

Neural Network with Backpropagation

Axel Ind (4563539)

November 6, 2017

Abstract

This document details work to create a neural network that successfully classifies the MNIST database of images into the numbers which each image represents. In testing the Sigmoid, Relu, and Tanh activation functions were used. Testing error of 2.32% is achieved.

1 Activation Functions

Sigmoid

Equation

$$f(x) = \frac{1}{1 + e^{-x}} \quad (1)$$

Derivative

$$f'(x) = f(x)(1 - f(x)) \quad (2)$$

Tanh

Equation

$$g(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \quad (3)$$

Derivative

$$g'(x) = 1 - g(x)^2 \quad (4)$$

Relu

Equation

$$h(x) = \max(0, x) \quad (5)$$

Derivative

$$h'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases} \quad (6)$$

2 Dataset

The data set used is the MNIST dataset of hand-written numbers represented by $28\text{px} \times 28\text{px}$ grayscale images with intensity values in the range $(0, 256)$.

3 Network Design

Inputs

Stochastic Gradient Descent: $784 \times n$ (where n is the batch size).

Gradient Descent: $784 \times N$ (where N is the number of samples).

Hidden Layers

Some number of hidden layers ($|h| > 0$), with some number of weights (W_i) and biases (b_i) at each layer (h_i). Each layer makes use of one of the described activation functions.

Outputs

Some linear or unhot array that uniquely identifies one of 10 classes.

4 Training

All training results are achieved using backpropagation to update weights.

General Observations

- Changing the number of hidden layers did not consistently improve network performance.
- Adding more hidden nodes to the network produced more consistently good results.
- As expected, increasing the number of entries in the data set significantly improved performance.
- Significantly better results were achieved with a batch size of 2, than with any greater number.

4.1 General Training

Without Optimization

Without any form of regularization, the greatest results in the table were above the required accuracy threshold.

Activation	Num Layers	Hidden Nodes	Learning Rate	Epochs	Error
Relu	3	130, 170, 20	0.004	35	2.32%
Relu	4	All 130	0.005	130	2.96%

Table 1: Unoptimized training results for the MNIST Test Set

With Optimization

Optimizing Steps

- Consensus training: three networks were trained simultaneously on the training data and the class which two or more agreed on was used as the prediction when the testing data was used.
- This approach, though it seemed promising, simply served to significantly slow computation time and did not improve accuracy on the most successful unoptimised networks.
- The failure of this technique suggests that the network tends to become trapped in the same small set of local minima.

Activation	Num Layers	Hidden Nodes	Learning Rate	Epochs	Error
Relu	4	All 84	0.05	170	2.92%
Relu	4	All 98	0.11	45	2.85%

Table 2: Optimized training results for the MNIST Test Set

5 Conclusion

The created neural network is able to classify numbers in the MNIST dataset to a very high degree of accuracy. The addition of a consensus optimizing technique did not significantly reduce the observed error rate.

It also appears clear that the sigmoid and tanh functions are consistently outperformed by the Relu function.

Further, the randomness inherent in the neural network algorithm may have obscured configurations which were more optimal