# How we imagined the working of the brain - 90s Edition.

Ajitesh Gupta, Jean Choi

**Abstract**.  In this assignment we learned and implemented multilayer neural networks and learned about their working on the MNIST dataset.

For part 1, where we used one hidden layer, we achieved **92.31%** for using batch process. **400** hidden layer units were used throughout the experiments. $50,000$ training images, $10,000$ validation images and $10,000$ testing images were used.

For part 2, the accuracies improved progressively as **92.31%**, **94.15%**, **94.45%**, **97.7%** as we use mini-batch and weight normalization, tanh activation, regulated weight initialization, momentum respectively.

For part 3, we get accuracies of **97.72%**, **85.43%** for networks with **800** and **5** neurons respectively. It also gets an accuracy of **97.7%** for a double hidden layer network with **300** and **270** neurons each.

## 1   Introduction

### 1.1   Classification

Here we use a neural network with a single hidden layer to perform classification on the MNIST dataset. We have varied the size of the hidden layer between 300 and 800 neurons and report the results for the same. For this part we used sigmoid activation in the hidden layers and softmax activation in the output layer. Since nothing is specified and we have to compare with mini-batch performance, here we train using stochastic gradient descent.

The implementation is based on the following notations and equations:

Denote $J$ as the number of units in hidden layer, $C$ as the number of classes, and $d$ as the dimension of the input, $X$.

For output layer, the output is given by

$$y_k = \frac{e^{sum_k}}{\sum_{k'} e^{sum_{k'}}} \tag{1}$$

and

$$E = -\sum_{l=1}^{C} t_l \log y_l \tag{2}$$

For hidden layer, the activation value of each unit is given by

$$z_j = \frac{1}{1 + e^{sum_j}} \tag{3}$$

Then

$$\delta_k = t_k - y_k \tag{4}$$

and

$$\delta_j = (1 - z_j) z_j \sum_{k=1}^{C} \delta_k w_{jk} \tag{5}$$

The weight update rules are:

$$w_{jk} = w_{jk} + \eta \delta_k z_j \tag{6}$$

$$w_{ij} = w_{ij} + \eta \delta_j x_i \tag{7}$$

## 1.2 "Tricks of the Trade"

Here we again conduct experiments using a standard 2 layer neural network architecture. But this time training is done in a mini-batch manner, in order to provide a speed boost. Also various optimizations are applied in the form of regulating the initial weights, using tanh as the activation function in the hidden layer, averaging the change in the weights over number of samples.

Since we use tanh as the activation function now, the corresponding $\delta_j$ is as follows :

$$\delta_j = 1.154 * (1 - \tanh^2(\frac{2}{3} z_j)) \sum_{k=1}^{C} \delta_k w_{jk} \tag{8}$$

## 1.3 Experiments with network topology

Here we saw the change in performance of the network by conducting 2 experiments with the structure of the network itself.

### 1.3.1 Doubling neurons

We double the number of neurons in the hidden layer i.e. double the number of parameters to be learned. We want to see whether increasing the number of parameters learned by the network increases the accuracy of the network and also what are its effects on the speed of the learning process itself.

### 1.3.2 Doubling hidden layers

We double the number of hidden layers while keeping similar number of parameters. We again want to see whether the increased complexity obtained by adding a second hidden layer makes the network learn to predict better and also again what is its effect on the learning speed. Using assumptions used in previous derivaitons and using $i, j, k, l$ to denote input, hidden layer 1, hidden layer 2 and output layer and $H2$ to denote number of units in hidden layer 2 we get :

$$y_l = \frac{e^{sum_l}}{\sum_{l'} e^{sum_{l'}}} \tag{9}$$

$$z_k = 1.732 * \tanh\left(\frac{2}{3} sum_k\right) \tag{10}$$

$$z_j = 1.732 * \tanh\left(\frac{2}{3} sum_j\right) \tag{11}$$

then

$$\delta_l = t_l - y_l \tag{12}$$

$$\delta_k = 1.154 * (1 - \tanh^2(\frac{2}{3} z_k)) \sum_{l=1}^{C} \delta_l w_{kl} \tag{13}$$

$$\delta_j = 1.154 * (1 - \tanh^2(\frac{2}{3} z_j)) \sum_{k=1}^{H2} \delta_k w_{jk} \tag{14}$$

The weight update rules are:

$$w_{kl} = w_{kl} + \eta \delta_l z_k \tag{15}$$

$$w_{jk} = w_{jk} + \eta \delta_k x_j \tag{16}$$

$$w_{ij} = w_{ij} + \eta \delta_j x_i \tag{17}$$

# 2 Methods

## 2.1 Classification

The images were read in and normalized to a range of [0,1] by element-wise division by 255. Then they were mean subtracted in order to center the data.

As mentioned before, logistic activation was used for hidden layer while the output layer used softmax activation.

The parameters used in gradient descent such as number of epochs, $m$, initial learning rate, $nu$, and metaparameter, $T$, were given values based on experiments.

The accuracies were continuously monitored at each iteration and each iteration was timed in order to calculate the training time.

The following was the training procedure :

1. Normalize each sample by dividing by 255.

2. Center each image by subtracting mean.

3. Initialize weights by sampling randomly from a normal distribution with mean 0 and standard deviation of 1.

4. In each epoch for each sample do a forward pass and store the error.

5. At the end of the epoch, update the weights using the accumulated error.

6. Calculate the classification accuracy once the weights are updated.

7. Repeat till convergence.

Also to evaluate the correctness of the backprop we compare backprop updates with the ones obtained using the equation given in the question. For this we iterated through each weight in the network and did a forward pass with its original value say x, then with value x+$\epsilon$ and then with value x-$\epsilon$. We used the normal weight to obtain backprop update and the results from the modified weights to calculate the numerical gradient. For each weight we keep summing the difference between the backprop update and the numerical gradient. Finally we average the sum over the number of weights to get the average difference.

## 2.2 "Tricks of the Trade"

In this part we incorporated various optimizations based on the reading of the paper by Yann LeCun.

- Mini Batch
  Since batch learning is not so accurate and stochastic learning takes time to iterate over the whole data, we use an in-between solution of mini-batches.

This helps in accurate and speedy learning. For the purpose of this assignment, we used mini-batches of size 125. We also averaged the weight update over the number samples in the batch.

- tanh Activation
  We used tanh activation in the hidden layer. This lead to a modified update equation as $\frac{d}{dx}\tanh(x)$ is $1 - \tanh^2(x)$.

- (Weight initialization)
  For each unit in the hidden layer and the output layer, we sample the initial weights from a normal distribution with mean of 0 and standard deviation of $\frac{1}{\sqrt{fan-in}}$ where fan-in stands for the number of inputs to that unit.

- Momentum
  Momentum term $\alpha$ used similar to its meaning in the physics world. Momentum involves using part of the previous update along with current update at any given iteration. It helps to avoid getting stuck in plateaus.

$$w(t+1) = w(t) + \alpha\Delta w(t-1) + \eta\Delta w(t) \tag{18}$$

So the final training procedure is as follows :

1. Normalize each sample by dividing by 255.

2. Center each image by subtracting mean.

3. Shuffle the examples.

4. Initialize weights by sampling randomly from a normal distribution with mean 0 and standard deviation of $\frac{1}{\sqrt{fan-in}}$.

5. For each mini-batch do a forward pass and calculate the error.

6. Normalize the error using the mini-batch size.

7. Update the weights using the calculated error and the also the previous error weighed by the momentum.

8. Calculate the classification accuracy after each epoch.

9. Repeat till convergence.

## 2.3 Experiments with Network Topology

- Doubling Neurons
  In this experiment we see the effect of doubling the number of neurons in the hidden layer of the network. In the previous networks we were using 400 neurons in the hidden layer so in this case we use 800 neurons. We

measure the time taken in each epoch along with the classification accuracy.

- Doubling the number of hidden layers
  Unlike normal 2 layer neural networks here we add another hidden layer. Again since we were using 400 neurons in the hidden layer leading to 318000 parameters, here we use 300 and 270 neurons in the two hidden layers leading to a very similar 319200 parameters. Again we measured both accuracy and time taken in each epoch.

  Owing to the added layer the new update equations are as follows :

## 3 Results and Discussions

### 3.1 Classification

For the full batch training we used $\eta = 0.1$ was used. The graph below shows the training accuracies for training, validation and testing sets as a function of number of epochs. The final accuracy achieved after 240 epochs was 92.31%. Also each epoch took around 1.5 seconds.
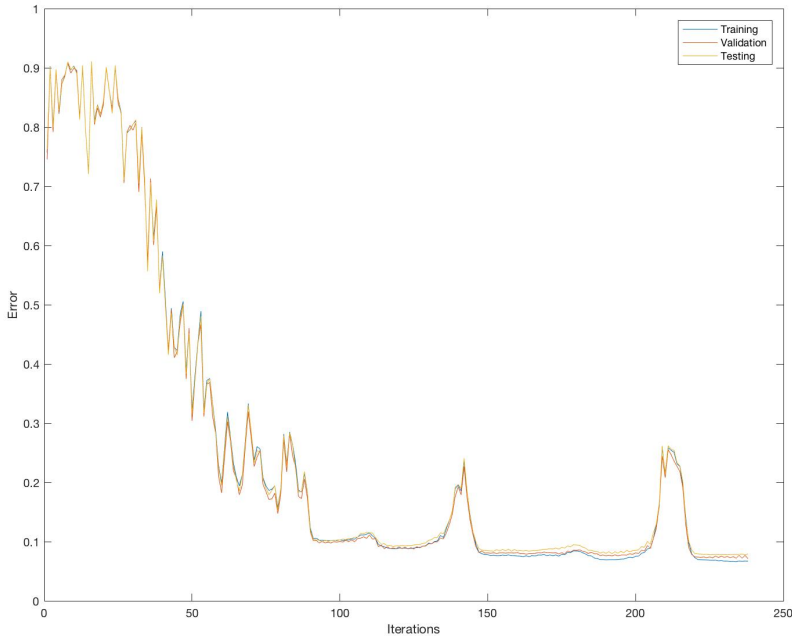


Figure 1: Standard full batch learning

As we can see the learning is very erratic owing to no normalization of the weight updates. Also the accuracy is only 5% better than the simple softmax regressor we used in the previous assignment.

Also regarding the correctness of the backprop update, we see that the average difference comes out to nearly zero with $\epsilon = 0.1$ and $\epsilon = 0.01$ and so on (ex. with $\epsilon = 0.01$ sum of differences comes out to 0.681 and averaging that over $3 * 10^6$ gives $2 * 10^{-7}$). Thus it is within the big-O of $\epsilon^2$. Hence our backprop implementation is correct.

## 3.2 "Tricks of the Trade"

(a) **Mini-batches and update normalization**

We used mini-batches of size 125 for convenience instead of 128. The error vs epoch graph for this method is as shown below.
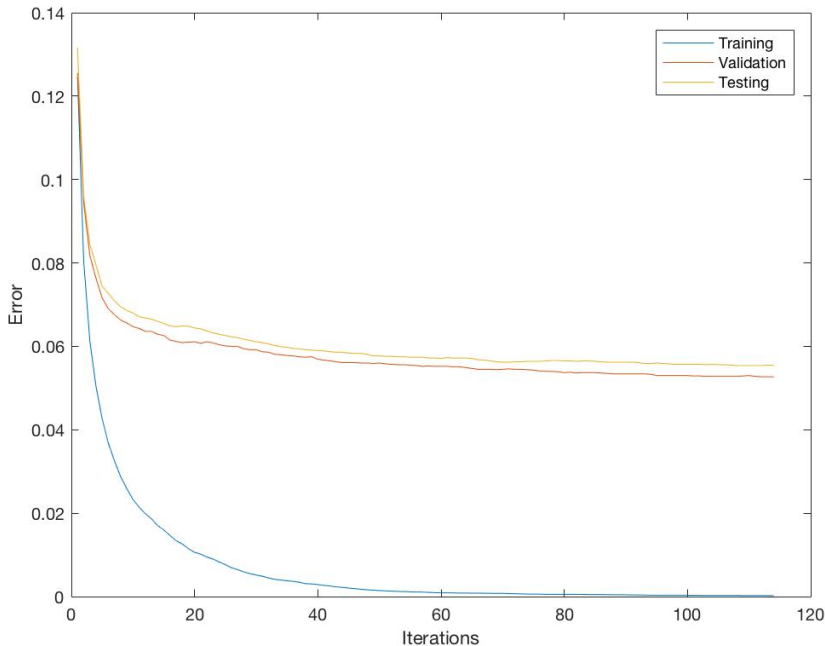


Figure 2: Mini batch learning

We can clearly see a marked improvement in the learning process. It is much more smoother than the previous batch method. We converge to 94.15% accuracy in just around 120 epochs. Also each epoch only takes 1.33 seconds. Thus mini-batches produce more accurate results than full batch while also being faster.

This is because the mini-batches retain qualities of both stochastic and batch gradient descent. The batch size being small makes updates faster and hence faster learning. Batches also means the updates are meaningful and not affected by random samples. Shuffling the samples also means that with high probability all successive samples belong to different classes.

(b) tanh **Activation**

Retaining the above modifications we further replaced the logistic activations in the hidden layer with $tanh$ activation function. Below we see the error vs epochs graph of the same.
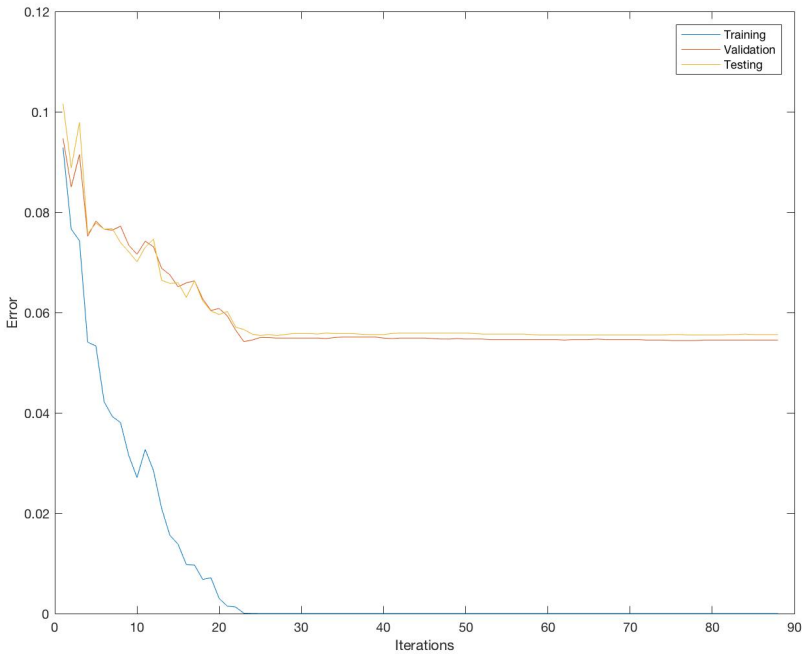


Figure 3: Learning using tanh activation

Again we can see the learning much stabler than the original full batch learning. It achieves a final accuracy of 94.45% and such the tanh function is better than the logistic activation function. Also although each epoch takes a little longer at 1.91 seconds seconds per epoch, it converges in around 80 epochs. As such it is still faster than the batch gradient descent.

(c) **Regulating initial weights**

Further we initialized the weights using a normal distribution of mean 0

and standard deviation equal to the inverse square root of the number of inputs in the corresponding neuron. Below is the corresponding error vs epoch graph.



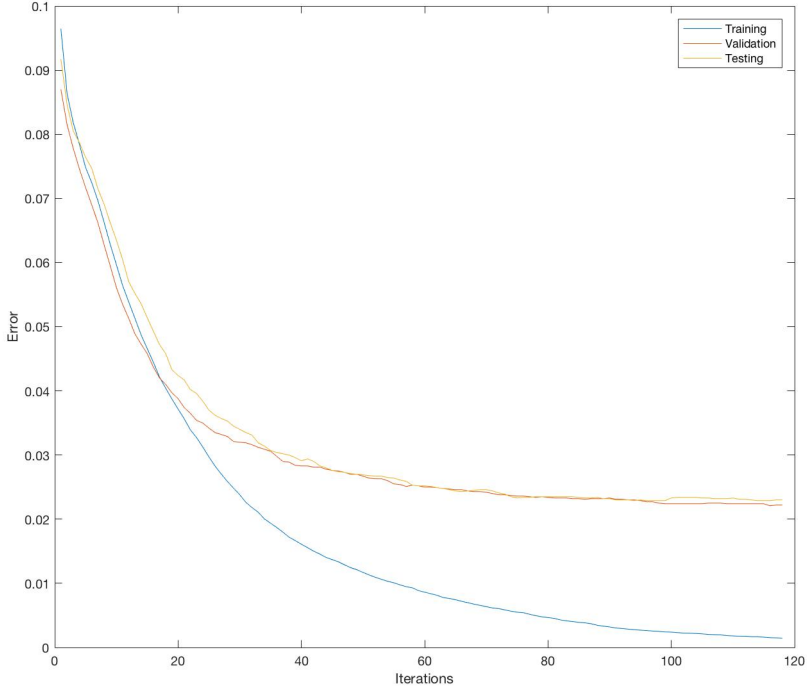Figure 4: Effect of regulating initial weights

This method reaches a final accuracy of 97.7% in around 120 epochs and each epoch takes around 1.8 seconds. Thus the regulating the initial weights drastically improves the accuracy while also maintaining training speed.

(d) **Momentum**

For the experiments we used a momentum of 0.7. Below we see the graph of error vs epochs while using momentum.
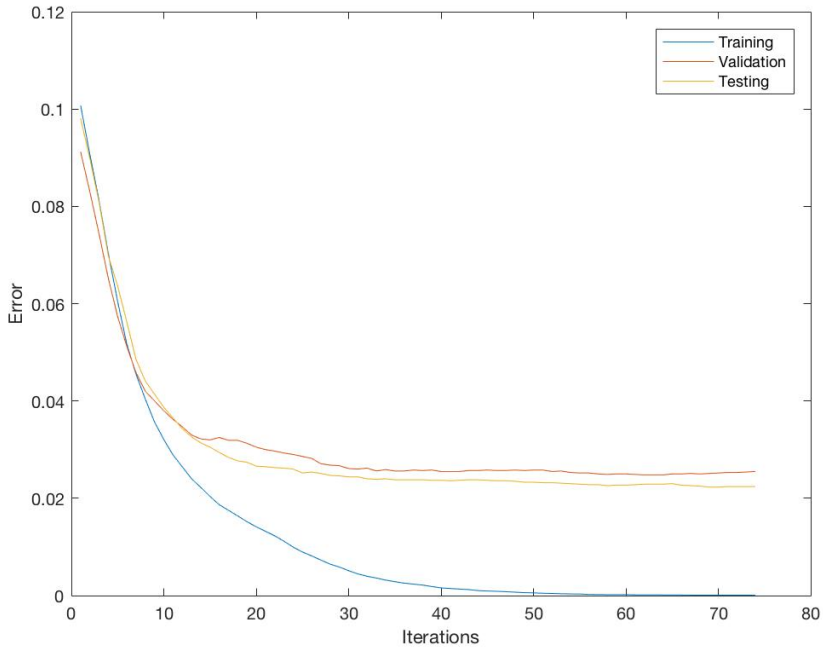
Figure 5: Learning with momentum

We see as expected the graph has a steeper slope meaning that the error goes down much faster. It reaches an accuracy of 97.77% in a much lesser 80 epochs. Thus again it generates better accuracy in faster time. The time taken per iteration is also around 2.07 seconds.

## 3.3 Experiments with network topology

### (a) Doubling number of neurons in hidden layer

From the error vs epochs graph we can see that increasing the number of neurons did not have a profound effect on the training. The final training accuracy was again around 97.72% after around 100 iterations. Each epoch took around 3.23 seconds to process. As such it only leads to increased training time without much change in accuracy.

Figure 6: Learning with double the neurons

In another case we took 5 neurons to check effect of too low number of neurons. In this case each epoch was much faster at 0.45 seconds per epoch, but the accuracy took a hit. The accuracy achieved in this case was just 85.43%.

Thus in order to maintain a balance between accuracy and speed, we need to pick an optimum network size which is neither too big nor too small. Big networks lead to larger training times whereas smaller ones give lower accuracy.

Figure 7: Learning with 5 neurons

(b) **Doubling hidden layers**

As we can see from the figure below, doubling the number of hidden layers again lead to similar accuracies of around 97.7% as achieved before. But it also takes around 150 epochs and each epoch takes around 3.17 seconds. Thus this is also unnecessary as although it might lead to a slightly better accuracy, it'll take much longer time to reach that. Only very complex classification problems can as such make use of the added complexity provided by an added hidden layer.

Figure 8: Learning with 2 hidden layers

## 4   Summary

(a) We see that neural networks work much better than the regression techniques we learnt in the last assignment.

(b) In particular we see that stochastic gradient descent leads to better learning than batch learning. A mid-way between the two can be in the form of mini-batches.

(c) Batch learning leads to low accuracy and the learning itself is also erratic.

(d) Mini-batches lead to a balance between fast learning and accurate learning.

(e) Normalizing the weight updates leads to faster learning and more accurate networks.

(f) A lot also depends on how we initialize our networks. Initializing weights to values, regulated by the number of their inputs leads to faster and more accurate learning.

(g) $1.732\tanh\left(\frac{2}{3}x\right)$ acts as a better activation function than regular logistic function as it is symmetrical around and centered at 0 whereas, the logistic function is always positive.

(h) Using momentum further speeds up the learning process and helps avoid plateaus.

(i) Needlessly increasing the number of neurons will just add to the overall training time whereas keeping the number too low will lead to lower learning.

(j) Increasing the number of layers might increase the accuracy in cases where the added level of complexity is required, but that is rare. So adding the extra layer again adds to training time while the accuracy remains nearly same as the network with single hidden layer.

## 5 Contributions

I, Jean Choi, did 3.
I, Ajitesh Gupta, did 4 and 5.

## Appendix

```matlab
% Change the filenames if you've saved the files under different
    names
% On some platforms, the files might be saved as
% train-images.idx3-ubyte / train-labels.idx1-ubyte
images_train = loadMNISTImages('train-images-idx3-ubyte')';
labels_train = loadMNISTLabels('train-labels-idx1-ubyte');
y_train = labels_train;
labels_train = onehot(labels_train,10);

for it=1:60000
    m = mean(images_train(it,:));
    images_train(it,:) = images_train(it,:)-m;
end

images_test = loadMNISTImages('t10k-images-idx3-ubyte')';
labels_test = loadMNISTLabels('t10k-labels-idx1-ubyte');
y_test = labels_test;
labels_test = onehot(labels_test,10);


for it=1:10000
    m = mean(images_test(it,:));
    images_test(it,:) = images_test(it,:)-m;
```

```matlab
23 end
24
25 ntrain = 50000;
26 nval = 10000;
27 ntest = 10000;
28
29 train_x = [ones(ntrain,1),images_train(1:ntrain,:)];
30 train_y = labels_train(1:ntrain,:)';
31
32 val_x = [ones(nval,1),images_train(ntrain+1:ntrain+nval,:)];
33 val_y = labels_train(ntrain+1:ntrain+nval,:)';
34
35 test_x = [ones(ntest,1),images_test(1:ntest,:)];
36 test_y = labels_test(1:ntest,:)';
37
38 % Setting hyperparameters
39 n_iters = 200;
40 nh = 400;
41 ni = 784;
42 no = 10;
43 eta0 = 0.3;
44 eta = eta0;
45 T = 10000;
46 mom = 0.7;
47 trainsize = 125;
48
49 % Initialization considering fan-in
50 w1 = normrnd(0,1/sqrt(ni),[ni+1,nh]);
51 w1(end,:) = normrnd(0,1,[1,nh]);
52 w2 = normrnd(0,1/sqrt(nh),[nh+1,no]);
53 w2(end,:) = normrnd(0,1,[1,no]);
54
55 % Initialization without fan-in
56 % w1 = normrnd(0,1,[ni+1,nh]);
57 % w1(end,:) = normrnd(0,1,[1,nh]);
58 % w2 = normrnd(0,1,[nh+1,no]);
59 % w2(end,:) = normrnd(0,1,[1,no]);
60
61 % Shuffling data
62 r = randi([1 ntrain],1,ntrain);
63 train_x(1:ntrain,:) = train_x(r,:);
64 train_y(:,1:ntrain) = train_y(:,r);
65
66 tr_ac = zeros([1,n_iters]);
67 val_ac = zeros([1,n_iters]);
68 te_ac = zeros([1,n_iters]);
69 time = 0.0;
70 for i=1:n_iters
71     tic;
72     % Running mini-batch
```

```
73      for k=1:ntrain/trainsize
74          updatek = zeros([nh+1,no]);
75          updatej = zeros([ni+1,nh]);
76          start = (k-1)*trainsize;
77          %forward propagate
78          netj = train_x(start+1:start+trainsize,:)*w1;
79          fnetj = [logistic(netj) ones([trainsize,1])];
80          %fnetj = [1.732.*tanh(2/3.*netj) ones([trainsize,1])];
81          netk = fnetj*w2;
82          fnetk = softmax(netk);
83
84          %error calc
85          dk = (train_y(:,start+1:start+trainsize)-fnetk');
86          u1 = updatek;
87          updatek = updatek + (fnetj'*dk');
88          dj = (dk'*w2').*fnetj.*(1-fnetj);
89          %dj = (dk'*w2').*1.154.*(1-tanh(2/3.*fnetj).*tanh(2/3.*
        fnetj));
90          u2 = updatej;
91          updatej = updatej + train_x(start+1:start+trainsize,:)'*dj
        (:,1:end-1);
92
93          w2 = w2 + mom.*u1 + eta.*updatek./trainsize;
94          w1 = w1 + mom.*u2 + eta.*updatej./trainsize;
95          eta = eta0/(1+i/T);
96      end
97      time = time+toc;
98      % Checking accuracies
99      netj = train_x(1:ntrain,:)*w1;
100     fnetj = [logistic(netj) ones([ntrain,1])];
101     %fnetj = [1.732.*tanh(2/3.*netj) ones([ntrain,1])];
102     netk = fnetj*w2;
103     fnetk = softmax(netk);
104     [~, Idx] = max(fnetk,[],2);
105     [~, Idx2] = max(train_y(:,1:ntrain),[],1);
106     %for x=1:nval
107     %    fprintf('%d %d %d\n',Idx(x),Idx2(x),y_train(ntrain+x));
108     %end
109     count0 = nnz(Idx==Idx2');
110
111     netj = val_x(1:nval,:)*w1;
112     fnetj = [logistic(netj) ones([nval,1])];
113     %fnetj = [1.732.*tanh(2/3.*netj) ones([nval,1])];
114     netk = fnetj*w2;
115     fnetk = softmax(netk);
116     [~, Idx] = max(fnetk,[],2);
117     [~, Idx2] = max(val_y(:,1:nval),[],1);
118     %for x=1:nval
119     %    fprintf('%d %d %d\n',Idx(x),Idx2(x),y_train(ntrain+x));
120     %end
```

```
121      count1  =  nnz ( Idx == Idx2 ' ) ;
122
123      netj  =  test_x ( 1 : ntest , : ) * w1 ;
124      fnetj  =  [ logistic ( netj )  ones ( [ ntest , 1 ] ) ] ;
125      %fnetj  =  [ 1.732.* tanh ( 2 / 3.* netj )  ones ( [ ntest , 1 ] ) ] ;
126      netk  =  fnetj * w2 ;
127      fnetk  =  softmax ( netk ) ;
128      [ ~ ,  Idx ]  =  max ( fnetk , [ ] , 2 ) ;
129      [ ~ ,  Idx2 ]  =  max ( test_y ( : , 1 : ntest ) , [ ] , 1 ) ;
130      %fprintf ( '%d %d %d\'n , Idx , Idx2 , y_test ) ;
131      count2  =  nnz ( Idx == Idx2 ' ) ;
132      tr_ac ( i )  =  count0 / ntrain ;
133      val_ac ( i )  =  count1 / nval ;
134      te_ac ( i )  =  count2 / ntest ;
135      fprintf ( 'Iteration : %d Train Accuracy : %0.4 f Val Accuracy :
         %0.4 f Test Accuracy : %0.4 f Time : %0.2 f \n ' , i , count0 / ntrain ,
         count1 / nval , count2 / ntest , time / i ) ;
136 end
137
138 function  [ g ]  =  logistic ( z )
139 par_res = 1 + exp ( − z ) ;
140 g = 1./ par_res ;
141 end
142
143 function  [ g ]  =  tanhh ( z )
144 g  =  2 / 3.* tanh ( z ) ;
145 end
146
147 function  [ softmaxA ]  =  softmax ( A )
148 dim  =  2 ;
149 s  =  ones ( 1 ,  ndims ( A ) ) ;
150 s ( dim )  =  size ( A ,  dim ) ;
151 maxA  =  max ( A ,  [ ] ,  dim ) ;
152 expA  =  exp ( A − repmat ( maxA ,  s ) ) ;
153 softmaxA  =  expA  ./  repmat ( sum ( expA , dim ) ,  s ) ;
154 end
155
156 function  [ oh ]  =  onehot ( v , c )
157 oh  =  zeros ( length ( v ) ,  c ) ;
158 for  i = 1 : length ( v )
159      oh ( i , v ( i ) + 1 )  =  1 ;
160 end
161 end
```

Listing 1: MATLAB code to run a 2 layer neural network

```
1 % Change the filenames if you've saved the files under different
     names
2 % On some platforms , the files might be saved as
3 % train − images . idx3 − ubyte  /  train − labels . idx1 − ubyte
```

```
4  images_train = loadMNISTImages ('train-images-idx3-ubyte ') ';
5  labels_train = loadMNISTLabels ('train-labels-idx1-ubyte ');
6  y_train = labels_train ;
7  labels_train = onehot (labels_train ,10) ;
8
9  for it =1:60000
10     m = mean (images_train (it ,:) );
11     images_train (it ,:) = images_train (it ,:)-m;
12  end
13
14  images_test = loadMNISTImages ('t10k-images-idx3-ubyte ') ';
15  labels_test = loadMNISTLabels ('t10k-labels-idx1-ubyte ');
16  y_test = labels_test ;
17  labels_test = onehot (labels_test ,10) ;
18
19
20  for it =1:10000
21     m = mean (images_test (it ,:) );
22     images_test (it ,:) = images_test (it ,:)-m;
23  end
24
25  ntrain = 50000;
26  nval = 10000;
27  ntest = 10000;
28
29  train_x = [ ones (ntrain ,1) ,images_train (1: ntrain ,:) ];
30  train_y = labels_train (1: ntrain ,:) ';
31
32  val_x = [ ones (nval ,1) ,images_train (ntrain +1: ntrain +nval ,:) ];
33  val_y = labels_train (ntrain +1: ntrain +nval ,:) ';
34
35  test_x = [ ones (ntest ,1) ,images_test (1: ntest ,:) ];
36  test_y = labels_test (1: ntest ,:) ';
37
38  % Setting hyperparameters
39  n_iters = 20000;
40  nh1 = 300;
41  nh2 = 270;
42  ni = 784;
43  no = 10;
44  eta0 = 0.3;
45  eta = eta0 ;
46  T = 1000000000;
47  trainsize = 50;
48  mom = 0.7;
49
50  % Initialization considering fan-in
51  w1 = normrnd (0 ,1/ sqrt (ni) ,[ ni +1,nh1 ]);
52  w1(end ,:) = normrnd (0 ,1 ,[1 ,nh1 ]);
53  w2 = normrnd (0 ,1/ sqrt (nh1) ,[ nh1 +1,nh2 ]);
```

```matlab
54  w2(end ,:)  =  normrnd(0 ,1 ,[1 ,nh2]);
55  w3  =  normrnd(0 ,1/ sqrt (nh2) ,[nh2+1,no]);
56  w3(end ,:)  =  normrnd(0 ,1 ,[1 ,no]);
57
58  tr_ac  =  zeros([1 ,n_iters]);
59  val_ac  =  zeros([1 ,n_iters]);
60  te_ac  =  zeros([1 ,n_iters]);
61
62  for  i =1:n_iters
63      for  k=1:ntrain / trainsize
64          updatek  =  zeros([nh2+1,no]);
65          updatej  =  zeros([nh1+1,nh2]);
66          updatei  =  zeros([ni+1,nh1]);
67          start  =  (k−1)* trainsize ;
68
69          %forward  propagate
70          neti  =  train_x(start +1: start + trainsize ,:) *w1;
71          %fneti  =  [sigmoid(neti)  ones([trainsize ,1])];
72          fneti  =  [tanh(1.732.*2/3.* neti)  ones([trainsize ,1])];
73          netj  =  fneti *w2;
74          %fnetj  =  [sigmoid(netj)  ones([trainsize ,1])];
75          fnetj  =  [1.732.* tanh(2/3.* netj)  ones([trainsize ,1])];
76          netk  =  fnetj *w3;
77          fnetk  =  softmax(netk);
78
79          %error  calc
80          dk  =  (train_y (:, start +1: start + trainsize )−fnetk ');
81          u1  =  mom.* updatek ;
82          updatek  =  updatek  +  (fnetj '* dk ');
83          %dj  =  (dk '* w3 ') .* fnetj .*(1− fnetj );
84          dj  =  (dk '* w3 ') .*(1− tanh( fnetj ) .* tanh( fnetj ));
85          u2  =  mom.* updatej ;
86          updatej  =  updatej  +  fneti '* dj (:, 1: end −1);
87          %di  =  (dj (:, 1: end −1)* w2 ') .* fneti .*(1− fneti );
88          di  =  (dj (:, 1: end −1)* w2 ') .* 1.154.*(1− tanh(2/3.* fneti ) .* tanh
          (2/3.* fneti ));
89          u3  =  mom.* updatei ;
90          updatei  =  updatei  +  train_x(start +1: start + trainsize ,:) '* di
          (:, 1: end −1);
91          w3  =  w3  +  u1+eta .* updatek ./ trainsize ;
92          w2  =  w2  +  u2+eta .* updatej ./ trainsize ;
93          w1  =  w1  +  u3+eta .* updatei ./ trainsize ;
94          eta  =  eta0 /(1+ i /T);
95
96      end
97      % Calculate  accuracies
98      neti  =  train_x(1:ntrain ,:) *w1;
99      %fneti  =  [sigmoid(neti)  ones([ntrain ,1])];
100     fneti  =  [1.732.* tanh(2/3.* neti)  ones([ntrain ,1])];
101     netj  =  fneti *w2;
```

```matlab
102       %fnetj = [sigmoid(netj) ones([ntrain,1])];
103       fnetj = [1.732.*tanh(2/3.*netj) ones([ntrain,1])];
104       netk = fnetj*w3;
105       fnetk = softmax(netk);
106       [~, Idx] = max(fnetk,[],2);
107       [~, Idx2] = max(train_y(:,1:ntrain),[],1);
108       count0 = nnz(Idx==Idx2');
109
110       neti = val_x(1:nval,:)*w1;
111       %fneti = [sigmoid(neti) ones([nval,1])];
112       fneti = [1.732.*tanh(2/3.*neti) ones([nval,1])];
113       netj = fneti*w2;
114       %fnetj = [sigmoid(netj) ones([nval,1])];
115       fnetj = [1.732.*tanh(2/3.*netj) ones([nval,1])];
116       netk = fnetj*w3;
117       fnetk = softmax(netk);
118       [~, Idx] = max(fnetk,[],2);
119       [~, Idx2] = max(val_y(:,1:nval),[],1);
120       count1 = nnz(Idx==Idx2');
121
122       neti = test_x(1:ntest,:)*w1;
123       %fneti = [sigmoid(neti) ones([ntest,1])];
124       fneti = [1.732.*tanh(2/3.*neti) ones([ntest,1])];
125       netj = fneti*w2;
126       %fnetj = [sigmoid(netj) ones([ntest,1])];
127       fnetj = [1.732.*tanh(2/3.*netj) ones([ntest,1])];
128       netk = fnetj*w3;
129       fnetk = softmax(netk);
130       [~, Idx] = max(fnetk,[],2);
131       [~, Idx2] = max(test_y(:,1:ntest),[],1);
132       count2 = nnz(Idx==Idx2');
133       tr_ac(i) = count0/ntrain;
134       val_ac(i) = count1/nval;
135       te_ac(i) = count2/ntest;
136       fprintf('Iteration : %d      Train Accuracy : %0.4f Val
          Accuracy : %0.4f      Test Accuracy : %0.4f\n',i,count0/ntrain,
          count1/nval,count2/ntest);
137 end
138
139 function [g] = sigmoid(z)
140 par_res=1+exp(-z);
141 g=1./par_res;
142 end
143
144 function [softmaxA] = softmax(A)
145 dim = 2;
146 s = ones(1, ndims(A));
147 s(dim) = size(A, dim);
148 maxA = max(A, [], dim);
149 expA = exp(A-repmat(maxA, s));
```

```matlab
150  softmaxA = expA ./ repmat(sum(expA,dim), s);
151  end
152
153  function [oh] = onehot(v,c)
154  oh = zeros(length(v), c);
155  for i=1:length(v)
156      oh(i,v(i)+1) = 1;
157  end
158  end
```

Listing 2: MATLAB code to run a 3 layer neural network

```matlab
1   function images = loadMNISTImages(filename)
2   %loadMNISTImages returns a 28x28x[number of MNIST images] matrix
        containing
3   %the raw MNIST images
4
5   fp = fopen(filename, 'rb');
6   assert(fp ~= -1, ['Could not open ', filename, '']);
7
8   magic = fread(fp, 1, 'int32', 0, 'ieee-be');
9   assert(magic == 2051, ['Bad magic number in ', filename, '']);
10
11  numImages = fread(fp, 1, 'int32', 0, 'ieee-be');
12  numRows = fread(fp, 1, 'int32', 0, 'ieee-be');
13  numCols = fread(fp, 1, 'int32', 0, 'ieee-be');
14
15  images = fread(fp, inf, 'unsigned char');
16  images = reshape(images, numCols, numRows, numImages);
17  images = permute(images,[2 1 3]);
18
19  fclose(fp);
20
21  % Reshape to #pixels x #examples
22  images = reshape(images, size(images, 1) * size(images, 2), size(
        images, 3));
23  % Convert to double and rescale to [0,1]
24  images = double(images) / 255;
25
26  end
```

Listing 3: MATLAB code to load images from data

```matlab
1   function labels = loadMNISTLabels(filename)
2   %loadMNISTLabels returns a [number of MNIST images]x1 matrix
        containing
3   %the labels for the MNIST images
4
5   fp = fopen(filename, 'rb');
6   assert(fp ~= -1, ['Could not open ', filename, '']);
```

```
7
8  magic = fread(fp, 1, 'int32', 0, 'ieee-be');
9  assert(magic == 2049, ['Bad magic number in ', filename, '']);
10
11 numLabels = fread(fp, 1, 'int32', 0, 'ieee-be');
12
13 labels = fread(fp, inf, 'unsigned char');
14
15 assert(size(labels,1) == numLabels, 'Mismatch in label count');
16
17 fclose(fp);
18
19 end
```

Listing 4: MATLAB code to load labels from data

```
1  % Change the filenames if you've saved the files under different
       names
2  % On some platforms, the files might be saved as
3  % train-images.idx3-ubyte / train-labels.idx1-ubyte
4  images_train = loadMNISTImages('train-images-idx3-ubyte')';
5  labels_train = loadMNISTLabels('train-labels-idx1-ubyte');
6  y_train = labels_train;
7  labels_train = onehot(labels_train,10);
8
9  for it=1:60000
10     m = mean(images_train(it,:));
11     images_train(it,:) = images_train(it,:)-m;
12 end
13
14 images_test = loadMNISTImages('t10k-images-idx3-ubyte')';
15 labels_test = loadMNISTLabels('t10k-labels-idx1-ubyte');
16 y_test = labels_test;
17 labels_test = onehot(labels_test,10);
18
19
20 for it=1:10000
21     m = mean(images_test(it,:));
22     images_test(it,:) = images_test(it,:)-m;
23 end
24
25 ntrain = 50000;
26 nval = 10000;
27 ntest = 10000;
28
29 train_x = [ones(ntrain,1),images_train(1:ntrain,:)];
30 train_y = labels_train(1:ntrain,:)';
31
32 val_x = [ones(nval,1),images_train(ntrain+1:ntrain+nval,:)];
33 val_y = labels_train(ntrain+1:ntrain+nval,:)';
```

```matlab
34
35  test_x = [ ones ( ntest ,1) , images_test (1: ntest ,:) ];
36  test_y = labels_test (1: ntest ,:) ';
37
38  load weights.mat;
39  eta = 0.001;
40  dnum = 0.0;
41  dcalc = 0.0;
42  count = 0;
43  epsilon = 0.01;
44  nlayers = length (w);
45  n = 100000;
46  delta = 0.0;
47  s = 1;
48  for i =1: nlayers
49      for j =1: size (w{ i } ,1)
50          for k =1: size (w{ i } ,2)
51              %forward propagate
52              wp = w;
53              wn = w;
54              wp{ i }( j ,k) = wp{ i }( j ,k)+epsilon ;
55              wn{ i }( j ,k) = wn{ i }( j ,k)−epsilon ;
56
57              netj = train_x (s ,:) *w{1};
58              fnetj = [ sigmoid ( netj ) 1];
59              netk = fnetj *w{2};
60              y = softmax ( netk );
61
62              netj = train_x (s ,:) *wp{1};
63              fnetj = [ sigmoid ( netj ) 1];
64              netk = fnetj *wp{2};
65              yp = softmax ( netk );
66
67              netj = train_x (s ,:) *wn{1};
68              fnetj = [ sigmoid ( netj ) 1];
69              netk = fnetj *wn{2};
70              yn = softmax ( netk );
71
72              dnum = calcNumGradient ( yp , yn , train_y (: , s ) , epsilon );
73
74              %error calc
75              update = w;
76              dk = ( train_y (: , s )−y ') ;
77              update {2} =( fnetj '*dk ') ;
78              dj = ( dk '*w{2} ') .* fnetj .*(1− fnetj );
79              update {1} = train_x (s ,:) '* dj (: ,1: end −1);
80              dcalc = −update { i }( j ,k );
81              delta = delta + abs ( dcalc −dnum );
82              count = count +1;
83          end
```

```matlab
84          end
85   end
86
87   fprintf('Avg Difference in Gradient : %0.4f        Epsilon square :
         %0.4f\n', delta/count, epsilon^2);
88
89   function [g] = sigmoid(z)
90   par_res =1+exp(−z);
91   g=1./par_res;
92   end
93
94   function [softmaxA] = softmax(A)
95   dim = 2;
96   s = ones(1, ndims(A));
97   s(dim) = size(A, dim);
98   maxA = max(A, [], dim);
99   expA = exp(A−repmat(maxA, s));
100  softmaxA = expA ./ repmat(sum(expA,dim), s);
101  end
102
103  function [oh] = onehot(v,c)
104  oh = zeros(length(v), c);
105  for i=1:length(v)
106      oh(i,v(i)+1) = 1;
107  end
108  end
109
110  function [numE] = calcNumGradient(yp,yn,t,ep)
111  yp =log(yp);
112  yn = log(yn);
113  Ep = −yp*t;
114  En = −yn*t;
115  numE = (Ep−En)/(2*ep);
116  end
```

Listing 5: MATLAB code to verify the correctness of the backprop implementation