
Classifying Handwritten Digits with a Convolutional Deep Neural Network

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Neural networks excel at solving classification problems. A simple implementation
2 of a deep neural network with convolutional and pooling layers saw success at
3 classifying members of the MNIST dataset: labeled images of handwritten digits.

4 1 Introduction

5 In order to classify the digits in the MNIST dataset, a deep convolutional neural network was used,
6 but what exactly does that mean? The network contained convolutional, pooling, and fully connected
7 layers, each serving either to increase efficiency of training or accuracy of classification. Training was
8 performed on randomly selected batches of 100 digits. 5000 of these batches were used for training.
9 This process was repeated twice more in order to achieve a clearer idea of the networks performance.

10 2 The MNIST Dataset

11 The MNIST dataset consists of 70,000 images of handwritten digits, and labels defining which digit
12 they represent. There are 60,000 in the training set, and 10,000 in the testing set. Each image has a
13 height and width of 28 pixels on a single black/white color channel. See figure 1.

14 3 Network Structure

15 The network consisted of of two convolutional layers, two pooling layers, and one fully connected
16 layer. See figure 2.

17 3.1 Convolutional Layers

18 Each convolutional layer applied a number of convolution filters to 5x5 subregions of the image. This
19 forced the neural network to consider regions of the image in the same way a human might from its
20 start, instead of allowing it to learn that groups of nearby pixels have significance over many training
21 steps.

22 The first convolutional layer applied 16 filters, which multiplied the size of the data by 16. The
23 second applied 32 filters, which similarly increased the size by 32.

24 The second convolutional layer was acting not on pixels, but on the data generated by the first
25 convolutional layer. Essentially, it drew inferences about the input not from relationships between
26 groups of pixels, but from relationships between groups of those groups.

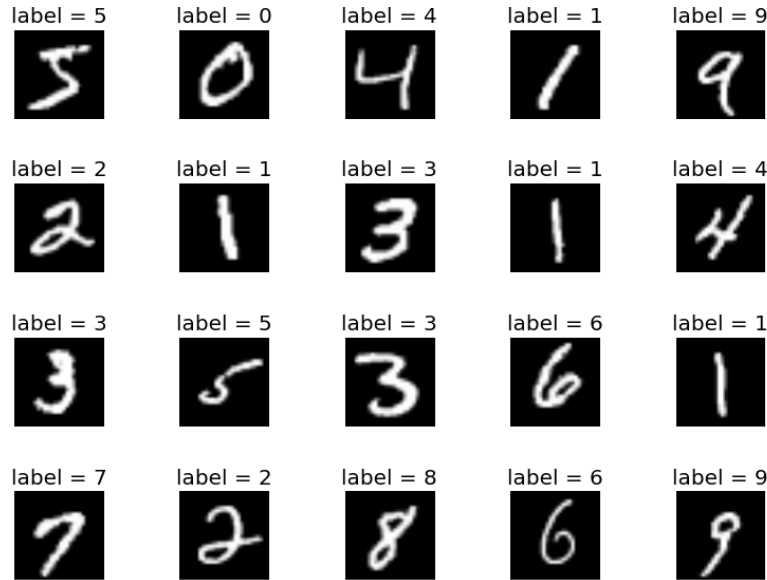


Figure 1: Sample images for the MNIST dataset.

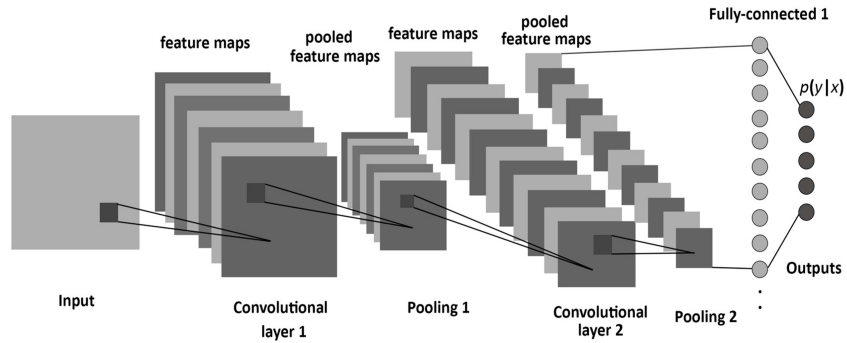


Figure 2: A diagram of the basic structure of the network.

3.2 Pooling Layers

Between the two convolutional layers, the amount of data being processed grows massively, which greatly increases the amount of processing time required for training. This was remedied in part by the inclusion of pooling layers after each convolutional layer, which reduced the resolution of the data by a factor of four. This comes at a cost to quality of classification, but with such a great decrease in processing time, its inclusion is well worth it.

The specific method of pooling used was max pooling over 2x2 regions with a step size of 2. This means that the data entering the layer was broken into discrete 2x2 regions, and each was reduced to its single highest value. Another method, such as taking the average of the values within each region, or taking the most extreme value could have been used. This process can be thought of in much the same way one can think of reducing the resolution on any image. Quality is indisputably lost, but without massive reduction in resolution, the image is still much the same.

3.3 Dense Layer

Finally, after all convolution and pooling is complete, the data passes through one dense, or fully connected layer, which simply draws patterns from the data passed in, which is finally reduced to

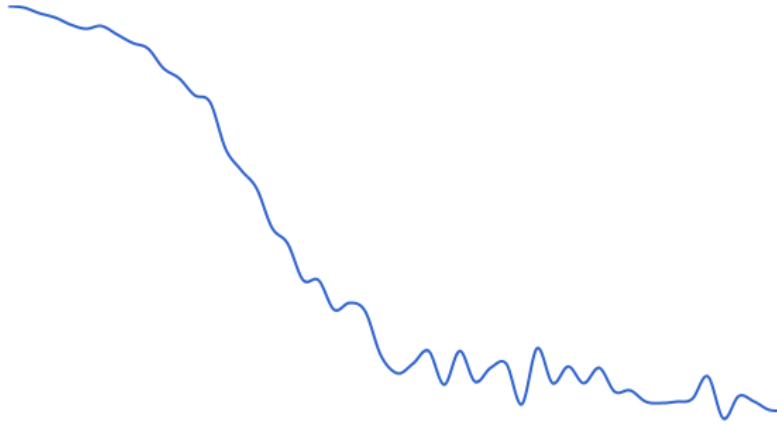


Figure 3: Change in loss over the course of the training.

ten values, one for each possible label. The greatest output value acts as the network’s prediction. The specific dense layer used here contained 512 nodes.

This layer had a dropout rate of 40%, which means that 40% of values, randomly selected, are simply not passed forwards. This encourages a robustness in the network that should allow it to succeed at classification even missing 40% of data, and should help to avoid total reliance on any one feature.

3.4 Activation Functions

Both the convolutional and dense layers use the ReLU activation function:

$$f(x) = \max(0, x)$$

This serves to normalize the outputs of each layer. This process has strong biological motivations and mathematical justifications, and has been shown to enable better training of neural networks (Hahnloser, et al. (2013)).

4 Training

In order to train the network, a batch of randomly selected training samples is selected and run through the network. Outputs are generated by the network for each member of the batch, and loss is calculated. Loss is simply the difference between the desired output and the actual output of the network. In this case, the desired output is a zero for each label except for the correct one, which should have a value of one.

This loss is used to perform gradient descent on the network, by which weight values are changed from the output layer working back towards the beginning, such that, should the same input arrive again, the actual output will be closer to the desired output. The learning rate for this process was selected as 0.001. This value was selected arbitrarily, and could very well be changed for better results. Empirical testing is required to find the best learning rate.

As can be seen in figure 3, reduction in loss started to diminish near the end of training. This indicates that simply performing additional training would see diminishing returns, and that, to improve the accuracy of the network, there would have to be changes in design.

This process was repeated 5000 times, and then the network was reset, and learning was done again from scratch twice more.

68 4.1 Stochastic Gradient Descent

69 Stochastic gradient descent is the process by which weights between nodes within the network are
70 modified to improve the accuracy of the network. Starting at the last layer of the network, and
71 working towards the front, weights are changed based on the following equation:

$$w = w - \eta \nabla Q(w) = w - \eta \sum_{i=1}^n \nabla Q_i(w) / n$$

72 Where w is the weight of a node, η is the learning rate, and $Q(w)$ is the error to be minimized.

73 5 Results

74 Training steps took on average 27.3 seconds each. The average accuracy after training was 93.55%.
75 Compared to some attempts at solving the same problem with neural networks, this was a bit of
76 a bust. As of 2013, some attempts were seeing accuracy as high as 99.79% (Wan, et al. (2013)).
77 Very impressive. With much more training time and number fiddling based on empirical testing, the
78 implimentation described here could possibly have reached similar quality.

79 References

- 80 [1] Wan, Li & Zeiler, Matthew & Zhang, Sixin & LeCun, Yann & Fergus, Rob (2013) *Regular-*
81 *ization of Neural Networks using DropConnect*. Dept. of Computer Science, Courant Institute of
82 Mathematical Science, New York University.
- 83 [2] R Hahnloser & R. Sarpeshkar & M A Mahowald & R.J. Douglas & H.S. Seung (2000) *Digital*
84 *selection and analogue amplification coexist in a cortex-inspired silicon circuit*. *Nature*. 405. pp.
85 947-951.