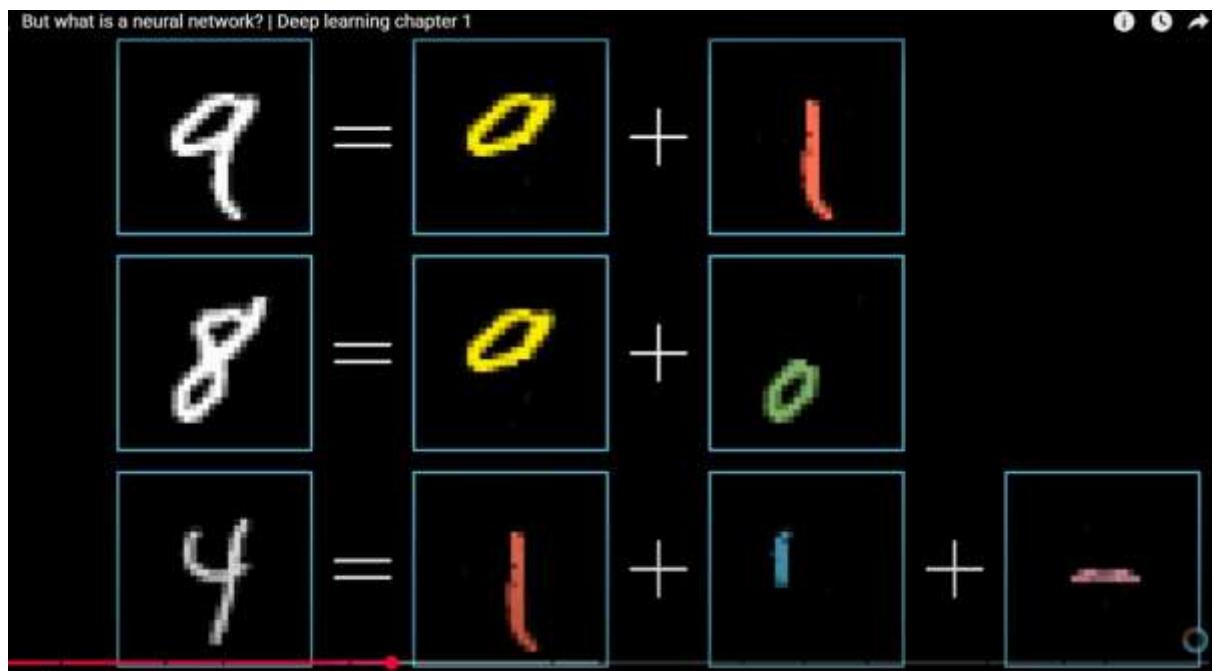
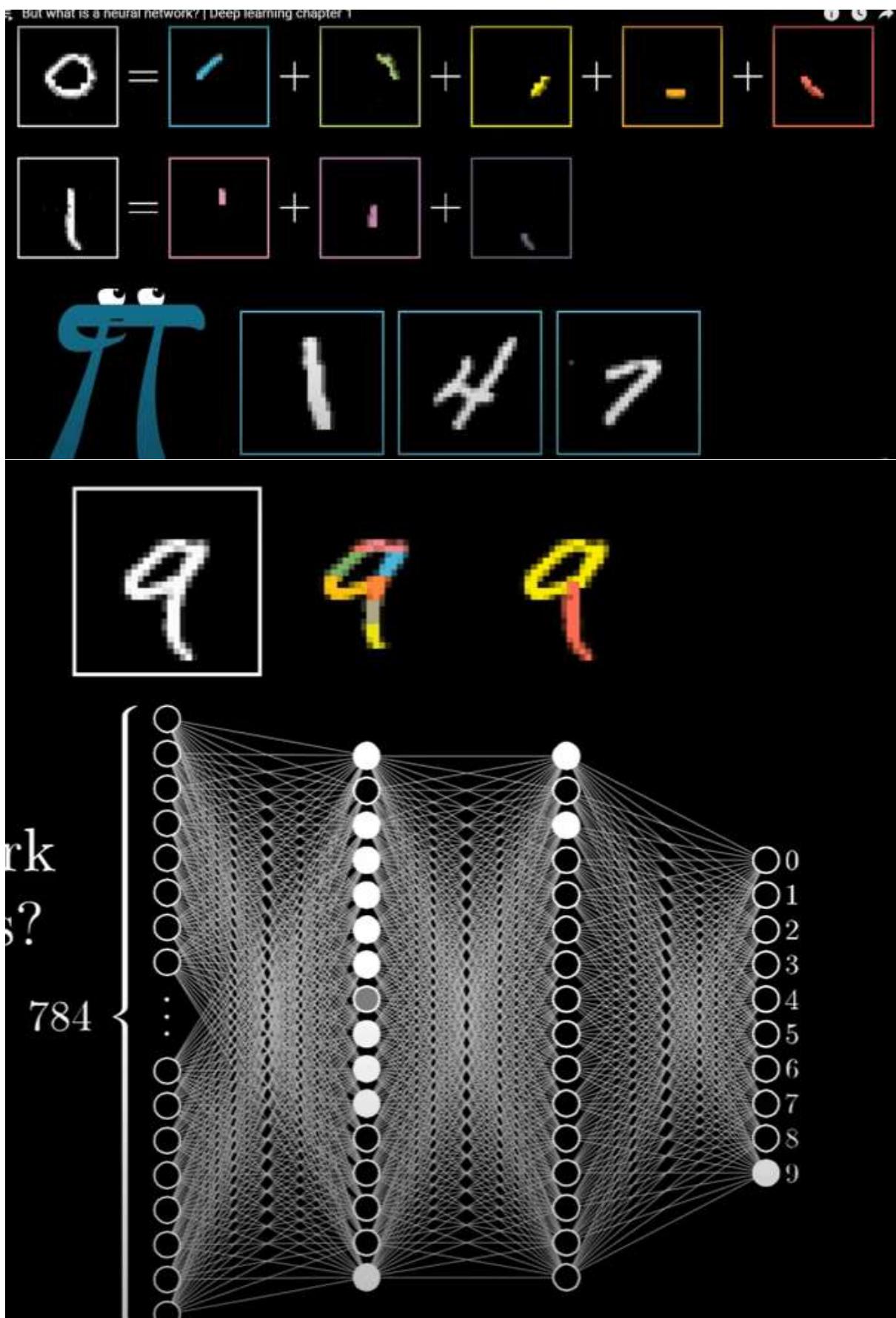
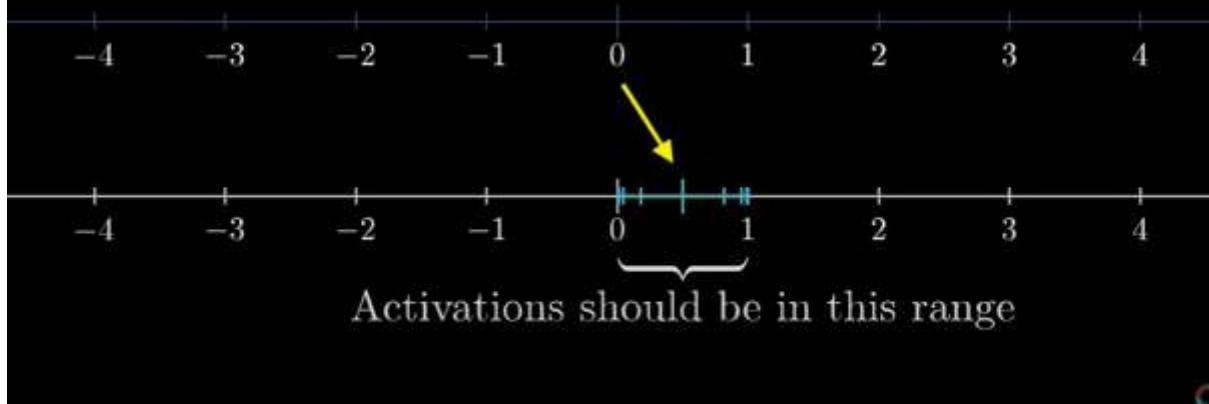


number inside neuron is called activation





$$w_1a_1 + w_2a_2 + w_3a_3 + w_4a_4 + \dots + w_na_n$$



we want to recognize 8 suppose

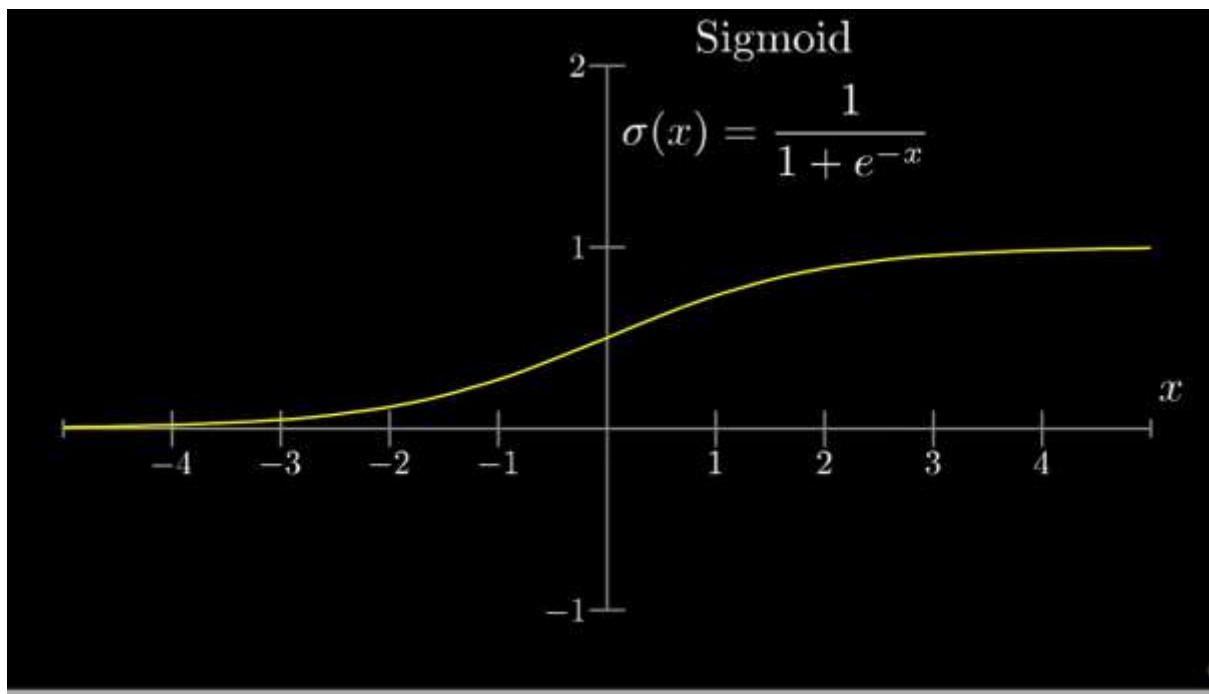
so we multiple the input image pattern with 8's pattern

and it will light up only neurons that match in both grid

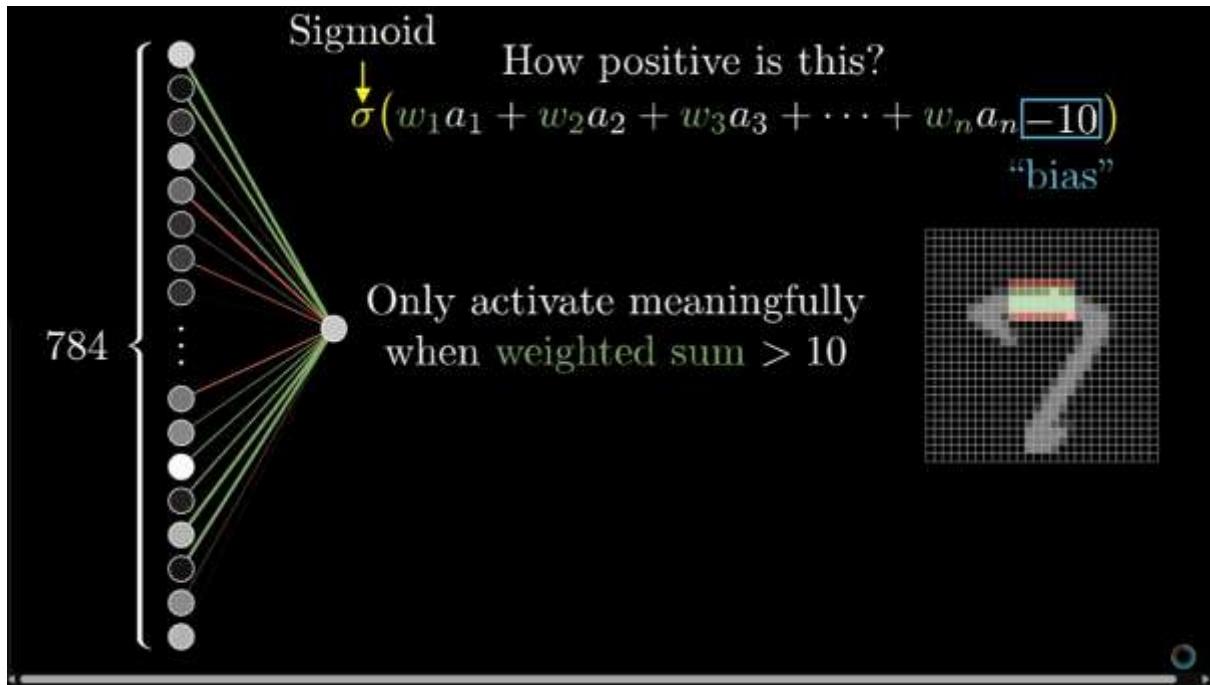
so only those will have higer number

after multiplication value can be anything

but what we want is vlaue between 0 and 1 so for this we use sigmoid

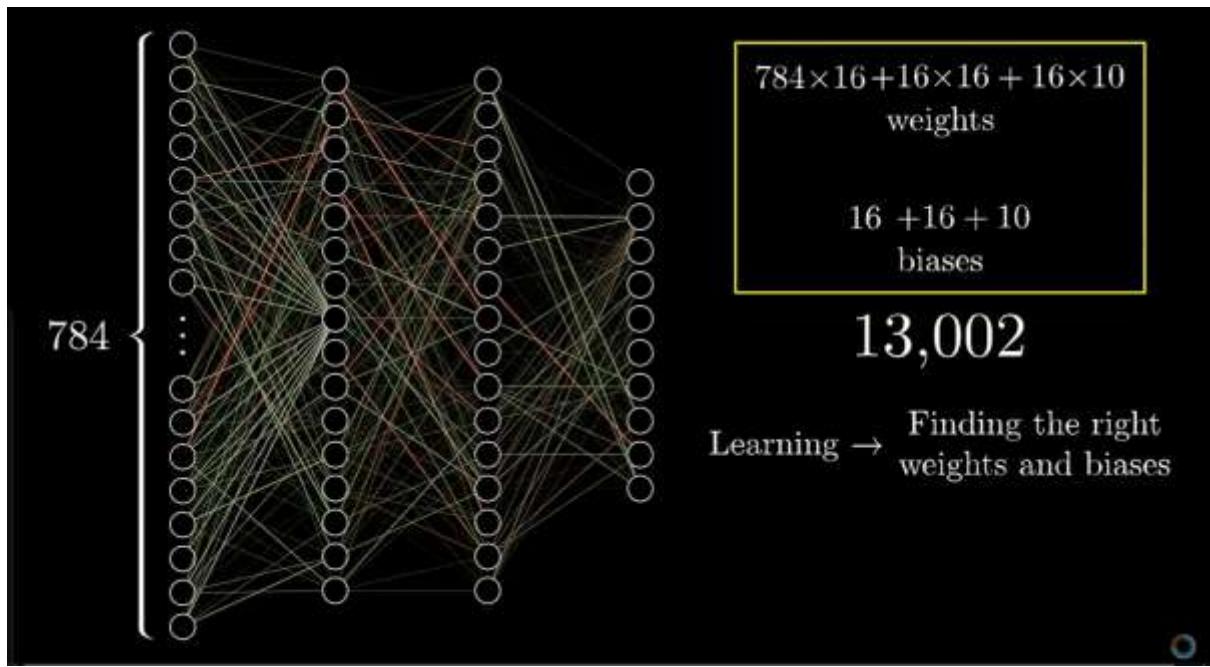


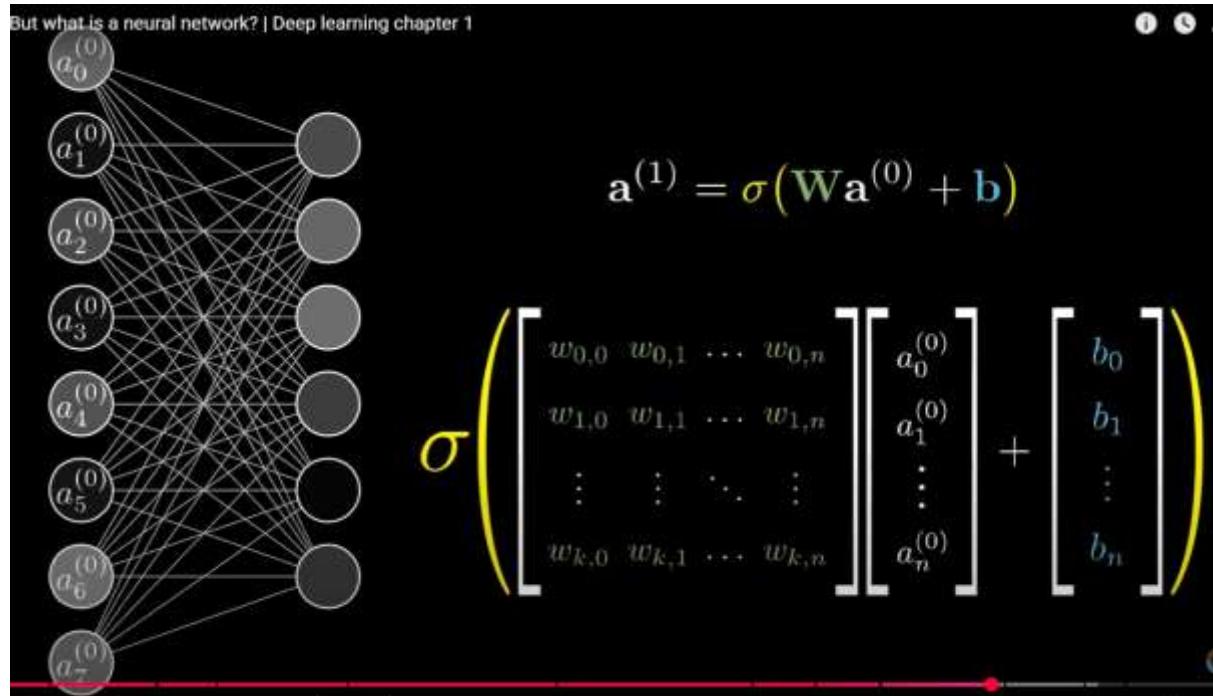
logistic curve



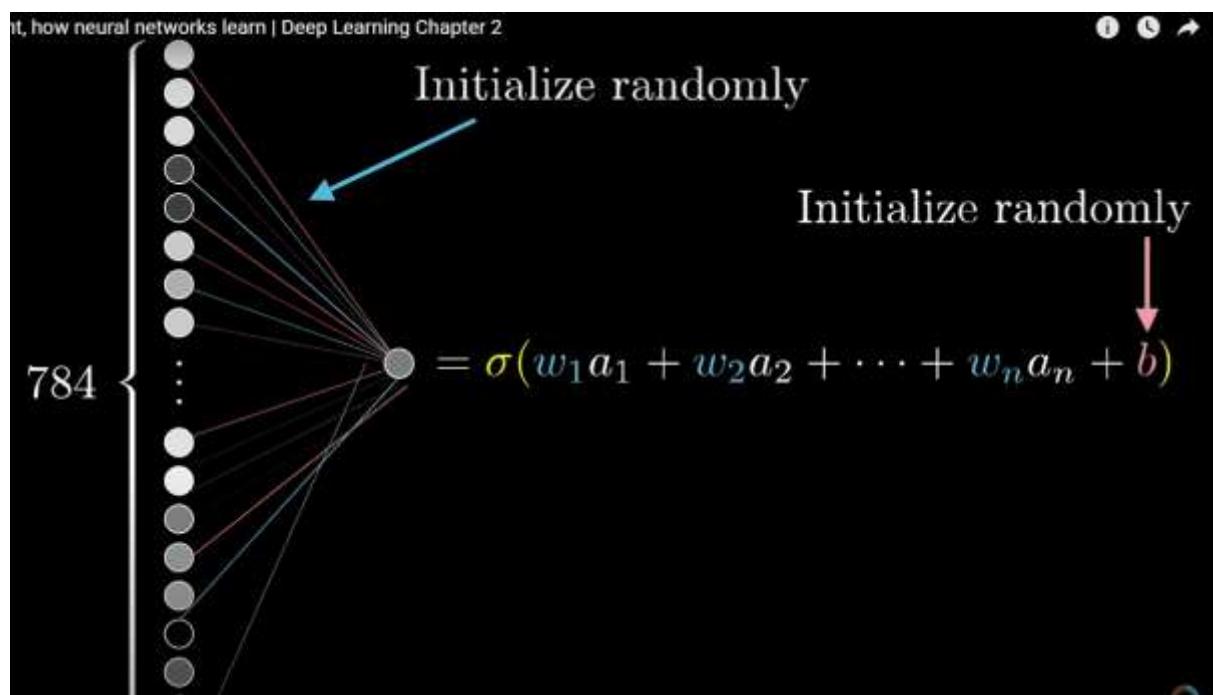
10 mean it is meaningfully active

so we keep that bias

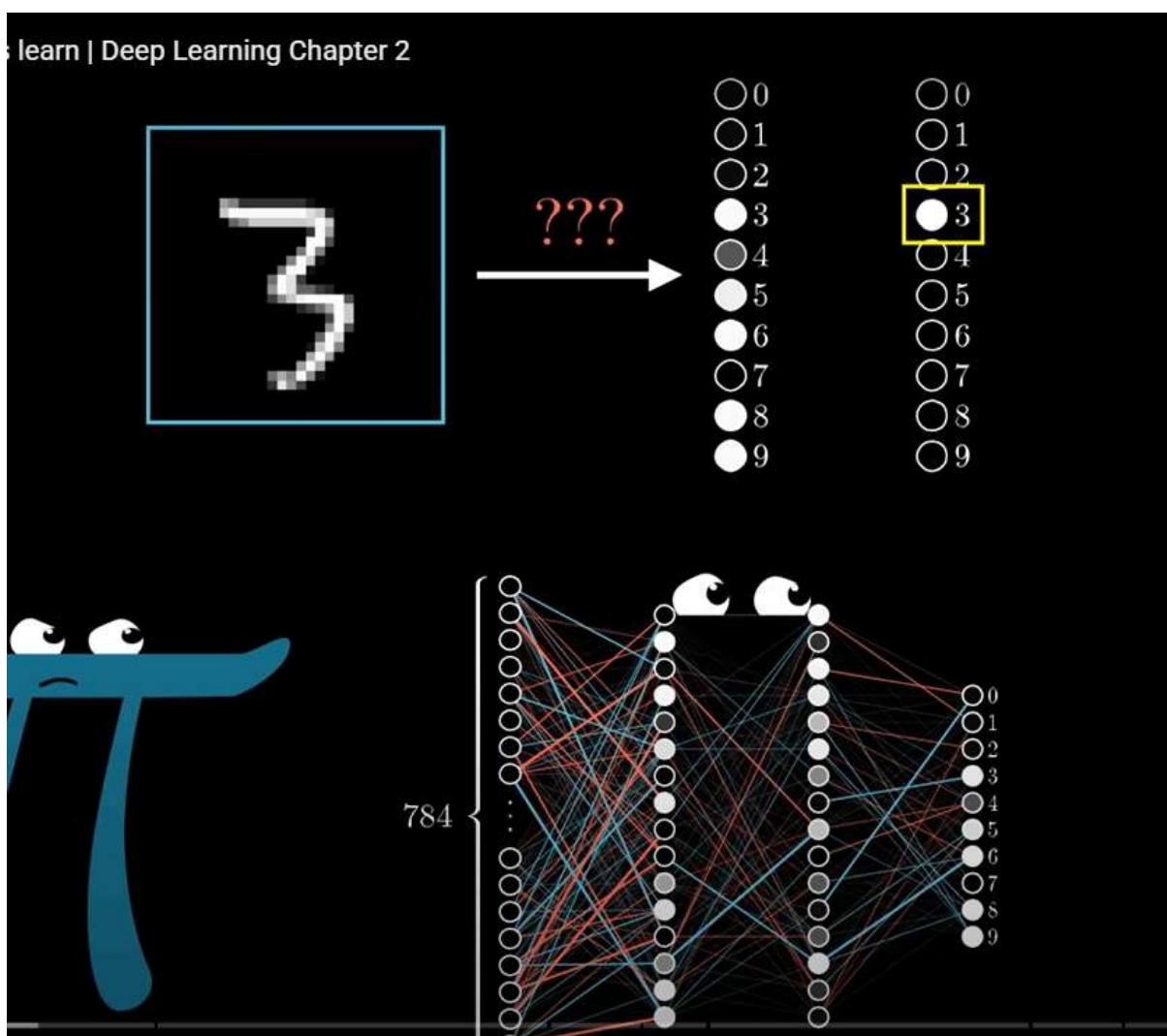




but how does neural network find pattern and learn it ??



learn | Deep Learning Chapter 2



to the left is what network gave with 3,5,6,8,9 ON and what i want is only 3 to be 1 and all other 0

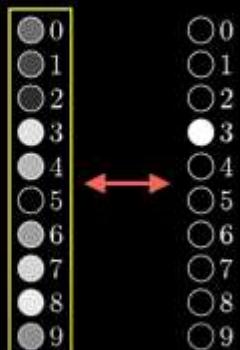
so i say network - bad boy u gave me wring output

so to say this mathematically that what u gave is not what i desired is by MSE

Cost of 3

$$3.37 \left\{ \begin{array}{l} 0.1863 \leftarrow (0.43 - 0.00)^2 + \\ 0.0809 \leftarrow (0.28 - 0.00)^2 + \\ 0.0357 \leftarrow (0.19 - 0.00)^2 + \\ 0.0138 \leftarrow (0.88 - 1.00)^2 + \\ 0.5242 \leftarrow (0.72 - 0.00)^2 + \\ 0.0001 \leftarrow (0.01 - 0.00)^2 + \\ 0.4079 \leftarrow (0.64 - 0.00)^2 + \\ 0.7388 \leftarrow (0.86 - 0.00)^2 + \\ 0.9817 \leftarrow (0.99 - 0.00)^2 + \\ 0.3998 \leftarrow (0.63 - 0.00)^2 \end{array} \right.$$

What's the "cost" of this difference?



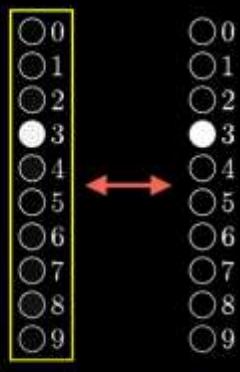
Utter trash

this sum is big so bad boy

Cost of 3

$$0.03 \left\{ \begin{array}{l} 0.0006 \leftarrow (0.02 - 0.00)^2 + \\ 0.0007 \leftarrow (0.03 - 0.00)^2 + \\ 0.0039 \leftarrow (0.06 - 0.00)^2 + \\ 0.0009 \leftarrow (0.97 - 1.00)^2 + \\ 0.0055 \leftarrow (0.07 - 0.00)^2 + \\ 0.0004 \leftarrow (0.02 - 0.00)^2 + \\ 0.0022 \leftarrow (0.05 - 0.00)^2 + \\ 0.0033 \leftarrow (0.06 - 0.00)^2 + \\ 0.0072 \leftarrow (0.08 - 0.00)^2 + \\ 0.0018 \leftarrow (0.04 - 0.00)^2 \end{array} \right.$$

What's the "cost" of this difference?

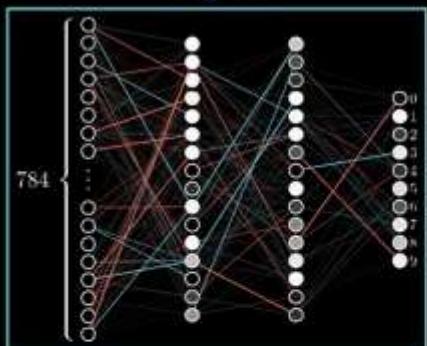


Utter trash

this sum is small so good boy

then we find avg cost of all training data

## Input



Cost: 5.4

## Cost function

Input: 13,002 weights/biases

Output: 1 number (the cost)

Parameters: Many, many, many training examples

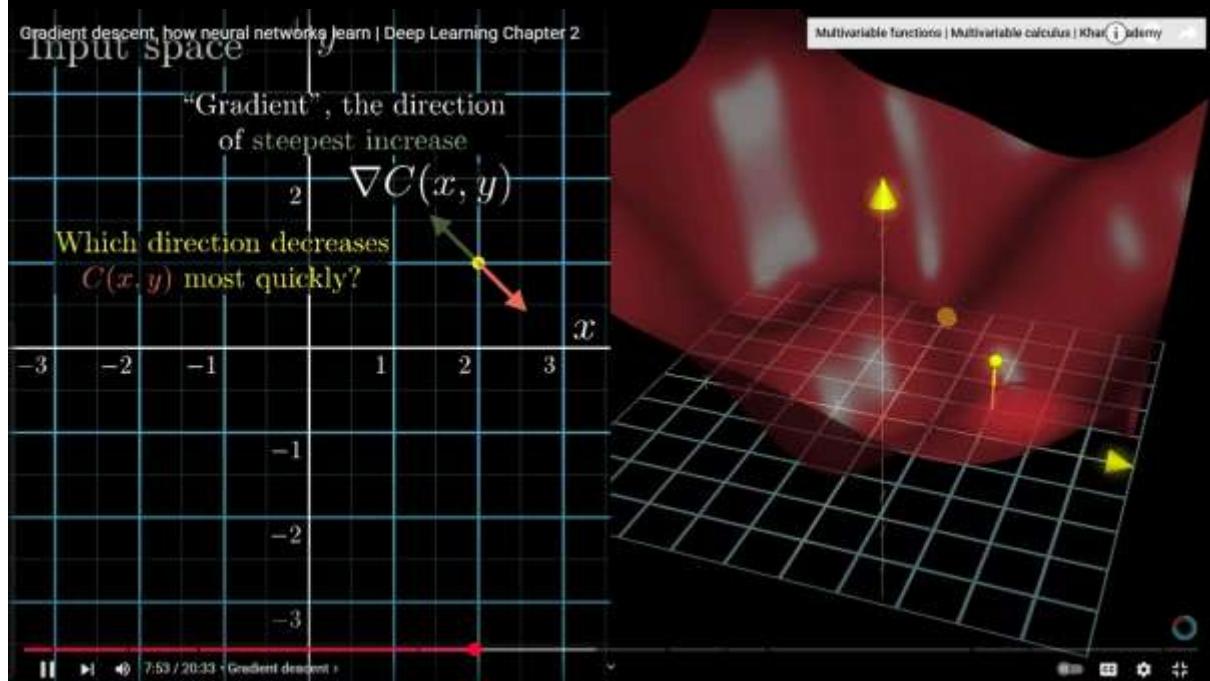
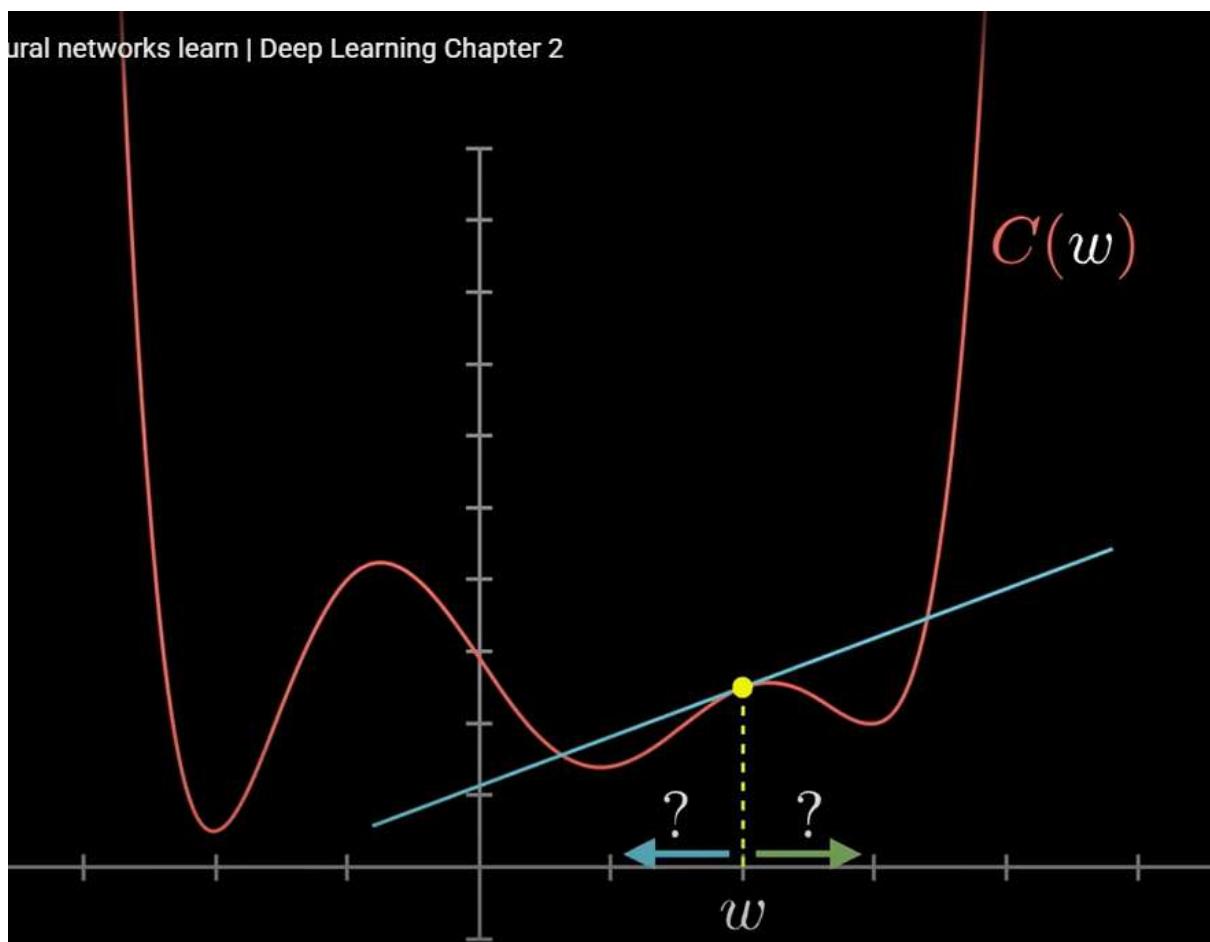
$$([\varphi], 9)$$

$C(w)$   
Single input

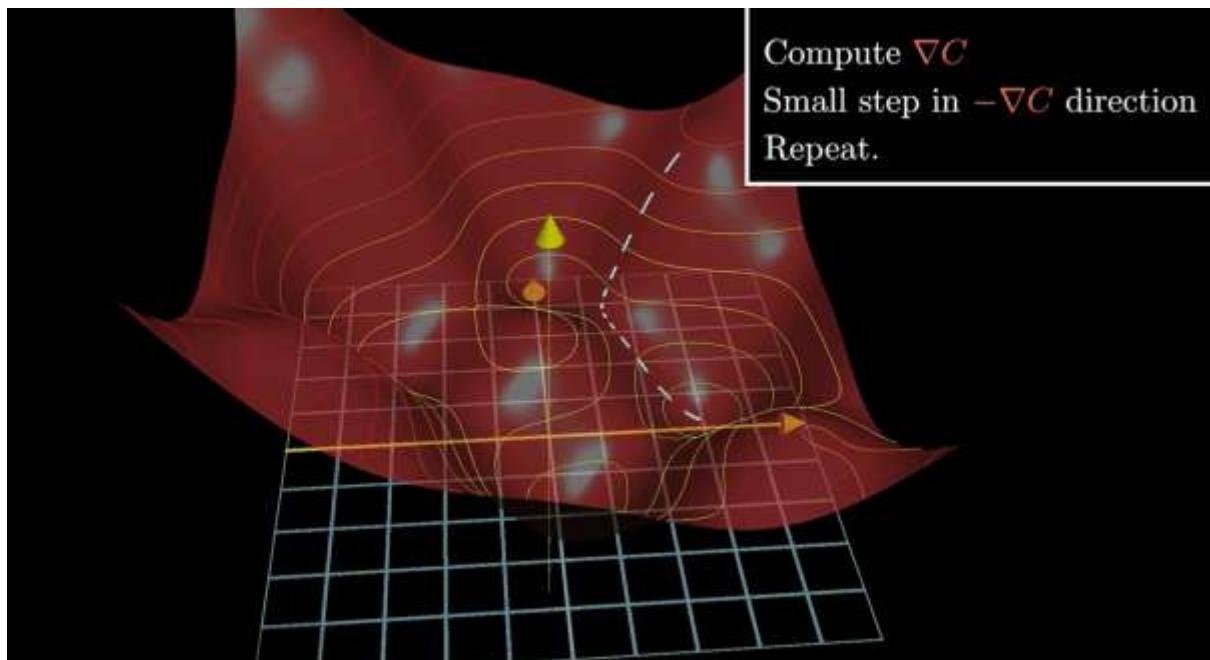
$$\frac{dC}{dw}(w) = 0$$

$w_{\min}$

but this is not feasible for very big values



find vector that tells the downhill direction and how steep that is



It's the same basic idea for a function that has 13,000 inputs instead of 2 inputs.

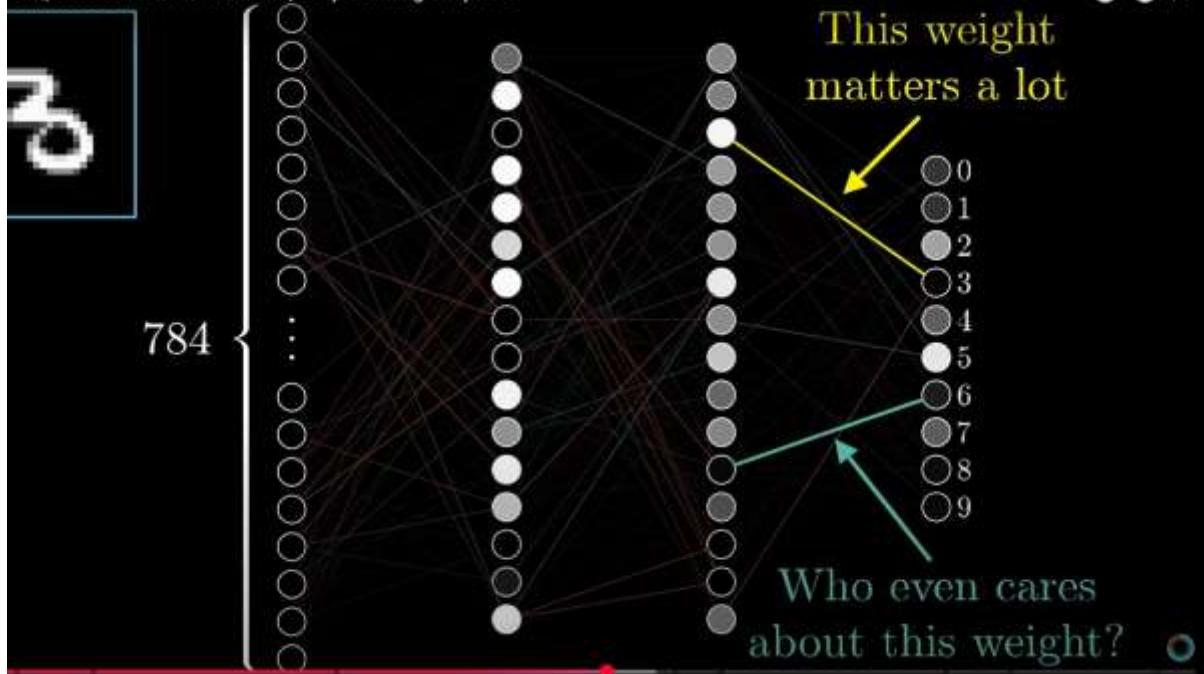
Gradient descent, how neural networks learn | Deep Learning Chapter 2

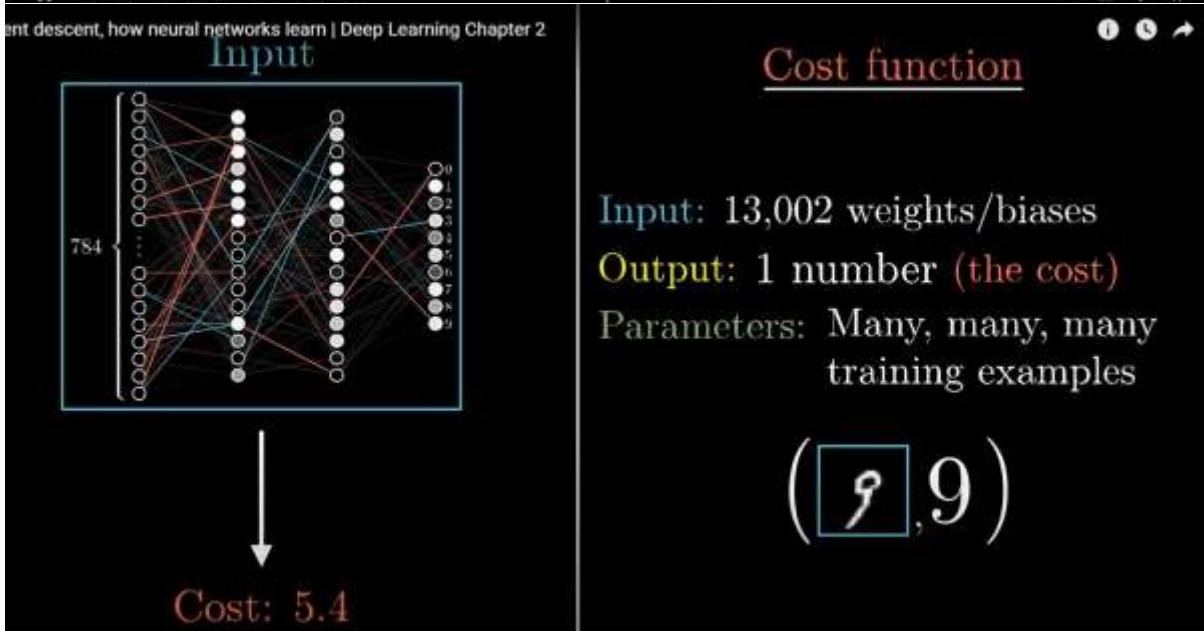
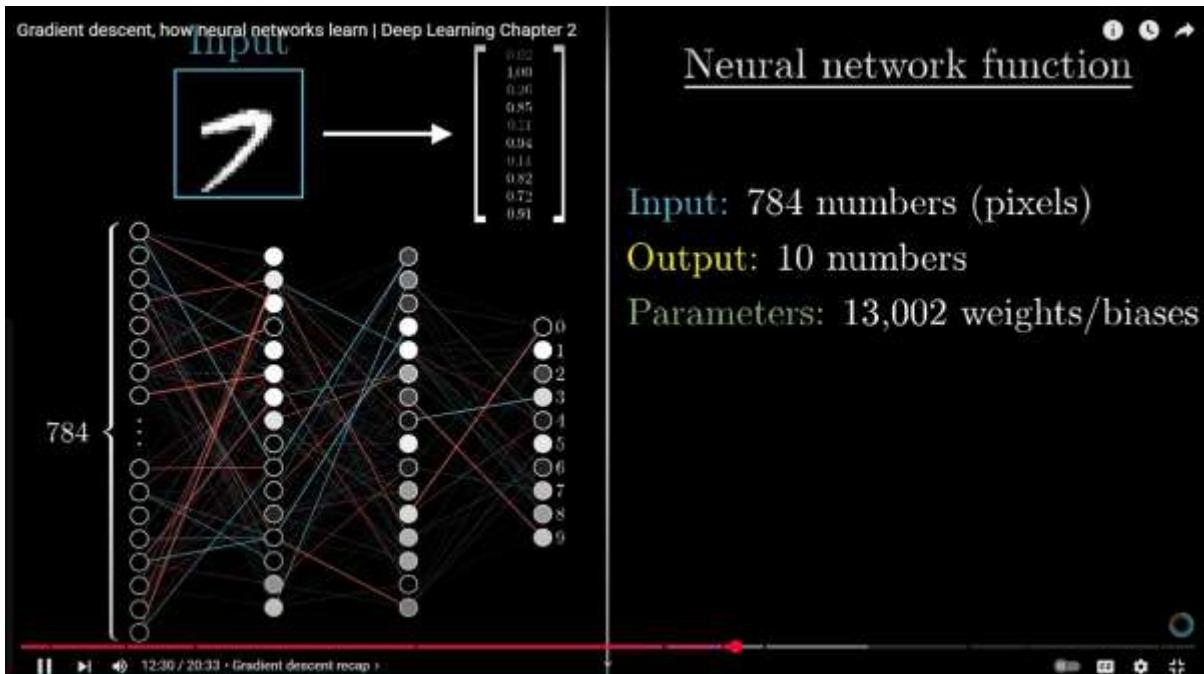


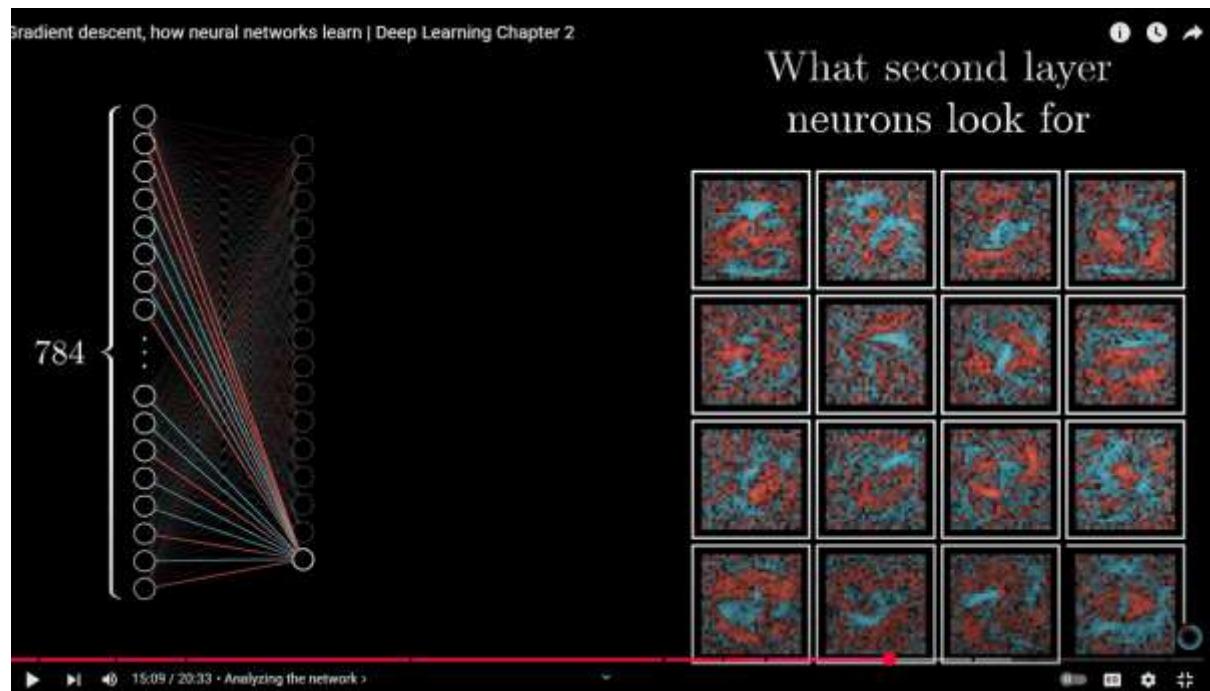
talk about a network learning is that it's just minimizing a cost function.

$$\vec{\mathbf{W}} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_{13,000} \\ w_{13,001} \\ w_{13,002} \end{bmatrix}$$

$$-\nabla C(\vec{\mathbf{W}}) = \begin{bmatrix} 0.31 \\ 0.03 \\ -1.25 \\ \vdots \\ 0.78 \\ -0.37 \\ 0.16 \end{bmatrix} \quad \begin{array}{l} w_0 \text{ should increase somewhat} \\ w_1 \text{ should increase a little} \\ w_2 \text{ should decrease a lot} \\ \vdots \\ w_{13,000} \text{ should increase a lot} \\ w_{13,001} \text{ should decrease somewhat} \\ w_{13,002} \text{ should increase a little} \end{array}$$

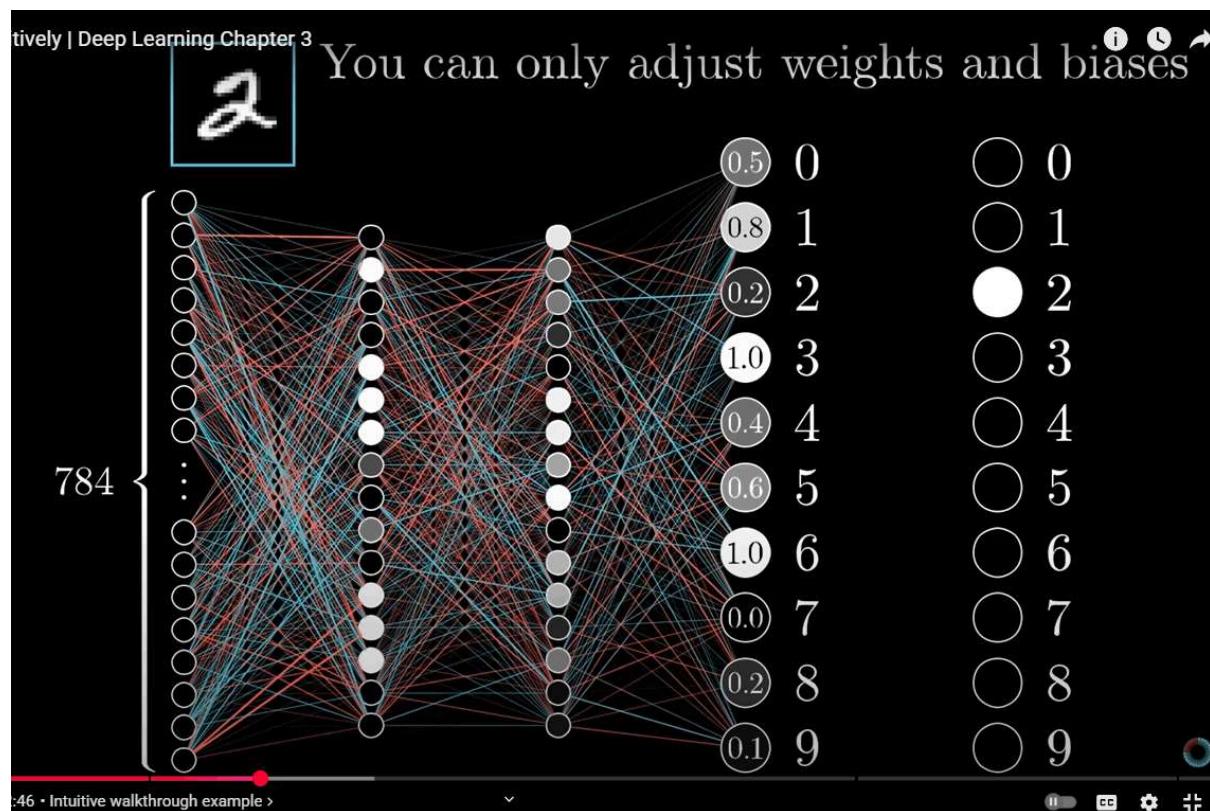


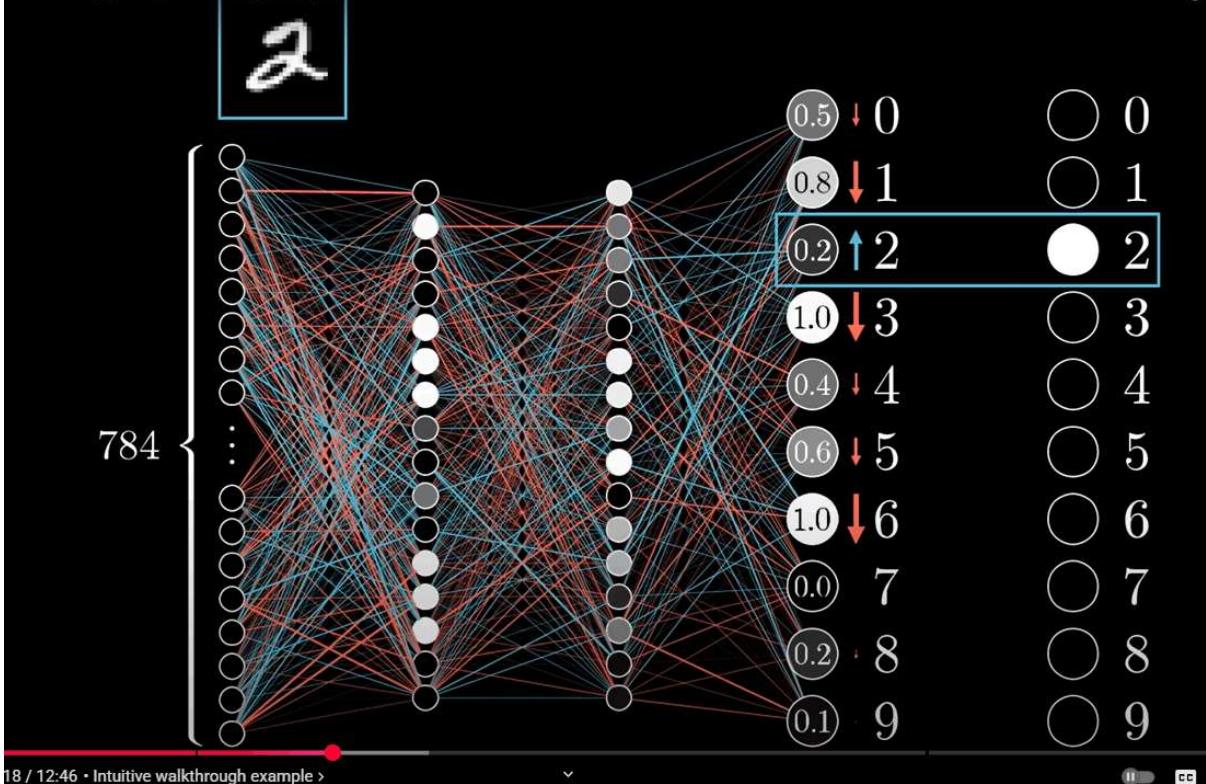




initially

bad network

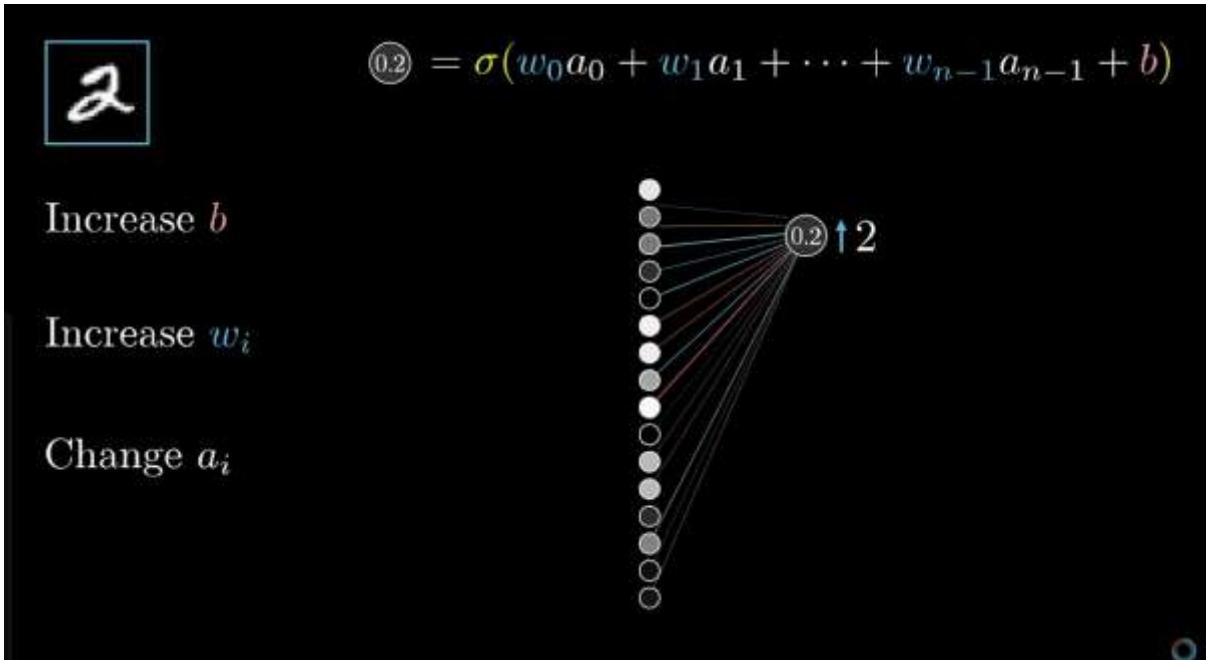


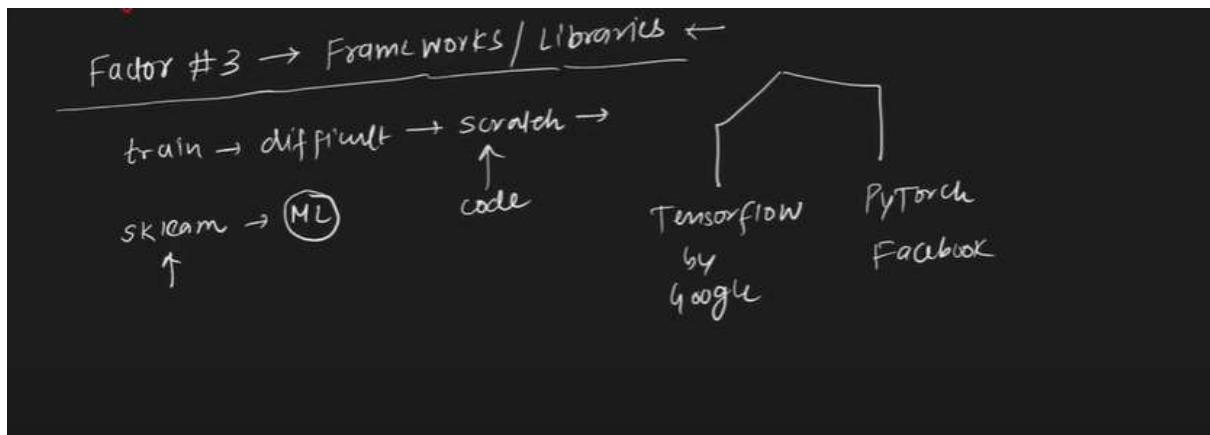


18 / 12:46 · Intuitive walkthrough example &gt;

CC

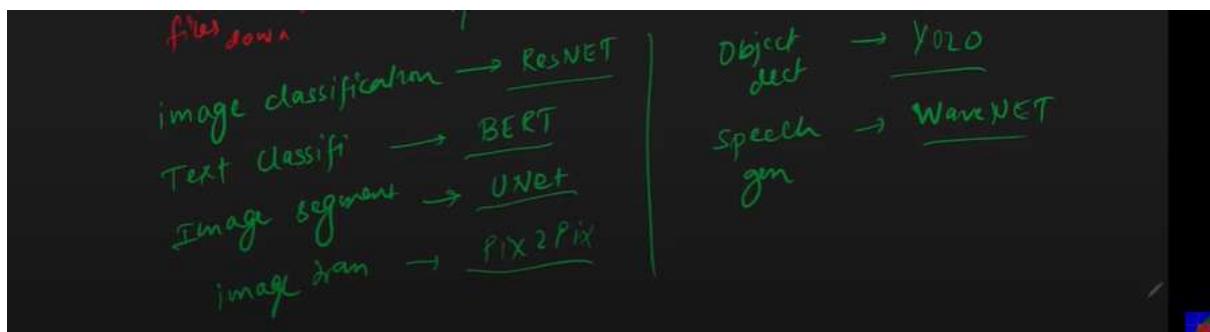
the sizes of UP and DOWN should be proportionate





tensorflow was very tough so they launches tensorflow + keras

all these are famous neural networks pretrained by researches that u can download and use



- type 1 - MLP - regression / classification problems
- type 2 - CNN - image processing
- type 3 - RNN / LSTM - NLP application
- type 4 - auto encoders - compression without losing quality
- type 5 - GAN - image / music generation

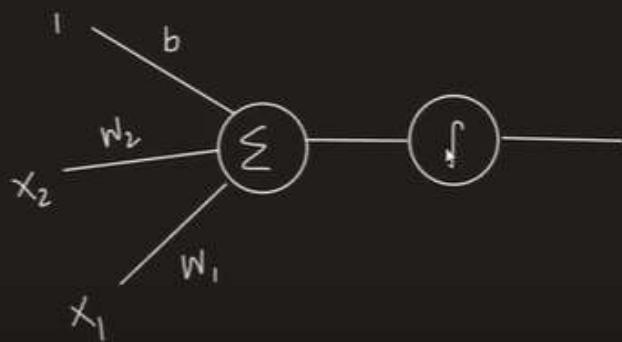
Type of Neural Network	Key Characteristics	Applications
Convolutional Neural Network (CNN)	Designed for grid-like data like images; uses convolutional layers.	Used in Image Recognition and medical image analysis.
Recurrent Neural Network (RNN)	Suitable for sequential data such as text and speech; has internal memory.	Applied in Speech recognition and language translation.
Long Short-Term Memory (LSTM)	A type of RNN for long-term dependencies; uses memory cells.	Utilized for language translation and sentiment analysis.
Transformer Networks	Processes sequences using a self-attention mechanism.	Commonly used for Natural Language Processing (NLP) and machine translation.
Generative Adversarial Networks (GANs)	Involves competing generator and discriminator networks.	Used for image generation and data augmentation.
Artificial Neural Networks (ANNs)	A foundational type with input, hidden, and output layers.	Applied in Image classification and predictive analytics.
Multilayer Perceptrons (MLPs)	Feedforward ANNs with multiple layers.	Used in Facial recognition and classification analysis.

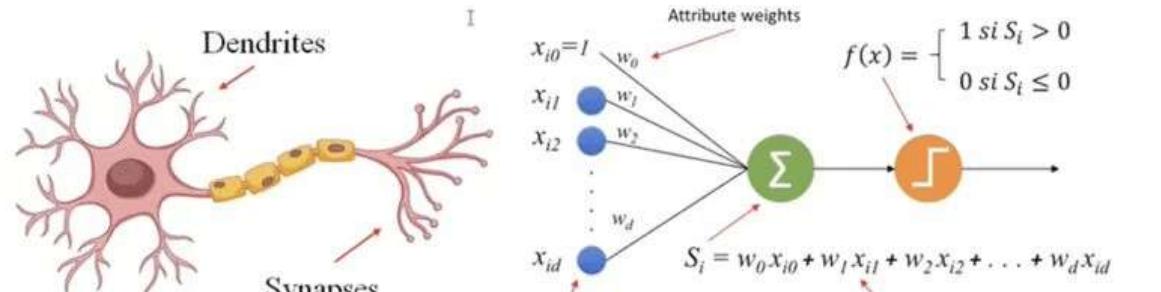
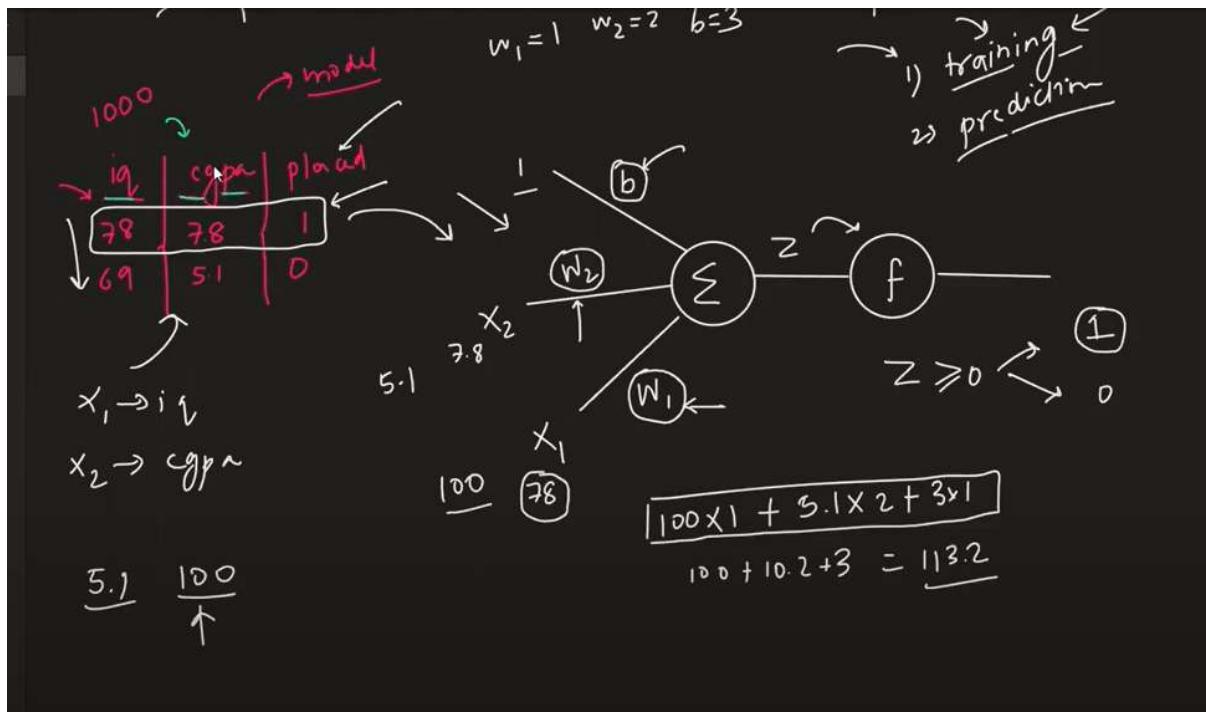
### What is a Perceptron?

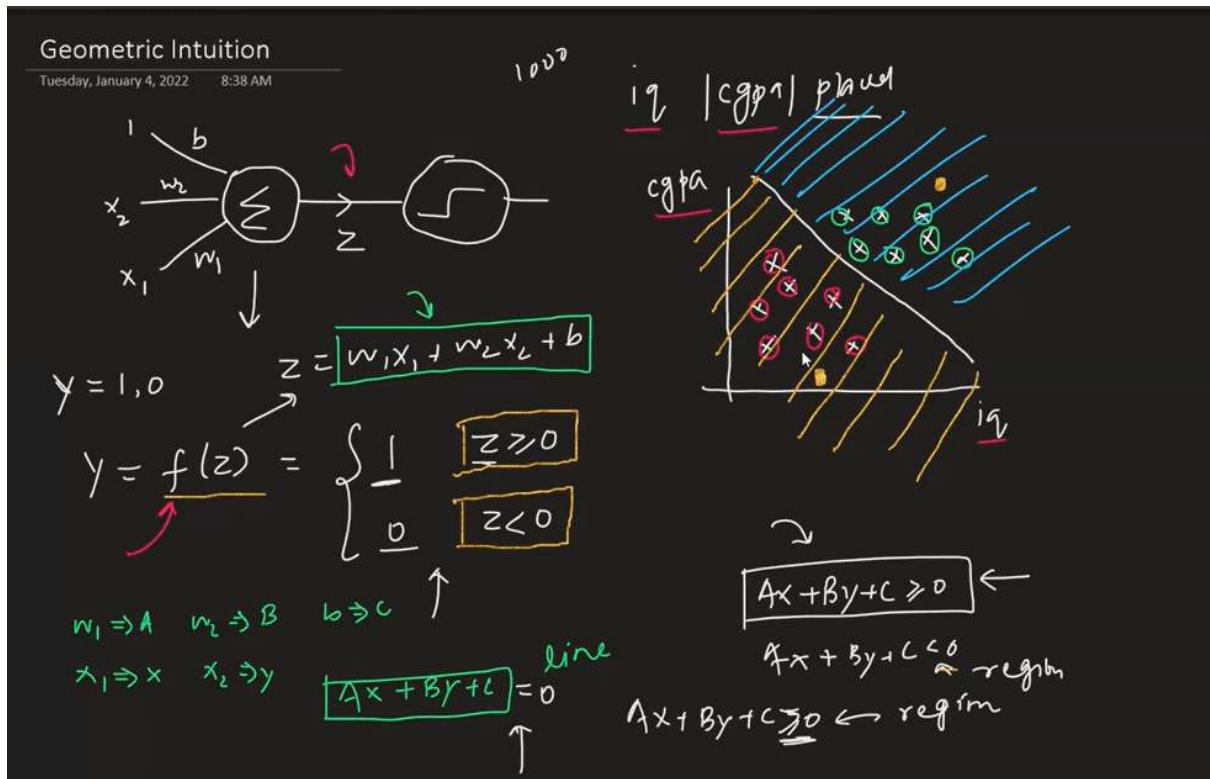
Tuesday, January 4, 2022 6:45 AM

Algorithm → supervised ml  
 ↓ design

- mathematical model
- function

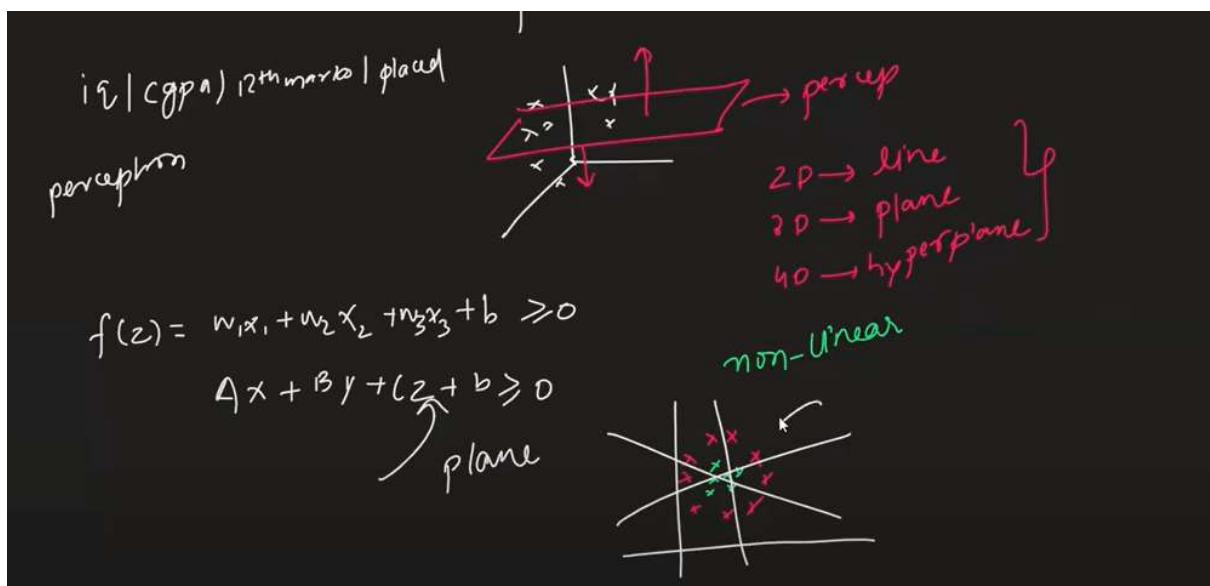






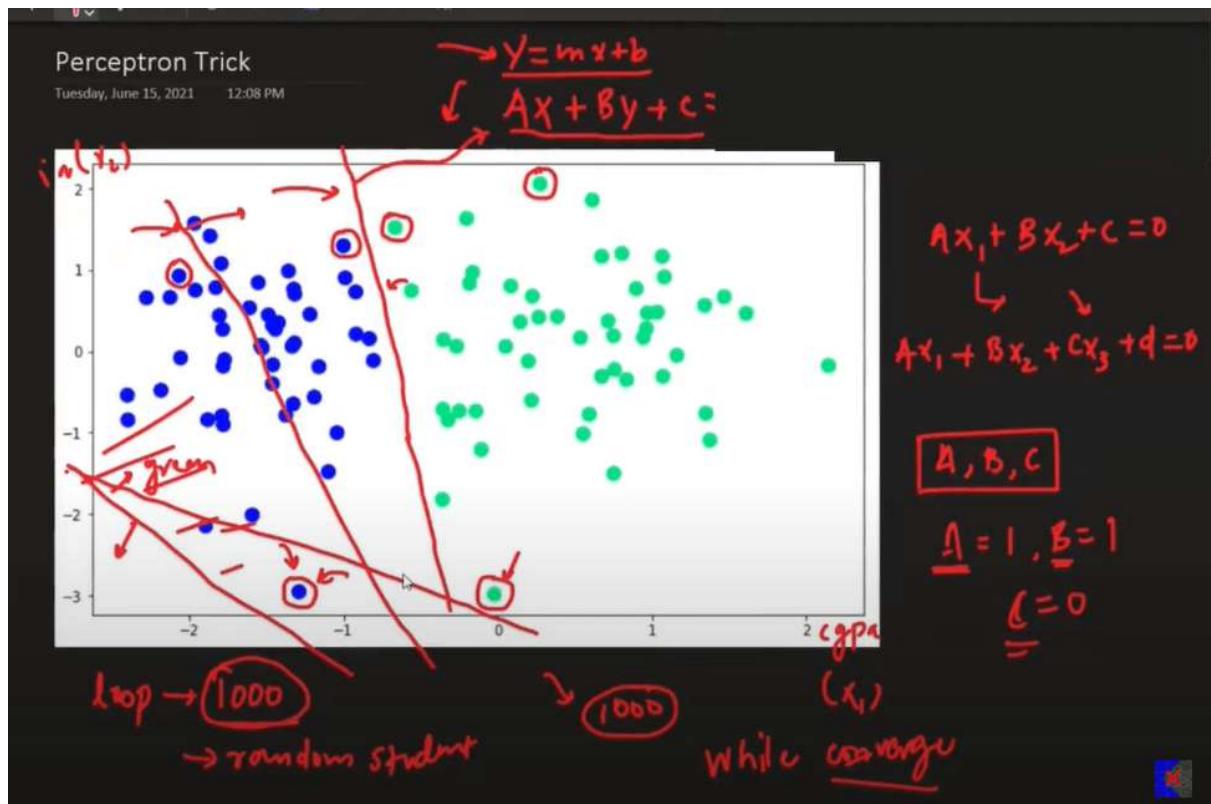
perceptron is a binary classifier

perceptron fails with non linear dataset



[scikit-learn.org/stable/modules/generated/sklearn.linear\\_model.Perceptron.html](https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.Perceptron.html)

[rasbt.github.io/mlxtend/user\\_guide/plotting/plot\\_decision\\_regions/](https://rasbt.github.io/mlxtend/user_guide/plotting/plot_decision_regions/)

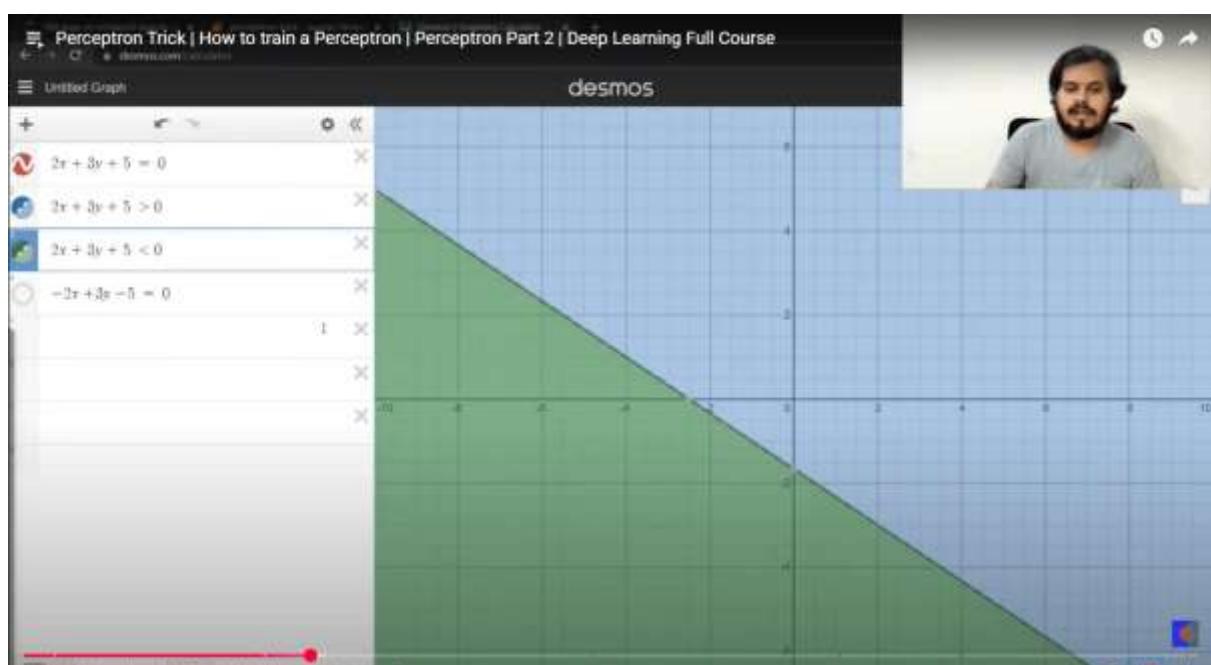


we start with random ABC values

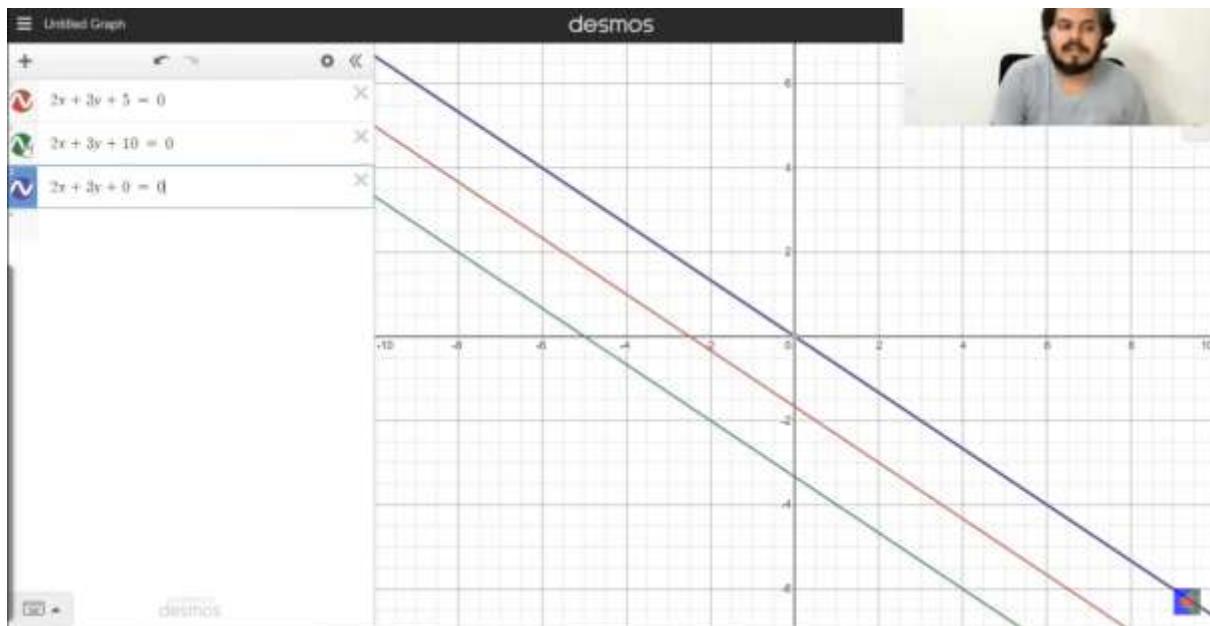
and then based on the line we ask random points whether they are guessed right or not

if they are guessed right no changes but if wrong then we shift line such that that selected point falls in correct region

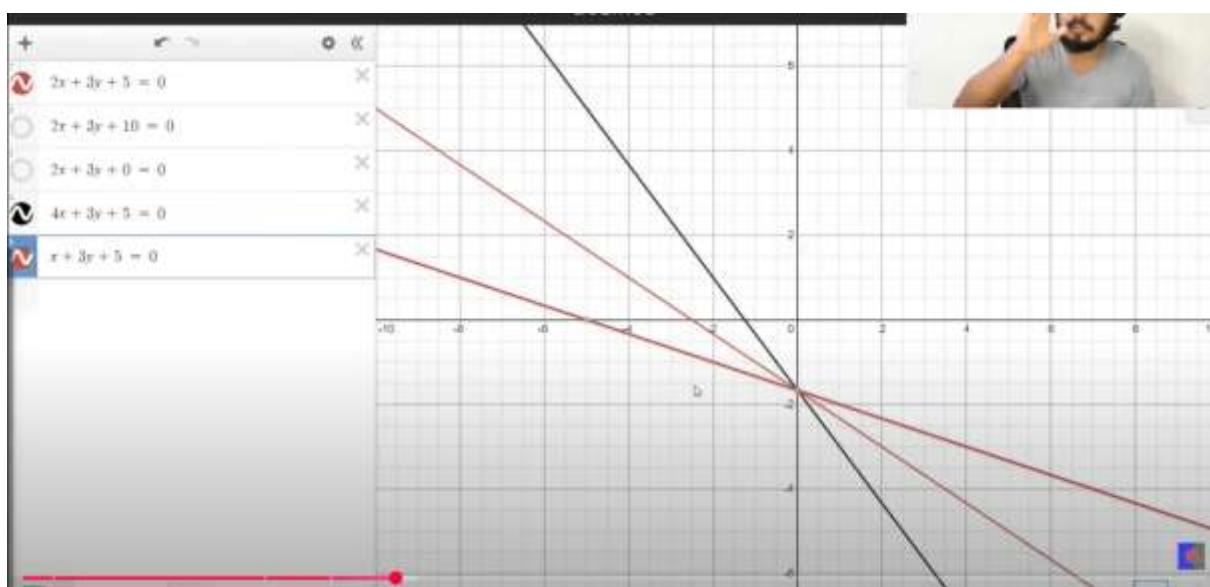
we can do this till convergence

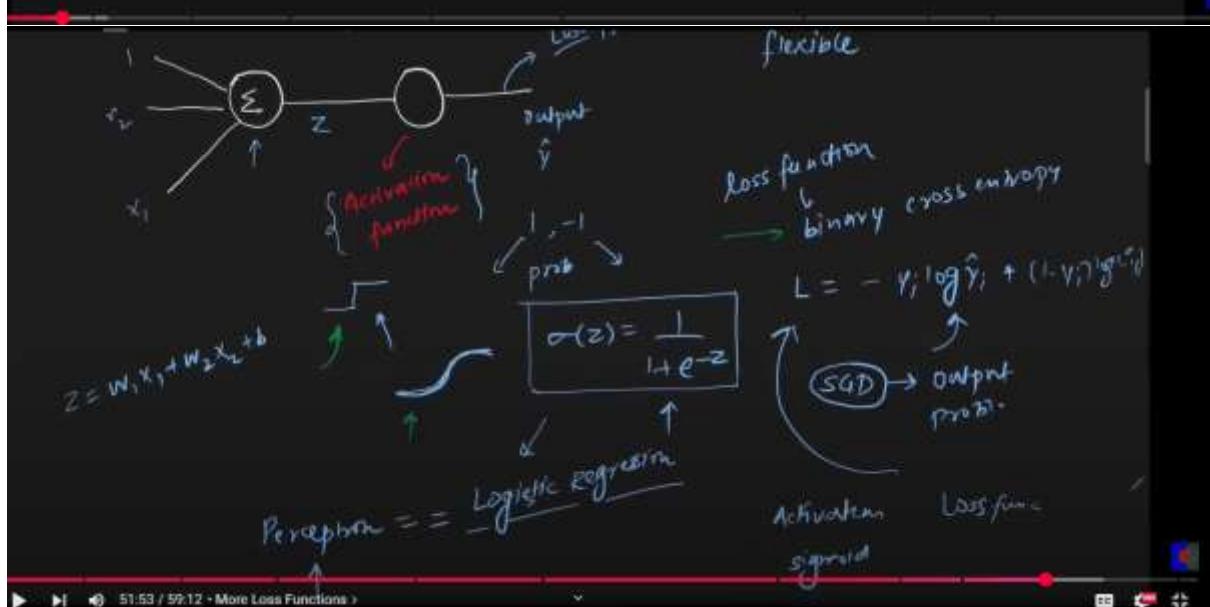
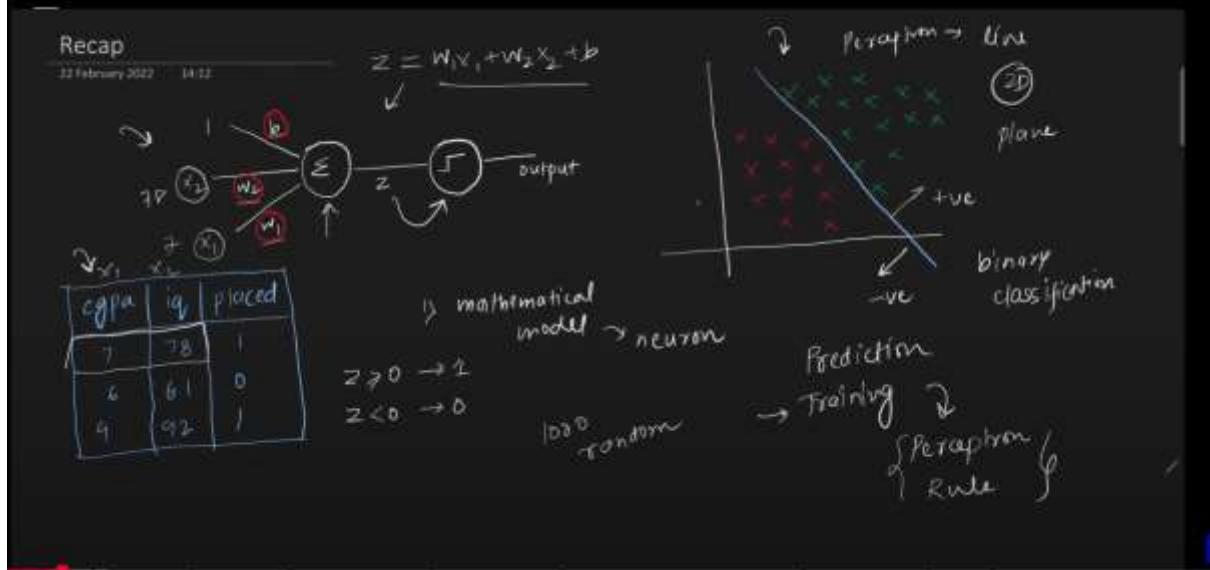
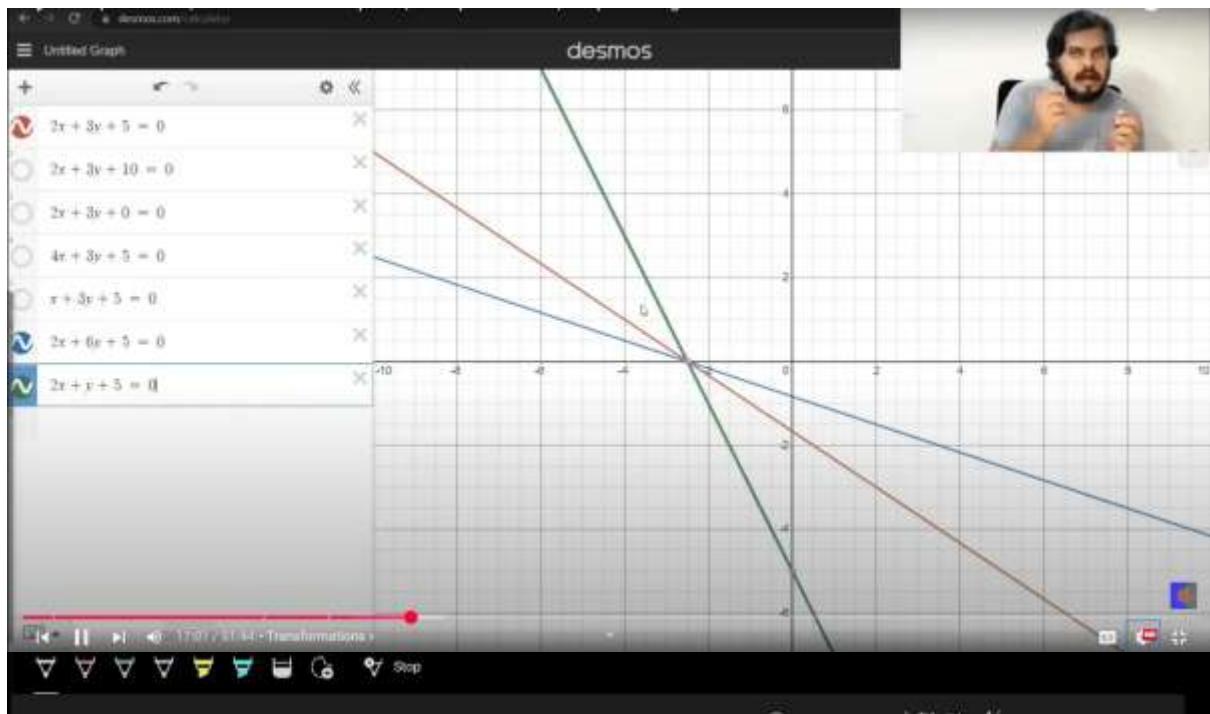


C is moving line up or down



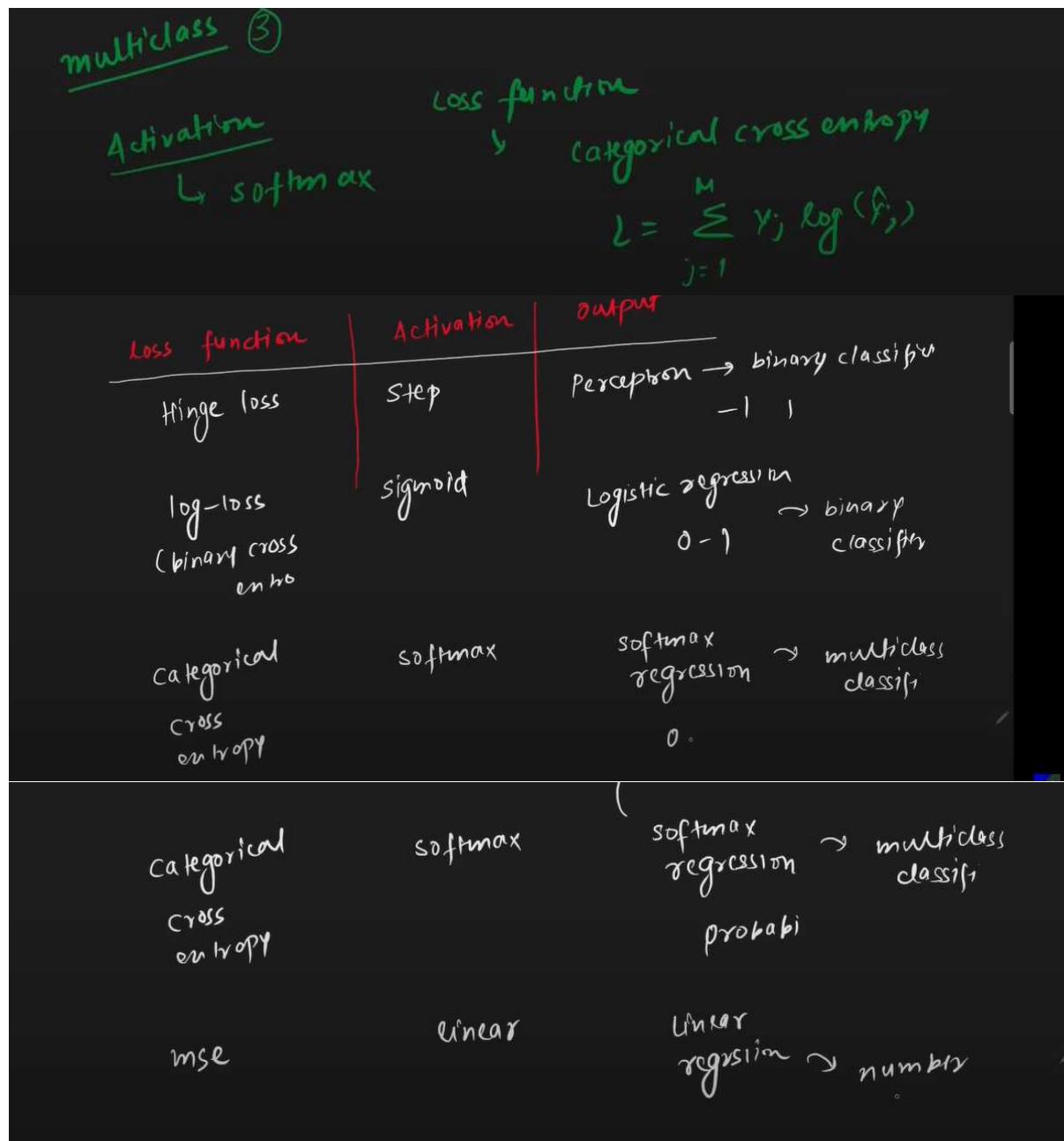
A , B is rotating line

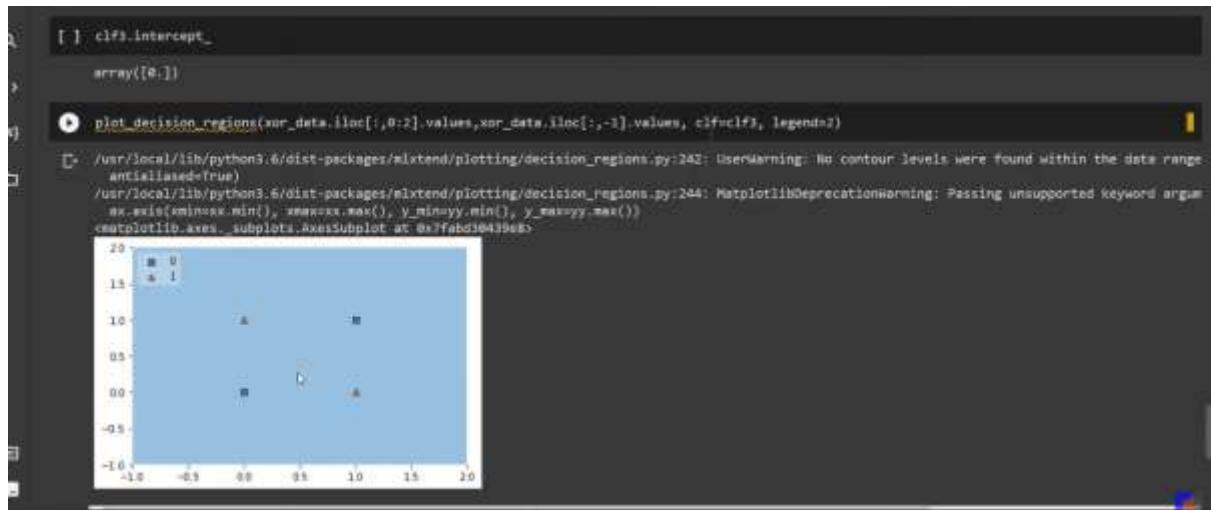




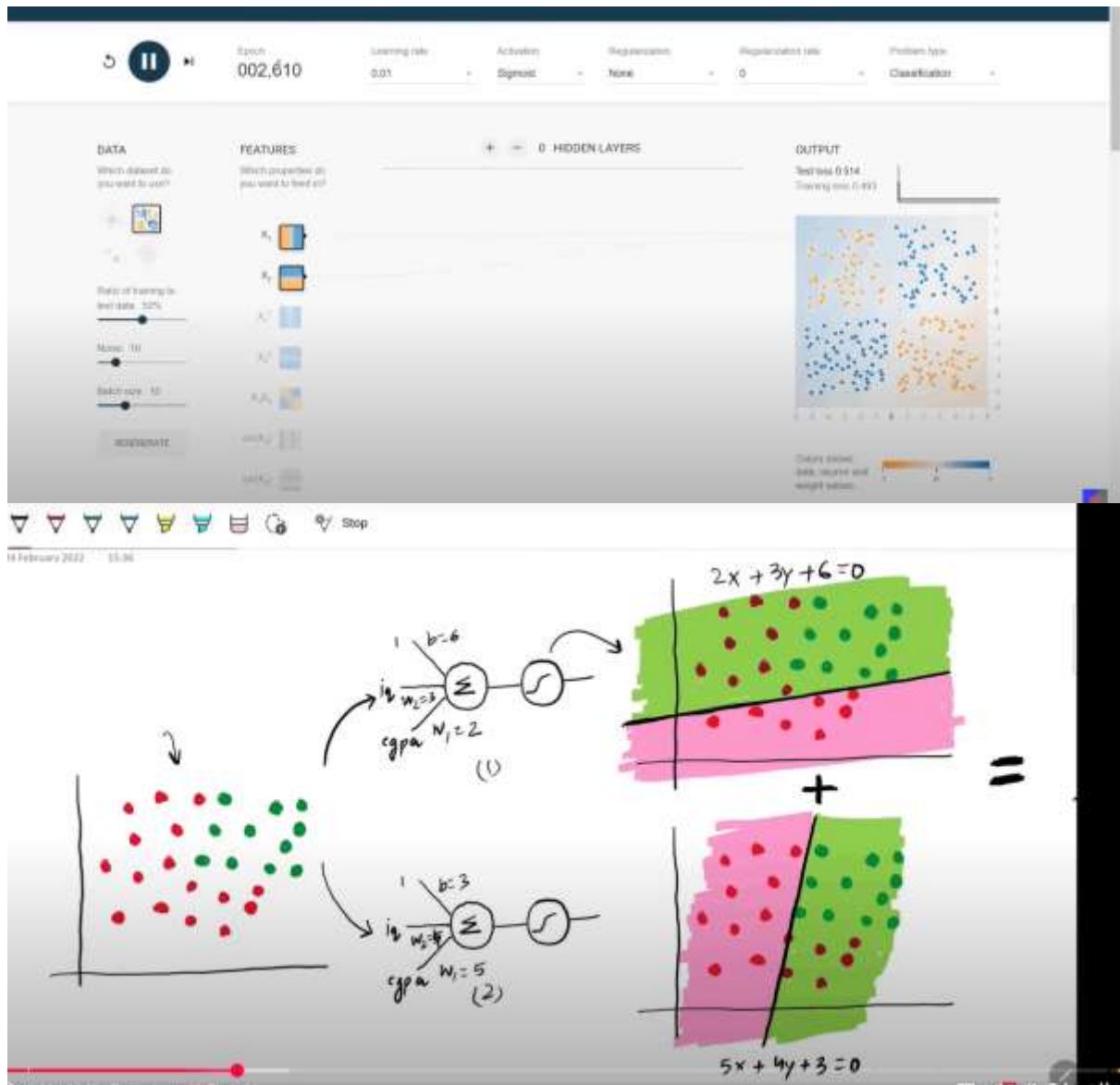
logistic regression has

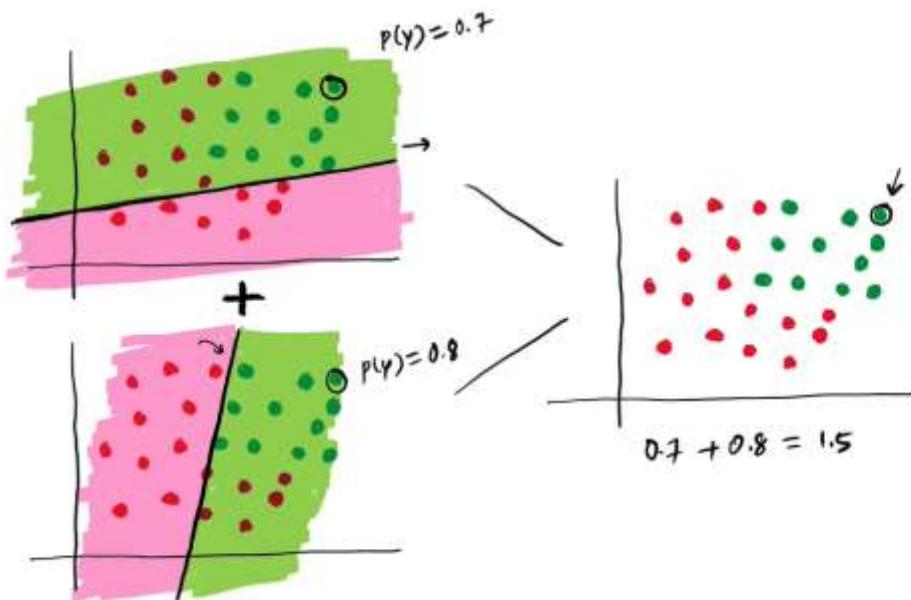
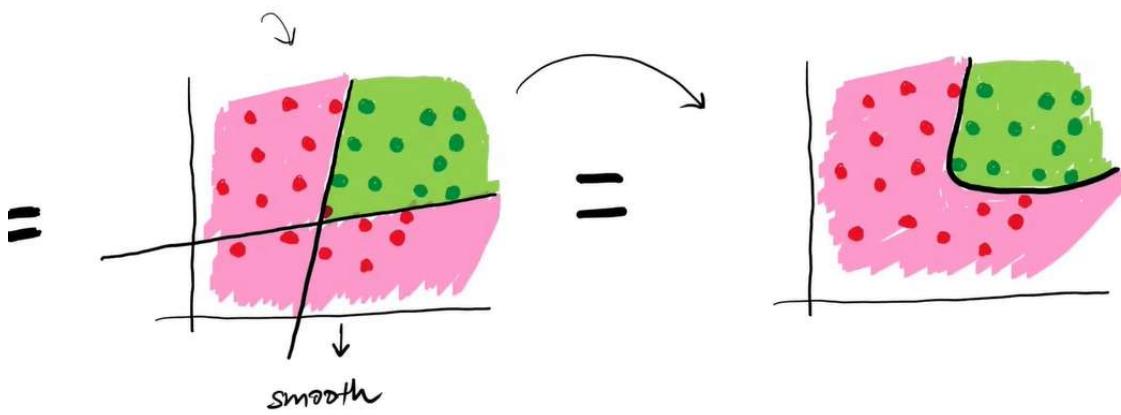
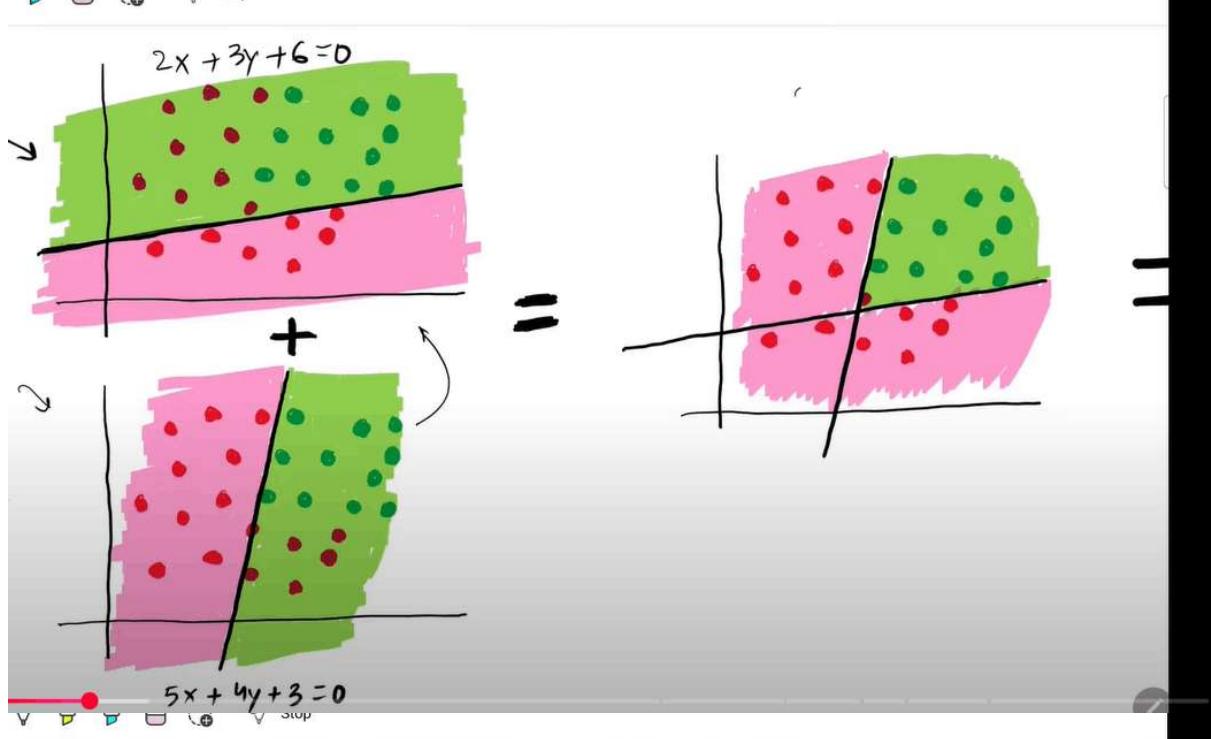
- activation func - sigmoid
- loss function - binary cross entropy

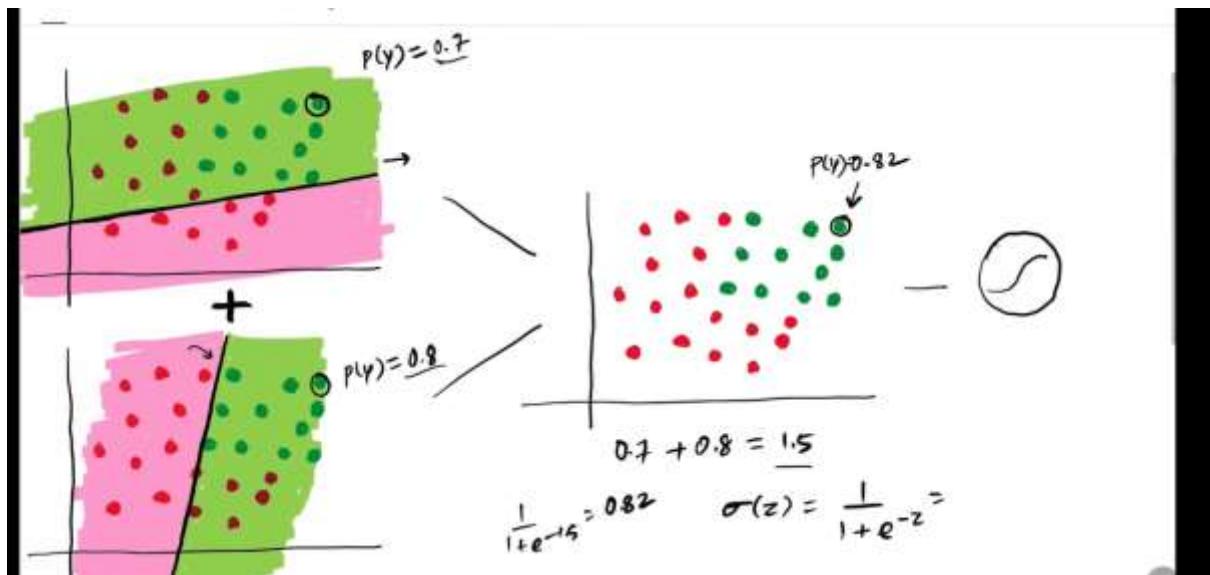




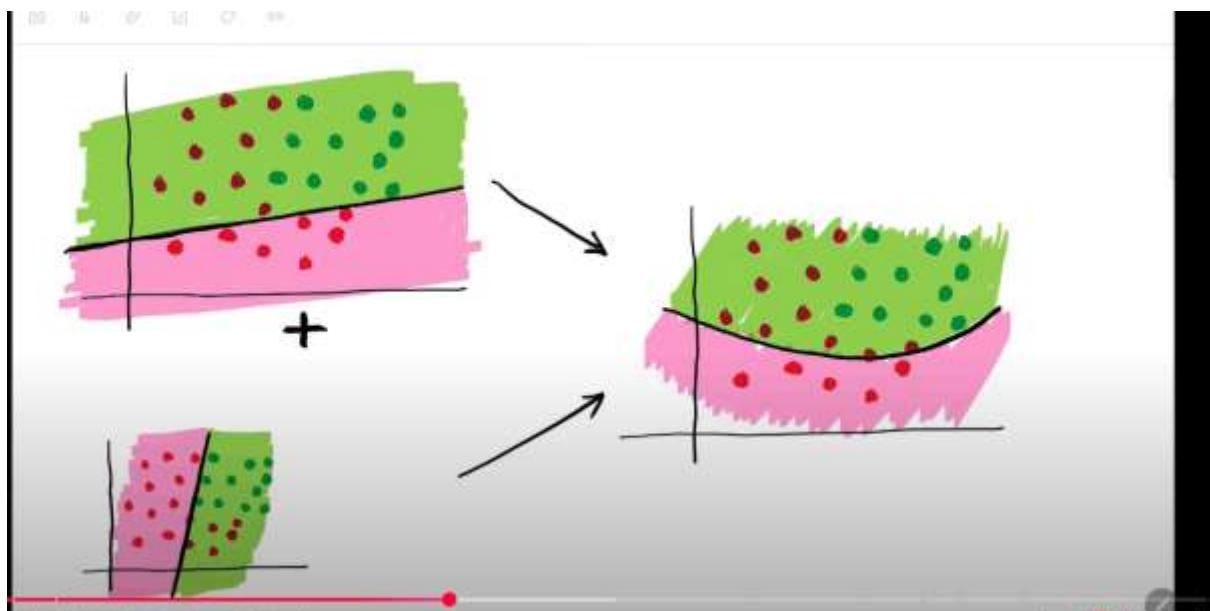
no line possible this is XOR data







now suppose i want below type of curve

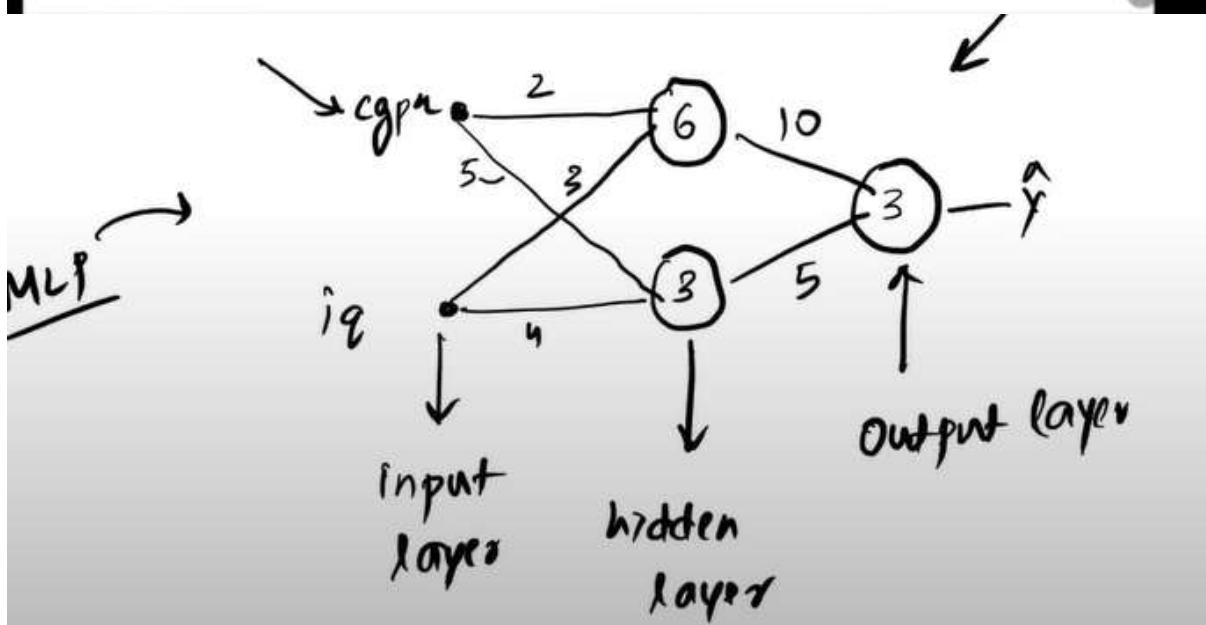
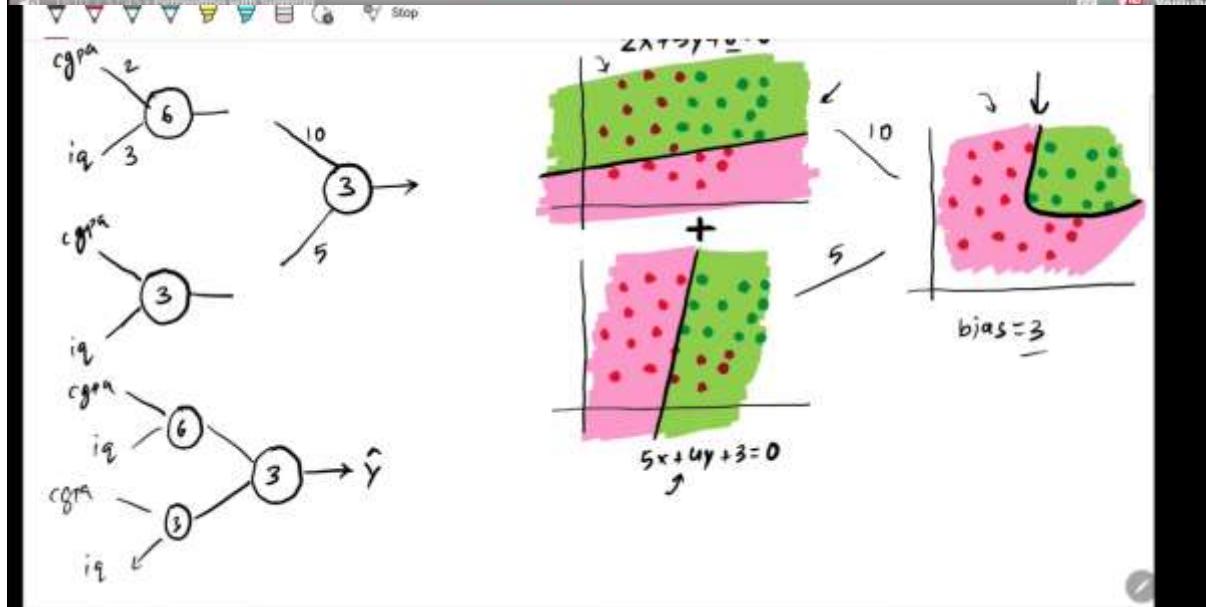
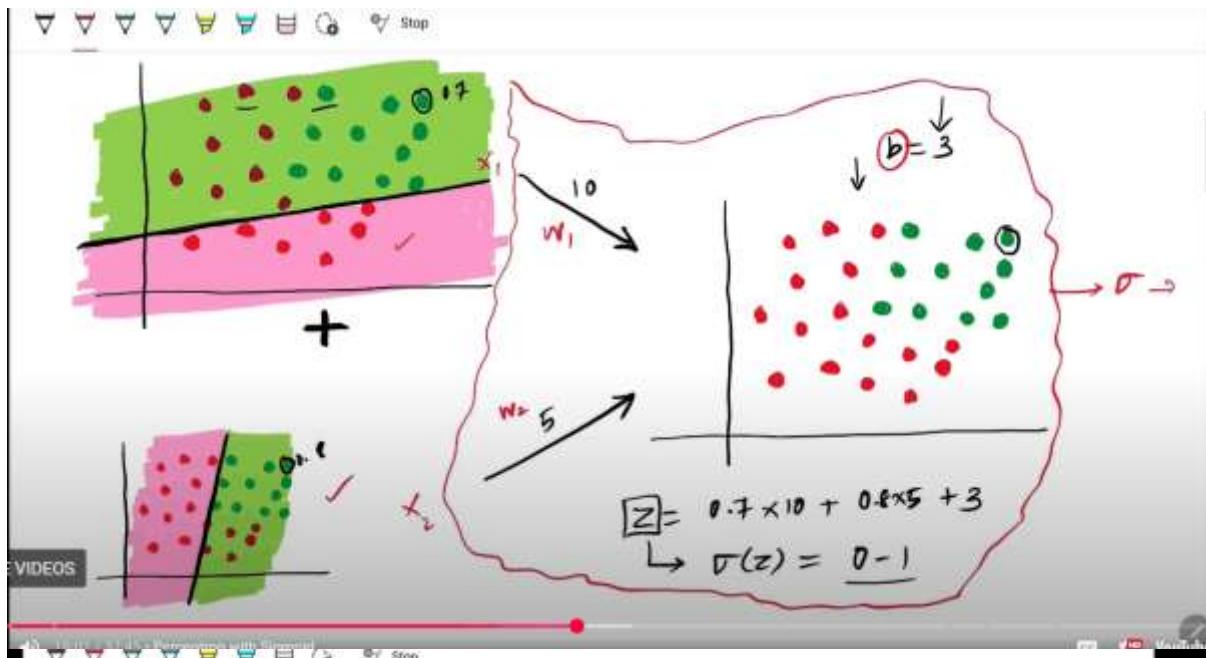


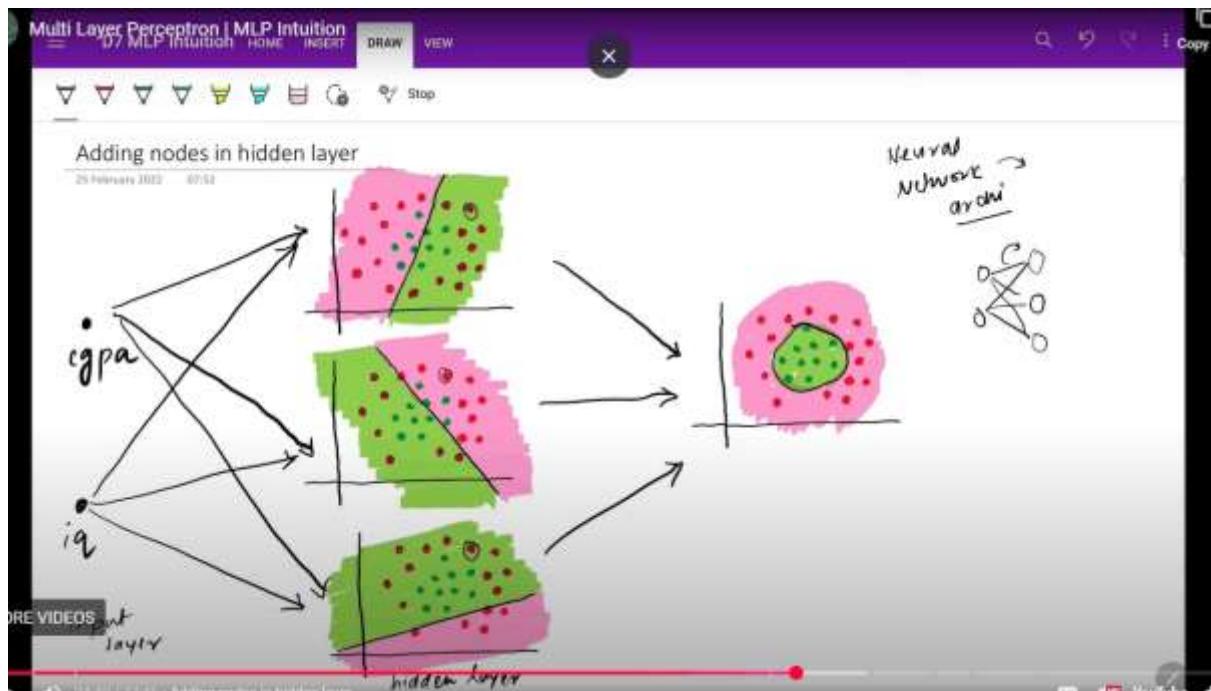
here perceptron 1 has more impact than perceptron 2

to do so we add weights to perceptron 1 and also bias

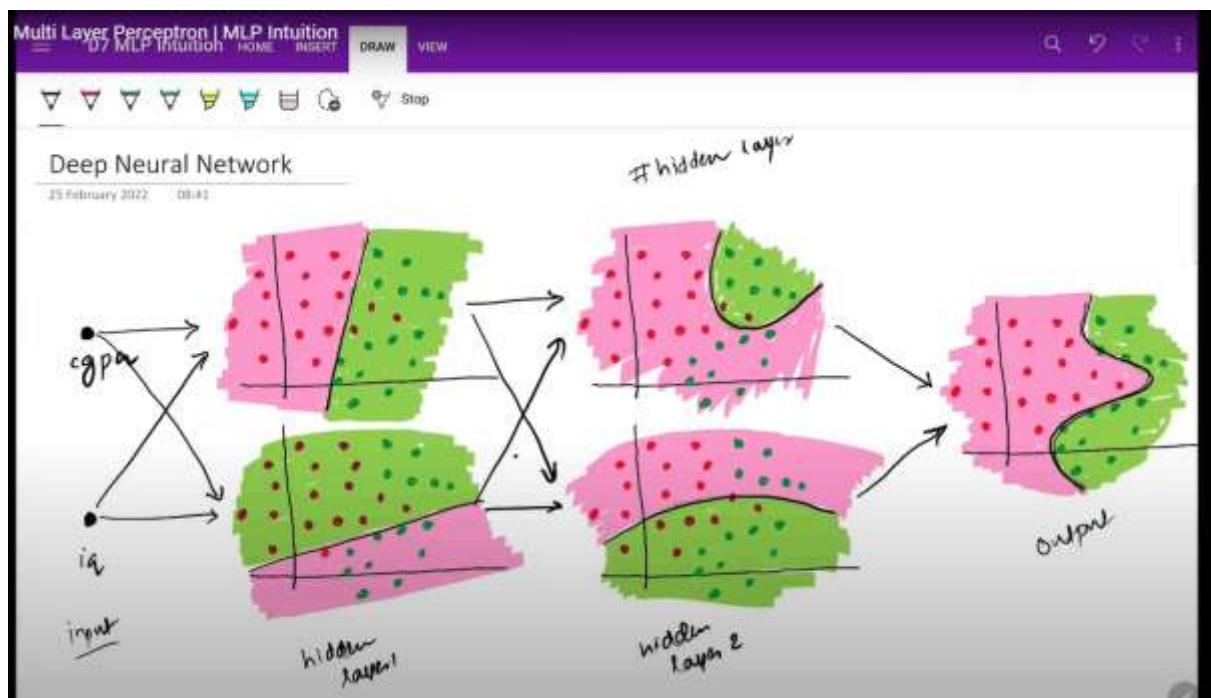
and then pass to sigmoid

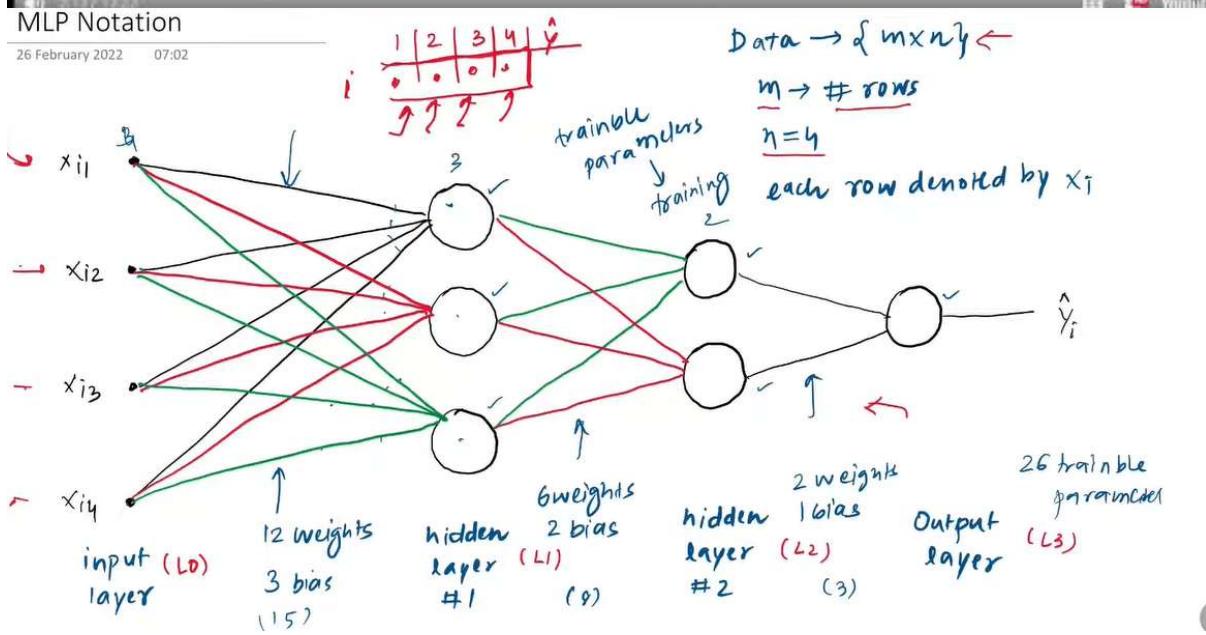
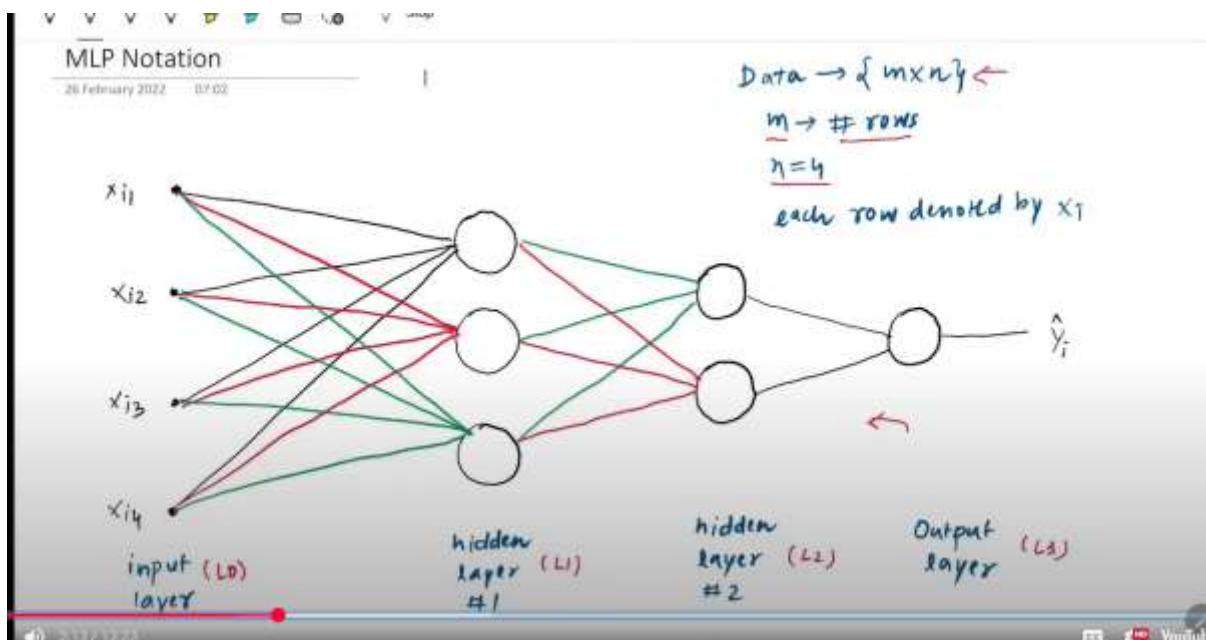
so in short the area marked in red below became a new perceptron

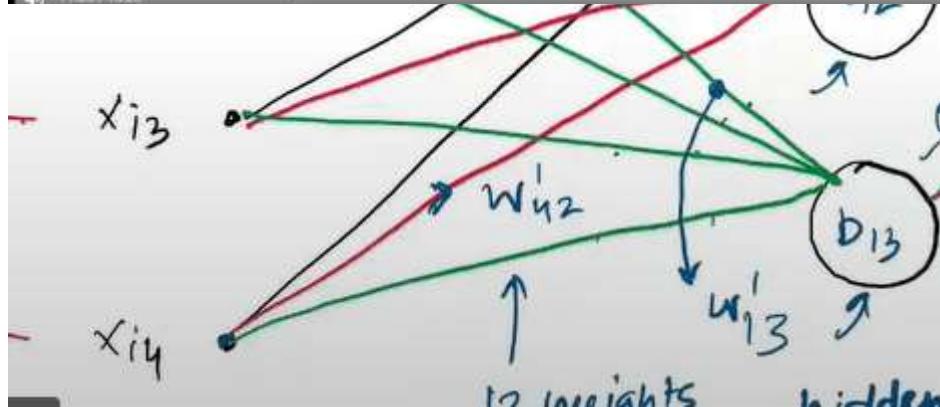
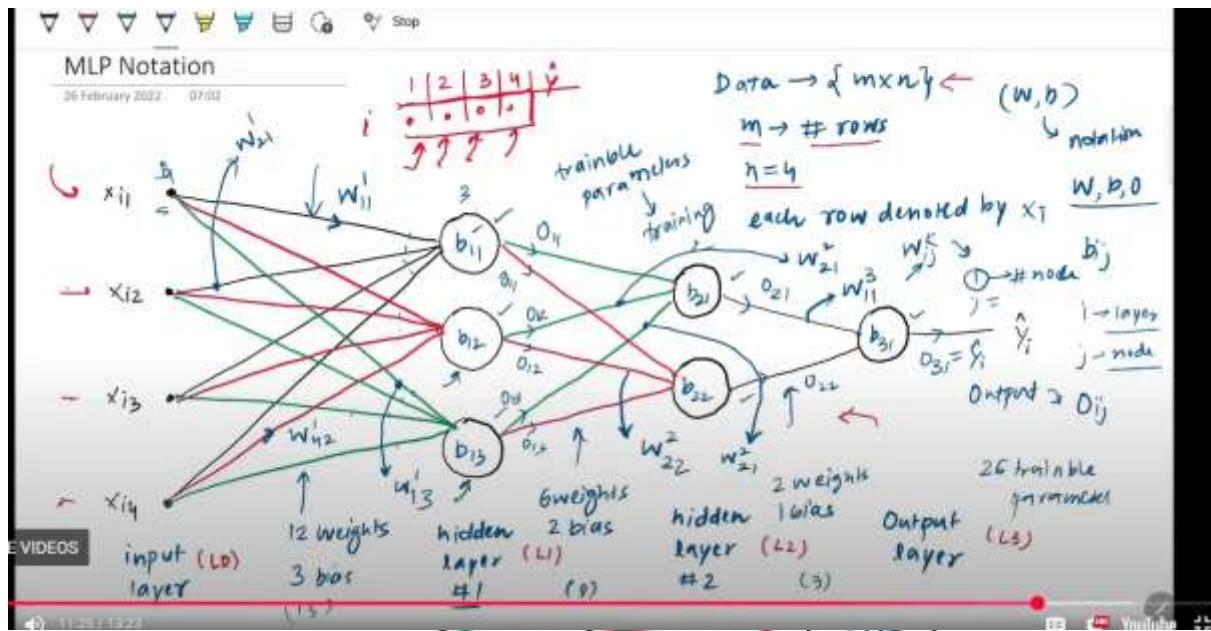




complex non linear data







here  $W_{142}$  means

1 -> going to layer 1

4 -> coming from node 4

2 -> going to node 2 g

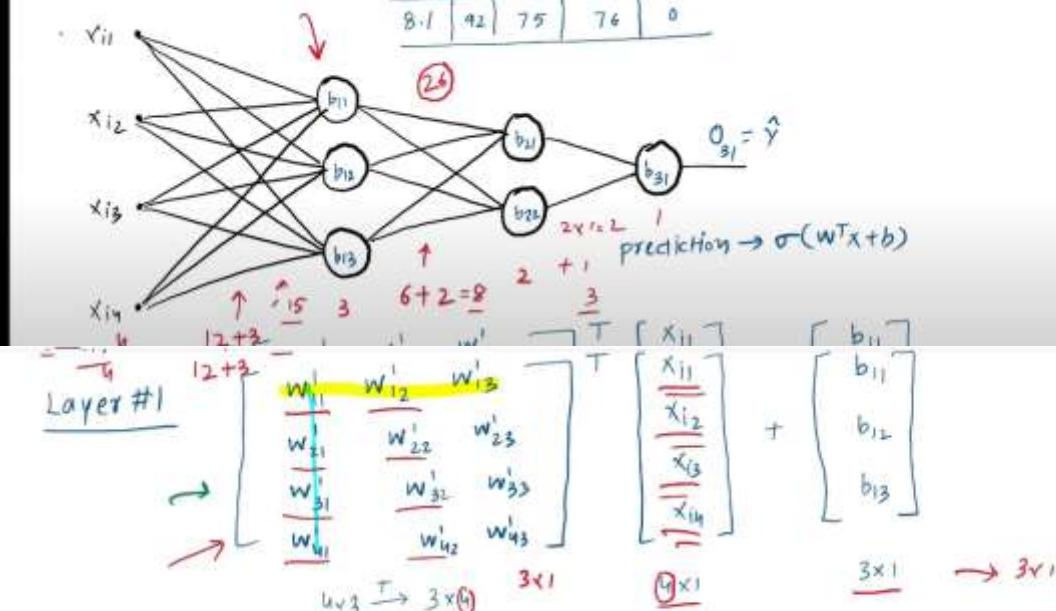
## Forward Propogartion

### Forward Propagation

03 March 2023 06:07

Cyphar	i9	10 <sup>th</sup> m	12 <sup>th</sup> m	placed
7.1	72	69	81	1
8.1	42	75	76	0

→ # of trainable parameters



$$= \begin{bmatrix} W_{11}^T x_{i1} + W_{12}^T x_{i2} + W_{13}^T x_{i3} + W_{14}^T x_{i4} \\ W_{21}^T x_{i1} + W_{22}^T x_{i2} + W_{23}^T x_{i3} + W_{24}^T x_{i4} \\ W_{31}^T x_{i1} + W_{32}^T x_{i2} + W_{33}^T x_{i3} + W_{34}^T x_{i4} \end{bmatrix} + \begin{bmatrix} b_{11} \\ b_{12} \\ b_{13} \end{bmatrix}$$

$$= \left( \begin{bmatrix} W_{11}^T x_{i1} + W_{12}^T x_{i2} + W_{13}^T x_{i3} + W_{14}^T x_{i4} + b_{11} \\ W_{21}^T x_{i1} + W_{22}^T x_{i2} + W_{23}^T x_{i3} + W_{24}^T x_{i4} + b_{12} \\ W_{31}^T x_{i1} + W_{32}^T x_{i2} + W_{33}^T x_{i3} + W_{34}^T x_{i4} + b_{13} \end{bmatrix} \right)$$

$$= \begin{bmatrix} o_{11} \\ o_{12} \\ o_{13} \end{bmatrix}$$

Layer #2

$$\begin{aligned}
 & \left[ \begin{array}{cc} w_{11}^2 & w_{12}^2 \\ w_{21}^2 & w_{22}^2 \\ w_{31}^2 & w_{32}^2 \end{array} \right]^T \left[ \begin{array}{c} o_{11} \\ o_{12} \\ o_{13} \end{array} \right] + \left[ \begin{array}{c} b_{21} \\ b_{22} \end{array} \right] \rightarrow \\
 & = \sigma \left( \left[ \begin{array}{c} \frac{w_{11}^2 o_{11} + w_{21}^2 o_{12} + w_{31}^2 o_{13} + b_{21}}{w_{12}^2 o_{11} + w_{22}^2 o_{12} + w_{32}^2 o_{13} + b_{22}} \end{array} \right] \right) = \left[ \begin{array}{c} o_{21} \\ o_{22} \end{array} \right]
 \end{aligned}$$

Layer #3

$$\begin{aligned}
 & \rightarrow \left[ \begin{array}{c} w_{11}^3 \\ w_{21}^3 \end{array} \right]^T \left[ \begin{array}{c} o_{21} \\ o_{22} \end{array} \right] + \left[ \begin{array}{c} b_{31} \end{array} \right] \\
 & = \sigma \left( \left[ \begin{array}{c} w_{11}^3 o_{21} + w_{21}^3 o_{22} + b_{31} \end{array} \right] \right) = \hat{y}_i
 \end{aligned}$$

$a^{[1]}$  =  $\sigma(a^{[0]} w^{[1]} + b^{[1]})$

$a^{[2]}$  =  $\sigma(a^{[1]} w^{[2]} + b^{[2]})$

$a^{[3]}$  =  $\sigma(a^{[2]} w^{[3]} + b^{[3]})$

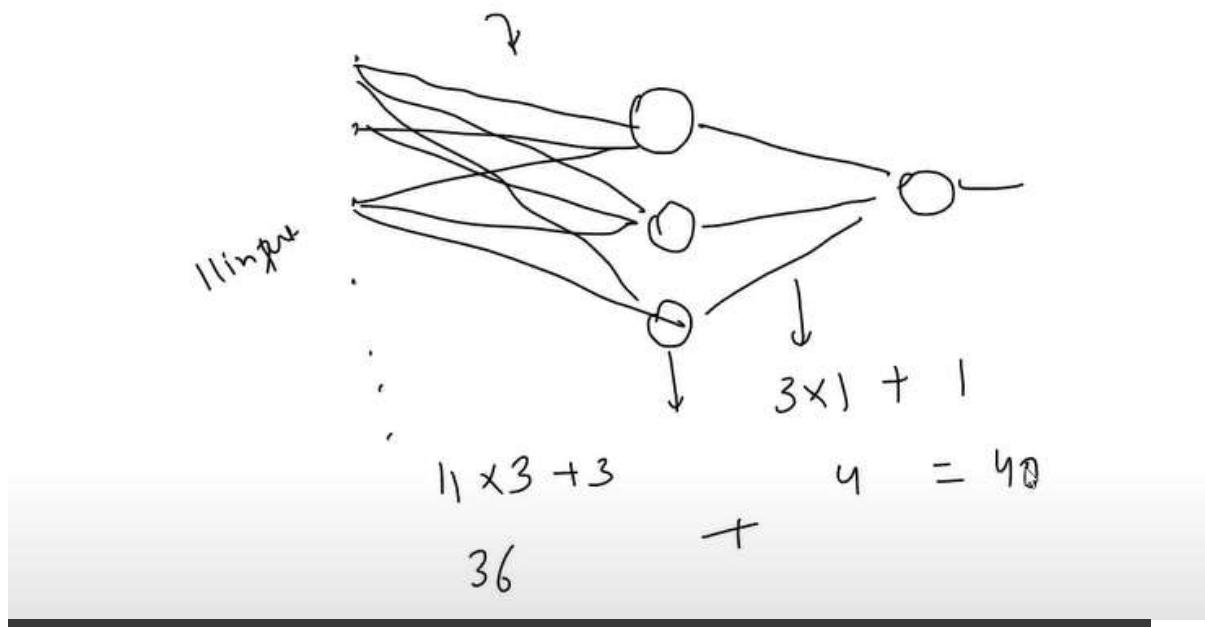
induction  $\rightarrow \sigma(w^T x + b)$

$$\sigma \left( \sigma \left( \left( \sigma \left( a^{[0]} w^{[1]} + b^{[1]} \right) \right) w^{[2]} + b^{[2]} \right) w^{[3]} + b^{[3]} \right)$$

$a^{[1]}$        $a^{[2]}$        $a^{[3]}$

EOS

$a[3]$  is final output



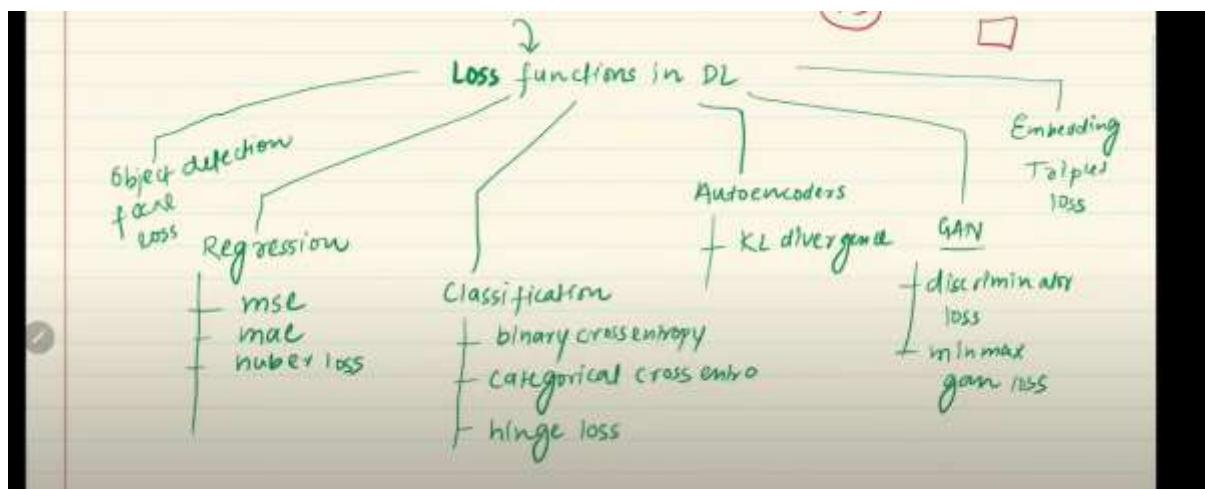
```
model.summary()
```

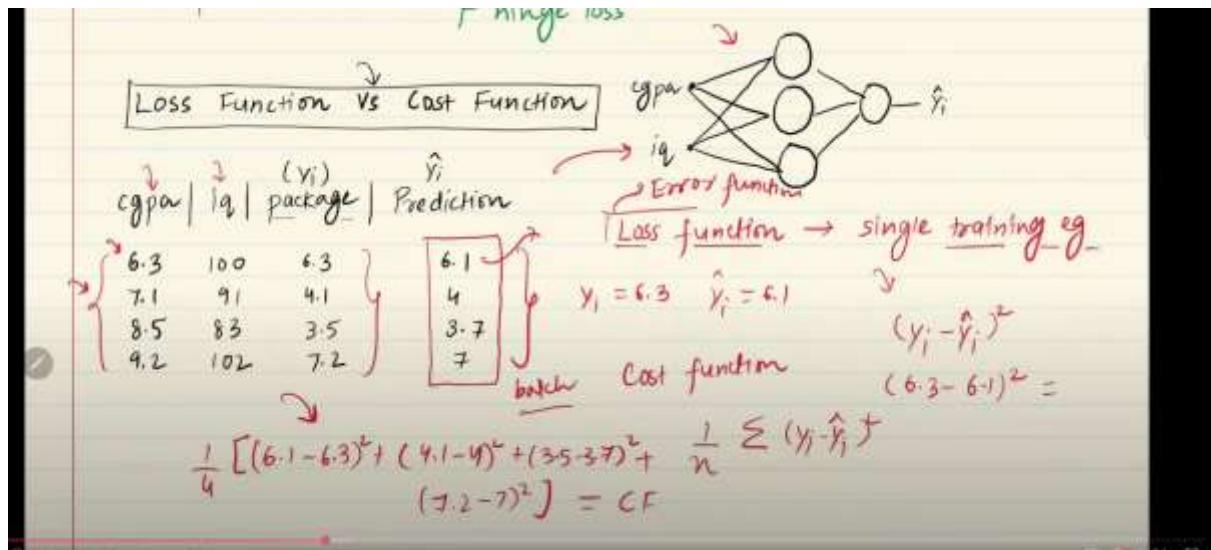
Model: "sequential"

Layer (type)	Output Shape	Param #
dense (Dense)	(None, 3)	36
dense_1 (Dense)	(None, 1)	4

Total params: 40 (160.00 B)  
Trainable params: 40 (160.00 B)  
Non-trainable params: 0 (0.00 B)

Loss functions based on problem

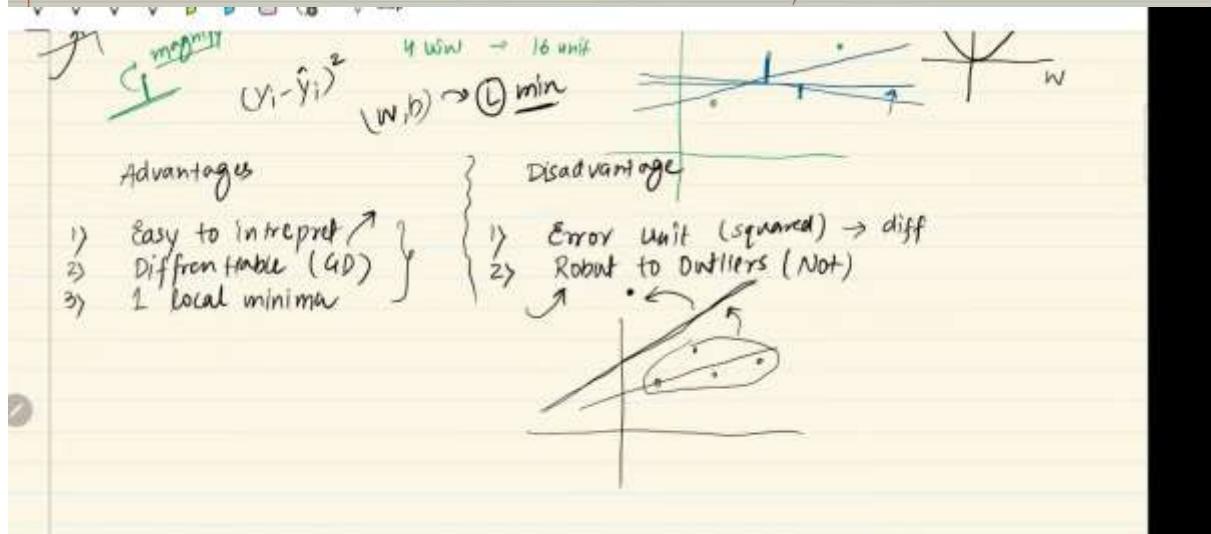
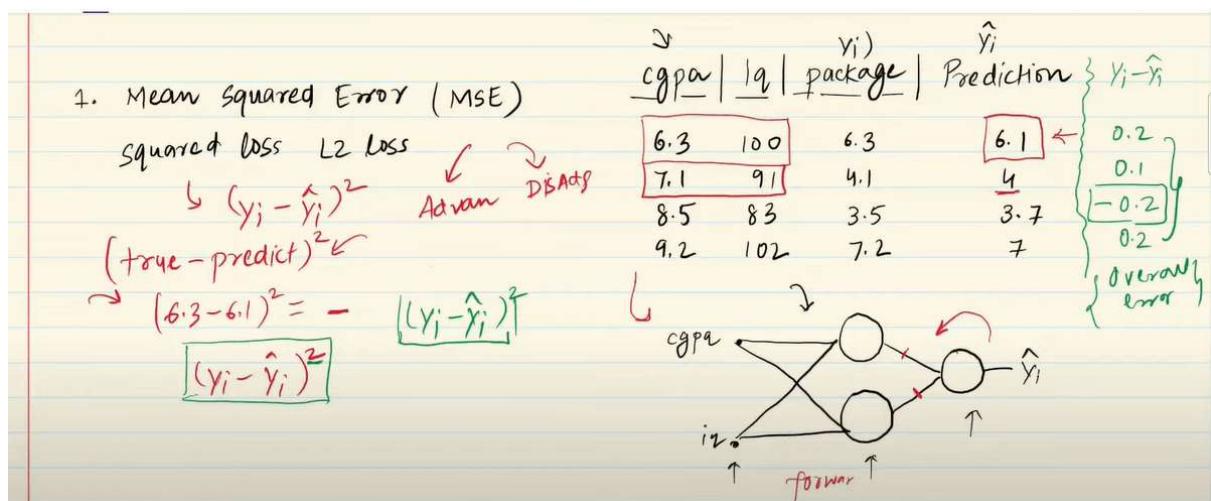




loss function is for single training example

cost function is for entire batch

so sum of all loss functions



for MSE last neuron should be linear type

```
► model = Sequential()  
  
model.add(Dense(7,activation='relu',input_dim=7))  
model.add(Dense(7,activation='relu'))  
model.add(Dense(1,activation='linear'))
```

+ Code

+ Markdown

[1]:

```
model.summary()
```

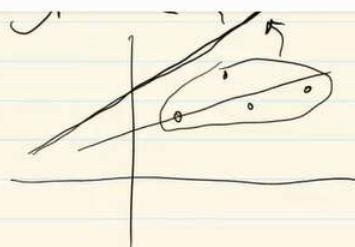
+ Code

+ Markdown

[1]:

```
model.compile(loss='mean_squared_error',optimizer='Adam')
```

CF =  $\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$



2. Mean Absolute Error (MAE) → L1 loss

$$L = |y_i - \hat{y}_i|$$

$$C = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

2. Mean Absolute Error (MAE) → L1 loss

$$L = |y_i - \hat{y}_i|$$

$$C = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|^2$$

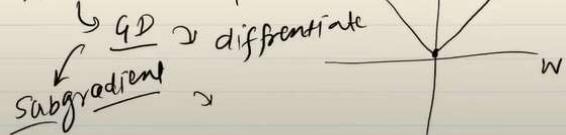
Advantages

- 1) Intuitive and easy
- 2) Unit → same → y
- 3) Robust to outliers

Disadvantages

- 1) Not differentiable

Subgradient



MSE doesn't handle outliers and MAE goes slowly towards outliers

Subgradiente ~

mae huber mse

3. Huber Loss ✓

$$L = \begin{cases} \frac{1}{2} (y - \hat{y})^2 & \text{for } |y - \hat{y}| \leq \delta \\ \delta |y - \hat{y}| - \frac{1}{2} \delta^2 & \text{otherwise} \end{cases}$$

hyper

The diagram illustrates the Huber loss function as a smooth approximation of the absolute difference (MAE) and the squared difference (MSE). It shows three regions: a central parabolic region where the loss is quadratic, a transition region where it follows the linear slope of the absolute value function, and an outer region where it becomes constant. Arrows point from the labels 'mae' and 'mse' to their respective regions in the graph.

Binary cross entropy also called logloss is used when

- classification problem
- output neuron is YES NO , 0 1
- output function should be sigmoid

4. Binary Cross Entropy

→ classification  
→ Two classes

	0	1
0	0	1
1	1	0
2	0	0

Loss function =  $-y \log(\hat{y}) - (1-y) \log(1-\hat{y})$   
 $y \rightarrow$  actual value / target  
 $\hat{y} \rightarrow$  NN prediction

cost function =  $-\frac{1}{n} \left[ \sum_{i=1}^n y_i \log \hat{y}_i + (1-y_i) \log(1-\hat{y}_i) \right]$

The diagram shows a neural network with two layers: hidden and output. The hidden layer has three neurons, and the output layer has one neuron. Arrows indicate the flow of activation from the hidden layer to the output layer. Above the output neuron, the text 'sigmoid' is written with an arrow pointing to it, indicating the activation function used for the output layer.

categorical cross entropy is used when

- multiclass classification problem
- output neuron is YES NO maybe
- output function should be softmax

```
▶ model = Sequential()
```

```
model.add(Flatten(input_shape=(28,28)))
model.add(Dense(128, activation='relu'))
model.add(Dense(32, activation='relu'))
model.add(Dense(10, activation='softmax'))
```

```
[ ] model.summary()
```

5. Categorical Cross Entropy [used in Softmax Regression]

→ Multi-class classification

1 point

$$L = - \sum_{j=1}^k y_j \log(\hat{y}_j)$$

where  $k$  is # classes in the data

cgpa	iq	placed?
8	80	Yes 1
6	60	No 2
7	70	Maybe 3

Yes	No	Maybe
1	0	0
0	1	0
0	0	1

$f(z) = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$

Activation neurons

Class ~ categories

→ Multi-class classification

1 point

$$L = - \sum_{j=1}^k y_j \log(\hat{y}_j)$$

where  $k$  is # classes in the data

1 point

$$L = - y_1 \log(\hat{y}_1) - y_2 \log(\hat{y}_2) - y_3 \log(\hat{y}_3)$$

ui	v	w
8	80	1
6	60	0
7	70	0

Yes	No	Maybe
1	0	0
0	1	0
0	0	1

$f(z) = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$

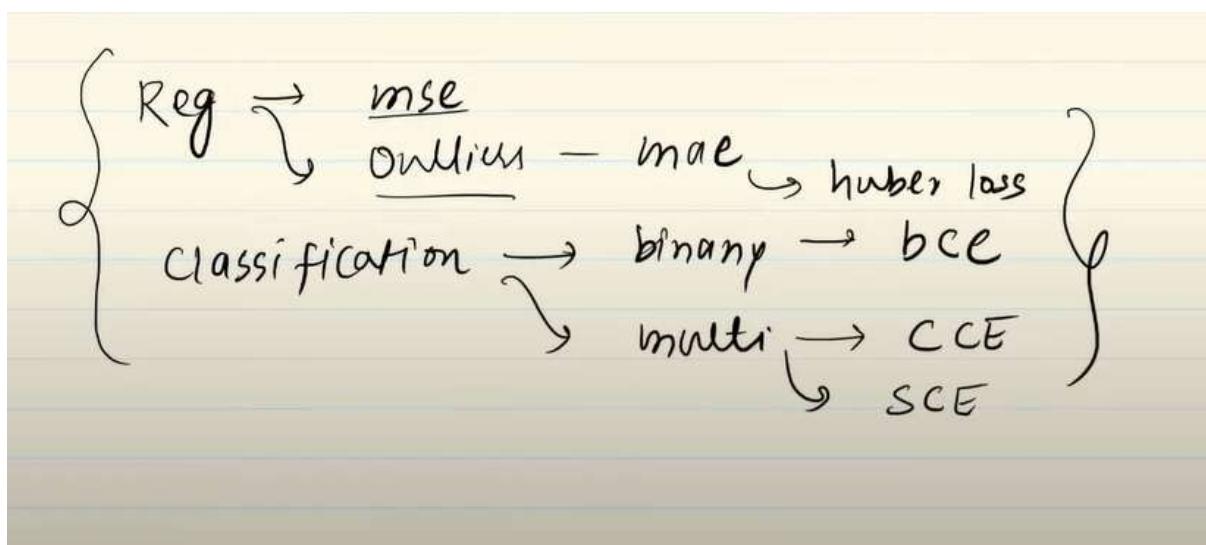
forward prop

Activation neurons

Class ~ categories

Sparse Categorical Cross Entropy			cgpa   iq	placed	DHE
④	1/2/3	[ 7 70 ]	Yes	1	
L ↗	[ 0.1 0.4 0.5 ] ↤	[ 8 80 ]	No	②	
	-1 × log(0.1)	[ 6 60 ]	Maybe	③	
		[ 0.1 0.4 0.6 ]			
SCI		-1 × log(0.4)			[ 0.6 0.2 0.2 ]
		-1 × log(0.2)			

if there are too many classes we should use sparse categorical cross entropy



**Backpropagation**, short for "backward propagation of errors", is a fundamental algorithm used to train artificial neural networks (ANNs). It plays a crucial role in enabling neural networks to learn from data and improve their predictive accuracy

## Training means finding correct values of weights and biases

1. **Forward Pass:** Input data is fed through the neural network, generating an output.
2. **Error Calculation:** The network's output is compared to the desired output, and a loss function calculates the error or difference between them.
3. **Backward Pass (Backpropagation):** The error is then propagated backward through the network, from the output layer towards the input layer. Using the chain rule of calculus, backpropagation calculates the gradient of the loss function with respect to

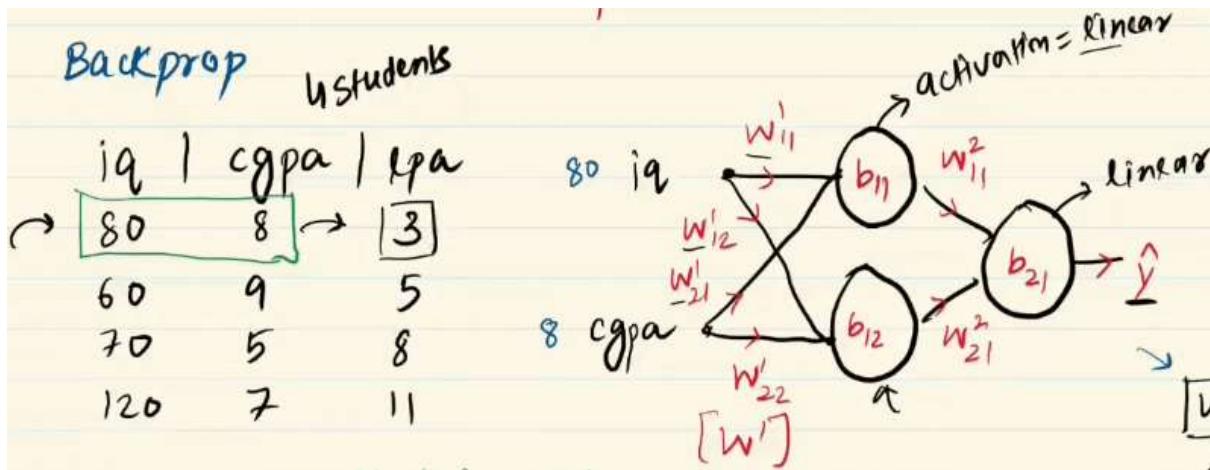
each weight and bias in the network. This essentially tells us how much each weight and bias contributes to the overall error.

4. **Weight and Bias Update:** Optimization algorithms like gradient descent or stochastic gradient descent use these gradients to adjust the weights and biases of the network. The goal is to minimize the loss function, thereby reducing the error and improving the network's performance.
5. **Iteration:** Steps 1-4 are repeated iteratively with multiple passes over the training data until the network learns to make accurate predictions

step 1 - initialize random values for weights and biases

u can also initialize 1 for weights and 0 for biases

step 2 - select 1 data entry and give neurons values as inputs



step 3 - do initial prediction with inputs , weights = 1 , biases = 0

for now assume prediction came out as 18LPA --- but it should be 3 LPA

step 4 - choose 1 loss function and calculate error

$$\mathcal{L} = (y - \hat{y})^2$$

$$(3 - 18)^2 = 225$$

error ↗

now to minimize loss function output ---> we need smaller  $\hat{y}$  value --->  $\hat{y}$  value is O21

---> O21 value is dependent on 5 variables ---> and tree grows into branches

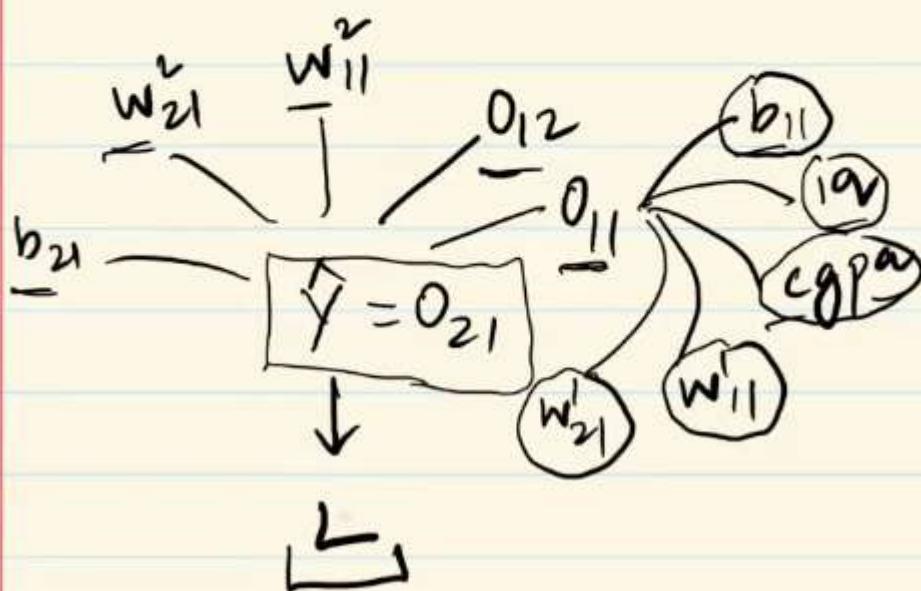
$$\hat{y} \quad \text{?}$$

$$O_{21} = w_{11}^2 O_{11} + w_{21}^2 O_{12} + b_{21}$$

Hierachony  
↓

2) Predict

3) Choose



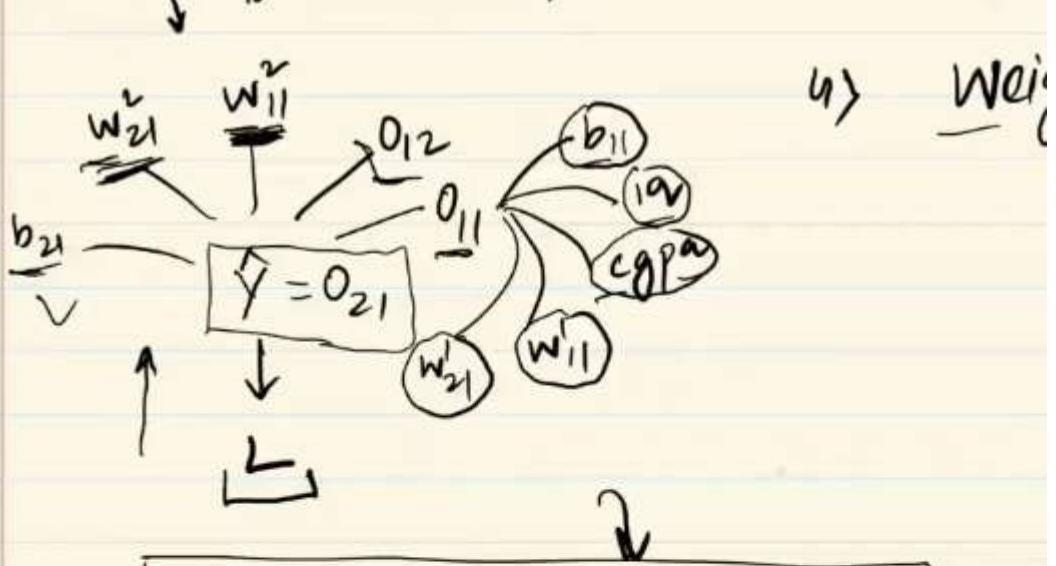
step 5 - update weights and biases using gradient descend

4) Weights and bias update ✓

↳ Gradient descent

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

$$b_{\text{new}} = b_{\text{old}} - \eta \frac{\partial L}{\partial b_{\text{old}}}$$



$$w_{11}^2_{\text{new}} = w_{11}^2_{\text{old}} - \eta \frac{\partial L}{\partial w_{11}^2}$$

$$w_{21}^2_{\text{new}} = w_{21}^2_{\text{old}} - \eta \frac{\partial L}{\partial w_{21}^2}$$

$$b_{21}_{\text{new}} = b_{21}_{\text{old}} - \eta \frac{\partial L}{\partial b_{21}}$$

derivative of loss wrt weight

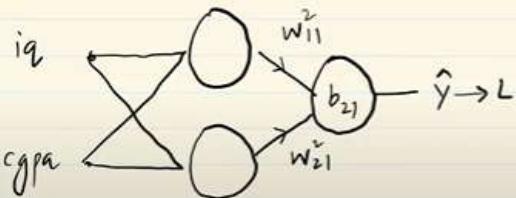
$$\frac{\partial L}{\partial w_{11}^2}$$

partial derivative

this is derivative of loss wrt weights

for example below is 9 derivatives ( 3 for each perceptron )

$$\frac{\partial L}{\partial w_{11}^2}, \frac{\partial L}{\partial w_{21}^2}, \frac{\partial L}{\partial b_{21}} \quad \left| \quad \frac{\partial L}{\partial w_{11}^1}, \frac{\partial L}{\partial w_{21}^1}, \frac{\partial L}{\partial b_{11}} \right| \quad \left| \quad \frac{\partial L}{\partial w_{12}^1}, \frac{\partial L}{\partial w_{22}^1}, \frac{\partial L}{\partial b_{12}} \right.$$

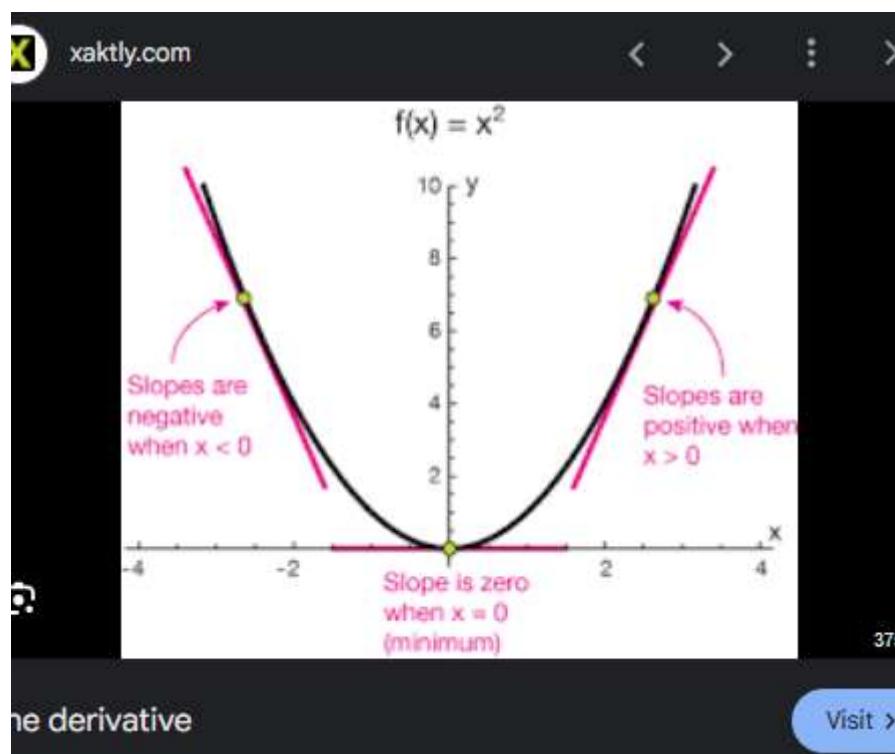


$\square \quad \square = 1$

$$w_{11}^2_{new} = w_{11}^2_{old} - \eta \frac{\partial L}{\partial w_{11}^2}$$

↑  
derivative

dy/dx says that upon a small change in X what changes in Y



so accordingly below derivative means that when weights are changed what is its impact on L (loss)

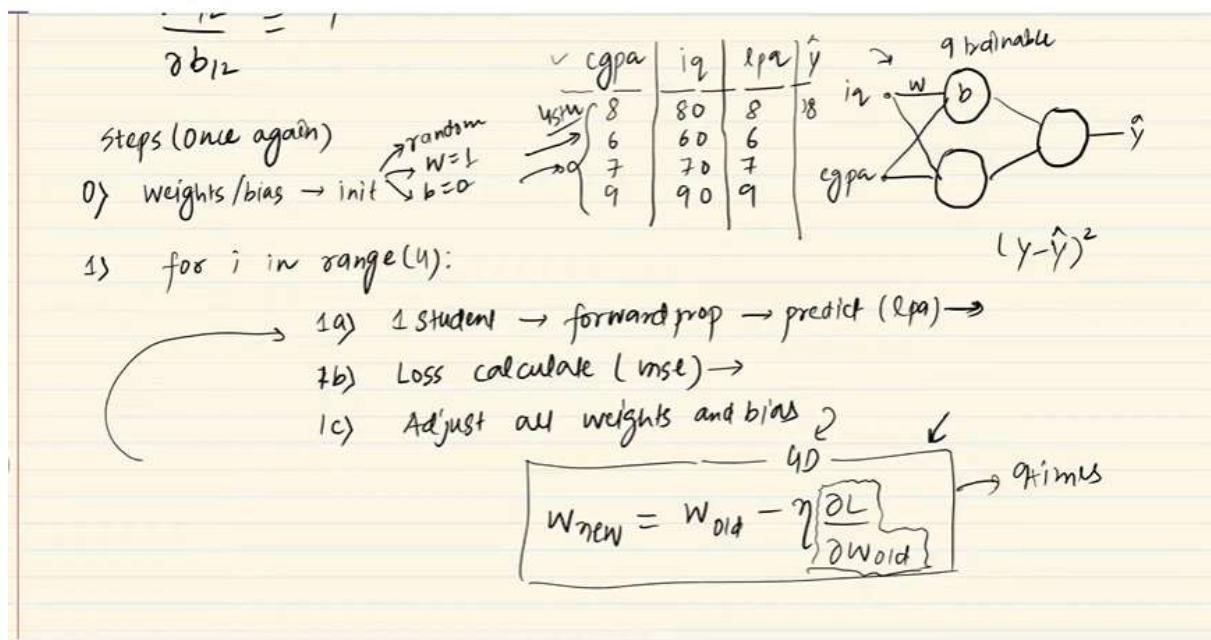
$$\frac{\partial L}{\partial w_{11}^2}$$

this above is written as

$$\frac{\partial L}{\partial w_{11}^2} = \left[ \frac{\partial L}{\partial \hat{y}} \right] \times \left[ \frac{\partial \hat{y}}{\partial w_{11}^2} \right] \rightarrow \text{chain rule of diffren}$$

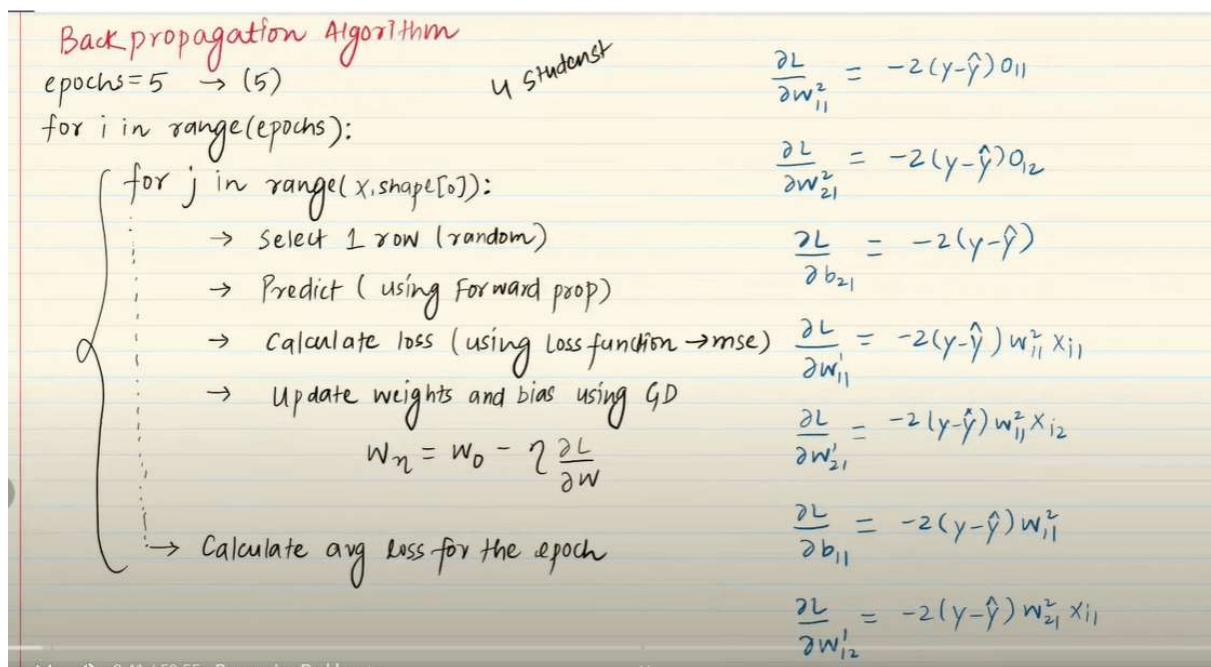
so weights change -> impact on  $\hat{y}$  ->  $\hat{y}$  changes -> impact on L loss

$$\begin{aligned}
 1. \quad & \frac{\partial L}{\partial w_{11}^2} = -2(y - \hat{y}) o_{11} \\
 2. \quad & \frac{\partial L}{\partial w_{21}^2} = -2(y - \hat{y}) o_{12} \\
 3. \quad & \frac{\partial L}{\partial b_{21}} = -2(y - \hat{y}) \\
 4. \quad & \frac{\partial L}{\partial w_{11}^2} = \left[ \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_{11}} \frac{\partial o_{11}}{\partial w_{11}^2} \right] \rightarrow -2(y - \hat{y}) w_{11}^2 x_{i1} \\
 5. \quad & \frac{\partial L}{\partial w_{21}^2} = \left[ \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_{12}} \frac{\partial o_{12}}{\partial w_{21}^2} \right] \rightarrow -2(y - \hat{y}) w_{21}^2 x_{i2} \\
 6. \quad & \frac{\partial L}{\partial b_{11}} = \left[ \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_{11}} \frac{\partial o_{11}}{\partial b_{11}} \right] \rightarrow -2(y - \hat{y}) w_{11}^2 \\
 7. \quad & \frac{\partial L}{\partial w_{12}^2} = \left[ \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_{12}} \frac{\partial o_{12}}{\partial w_{12}^2} \right] \rightarrow -2(y - \hat{y}) w_{12}^2 x_{i1} \\
 8. \quad & \frac{\partial L}{\partial w_{22}^2} = \left[ \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_{12}} \frac{\partial o_{12}}{\partial w_{22}^2} \right] \rightarrow -2(y - \hat{y}) w_{22}^2 x_{i2} \\
 9. \quad & \frac{\partial L}{\partial b_{12}} = \left[ \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_{12}} \frac{\partial o_{12}}{\partial b_{12}} \right] \rightarrow -2(y - \hat{y}) w_{12}^2
 \end{aligned}$$



step 1 will be done 100/1000 times till convergence is not met

this is called epochs



```

[22] # epochs implementation

parameters = initialize_parameters([2,2,1])
epochs = 5

for i in range(epochs):

    Loss = []

    for j in range(df.shape[0]):

        X = df[['cgpa', 'profile_score']].values[j].reshape(2,1) # Shape(no of features, no. of
        training example)
        y = df[['lpa']].values[j][0]

        # Parameter initialization

        y_hat,A1 = l_layer_forward(X,parameters)
        y_hat = y_hat[0][0]

        update_parameters(parameters,y,y_hat,A1,X)

        Loss.append((y-y_hat)**2)

    print('Epoch - ',i+1,'Loss - ',np.array(Loss).mean())

parameters

```

→ Concept of Gradient → Gradient descent  
 fancy word for derivative

$$Y = f(x) = x^2 + x$$

$$\frac{dy}{dx} = \frac{d}{dx}(fx) = \frac{d}{dx}(x^2 + x) = 2x + 1$$

derivative  $y \rightarrow x \rightarrow$  derivative  $\frac{d}{dx}$

$$z = f(x,y) = \sqrt{x^2 + y^2}$$

$$\frac{\partial z}{\partial x} = 2x \quad \frac{\partial z}{\partial y} = 2y$$

gradient  $\frac{\partial}{\partial x}$

derivative  $y \rightarrow x \rightarrow$  derivative  $\frac{d}{dx}$

gradient complex  $\frac{\partial L}{\partial w}, \frac{\partial L}{\partial b}$

$L(w_1^1, w_2^1, \dots, w_1^2, w_2^2, b_1, b_2)$

9D function → 9 diff slopes wrt each dim

3D →  $z = x^2 + y^2 \quad z = f(x,y)$

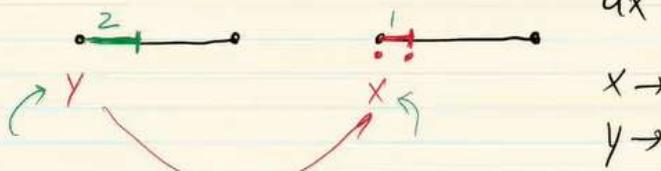


→ Concept of Derivative → Derivative at a point  
 ↑ intuition

$\frac{dy}{dx}$  (y) func x (x)  
 rate of change

$\frac{dy}{dx} = 2$  (+ve)  
 mag sign

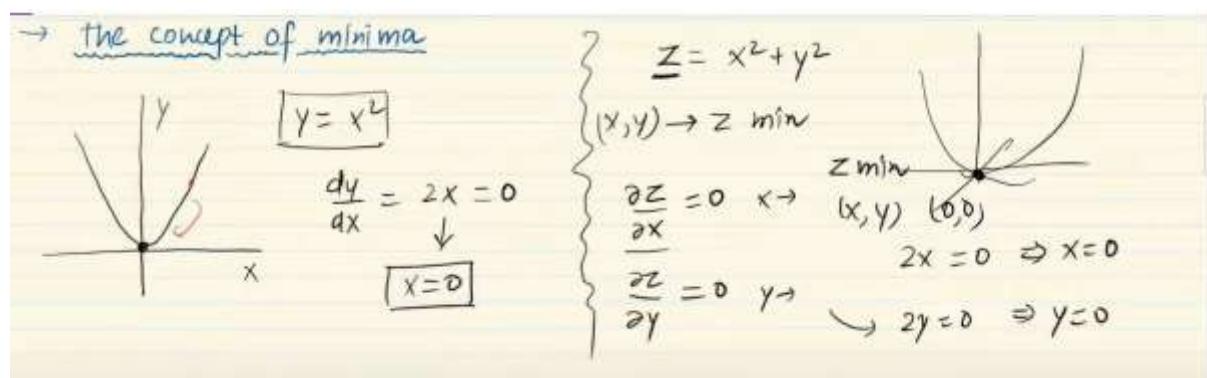
$\frac{dy}{dx} = -2$  (-ve)



suppose  $dy/dx = 2$  ..... so on 1 times change of x there will be 2 times change in y

2 is magnitude

so 2 is magnitude / slope / rate of change



to find minima we just say =0

so for  $y=x^2$

we do  $dy/dx = 2x = 0$

so  $x=0$

$L$  (q param)  $L \downarrow$   $\frac{\partial L}{\partial w_{11}} \dots \frac{\partial L}{\partial b_{12}} = 0$  qdim 2 minima

for example assume we are changing b21

The diagram illustrates the sensitivity of the cost function  $L$  to changes in parameter  $b_{21}$ . A graph shows  $L$  as a parabola opening upwards, with the minimum labeled  $(b_{21})$  and  $L \text{ min}$ .

Derivative calculation:

$$\frac{\partial L}{\partial b_{21}} = \boxed{b_2} \quad \boxed{\frac{\partial L}{\partial b_{21}}} \quad b_{21} = 5 \times$$

Sign analysis:

- $\frac{\partial L}{\partial b_{21}}$  is positive (+ve) when  $b_{21}$  is positive.
- $\frac{\partial L}{\partial b_{21}}$  is negative (-ve) when  $b_{21}$  is negative.

Effect of change in  $b_{21}$ :

Smart L derivative w.r.t  $b_{21}$

Unit change in  $b_{21}$  leads to a change in  $L$  (indicated by arrows). The sign of the derivative determines the direction of change:

- If  $\frac{\partial L}{\partial b_{21}} < 0$ , then  $L \downarrow$  when  $b_{21} \downarrow$  and  $b_{21} \uparrow$ .
- If  $\frac{\partial L}{\partial b_{21}} > 0$ , then  $L \uparrow$  when  $b_{21} \uparrow$  and  $b_{21} \downarrow$ .

Equations for partial derivatives:

$$\frac{\partial L}{\partial b_{21}} = +ve \quad \frac{\partial L}{\partial b_{21}} = -ve$$

$$\frac{\partial L}{\partial b_{21}} = -ve \quad \frac{\partial L}{\partial b_{21}} = +ve$$

if derivative of  $L/b_{21}$  is +ve then we decrease  $b_{21}$  value

and then subtract it from old b21 value

if derivative of  $L/b_{21}$  is -ve then we increase  $b_{21}$  value

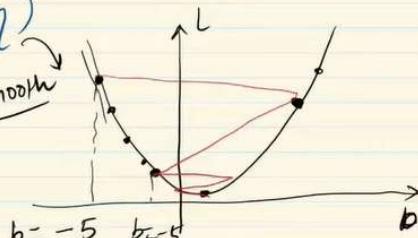
and then subtract it from old b21 value ( bcz - + is - anyways )

## Effect of Learning Rate ( $\eta$ )

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

$$\frac{\partial L}{\partial b} = -50$$

$$b = -5 - (-50) = 45$$



$$\frac{\partial L}{\partial b} = 50$$

$$b = 45$$

$$b = 45 - 50$$

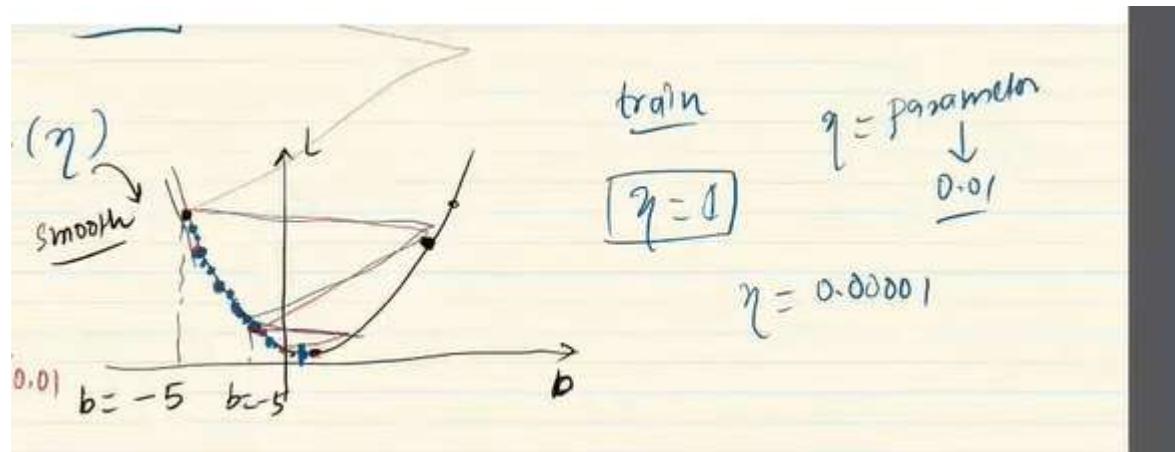
if formula does not have learning rate

then first  $b = -5$

then  $b = 45$

then  $b = 5$

so its very zig zag



▽ ▷ ▲ △ ▵ ▶ ▶ Stop

→ What is convergence?

$$w_n = w_0 - \eta \frac{\partial L}{\partial w} \approx 0$$

while → con → for 1000



[epoch - 100]  
for 1000

convergence is point where wold and wnew has not much difference

now we dont know exactly how many loops will it take to reach that

so we run it 100/1000 epochs time

[developers-dot-devsite-v2-prod.appspot.com/machine-learning/crash-course/backprop-scroll](https://developers-dot-devsite-v2-prod.appspot.com/machine-learning/crash-course/backprop-scroll)

## Gradient Descent

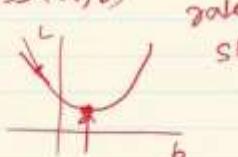
Thursday, April 7, 2023 7:44 AM

Optimization

Gradient descent is one of the most popular algorithms to perform optimization and by far the most common way to optimize neural networks.

$J(\theta) = \text{LOSS function}$  Neural Loss ( $w, b$ ) rate SLP

Gradient descent is a way to minimize an objective function  $J(\theta)$  parameterized by a model's parameters  $\theta \in \mathbb{R}^d$  by updating the parameters in the opposite direction of the gradient of the objective function  $\nabla_{\theta} J(\theta)$  w.r.t. to the parameters. The learning rate  $\eta$  determines the size of the steps we take to reach a (local) minimum. In other words, we follow the direction of the slope of the surface created by the objective function downhill until we reach a valley.



## Back propagation Algorithm

epochs = 5

for i in range(epochs):

    for j in range(x.shape[0]):

        → Select 1 row (random)

        → Predict (using Forward prop)

        → Calculate loss (using loss function → mse)

        → Update weights and bias using GD

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

    → Calculate avg loss for the epoch  
 $L_{avg}$

0

$\{ \text{Batch} / \text{Stochastic} / \text{mini batch} \}$

Batch GD  
nse) epoch = 5 → current weight

↳ 50 points → predict ↓  
dot product

$y_{\text{hat}} = \text{np.dot}(X, w) + b$   
↳ 50 predict

1 single  $y = 50$  actual val.

( $w, b$ )

$y$   $y_{\text{hat}}$   $\rightarrow$  loss

off

$$f = -\sum_{i=1}^{50}$$

i have curr weights and biases

in a single epoch using current weight and biases all predictions are made for entire dataset

then using 50 new predictions 1 single time weights and biases are updated

so no. of epochs = no. of updates

5 epoch = 5 updates

i have curr weights and biases

in a single epoch using current weight and biases 50 times loop runs everytimes selects 1 random point -> calculates loss and updates w and b

so suppose 10 epoch and 50 rows

then  $10 \times 50 = 500$  weights , biases updates

suppose i have 50 rows

then in code

batch size = 50 means batch gradient descendent

batch size = 1 means SGD

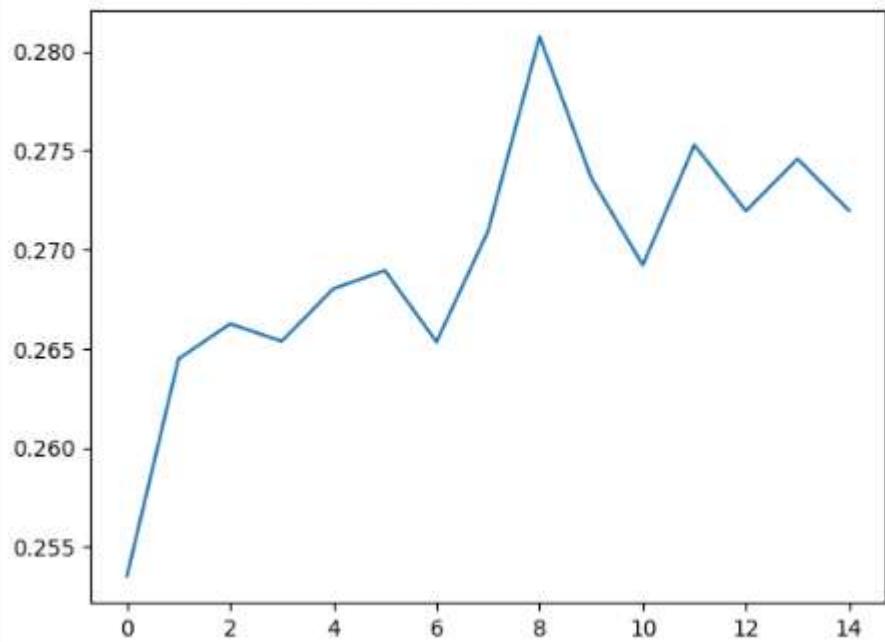
SGD

```
18s  model.compile(loss='binary_crossentropy',metrics=['accuracy'])
      start = time.time()
      history = model.fit(X_scaled,y,epochs=15,batch_size=1,validation_split=0.2)
      print(time.time() - start)

      Epoch 1/15
320/320    2s 4ms/step - accuracy: 0.8601 - loss: 0.2756 - val_accuracy: 0.9750 - val_loss: 0.1789
Epoch 2/15
320/320    1s 4ms/step - accuracy: 0.8688 - loss: 0.3250 - val_accuracy: 0.9750 - val_loss: 0.1809
Epoch 3/15
320/320    1s 3ms/step - accuracy: 0.8676 - loss: 0.2654 - val_accuracy: 0.9750 - val_loss: 0.1705
Epoch 4/15
320/320    1s 2ms/step - accuracy: 0.8794 - loss: 0.2818 - val_accuracy: 0.9750 - val_loss: 0.1657
Epoch 5/15
320/320    1s 2ms/step - accuracy: 0.8969 - loss: 0.2699 - val_accuracy: 0.9750 - val_loss: 0.1653
Epoch 6/15
320/320    1s 2ms/step - accuracy: 0.9141 - loss: 0.2210 - val_accuracy: 0.9750 - val_loss: 0.1637
Epoch 7/15
320/320    1s 2ms/step - accuracy: 0.8863 - loss: 0.2333 - val_accuracy: 0.9875 - val_loss: 0.1672
Epoch 8/15
320/320    1s 2ms/step - accuracy: 0.8721 - loss: 0.2865 - val_accuracy: 0.9875 - val_loss: 0.1621
Epoch 9/15
320/320    1s 2ms/step - accuracy: 0.8873 - loss: 0.2841 - val_accuracy: 0.9875 - val_loss: 0.1596
Epoch 10/15
320/320   1s 2ms/step - accuracy: 0.8693 - loss: 0.2953 - val_accuracy: 0.9875 - val_loss: 0.1541
Epoch 11/15
320/320   1s 2ms/step - accuracy: 0.8900 - loss: 0.2774 - val_accuracy: 0.9875 - val_loss: 0.1559
Epoch 12/15
320/320   1s 3ms/step - accuracy: 0.8770 - loss: 0.2821 - val_accuracy: 0.9875 - val_loss: 0.1524
Epoch 13/15
320/320   2s 5ms/step - accuracy: 0.8711 - loss: 0.3325 - val_accuracy: 0.9875 - val_loss: 0.1506
Epoch 14/15
320/320   2s 2ms/step - accuracy: 0.9048 - loss: 0.2560 - val_accuracy: 0.9875 - val_loss: 0.1516
Epoch 15/15
320/320   1s 2ms/step - accuracy: 0.8960 - loss: 0.2415 - val_accuracy: 0.9875 - val_loss: 0.1464
18.50076150894165
```

```
import matplotlib.pyplot as plt
plt.plot(history.history['loss'])
```

```
[<matplotlib.lines.Line2D at 0x7819561569d0>]
```



# Batch

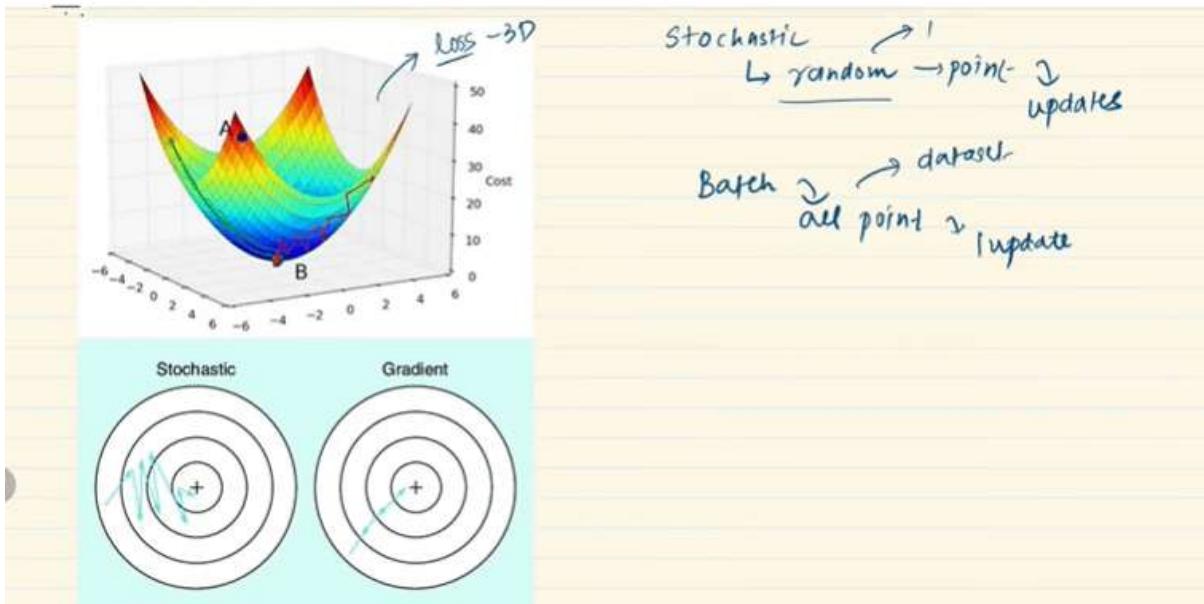
```
✓ 3s  model.compile(loss='binary_crossentropy',metrics=['accuracy'])
    start = time.time()
    history = model.fit(X_scaled,y,epochs=15,batch_size=400,validation_split=0.2)
    print(time.time() - start)

    Epoch 1/15
1/1    1s 1s/step - accuracy: 0.8969 - loss: 0.2646 - val_accuracy: 0.9875 - val_loss: 0.1507
Epoch 2/15
1/1    0s 102ms/step - accuracy: 0.8969 - loss: 0.2615 - val_accuracy: 0.9875 - val_loss: 0.1537
Epoch 3/15
1/1    0s 99ms/step - accuracy: 0.8938 - loss: 0.2597 - val_accuracy: 0.9875 - val_loss: 0.1563
Epoch 4/15
1/1    0s 157ms/step - accuracy: 0.8938 - loss: 0.2582 - val_accuracy: 0.9875 - val_loss: 0.1586
Epoch 5/15
1/1    0s 112ms/step - accuracy: 0.8938 - loss: 0.2570 - val_accuracy: 0.9875 - val_loss: 0.1607
Epoch 6/15
1/1    0s 127ms/step - accuracy: 0.8938 - loss: 0.2559 - val_accuracy: 0.9875 - val_loss: 0.1626
Epoch 7/15
1/1    0s 97ms/step - accuracy: 0.8938 - loss: 0.2550 - val_accuracy: 0.9875 - val_loss: 0.1644
Epoch 8/15
1/1    0s 85ms/step - accuracy: 0.8938 - loss: 0.2542 - val_accuracy: 0.9875 - val_loss: 0.1661
Epoch 9/15
1/1    0s 86ms/step - accuracy: 0.8969 - loss: 0.2534 - val_accuracy: 0.9875 - val_loss: 0.1678
Epoch 10/15
1/1    0s 138ms/step - accuracy: 0.8969 - loss: 0.2526 - val_accuracy: 0.9875 - val_loss: 0.1694
Epoch 11/15
1/1    0s 87ms/step - accuracy: 0.8969 - loss: 0.2520 - val_accuracy: 0.9875 - val_loss: 0.1709
Epoch 12/15
1/1    0s 140ms/step - accuracy: 0.8969 - loss: 0.2514 - val_accuracy: 0.9875 - val_loss: 0.1723
Epoch 13/15
1/1    0s 135ms/step - accuracy: 0.8938 - loss: 0.2508 - val_accuracy: 0.9875 - val_loss: 0.1737
Epoch 14/15
1/1    0s 141ms/step - accuracy: 0.8938 - loss: 0.2502 - val_accuracy: 0.9875 - val_loss: 0.1751
Epoch 15/15
1/1    0s 140ms/step - accuracy: 0.8938 - loss: 0.2497 - val_accuracy: 0.9875 - val_loss: 0.1764
3.0284104347229004
3.0284104347229004
```

```
✓ 0s  import matplotlib.pyplot as plt
    plt.plot(history.history['loss'])

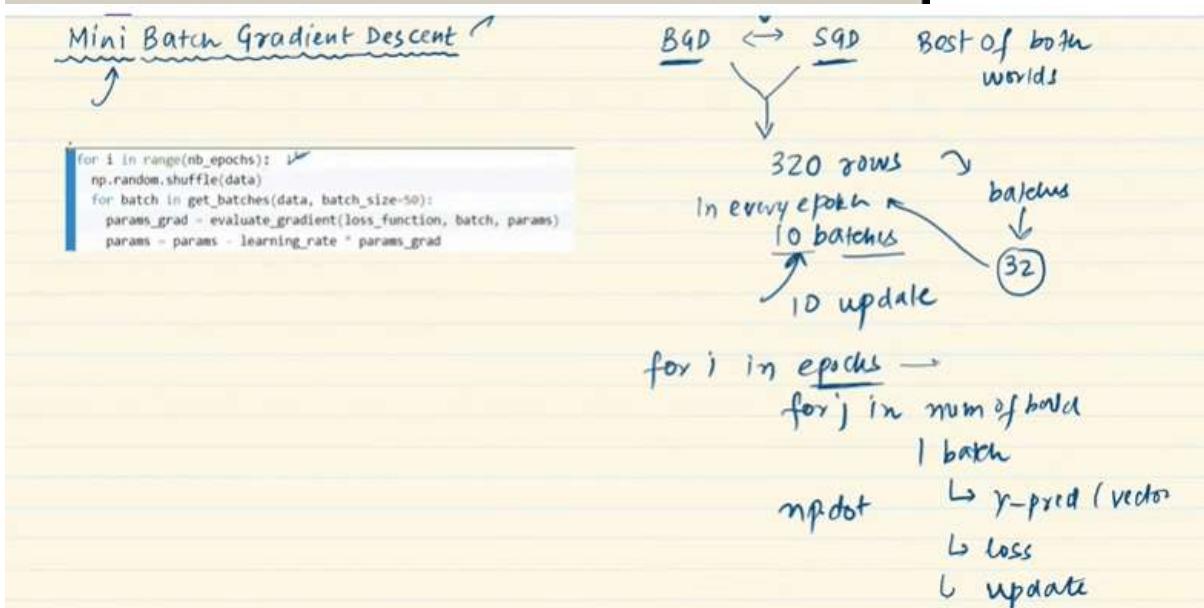
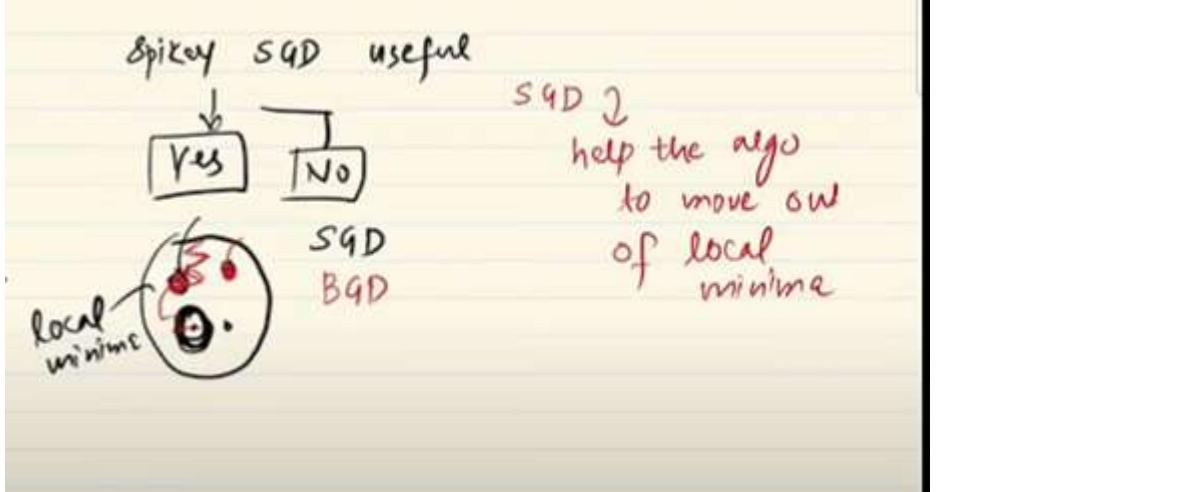
    []
```

Epoch	Loss
0	0.2645
2	0.2605
4	0.2565
6	0.2545
8	0.2525
10	0.2510
12	0.2505
14	0.2500



Stochastic  
random → point updates

Batch  
all point → update



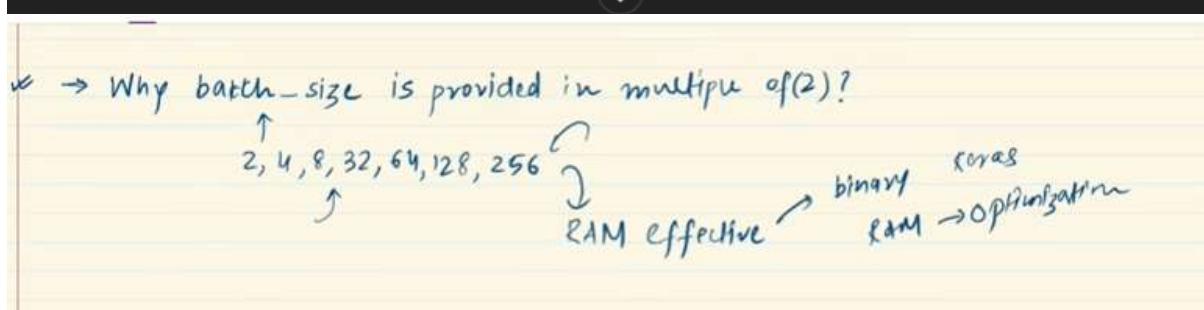
Let's say you have:

- A dataset with 1000 samples
- A batch size = 50

**Step-by-step:**

1. Initialize weights and biases.
2. Epoch 1 starts:
  - o Shuffle the data (optional but common).
  - o Split it into mini-batches of size 50 → you'll get  $1000 / 50 = 20$  mini-batches.
  - o For each mini-batch (i.e., 50 samples at a time):
    - Make predictions using current weights and biases.
    - Calculate the loss and gradients.
    - **Update weights and biases ONCE using that mini-batch.**
  - o After all 20 mini-batches, epoch 1 ends

Aspect	Batch GD	Stochastic GD (SGD)	Mini-Batch GD
Time per Epoch	● Slowest (processes all data once)	● Fastest (1 sample = 1 update)	● Moderate (processes batch size)
Updates per Epoch	1	= Number of samples	= Number of samples / batch size
Convergence Speed	● Slow (fewer updates)	● Fast per update (high frequency)	● Balanced
Convergence Stability	● Stable (smooth updates)	● Unstable (noisy, high variance)	● Moderate stability
Final Accuracy	● Often good	● May fluctuate, may not settle well	● Often good and practical
Memory Usage	● High (entire dataset in memory)	● Low (1 sample at a time)	● Medium (batch size in memory)
GPU Efficiency	● Poor	● Poor	● Excellent (vectorized operations)
Practical Use	Rare in large data	Rare in pure form	● Standard in practice



## Vanishing Gradient Problem

Thursday, April 7, 2022 7:45 AM

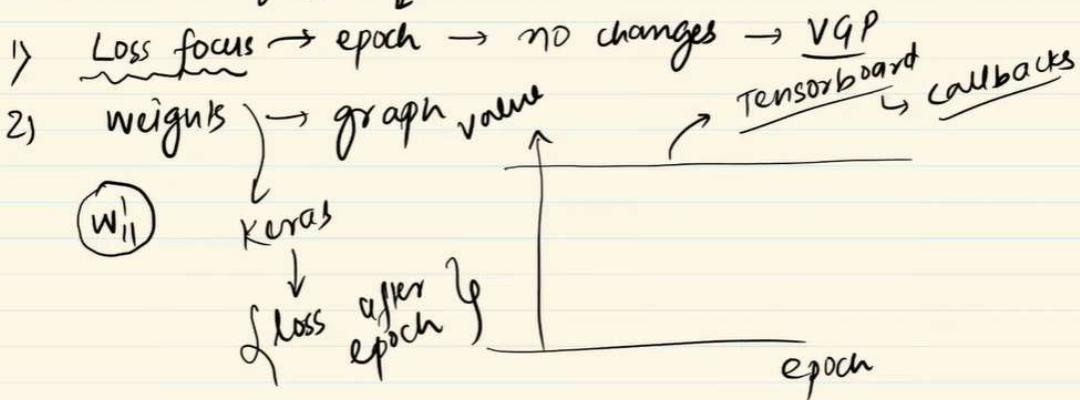
In machine learning, the **vanishing gradient problem** is encountered when training artificial neural networks with gradient-based learning methods and backpropagation. In such methods, during each iteration of training each of the neural network's weights receives an update proportional to the partial derivative of the error function with respect to the current weight. The problem is that in some cases, the gradient will be vanishingly small, effectively preventing the weight from changing its value. In the worst case, this may completely stop the neural network from further training. ~~An example of the problem cause traditional solution~~

1)  $0.1 \times 0.1 \times 0.1 \times 0.1 = [0.0001] \rightarrow \text{VGP}$

2) Deep NN  $\rightarrow \square \square \square \square \square$

3) Sigmoid / tanh  $\rightarrow Af \rightarrow$

How to recognize?



```
[75] old_weights
```

```
array([[ 0.15133941, -0.5475049 ,  0.20439553, -0.5383658 ,  0.01435351,
       0.6393233 , -0.4457206 , -0.7024723 ,  0.7043677 ,  0.26854134],
       [ 0.19128346,  0.2853908 , -0.42532492,  0.23361963,  0.48123974,
      -0.3359478 , -0.44716895, -0.57663697,  0.2556631 , -0.6959247 ]],  
dtype=float32)
```

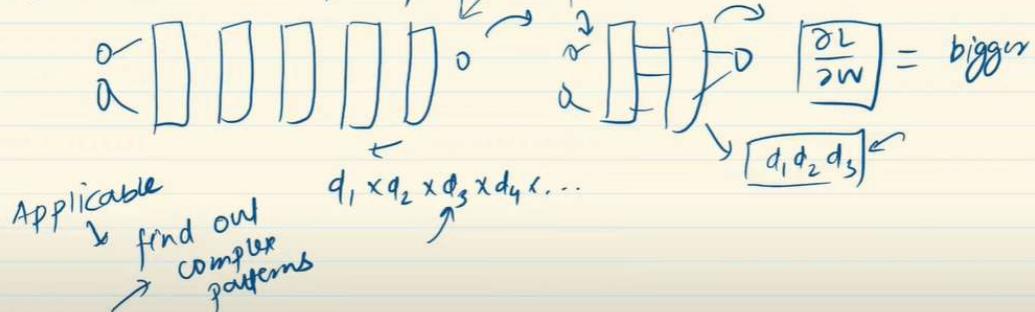
```
new_weights
```

```
array([[ 0.15132847, -0.5475164 ,  0.20438164, -0.5383458 ,  0.01440594,
       0.639327 , -0.4457354 , -0.7025074 ,  0.7043448 ,  0.26854628],
       [ 0.19130446,  0.28541353, -0.42529657,  0.23357487,  0.48114753,
      -0.3359549 , -0.44714212, -0.5765698 ,  0.25571352, -0.6959336 ]],  
dtype=float32)
```

we can see after epochs there is no change in weights --- almost negligible

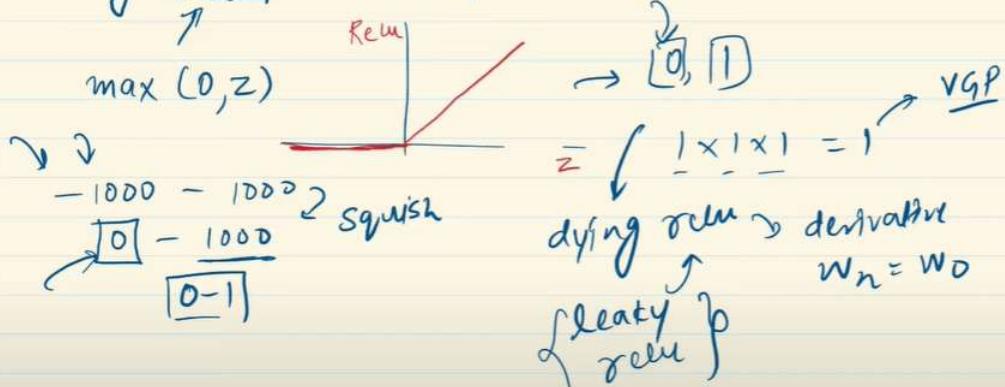
How to handle Vanishing Gradient Problem →

1) Reduce model complexity ↗



not applicable bcz nn is used for complex patterns and reducing layers would be wrong

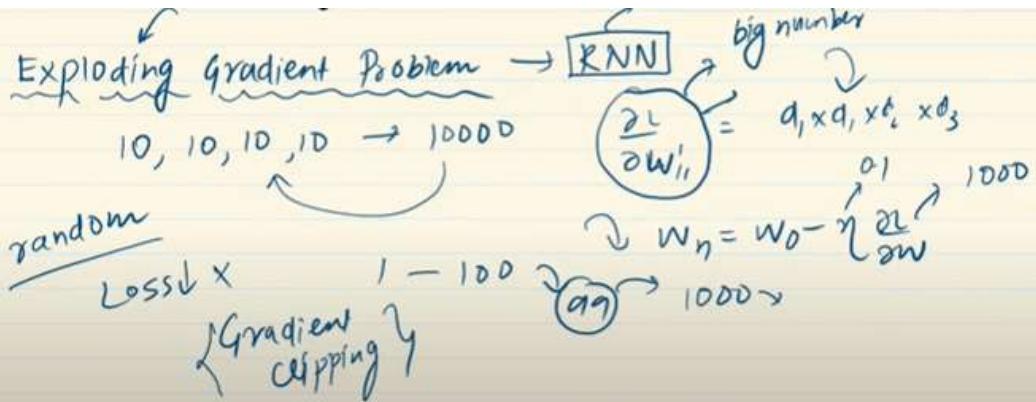
2) Using ReLU Activation functions



3) Proper weight init ↗ Glorot ↗ Xavier

4) Batch norm → layer →

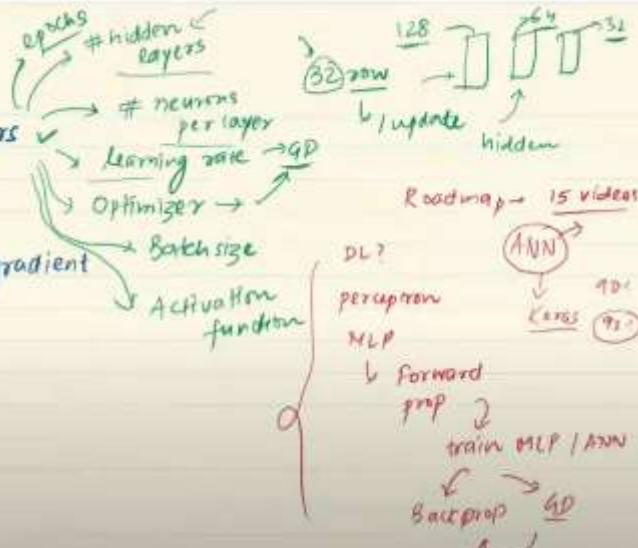
5) Residual Network ↗ CNN → ResNET  
↳ building block → ANN



### How to improve a neural network

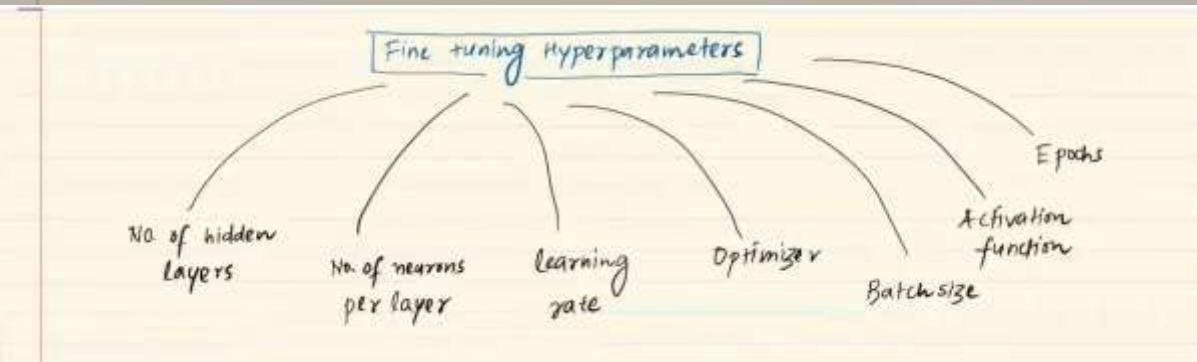
29 April 2022 18:51

#### 1. Fine tuning NN hyperparameters



#### 2. By solving problems:

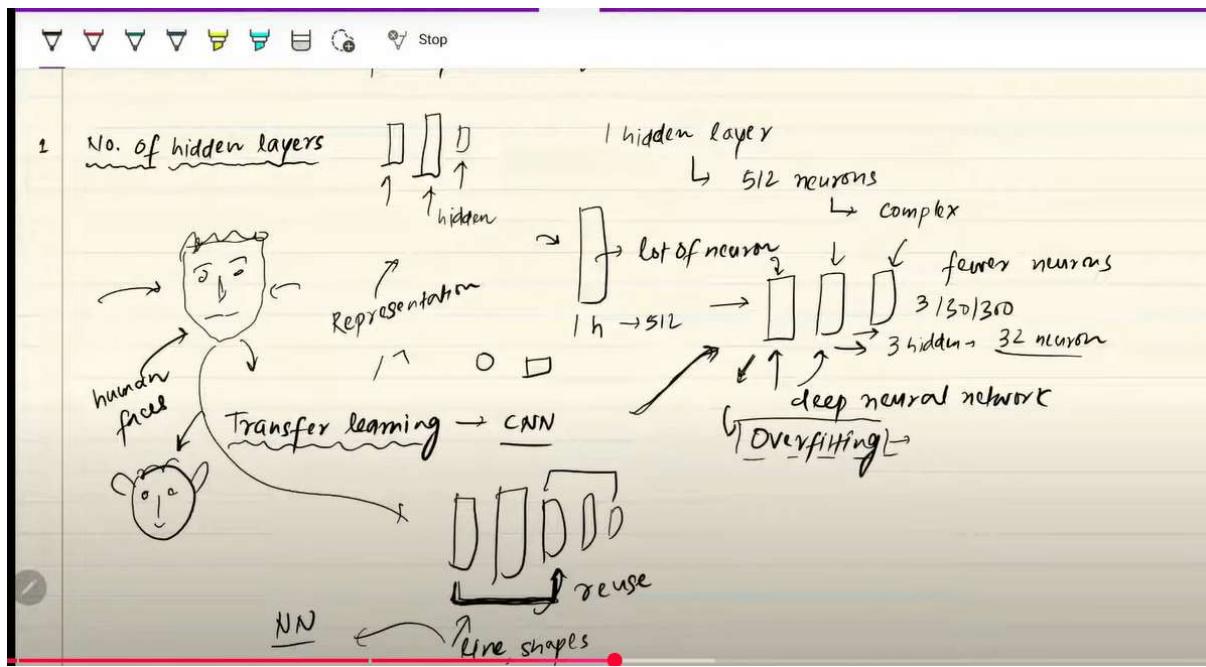
- Vanishing / Exploding gradient
- Not enough data
- Slow training
- Overfitting.



#### 1. no of hidden layers

prefer increasing layers rather than adding more neurons in one layer

keep increasing layers till overfitting error does not occur

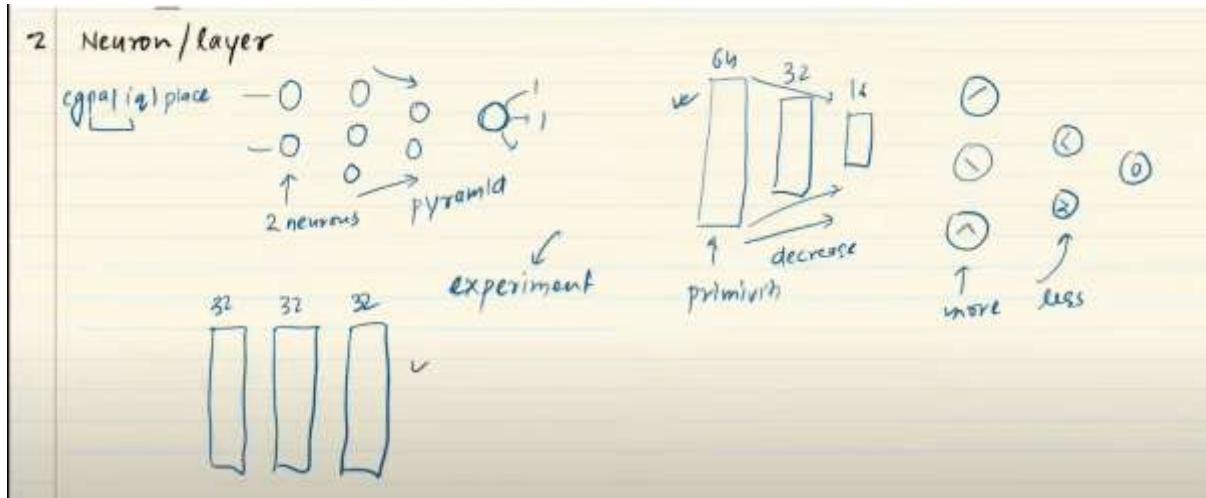


## 2. NO of neurons

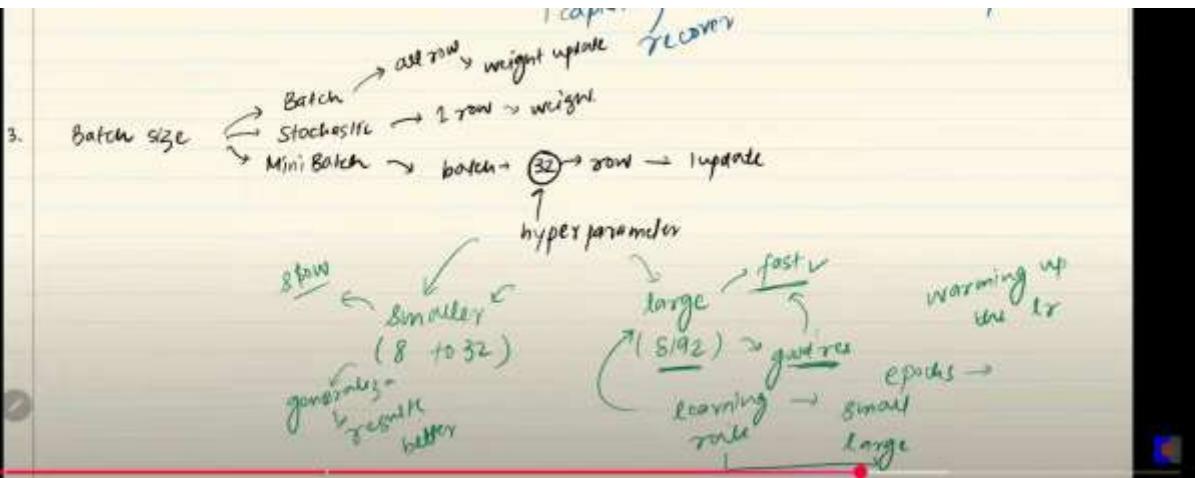
input layer => no. of neurons = no. of inputs

output layer => type of output

- regression - 1
- classification - 1
- multiclass classification - depends on no. of classes to be guessed



## 3. batch size



#### 4. epochs



#### Problems with Neural Networks

Vanishing and exploding gradients

Not enough data

Slow training

Overfitting

Vanishing and exploding gradients

- Weight Init →
- Activation function →
- Batch Norm →
- Gradient clipping →

Not enough data

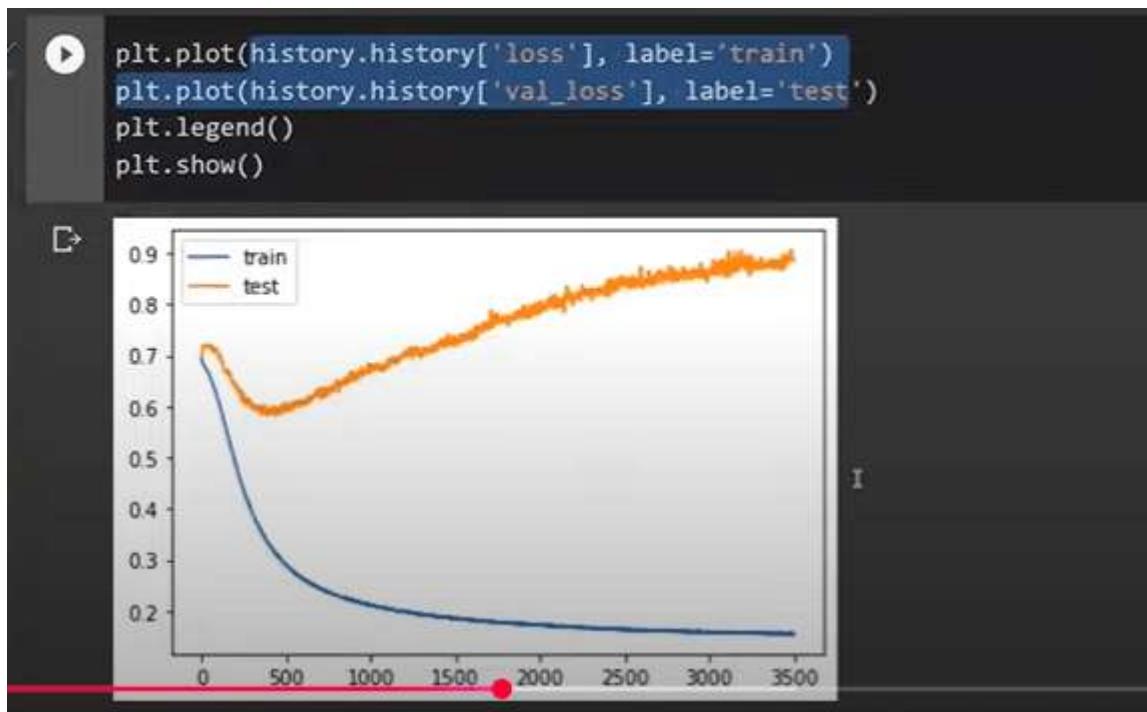
- Transfer learning →
- Unsupervised pretraining →

Slow training

- Optimizers → Adam
- learning rate scheduler

Overfitting

- $L_1$  and  $L_2$  reg
- Dropouts



◆ **Blue Curve – `train` (Training Loss)**

- Plotted using `history.history['loss']`
- Shows how well the model is fitting the **training data**
- **Decreases over time**, which is expected as the model learns

◆ **Orange Curve – `test` (Validation Loss)**

- Plotted using `history.history['val_loss']`
- Shows how well the model performs on **unseen validation data**
- **Initially decreases**, then **starts increasing** — this is a classic sign of **overfitting**

[keras.io/api/callbacks/early\\_stopping/](https://keras.io/api/callbacks/early_stopping/)

Commands | Code | Text | Run All | Copy to Drive

### Early Stopping

```
[24] model = Sequential()
model.add(Dense(256, input_dim=2, activation='relu'))
model.add(Dense(1, activation='sigmoid'))

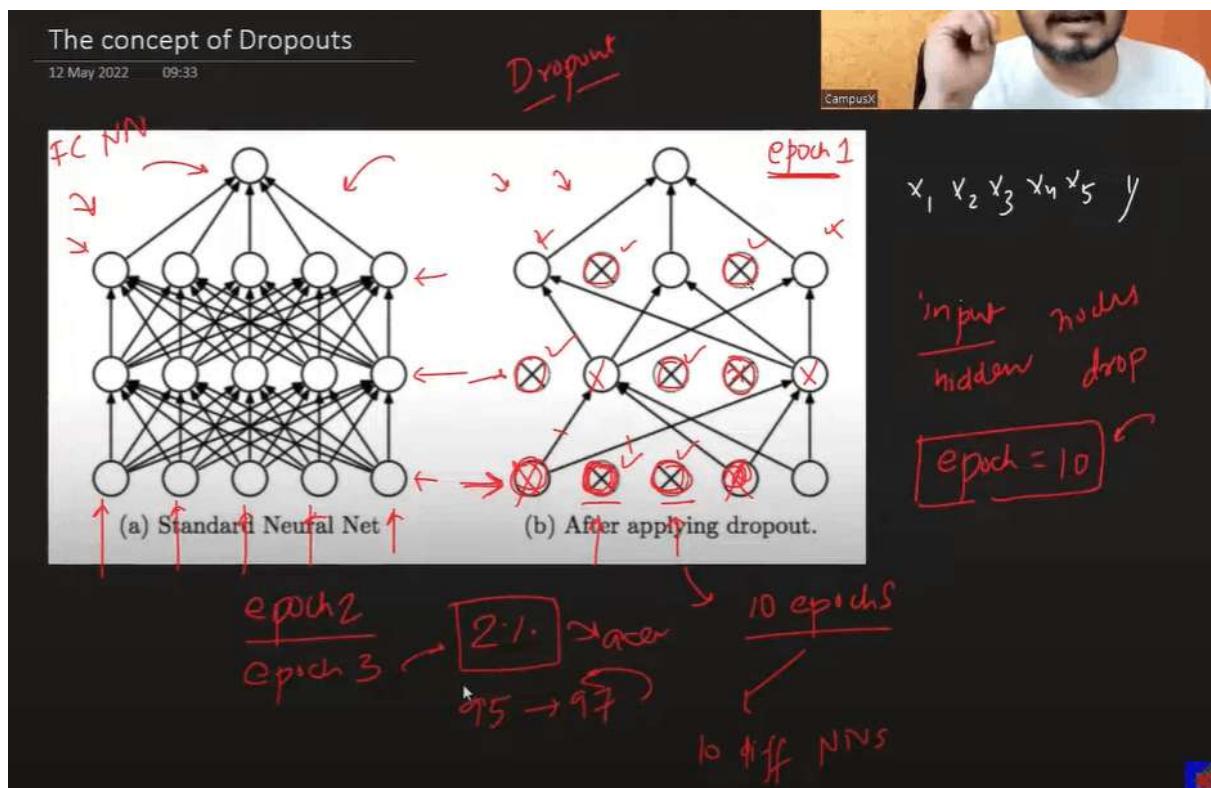
/usr/local/lib/python3.11/dist-packages/keras/src/layers/core/dense.py:87: UserWarning: Do not pass an `input_shape`/`input_dim` super().__init__(activity_regularizer, **kwargs)

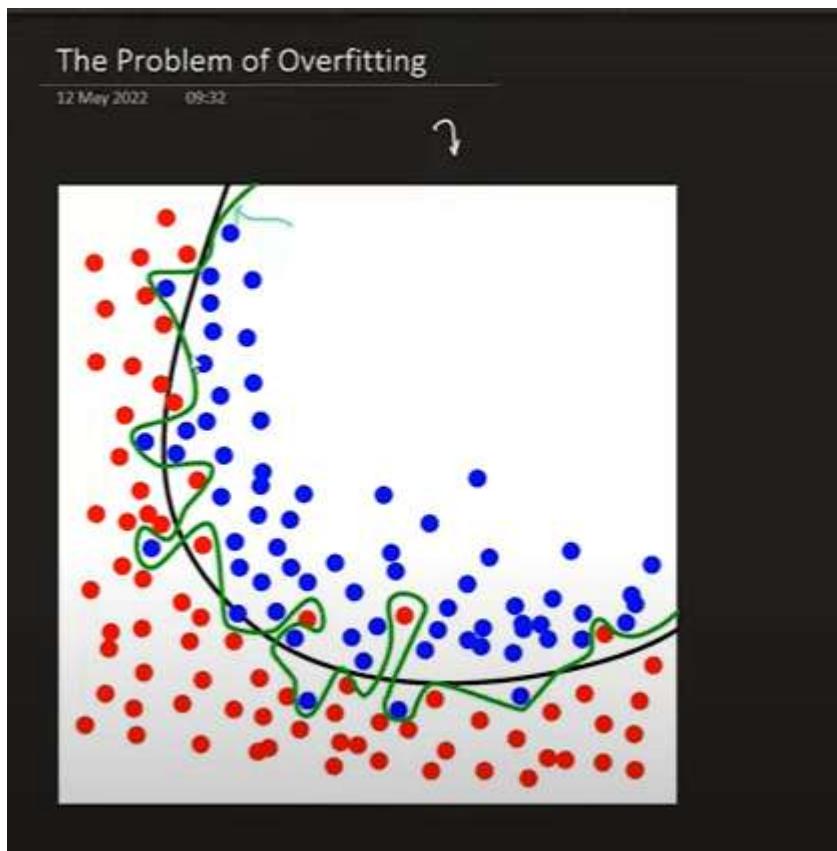
[25] model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])

callback = EarlyStopping(
    monitor="val_loss", #Quantity to be monitored. Defaults to "val_loss".
    min_delta=0.0001, #Minimum change in monitored metric to qualify as an actual improvement.
    patience=20, #Number of epochs with no improvement before training is stopped.
    verbose=1, #Controls whether to display early stopping messages (0 = silent, 1 = show).
    mode="auto",#Defines improvement direction - 'min' for decreasing, 'max' for increasing, 'auto' infers from metric name.
    baseline=None,# Minimum baseline value; training stops if no improvement beyond it is achieved
    restore_best_weights=True #If True, restores model weights from the epoch with best monitored
    value.
)

[26] history = model.fit(X_train, y_train, validation_data=(X_test, y_test), epochs=3500,
callbacks=callback)
```

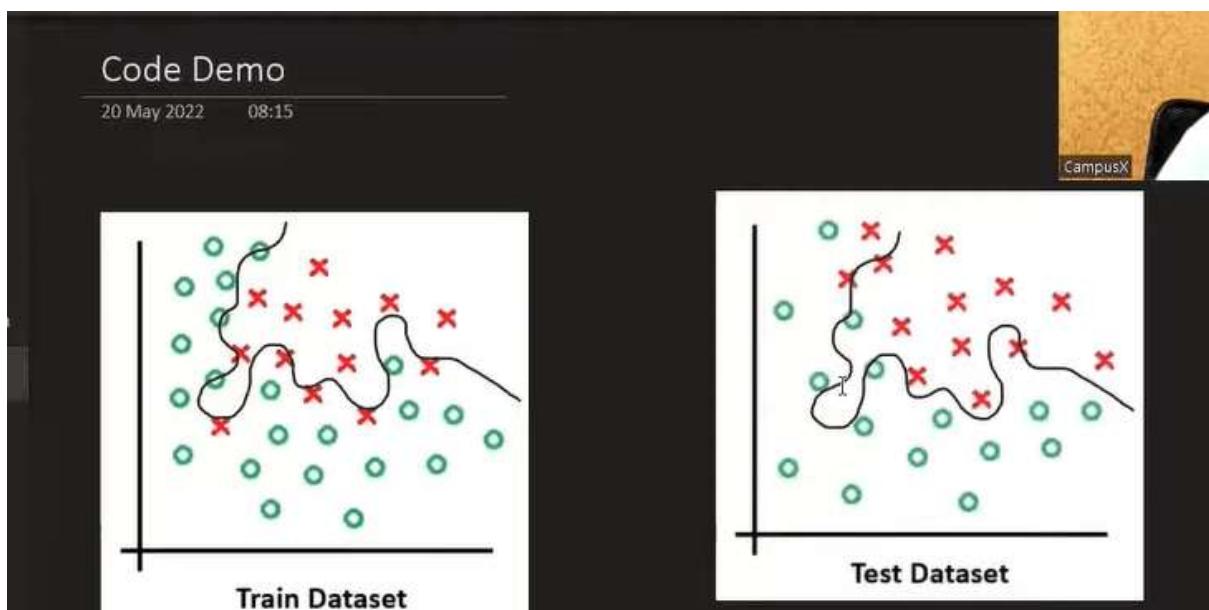
# OVERFITTING

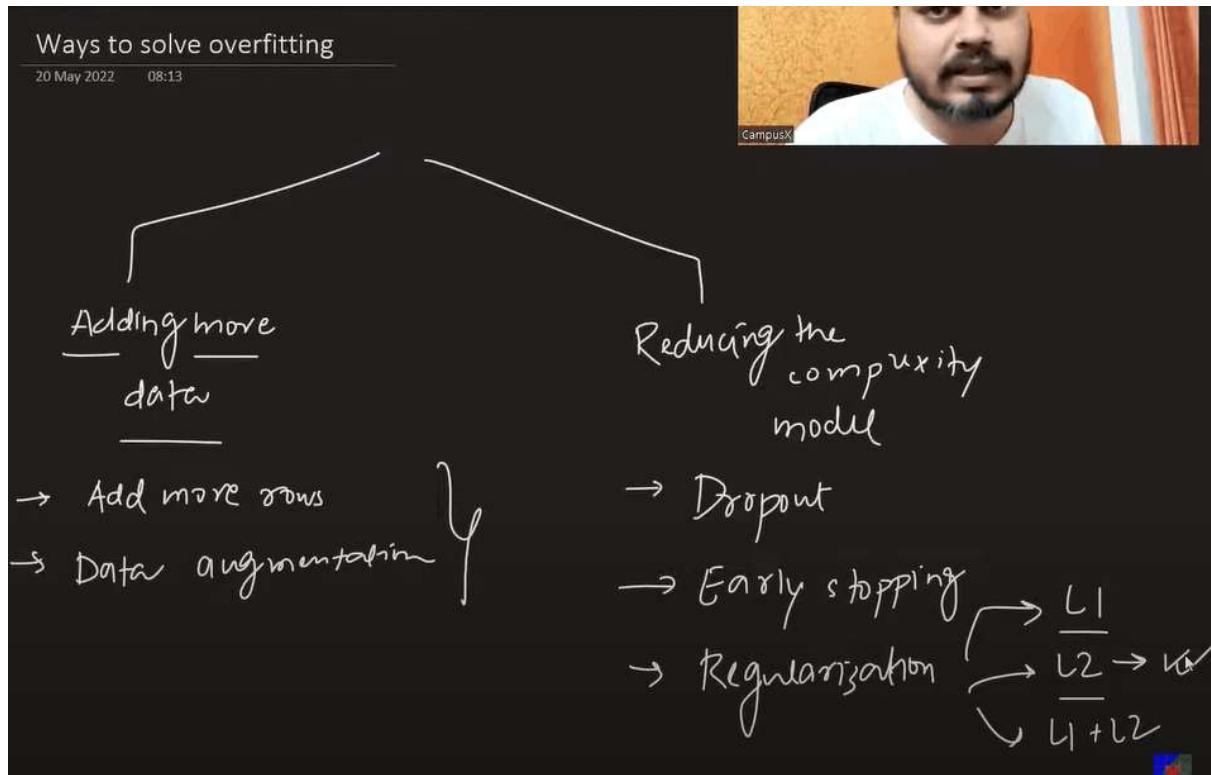




here green NN is trying to memorize and is very specific from training data perspective and performs bad with testing data

whereas the black NN is perfect

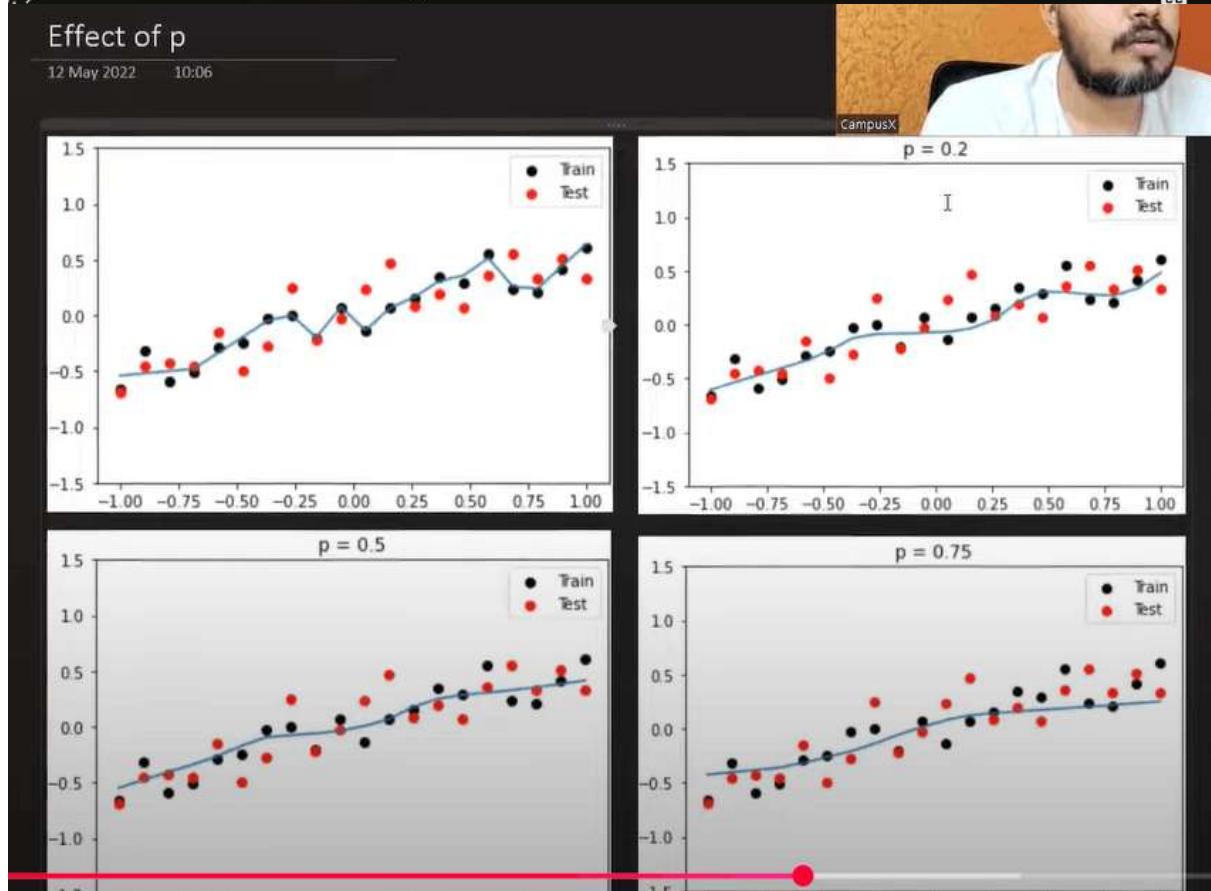
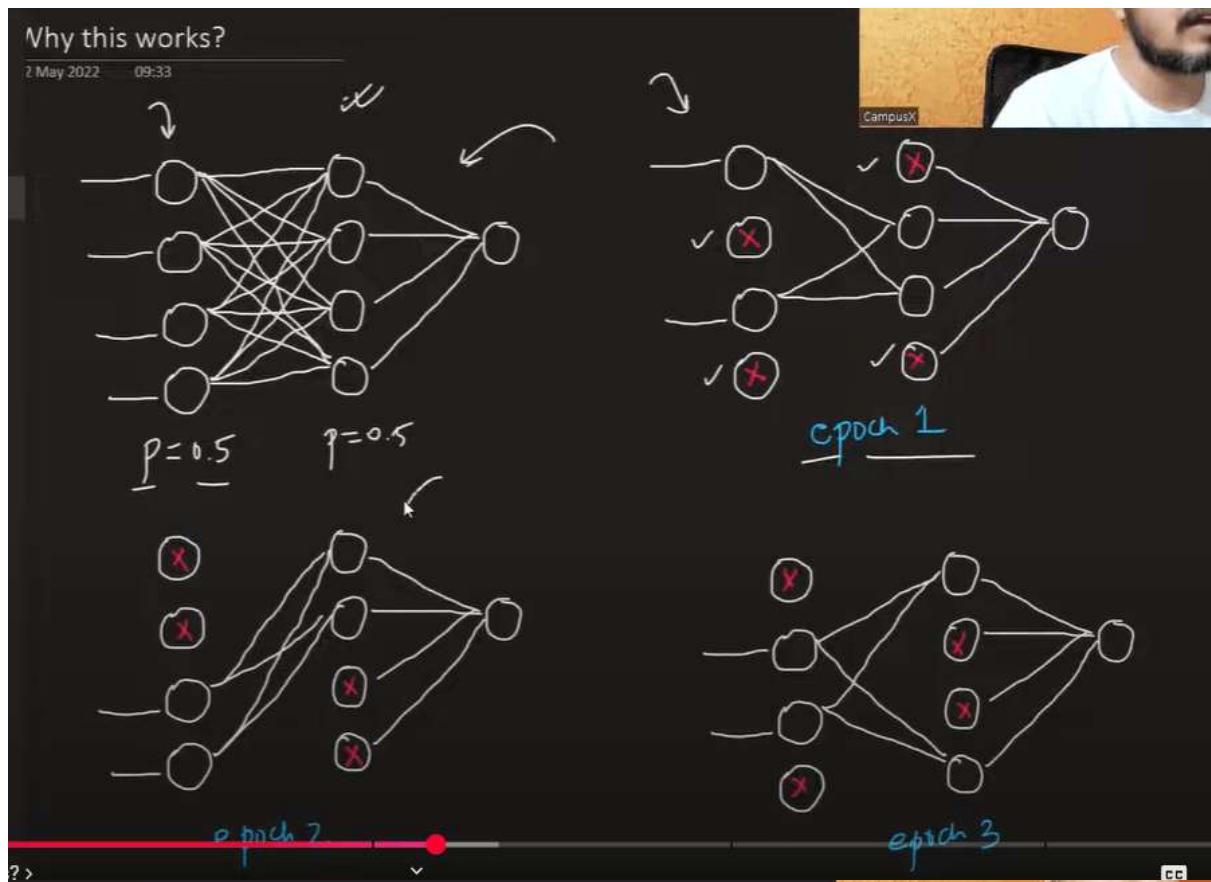


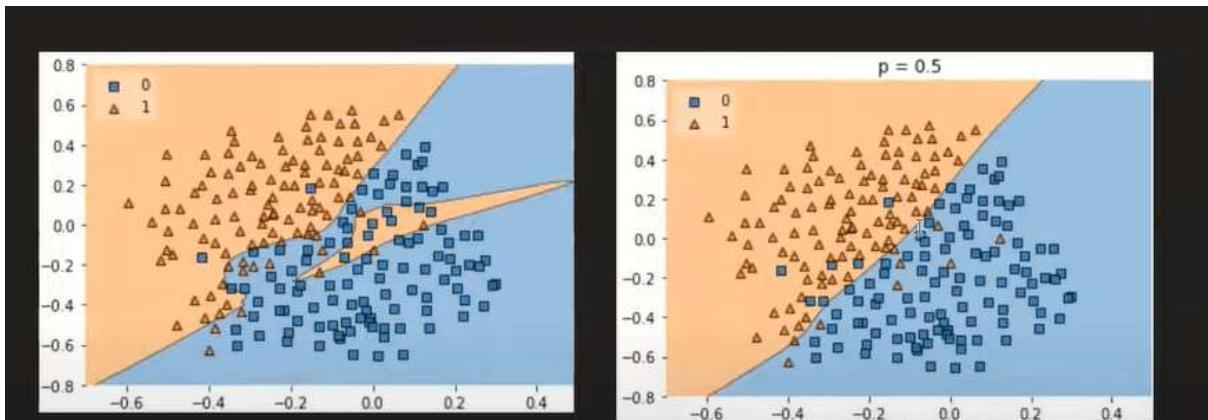


## DROPOUTS

10 epoch - 10 diff NN

randomly few neurons are shut down





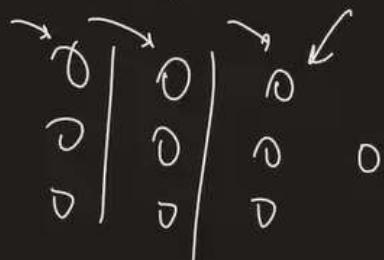
## Practical Tips and Tricks

12 May 2022 09:34



1) Overfitting  $P \uparrow$ , underfitting  $P \downarrow$

2) Last layer  $\rightarrow$  dropout



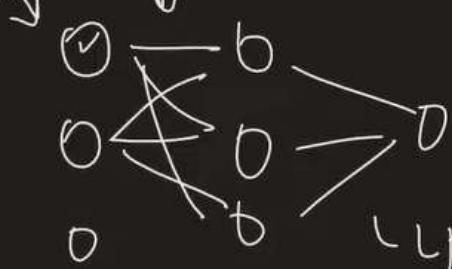
3) CNN  $\rightarrow$   $40-50 \cdot 1$  (P)  $\rightarrow \checkmark$

20-30 RNN

$\boxed{50 >}$

ANIN  $\rightarrow$   $\boxed{10-50}$



- 1) convergence  $\rightarrow$  delay
- 2) 
- LF - has to change

## L1 L2 regularization

$\hookrightarrow \min \text{ Lossfunction}$

$L = \text{mse}$   
 $\downarrow$  binary

$C = \left[ \frac{1}{n} \sum_{i=1}^n L(y_i, \hat{y}_i) \right] + \text{penalty term}$

$C = L + \left[ \frac{\lambda}{2n} \sum_{i=1}^k \|w_i\|^2 \right] \rightarrow \text{weightage}$

$w_i \rightarrow w_{i,0}$

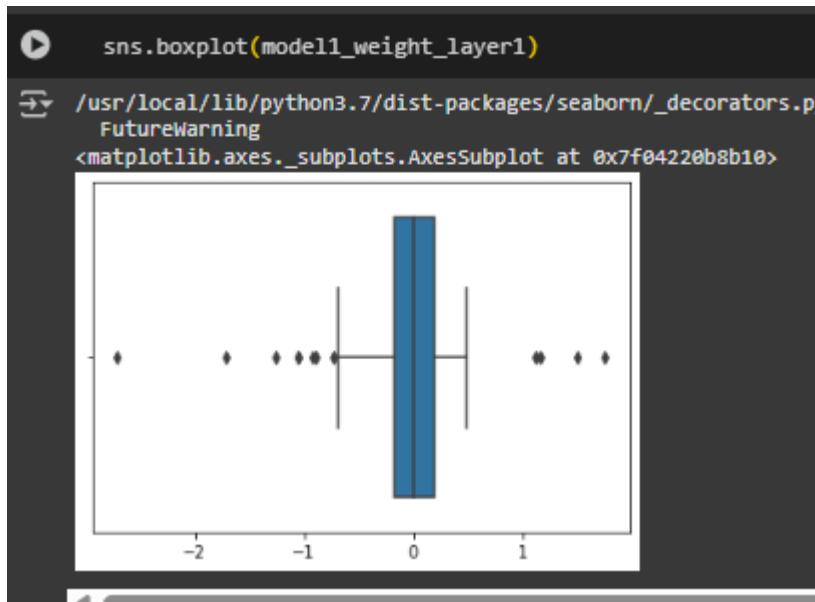
$\left( \frac{\lambda}{2n} \right) [ w_1^2 + w_2^2 + \dots + w_k^2 ]$

$\uparrow \lambda = \text{hyperparameter} \uparrow$

$\boxed{\lambda=0}$

$$\sum_{i=1}^n \left( y_i - \underbrace{\sum_{j=1}^p x_{ij} b_j}_{\text{Loss Term}} \right)^2 + \underbrace{\lambda \cdot \sum_{j=1}^p |b_j|}_{\text{Regularizer term}}$$

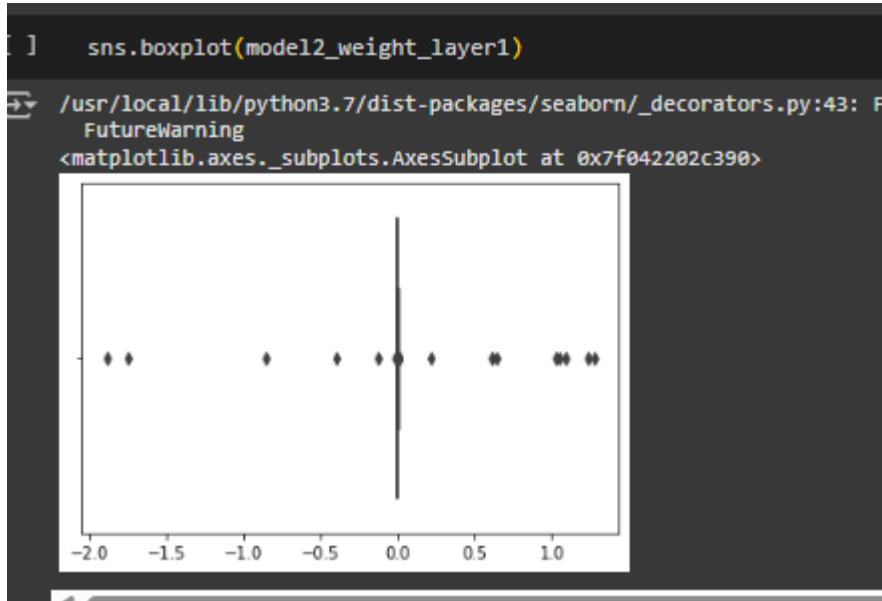
without regularization model 1



we can see weights are mainly between -0.2 and 0.2 with outliers till -2.5

whereas regularization will decay weights and bring them smaller

so in model2 all weights are in between 0 range and outliers go till -2.0 max



```

} [48] model1_weight_layer1.min()
-2.7185578

[49] model2_weight_layer1.min()
-0.57482177

```

## ACTIVATION FUNCTION

activation function is a gate which tell whether any input given to a neuron will pass or not

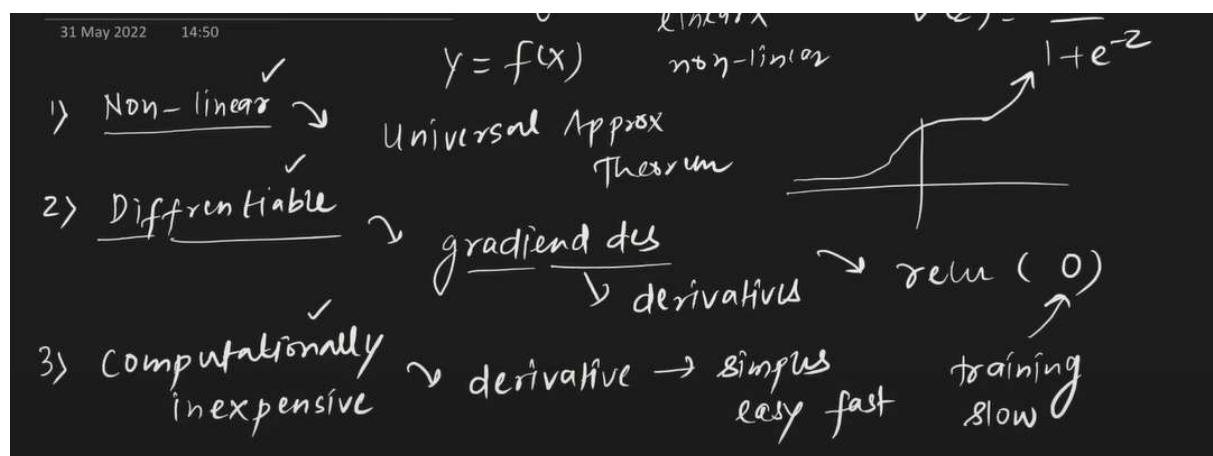
and if it passes

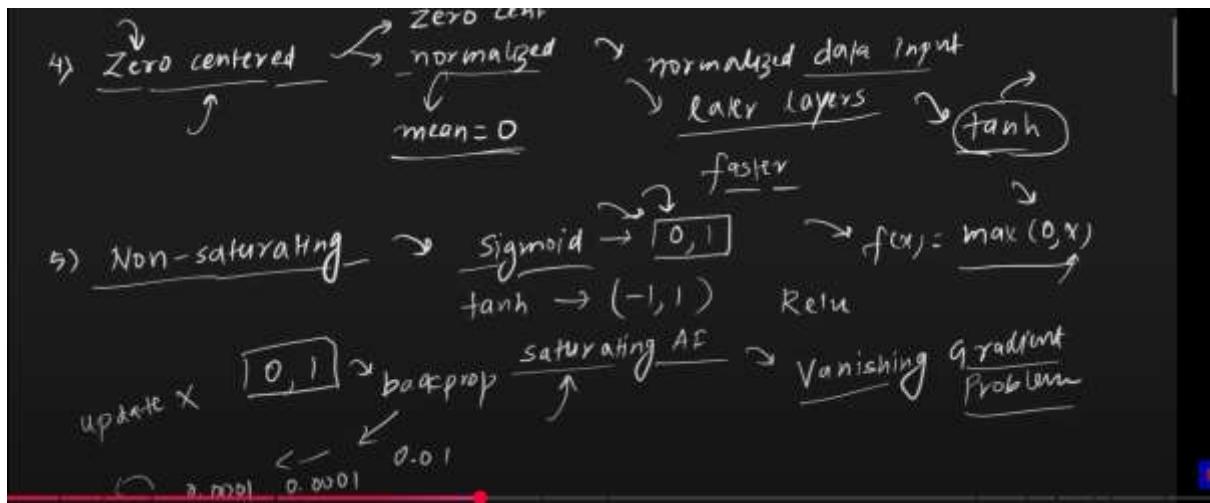
it decides

- whether a neuron will be activated or not
- if activated then how much activated will it be

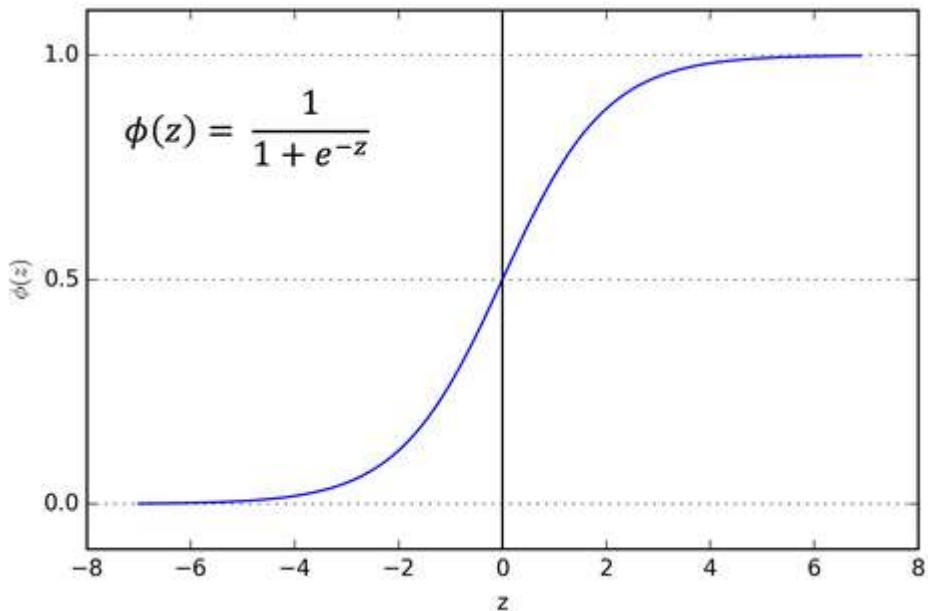
without activation function ANN will work like linear regg and will not catch no linearity

qualities of a activation func

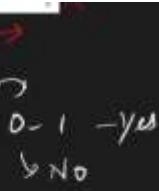




## sigmoid



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

Yes →  sigmoid  
 No → 

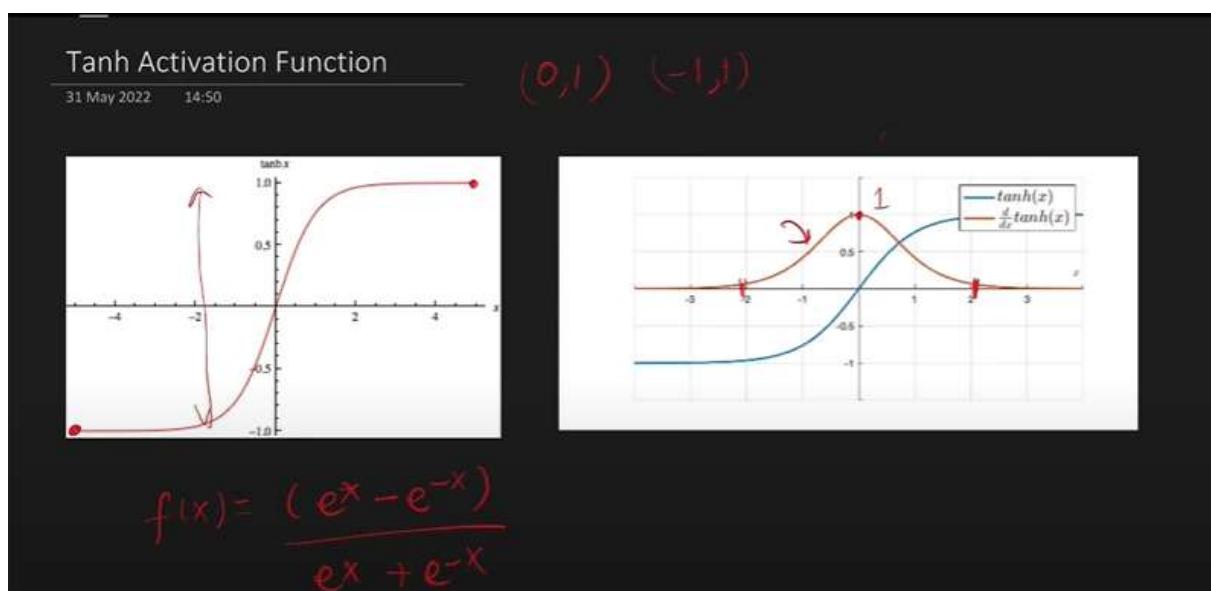
Advantages

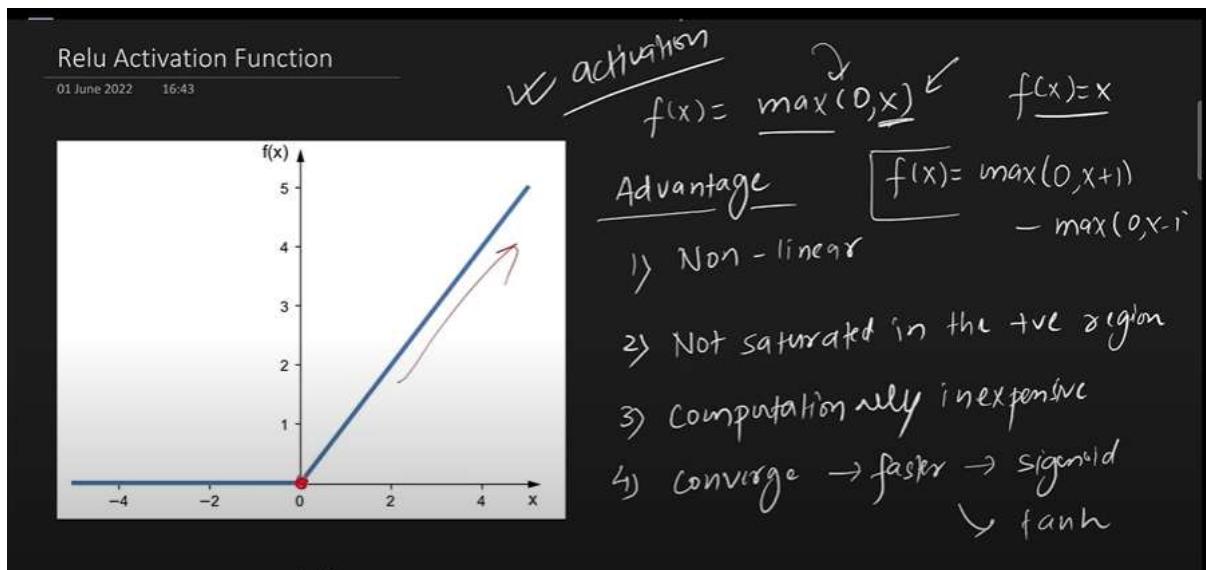
- 1)  $[0, 1] \rightarrow \text{probability} \rightarrow \text{output layer} \rightarrow \text{Binary classification}$
- 2) Non-Linear  $\rightarrow$   $\boxed{\text{Non-Linear data}} \rightarrow \text{good option}$
- 3) Differentiable  $\rightarrow$  Backprop  $\rightarrow \boxed{\frac{\partial L}{\partial w}}$

sigmoid is not used in hidden layers

used in output layer in case of binary classification problem

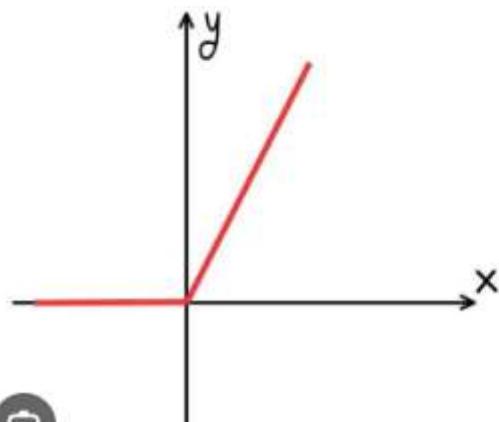
## Tanh





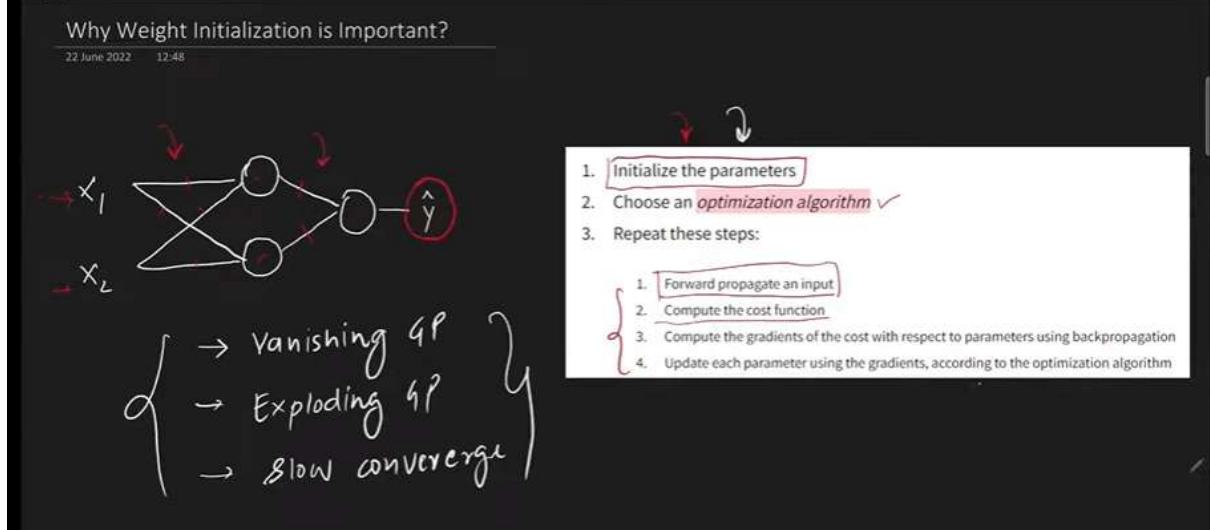
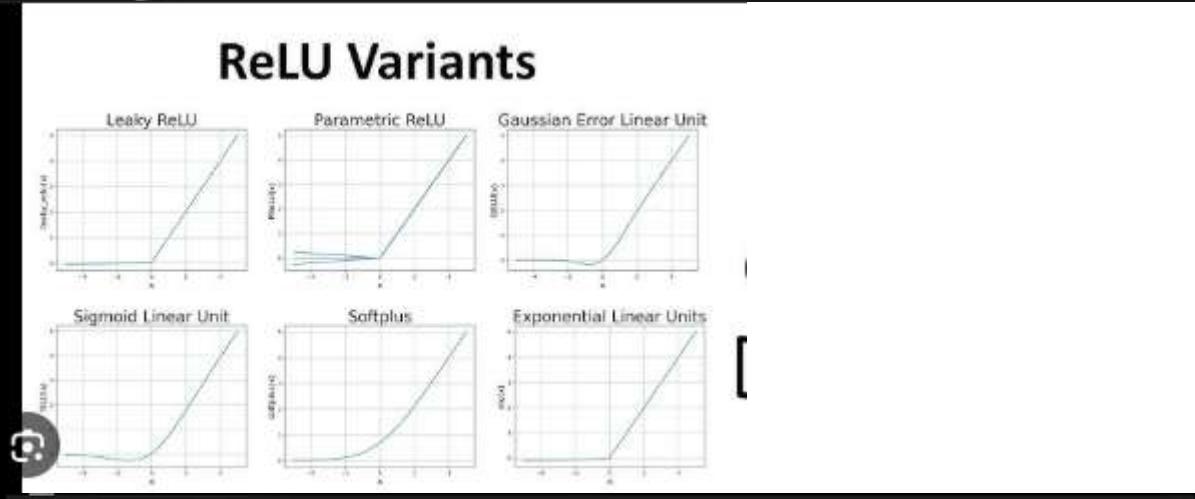
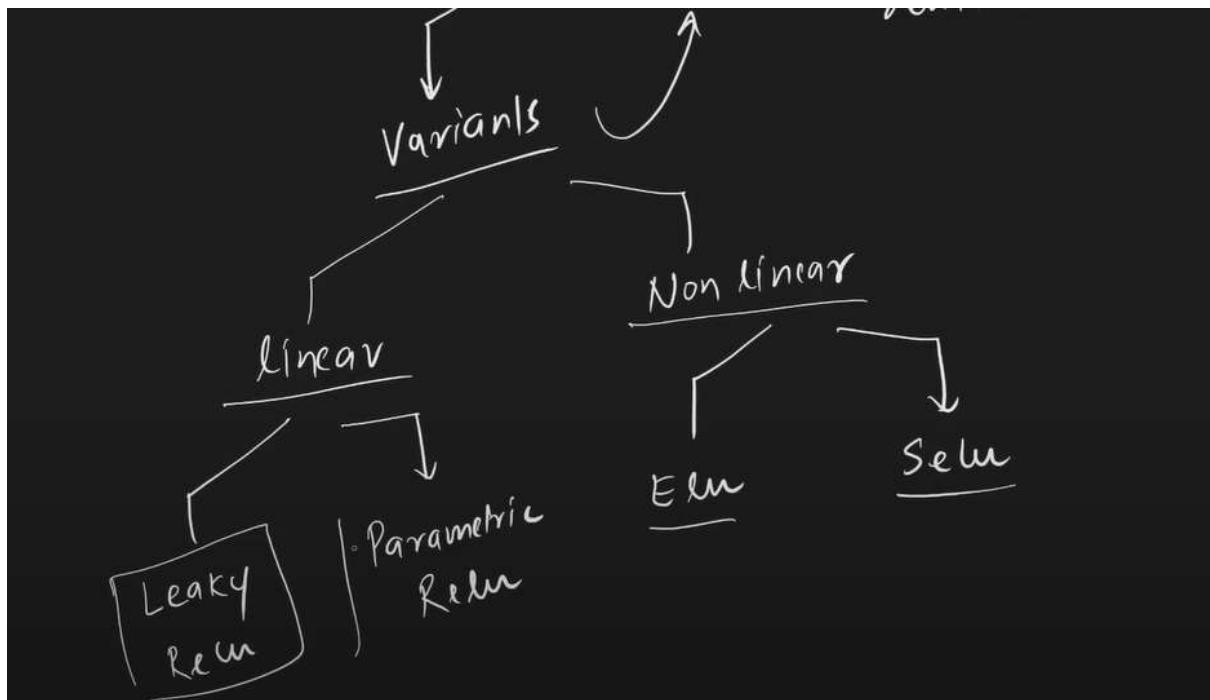
## ReLU Activation Function

CLEARLY EXPLAINED..



dying relu

- Solutions
- Set low learning rate
  - bias → +ve value →  $0.01$
  - Don't use relu → variants



As number of hidden layers grow, gradient becomes very small and weights will hardly change . This will hamper the learning process

## Vanishing Gradients

### WEIGHT INITIALIZATION TECHNIQUES

Why Weight Initialization is Important?  
22 June 2022 12:48

2000's Deep learning

1. Initialize the parameters  
2. Choose an *optimization algorithm* ✓  
3. Repeat these steps:  
1. Forward propagate an input  
2. Compute the cost function  
3. Compute the gradients of the cost with respect to parameters using backpropagation  
4. Update each parameter using the gradients, according to the optimization algorithm

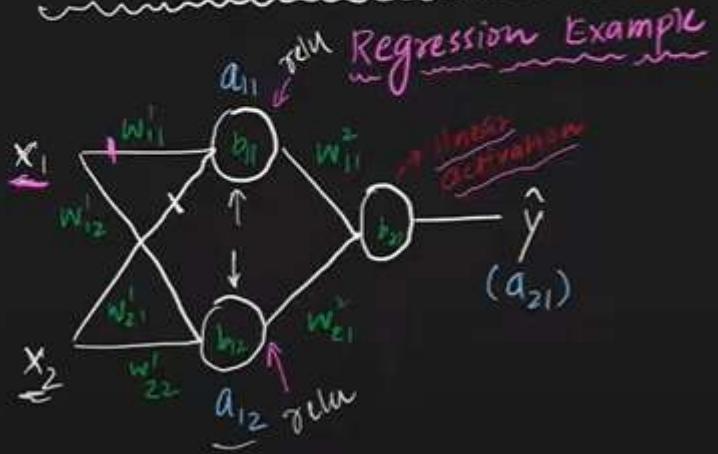
Sigmoid Activ

wrong weight init

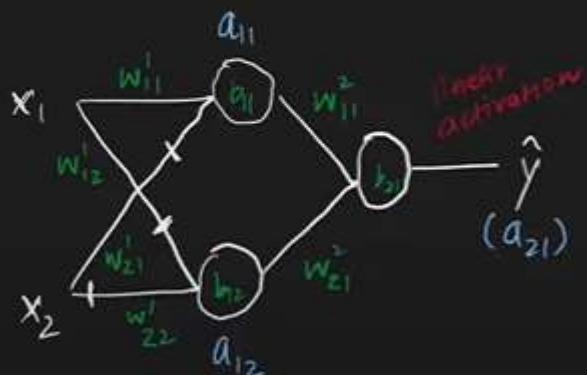
What not to do?

22 June 2022 12:49

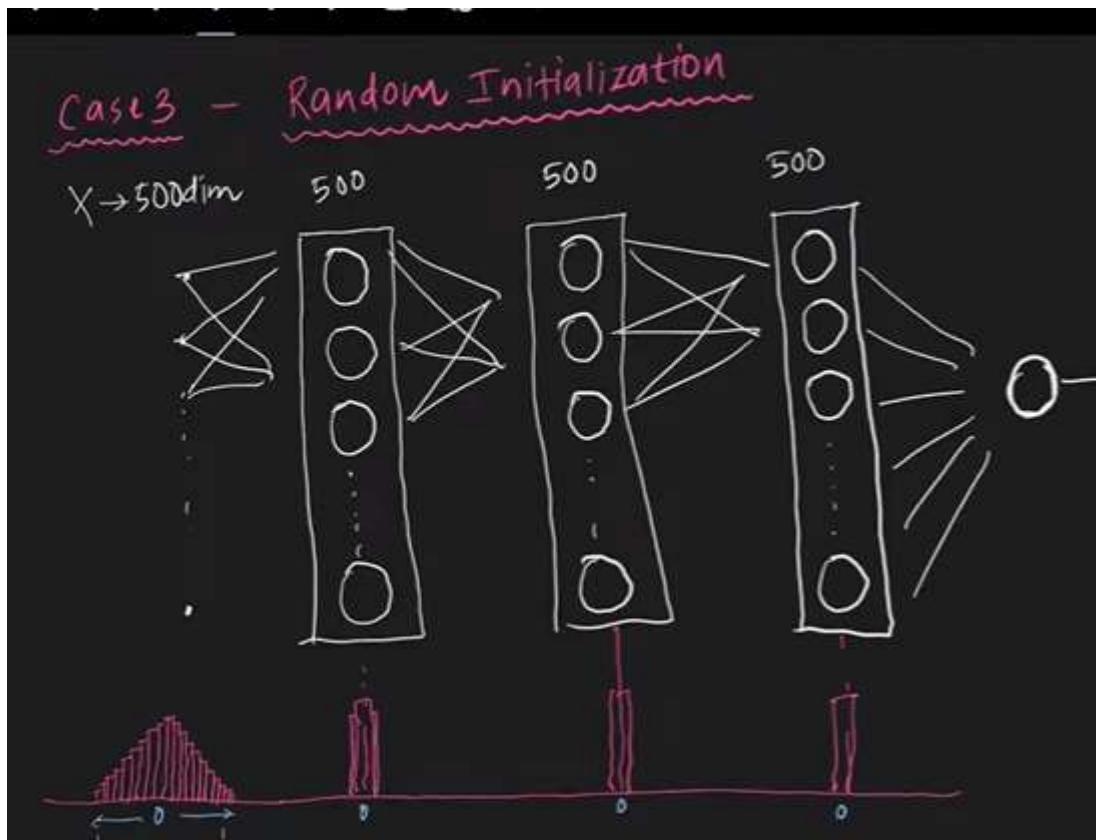
Case1 → Zero Initialization



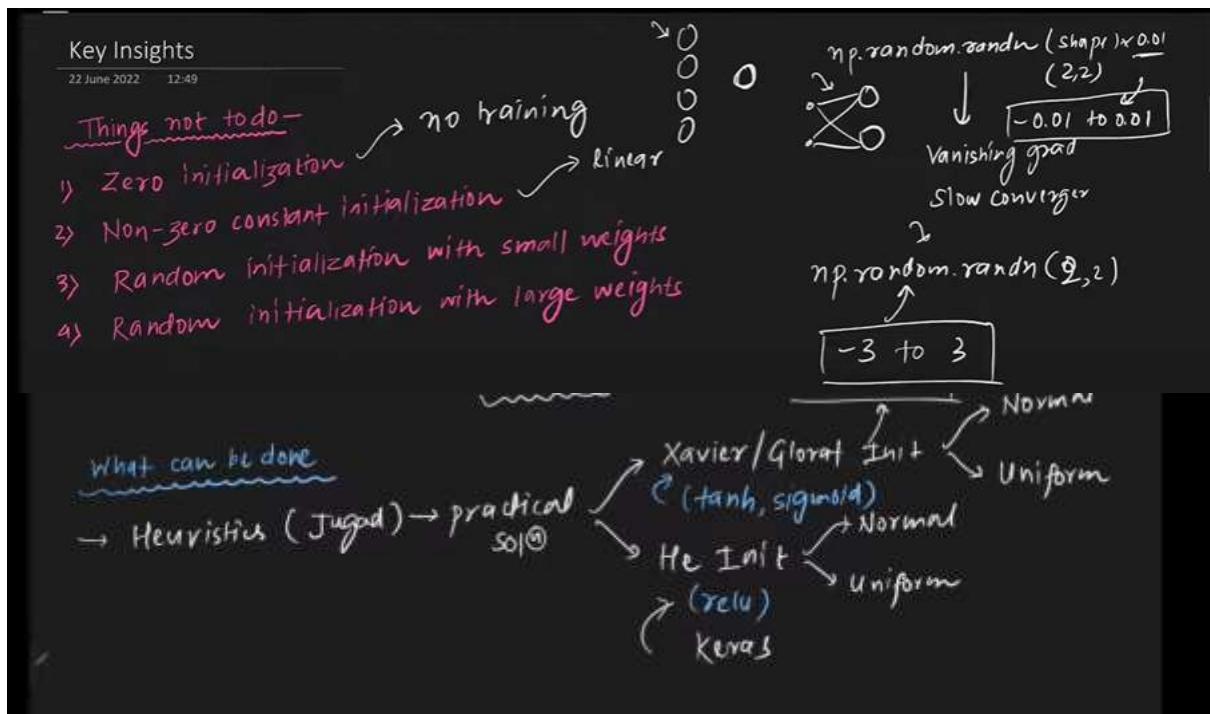
Case2 → Non-0 constant value



all weights and biases given same value



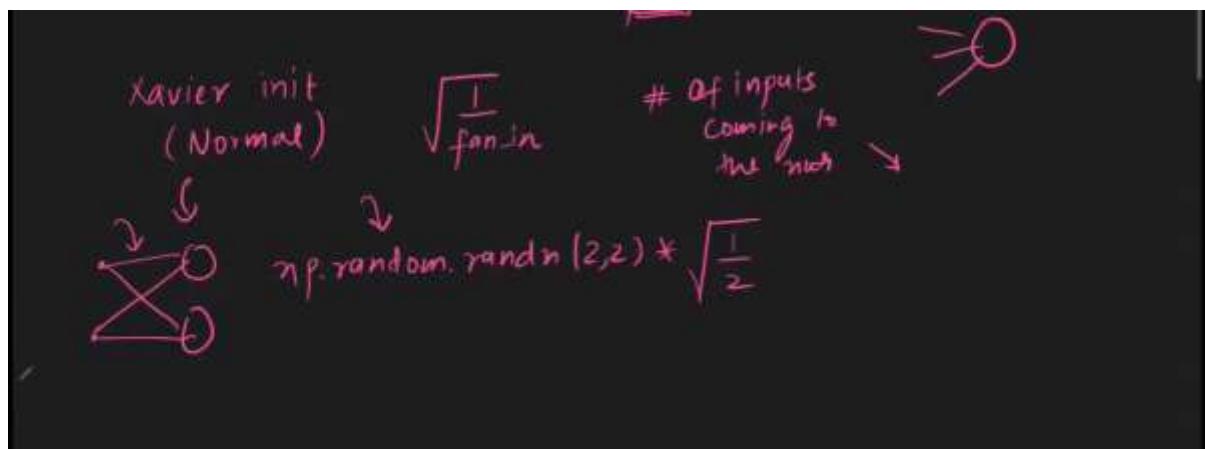
all below 4 techniques are wrong



we know weights can only be initialized by random values

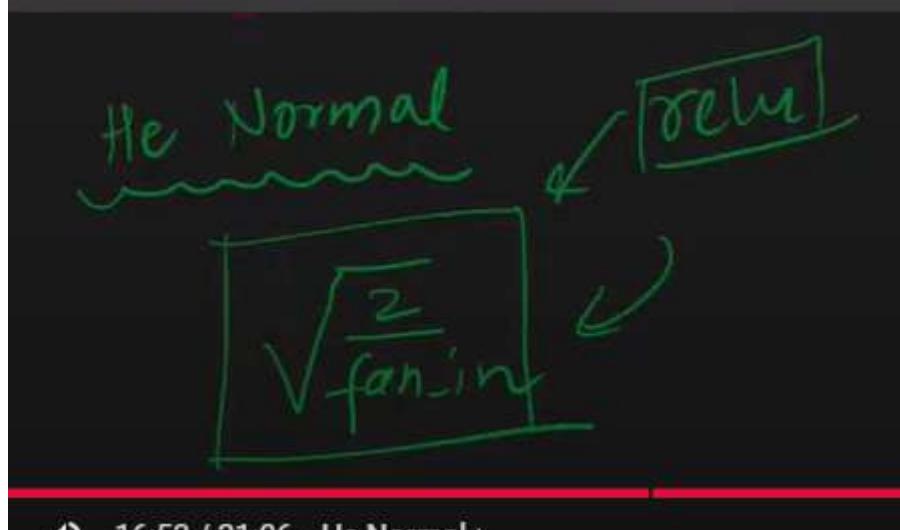
but problem is that they cant be very small values nor it can be very big values

so we should have some appropriate range of values from which random w and b are picked



```
initial_weights[0] = np.random.randn(2,10)*np.sqrt(1/2)
initial_weights[1] = np.zeros(model.get_weights()[1].shape)
initial_weights[2] = np.random.randn(10,10)*np.sqrt(1/10)
initial_weights[3] = np.zeros(model.get_weights()[3].shape)
initial_weights[4] = np.random.randn(10,10)*np.sqrt(1/10)
initial_weights[5] = np.zeros(model.get_weights()[5].shape)
initial_weights[6] = np.random.randn(10,10)*np.sqrt(1/10)
initial_weights[7] = np.zeros(model.get_weights()[7].shape)
initial_weights[8] = np.random.randn(10,1)*np.sqrt(1/10)
initial_weights[9] = np.zeros(model.get_weights()[9].shape)
```

[48] model.set\_weights(initial\_weights)



16:52 / 21:06 • He Normal

## Uniform Distribution

Xavier Uniform

$[-\text{limit}, \text{limit}]$



$$\text{limit} = \sqrt{\frac{6}{(\text{fan-in} + \text{fan-out})}}$$

## He Uniform

$[-\text{limit}, \text{limit}]$

$$\text{limit} = \sqrt{\frac{6}{\text{fan-in}}}$$

`class GlorotNormal`: The Glorot normal initializer, also called Xavier normal initializer.

`class GlorotUniform`: The Glorot uniform initializer, also called Xavier uniform initializer.

`class HeNormal`: He normal initializer.

`class HeUniform`: He uniform variance scaling initializer.

```
model = Sequential()

model.add(Dense(10, activation='relu', input_dim=2, kernel_initializer='he_uniform'))
model.add(Dense(10, activation='relu', kernel_initializer='he_normal'))
model.add(Dense(10, activation='relu', kernel_initializer='he_normal'))
model.add(Dense(10, activation='relu', kernel_initializer='he_normal'))
model.add(Dense(1, activation='sigmoid'))
```

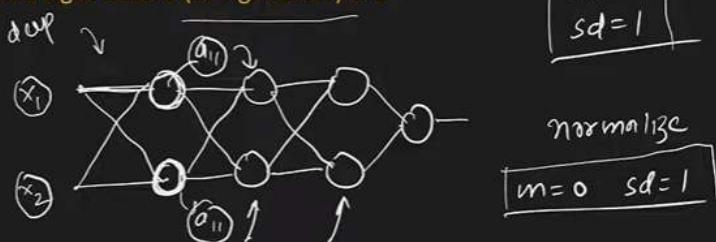
model.summary()

## What is Batch Norm?

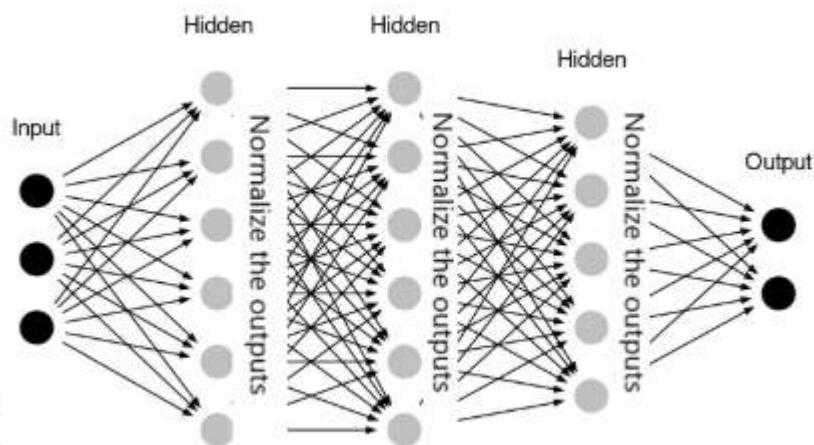
27 June 2022 11:00

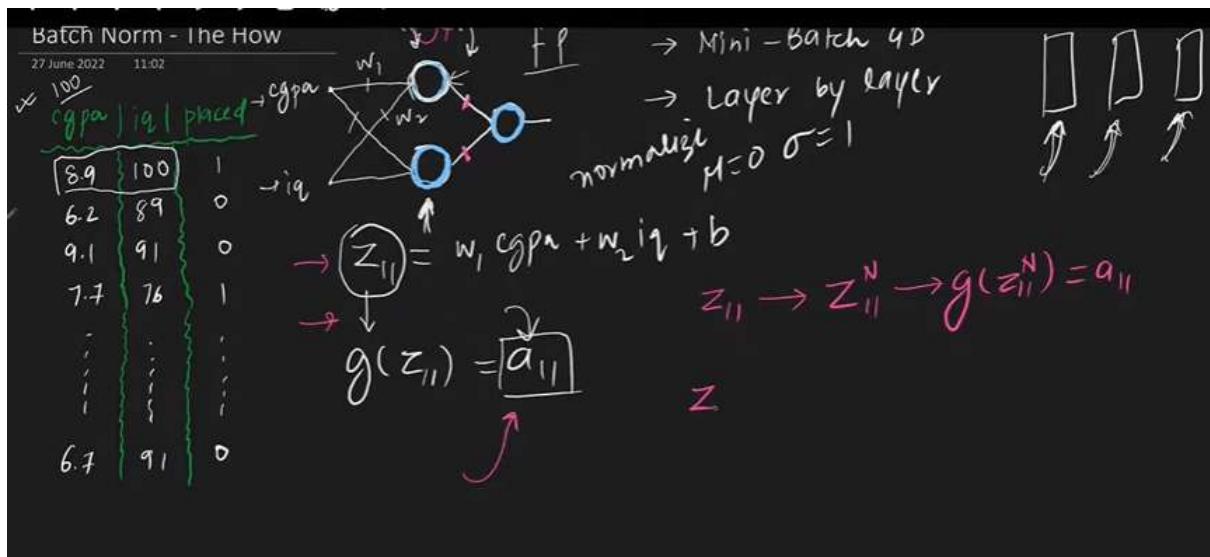
Batch-Normalization (BN) is an algorithmic method which makes the training of Deep Neural Networks (DNN) faster and more stable.

- It consists of normalizing activation vectors from hidden layers using the mean and variance of the current batch. This normalization step is applied right before (or right after) the nonlinear function.



Analytics India Magazine





$$\frac{z_{11} - \mu}{\sigma} = z_{11}^N$$

- 27 June 2022 11:02
- 1) stable → hyperpara → wider range of values
  - 2) faster → learning rate (higher)
  - 3) Regularizer → batch randomness → mix overfitting
    - ↓
    - Dropout
  - 4) weight init impact reduce
- 

```

model = Sequential()

model.add(Dense(3,activation='relu',input_dim=2))
model.add(BatchNormalization())
model.add(Dense(2,activation='relu'))
model.add(BatchNormalization())
model.add(Dense(1,activation='sigmoid'))

model.summary()

```

# OPTIMIZERS

10:35 05.7.22 10:32

D27 Optimizers HOME INSERT DRAW VIEW

Types of Optimizers

05 July 2022 10:02

$\begin{cases} \rightarrow \text{Batch GD} \\ \rightarrow \text{Stochastic GD} \\ \rightarrow \text{Mini batch GD} \end{cases}$

$(10)$

$\# \text{ epochs} = (10)$

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

500 rows  $\rightarrow$  pred

loss  $\rightarrow$  weight update

$(10 \times 500) \rightarrow 5000$

batch size = 100

5 batch  $\rightarrow$   $10 \times 5$  hours

GRADIENT DESCENT IN NEURAL NETWORKS

Gradient Descent in Neural Networks | Batch vs...

1.6K views • 2 months ago • 37:53

challengers with GD

Challenges

05 July 2022 10:02

1) learning rate

2) learning rate scheduling  $\rightarrow$  pre define

3)

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

$(\eta)$  same both the directions

$10-D \rightarrow 10 \text{ directions}$

weights and bias learn

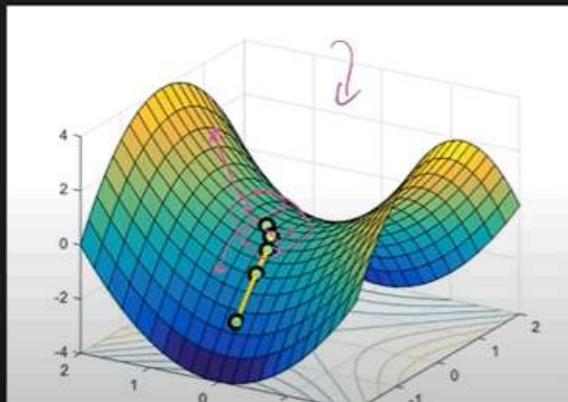
complex minima

multiple sub optimum

us local minima

5) saddle point

$$\frac{\partial L}{\partial w} = 0$$



What next?

05 July 2022 10:02

1) Momentum

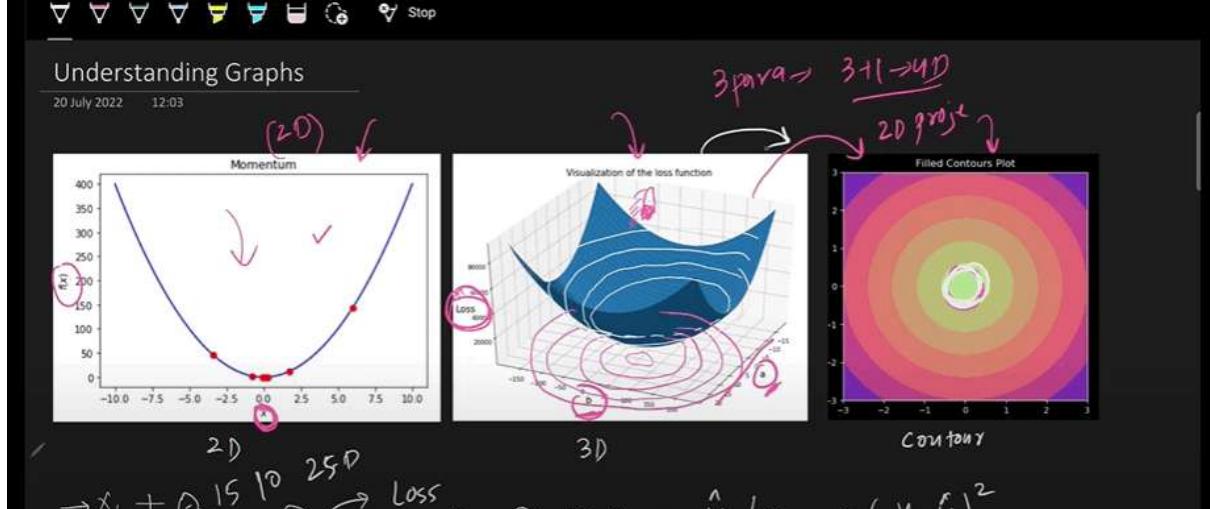
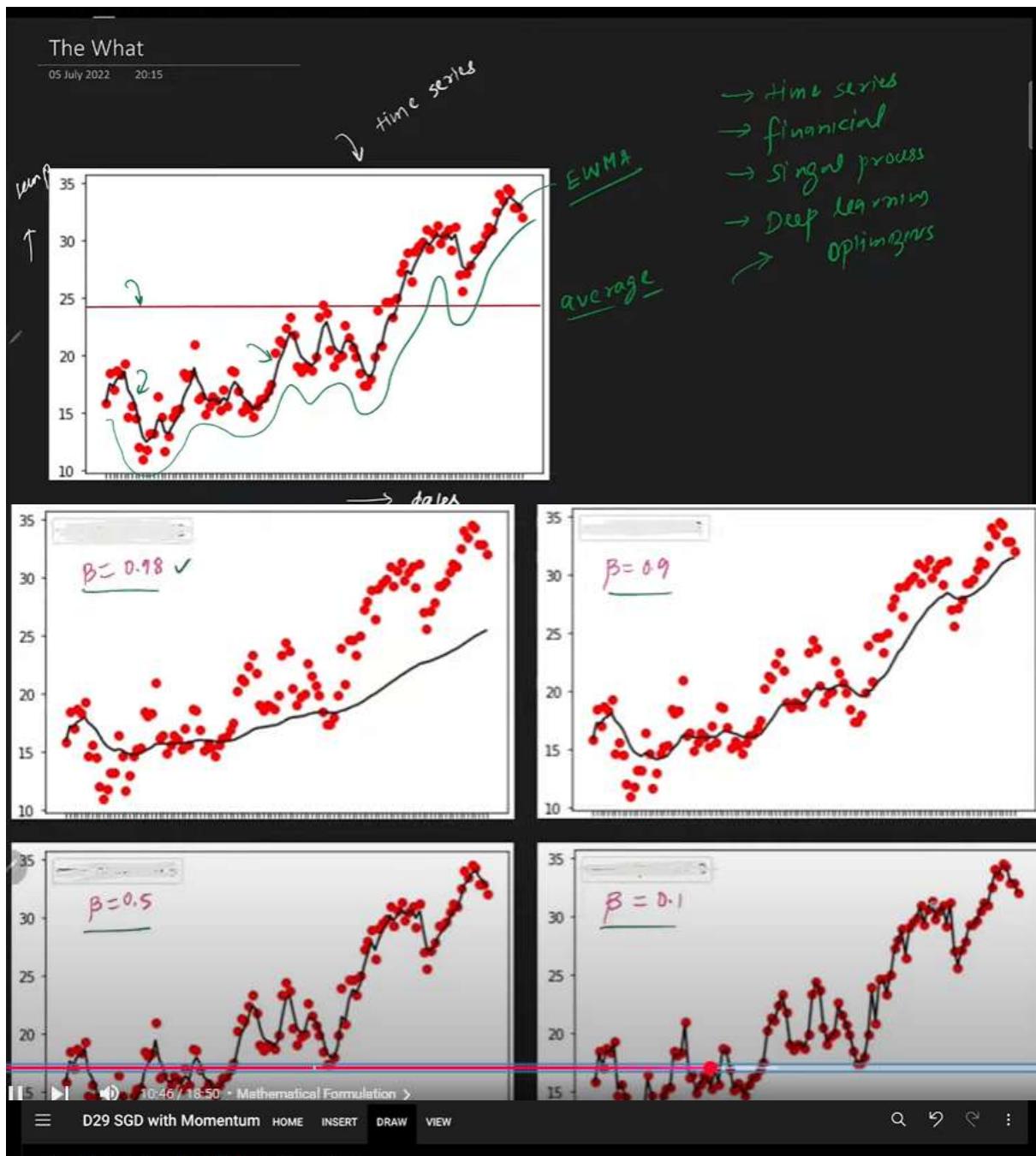
2) Adagrad

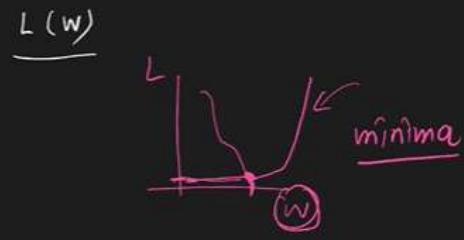
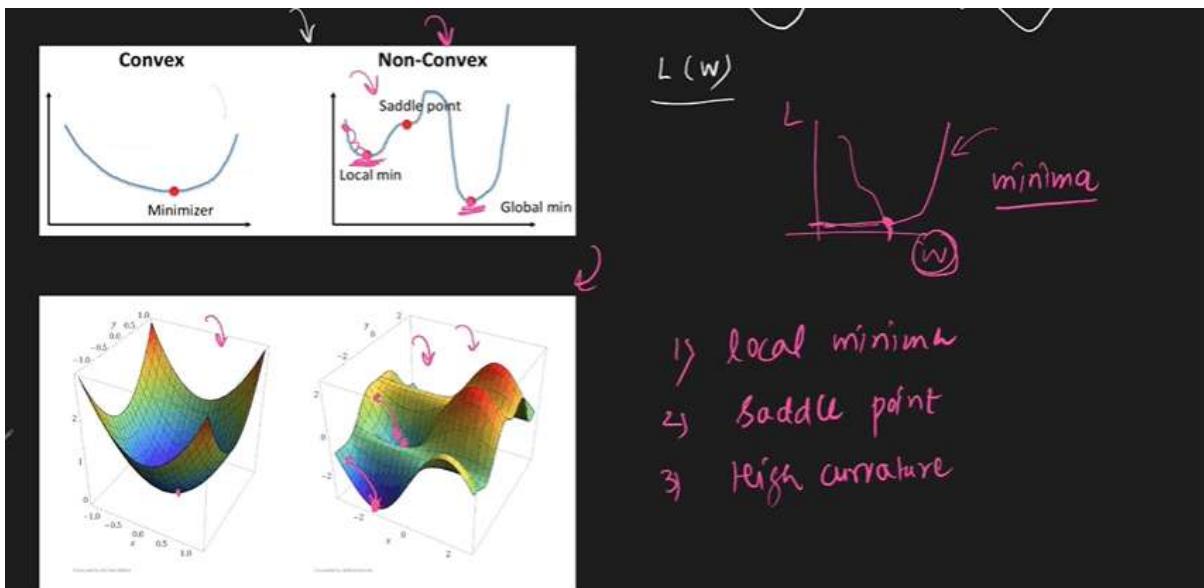
3) NAG

us RMSprop

5) Adam

EWMA Exponentially Weighted Moving Average or Exponential Weighted Average

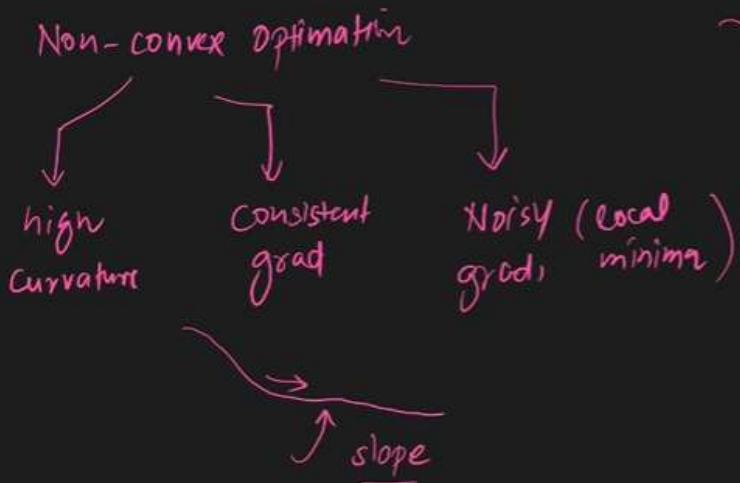




- 1) local minimum
- 2) saddle point
- 3) high curvature

## Momentum Optimization - The Why?

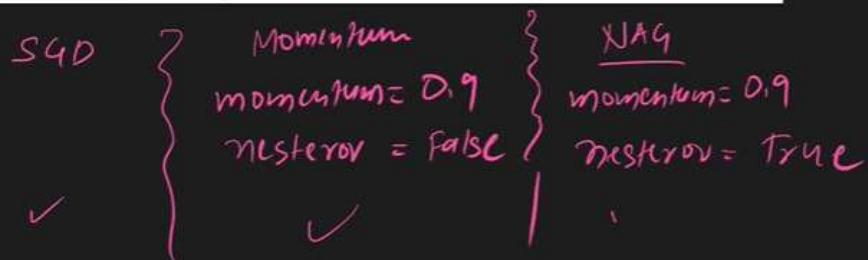
20 July 2022 14:28



## Keras Code

24 July 2022 13:18

```
tf.keras.optimizers.SGD()
    learning_rate=0.01, momentum=0.0, nesterov=False, name="SGD", **kwargs
)
```



batch gradient descent problem is speed

momentum uses past gradients for current update , but problem is oscillations

NAG overcame oscillations

but still problem was with sparse column

Adagrad solve sparse problem by decaying learning rate , but problem is it was not able to reach minima

RMS prop solved everything , no disv

so overall two optimzation mindsets

either momentum or LR decay , ADAM merged both

Mathematical Formulation

14 August 2022 10:45

$$w_{t+1} = w_t - \frac{\eta}{\sqrt{v_t + \epsilon}} * m_t \rightarrow \hat{m}_t = \frac{m_t}{1 - \beta_1^t} \rightarrow \hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

$\beta_1 = 0.9$        $\beta_2 = 0.99 \rightarrow \text{keras}$

where

$$\rightarrow m_t = \beta_1 m_{t-1} + (1 - \beta_1) \nabla w_t \rightarrow \text{momentum}$$
$$\rightarrow v_t = \beta_2 v_{t-1} + (1 - \beta_2) (\nabla w_t)^2 \rightarrow \text{Adagrad}$$

→ 1) momentum → merge  
→ 2) learning decay

ADAM IS BEST

THEN RMS PROP

## KERAS TUNER

```
●  from keras.layers import Dropout
def build_model(hp):
    model = Sequential()

    # Input layer + First hidden layer: Tune the number of neurons
    model.add(Dense(units=hp.Int('input_units', min_value=8, max_value=128, step=8),
                    activation='relu', input_dim=8))

    # Add a tunable dropout layer after the first hidden layer
    model.add(Dropout(hp.Choice('dropout_input', values=[0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9])))

    # Add a tunable number of hidden layers between 0 and 10
    for i in range(hp.Int('num_layers', min_value=0, max_value=10)):
        # Tune the number of neurons in each hidden layer
        model.add(Dense(units=hp.Int(f'layer_{i}_units', min_value=8, max_value=128, step=8),
                        activation='relu'))
        # Add a tunable dropout layer after each hidden layer
        model.add(Dropout(hp.Choice('dropout' + str(i), values=[0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9])))

    # Output layer
    model.add(Dense(1, activation='sigmoid'))

    # Tune the optimizer
    optimizer = hp.Choice('optimizer', values=['adam', 'sgd', 'rmsprop', 'adadelta'])

    model.compile(
        optimizer=optimizer,
        loss='binary_crossentropy',
        metrics=['accuracy']
    )

return model
```