# MNIST Neural Network



28 Pixel

784 Pixels

28 Pixel

$$X = \begin{bmatrix} \underline{\quad} x^{(1)} \underline{\quad} \\ \underline{\quad} x^{(2)} \underline{\quad} \\ \vdots \\ \underline{\quad} x^{(n)} \underline{\quad} \end{bmatrix}^{T} = \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \cdots & x^{(n)} \\ | & | & & | \end{bmatrix}$$

Each row is an example, and each row is 784 Colums long each Corresponding to a pixel in the Image

After transposing we have 784 rows Corresponding to each pixel

We will build a neural network to predict a digit from 0 to 9 given an Image

Input layer

Output layer



$\hat{Y}$

Prediction

Our Neural Net will have 3 layers

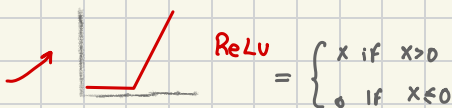784

Parameters

hidden layer

10

10

# Forward Prop

1. $A^{[0]} = X$ $(784 \times m)$ This is our input layer

2. $Z^{[1]} = W^{[1]} A^{[0]} + b^{[1]}$

   ↳ Unactivated first layer obtained by multiplying (dot product) of our input layer and a weight and adding a bias term (constant)

3. $A^{[1]} = g(Z^{[1]}) = ReLu(Z^{[1]})$

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

   ↳ In this step we apply an activation function (ReLu) so the layers are not linear combinations of each other.
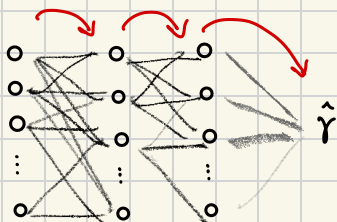
4. $Z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$

   ↳ $Z^{[2]}$ is our unactivated second layer is equal to a second weight parameter times our (activated first layer) $A^{[1]}$ plus a bias term

5. $A^{[2]} = SoftMax(Z^{[2]})$

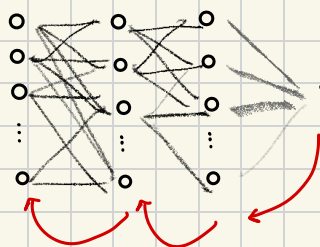   ↳ Finaly we use softmax to convert the output layer to probability

| Output layer | SoftMax | Probabilities |
|---|---|---|
| $\begin{bmatrix} 1.3 \\ 5.1 \\ 2.2 \\ 0.7 \end{bmatrix}$ | $\dfrac{e^{z_i}}{\sum_{j=1}^{k} e^{z_j}}$ | $\begin{bmatrix} 0.02 \\ 0.9 \\ 0.05 \\ 0.01 \end{bmatrix}$ |



$\hat{\gamma}$

So in forward propagation we navegate through our neural net so we can get to our predictor.

# Back Prop

We use back propagation to adjust the weights and biases to make good predictions.



In this case we will start at the prediction and we will see by how much did it deviated from the actual label

That will give us an error. Then we will see how much the previous weights and biases contribute to the error. So we can adjust the weights.

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

So $dz^{[2]}$ Represents the error of the second layer. So its our predictions $(A^{[2]})$ and subtract the actual errors $(y)$.

↳ This subtraction is made by 1-hot encoding

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

$dw^{[2]}$ is the derivative of the loss function with respect to the weights in layer 2

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

This is the averedge of the absolut error. this is how much the second layer was off by from the prediction

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

Here we take the error from the second layer and apply the weights in reverse to get to the errors in the first layer. And multiply the derivatie of g to undo the activation function

5. $\begin{cases} dw^{[1]} = \frac{1}{m} dz^{[1]} x^T \\ db^{[1]} = \frac{1}{m} \sum dz^{[1]} \end{cases}$

And we do the same thing that we did with the second layer to calculate how much $w^{[1]}$ and $b^{[1]}$ contributed to the error in the first layer.

After finding how much each weight term and each bias term contributed to the error we update our parameters

$$W^{[1]} = W^{[1]} - \alpha \, dw^{[1]}$$
$$b^{[1]} = b^{[1]} - \alpha \, db^{[1]} \quad$$ <span style="color:red">$\alpha$ is the learning rate</span>

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

After our parameters are updated we go through Forward-Prop again