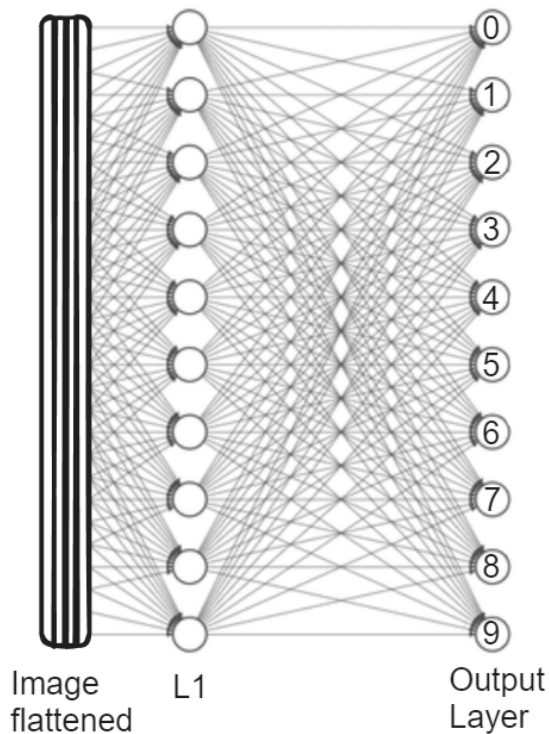


Neuron Understanding



This is a map of the simplest neural network created for Mnist dataset

This Neural Network flattens the 28×28 image into a long array to be sent into the network

It has 3 layers -

- Input layer (784 neurons)
- Layer 1 (10 neurons)
- Output Layer (10 neurons) - (Probabilities)

This paper targets to understand - In the condition of large data compression does a neuron depict significant understandable variables

(Spoiler - yes)

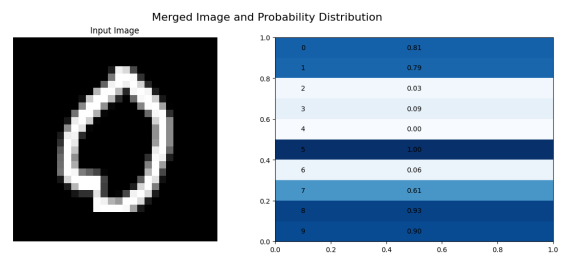
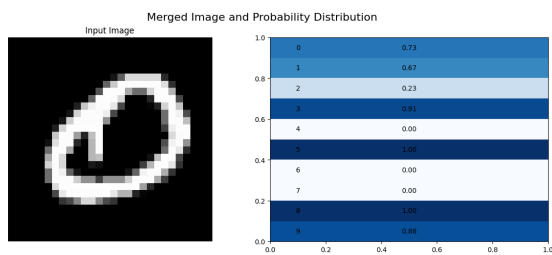
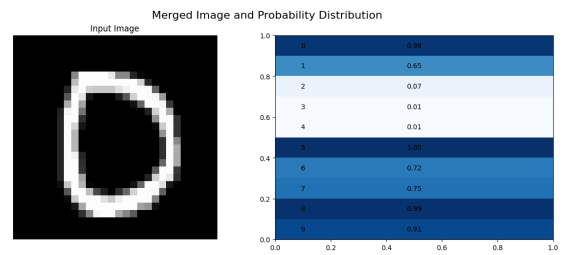
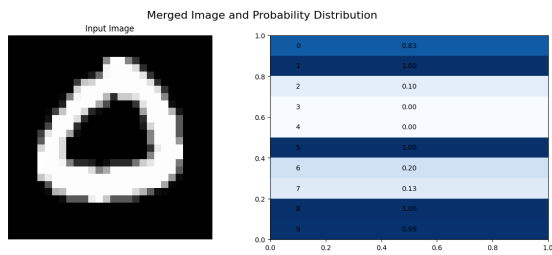
First of all this neural network is trained and gets up to 93% accuracy in test dataset

Secondly the entire code is uploaded

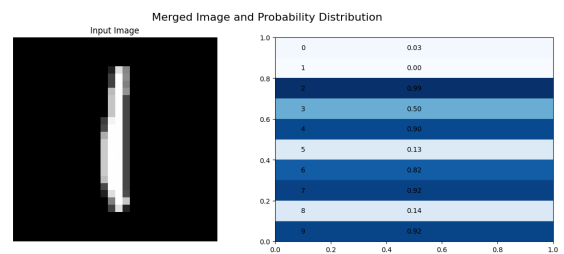
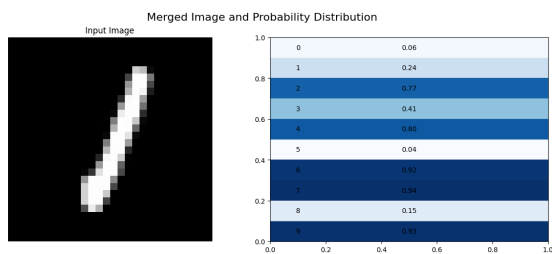
The main target to be understood is the activation of layer 1 (or the middle layer)

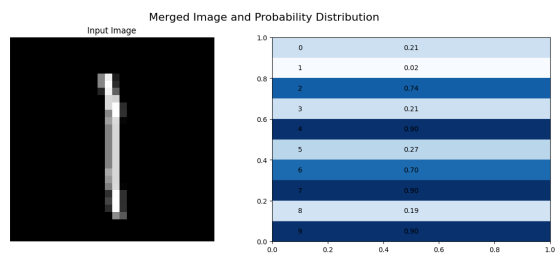
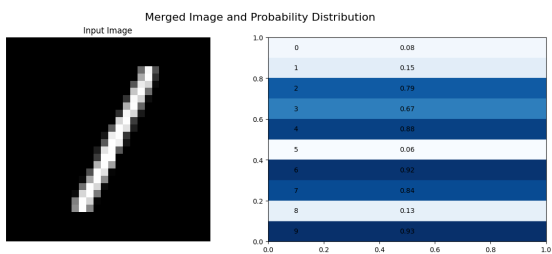
On the following examples, these were the activations of L1 layer neurons -

0

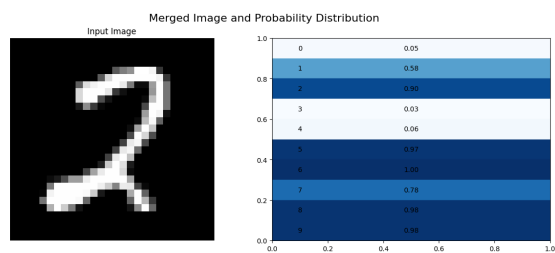
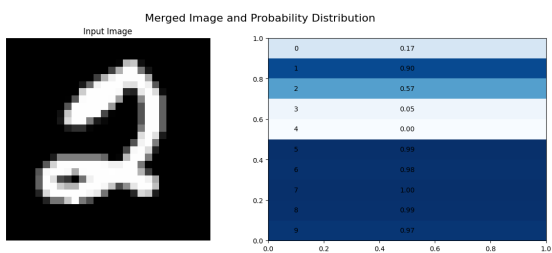
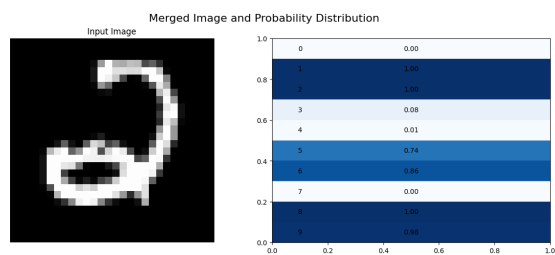
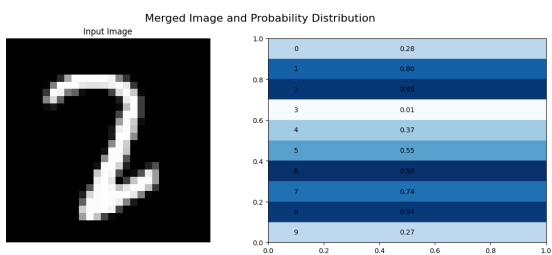


1

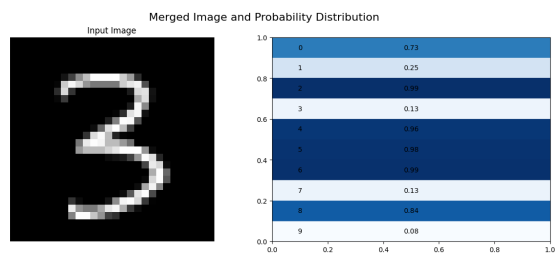
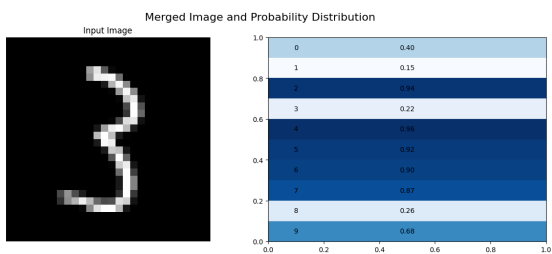
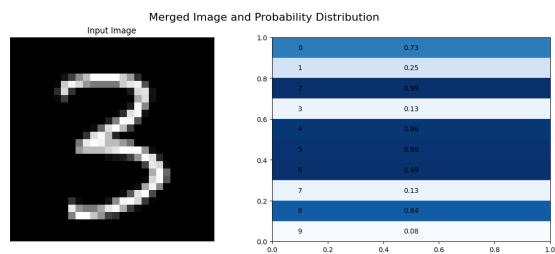
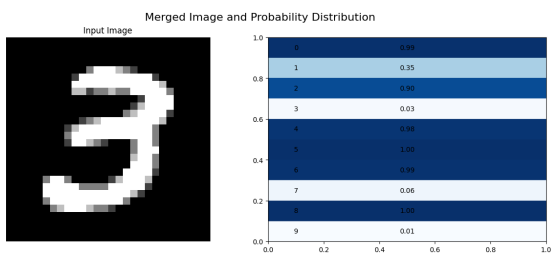




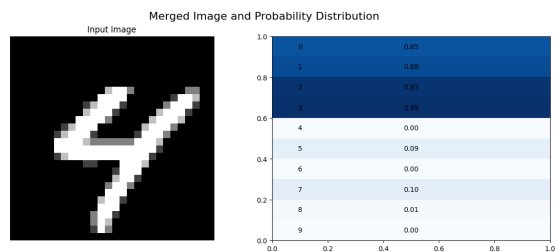
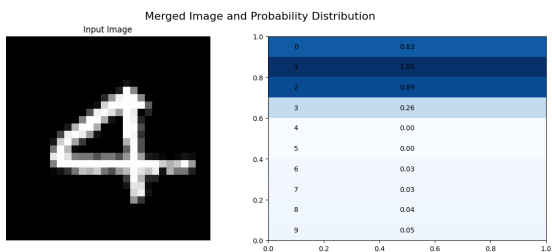
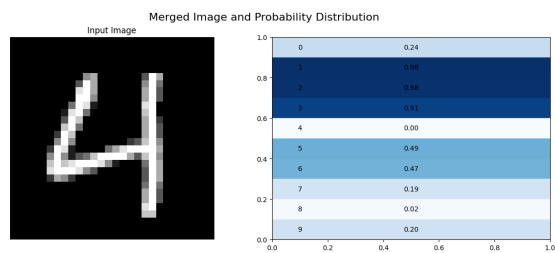
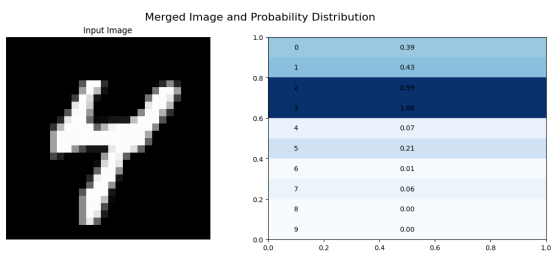
2



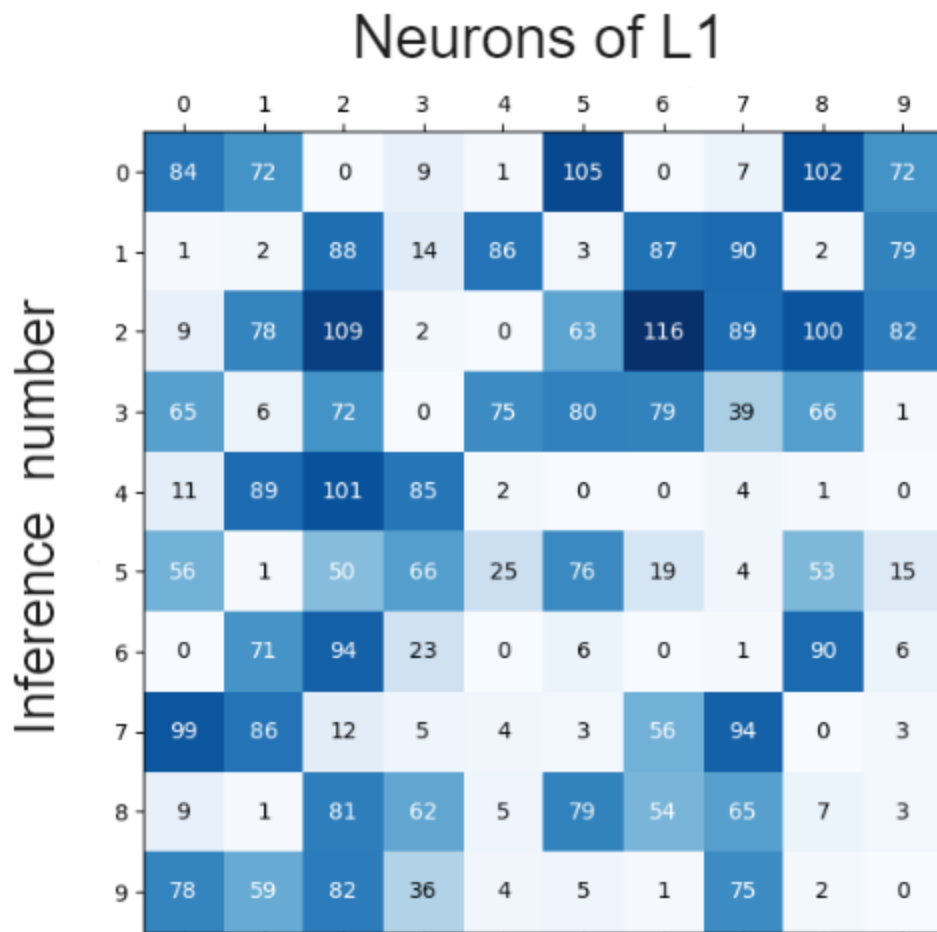
3



4



You Get the point, these were the distribution for the dominance of a specific neuron used in that number



The dominance of a specific neuron in recognizing a number indicates that each neuron in the layer has learned to recognize certain features. This finding supports the hypothesis that neurons in a neural network learn to identify and respond to specific patterns in the input data.