

# Neural Nets from Scratch

7

$\Rightarrow$  'm' no. of such training images

Image -  $28 \times 28$  pixels = 784 pixels

784  
 $28 \times 28$

$\Rightarrow$  0, 1, 2 - ... 9  
10 classes

each class is for  
each digit 0-9

↳ each of these 784 pixels is a value b/w 0 & 255  
 ↓  
 totally black      totally white

So, if we make an array out of the pixels of each digit and map the value b/w 0 and 255.

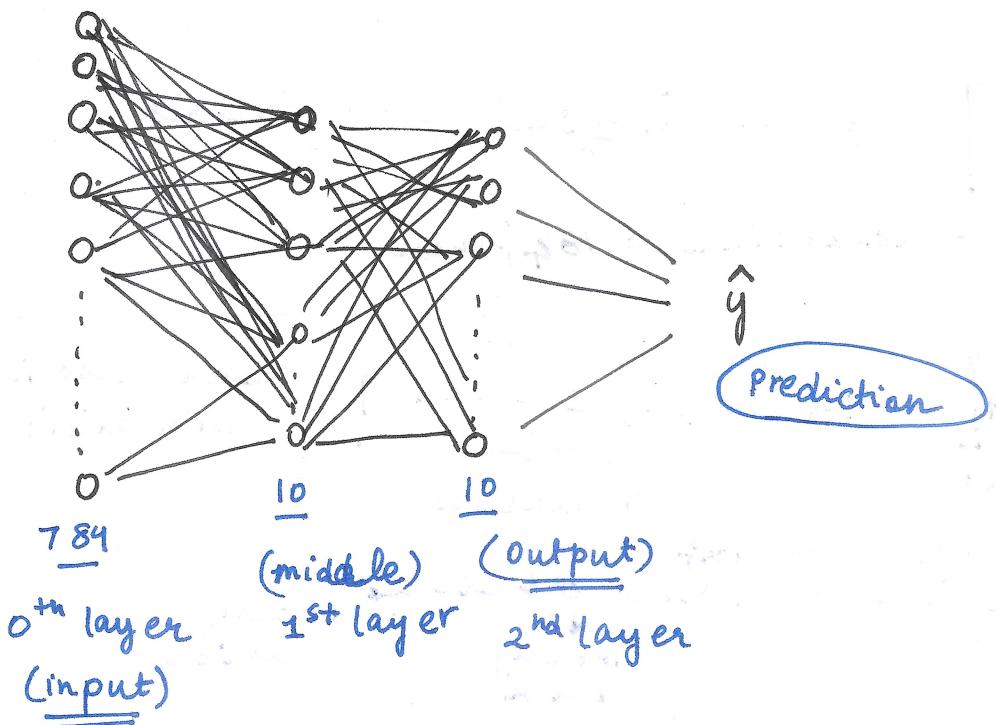
$$\text{we get an array } X = \begin{bmatrix} \xrightarrow{\hspace{1cm}} X^{(1)} \\ \xrightarrow{\hspace{1cm}} X^{(2)} \\ \vdots \\ \xrightarrow{\hspace{1cm}} X^{(m)} \end{bmatrix}$$

where, Each row = example image of digit  
 & 784 columns = corresponding pixels of that in one row one sample image.

If I were to transpose the matrix, each column is an example digit and rows correspond to pixels.

$$X = \begin{bmatrix} \xrightarrow{\hspace{1cm}} X^{(1)} \\ \xrightarrow{\hspace{1cm}} X^{(2)} \\ \vdots \\ \xrightarrow{\hspace{1cm}} X^{(m)} \end{bmatrix} \overset{T}{=} \begin{bmatrix} | & | & \dots & | \\ X^{(1)} & X^{(2)} & \dots & X^{(m)} \end{bmatrix}$$

## structure of the Net



## Setup of the Neural net -

### Step 1 # - FORWARD PROPAGATION

- 1  $A^{[0]} = X \quad (784 \times m)$  i.e. unactivated first layer which is input from dataset itself.
- 2  $g(\cdot)$  = activation function — it serves to have multiple connections to a node in all layers.

① In the absence of an activation function, each layer's nodes are just a combination of previous nodes, and we get just linear combinations.

In other words if, no multiple node connection to each layer  $\Rightarrow$

Fancy Linear regression

we add weights & a bias -

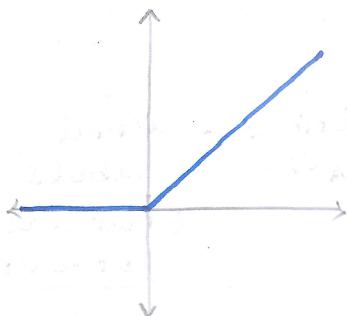
$$Z^{[1]} = W^{[1]} \cdot A^{[0]} + b^{[1]} \rightarrow Z^{[1]} = \text{term } A^{[0]} \text{ with,}$$

$\downarrow \qquad \downarrow \qquad \downarrow$   
 $10 \times m \quad 784 \times 10 \quad 10 \times m$

$\text{Weight} = W$   
 $\& \text{bias} = b$

Activation function - any nonlinear function such as tanh or sigmoid or ReLU

using ReLU( $n$ ), (rectified linear unit)



$$\text{ReLU}(n) = \begin{cases} n, & n > 0 \\ 0, & n \leq 0 \end{cases}$$

so, my  $A^{[1]}$  becomes the output get after I perform  $\text{ReLU}(z^{[1]})$

$$\Rightarrow A^{[1]} = \text{ReLU}(z^{[1]})$$

This is for the "Hidden" middle layer

For the output layer however, we need each of the 10 nodes to correspond to each of the 10 digits that can be recognised

$\Rightarrow$  Activation function we use is  $\rightarrow \text{SOFTMAX}(n)$

$$\text{i.e. } \text{SOFTMAX}(n) = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

SOFTMAX( $n$ ) makes the input vector ' $z$ ' of  $K$  real no.s and normalizes it into a probability distribution consisting of ' $K$ ' probabilities  $\propto$  input of number's exponent



We need each value of each of the output nodes to be b/w 0 & 1 to determine which is the given image closest to.

0 = no chance at all

1 = absolute certainty.

so,

$$A^{[2]} = \text{softmax}(z^{[2]})$$

## Step 2 # - BACKWARD PROPAGATION

- Essentially we start with our prediction, find out how much our weights & biases offset us from the correct ones
- Find and adjust weights & biases accordingly to minimize the loss function.

$$\left[ \begin{array}{c} \text{Error of} \\ \text{OUTPUT (2nd)} \\ \text{Layer} \end{array} \right]_{10 \times m} dZ^{[2]} = A^{[2]} - y \rightarrow \begin{array}{l} \text{Predictions - actual} \\ (A^{[2]}) \text{ labels} \\ (\text{that are} \\ \underline{\text{hot-encoded}}) \end{array}$$

$\downarrow$

derivative of the LOSS FUNCN w.r.t Weights  $dW^{[2]}_{10 \times 10} = \frac{1}{m} dZ^{[2]}_{10 \times m} \cdot A^{[1] T}_{m \times 10}$

$\downarrow$

Nothing but the Avg. Absolute Error in output layer

derivative of LOSS FUNCN w.r.t biases  $db^{[2]}_{10 \times 1} = \frac{1}{m} \sum dZ^{[2]}_{10 \times 1}$

$\downarrow$

i.e.  $[0, 0, 0, 0, 1, 0, 0, \dots, 0]_{m=1 \dots m}$  pos 4

similarly, for the middle "hidden" layer we compute the errors and also account for the error of the output layer & apply the corresponding ~~errors~~ weights in reverse.

Intuition →

error from OUTPUT (2nd) Layer  $\Rightarrow$  apply weights to it in reverse  $\Rightarrow$  get to the errors of the first (hidden) layer  $\Rightarrow$  multiply with derivative of the activation funcn

so,

$$\left[ \begin{array}{c} \text{Error of} \\ \text{Output of} \\ \text{HIDDEN 1st} \\ \text{Layer} \end{array} \right]_{10 \times m} dZ^{[1]} = W^{[2] T}_{10 \times 10} dZ^{[2]}_{10 \times m} \cdot g'(dZ^{[2]})_{10 \times m}$$

$\downarrow$

Transpose, to apply weights to Errors in 2nd layer to get to Errors in 1st layer (HIDDEN)

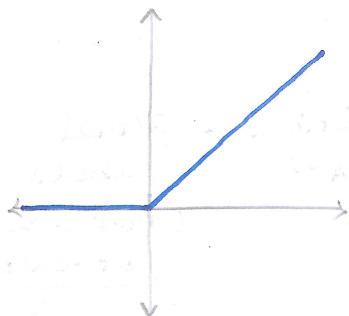
$\downarrow$

derivative of the Activation funcn  $\downarrow$

to undo the activation funcn to get proper error

Activation function - any nonlinear function such as tanh or sigmoid or ReLU

using ReLU( $n$ ), (rectified linear unit)



$$\text{ReLU}(n) = \begin{cases} n, & n > 0 \\ 0, & n \leq 0 \end{cases}$$

so, my  $A^{[1]}$  becomes the output get after I perform  $\text{ReLU}(z^{[1]})$

$$\Rightarrow A^{[1]} = \text{ReLU}(z^{[1]})$$

This is for the "Hidden" middle layer

For the output layer however, we need each of the 10 nodes to correspond to each of the 10 digits that can be recognised

⇒ Activation function we use is → SOFTMAX( $n$ )

$$\text{i.e. } \text{SOFTMAX}(n) = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

SOFTMAX( $n$ ) makes the input vector ' $z$ ' of  $K$  real no.s and normalizes it into a probability distribution consisting of ' $K$ ' probabilities  $\propto$  input of number's exponent



we need each value of each of the output nodes to be b/w 0 & 1 to determine which is the given image closest to.

0 = no chance at all

1 = absolute certainty.

so,

$$A^{[2]} = \text{softmax}(z^{[2]})$$

How much  
 did the  
 weights  
 &  
 biases  
 contribute  
 to  
 errors  
 in the "Hidden" layer

$$dW^{[1]} = \frac{1}{m} \sum_{i=1}^m dz^{[1]}_i \cdot X^{[T]}_{i \times 784}$$

$$db^{[1]} = \frac{1}{m} \sum_{i=1}^m dz^{[1]}_{i \times 1}$$

### Step 3 # UPDATE PARAMETERS

with ' $\alpha$ ' (learning rate  $\approx 0.01$ )

low  
value  
generally  
better

- set by user during  
gradient decent  
run on the loss function

- Once weights' and biases' losses are found, we can update parameters accordingly -

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

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

&

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

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