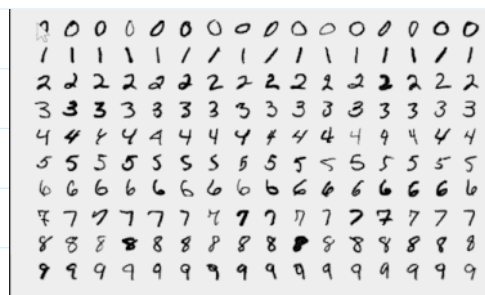


→ mnist data set



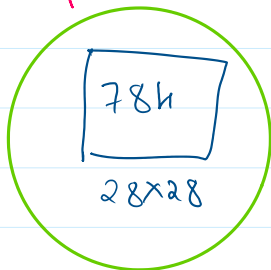
↑  
28

$\times m$

is the volume of dataset

← 28 pixd →

⇒ implementation



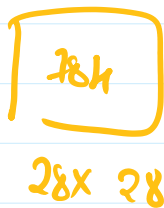
time's  $m$  ; training images

$x \times m$  is the training dataset

$x$

$$x = \begin{bmatrix} - & x^1 & - \\ - & x^2 & - \\ - & x^3 & - \\ - & x^m & - \end{bmatrix}^T = \begin{bmatrix} 1 & 1 & 1 & 1 & \dots & 1 \\ x^1 & x^2 & x^3 & x^4 & \dots & x^m \\ 1 & 1 & 1 & 1 & \dots & 1 \end{bmatrix}$$

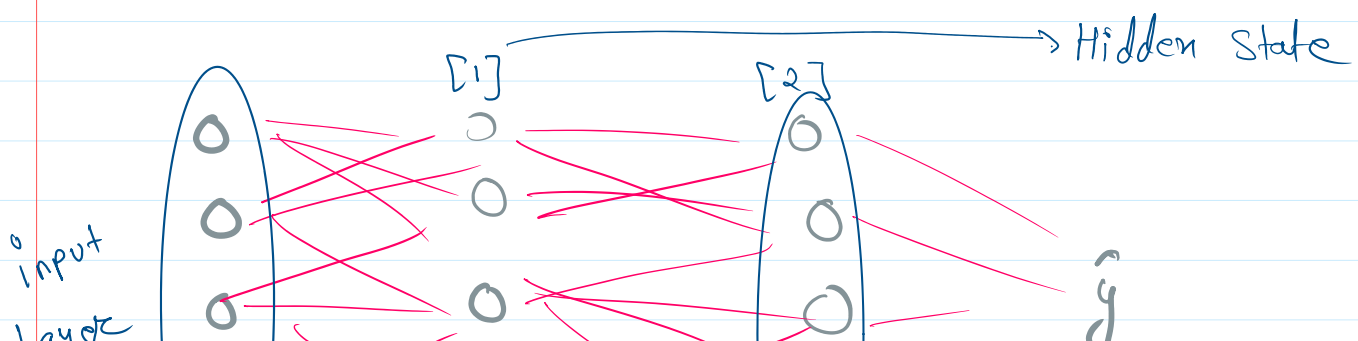
total dataset

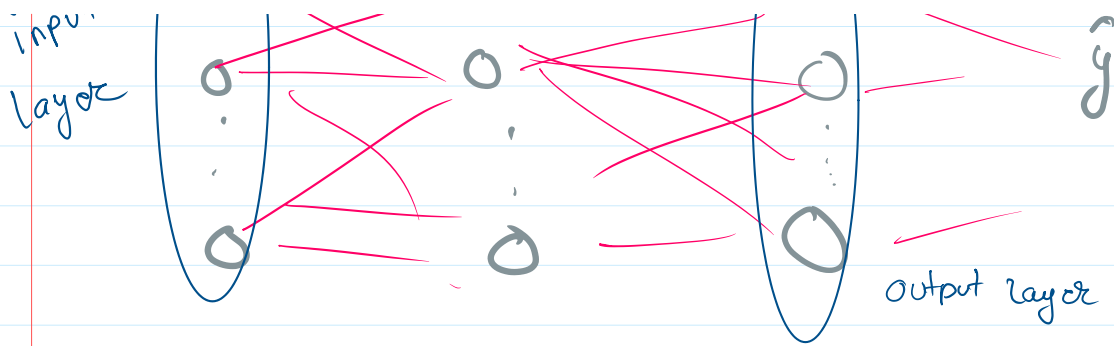


⇒

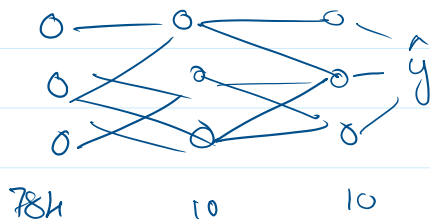
0, 1, 2, ..., 9

10 classes





## ① Forward Propagation



$$A^{[0]} = X \quad \left\{ 784 \times m \right\}$$

Unactivated  
1st layer

$$Z^{[1]} = W^{[1]} \cdot A^{[0]} + b^{[1]}$$

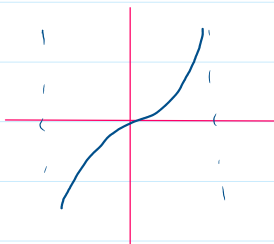
$10 \times m$        $10 \times 784$      $784 \times m$        $10 \times 1 \rightarrow 10 \times m$

$\left\{ \begin{array}{l} w \rightarrow \text{weight} \\ b \rightarrow \text{bias} \end{array} \right\}$

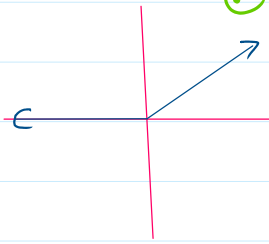
Activation  
Function

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

↳ without it, it would be a boring Linear Regression



$\tanh(h)$



rectified linear

$$\text{ReLU} = \begin{cases} x & x > 0 \\ 0 & x \leq 0 \end{cases}$$

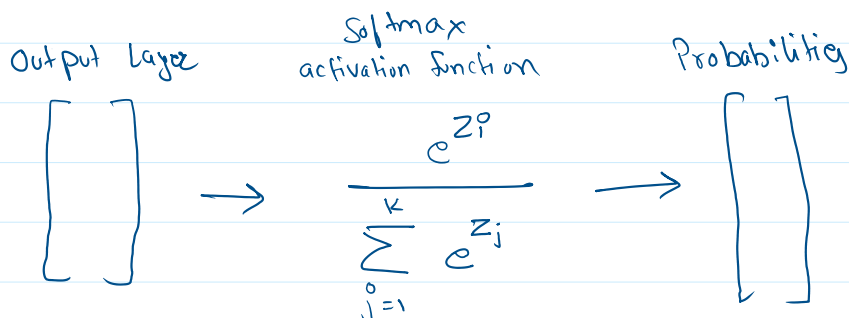
Unactivated  
2nd layer

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

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

Activation  
function

$$A^{[2]} = \text{Softmax}(Z^{[2]})$$



## (2) Backward Propagation

$$dz^{[2]} = A^{[2]} - y$$

$10 \times m$      $10 \times m$      $10 \times m$

$\left\{ \begin{array}{l} dz^{[2]} \leftarrow \text{error of 2nd layer} \\ y \rightarrow \text{one-hot code} \end{array} \right\}$

$$dw^{[2]} = \frac{1}{m} dz^{[2]} A^{[1]T}$$

$10 \times 10$      $10 \times m$      $m \times 10$

$\left\{ dw^{[2]} = \text{derivative of Cost function} \right\}$

$$db^{[2]} = \frac{1}{m} \sum dz^{[2]}$$

$10 \times 1$      $10 \times 1$

$\left\{ db^{[2]} = \text{avg of the absolute error} \right\}$

error for 1st layer

$$dz^{[1]} = W^{[2]T} dz^{[2]} * g'(z^{[1]})$$

$10 \times m$      $10 \times 10$      $10 \times m$      $10 \times m$

$\rightarrow \text{derivative of the activation function}$

Contribution of  $w^{[1]}$  to the error

$$dw^{[1]} = \frac{1}{m} dz^{[1]} x^T$$

$10 \times 784$      $10 \times m$      $m \times 784$

of  $w^{[1]}$

to the error

$10 \times 784$

$m$

$10 \times m$

$m \times 784$

Contribution

of  $b^{[1]}$

to the error

$10 \times 1$

$10 \times 1$

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

### ③ Parameters updation

$$\begin{aligned} w^{[1]} &:= w^{[1]} - \alpha dw^{[1]} \\ b^{[1]} &:= b^{[1]} - \alpha db^{[1]} \\ w^{[2]} &:= w^{[2]} - \alpha dw^{[2]} \\ b^{[2]} &:= b^{[2]} - \alpha db^{[2]} \end{aligned}$$

learning rate

