

Objective: Build a neural network for hand-written number detection.

Theory

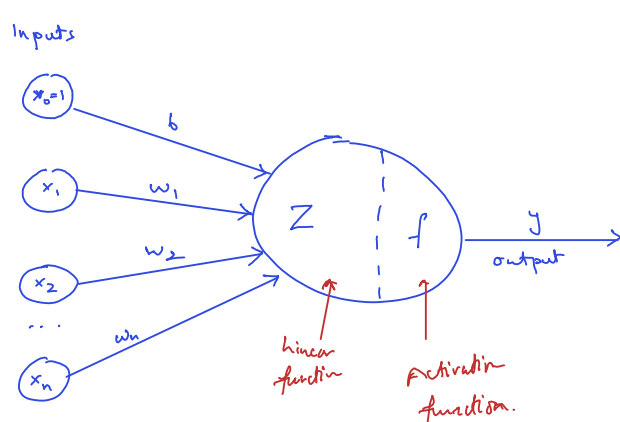
Definition of a neural network

They are machine learning models that mimic the complex functions of the human brain.

These models consist of interconnected nodes or neurons that process data, learn patterns, and enable tasks such as pattern recognition & decision-making.

It follows a three-stage process:

- ① Input Computation: Data is fed into the network.
- ② Output Generation: Based on the current parameters, the network generates an output.
- ③ Iterative Refinement: The network refines its output by adjusting weights & biases, gradually improving its performance on diverse tasks.



Working of a neural network

Forward Propagation

input layer \rightarrow hidden layer \rightarrow output layer

This process is known as forward propagation.

During this phase:

- ① Linear Transformation: Each neuron in a layer receives inputs, which are multiplied by the weights associated w/ the connections. These products are summed together, & a bias is added to the sum.

Representation: $Z = w_1x_1 + w_2x_2 + \dots + w_nx_n + b$

w = weights,
 x = input values,
 b = bias.

- ② Activation: The result of linear transformation (Z) is then passed through an activation fn. The activation fn. is crucial because it introduces non-linearity into the system, enabling the network to learn more complex patterns.

Ex: ReLU, sigmoid, and tanh

Backward Propagation

After forward propagation, the network evaluates its performance using a loss fn, which measures the difference b/w the actual output & the predicted output. The goal of training is to minimize this loss. This is where backward propagation comes to play.

- ① Loss Calculations: The network calculates the loss, which provides a measure of error in the prediction. The loss fn. could vary: common choices are mean squared error for regression tasks or cross-entropy loss for classification.

- ② Gradient Calculation: The network computes the gradients of the loss fn. w/ respect to each weight & bias in the network. This involves applying the chain rule of calculus to find out how much each part of the output error can be attributed to each weight & bias.

- ③ Weight Update: Once the gradients are calculated, the weights & biases are updated using an optimization algorithm such as Stochastic Gradient Descent (SGD). The weights are adjusted in the opposite direction of the gradient to minimize the loss. The size of the step taken in each update is determined by the learning rate.

A neural network can be supervised or unsupervised,

In supervised learning, labeled data is used to train a neural network so that it may learn to map inputs to matching outputs.

Unsupervised learning works w/ unlabeled data & looks for structures or patterns in the data.

Tools to be Used

- Python,
- numpy, — for numerical computation.
- matplotlib, — for visualizing the data.
- MNIST dataset. — ready to use labelled data
- pandas. — for reading & loading the data from the MNIST dataset.

Steps for the Project

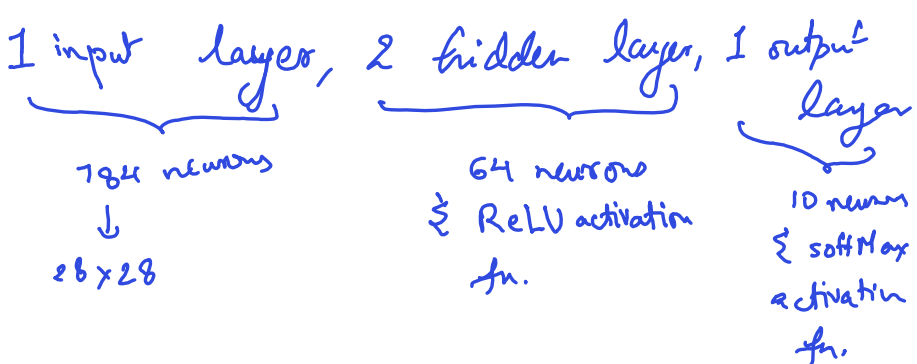
① Load & Preprocess the data.

MNIST dataset,

28x28 pixels in size (grayscale)

We use pandas to load & process.

② Defining the Neural Network Architecture



③ Compute the gradients.

Repeat weight adjustment until algorithm converges on a set of weights that produces the lowest possible error.

Error % $\neq 0$, this means the model has stopped learning.

④ Create a neural network object.

Initialize a neural network w/ 784 input features, 64 hidden nodes & 10 output neurons. And run the backpropagation algorithm for 10 epochs. w/ a learning rate of 0.1.

⑤ Train the neural network.

Once the neural network object is created, it can now be trained using the dataset & then used to make predictions on new data.

⑥ Evaluating the neural network.

The performance of the network will be evaluated by comparing the predicted o/p w/ the actual o/p & calculating a metric such as accuracy.

Additionally, analyze any patterns or trends in the errors made by the network to identify areas for improvement.

784
28x28 px

m training images

can be represented as matrix

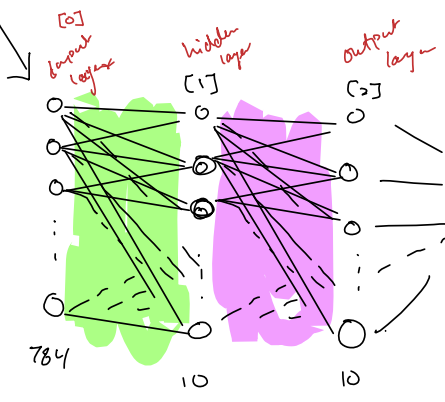
$$X = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(m)} \end{bmatrix} = \begin{bmatrix} | & | & \dots & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & \dots & | \end{bmatrix}$$

Transpose for calc. purpose

Each row is 784 cols wide, each col represents a pixel from the image the value stored in each cell $\in [0, 255]$

now each column is an example.

color range ex: 0 = white 255 = black



Forward Propagation Calc.

Basically taking an image from matrix and running thru the network

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

X is any input from

Scaling w matrices $w^{[1]}$ = weights matrix for hidden layer

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

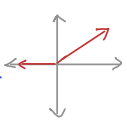
$b^{[1]}$ = bias values matrix for hidden layer

Now, we apply an activation function to hidden layer matrix. Why? \rightarrow The hidden layer is a linear combination of the input layer, & the output layer is a linear combination of the hidden layer. w/o the activation fn, there is no change so it is as if the o/p layer is the linear combination of the i/p layer. Obviously this is not desired. Also makes the model more complex & powerful.

Here, we use Rectified Linear Unit. (ReLU)

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

very simple fn. makes our graph look like



So,

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

Example of softmax act. fn.

$$softmax(x)_i = \frac{e^{x_i}}{\sum_{j=1}^N e^{x_j}}$$

Softmax always delivers in a range $\in (0, 1)$

Now, our unactivated second layer.

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

O/p layer

1.3
5.1
2.2
0.7
1.1

$\rightarrow softmax(x)_i \rightarrow$

0.02
0.90
0.05
0.01
0.02

softmax(1.3) = $\frac{e^{1.3}}{e^{1.3} + e^{5.1} + e^{2.2} + e^{0.7} + e^{1.1}} = \frac{3.66529}{181.73413} = 0.02019... \approx 0.02$

Probabilities

This time we apply another activating fn.

softmax(x)

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

now we get a prediction based on the image input :)

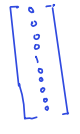
why? Since now we deal w/ the o/p layer, where each of the 10 nodes correspond to a number that represents the predicted output, we want each of them to have a probability a value $\in (0, 1)$. So, to get that we use softmax.

Backward Propagation

We use this algorithm to optimize the weights & biases

$$dz^{[2]} = A^{[2]} - y \rightarrow \text{one-hot encode}$$

Ex, if $y = 4$



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

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

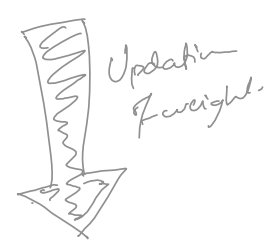
average formula Avg. of absolute error

for 1st layer

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

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

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



$$w^{[1]} = w^{[1]} - \alpha dw^{[1]}$$

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

$$w^{[2]} = w^{[2]} - \alpha dw^{[2]}$$

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

α = learning rate

α is a hyperparameter, its set by us.

Repetition for set 'epochs'.