# ANN - DIGIT CLASSIFICATION.
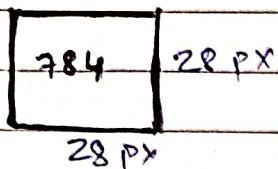
### Part 1 - Problem Statement :

A Digit classifying model of MNIST Dataset (kaggle)

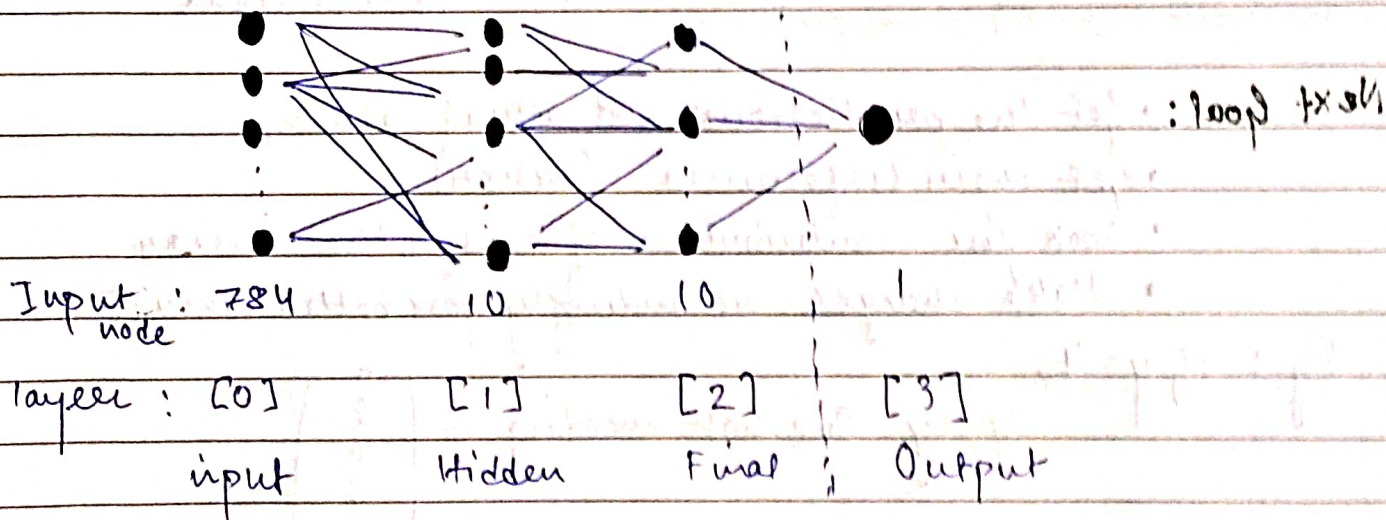### Part 2 - Math

pixel value :  $0 \longrightarrow 255$

black      white

| 784 | 28 px |

28 px

$I/P \Rightarrow X =$

$$\begin{bmatrix} \text{---} x^{[1]} \text{---} \\ \text{---} x^{[2]} \text{---} \\ \vdots \\ \text{---} x^{[m]} \text{---} \end{bmatrix}^T = \begin{bmatrix} \Big| & \Big| & & \Big| \\ x^{[1]} & x^{[2]} & \cdots & x^{[m]} \\ \Big| & \Big| & & \Big| \end{bmatrix}$$

m examples

$x^{[i]}$ contains 784 rows



: Next layer

Input : 784     10     10     1
node

layer : [0]     [1]     [2]     [3]

input     Hidden     Final     Output

### Forward Propagation :

1st layer
$$A^{[0]} = X \;(784 \times m)$$

$$Z^{[1]} = \underset{\underset{(10 \times 784)}{\text{weight}}}{w^{[1]}} \cdot \underset{\underset{(784 \times m)}{\text{I/P}}}{A^{[0]}} + \underset{\underset{(10 \times 1)}{\text{bias}}}{b^{[1]}}$$

$(10 \times m)$

$$A^{[1]} = g(z^{[1]}) = Relu(z^{[1]})$$

2nd layer
$$Z^{[2]} = \underset{(10 \times 10)}{w^{[2]}} \; \underset{(10 \times m)}{A^{[1]}} + \underset{(10 \times 1)}{b^{[2]}}$$

$(10 \times m)$

Camlin

Ⓐ Activation functions:



tanh

rectified linear unit (ReLu)

$$\begin{cases} \text{if } x > 0, & x \\ \text{if } x \leq 0, & 0 \end{cases}$$

→ to avoid nodes being the linear combination of previous layer.

3rd Layer $\{$ $[A^{[3]} = \text{softmax}(Z^{[2]})]$

Output layer $\begin{bmatrix} 1.3 \\ 5.1 \\ 2.2 \\ \vdots \\ 1.1 \end{bmatrix} \Rightarrow \begin{bmatrix} \dfrac{e^{z_i}}{\sum_{j=0}^{*} e^{z_j}} \end{bmatrix} \Rightarrow \begin{bmatrix} 0.02 \\ 0.91 \\ 0.03 \\ \vdots \\ 0.1 \end{bmatrix} \checkmark$

softmax activation

Predictions

Next Goal:
- Get the predictions and actual values,
- Get their differences (error)
- Look for contribution of bias to the error
- Make changes accordingly for better result.

Eg: if $y = 4$,

apply One hot encoding ⇒ $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ \vdots \end{bmatrix}$

apply backward propagation,

layer 2 $\begin{cases} dz^{[2]} = A^{[2]} - y \\ dw^{[2]} = \dfrac{1}{m} dz^{[2]} A^{[1]T} \\ \quad (10 \times 10) \qquad (10 \times m) \ (m \times 10) \\ db^{[2]} = \dfrac{1}{m} \sum dz^{[2]} \\ \quad (10 \times 1) \qquad (10 \times 1) \end{cases}$

{error in 2nd layer}

$$dz^{[1]} = w^{[2]T} \, dz^{[2]} \, .* \, g'(z^{[1]})$$

layer 1

(10×m)        (10×10) (10×m)      (→ derivative of activation fn to w.r.t.it.

$$dw = \frac{1}{m} \, dz \quad x^T$$

(10×728)     $\frac{1}{m}$ (10×m) (m×728)

$$db = \frac{1}{m} \sum dz^{[1]}$$

(10×1)      $\frac{1}{m}$ (10×1)

## Error Reduction (Gradient Descent):
### Repeat until Convergence {

$$w^{[1]} = w^{[1]} - \alpha \, dw^{[1]}$$

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

$$w^{[2]} = w^{[2]} - \alpha \, dw^{[2]}$$

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

} $\alpha$ = learning rate (hyper parameter)