# Implementing a Neural Network for MNIST digit classification

→ The neural network to implement will allow MNIST digit classification

↳ Training images are 28×28 pixels $\boxed{784}$ ²⁸

↳ Each pixel has a pixel value [0,256], where 0 = Black and 255 = white

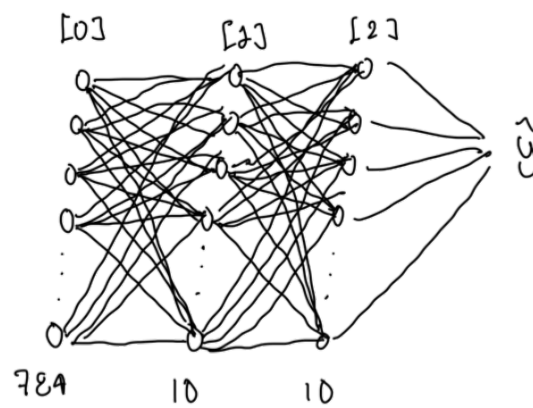→ The images can be represented as a matrix

↳ Normally the matrix would be something like this: $A = \begin{bmatrix} — x^{(1)} — \\ — x^{(2)} — \\ — x^{(3)} — \\ — x^{(n)} — \end{bmatrix}$ ← each row constitutes an example

however, I want to work with columns, so I will transpose it $A^T = \begin{bmatrix} | & | & | & | \\ x^{(1)} & x^{(2)} & x^{(3)} & x^{(n)} \\ | & | & | & | \end{bmatrix}$

→ Each column is an example, so each column has 784 rows corresponding to each pixel.

→ The neural network will have 2 layers

↳ 0th : 784 nodes (input layer)
↳ 1st : 10 units (hidden layer)
↳ 2nd : 10 units (output layer) → digit



[0]    [1]    [2]

ŷ

784   10   10

## 1) Forward Propagation

→ Taking an image and running it through the network and computing the output

$A^{[0]} = X$   $(784 \times m)$ → Input layer

Unactivated 1st layer ← $Z^{[1]} = W^{[1]} A^{[0]} + b^{[1]}$ { This basically means to multiply the input by the weight of each connection to the 1st layer. and add a bias

       10×m    10×784   784×m    10×1

       Weight      Bias

Now we need to apply an activation function, otherwise each node would just be a linear combination of the nodes before it plus a bias, without a activation function it would just be a linear regression.

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

          ↓ (Rectified Linear Unit)



$ReLU(x) \begin{cases} x \text{ if } x>0 \\ 0 \text{ if } x\leq0 \end{cases}$

Now to get from layer 1 to layer 2 first we add the weights on the bias

$Z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$ { multiply the 1st layer by the weight of each connection to the 2nd layer and add a bias

→ Now the second activation function will be Softmax

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

       ↳ This will give a probability (between 0 and 1)

$\dfrac{e^{z_i}}{\sum_{j=1}^{K} e^{z_j}}$

## 2) Backwards Propagation (Back Prop.)

→ Backwards propagation serves as a way to adjust the weights and biases.

→ It is the opposite process to forward propagation, so it starts with a prediction and finds how much the prediction deviated from the actual label (so instead of giving a success probability it gives an epsilon (error rate)) so it is possible to see how much did the previous weights and biases contributed to the actual error, and then adjust them accordingly.

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

10×m    10×m   10×m

↳ Error of the 2nd layer   ↳ Predictions   ↳ Actual label

→ this y will need to be one-hot encoded

$dW^{[2]} = \dfrac{1}{m} dZ^{[2]} A^{[1]T}$

10×10       10×m   m×10

↳ Average of the absolute error, how much the 2nd layer was off by from the prediction, and how much should the weights of the 2nd layer be adjusted and bias

$db^{[2]} = \dfrac{1}{m} \sum dZ^{[2]}$

10×1        10×1

$dZ^{[1]} = W^{[2]T} dZ^{[2]} \cdot g'(Z^{[1]})$

10×m    10×10   10×m    10×m

↳ How much was the hidden layer off by, does propagation in reverse → error from the second layer and aplaying the weights to it in reverse to get to the 1st layer, and then undoing the activation function (by its derivative) to get the proper error for the 1st layer

$dW^{[1]} = \dfrac{1}{m} dZ^{[1]} X^T$

10×784      10×m   m×784

↳ // //

$db^{[1]} = \dfrac{1}{m} \sum dZ^{[1]}$

10×1

## 3) Then, the parameters are updated:

$W^{[1]} = W^{[1]} - \alpha \, dW^{[1]}$
$b^{[1]} = b^{[1]} - \alpha \, db^{[1]}$
$W^{[2]} = W^{[2]} - \alpha \, dW^{[2]}$
$b^{[2]} = b^{[2]} - \alpha \, db^{[2]}$

{ Where alpha (α) is the learning rate and is a hyper parameter → Means it is not trained by the model, when the cicle is executed (when gradient descent is applied) alpha is set by whoever excecutes the code.

## 4) go back to forward propagation