Training computer to recognize handwritten digits.

- Using MINST data base of 28×28 pixel digits

28 | 784 | } training image dimensions | each pixel has 0...255 weighting for intensity
28

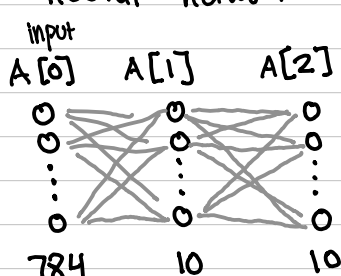We have $m$ training images, put into a row matrix (csv):

$$X = \begin{bmatrix} \underline{\quad x^{[1]} \quad} \\ \underline{\quad x^{[2]} \quad} \\ \vdots \\ \underline{\quad x^m \quad} \end{bmatrix}^T \Rightarrow \begin{bmatrix} | & | & & | \\ x^{[1]} & x^{[2]} & \cdots & x^{[n]} \\ | & | & & | \end{bmatrix}$$

Goal: Take the 28×28 or 784 pixel image and narrow to 10 digit options (0, 1, ..., 9)

Two layer neural network:   [1]: hidden layer
                            [2]: output layer

input
A[0]    A[1]    A[2]



784     10      10

## Training the Network:         ★ For this Neural net   $m=1$

### 1. Forward Propagation: running through our network [0]→[2]

\* $A^{[0]} = X$  [784 × m]          \* input layer

\*\* $Z^{[1]} = W^{[1]} A^{[0]} + b^{[1]}$      \*\* unactivated A[1] layer
[10 × m] [10×784][784×m] [10×m]
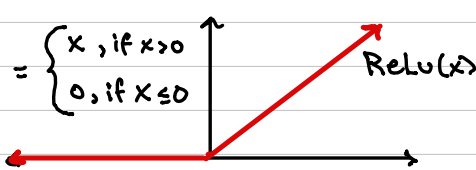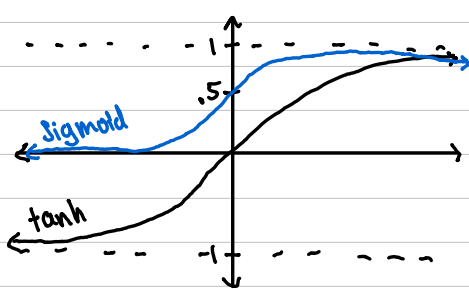
$W^{[1]}$ corresponds to a weight for each connection between the neurons in our $A^{[0]}$ to $A^{[1]}$ layer.

\*\*\* $A^{[1]} = g(Z^{[1]}) = ReLu(Z^{[1]})$        \*\*\* Normalized layer using activation funct.
        activation function

### Different Activation Functions:



Sigmoid
tanh

$= \begin{cases} x, \text{ if } x > 0 \\ 0, \text{ if } x \leq 0 \end{cases}$      $ReLu(x)$

- Using ReLu
- $\sigma$ and tanh are more complex

Going from layer 1 to layer 2:

$$Z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$$
  10×m   10×10  10×m    10×m

$$A^{[2]} = softmax(Z^{[2]})$$

$$softmax = \frac{e^{z_i}}{\sum_{j=1}^{K} e^{z_j}}$$   • We use this to derive a column of values that add to 1

### 2. Backwards Propagation: Optimization of weightings/biases

## Second Layer

$$dZ^{[2]} = A^{[2]} - Y$$
  10×m    10×m   10×m

- taking prediction and subtracting the one hot encoding
- One hot Y puts a 1 at the index which is encoding
- shows the error of the second layer

$$dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1]T}$$
  10×10      10×m   m×10

$$db^{[2]} = \frac{1}{m} \sum dz^{[2]}$$   • Average of absolute error
  10×1

## First Layer

$$dz^{[1]} = W^{[2]T} dz^{[2]} \cdot * g'(Z^{[1]})$$   • derivative of activation
  10×m   10×10   10×m          10×m

$$dW^{[1]} = \frac{1}{m} dz^{[1]} X^T$$
  10×784     10×m   m×784

$$db^{[1]} = \frac{1}{m} \sum dz^{[1]}$$
  10×1        10×1

### 3. Updating Parameters — $\alpha$: learning rate
                            - independent of model

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

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

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

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