

BUILDING A NEURAL MODEL FROM SCRATCH

NO TENSORFLOW &
PYTORCH

JUST MATHS 😊

— ASHISH
KISHORE

28x28

greyscale images

784
Pixels

m
training
images

Range from

0 - 255

0 = Black & 255 =

white

$$X = \begin{bmatrix} \text{---} & x^{(1)} & \text{---} \\ \text{---} & x^{(2)} & \text{---} \\ & \vdots & \\ \text{---} & x^{(m)} & \text{---} \end{bmatrix}^T = \begin{bmatrix} | & | & & | \\ & x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ & | & | & & | \end{bmatrix}$$

→ Each row is an example

→ Each row is going to
784 columns long

→ 784 rows each
column, corresponding
to pixel value

→ m columns
[Example]

784

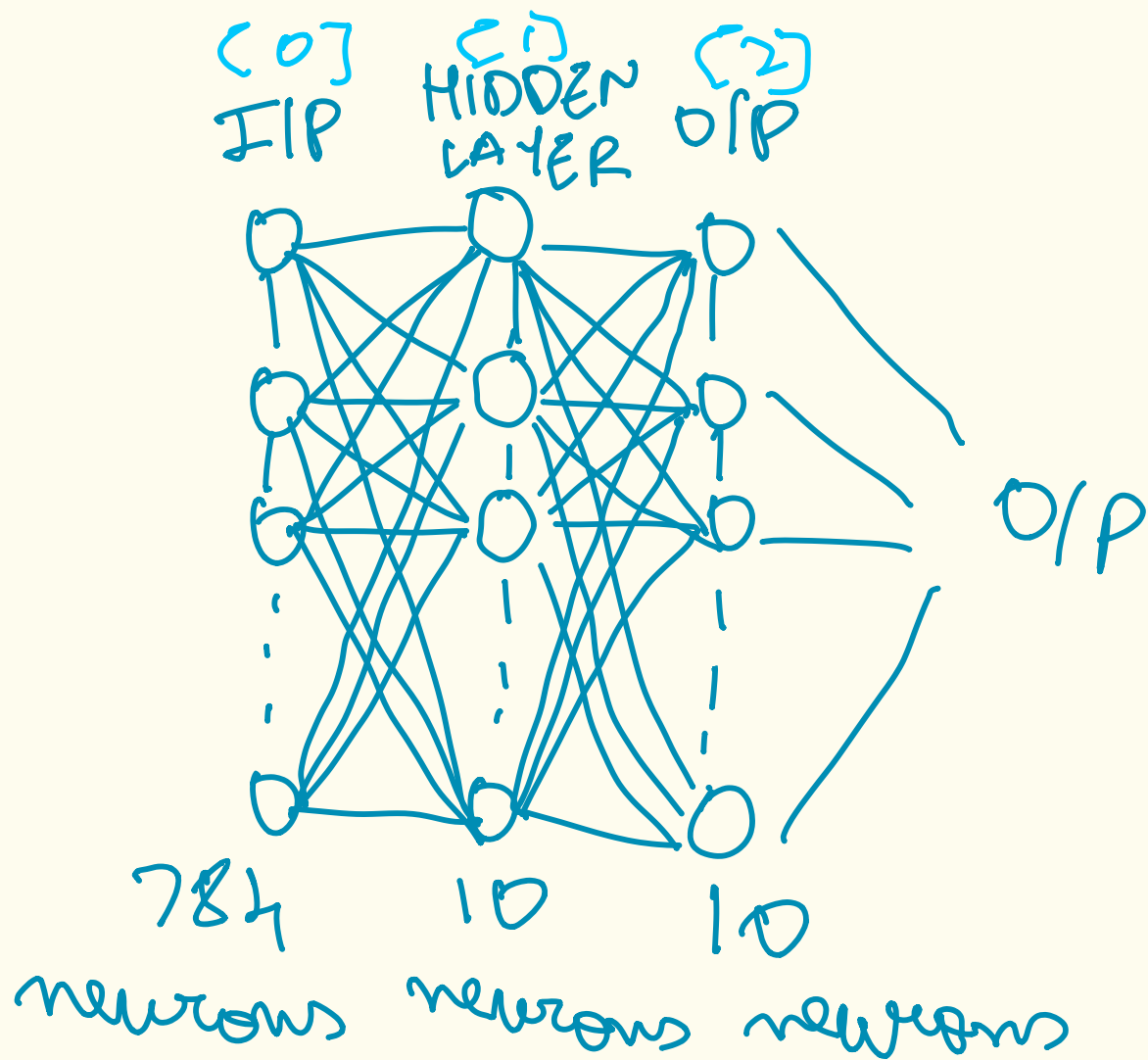
28x28

\Rightarrow

0, 1, 2, ..., 9

└──────────┘

10 classes



TRAINING

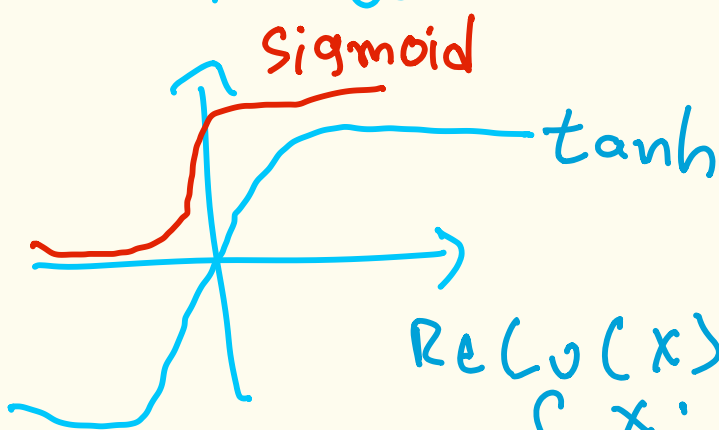
A) Forward Propagation: Take the image & run it through the NW

$$\begin{aligned}
 A^{[0]} &= [X] \quad (784 \times m) \\
 Z^{[1]} &= W^{[1]} A^{[0]} + b^{[1]} \rightarrow \text{Unactivated bias} \\
 10 \times m &= 10 \times 784 \quad 784 \times m \quad 10 \times 1 \Rightarrow 10 \times m
 \end{aligned}$$

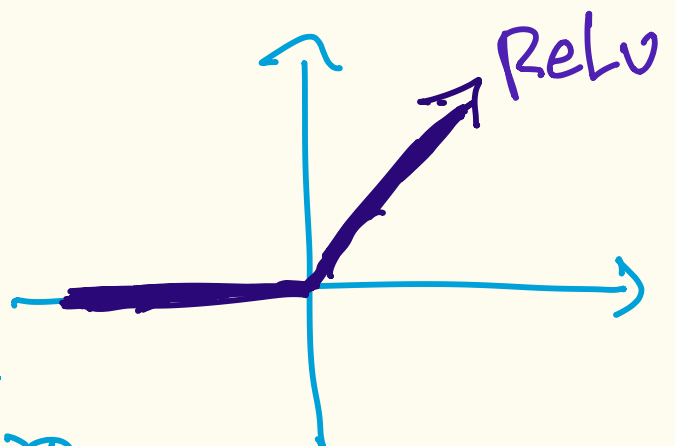
Apply activation function

- If there is no activation, each neuron/node will be a linear combination of previous nodes with bias term.
- The 2nd layer is gonna be the combination of nodes in 1st layer but 1st layer is just linear combination of IP layer. If there is only linear combination, with weights & biases there will be no interesting function but a complete linear regression.

So use of ACTIVATION FUNCTION is required.



$$\text{Relu}(x) = \begin{cases} x; & x > 0 \\ 0; & x < 0 \end{cases}$$



→ If any activation function is added layers turn into complex from linearity, layers become complex with preceding & becomes Non linear.

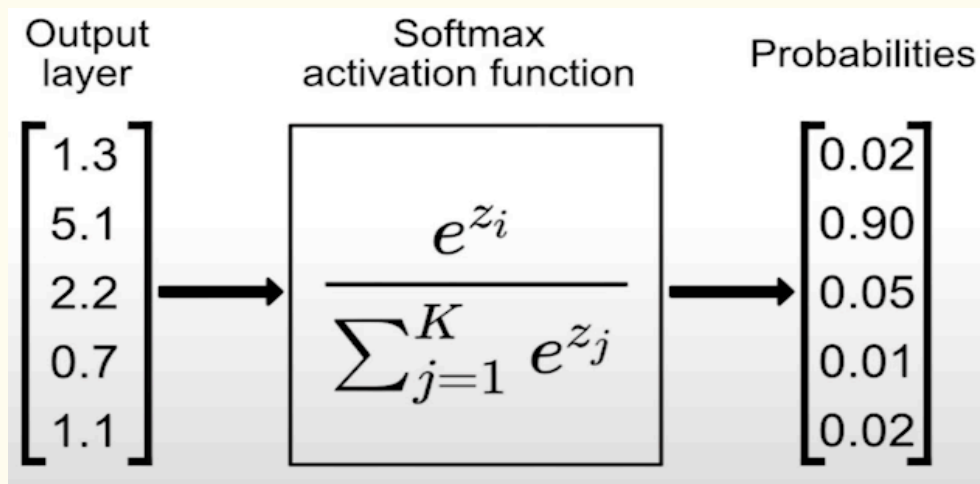
$$A^{[1]} = g(Z^{[1]}) = \text{ReLU}(Z^{[1]})$$

$$Z^{[1]} = W^{[2]} A^{[1]} + b^{[2]}$$

$$10 \times m \quad 10 \times 10 \quad 10 \times m \quad 10 \times 1 \Rightarrow 10 \times m$$

↓
Activated
1st Layer

Software Activation function ÷ It is applied to O/P layer, each 10 neurons corresponds to each of 10 digits, there must be probability for each digit b/w 0 & 1 with 1 absolute certainty & 0 no chance at all.



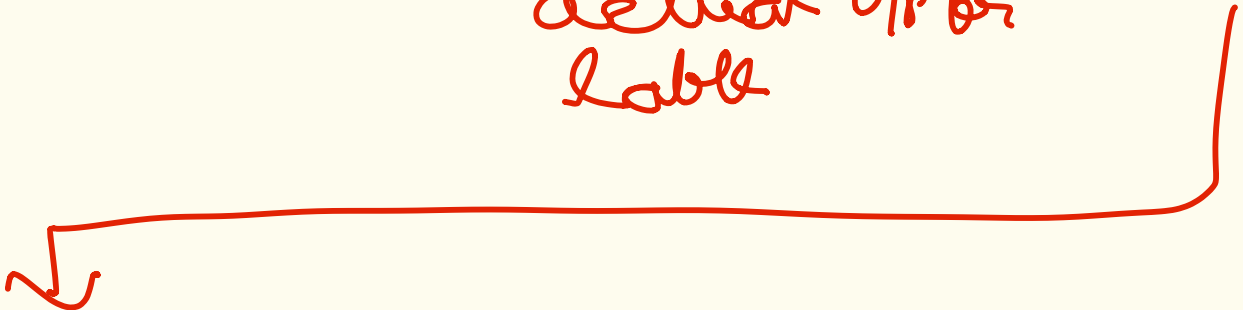
$e^{\text{that node}}$
Sum of $e^{\text{all nodes}}$
 O/P is b/w
 0-1

B) Back Propagation:

To get good weights & biases to make better predictions we use Back Propagation.

→ We go in the opposite way,

Prediction → Find out by how much prediction got deviated by actual O/P or Label → This gives us error



How much
each of these
weights & biases

contributed to that \rightarrow
error

Adjust
weights &
biases

$$\underset{10 \times m}{dZ^{[2]}} = \underset{10 \times m}{A^{[2]}} - \underset{10 \times m}{Y}$$

↓
Probabilities
of all digits

↓
One hot
Encoding of digits

$$\underset{10 \times 10}{dW^{[2]}} = \underset{10 \times m}{1/m} \underset{m \times 10}{dZ^{[2]} A^{[3]T}}$$

$$\underset{10 \times 1}{dB^{[2]}} = \underset{10 \times 1}{1/m} \sum \underset{10 \times 1}{dZ^{[2]}}$$

$$\underset{10 \times m}{dZ^{[1]}} = \underset{10 \times 10}{W^{[2]T}} \underset{10 \times m}{dZ^{[2]}} \cdot * \underset{10 \times m}{g'(Z^{[1]})}$$

$g' =$ Derivative of activation func

$$\underset{10 \times 784}{dW^{[1]}} = \underset{10 \times m}{1/m} \underset{m \times 784}{dZ^{[1]}} X^T$$

$$\underset{10 \times 1}{dB^{[1]}} = \underset{10 \times 1}{1/m} \sum dZ^{[1]}$$

c) Update the weights

$$W^{[1]} = W^{[1]} - \alpha dW^{[1]}$$

$$b^{[1]} = b^{[1]} - \alpha db^{[1]}$$

$$W^{[2]} = W^{[2]} - \alpha dW^{[2]}$$

$$b^{[2]} = b^{[2]} - \alpha db^{[2]}$$

$\alpha =$ learning rate,
hyperparameter

Cycle of Model

