

Neural Network Assignment 6 Report

Sreeraj Rimmalapudi

Neural network architecture - Neural network with 3 layers of 1000, 250, 10 neurons achieved best performance on both the validation and test with test accuracy reaching 97.79% on average. Increasing the number of layers to 4 - 1000, 500, 100, 10 neurons, gave 97.25% accuracy on the test set. On reducing the network to 1 hidden layer with 1000 neurons and 10 output neurons, the performance dropped to 97.53%.

Inputs - The inputs were normalized as suggested in the paper[4]. The network was more stable with these inputs and the accuracy on all sets increased by 0.1%.

Labels - The output labels were encoded as vectors of size 10 with [1 0 0 0 0 0 0 0 0 0] representing 0, [0 1 0 0 0 0 0 0 0 0] representing 1 and so on.

Initial Weights - The weights were uniformly chosen at random between [-0.001, 0.001]. Weights 10 times larger than this caused overflow in exponential term and weights 10 times smaller caused gradient to be very small and hence no update.

Activation Function - The activation function for each neuron was ReLU, defined as $f(x) = \max(0, x)$. It gave nearly 2% improvement on accuracy and faster convergence than hyperbolic tangent function. Softmax function given by $\sigma(z)_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}}$ was used in the final layer to treat the outputs as probabilities where z is the weighted input to the neuron.

Loss Function - Cross entropy or log loss was used as the loss function. The loss function is given as $\frac{1}{n} \sum_{i=1}^n -d_i \log(y_i) - (1-d_i) (\log(1 - y_i))$ where d is the desired output and y is the output of the network.

Mini batch update - Mini batches were used[3]. The batch size was set to 200.

Learning rate - Learning rate was set to 0.001. It was constant for each of the weights throughout the entire training process. Increasing the learning rate to 0.01 caused the network to sometimes diverge with no increase in accuracy. There was a plan initially to decay learning rate after certain epochs[3] but it performed well even without the decay.

Adaptive Momentum Estimation (Adam) Update - Adam update was used. The parameters for the update were set as recommended in the paper[1].

Epochs - Epochs were initially set to 100. On closer inspection of the loss function, increasing the epoch would have had no effect. It was decreased to 25 the point after which no drop in error rate was observed.

Dropout - Dropout[2] was used with drop probability set to 0.5 for each neuron. It was not used in the input and output layers.

Regularization - L2 regularization or weight decay was used to limit the weights from growing large. The regularization parameter was set to 0.001. Lower values caused problems with logarithm.

References

- [1] Adam: a Method for Stochastic Optimization. Kingma, D. P., Ba, J. L. (2015). International Conference on Learning Representations, 1–13.
- [2] Dropout: A simple way to prevent neural networks from overfitting. N. Srivastava, G.E. Hinton, A. Krizhevsky, I. Sutskever, and R. R. Salakhutdinov. The Journal of Machine Learning Research, 15(1):1929–1958, 2014.
- [3] Practical Recommendations for Gradient-Based Training of Deep Architectures Yoshua Bengio. In Neural Networks: Tricks of the Trade, pg. 437–478. Springer, 2012
- [4] Efficient BackProp, LeCun, Y., Bottou, L., Orr, G. B., and Muller, K., Tricks of the Trade (1998a).

Description of Algorithm - The algorithm takes a modular approach and performs the following operation -

- 1) Perform Forward Propagation and store values for backprop
- 2) Get loss
- 3) Perform Back Propagation and get the gradients
- 4) Update the weights using the gradients

The forward and backward propagation is performed layer by layer in forward and backward step function respectively.

The algorithm mainly uses mini batch learning with fixed batch size. It utilizes regularization, dropout, early stopping and weights updating using Adam method techniques to improve performance.

Avg. Test accuracy - 97.79%

Best Test accuracy - 97.9%

Time to train - Approximately 20 mins

Algorithm 1 Neural Network Mnist Classification

```
1: function MAIN
2:   Get the inputs and divide the data into training, validation and testing set
3:   Initialize weights, biases and set the hyper parameters
4:   set patience = 0
5:   for  $epoch = 1$  to 25 do
6:     Divide the training data into batches each of size 200 ▷ Mini batches
7:     for each batch  $b_i$  do ▷ Training
8:       Get batch output of  $b_i$  by performing forward propagation
9:       Calculate cross entropy loss
10:      Back propagate the gradients
11:      Update the weights using Adam update
12:      Forward propagate the entire training sample and validation sample and get out-
      put  $y$  and  $v$ 
13:      Calculate training, validation loss and accuracy ▷ Validation
14:      if validation accuracy increases then
15:        Store the parameters
16:      else
17:        Decrease the patience
18:      if patience = 0 then
19:        End training ▷ Early stopping

20: function FORWARD PROPAGATION
21:   for each layer  $l$  do
22:     Let  $x$  be the input of size  $N \times D$ , where  $N$  is the batch size and  $D$  is the number of
     neurons in previous layer
23:     Let  $w$  be the weights of size  $D \times H$  ( $H$  is the number of hidden neurons of layer  $l$ )
24:     Let  $b$  be the bias of size  $1 \times H$ 
25:     if  $l$  is not the last layer then
26:       Perform Forward step with  $x, w$  and  $b$ 
27:     else
28:       Perform softmax transformation with  $x, w$  and  $b$ 
29:       Store values for backward step

30: function FORWARD STEP( $x, w, b$ )
31:   Perform affine transformation,  $out = xw + b$ 
32:   Use ReLU transformation,  $reluout = \max(0, x)$ 
33:   With probability of 0.5 drop the output, ie, make it zero
34:   Store  $x, w, out$ 

35: function SOFTMAX TRANSFORMATION( $x, w, b$ )
36:   Perform affine transformation,  $z = xw + b$ 
37:   for  $j = 1..10$  do
38:      $out_j = \frac{e^{z_j}}{\sum_{k=1}^{10} e^{z_k}}$ 
39:   out is the network output of size  $N \times 10$  giving probabilities of each class
```

Algorithm 1 Neural Network Mnist Classification (continued)

```
40: function BACKWARD PROPAGATION
41:   for each layer l in reversed order do
42:     if l is the last layer then
43:       Perform backprop of softmax function
44:     else
45:       Perform backstep to calculate gradient

46: function SOFTMAX BACKPROP
47:   Let  $t$  be the desired label with one hot encoding
48:   Let  $out$  be the network output
49:    $delta = out_i - t_i$  ▷ Softmax derivative
50:    $dw = x^T delta$  ▷ Derivative of loss wrt weight w
51:   for each hidden neuron k do
52:      $db_k = \sum_{i=1}^n delta_{i,k}$  ▷ Derivative of loss wrt bias k
53:   delta for previous layer =  $delta w^T$ 

54: function BACKWARD STEP
55:   Get the stored values from forward propagation
56:   if neuron was dropped then
57:     Gradients of the weights connected to that neuron is 0
58:   else
59:     If the input value before applying ReLU in forward prop was less than zero, kill
        the gradient at that location ▷ ReLU derivative
60:      $dw = x^T delta$  ▷ Derivative of loss wrt weight w
61:     for each hidden neuron k do
62:        $db_k = \sum_{i=1}^n delta_{i,k}$  ▷ Derivative of loss wrt bias k
63:     delta for previous layer =  $delta w^T$ 

64: function CROSS ENTROPY LOSS
65:   Let  $out$  of size NX10 denote the network output representing probability of each class
66:   Let  $t_i$  denoted the desired label for each input of the batch
67:   Let  $W$  denote the network weights
68:    $loss = - \sum_{i=1}^N \log(out_{i t_i}) + \frac{\lambda}{2} * \|W\|^2$ 
```
