

# Neural Network Notes

The Problem: Digit Classification (MNIST Databases)

784  $\Rightarrow$  Output what digit that image represents:  $\underbrace{0, 1, 2, \dots, 9}_{10 \text{ classes}}$   
 $28 \times 28$  pixel training images

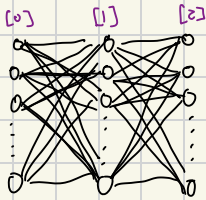
m training images

each row is an example:

$$X = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(m)} \end{bmatrix}^T \rightarrow \begin{bmatrix} \downarrow 784 \text{ rows per pixel} \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ \downarrow & \downarrow & & \downarrow \end{bmatrix}$$

## 2-Layer Neural Network

$\hookrightarrow$  each of 784 pixels maps to a node



0 layer  $\rightarrow$  just input

1st layer  $\rightarrow$  1st hidden layer

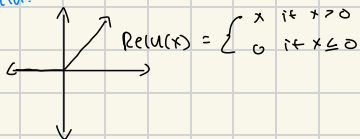
2nd layer  $\rightarrow$  output layer

## 1. Forward Propagation

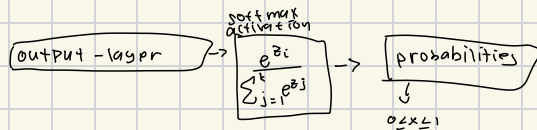
$\hookrightarrow$  getting output by passing image through network

$$\begin{aligned} A^{[0]} &= X \\ Z^{[1]} &= W^{[1]} A^{[0]} + b^{[1]} \rightarrow \text{unactivated first layer (weights and bias)} \\ A^{[1]} &= g(Z^{[1]}) = \text{ReLU}(Z^{[1]}) \rightarrow \text{activation function to make it a hidden layer (ReLU as func.)} \end{aligned}$$

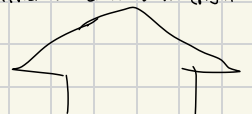
ReLU:



$$\begin{aligned} Z^{[2]} &= W^{[2]} A^{[1]} + b^{[2]} \rightarrow \text{unactivated 2nd layer} \\ A^{[2]} &= \text{softmax}(Z^{[2]}) \end{aligned}$$



The probabilities represent the likelihood that a pixel could represent a certain digit.



## 2. Backward Propagation

↳ going backwards and comparing result w/ actual value

Error of output: how much output is off by

$$\underset{\text{loss}}{\delta z^{[2]}} = \underset{\text{result}}{A^{[2]}} - \underset{\text{acc. value}}{Y} \rightarrow \text{will one-hot encode (i.e. } y=4: \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix})$$

How much weight/bias contributed to error in output layer:

$$dW^{[2]} = \frac{1}{m} \delta z^{[2]} A^{[1]T} \rightarrow$$

$$db^{[2]} = \frac{1}{m} \sum \delta z^{[2]} \rightarrow \text{average of absolute error}$$

Error of hidden layer 1: how much hidden layer 1 was off by

$$\delta z^{[1]} = W^{[2]T} \underset{\substack{\text{Applying weights to error} \\ \text{from 2nd layer}}}{\delta z^{[2]}} \cdot \underset{\substack{\text{derivative of} \\ \text{Activation func.} \\ \text{(Undoing to get right error)}}}{g'(z)}$$

How much weight/bias contributed to error in first layer:

$$dW^{[1]} = \frac{1}{m} \delta z^{[1]}$$

## 3. Update parameters

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

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

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

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

$\alpha$ : learning rate  $\rightarrow$  hyperparameter  
not trained by model

