
Using Neural Networks to Effectively Classify Hand-Written Digits of the MNIST Dataset

Chetan Gandotra	Rishabh Misra
UC San Diego	UC San Diego
9500 Gilman Drive	9500 Gilman Drive
cgandotr@ucsd.edu	rimisra@ucsd.edu

Abstract

1 This report discusses the second programming assignment of our course CSE 253:
2 Neural Networks and Pattern Recognition, its solutions and the inferences we drew.
3 A variable layer neural network was implemented from the scratch and observations
4 were made based on various parameters and before/after adding certain trades of
5 tricks as discussed in Yann LeCun's famous paper "Efficient BackProp". The
6 data-set used was the famous MNIST data-set and a ten-way classification was
7 performed on it. A test data-set accuracy in excess of 97% was achieved using
8 various mechanisms and tricks, which is almost at par with the accuracy reported
9 by LeCun on his website.

10 1 Task 3: Implementing Neural Network and Gradient Calculation

11 1.1 Introduction

12 This problem asks us to build a neural network from scratch using our derivations from Problem 2.
13 We are required to read the MNIST data-set and normalize the data by dividing each pixel value by
14 255, followed by subtracting the mean over all of the pixels in each image. We have been asked to
15 use the softmax activation function for output layer and logistic activation function for hidden layers.
16 We need to check our code for gradient computation by using the following numerical approximation
17 formula for each weight:

$$\frac{d}{dw} E^n(w) \approx \frac{E^n(w + \epsilon) - E^n(w - \epsilon)}{2\epsilon} \quad (1)$$

18 Finally, using the update rule from 2(c), we will perform gradient descent and report our accuracy on
19 training and test data-sets.

20 1.2 Methodology

21 The MNIST data-set was downloaded and extracted into local variables using existing libraries
22 available online on GitHub. [Note that we used the same off-the-shelf GitHub library for PA1 as
23 well]. Once we had the data, we normalized it using the instructions mentioned. Each value was
24 divided by 255, and the row-wise mean was subtracted from each element in that row.

25 Then, a neural network was implemented that included forward propagation, back-propagation and
26 weight update steps. The activation functions used were logistic (for hidden layer) and softmax (for
27 output layer). Thus, our neural network had three layers - one input layer of size 784, one hidden
28 layer of size 300 and lastly, the output layer of size 10. The hidden layer was chosen to have 300
29 units as many of the experiments on LeCun's website had 300 hidden units in them. Note that the
30 output was converted into a "one-hot" representation form. The learning rate was set to be 0.0000044
31 (obtained using - 0.22/50k examples) after much experimentation and inferring the best learning rate.

The full batch of 50k examples was given to the neural network to learn, and a stopping condition was enforced using a neural network. If at any point 5 or more continuous validation error (mis-classification error) values increased, we stopped training our model and exited. All the previous weights were stored and the weights having least error on validation set were selected to be the optimal weights. Note that we also used cross entropy error on validation set to enforce a stopping condition and it yielded a similar learning pattern. Due to this, we decided to go with either one of the two error calculations.

To sum: $\eta = 0.0000044$, examples = 50,000 (full batch used), convergence criteria = Validation set mis-classification error

For 3(d), we computed the slope with respect to one weight using the numerical approximation in equation (29). The value of ϵ was taken to be 0.01 and the gradient for weights of each layer was compared to our gradient calculation method. Both the gradients were calculated and stored in lists and then subtracted. Finally, the subtracted result was used to draw inferences.

1.3 Results

The gradient difference findings are as under:

Maximum Absolute Difference in W_{ij} : $5.01416981156e-07$

Maximum Absolute Difference in W_{jk} : 0.000900066795115

Mean Difference for W_{ij} : $3.80757313778e-08$

Mean Difference for W_{jk} : $9.02485484117e-05$

The detailed differences for every weight in W_{ij} and W_{jk} are given in attached files - GradientDifferenceIJ.csv and GradientDifferenceJK.csv, respectively.

Using gradient descent on full training set (minus the validation set) resulted in an accuracy of 91.8% on the training set and 91.7% on the test-set in 527 iterations. Note that validation error was used as a means of stopping criteria.

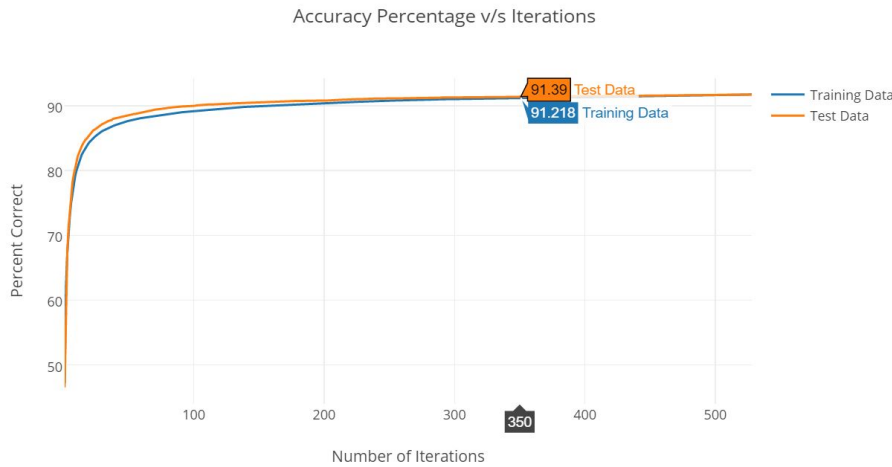


Figure 1: Accuracy Percentage v/s Iterations for Full Batch in 2-layer Neural Network with 300 Hidden Units

55

Graph 1 shows the percent correct plots for training and test data over 500+ iterations. As shown in the graph, at 350th iteration, we had a training accuracy of 91.39% and test accuracy of 91.218%.

1.4 Discussion

The gradient calculation from part 3(d) performed as expected. With an ϵ value 0.01, the difference of gradients agreed in $O(\epsilon^2)$, i.e., 10^{-4} . This tells us that our gradient computation is reliable and can be relied upon for the purpose of this assignment.

62 The results of 10-way classification using a two-layer neural network on MNIST data-set were
63 satisfactory. The convergence slows down a lot as number of iterations increase and more than 500
64 iterations were required to achieve an accuracy of 91.7% on the test set. As you can see in graph 1,
65 there was little improvement on the training and test accuracies after 200 iterations and the graph was
66 this point was almost a straight line parallel to the iteration axis. This gives us a motivation to use
67 and try out the various tricks of trades mentioned by Yann LeCun et al in the famous paper "Efficient
68 BackProp".

69 2 Task 4: Adding the "Tricks of the Trade"

70 2.1 Introduction

71 In this section, we would apply certain tricks from a famous paper by Yann Lecun and analyze their
72 effect on the performance of our neural network. Particularly, we would be doing following things:

- 73 1. We would shuffle the data and apply stochastic gradient descent over a mini-batch of varying
74 size to analyze the effect of mini-batch size on performance.
- 75 2. We would use a different sigmoid function ($1.7159 \tanh(\frac{2x}{3})$) as an activation function for
76 the hidden layer and observe its effect.
- 77 3. We would observe the effects when weights are initialized in a specific way suggested in the
78 paper.
- 79 4. We would add a momentum term in the update rules and would observe its effects.

80 2.2 Methodology

81 We will be applying all the above task incrementally and observe their effects. For the first part, we
82 would randomly shuffle the data using numpy library's inbuilt shuffle function. Then we consider
83 mini-batches of size 50, 128 and 500 and analyze their performance on our data. Learning rate is set
84 to 0.1 and number of units in hidden layer is 300.

85 For the second part, we choose mini-batch size of 128 and use the sigmoid function ($1.7159 \tanh(\frac{2x}{3})$)
86 as an activation function for the hidden layer. The value of the hyperparameters is same as that of
87 previous part. Derivative of this function could be found as following.

$$\begin{aligned} \frac{d}{dx} 1.7159 \tanh\left(\frac{2x}{3}\right) &= 1.7159 (1 - \tanh^2\left(\frac{2x}{3}\right)) \frac{2}{3} \\ &= 1.1439 (1 - \tanh^2\left(\frac{2x}{3}\right)) \end{aligned}$$

88
89 For the next part, we initialize the weights by random samples drawn from normal distribution with
90 $\mu = 0$ and $\sigma = \sqrt{\frac{1}{fan-in}}$ where the *fan-in* is the number of units in the respective layer. The
91 value of the hyperparameters is same as that of previous part.

92 For the last part, we add a momentum term in our update rule. Now the update rules for the weights
93 would look like following:

$$\begin{aligned} w_{jk}^{t+1} &= w_{jk}^t - \eta \frac{\partial E}{\partial w_{jk}^{t+1}} + \beta(w_{jk}^t - w_{jk}^{t-1}) \\ w_{ij}^{t+1} &= w_{ij}^t - \eta \frac{\partial E}{\partial w_{ij}^{t+1}} + \beta(w_{ij}^t - w_{ij}^{t-1}) \end{aligned}$$

94
95 where β controls how much weight the momentum term should have. For our experiments, we have
96 considered the value of β as 0.9. The value of other hyperparameters is same as that of previous part.

97 We'll run all the experiments on 1000 iterations with early stopping horizon value of 5. However,
98 since our empirical results showed no significant improvement after 100 iterations, we will show the
99 results only on first 100 iterations for our comparative study (otherwise the shape of curve would not
100 be visible clearly). Another motivation for taking first 100 iterations for plotting is that any setting
101 could lead to best possible value when allowed to run sufficiently large time, but that wouldn't give
102 us a clear picture of effects of the tricks we are applying.

103 2.3 Results

104 Results for each of the part is shown below:

105 1 - Results for shuffling of data and applying stochastic gradient descent over a mini-batch of varying
106 size are shown below.

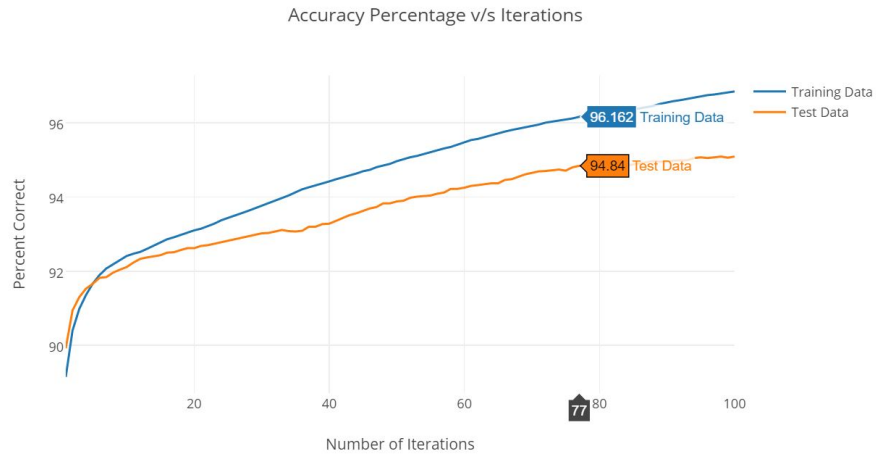


Figure 2: Train and test accuracy for mini-batch size 50

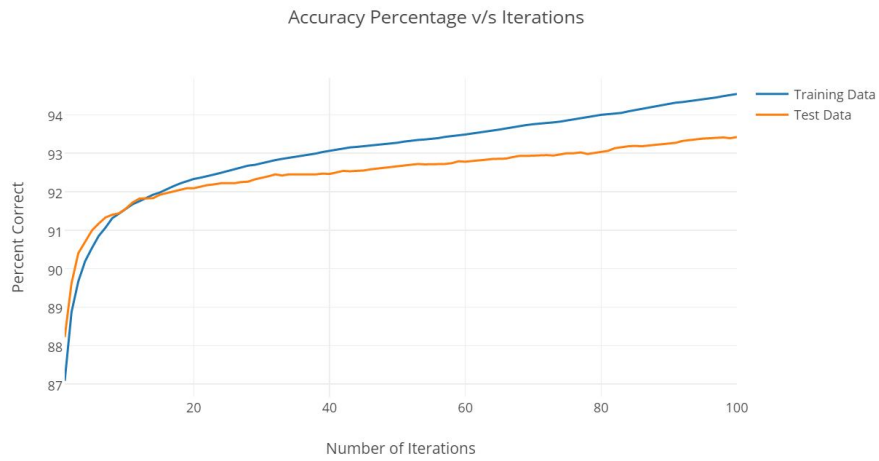


Figure 3: Train and test accuracy for mini-batch size 128

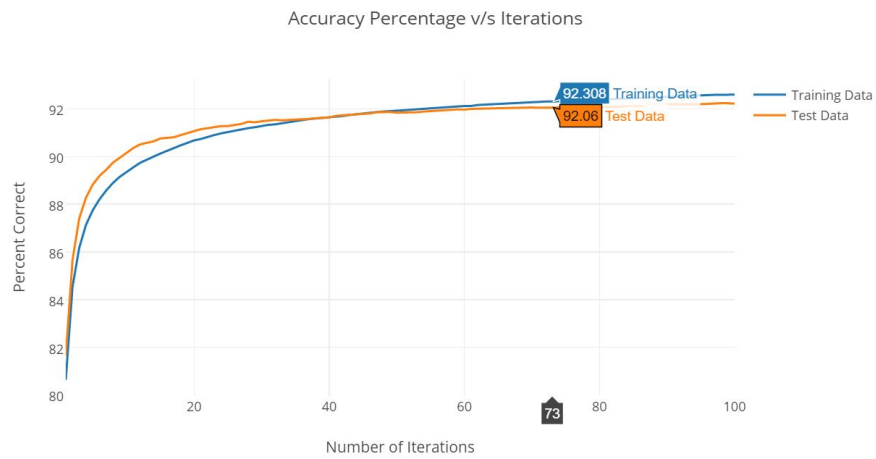


Figure 4: Train and test accuracy for mini-batch size 500

107 2 - Results when using function $1.7159 \tanh(\frac{2x}{3})$ as an activation function for the hidden layer are
 108 shown below:

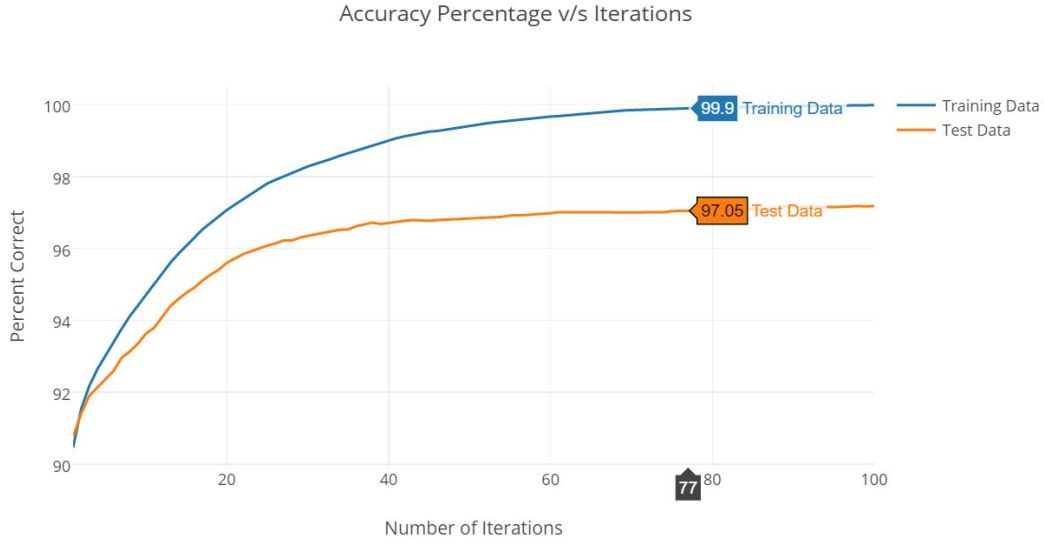


Figure 5: Train and test accuracy for tanh activation function with mini-batch size 128

109 3 - Results when weights are sampled from a normal distribution with $\mu = 0$ and $\sigma = \sqrt{\frac{1}{fan-in}}$

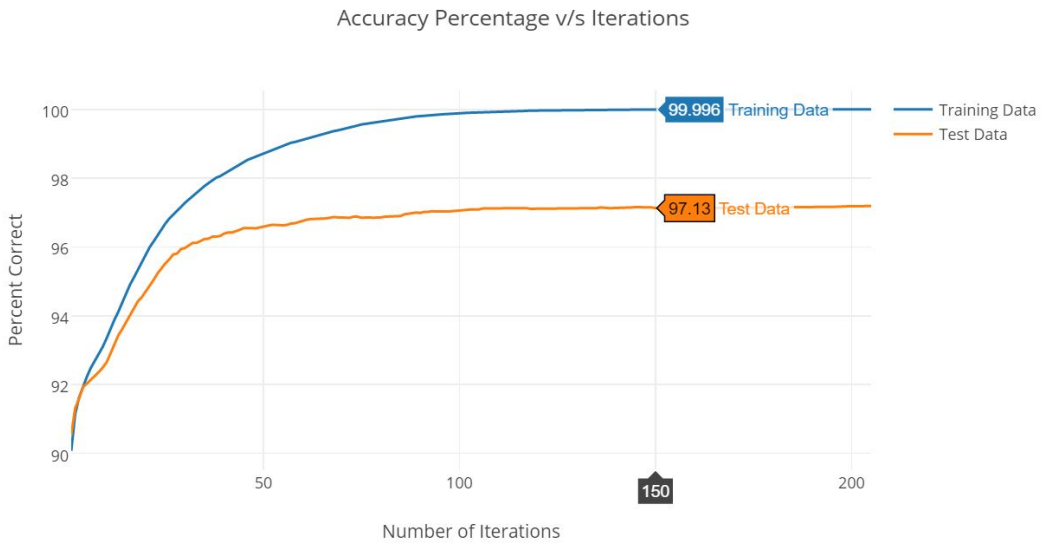


Figure 6: Train and test accuracy when weights sampled from a normal distribution with $\mu = 0$ and $\sigma = \sqrt{\frac{1}{fan-in}}$

110 4 - Results when momentum term is added in our update rule.

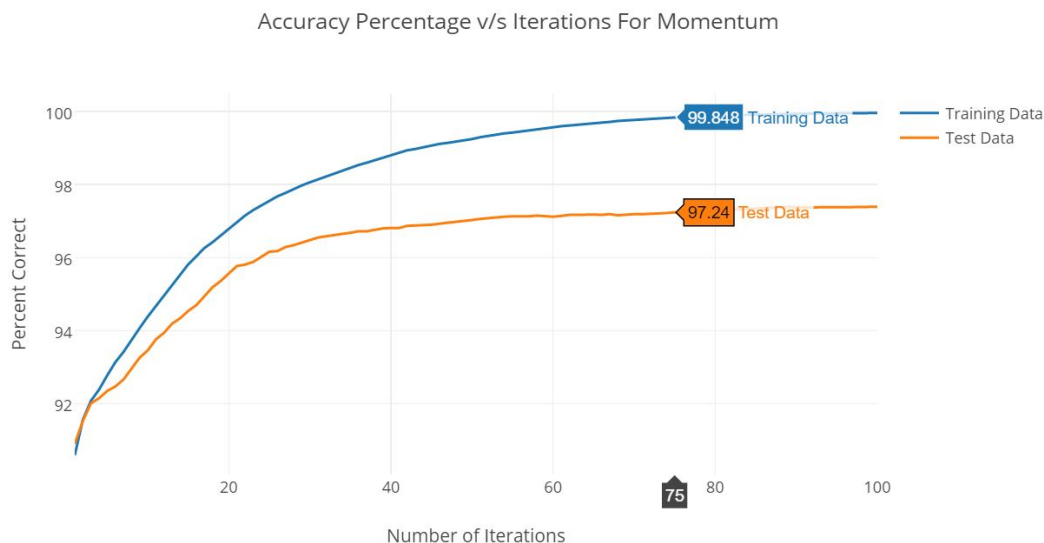


Figure 7: Train and test accuracy when we consider momentum with weight 0.9

2.4 Discussion

For the first part, where we introduced the shuffling and mini-batching, the precision when compared with the full batch learning is better by at least 5% on the test data at 100 epochs (2 and 1). Also, the precision for full batch learning starts from 10% where as for mini-batches, it starts from 90%. Mini-batch learning converges faster to a better precision when compared with full-batch learning. This is because mini-batch learning uses noisy gradient to move towards the optimum which results in faster convergence. It might have happened that full-batch learning got stuck at local minima because the precision saturates around 90% where as for mini-batch learning precision reaches around 95%. As the number of batches are increased, we see that performance getting a little bit subdued. This is because the gradient become less and less noisy and mini-batch starts having the side effects of full-batch learning. The described behavior is evident from Graphs 2, 3, and 4.

For the second part, where we chose different sigmoid function for hidden layer, the precision when compared with the case where logistic function was can be clearly seen in graphs 1 and 5. There's an improvement of about 2% within 100 epochs. This is because the tanh activation function is symmetric about origin and is more likely to produce output whose mean is close to zero unlike logistic function. Also, we notice that in case of tanh, the convergence is better than that of logistic function.

For the next part, where we initialize the weights using normal distribution with $\mu = 0$ and $\sigma = \sqrt{\frac{1}{fan-in}}$, we see a slight improvement in the performance when comparing graphs 5 and 6. The difference is not much because prior to this we were initializing the weights using normal distribution with $\mu = 0$ and $\sigma = 0.1$. Nevertheless, we can see from the graphs that the convergence speed is slightly improved. This is expected because if the weights are large or small the learning speed would be slow. Our initialization method here ensures that weights are neither too large nor too small.

For the last part, where we will be using a momentum term to update the weights, we see a slight improvement in the convergence rate when comparing graphs 7 and 5. This is expected because the momentum term forces the weights to move into the direction of previous update. If the weights are diverging, momentum would add a positive term to subdue this divergence. Similarly, if weights are converging, momentum would improve the gradient so that the convergence is sped up.

Thus, after applying all the tricks, we can see a lot of improvement in the performance and accuracy from our neural network. This is evident when we compare graphs 1 and 7.

141 3 Task 5: Experiment with Network Topology

142 3.1 Introduction

143 In this section, we would experiment with the network architecture and observe the changes it has on
144 the performance on our dataset.

145 First, we would change the number of hidden units in our network and analyze the difference. We
146 will try four configuration where the number of hidden units would be 150 (half), 600 (double), 10
147 (very small) and 1000 (very large) respectively.

148 Next, we would introduce an additional hidden layer (such that the number of weights parameters are
149 approximately the same as that of single layer) and would analyze the effect it has on the performance
150 on our dataset.

151 3.2 Methodology

152 For the first part, since we had taken 300 hidden units in the section 4, we would half (150) and
153 double (600) the number of hidden units here to see the effect on performance. All the parameters
154 like learning rate, momentum, batch size, early stopping horizon etc. are same as that of previous
155 part (0.1, 0.9, 128, 5 respectively).

156 For the second part, we have to introduce one additional hidden layer such that the number of
157 parameters are approximately the same. This means we have to find the solution of following
158 equation:

$$785x + (x + 1) * x + (x + 1) * 10 = 785 * 300 + 301 * 10$$

$$796x + x^2 + 10 = 235500 + 3010$$

$$x^2 + 796x - 238500 = 0$$

$$x = \frac{-796 \pm \sqrt{796 * 796 + 4 * 238500}}{2}$$

$$x \approx \frac{-796 \pm 1260}{2}$$

$$x \approx \frac{-796 \pm 1260}{2}$$

$$x \approx \frac{464}{2} (\text{negative } x \text{ not possible})$$

$$x \approx 232$$

159 Hence, we would consider 232 units (plus 1 bias unit) for each of the hidden layers. All the hyperpa-
160 rameters remain the same as that of previous part. We also consider one additional configuration of
161 300 and 150 hidden units to see the affect of increasing the number of parameters on 2 hidden layer
162 network architecture.

163 3.3 Results

164 For the part where we have to increase/decrease the number of hidden units in the single layer, the
165 results are shown below:

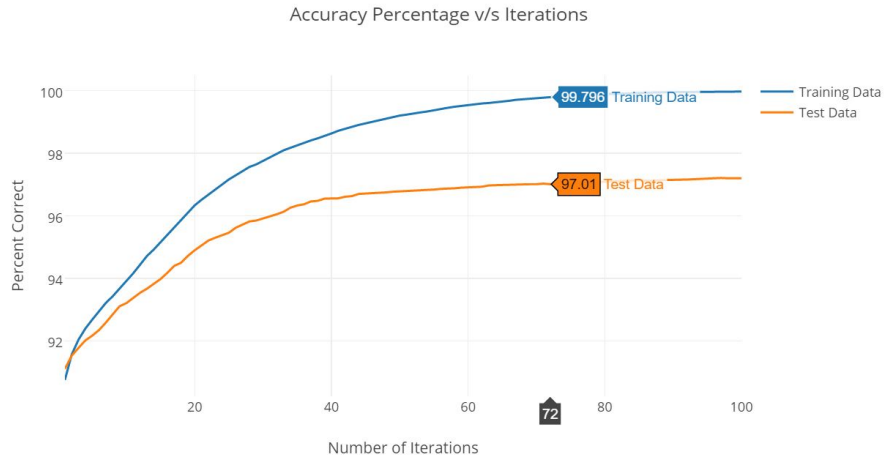


Figure 8: Train and test accuracy for 1 hidden layer with 600 units

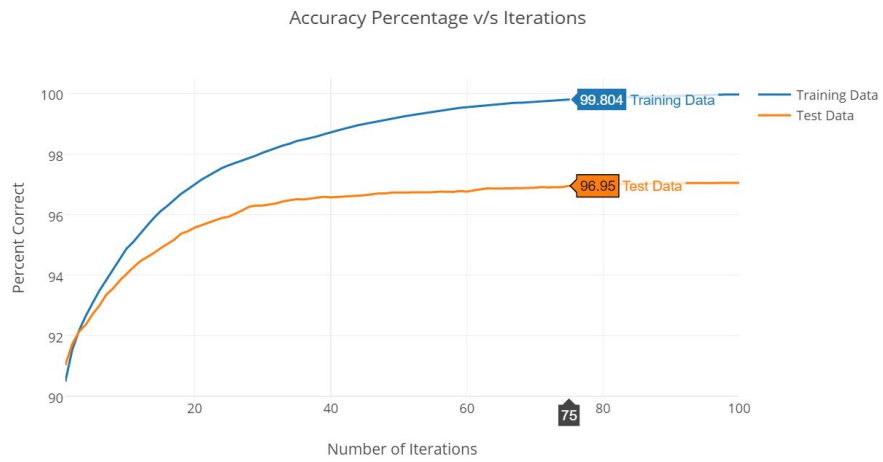


Figure 9: Train and test accuracy for 1 hidden layer with 150 units

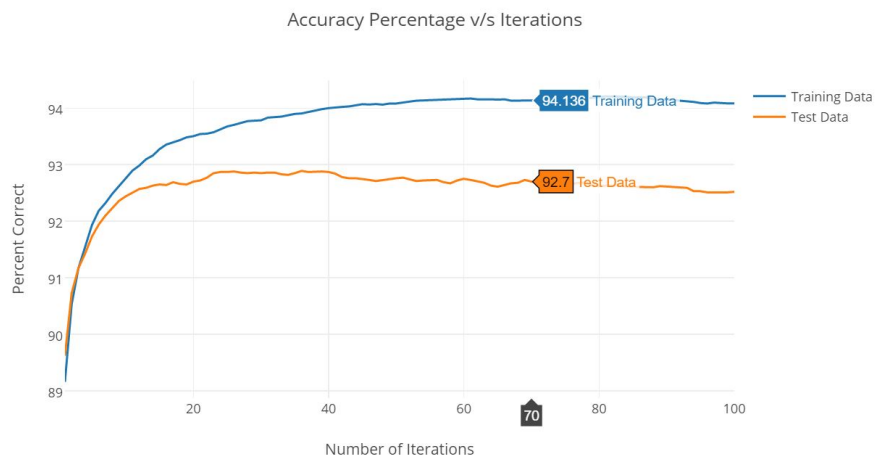


Figure 10: Train and test accuracy for 1 hidden layer with 10 units

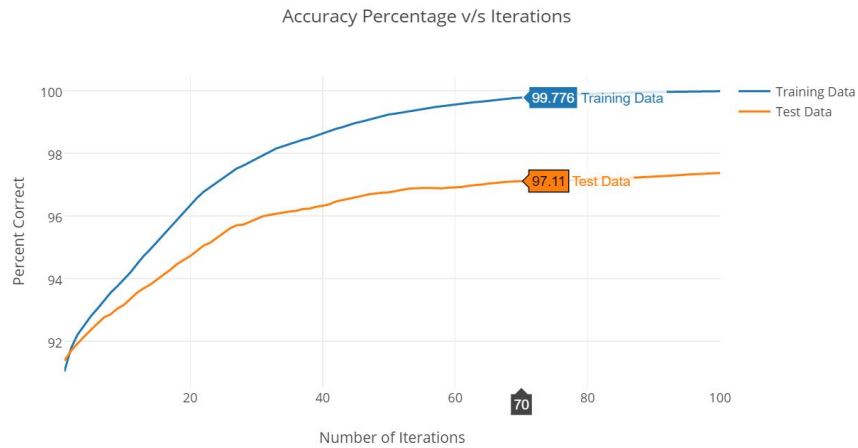


Figure 11: Train and test accuracy for 1 hidden layer with 1000 units

166 For the part where we have to introduce additional hidden layer, the results are shown below:

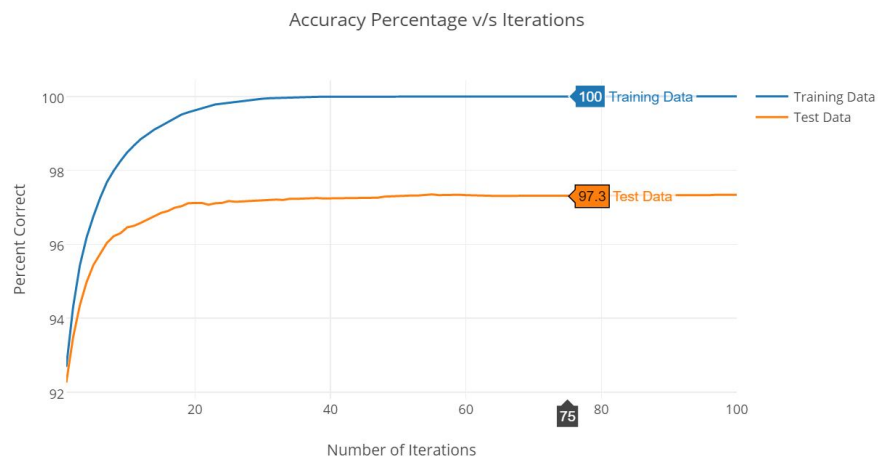


Figure 12: Train and test accuracy for 2 hidden layers with 232 units each

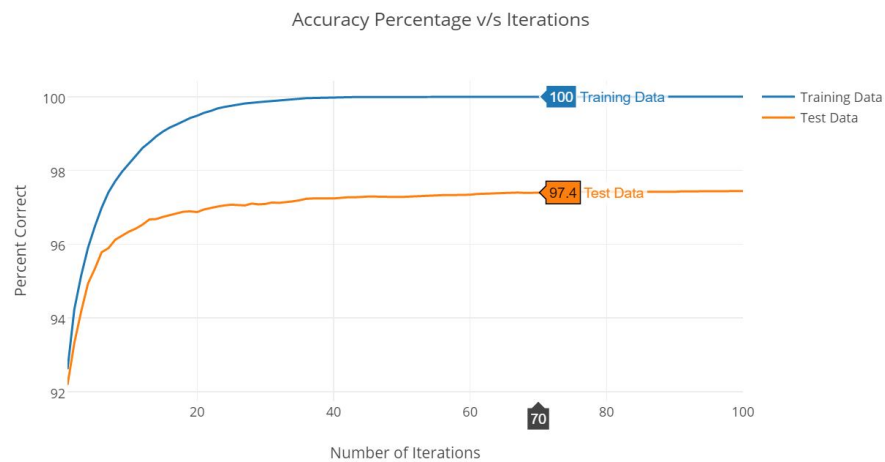


Figure 13: Train and test accuracy for 2 hidden layers with 300 and 150 units respectively

3.4 Discussion

Though, the experiments were allowed to run more than 100 epochs with early stopping condition, we are reporting the results of first 100 iteration, where the result has almost converged, to have a better comparative study.

From the Graphs 8, 9, 10 and 11, we infer the following things. When the number of units are very less (10), our network has an overtly less accuracy on both test and train data as compared to the original case (300 units). This is expected as the function learned by our network would be very simplistic and thus the predictions are a bit off. As we increase the number of hidden units, the performance on training as well as test starts to improve. We can see the improvement in performance from graph 9 and 7 on both test and train data. However, when we increase the hidden units further, we do not see any big improvements in the results. This can be seen from the graph 10 and 11. The accuracy on test and train is approximately the same. Though, cases with more hidden units takes a bit more time to converge. This could be because more units take more time to get tuned to learn the same function. Also, as we increase the number of units, our method takes more time to perform the same number of epochs as the number of parameters to tune would increase. Hence, our takeaway from this is that after a certain number of units for the hidden layer, we don't see any significant improvement.

For the part where we have to introduce additional layer in our network, we infer the following things from graphs 12 and 13. For the first graph, the number of parameters are approximately the same as that of the network with 1 layer with 300 units. In the graph, we see slight improvement in accuracy for 2-layered network. This could be because our 2 hidden layered network is capable learning more complex function with arbitrary decision boundaries. One more thing to note is that the weights converge significantly quickly in 2 hidden layered network when compared to 1 hidden layered network. For 2-layered network where the number of parameters are more (300 and 150 units), we see a slight improvement in test data accuracy (train accuracy is already 100% in both) and convergence is relatively a bit slow. Again, this is because we have more parameters to optimize.

193 4 Results and Learnings

194 The best accuracy on testing data set was achieved to be 97.54% with two hidden layers, while more
195 than 97% accuracy was achieved easily using tricks of trades as mentioned by Yann LeCun. Using
196 these tricks improves the learning by a significant amount, saving the developer a lot of time.

197 Shuffling and mini-batching improved the performance by almost 5% and also increased the learning
198 speed. Same goes for the *tanh* sigmoid function, which seems to be way faster than the traditional
199 logistic sigmoid function. Initializing the weights using normal distribution with $\mu = 0$ and $\sigma =$
200 $\sqrt{\frac{1}{fan-in}}$, we see a slight improvement in the performance when comparing graphs 5 and 6. This
201 gives us a good idea as to how to initialize the weights while using neural networks in future. Lastly,
202 we learned the concept of momentum, which plays a great part in correcting the direction of gradient
203 allowing faster and better convergence.

204 Both of us agree that this programming assignment was definitely an eye-opener for us in the field
205 of neural networks. We had a good time implementing it and experimenting with various values
206 of the learning rate η . The "tricks of trade" mentioned improved the performance by a significant
207 amount and this is something that we never knew of. We definitely look forward to more of such
208 interactive and practical programming assignments which take more time in implementation and
209 focus on understanding the concepts.

210 5 Individual Contributions

211 Being roommates, it was extremely convenient and simple for both the authors to coordinate and
212 work in sync, while ensuring equal distribution of work and time spent on the assignment. Whenever
213 one of the authors got stuck at some point, the other was there to help him out and unblock instantly.
214 The authors started out with a plan of having perfectly modularized and re-usable code, so that for
215 question 5, the effort to add another hidden layer would be minimal. The work was taken up as under:

216 Chetan Gandotra implemented the part which involved reading of MNIST data, normalization of
217 data and forward propagation. Then, Rishabh Misra implemented back-propagation and weight
218 updates. Thereafter, Chetan worked with code for momentum and tanh activation function, and
219 Rishabh worked with the code for 3 (d). The debugging and parameter tuning task was divided
220 equally, with both of the authors testing performance on defined range of parameters. This allowed
221 them to come up with good hyper-parameters fairly quickly.

222 When it came to writing the report, we took up alternate question sub-parts on shared latex file online,
223 ensuring equal work distribution.

References

[1] LeCun, Yann - Efficient BackProp : <http://yann.lecun.com/exdb/publis/pdf/lecun-98b.pdf>

Appendix

The code consists of three files - LoadMNIST.py, PA2.py and 3d.py

LoadMNIST.py

```
import os, struct
from array import array as pyarray
from numpy import append, array, int8, uint8, zeros
```

```
def load_mnist(dataset="training", digits=None, path=None, asbytes=False, selection=None, return_labels=False, return_indices=False):
```

```
    """
```

```
    Loads MNIST files into a 3D numpy array.

    You have to download the data separately from [MNIST]_. It is recommended
    to set the environment variable 'MNIST' to point to the folder where you
    put the data, so that you don't have to select path. On a Linux+bash setup,
    this is done by adding the following to your '.bashrc':
```

```
        export MNIST=/path/to/mnist
```

Parameters

dataset : str

Either "training" or "testing", depending on which dataset you want to load.

digits : list

Integer list of digits to load. The entire database is loaded if set to 'None'. Default is 'None'.

path : str

Path to your MNIST datafiles. The default is 'None', which will try to take the path from your environment variable 'MNIST'. The data can be downloaded from <http://yann.lecun.com/exdb/mnist/>.

asbytes : bool

If True, returns data as 'numpy.uint8' in [0, 255] as opposed to 'numpy.float64' in [0.0, 1.0].

selection : slice

Using a 'slice' object, specify what subset of the dataset to load. An example is 'slice(0, 20, 2)', which would load every other digit until—but not including—the twentieth.

return_labels : bool

Specify whether or not labels should be returned. This is also a speed performance if digits are not specified, since then the labels file does not need to be read at all.

return_indices : bool

Specify whether or not to return the MNIST indices that were fetched. This is valuable only if digits is specified, because in that case it can be valuable to know how far in the database it reached.

Returns

images : ndarray

Image data of shape '(N, rows, cols)', where 'N' is the number of images. If neither

labels : ndarray

Array of size 'N' describing the labels. Returned only if 'return_labels' is 'True'

indices : ndarray

The indices in the database that were returned.

Examples

```

284     Assuming that you have downloaded the MNIST database and set the
285     environment variable '$MNIST' point to the folder, this will load all
286     images and labels from the training set:
287
288     >>> images, labels = ag.io.load_mnist('training') # doctest: +SKIP
289
290     Load 100 sevens from the testing set:
291
292     >>> sevens = ag.io.load_mnist('testing', digits=[7], selection=slice(0, 100), return_label
293     ""
294
295     # The files are assumed to have these names and should be found in 'path'
296     files = {
297         'training': ('train-images.idx3-ubyte', 'train-labels.idx1-ubyte'),
298         'testing': ('t10k-images.idx3-ubyte', 't10k-labels.idx1-ubyte'),
299     }
300
301     if path is None:
302         try:
303             path = 'C:\\Users\\Chetan\\Documents\\Python Scripts\\Way1'
304             #path = os.environ['MNIST']
305         except KeyError:
306             raise ValueError("Unspecified path requires environment variable $MNIST to be set")
307
308     try:
309         images_fname = os.path.join(path, files[dataset][0])
310         labels_fname = os.path.join(path, files[dataset][1])
311     except KeyError:
312         raise ValueError("Data set must be 'testing' or 'training'")
313
314     # We can skip the labels file only if digits aren't specified and labels aren't asked for
315     if return_labels or digits is not None:
316         flbl = open(labels_fname, 'rb')
317         magic_nr, size = struct.unpack(">II", flbl.read(8))
318         labels_raw = pyarray("b", flbl.read())
319         flbl.close()
320
321         fimg = open(images_fname, 'rb')
322         magic_nr, size, rows, cols = struct.unpack(">IIII", fimg.read(16))
323         images_raw = pyarray("B", fimg.read())
324         fimg.close()
325
326         if digits:
327             indices = [k for k in range(size) if labels_raw[k] in digits]
328         else:
329             indices = range(size)
330
331         if selection:
332             indices = indices[selection]
333     N = len(indices)
334
335     images = zeros((N, rows, cols), dtype=uint8)
336
337     if return_labels:
338         labels = zeros((N), dtype=int8)
339     for i, index in enumerate(indices):
340         images[i] = array(images_raw[ indices[i]*rows*cols : (indices[i]+1)*rows*cols ]).resha
341         if return_labels:
342             labels[i] = labels_raw[indices[i]]
343
344     #if not asbytes:
345     #    images = images.astype(float)/255.0
346
347     ret = (images,)
348

```

```

349     if return_labels:
350         ret += (labels,)
351     if return_indices:
352         ret += (indices,)
353     if len(ret) == 1:
354         return ret[0] # Don't return a tuple of one
355     else:
356         return ret

357 PA2.py

358 """
359 Variable layer Neural Networks code to Classify Hand-written digits
360 in MNIST Dataset
361 """
362 import numpy
363 import math
364 from LoadMNIST import load_mnist
365 from sklearn.utils import shuffle
366 import plotly.plotly as pyl
367 import plotly.graph_objs as go
368 #import bigfloat
369
370 #-----Utility functions-----
371
372 def get_data(N=60000, N_test=10000, validationReqd = True):
373     # Load MNIST data using libraries available
374     training_data, training_labels = load_mnist('training ')
375     test_data, test_labels = load_mnist('testing ')
376
377     # Training_data is N x 784 matrix
378     training_data = flatten(N, 784, training_data)
379     training_labels = training_labels[:N]
380     test_data = flatten(N_test, 784, test_data)
381     test_labels = test_labels[:N_test]
382
383     # Adding column of 1s for bias
384     training_data = addOnesColAtStart(training_data)
385     test_data = addOnesColAtStart(test_data)
386
387     if (validationReqd):
388         # Last 10% of training data size will be considered as the validation set
389         N_validation = int (N / 6.0)
390         validation_data = training_data[N-N_validation:N]
391         validation_labels = training_labels[N-N_validation:N]
392         N=N-N_validation
393     else:
394         validation_data = []
395         validation_labels = []
396
397     # Update training data to remove validation data
398     training_data = training_data[:N]
399     training_labels = training_labels[:N]
400
401     # Normalization of Data
402     training_data = training_data/255.0
403     test_data = test_data/255.0
404     validation_data = validation_data/255.0
405     training_data = training_data - numpy.mean(training_data, axis=0)[numpy.newaxis, :]
406     test_data = test_data - numpy.mean(test_data, axis=0)[numpy.newaxis, :]
407     validation_data = validation_data - numpy.mean(validation_data, axis=0)[numpy.newaxis, :]
408
409     return training_data, training_labels, test_data, test_labels, validation_data, validation_labels
410
411 # Convert from tuple form to Matrix form

```

```

412 def flatten(rows, cols, twoDArr):
413     flattened_arr = numpy.zeros(shape=(rows, cols))
414     for row in range(0, rows):
415         i=0
416         for element in twoDArr[row]:
417             for ell in element:
418                 flattened_arr[row][i] = ell
419                 i = i+1
420     return flattened_arr
421
422 def addOnesColAtStart(matrix):
423     Ones = numpy.ones(len(matrix))
424     newMatrix = numpy.c_[Ones, matrix]
425     return newMatrix
426
427 # Single method for calculating sigmoid (logisitic) activation
428 # and its derivative
429 def f_sigmoid(X, derivativeReqd=False):
430     if not derivativeReqd:
431         return 1.0 / (1 + numpy.exp(-X))
432         #return 1.0 / (1 + bigfloat.exp(-X, bigfloat.precision(100)))
433     return numpy.multiply(f_sigmoid(X), (1 - f_sigmoid(X)))
434
435 # Method to calculate the softmax activation
436 def f_softmax(X):
437     Z = numpy.sum(numpy.exp(X), axis=1)
438     Z = Z.reshape(Z.shape[0], 1)
439     return numpy.exp(X) / Z
440
441 # Single method for tanh sigmoid and its derivative
442 def f_tanh(x, a=0, derivativeReqd = False):
443     mul_factor = 1.7159
444     div_factor = 2.0/3.0
445     tanh_term = numpy.tanh(div_factor*x)
446     if not derivativeReqd:
447         return (mul_factor * tanh_term + a*x)
448     return (div_factor * mul_factor * (1 - (tanh_term*tanh_term))) + a
449
450 # Method to return batches from data and labels after shuffling
451 def get_batches(X, Y, batch_size):
452     N = X.shape[0]
453     batch_X = []
454     batch_Y = []
455     count = 0
456     X, Y = shuffle(X, Y, random_state=0)
457     while count + batch_size <= N:
458         batch_X.append(X[count:count+batch_size, :])
459         one_hot = get_one_hot_representation(Y[count:count+batch_size], 10)
460         batch_Y.append(one_hot)
461         count += batch_size
462     return batch_X, batch_Y
463
464 # Convert numerical y value (0-9) to a one-hot representation
465 def get_one_hot_representation(Y, C=10):
466     one_hot = numpy.zeros((Y.shape[0], C))
467     for i in range(Y.shape[0]):
468         one_hot[i, Y[i]] = 1.0
469     return one_hot
470
471 # Layer specific forward propogation
472 def forward_prop(W, Z_prev, layer_no, num_layers,
473                 batch_size, layer_config, activation=f_tanh):
474     #Fprime = numpy.zeros((layer_config[layer_no], batch_size))
475     Z = activation(numpy.dot(Z_prev, W))
476     if layer_no == num_layers - 1:

```



```

477         return Z, []
478     else:
479         # Hidden layers need to have their Fprime values computed
480         Fprime = activation(Z, derivativeReqd=True).T
481         # Add bias terms for the hidden layers
482         Z = numpy.append(numpy.ones((Z.shape[0], 1)), Z, axis=1)
483         return Z, Fprime
484
485 # Forward propagation begins
486 def forward_prop_for_all_layers(W, Z, train_data, num_layers,
487                                batch_size, layer_config):
488     Z[0] = train_data
489     Fprime = []
490     for i in range(1, num_layers - 1):
491         Z[i], Fprime1 = forward_prop(W[i - 1], Z[i - 1], i, num_layers,
492                                     batch_size, layer_config)
493         Fprime.append(Fprime1)
494     # Separate call to send f_softmax as the activation function parameter
495     Z[-1], Fprime1 = forward_prop(W[-1], Z[-2], num_layers - 1,
496                                   num_layers, batch_size, layer_config, f_softmax)
497     return Z, Fprime
498
499 # Backpropagation step
500 def backprop(y, t, num_layers, Fprime, delta, W):
501     delta[-1] = (t - y).T
502     for i in range(num_layers - 2, 0, -1):
503         # Remove the bias column before operating on W
504         W1 = W[i][1:, :]
505         temp = numpy.dot(W1, delta[i])
506         delta[i - 1] = numpy.multiply(temp, Fprime[i - 1])
507     return delta
508
509 # Update the weight vectors
510 def update_weights(learning_rate, num_layers, Z, delta, W,
511                   momentum, prev_del_w):
512     ret_val_W_prev = []
513     for i in range(0, num_layers - 1):
514         W_grad = (learning_rate * (numpy.dot(delta[i], Z[i])).T)
515         W_grad /= (len(Z[i]))
516         ret_val_W_prev.append(W[i])
517         W[i] += W_grad + momentum * (W[i] - prev_del_w[i])
518     return W, ret_val_W_prev
519
520 # Hard coded forward propagation for 2 layer NN - only for testing
521 def hard_code_forward_prop(W, Z, num_layers,
522                             batch_size, layer_config):
523     A = []
524     A.append(numpy.dot(Z[0], W[0]))
525     Z[1] = f_sigmoid(A[0])
526     Fprime = []
527     Fprime.append(f_sigmoid(Z[1], derivativeReqd=True).T)
528     Z[1] = numpy.append(numpy.ones((Z[1].shape[0], 1)), Z[1], axis=1)
529     A.append(numpy.dot(Z[1], W[1]))
530     Z[2] = f_softmax(A[1])
531     return Z, Fprime
532
533 # train the model
534 def fit(X, y, X_test, y_test, X_val, y_val, iterations, learning_rate,
535        num_layers, W, Z, Z_val, Z_test, delta,
536        batch_size, layer_config, batch_size_val, batch_size_test,
537        momentum, prev_del_W):
538
539     train_data_batches, train_label_batches = get_batches(X, y,
540                                                            batch_size)
541     val_data_batches, val_label_batches = get_batches(X_val, y_val,

```

```

542                                     batch_size_val)
543 test_data_batches , test_label_batches = get_batches(X_test , y_test ,
544                                                     batch_size_test)
545
546 percent_correct_train = []
547 percent_correct_test = []
548 weights_array= []
549 val_error_array = []
550 test_error = 0.0
551 stopping_threshold = 20
552 W_opt = W
553
554 for t in range(iterations):
555     # Train for particular iteration
556     for i in range(len(train_label_batches)):
557         batch_data = train_data_batches[i]
558         batch_labels = train_label_batches[i]
559         # Forward Propagation
560         Z_updated, Fprime = forward_prop_for_all_layers(W, Z,
561                                                         batch_data ,
562                                                         num_layers ,
563                                                         batch_size ,
564                                                         layer_config)
565         # Back-propagation
566         delta = backprop(Z_updated[-1], batch_labels , num_layers ,
567                         Fprime, delta , W)
568         # Update the weights
569         W, prev_del_W = update_weights(learning_rate , num_layers ,
570                                       Z_updated, delta , W, momentum, prev_del_W)
571
572     # Check error on training data for this iteration
573     # and add to plot array
574     train_error = find_misclassification_error(train_data_batches ,
575                                               train_label_batches ,
576                                               Z, W)
577     percent_correct_train.append(((1 - train_error/len(y))*100))
578
579     # Check error on validation data for this iteration
580     val_error = find_misclassification_error(val_data_batches ,
581                                             val_label_batches ,
582                                             Z_val, W)
583     print (" Validation error = " + str(val_error/len(y_val))
584           + " iteration number = " + str(t))
585
586     weights_array.append(W)
587     val_error_array.append(val_error)
588     W_opt = W
589
590     # Check error on Test Data and add to test plot array
591     test_error = find_misclassification_error(test_data_batches ,
592                                             test_label_batches ,
593                                             Z_test, W)
594     percent_correct_test.append(((1 - test_error/len(y_test))*100))
595
596     # Setting threshold of minimum 15 iterations before we abort
597     if (early_stop_reqd(t, stopping_threshold , val_error_array)):
598         W_opt = weights_array[numpy.argmin(val_error_array)]
599         break
600
601     # Check error on Test Data with final weight vector chosen
602     test_error = find_misclassification_error(test_data_batches ,
603                                             test_label_batches ,
604                                             Z_test, W_opt)
605     print (" Test error = " + str(test_error/len(y_test)))
606     plotly_graphs(percent_correct_train , percent_correct_test)

```

```

607
608 def early_stop_reqd(t, stopping_threshold, val_error_array):
609     if (t > stopping_threshold):
610         count = 0
611         for index in range(t - 1, t-stopping_threshold, -1):
612             if (count < stopping_threshold
613                 and val_error_array[index] >= val_error_array[index - 1]):
614                 count += 1
615             else:
616                 return False
617         if (count >= stopping_threshold - 1):
618             return True
619     return False
620
621 # Find the misclassification error given batches of labels and data
622 def find_misclassification_error(data_batches, label_batches, Z, W):
623     error = 0.0
624     for i in range(len(label_batches)):
625         batch_data = data_batches[i]
626         batch_labels = label_batches[i]
627         Z_updated, Fprime_test = forward_prop_for_all_layers(W,
628                                                                Z,
629                                                                batch_data,
630                                                                num_layers,
631                                                                batch_size_test,
632                                                                layer_config)
633         y_pred = numpy.argmax(Z_updated[-1], axis=1)
634         error += numpy.sum(1 - batch_labels[numpy.arange(len(batch_labels)),
635                                                                y_pred])
636     return error
637
638 # Plot Percent correct graphs for testing and training data using Plotly
639 def plotly_graphs(percent_correct_train, percent_correct_test):
640     pyl.sign_in('chetang', 'vi17vTAuCSWt2lEZvaH9')
641     trace = []
642     graph_y = []
643     graph_y.append(percent_correct_train)
644     graph_y.append(percent_correct_test)
645     for i in range(2):
646         name = "Training Data"
647         if i == 1:
648             name = "Test Data"
649         y1 = graph_y[i]
650         x1 = [j+1 for j in range(len(y1))]
651         trace1 = go.Scatter(
652             x=x1,
653             y=y1,
654             name = name,
655             connectgaps=True
656         )
657         trace.append(trace1)
658     data = trace
659     fig = dict(data=data)
660     pyl.iplot(fig, filename='5b_232_232HU_Point1_128BS_point9M')
661
662 # -----Main function-----
663
664 if __name__ == "__main__":
665     numpy.random.seed(0)
666     learning_rate = 0.22
667     momentum = 0.0
668     N = 60000
669     N_test = 10000
670     iteration = 2000
671     X, y, X_test, y_test, X_validation, y_validation = get_data(N,

```

```

672                                     N_test ,
673                                     True)
674     batch_size = 50000
675     batch_size_val = 10000
676     batch_size_test = 10000
677
678     # Layers - 784, 300, 10
679     layer_config = [784, 300, 10]
680     num_layers = len(layer_config)
681     W = []
682     Z = []
683     Z_test = []
684     Z_val = []
685     delta = []
686     prev_del_W = []
687
688     for i in range(num_layers):
689         # common for all layers except input layer
690         if i != 0:
691             Z.append(numpy.zeros((batch_size , layer_config[i])))
692             Z_test.append(numpy.zeros((batch_size_test , layer_config[i])))
693             Z_val.append(numpy.zeros((batch_size_val , layer_config[i])))
694             delta.append(numpy.zeros((batch_size , layer_config[i])))
695         else:
696             Z.append(numpy.zeros((batch_size , layer_config[i]+1)))
697             Z_test.append(numpy.zeros((batch_size_test , layer_config[i]+1)))
698             Z_val.append(numpy.zeros((batch_size_val , layer_config[i]+1)))
699         # For all layers except output
700         if i != num_layers - 1:
701             W.append(numpy.random.normal(size=[layer_config[i]+1,
702                                                layer_config[i+1]],
703                                           loc = 0.0,
704                                           #scale = 0.1))
705                                           scale = 1.0/math.sqrt(layer_config[i]))
706             prev_del_W.append(numpy.zeros((layer_config[i]+1,
707                                           layer_config[i+1])))
708
709     fit(X, y, X_test, y_test, X_validation, y_validation, iteration,
710         learning_rate, num_layers, W,
711         Z, Z_val, Z_test, delta,
712         batch_size, layer_config,
713         batch_size_val, batch_size_test, momentum, prev_del_W)
714
715 3d.py
716
717 import numpy
718 from LoadMNIST import load_mnist
719 from sklearn.utils import shuffle
720 import plotly.plotly as pyl
721 import plotly.graph_objs as go
722 import copy
723 #import bigfloat
724
725 #-----Utility functions-----
726
727 def get_data(N=60000, N_test=10000, validationReqd = True):
728     # Load MNIST data using libraries available
729     training_data, training_labels = load_mnist('training')
730     test_data, test_labels = load_mnist('testing')
731
732     # Training_data is N x 784 matrix
733     training_data = flattenAndNormalize(N, 784, training_data)
734     training_labels = training_labels[:N]
735     test_data = flattenAndNormalize(N_test, 784, test_data)

```

```

735     test_labels = test_labels[:N_test]
736
737     # Adding column of 1s for bias
738     training_data = addOnesColAtStart(training_data)
739     test_data = addOnesColAtStart(test_data)
740
741     if (validationReqd):
742         # Last 10% of training data size will be considered as the validation set
743         N_validation = int(N / 6.0)
744         validation_data = training_data[N-N_validation:N]
745         validation_labels = training_labels[N-N_validation:N]
746         N=N-N_validation
747     else:
748         validation_data = []
749         validation_labels = []
750     #update training data to remove validation data
751     training_data = training_data[:N]
752     training_labels = training_labels[:N]
753
754     return training_data, training_labels, test_data, test_labels, validation_data, validation_labels
755
756 def flattenAndNormalize(rows, cols, twoDArr):
757     flattened_arr = numpy.zeros(shape=(rows, cols))
758     for row in range(0, rows):
759         i=0
760         for element in twoDArr[row]:
761             for ell in element:
762                 flattened_arr[row][i] = ell
763                 i = i+1
764     return flattened_arr
765
766 def addOnesColAtStart(matrix):
767     Ones = numpy.ones(len(matrix))
768     newMatrix = numpy.c_[Ones, matrix]
769     return newMatrix
770
771 # Single method for calculating sigmoid (logisitic) activation
772 # and its derivative
773 def f_sigmoid(X, derivativeReqd=False):
774     if not derivativeReqd:
775         return 1.0 / (1 + numpy.exp(-X))
776         #return 1.0 / (1 + bigfloat.exp(-X, bigfloat.precision(100)))
777     return numpy.multiply(f_sigmoid(X), (1 - f_sigmoid(X)))
778
779 # Method to calculate the softmax activation
780 def f_softmax(X):
781     Z = numpy.sum(numpy.exp(X), axis=1)
782     Z = Z.reshape(Z.shape[0], 1)
783     return numpy.exp(X) / Z
784
785 # Single method for tanh sigmoid and its derivative
786 def f_tanh(x, a=0, derivativeReqd = False):
787     mul_factor = 1.7159
788     div_factor = 2.0/3.0
789     tanh_term = numpy.tanh(div_factor*x)
790     if not derivativeReqd:
791         return (mul_factor * tanh_term + a*x)
792     return (div_factor * mul_factor * (1 - (tanh_term*tanh_term))) + a
793
794 # Method to return batches from data and labels after shuffling
795 def get_batches(X, Y, batch_size):
796     N = len(X)#X.shape[0]
797     batch_X = []
798     batch_Y = []
799     count = 0

```

```

800     X, Y = shuffle(X, Y, random_state=0)
801     while count + batch_size <= N:
802         batch_X.append(X[count:count+batch_size, :])
803         one_hot = get_one_hot_representation(Y[count:count+batch_size], 10)
804         batch_Y.append(one_hot)
805         count += batch_size
806     return batch_X, batch_Y
807
808 # Convert numerical y value (0-9) to a one-hot representation
809 def get_one_hot_representation(Y, C=10):
810     one_hot = numpy.zeros((Y.shape[0], C))
811     for i in range(Y.shape[0]):
812         one_hot[i, Y[i]] = 1.0
813     return one_hot
814
815 # Layer specific forward propagation
816 def forward_prop(W, Z_prev, layer_no, num_layers,
817                 batch_size, layer_config, activation=f_sigmoid):
818     #Fprime = numpy.zeros((layer_config[layer_no], batch_size))
819     Z = activation(numpy.dot(Z_prev, W))
820     if layer_no == num_layers - 1:
821         return Z, []
822     else:
823         # Hidden layers need to have their Fprime values computed
824         Fprime = activation(Z, derivativeReqd=True).T
825         # Add bias terms for the hidden layers
826         Z = numpy.append(numpy.ones((Z.shape[0], 1)), Z, axis=1)
827         return Z, Fprime
828
829 # Forward propagation begins
830 def forward_prop_for_all_layers(W, Z, train_data, num_layers,
831                                batch_size, layer_config):
832     Z[0] = train_data
833     Fprime = []
834     for i in range(1, num_layers-1):
835         Z[i], Fprime1 = forward_prop(W[i-1], Z[i-1], i, num_layers,
836                                     batch_size, layer_config)
837         Fprime.append(Fprime1)
838     # Separate call to send f_softmax as the activation function parameter
839     Z[-1], Fprime1 = forward_prop(W[-1], Z[-2], num_layers-1,
840                                   num_layers, batch_size, layer_config, f_softmax)
841     return Z, Fprime
842
843 # Backpropagation step
844 def backprop(y, t, num_layers, Fprime, delta, W):
845     delta[-1] = (t - y).T
846     for i in range(num_layers-2, 0, -1):
847         # Remove the bias column before operating on W
848         W1 = W[i][1:, :]
849         temp = numpy.dot(W1, delta[i])
850         delta[i-1] = numpy.multiply(temp, Fprime[i-1])
851     return delta
852
853 # Update the weight vectors
854 def update_weights(learning_rate, num_layers, Z, delta, W,
855                   momentum, prev_del_w, size):
856     ret_val_W_prev = []
857     gradient_bp = []
858     for i in range(0, num_layers-1):
859         W_grad = (learning_rate*(numpy.dot(delta[i], Z[i])).T + momentum*prev_del_w[i])/float(size)
860         gradient_bp.append(W_grad)
861         W[i] += W_grad
862         ret_val_W_prev.append(-1*W_grad)
863     return W, ret_val_W_prev, gradient_bp
864

```

```

865 # Update the weight vectors
866 def update_weights_analytical_gradient(learning_rate, num_layers, Z, delta, W, t, activation_h
867     e = 0.01
868     gradient = []
869     W_save = copy.deepcopy(W)
870     for x in range(0, num_layers - 1):
871         gradient.append(numpy.zeros((W[x].shape[0], W[x].shape[1])))
872         for i in range(0, W[x].shape[0]):
873             for j in range(0, W[x].shape[1]):
874                 W_pos = copy.deepcopy(W_save)
875                 W_neg = copy.deepcopy(W_save)
876
877                 W_pos[x][i][j] = W_save[x][i][j] + e
878                 W_neg[x][i][j] = W_save[x][i][j] - e
879                 #print(W_save[x][i][j] - W_pos[x][i][j])
880                 Z_pos = copy.deepcopy(Z)
881                 Z_neg = copy.deepcopy(Z)
882                 for k in range(1, num_layers - 1):
883                     Z_pos[k] = activation_hidden(numpy.dot(Z_pos[k-1], W_pos[k-1]))
884                     Z_neg[k] = activation_hidden(numpy.dot(Z_neg[k-1], W_neg[k-1]))
885                     Z_pos[k] = numpy.append(numpy.ones((Z_pos[k].shape[0], 1)), Z_pos[k], axis
886                     Z_neg[k] = numpy.append(numpy.ones((Z_neg[k].shape[0], 1)), Z_neg[k], axis
887
888                 Z_pos[-1] = activation_output(numpy.dot(Z_pos[-2], W_pos[-1]))
889                 Z_neg[-1] = activation_output(numpy.dot(Z_neg[-2], W_neg[-1]))
890
891                 gradient[x][i][j] = ( learning_rate * (numpy.dot(t, numpy.log(Z_pos[-1].T)) -
892
893     return gradient
894
895 # Hard coded forward propagation for 2 layer NN - only for testing
896 def hard_code_forward_prop(W, Z, num_layers,
897     batch_size, layer_config):
898     A = []
899     A.append(numpy.dot(Z[0], W[0]))
900     Z[1] = f_sigmoid(A[0])
901     Fprime = []
902     Fprime.append(f_sigmoid(Z[1], derivativeReqd=True).T)
903     Z[1] = numpy.append(numpy.ones((Z[1].shape[0], 1)), Z[1], axis=1)
904     A.append(numpy.dot(Z[1], W[1]))
905     Z[2] = f_softmax(A[1])
906     return Z, Fprime
907
908 # train the model
909 def fit(X, y, X_test, y_test, X_val, y_val, iterations, learning_rate,
910     num_layers, W, Z, Z_val, Z_test, delta,
911     batch_size, layer_config, batch_size_val, batch_size_test,
912     momentum, prev_del_W):
913
914     train_data_batches, train_label_batches = get_batches(X, y,
915                                                         batch_size)
916     val_data_batches, val_label_batches = get_batches(X_val, y_val,
917                                                         batch_size_val)
918     test_data_batches, test_label_batches = get_batches(X_test, y_test,
919                                                         batch_size_test)
920
921     test_error = 0.0
922
923     for t in range(iterations):
924         # Train for particular iteration
925         for i in range(len(train_label_batches)):
926             batch_data = train_data_batches[i]
927             batch_labels = train_label_batches[i]
928             # Forward Propagation
929

```

```

930         Z_updated, Fprime = forward_prop_for_all_layers(W, Z,
931                                                         batch_data,
932                                                         num_layers,
933                                                         batch_size,
934                                                         layer_config)
935     # Back-propagation
936     delta = backprop(Z_updated[-1], batch_labels, num_layers,
937                     Fprime, delta, W)
938     # Update the weights
939     W_true, prev_del_W, gradient_bp = update_weights(learning_rate, num_layers,
940                                                         Z_updated, delta, W, momentum, prev_del_W, len(batch_labels))
941
942     # Update the weights
943     gradient_ana = update_weights_analytical_gradient(learning_rate, num_layers,
944                                                         Z_updated, delta, W, batch_labels)
945
946     gradient_bp = [gradient_bp[i] - gradient_ana[i] for i in range(0, len(gradient_bp))]
947     print(gradient_bp)
948
949     diff_grad_ij = numpy.fabs(gradient_bp[0] - gradient_ana[0])
950     diff_grad_jk = numpy.fabs(gradient_bp[1] - gradient_ana[1])
951
952     max_diff_ij = numpy.amax(diff_grad_ij)
953     max_diff_jk = numpy.amax(diff_grad_jk)
954     mean_diff_ij = numpy.mean(diff_grad_ij)
955     mean_diff_jk = numpy.mean(diff_grad_jk)
956     print("max_diff_ij: " + str(max_diff_ij))
957     print("max_diff_jk: " + str(max_diff_jk))
958
959     print("mean_diff_ij: " + str(mean_diff_ij))
960     print("mean_diff_jk: " + str(mean_diff_jk))
961
962     numpy.savetxt("Gradient_diff.csv", gradient_bp[0], delimiter=",")
963     numpy.savetxt("Gradient_diff2.csv", gradient_bp[1], delimiter=",")
964
965     print("Test error = " + str(test_error/len(y_test)))
966
967 def plotly_graphs(percent_correct_train, percent_correct_test):
968     pyl.sign_in('chetang', 'v1l7vTAuCSWt2lEZvaH9')
969     trace = []
970     graph_y = []
971     graph_y.append(percent_correct_train)
972     graph_y.append(percent_correct_test)
973     for i in range(2):
974         name = "Training Data"
975         if i == 1:
976             name = "Test Data"
977         y1 = graph_y[i]
978         x1 = [j+1 for j in range(len(y1))]
979         trace1 = go.Scatter(
980             x=x1,
981             y=y1,
982             name = name,
983             connectgaps=True
984         )
985         trace.append(trace1)
986     data = trace
987     fig = dict(data=data)
988     pyl.iplot(fig, filename='percentCorrect-connectgaps_1')
989
990 #-----Main function-----
991
992 if __name__ == "__main__":
993     numpy.random.seed(0)
994     learning_rate = 0.001

```



```

995     momentum = 0.0
996     N = 1
997     N_test = 1
998     iteration = 1
999     X, y, X_test, y_test, X_validation, y_validation = get_data(N,
1000                                                                N_test,
1001                                                                False)
1002     X = X/255.0
1003     #X_test = X_test/255.0
1004     #X_validation = X_validation/255.0
1005     #X = X - numpy.mean(X, axis=0)[numpy.newaxis, :]
1006     for r in range(0,X.shape[0]):
1007         X[r] = X[r] - sum(X[r])/X.shape[1]
1008     #X_test = X_test - numpy.mean(X_test, axis=0)[numpy.newaxis, :]
1009     #X_validation = X_validation - numpy.mean(X_validation, axis=0)[numpy.newaxis, :]
1010
1011     batch_size = 1
1012     batch_size_val = 1
1013     batch_size_test = 1
1014
1015     # Layers = 784, 100, 10
1016     layer_config = [784, 300, 10]
1017     num_layers = len(layer_config)
1018     W = []
1019     Z = []
1020     Z_test = []
1021     Z_val = []
1022     delta = []
1023     prev_del_W = []
1024
1025     for i in range(num_layers):
1026         # common for all layers except input layer
1027         if i != 0:
1028             Z.append(numpy.zeros((batch_size, layer_config[i])))
1029             Z_test.append(numpy.zeros((batch_size_test, layer_config[i])))
1030             Z_val.append(numpy.zeros((batch_size_val, layer_config[i])))
1031             delta.append(numpy.zeros((batch_size, layer_config[i])))
1032         else:
1033             Z.append(numpy.zeros((batch_size, layer_config[i]+1)))
1034             Z_test.append(numpy.zeros((batch_size_test, layer_config[i]+1)))
1035             Z_val.append(numpy.zeros((batch_size_val, layer_config[i]+1)))
1036         # For all layers except output
1037         if i != num_layers - 1:
1038             W.append(numpy.random.normal(size=[layer_config[i]+1,
1039                                                layer_config[i+1]],
1040                                          loc = 0.0,
1041                                          scale=0.001))
1042             prev_del_W.append(numpy.zeros((layer_config[i]+1,
1043                                           layer_config[i+1])))
1044
1045     fit(X, y, X_test, y_test, X_validation, y_validation, iteration,
1046         learning_rate, num_layers, W,
1047         Z, Z_val, Z_test, delta,
1048         batch_size, layer_config,
1049         batch_size_val, batch_size_test, momentum, prev_del_W)

```