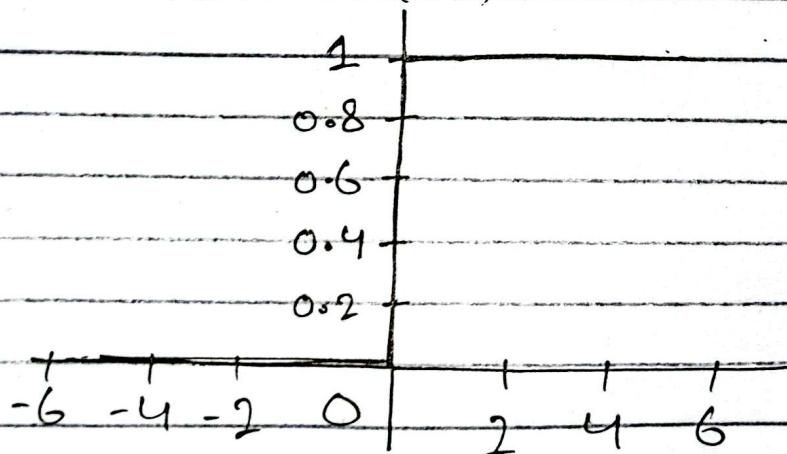
It's activation graph is 0 centered and its derivative of activation func graph is b.w 0 and 1, Due to which value of weight is efficiently updates, which can help with vanishing gradient problem, but still might not helpful in deep neural Network.

② ReLU (Rectified linear unit):

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



derivative (ReLU (x)) :



ReLU is faster than tanh and sigmoid function as they both need to calculate exponent which is slower.

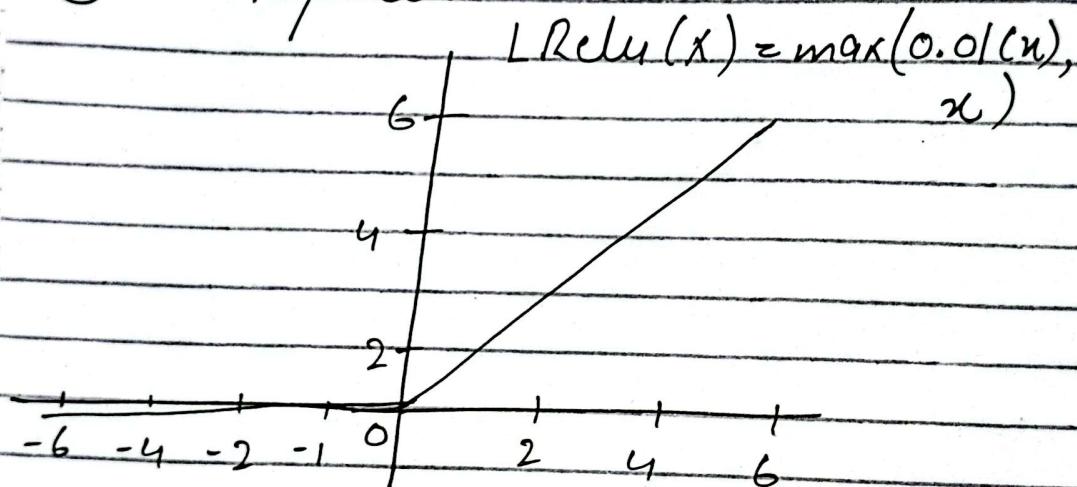
Disadvantage:

→ Don't handle $(-ve)$ value
↳ means ReLU will die once $(-ve)$ input is entered.

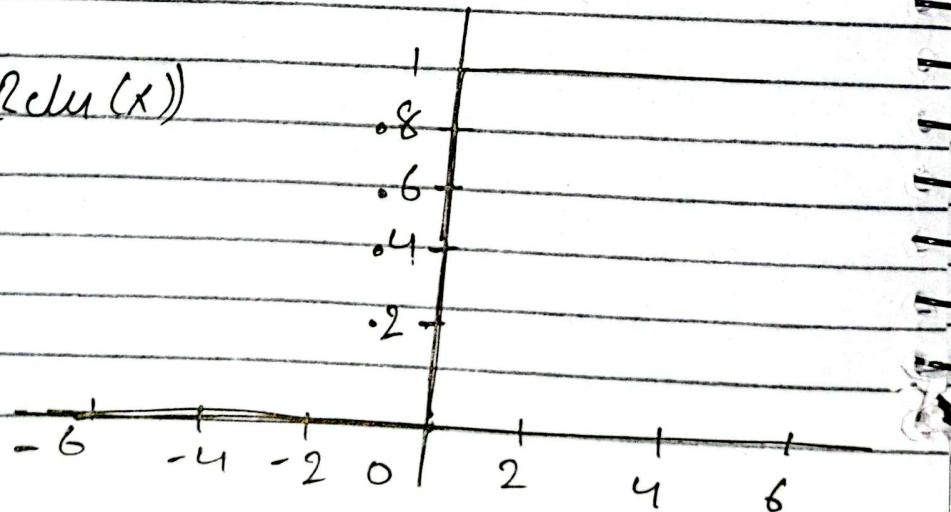
In order to solve dead ReLU or dead we'll use Leaky ReLU.

Neuron

③ Leaky ReLU:

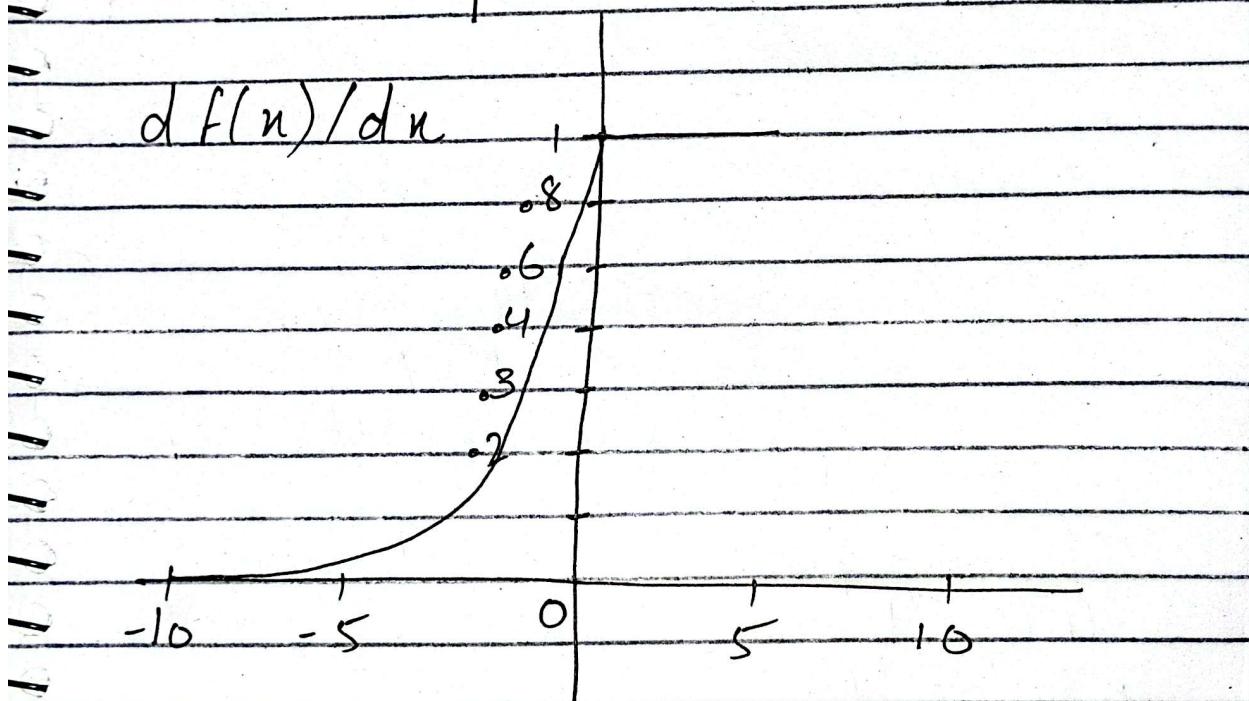
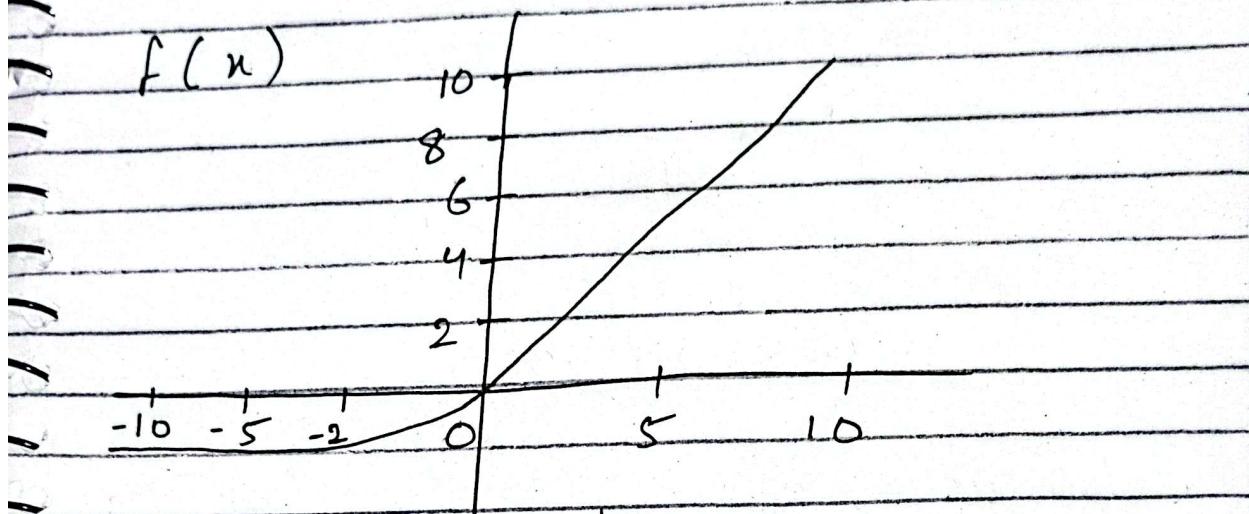


$$\frac{d(LReLU(x))}{dx}$$

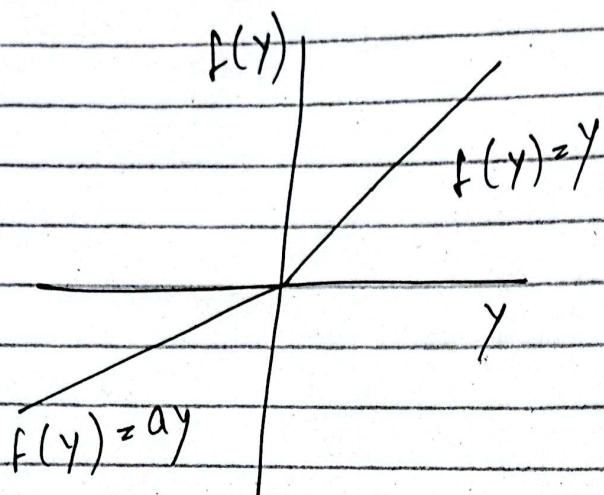


(ii) ELU (Exponential linear unit)

$$f(x) = \begin{cases} x, & \text{if } x > 0 \\ \alpha(e^x - 1), & \text{otherwise} \end{cases}$$



(5) PRelu:



$$f(y) = \begin{cases} y, & \text{if } y > 0 \\ ay, & \text{if } y \leq 0 \end{cases}$$

Technique \Rightarrow which Activation function to use?

\rightarrow For Binary Classification:

1 \rightarrow Use (ReLU)

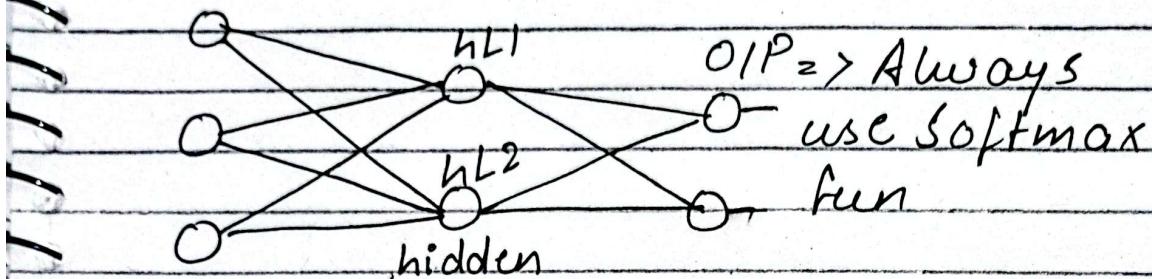
In hidden layer.

1 \rightarrow In output layer

use (Sigmoid.) only.

If convergence is not happening with (ReLU) use (PRelu), or (ELU) in hidden layer.

→ For Multiclass Classification:



In layers use ReLU or PReLU and ELU if convergence is not happening with (ReLU).

But try to use ReLU in hidden layer.

→ For Regression:

use linear Activation function in output layer and ReLU or any other ReLU in hidden layers.

Loss func:

$$L = \frac{1}{n} \sum (\hat{y} - y)^2 \quad | \quad \text{SVM models: } L = \max(0, 1 - y \cdot \hat{y})$$

Binary Cross-Entropy

$$L = -[y \log(\hat{y}) + (1-y) \log(1-\hat{y})]$$

Multi Class Classification

$$L = -\sum y_i \log(\hat{y}_i)$$

① Loss Functions:

Diff b/w Loss and Cost

$$\text{Loss} = \frac{1}{n} (y - \hat{y})^2 \quad \text{Data} = 100$$

$$\text{Cost Function} = \frac{1}{n} \sum_{i=1}^n (y - \hat{y})^2$$

→ For Regression:

→ MAE

→ MSE

→ RMSE

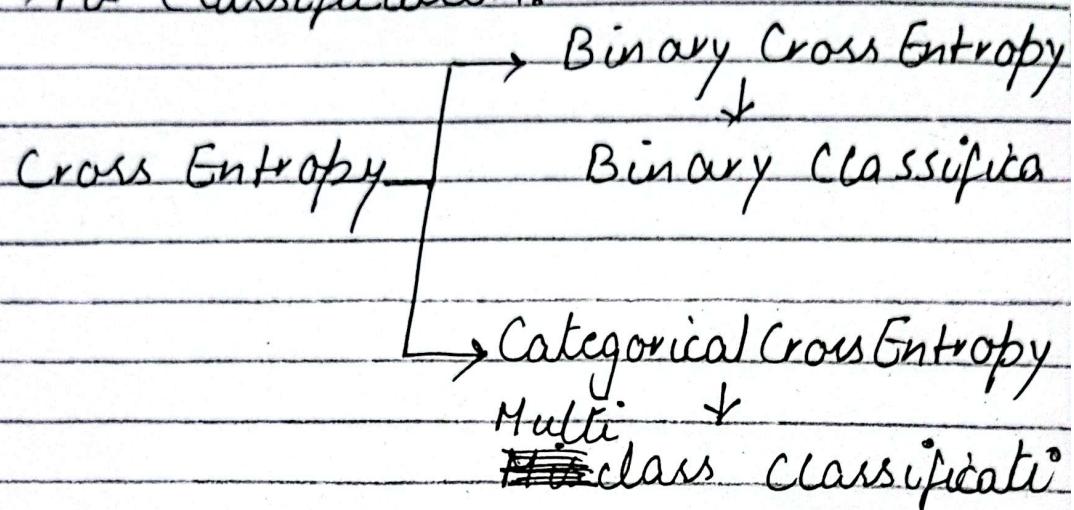
→ Huber Loss

MAE:

Disadvantage: Not Robust to Outliers

MAE → Robust to Outliers.

→ For Classification:



→ Binary Cross Entropy:

$$\text{Loss} = -y \log(\hat{y}) - (1-y) \log(1-\hat{y})$$

Logistic Regression

$$\text{Loss} = \begin{cases} -\log(1-\hat{y}) & \text{if } y=0 \\ -\log(\hat{y}) & \text{if } y=1 \end{cases}$$

$$\hat{y} = \frac{1}{1+e^{-x}}$$

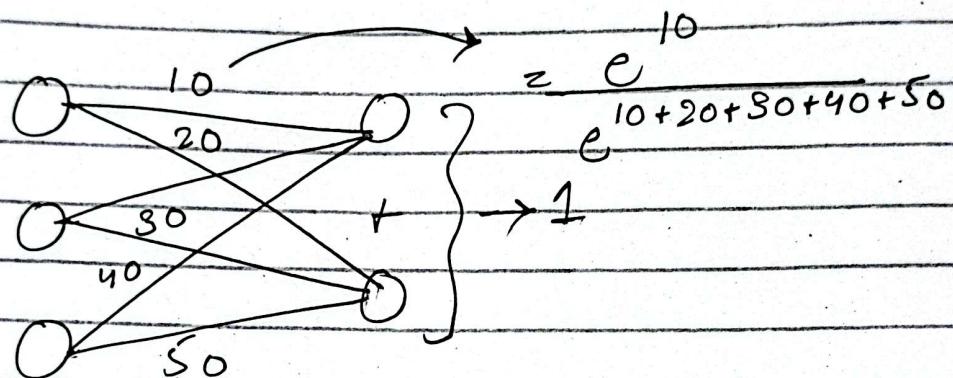
→ Categorical Cross Entropy:

$$L(y_i, \hat{y}_i) = -\sum_{j=1}^C y_{ij} \ln(\hat{y}_{ij})$$

y_{ij}^1 = Softmax Activation function.

$$\sigma(z) = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

} Softmax
} Activation



The output will be the one with higher probability in case of 2 or more outputs (Multiclass classification).

→ Conclusion: Categorical cross classification ↓ entropy

ReLU(hL), Softmax(OL) \Rightarrow Multiclass

ReLU(hL), Sigmoid(OL) \Rightarrow Binary

Regression Binary cross entropy

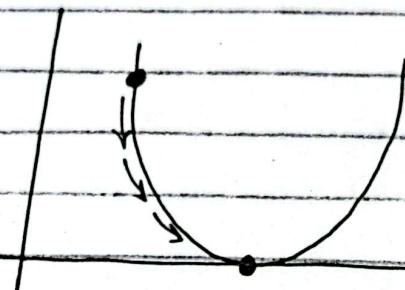
ReLU(hL), Linear act func(OL)

↓
Loss func: MSE, MAE, Huber loss

→ Types of Optimizers

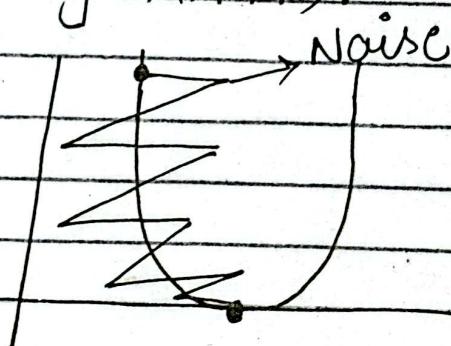
- ① Gradient Descent
- ② Stochastic Gradient Descent
- ③ Min Batch SGD
- ④ SGD with momentum.

GD =>



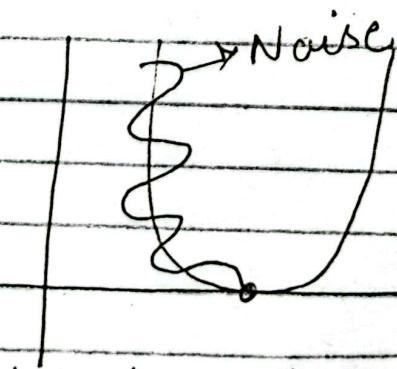
→ Requires Intense or Extensive Resource (Huge RAM).

SGD =>



Sends individual value from data in 1 Epoch. → Much slower

MB SGD =>



Sends a Batch from data in 1 epoch. → a little faster but still

have some noise.

How to Reduce Noise?

→ SGD with momentum.

As we know

$$w_t = w_{t-1} - \eta \frac{dL}{dw_{t-1}} \rightarrow ①$$

We use Exponential weighted Average

$$\begin{matrix} t_1, t_2, t_3, t_4, \dots, t_n \\ a_1, a_2, a_3, a_4, \dots, a_n \end{matrix}$$

$$v_{t_1} = a_1$$

$$v_{t_2} = \beta * v_{t_1} + (1-\beta) * a_2 \rightarrow ②$$

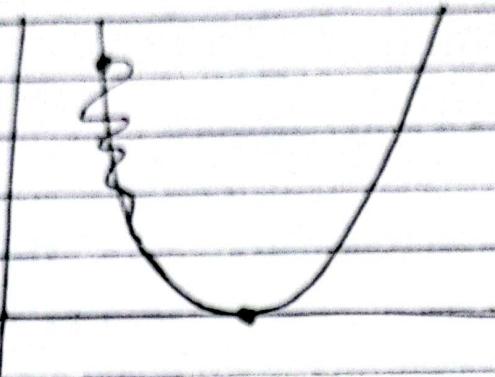
$\beta \Rightarrow$ Hyperparameter
→ ranges from 0 to 1.

$$v_{t_3} = \beta * v_{t_2} + (1-\beta) * a_3$$

By using β as hyperparameter we will decide which value to choose for updating weight as it will reduce noise.

$$\beta = 0.95 \Rightarrow \text{Put in } ②$$

$$v_{t_2} = \underbrace{(0.95) * v_{t_1}}_{\text{we'll choose this}} + \underbrace{(1-\beta(0.95)) * a_2}_{0.05}$$



Final E.W.A.:

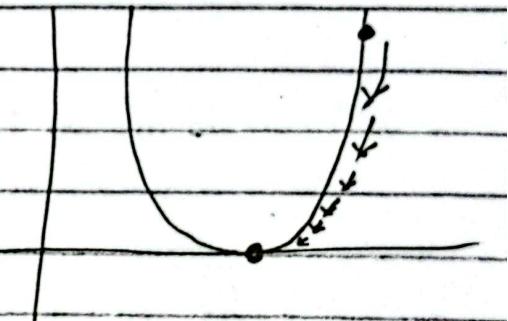
$$w_t = w_{t-1} - \eta Vdw_t$$

$$Vdw_t = \beta * Vdw_{t-1} + (1-\beta) * \frac{dL}{dw_{t-1}}$$

⑤ Adagrad \Rightarrow Adaptive Gradient Descent

$$w_t = w_{t-1} - \eta \frac{dL}{dt_{t-1}}$$

$$w_t = w_{t-1} - \eta' \frac{dL}{dt_{t-1}}$$



$$\eta' = \eta$$

$\sqrt{\alpha_t + \epsilon}$ \Rightarrow here value is decreasing as we reach Global minima,

$\alpha_t = \sum_{i=1}^t \left(\frac{dL}{dw_i} \right)^2$, due to which our learning rate is adaptive. So that the learning rate is not fixed.

what if αt increase greatly?
then: $\rightarrow w_t \approx w_{t-1}$
vanishing GrD Problem

(6) AdaDelta and RMSProp

$$\eta' = \frac{\eta}{\sqrt{Sdw_t + \epsilon}} \rightarrow \text{now it will increase by } \frac{1}{\sqrt{\epsilon}}$$

$$Sdw_t = \beta * Sdw_{t-1} + (1 - \beta) \left(\frac{\nabla L}{\|w_{t-1}\|} \right)^2$$

(7) Adam Optimizer: (Best Optimizer)

Momentum + RMSprop.

So Finally best option or formula
to update weight value:

$$w_t = w_{t-1} - \eta' Vdw$$

$$b_t = b_{t-1} - \eta' Vdb \rightarrow \text{to update bias value.}$$

So ~~What~~ what we did:

- Smoothing (Reducing Noise).
- Learning Rate (Making it Adaptive).

ANN Implementation

→ Tensorflow important importing
Code

from tensorflow.keras.models import
Sequential

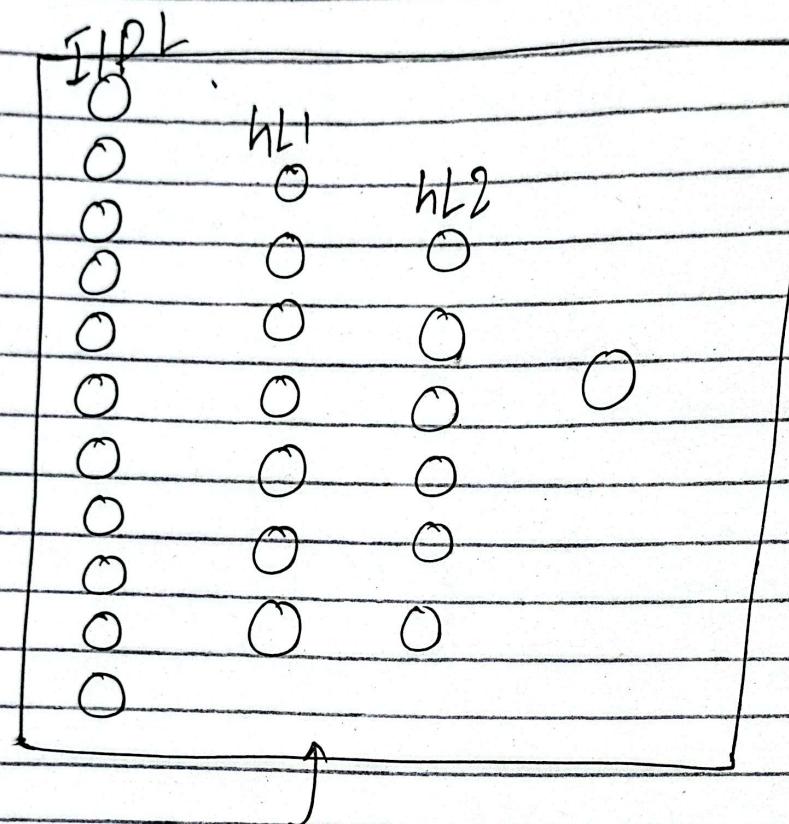
from tensorflow.keras.layers import
Dense, ReLU, Dropout → helps with overfitting
↓ ↓ uses it inside neurons or
helps to create hidden layers.
neurons or hidden
layer and also
I/P or O/P layers

classifier = Sequential()

Creates a space
to make neural
network

Now to add input, hidden and
Output layer we'll use Dense

Dense
classifier.add(Dense(unit=11, activation
= 'relu')).



classifier.add(Dense(unit=7, activation='relu'))

Creates 1 hidden layer with 7
units if also we make another layer
with 6 units then

CNN (Convolution Neural Network)

CNN is a type of DL model that is mainly used to analyze images.

→ It works by automatically detecting important features (like edges, shapes and textures) from images using convolutional layers, so it can recognize patterns and objects without manual feature extraction.



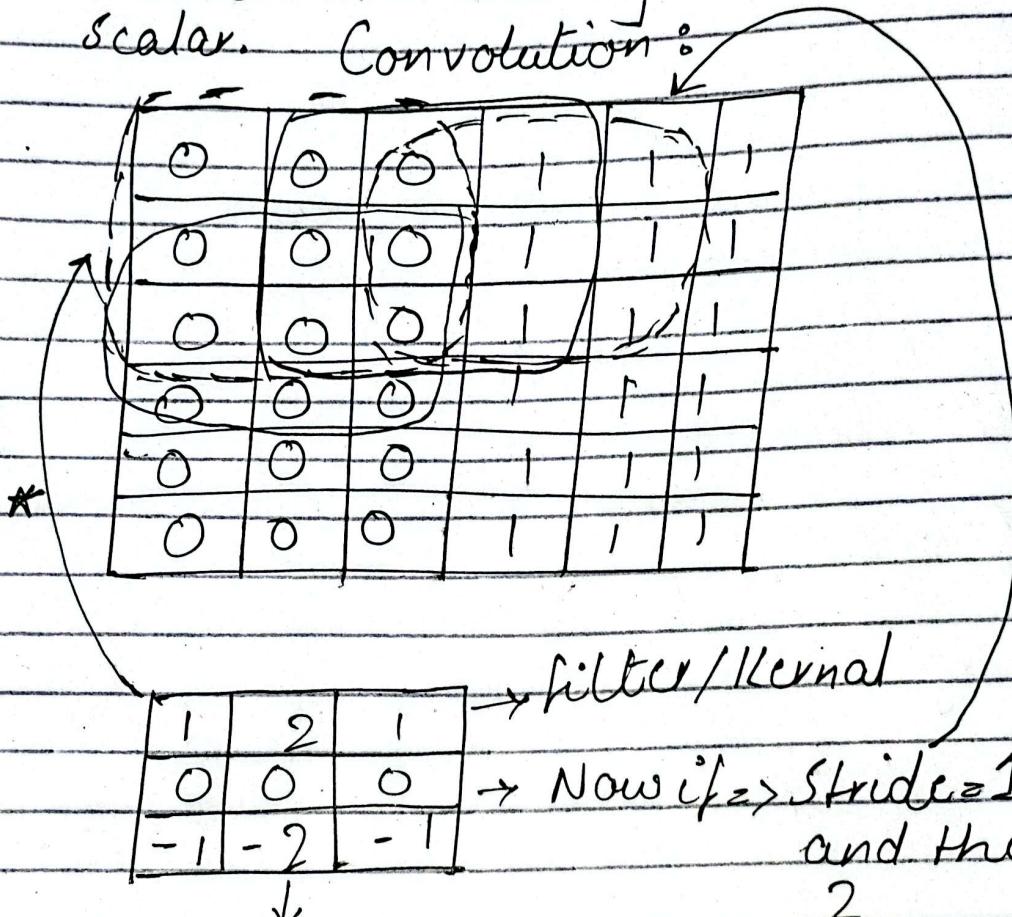
Images: each pixel ranges from 0 - 255

(5 x 5) pixel image
row column

0 → Black
255 → white

① Convolutional Operations

In first step of CNN we try to convert the pixels of an image into 0 and 1 using min, max scalar. Convolution:



$$1(0) + 2(0) + 1(0) + 0(0) + 0(0) + 0(0) \\ + (-1)(0) + (-2)(0) + (-1)(0) \geq 0$$

$\Rightarrow 0$

$$0 + 2 + 1 + 0 + 0 + 0 - 2 - 1 \\ \geq 0$$

\rightarrow Why we multiply filter with image pixels

1	0	-1
2	0	-2
1	0	-1

\Rightarrow vertical edge filter
(Sobel)

\rightarrow Filters / Kernel helps to detect features like: edges, textures and patterns.

$$\text{if } n = 6 ; f = 3$$

what will be the size of O/P image

$$= n - f + 1$$

$$n - f + 1 \Rightarrow 6 - 3 + 1 = 4 \times 4$$

255	0	0	255
255	0	0	255
255	0	0	255
255	0	0	255

$\rightarrow 4 \times 4$

As we see that the size of our image is decreased due to which it means that some of the information is lost to prevent that we use a concept \Rightarrow Padding.

Padding \Rightarrow Building a compound around images.

New formula to find size of O/P image after using padding:

$$\begin{aligned} &= n + 2p - f + 1 \\ &= 6 + 2(1) - 3 + 1 \\ n = 6 \quad &= 8 - 3 + 1 \\ p = 1 \quad &= 9 - 3 \\ f = 3 \quad &= 6 \end{aligned}$$

We need to update the values of filter on the base of input image by using backpropagation.

After getting output image we apply ReLU activation function on each value of image's pixel.

If we give the value of Stride > 1 then O/P image size formula will be:

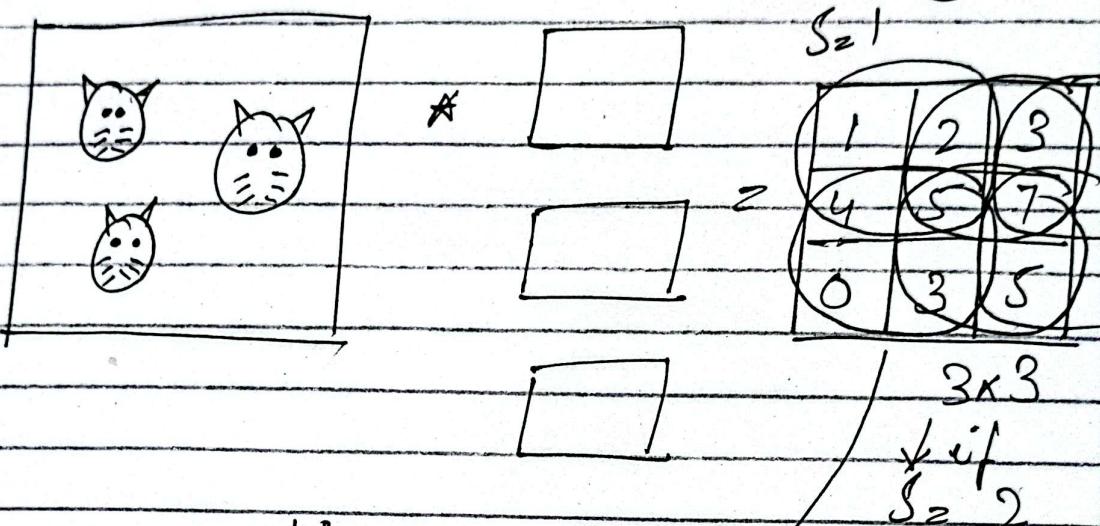
$$\frac{n + 2p - f + 1}{s}$$

② Max pooling:

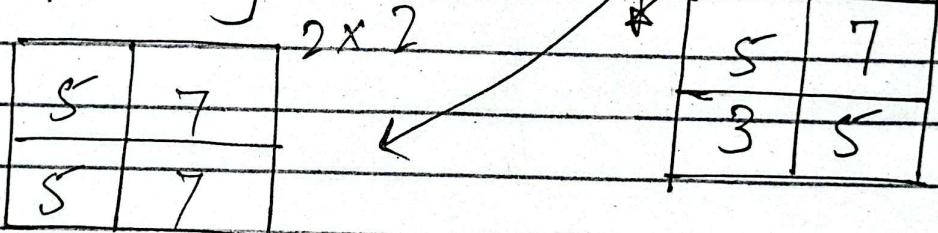
Other 2 types of Pooling:

→ Avg pooling

→ Min pooling



Max Pooling:



Max Pool will contain the highest values of pixels in order have clear data.

③ Flattening layer:

Just Converting the Pooling into 1D Data.

or combine all O/P. then we can throw this as I/P in NN

Conclusion:

