# Analysis on Logistic and Softmax Regression Using MNIST Dataset

**Chetan Gandotra**
UC San Diego
9500 Gilman Drive
cgandotr@ucsd.edu

**Rishabh Misra**
UC San Diego
9500 Gilman Drive
r1misra@ucsd.edu

## Abstract

This report discusses the first programming assignment of course CSE 253: Neural Networks and Pattern Recognition, its solutions and the inferences. MNIST dataset was used and the hand-written digits in it were classified using Logistic and Softmax Regressions. Under Logistic Regression, two-way classification was performed on specific digits (2's and 3's, 2's and 8's). An accuracy of more than 97% was achieved on both of these subsets of data using Logistic Regression. For Softmax Regression, we performed a ten-way classification (for all digits from 0 to 9) and achieved an accuracy of 87.65% on the test set.

## 1 Derivation of Gradient for Logistic Regression

### 1.1 Introduction

The problem statement here is to find the gradient of the cost function. The error function for the logistic regression follows from the negative log likelihood, which can be written as:

$$E(w) = -\sum_{n=1}^{N}\{t^n \ln y^n + (1 - t^n)\ln(1 - y^n)\}$$

### 1.2 Methodology

In this section, we will derive the gradient of cost function, which will be used in the later parts of this report. To find the optimal weight parameters, we need to take the partial derivative of the error function with respect to $w_j$ as follows:

$$\frac{\partial E(w)}{\partial w_j} = -\sum_{n=1}^{N}\left[t^n\frac{\partial \ln y^n}{\partial w_j} + (1 - t^n)\frac{\partial \ln(1 - y^n)}{\partial w_j}\right]$$

$$\frac{\partial E(w)}{\partial w_j} = -\sum_{n=1}^{N}\left[\frac{t^n}{y^n}\frac{\partial y^n}{\partial w_j} + \frac{(1 - t^n)}{1 - y^n}\frac{\partial(1 - y^n)}{\partial w_j}\right]$$

Since $y^n = \sigma(\mathbf{w}.\mathbf{x^n})$, the derivative can be written as:

$$\frac{\partial E(w)}{\partial w_j} = -\sum_{n=1}^{N}\left[\frac{t^n}{\sigma(\mathbf{w}.\mathbf{x^n})}\frac{\partial\sigma(\mathbf{w}.\mathbf{x^n})}{\partial w_j} + \frac{(1 - t^n)}{1 - \sigma(\mathbf{w}.\mathbf{x^n})}\frac{\partial(1 - \sigma(\mathbf{w}.\mathbf{x^n}))}{\partial w_j}\right]$$

12 Using the following properties of sigmoid function

$$\sigma(-\mathbf{x}) = 1 - \sigma(\mathbf{x}) \tag{1}$$

$$\frac{\partial\sigma(-\mathbf{x})}{\partial x} = \sigma(\mathbf{x})\sigma(-\mathbf{x}) \tag{2}$$

the derivative can be written as:

$$\frac{\partial E(w)}{\partial w_j} = -\sum_{n=1}^{N}\left[\frac{t^n}{\sigma(\mathbf{w}.\mathbf{x^n})}\sigma(\mathbf{w}.\mathbf{x^n})\sigma(-\mathbf{w}.\mathbf{x^n})x_j^n - \frac{(1-t^n)}{\sigma(-\mathbf{w}.\mathbf{x^n})}\sigma(-\mathbf{w}.\mathbf{x^n})\sigma(\mathbf{w}.\mathbf{x^n})x_j^n\right]$$

$$\frac{\partial E(w)}{\partial w_j} = -\sum_{n=1}^{N}x_j^n\left[t^n\sigma(-\mathbf{w}.\mathbf{x^n}) - (1-t^n)\sigma(\mathbf{w}.\mathbf{x^n})\right]$$

Using (1) we get,

$$\frac{\partial E(w)}{\partial w_j} = -\sum_{n=1}^{N}x_j^n\left[t^n(1-\sigma(\mathbf{w}.\mathbf{x^n})) - (1-t^n)\sigma(\mathbf{w}.\mathbf{x^n})\right]$$

Solving the above equation, we get:

$$\frac{\partial E(w)}{\partial w_j} = -\sum_{n=1}^{N}x_j^n\left[t^n - \sigma(\mathbf{w}.\mathbf{x^n})\right]$$

or

$$-\frac{\partial E(w)}{\partial w_j} = \sum_{n=1}^{N}(t^n - y^n)x_j^n$$

13 ## 1.3   Results

Hence, from the above derivation, it follows that for $n^{th}$ sample, the gradient of error can be written as:

$$-\frac{\partial E^n(w)}{\partial w_j} = (t^n - y^n)x_j^n$$

14 ## 1.4   Discussion

15 The expression takes the difference between true label and predicted label and weigh it by the input
16 data value. It makes sense because if there is a stark difference between the true and the predicted
17 labels, the gradient value would be large. Thus, the corresponding component of the weight vector
18 would be adjusted quickly in the direction of gradient to reduce the loss.

## 2 Derivation of Gradient for Softmax Regression

### 2.1 Introduction

In this section, the focus is to find the gradient of the loss function of Softmax Regression - *E(w)*. The error function for the softmax regression follows from the negative log likelihood, which can be written as:

$$E(w) = -\sum_{n=1}^{N} \sum_{k'=1}^{C} t_{k'}^{n} \ln y_{k'}^{n}$$

### 2.2 Methodology

To find the optimal weight parameters for each class, we need to take the partial derivative of the error function with respect to $w_{jk}$ as follows:

$$\frac{\partial E(w)}{\partial w_{jk}} = -\sum_{n=1}^{N} \sum_{k'=1}^{C} \left[ t_{k'}^{n} \frac{\partial \ln y_{k'}^{n}}{\partial w_{jk}} \right]$$

$$-\frac{\partial E(w)}{\partial w_{jk}} = \sum_{n=1}^{N} \left[ \frac{t_{k}^{n}}{y_{k}^{n}} \frac{\partial y_{k}^{n}}{\partial w_{jk}} \right] + \sum_{k' \neq k} \left[ \frac{t_{k'}^{n}}{y_{k'}^{n}} \frac{\partial y_{k'}^{n}}{\partial w_{jk}} \right] \tag{3}$$

Since,

$$y_{k}^{n} = \frac{e^{\mathbf{w_k} \cdot \mathbf{x^n}}}{\sum_{l=1}^{C} e^{\mathbf{w_l} \cdot \mathbf{x^n}}}$$

the derivatives $\frac{\partial y_{k}^{n}}{\partial w_{jk}}$ and $\frac{\partial y_{k'}^{n}}{\partial w_{jk}}$ can be written as:

$$\frac{\partial y_{k}^{n}}{\partial w_{jk}} = \frac{e^{\mathbf{w_k} \cdot \mathbf{x^n}}}{\sum_{l=1}^{C} e^{\mathbf{w_l} \cdot \mathbf{x^n}}} x_{j}^{n} - e^{\mathbf{w_k} \cdot \mathbf{x^n}} \left[ \frac{1}{(\sum_{l=1}^{C} e^{\mathbf{w_l} \cdot \mathbf{x^n}})^2} \right] e^{\mathbf{w_k} \cdot \mathbf{x^n}} x_{j}^{n}$$

$$\frac{\partial y_{k}^{n}}{\partial w_{jk}} = \left( y_{k}^{n} - (y_{k}^{n})^2 \right) x_{j}^{n} \tag{4}$$

$$\frac{\partial y_{k'}^{n}}{\partial w_{jk}} = -e^{\mathbf{w_{k'}} \cdot \mathbf{x^n}} \left[ \frac{1}{(\sum_{l=1}^{C} e^{\mathbf{w_l} \cdot \mathbf{x^n}})^2} \right] e^{\mathbf{w_{k'}} \cdot \mathbf{x^n}} x_{j}^{n}$$

$$\frac{\partial y_{k'}^{n}}{\partial w_{jk}} = -\left( y_{k'}^{n} \right)^2 x_{j}^{n} \tag{5}$$

Substituting (4) and (5) in (3), we get

$$-\frac{\partial E(w)}{\partial w_{jk}} = \sum_{n=1}^{N} \left[ \frac{t_{k}^{n}}{y_{k}^{n}} \left( y_{k}^{n} - (y_{k}^{n})^2 \right) x_{j}^{n} \right] - \sum_{k' \neq k} \left[ \frac{t_{k'}^{n}}{y_{k'}^{n}} \left( y_{k'}^{n} \right)^2 x_{j}^{n} \right]$$

$$-\frac{\partial E(w)}{\partial w_{jk}} = \sum_{n=1}^{N} \left[ t_{k}^{n} \left( 1 - y_{k}^{n} \right) x_{j}^{n} \right] - \sum_{k' \neq k} \left[ t_{k'}^{n} y_{k'}^{n} x_{j}^{n} \right] \tag{6}$$

Now, for any sample, only one of the C labels in $t^n$ would be 1, and all the others would be 0. This is because the label one would be the label set, and each training example can correspond to only one label. Thus, for any sample $a$ where $t_{k}^{n}$ is 1, the derivative would be:

$$-\frac{\partial E^a(w)}{\partial w_{jk}} = \left[ (1 - y_k^a) \, x_j^a \right]$$

or it could be written as:

$$-\frac{\partial E^a(w)}{\partial w_{jk}} = \left[ (t_k^n - y_k^a) \, x_j^a \right] \tag{7}$$

For any sample $b$ where one of $t_{k'}^n$ is 1 (where $k' \neq k$), the derivative would be:

$$-\frac{\partial E^b(w)}{\partial w_{jk}} = \left[ \left( -y_k^b \right) x_j^b \right]$$

or it could be written as:

$$-\frac{\partial E^b(w)}{\partial w_{jk}} = \left[ \left( t_k^n - y_k^b \right) x_j^b \right] \tag{8}$$

Using the results of (7) and (8), (6) could be written as:

$$-\frac{\partial E(w)}{\partial w_{jk}} = \sum_{n=1}^{N} (t_k^n - y_k^n) \, x_j^n$$

## 2.3   Results

Thus, using our findings above, we can say that for $n^{th}$ sample, the derivative can be written as:

$$-\frac{\partial E^n(w)}{\partial w_{jk}} = (t_k^n - y_k^n) \, x_j^n \tag{9}$$

## 2.4   Discussion

Interestingly, the expression of gradient looks similar to that of logistic regression. In this case, the derivative takes the difference between true label and predicted label for the $k^{th}$ class and weigh it by the input data value. Again, if the difference is big between the true and the predicted labels, the gradient value would be large. Thus, the corresponding component of the weight vector would be adjusted quickly in the direction of gradient to reduce the loss.

# 3 Read in Data

## 3.1 Introduction

As mentioned in the abstract, we deal with the "MNIST" dataset in this programming assignment. The first and foremost task before operating on the data was to load it.

## 3.2 Methodology

The MNIST data was downloaded from the website at *http://yann.lecun.com/exdb/mnist/* (the same link as given in PA1). To read the data, GitHub library at *https://github.com/akosiorek/CSE/blob/master/MLCV/* was used, which returns the training and testing data in matrix form, and labels as a vectors. Operations were then performed on this data to add a column of ones (bias term) at the beginning and to extract digit specific data (2-3's and 2-8's). Also, the data was restricted to first 20k entries, 10% of which was allocated to a hold-out set. The size of test data was kept as two thousand. This was done by picking the first 2k entries from the test data returned by the library.

## 3.3 Results

Using some existing libraries on Github, we were able to extract the data into variables in Python. This data, however, consisted of the full 60k training data points and 10k testing data points. We extracted the first 20k training data points and the first 2k testing data points. 10% of the training data was designated as a hold-out set.

## 3.4 Discussion

Computation on large data sets can often be time consuming. Due to this reason, we extracted the full data and restricted the size of training, validation and test sets. This allowed faster computations throughout the programming assignment. A hold-out set acts as a dummy test set, which we use so as to improve performance of our model by testing it on the hold-out set. Good accuracy on hold-out set leads to a good accuracy on test set in general.

# 4 Logistic Regression via gradient descent

## 4.1 Introduction

In this part, we are required to use logistic regression and classify a given hand-written digit as either 2 or 3. Since logistic regression uses binary output, we say that the target is 1 if the input is from the "2" category and 0 if it is from the other category. We are required to produce the following:

1. Plot of loss function ($E$) over the training set, test set and the hold-out set
2. Plot of percent correct classification over training for the training set, the hold-out set, and the test set
3. The above two plots for digits 2 and 8
4. Display weights as images for both the classifiers (2 vs. 3 and 2 vs. 8). Plot the difference between weights as well.

## 4.2 Methodology

To plot the first graph, loss function is put against the Y-axis, and the iteration number along the X-axis. The loss function used was:

$$E(w) = -\sum_{n=1}^{N}\{t^n \ln y^n + (1-t^n)\ln(1-y^n)\}$$

This was done for all the three sets - training, test and validation. Hence, for each set, the value of $N$ and the corresponding data/labels change. This procedure was repeated for the 2's and 8's data set.

In the next graph, we plotted "percent correct" against the iteration number. The value of percent correct can be inferred by going over all the examples as test set in a way, and seeing what our model predicts on it. For every correct classification, we add one to the number of data points classified correctly. Then, at the end, we find the corresponding percentage. This is repeated at each iteration for all the three sets - training, test and validation. Again, this procedure was repeated for the 2's and 8's data set.

To display the weight vectors for both the cases and their difference, the bias term in them was dropped. This reduced the weight vector to a 784 dimensional vector, which was re-shaped into a 28 x 28 matrix and then plotted using Python.
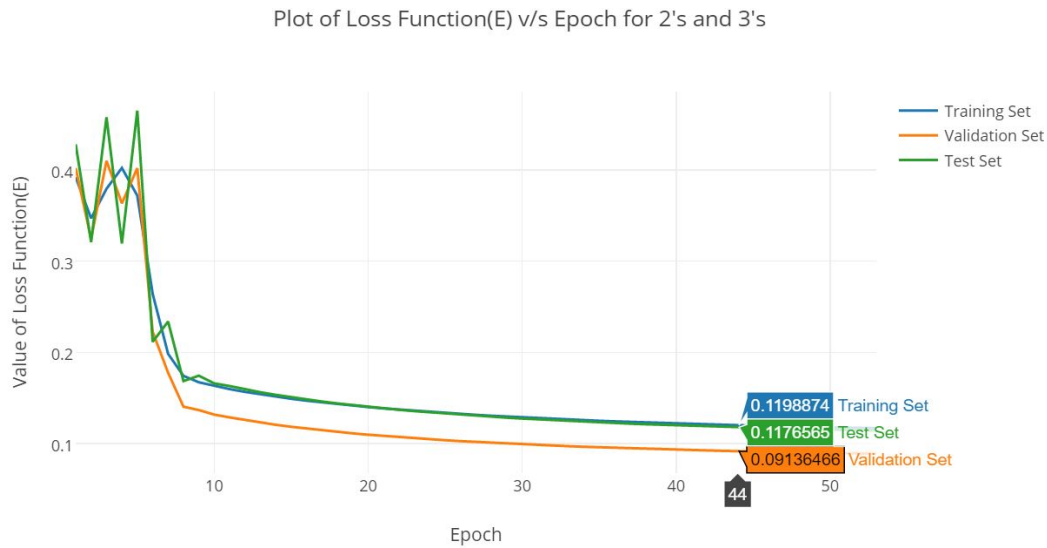
## 4.3 Results

Various results are plotted below:

Figure 1: Loss Function (E) vs Epoch for 2's and 3's

Figure 2: Loss Function calculated every 1/20 of Epoch vs Epoch for 2's and 3's
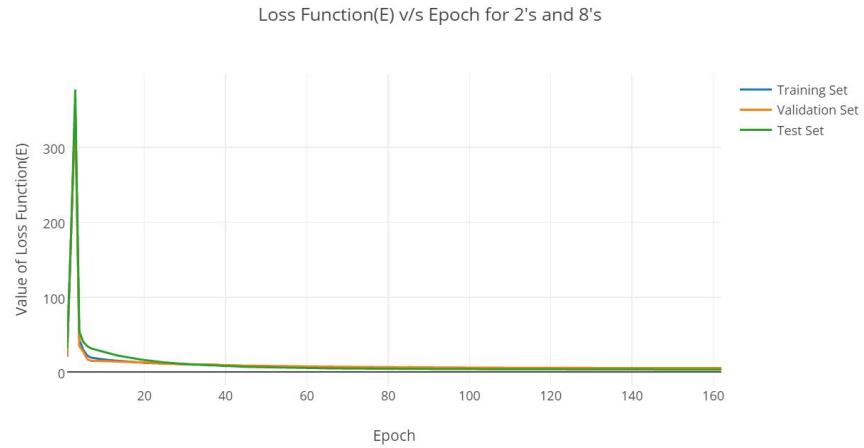
6

Loss Function(E) v/s Epoch for 2's and 8's

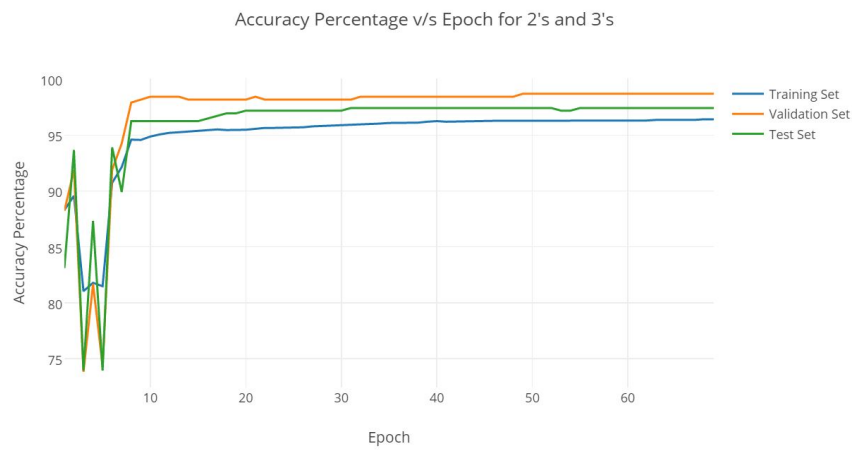Figure 3: Loss Function (E) vs Epoch for 2's and 8's

Accuracy Percentage v/s Epoch for 2's and 3's

Figure 4: Percent correct classification vs Epoch for 2's and 3's

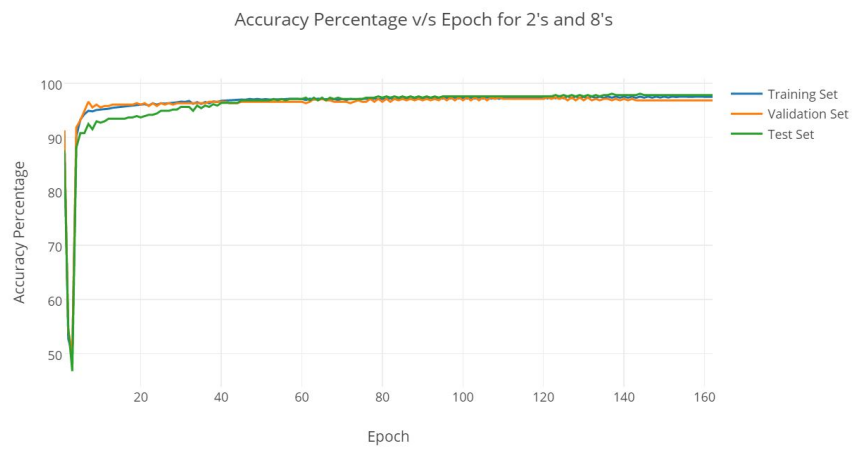Accuracy Percentage v/s Epoch for 2's and 8's

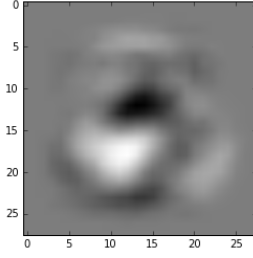Figure 5: Percent correct classification vs Epoch for 2's and 8's

Figure 6: Weights as an image for 2's and 3's
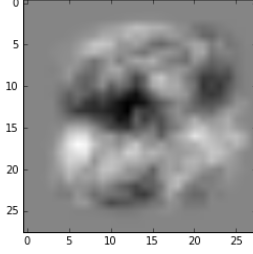


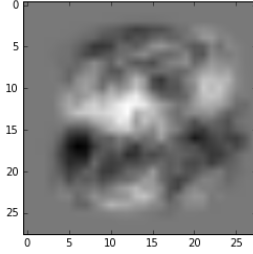Figure 7: Weights as an image for 2's and 8's



Figure 8: Difference of weights as an image for 2's and 3's and 2's and 8's

## 4.4 Discussion

Generally, we assume that the data for training and testing are generated from same underlying distribution. Thus, they have similar underlying properties. However, since test data is not available to us in the real world, we extract a portion of training data to test our model and call it a validation set. And since the validation set and the test set are generated from the same underlying distribution, the performance of our model should not vary much on validation and test set in ideal case (it would very a little bit in our case as validation set is small as compared to test set). The aforementioned behavior is reaffirmed by our experiments. From figure **??** and figure **??**, we see that curves for validation error and test error are similar and validation error hovers around the test error. Same is the case for classification of 2's and 8's in figure **??**. Thus, we can safely say that validation set closely captures how our model would perform on test data. One interesting thing to note is that in the case of batching (contained 1/20 of all points), the curve of error function is smooth . This is because the gradient is calculated only on few points at a time, which makes the gradient increase gradually towards the optimum.

We have chosen the value of the hyperparameter T as 2000 empirically. Also, for the early stopping of our algorithm, we have checked whether the error is non-decreasing on 15 iterations. The reason for choosing a bit large value is that the error on the validation set used to change in the steps of 5-10 iterations. Thus, just to insure that the error does not decrease after early stopping, we took added some buffer iterations.

In figure **??** and figure **??**, we see that the accuracy percentage increases over time. This is expected, as with more number of iterations, we minimize loss and fit the data. This means that we will be able

to classify more number of points correctly leading to higher accuracy. Also, note that the graph of 2's and 8's is more overlapping. This can be because the validation/test data is very similar to the training data.

Next, we plot the weight vectors for three cases - 2's and 3's, 2's and 8's, and the difference of these two in figure **??**, **??** and**??** respectively. For both 2's and 3's case and the 2's and 8's case, the weight vector seems like images of 2 and 3/8 have been superimposed. This seems intuitively correct as the weight vector needs to predict either of these. For the difference case, the image looks like a mirror-image of 3. This is because the 2's component is common in both the cases and must get cancelled out. Thereafter, we are left with 8's and 3's, in which superimposition cancels out too. Hence, what we are left with is the portion of 8 that was not covered by 3 and thus looks like a mirror image of 3.

## 5    Regularization

### 5.1    Introduction

Regularization is a commonly used technique to improve the model generalization. We write the regularized loss function *J(w)* as:

$$J(w) = E(w) + \lambda C(w)$$

where *C(w)* is the complexity penalty and $\lambda$ is the strength of regularization. For $\lambda = 0$, *J* reduces to *E*. Considering $L_2$ norm as the complexity penalty, we have:

$$J(w) = E(w) + \lambda ||w||^2$$

For $L_1$ norm, we have:

$$J(w) = E(w) + \lambda |w|$$

In the first part, we are expected to derive the update term for both $L_1$ and $L_2$ penalties.

Next, we are expected to plot the percent correct v/s iterations graph for different $\lambda$ values. This is followed by plotting length of weight vector v/s iterations for different $\lambda$ values. Then, we plot the final test error with each of the $\lambda$. Finally, we are expected to plot the weights as images.

### 5.2    Methodology

To derive the update term, we take derivative of this function with respect to *w*, the weight vector. Hence, we have:

$$\frac{\partial J(w)}{\partial w} = \frac{\partial E(w)}{\partial w} + \lambda \frac{\partial C(w)}{\partial w}$$

We have already calculate the first part of the equation - $\frac{\partial E(w)}{\partial w}$ in Question 1. Hence, according to the question, solving for $\frac{\partial C(w)}{\partial w}$, we have:

$$\frac{\partial C(w)}{\partial w} = \frac{||w||^2}{\partial w} = 2w$$

Therefore, we have:

$$\frac{\partial J(w)}{\partial w} = \frac{\partial E(w)}{\partial w} + 2\lambda w$$

Similarly, for $L_1$ norm as the complexity penalty, we have

$$J(w) = E(w) + \lambda |w|$$

Therefore,

$$\frac{\partial J(w)}{\partial w} = \frac{\partial E(w)}{\partial w} + \frac{\partial \lambda |w|}{\partial w}$$

$$\frac{\partial J(w)}{\partial w} = \frac{\partial E(w)}{\partial w} + \lambda \frac{\partial |w|}{\partial w}$$

Now, the entry at index *j* of partial derivative of |w| can be written as:

$$\frac{\partial |w|}{\partial w_j} = \left\{ \begin{array}{l} 1, if w_j \geq 0 \\ -1, otherwise \end{array} \right.$$

136 Thus, the value of $\frac{\partial |w|}{\partial w}$ is A vector of all one's or minues one's depending upon the sign of entries in
137 (w) vector and has the same number of elements as in *w*.

138 To plot the graphs, a similar approach as in Section 4 was undertaken. The only difference was that
139 earlier we did it for different data sets - training, test and validation. In these graphs, we always take
140 the training set and calculate percent error and length of weight vector at that iteration. This is done
141 multiple times by changing $\lambda$ values. Then, we plot the final test error for each $\lambda$ value, keeping the
142 learning rate fixed. This plot is made as a bar graph, with one bar for each $\lambda$. Finally, we plot the
143 weights as images like we did in Section 4. To display the weight vectors, the bias term was dropped.
144 This reduced the weight vector to a 784 dimensional vector, which was re-shaped into a 28 x 28
145 matrix and then plotted using Python.

## 5.3 Results

The partial derivative of loss function with $L_1$ norm regularization is:

$$\frac{\partial J(w)}{\partial w} = \frac{\partial E(w)}{\partial w} + \lambda \frac{\partial |w|}{\partial w}$$

147 where the partial derivative of |w| can be written as:

$$\frac{\partial |w|}{\partial w_j} = \left\{ \begin{array}{l} 1, if w_j \geq 0 \\ -1, otherwise \end{array} \right.$$

148 and the partial derivative in case of $L_2$ norm regularization is:

$$\frac{\partial J(w)}{\partial w} = \frac{\partial E(w)}{\partial w} + 2\lambda w$$

149 The "percent correct" value was plotted over the number of training iterations for the training set, for
150 different lambda values keeping the other hyper-parameters same.

151 Parameters used: **Penalty** = $L_2$ norm, **Learning rate** $\eta = 0.0001$



Figure 9: Accuracy Percentage v/s Epoch for $L_2$ norm

152 Now for penalty as $L_1$ norm and Learning rate $\eta = 0.0001$, we have:

Accuracy Percentage v/s Epoch for L-1 Norm



Figure 10: Accuracy Percentage v/s Epoch for $L_1$ norm

The length of weight vector against training iterations produced a graph as follows for the $L_2$ norm case:

Weight Vector Length v/s Epoch for L-2 Norm



Figure 11: Weight Vector Length v/s Epoch with $L_2$ norm

For $L_1$ norm, it becomes:

Weight Vector Length v/s Epoch for L-1 Norm



Figure 12: Weight Vector Length v/s Epoch with $L_1$ norm

156 Note that the learning rate $\eta$ used was 0.0001 in both the cases.
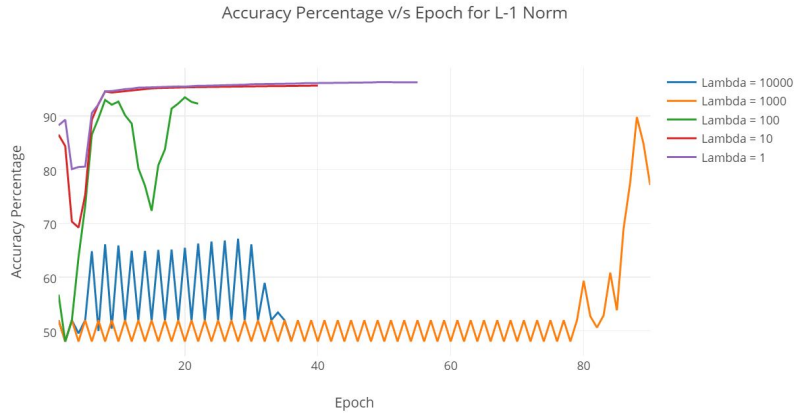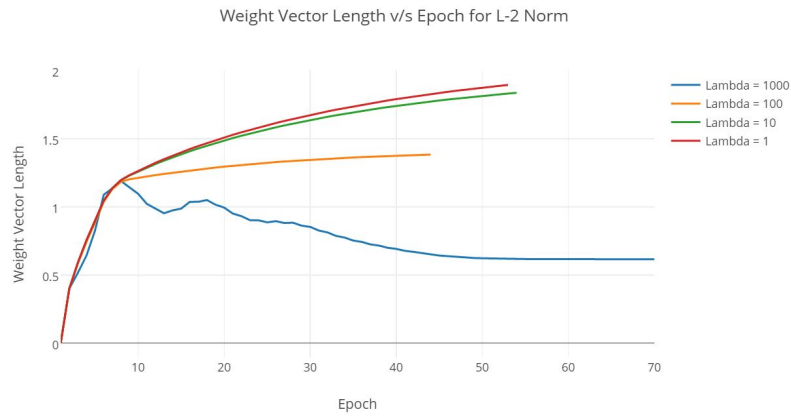
157 The plot of final test error for various $\lambda$ values with $L_2$ norm penalty is as follows:

Final Test Error for Different log(Lambda) Values



Figure 13: Final Test Error v/s Log($\lambda$) $L_2$ norm

158 The plot of final test error for various $\lambda$ values with $L_1$ norm penalty is as follows:

Final Test Error v/s Log(Lambda) for L-1 Norm



Figure 14: Final Test Error v/s Log($\lambda$) $L_1$ norm

159 Note that the learning rate $\eta$ used was 0.0001 in both the cases.

160 For **L1** norm: using learning rate $\eta$ as 0.0001 in both the cases and different $\lambda$ values, the final weight
161 vectors were plotted as images. Here are the findings:

13

Figure 15: Weight Vector Image for 2's and 3's with $L_1$ norm and $\lambda = 1$



Figure 16: Weight Vector Image for 2's and 3's with $L_1$ norm and $\lambda = 10$



Figure 17: Weight Vector Image for 2's and 3's with $L_1$ norm and $\lambda = 100$



Figure 18: Weight Vector Image for 2's and 3's with $L_1$ norm and $\lambda = 1000$

For **L1** norm: Using learning rate $\eta$ as 0.0001 in both the cases and optimal $\lambda$ value of 0.0001, the image looks like following:

14

Figure 19: Weight Vector Image for 2's and 3's with $L_1$ norm and $\lambda = 0.0001$

For **L2** norm: using learning rate $\eta$ as 0.0001 in both the cases and different $\lambda$ values, the final weight vectors were plotted as images. Here are the findings:



Figure 20: Weight Vector Image for 2's and 3's with $L_2$ norm and $\lambda = 1$



Figure 21: Weight Vector Image for 2's and 3's with $L_2$ norm and $\lambda = 10$
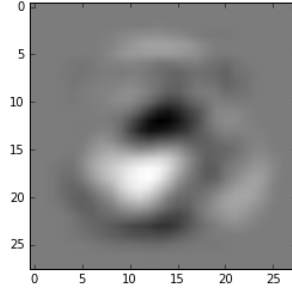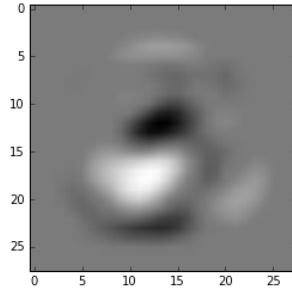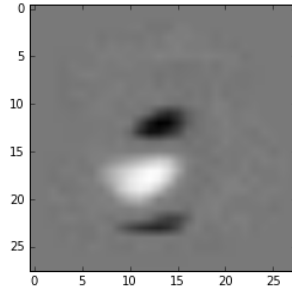


Figure 22: Weight Vector Image for 2's and 3's with $L_2$ norm and $\lambda = 100$
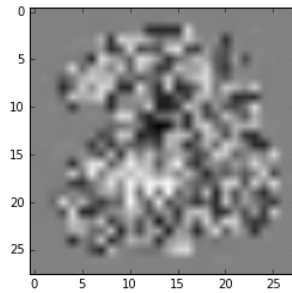
Figure 23: Weight Vector Image for 2's and 3's with $L_2$ norm and $\lambda = 1000$

For **L2** norm: Using learning rate $\eta$ as 0.0001 in both the cases and optimal $\lambda$ value of 0.0001, the image looks like following:
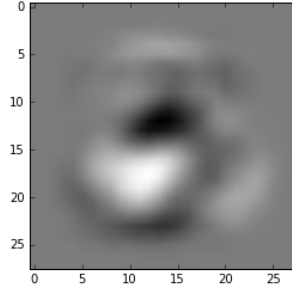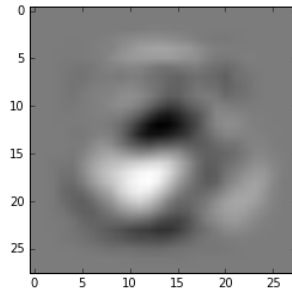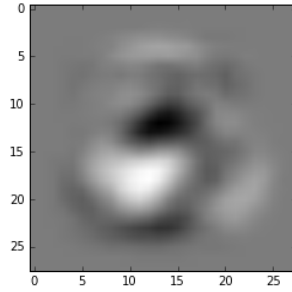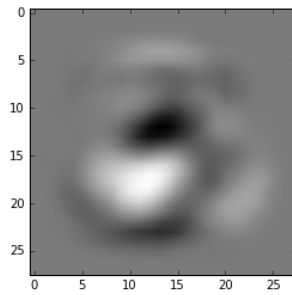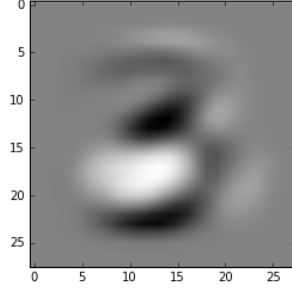


Figure 24: Weight Vector Image for 2's and 3's with $L_2$ norm and $\lambda = 0.0001$

## 5.4 Discussion

Note that the $L_1$ norm is not differentiable at 0. However, all that matters is that we can compute a subgradient / subderivative. Since it's differentiable everywhere else, we can just fill in any "reasonable" value (such as -1 or 1; we have chosen -1) for the gradient at 0. So, for the cases where the weights are negative, we would have positive regularization term that would drive them towards zero and for the cases where the weights are positive, the regularization term would be negative, thus decreasing the weights. Hence, $L_1$ tries to limit the coefficient which restricts overfitting.

Graph **??** shows how the accuracy on training data varies according to the iteration number (epoch) for several $\lambda$ values. $\lambda$ is added to avoid over-fitting. Thus, higher the $\lambda$ value, lower will be the accuracy on training set as more emphasis is given to the complexity function rather than the original loss function *E(w)*.

Graph **??** shows how the accuracy on training data varies according to the iteration number (epoch) for several $\lambda$ values. The same generalization as above holds here as well. Note that the graph line for $\lambda = 10000$ overlaps with that of $\lambda = 1000$ and is not visible clearly.

Next, we plot the length of weight vectors against training iterations for different $\lambda$ values. In both the cases, as $\lambda$ increases, the corresponding weight vector length value decreases. Both these values are inversely related. This is due because as $\lambda$ increases, the overall regularized loss function value tends to increase. However, since our goal is to minimize loss, the weight vector balances this increase in $\lambda$ by decreasing itself. A similar argument is valid for decreasing $\lambda$ values as well.

Then we have the plots of final test error for different $\log(\lambda)$ values. In both the above graphs, as $\lambda$ (or $\log(\lambda)$) increases, the final test error increases. This is because for large $\lambda$ values, the model complexity is low. Beyond a certain level, the complexity may become so low that it no longer fits the data well leading to mis-classification on a large number of points. Increasing $\lambda$ value only decreases the complexity further, leading to even further decline in test accuracy. Similarly, if the $\lambda$ value is too low, the model may become highly complex, so much so that over-fitting happens. Test error in this case will again be large, as the model won't generalize well for points in test set.

16

Finally, we have the weight vector images. Since, regularization limits the weights from being too high, the images of weights are little bit softer in nature.

# 6 Softmax Regression Via Gradient Descent

## 6.1 Introduction

The task at hand is to perform Softmax Regression on the MNIST data set and come up with the best parameters that may perform well on the test data, without actually looking at test data. Then, we need to plot the loss function values over number of training iterations for training, hold-out and test data sets. Finally, we are expected to plot the percent correct values over training iterations the three data set parts.

## 6.2 Methodology

Softmax Regression was performed on the first 20,000 training data points and was used to do a 10-way classification of the hand-written digits using a hold-out set and regularization parameter $\lambda$. The size of hold-out set was again set to 10% of the size of training data. To figure out the best hyper-parameter values, a hold-out set was used. The parameters performing the best on this hold-out set were chosen to be the final parameters. The loss function(E) was plotted over the number of training iterations for the training, hold-out and test sets.

Parameters used: Penalty = $L_2$ norm; Regularization parameter $\lambda = 0.0001$; Learning rate $\eta = 0.0001$

Next, the "percent correct" (or accuracy percentage) was plotted over the number of training iterations for the training, hold-out and test sets. Same parameters as above were used.

## 6.3 Results

The percentage error values recorded on the hold-out set for different values of hyper-parameters are as follows:

Table 1: Error on Hold-out Set for Different Hyper-parameters

| $\eta$ | $\lambda$ | Norm | Error % |
|--------|-----------|--------|---------|
| 0.0001 | 0.01 | L-1/L-2 | 8.1 |
| 0.0001 | 0.1 | L-1/L-2 | 8.1 |
| 0.0001 | 0.0 | L-1/L-2 | 8.1 |
| 0.001 | 0.0 | L-1/L-2 | 8.15 |
| 0.01 | 0.0 | L-1/L-2 | 8.2 |
| 0.1 | 0.0 | L-1/L-2 | 88.2 |

All the graphs produced are plotted and reported herein.



Figure 25: Loss Function(E) v/s Epoch

18

Accuracy Percentage v/s Epoch



Figure 26: Accuracy Percentage v/s Epoch

## 6.4 Discussion

As evident from the given table, the optimum value of error on hold-out (validation) set was obtained at more than one pair of values of $\eta$ and $\lambda$. We decided to go with the highest $\eta$ and the largest non-zero $\lambda$ that gave us the optimum error value on hold-out set. This is because we wanted $\lambda$ as large as possible to help generalization and avoid over-fitting, and at the same time being a good representative of the data. Also, higher learning rate was preferred for faster convergence. Since the penalties at both $L_2$ and $L_1$ norm did not seem to affect the error values, we decided to go forward with $L_2$ norm, as it gives a better measure of loss and is convex everywhere.

Thus, the hyper-parameters chosen were": Penalty = $L_2$ norm Regularization parameter $\lambda = 0.1$ Learning rate $\eta = 0.0001$

Validation Error = 8.1

Test Error obtained = 12.35

Note that early stopping was used to make sure that we do not over-fit the data.

Graph **??** shows how the training, validation and test errors vary according to the iteration number (epoch) for the same values of hyper-parameters. As evident from the graph, the loss function decreases and stabilizes over time, which corresponds to convergence. As an indicative measure, the values of test, validation and training errors have been displayed for a particular value (93) of epoch.

Graph **??** shows how the training, validation and test percent correct values vary according to the iteration number (epoch) for the same values of hyper-parameters. As evident from the graph, the accuracy saturates over time and does not improve. The weights learned give the best performance on validation set, followed by training and test set. It is interesting to note that all the values are fairly close to each other and similar in shape, which means that the training set is a good representative of the hold-out and test sets.

19

# 7  Results and Learnings

The best accuracy achieved using Logistic Regression was more than 97% for both the subsets of data - having digits 2 and 3, and having digits 2 and 8. The best learning rate was $\eta = 0.0001$ and the best regularization parameter $\lambda$ was 0.0001 for 2's and 3's, and 0.1 for 2's and 8's. Using Softmax Regression, our accuracy was around 87.65% on test data. The $\lambda$ used in this case was 0.001 while the learning rate was kept to be 0.0001. Annealing of learning rate per iteration helped us avoid over-fitting of the data, even when number of iterations was huge. This annealing parameter was set to a convenient value (2000) to ensure a gradual yet sufficient decrease in learning rate. The early stopping margin for number of iterations was set to 15 for the both the cases to ensure that the algorithm was not stopping at some sub-optimal value, which was the case when this value was small.

Having taken courses like 250A and 250B, we knew the working and mechanism of Logistic and Softmax Regression, but we had never had a chance to perform their in-depth analysis by ourselves. Deriving the expressions and implementing these two algorithms along with regularization gave us a new insight and deep understanding of the working of these methods. Questions involving plotting of loss functions, weight vectors as images and weight vector lengths showed us how these values vary with different regularization parameters $\lambda$ and the iteration number. We had never looked at the MNIST digit classification problem from this perspective and now have a clearer idea as to how the various hyper-parameters are related to each other. The impact of regularization on the model and weights is now in front of us, while previously we only looked at it theoretically. The concepts of annealing and early stopping were entirely new to us, as these concepts were never visible when using *scikit* for performing these computations.

# 8  Individual Contributions

Being roommates, it was extremely convenient and simple for both the authors to coordinate and work in sync, while ensuring equal distribution of work and time spent on the assignment. Whenever one of the authors got stuck at some point, the other was there to help him out and unblock instantly. The process was initially started with both of us sitting down and solving on a white board the various derivations involved. The work thereafter was taken up as under:

Chetan Gandotra implemented the part which involved reading of MNIST data and implemented Softmax Regression. Then, debugging and graph plotting of Logistic Regression was taken up by him.

Rishabh Misra took up the implementation of Logistic Regression after extracting digit specific data. Thereafter, debugging and graph plotting of Softmax Regression was taken up by him.

The implementation and graphs for regularization question (Q5) were divided equally, with 5 (b) and 5 (c) being taken up by Chetan, and the others by Rishabh. For parameter tuning (values of $\lambda$, $\eta$, T, early stopping iteration number etc.), we made an Excel sheet with possible list of values and divided them equally amongst us. We then ran the code for these values for parameters on our respective systems. When it came to writing the report, we took alternate question parts, with Rishabh taking odd questions and Chetan taking up even questions.

278 **References**

279 [1] https://github.com/akosiorek/CSE/tree/master/MLCV

280 **Appendix**

281 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=True, 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_indicies : 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 nor inices are returned, then this is returned directly, and not inside
    labels : ndarray
        Array of size ``N`` describing the labels. Returned only if ``return_labels`` is
        'True', which is default.
    indices : ndarray
        The indices in the database that were returned.
```

```
338         Examples
339         ————
340         Assuming that you have downloaded the MNIST database and set the
341         environment variable ``$MNIST`` point to the folder, this will load all
342         images and labels from the training set:
343
344         >>> images, labels = ag.io.load_mnist('training') # doctest: +SKIP
345
346         Load 100 sevens from the testing set:
347
348         >>> sevens = ag.io.load_mnist('testing', digits=[7], selection=slice(0, 100),
349          return_labels=False) # doctest: +SKIP
350
351         """
352
353         # The files are assumed to have these names and should be found in 'path'
354         files = {
355             'training': ('train-images.idx3-ubyte', 'train-labels.idx1-ubyte'),
356             'testing': ('t10k-images.idx3-ubyte', 't10k-labels.idx1-ubyte'),
357         }
358
359         if path is None:
360             try:
361                 path = 'C:\\Users\\Chetan\\Documents\\Python Scripts\\Way1'
362                 #path = os.environ['MNIST']
363             except KeyError:
364                 raise ValueError("Unspecified path requires environment variable $MNIST
365                 to be set")
366
367         try:
368             images_fname = os.path.join(path, files[dataset][0])
369             labels_fname = os.path.join(path, files[dataset][1])
370         except KeyError:
371             raise ValueError("Data set must be 'testing' or 'training'")
372
373         # We can skip the labels file only if digits aren't specified and labels aren't
374         asked for
375         if return_labels or digits is not None:
376             flbl = open(labels_fname, 'rb')
377             magic_nr, size = struct.unpack(">II", flbl.read(8))
378             labels_raw = pyarray("b", flbl.read())
379             flbl.close()
380
381         fimg = open(images_fname, 'rb')
382         magic_nr, size, rows, cols = struct.unpack(">IIII", fimg.read(16))
383         images_raw = pyarray("B", fimg.read())
384         fimg.close()
385
386         if digits:
387             indices = [k for k in range(size) if labels_raw[k] in digits]
388         else:
389             indices = range(size)
390
391         if selection:
392             indices = indices[selection]
393         N = len(indices)
394
395         images = zeros((N, rows, cols), dtype=uint8)
396
397         if return_labels:
398             labels = zeros((N), dtype=int8)
399         for i, index in enumerate(indices):
400             images[i] = array(images_raw[ indices[i]*rows*cols : (indices[i]+1)*rows*cols
401             ]).reshape((rows, cols))
402             if return_labels:
```

```
403                     labels[i] = labels_raw[indices[i]]

404

405        if not asbytes:
406            images = images.astype(float)/255.0

407

408        ret = (images,)
409        if return_labels:
410            ret += (labels,)
411        if return_indices:
412            ret += (indices,)
413        if len(ret) == 1:
414            return ret[0] # Don't return a tuple of one
415        else:
416            return ret
```

417 ——————————————————————————————————————————————————————

418  Logistic_Regression_via_Gradient_Descent.py

419

```
420   """
421   CSE 253: Neural Networks and Pattern Recognition
422   Logistic Regression With and Without Gradient Descent
423   This file contains code for questions 4 and 5, including all graph plots
424   """
425   import numpy
426   import math
427   import matplotlib.pyplot as plt
428   import plotly.plotly as py1
429   import plotly.graph_objs as go

430

431   from LoadMNIST import load_mnist
```

432   #————————————————————————————————Utility functions————————————————————————————————

```
433   def get_data(N, N_test):
434       #load MNIST data using libraries available
435       training_data, training_labels = load_mnist('training')
436       test_data, test_labels = load_mnist('testing')

437

438       training_data = flatArray(N, 784, training_data) #training_data is N x 784 matrix
439       training_labels = training_labels[:N]
440       test_data = flatArray(N_test, 784, test_data)
441       test_labels = test_labels[:N_test]

442

443       # adding column of 1s for bias
444       training_data = addOnesColAtStart(training_data)
445       test_data = addOnesColAtStart(test_data)

446

447       # Last 10% of training data size will be considered as the validation set
448       N_validation = int (N / 10)
449       validation_data = training_data[N-N_validation:N]
450       validation_labels = training_labels[N-N_validation:N]
451       N=N-N_validation
452       #update training data to remove validation data
453       training_data = training_data[:N]
454       training_labels = training_labels[:N]

455

456       return training_data, training_labels, test_data, test_labels, validation_data,
457       validation_labels

458

459   def flatArray(rows, cols, twoDArr):
460       flattened_arr = numpy.zeros(shape=(rows, cols))
461       for row in range(0, rows):
462           i=0
463           for element in twoDArr[row]:
464               for el1 in element:
465                   flattened_arr[row][i] = el1
466                   i = i+1
467       return flattened_arr
```

```
468
469   def addOnesColAtStart(matrix):
470       Ones = numpy.ones(len(matrix))
471       newMatrix = numpy.c_[Ones, matrix]
472       return newMatrix
473
474   # custom sigmoid function; if -x is too large, return value 0
475   def sigmoid(x):
476       if(-x < 709):
477           return 1 / (1 + math.exp(-x))
478       else:
479           return 1 / (1 + math.exp(708))
480
481   def extract_digit_specific_data(digits, data, label):
482       pruned_data = numpy.zeros(shape = (1, len(data[0])))
483       pruned_labels = []
484       cnt = 0
485       for i in range(0, len(label)):
486           if label[i] in digits:
487               if (cnt == 0):
488                   for j in range(len(data[0])):
489                       pruned_data[0][j] = data[i][j]
490               else:
491                   pruned_data = addRowToMatrix(pruned_data, data[i])
492               pruned_labels.append(digits.get(label[i]))
493               cnt = cnt + 1
494       return pruned_data, pruned_labels
495
496   def addRowToMatrix(matrix, row):
497       newMatrix = numpy.zeros(shape = ((len(matrix) + 1), len(matrix[0])))
498       for i in range(len(matrix)):
499           newMatrix[i] = matrix[i]
500       newMatrix[i+1] = row
501       return newMatrix
502
503   def calculate_log_liikelihood(data, label, weights, t):
504       log_likelihood = 0.0
505       for j in range(0, t):
506           #print(numpy.log(sigmoid(numpy.dot(weights, data[j]))))
507           log_likelihood += (label[j]*numpy.log(sigmoid(numpy.dot(weights, data[j])))) +
508           ((1-label[j])*numpy.log(sigmoid(-1*numpy.dot(weights, data[j]))))
509
510       return -1*log_likelihood/t
511
512   # 5.b, 5.c - Plot of Percent correct on training data v/s iterations,
513   #length of weight vector v/s lambda
514   def plotlyGraphsRegularization(error_plot_array, lamda_vals, graph_name, min_index = -1):
515       py1.sign_in('chetang', 'vil7vTAuCSWt2lEZvaH9')
516
517       trace = []
518
519       for i in range(len(lamda_vals)):
520           y1 = error_plot_array[i]
521           #y1 = y1[:min_index[i]]
522           x1 = [j+1 for j in range(len(y1))]
523
524           trace1 = go.Scatter(
525               x=x1,
526               y=y1,
527               name = 'lambda = ' + (str)(lamda_vals[i]), # Style name/legend entry with
528               html tags
529               connectgaps=True
530           )
531
532           trace.append(trace1)
```

24

```
533        data = trace

535        fig = dict(data=data)
536        py1.iplot(fig, filename=graph_name)

538    def plotlyGraphs(error_plot_array, labels, name):
539        py1.sign_in('cgandotr', '3c9fho4498')
540        trace = []

542        for i in range(len(labels)):
543            y1 = error_plot_array[i]
544            y1 = [k for k in y1]
545            x1 = [(j+1) for j in range(len(y1))]

547            trace1 = go.Scatter(
548                x=x1,
549                y=y1,
550                name = str(labels[i]), # Style name/legend entry with html tags
551                connectgaps=True
552            )
553            trace.append(trace1)
554        data = trace
555        fig = dict(data=data)
556        py1.iplot(fig, filename=name)

558    def early_stopping(early_stopping_horizon, accuracy, i):
559        if i > early_stopping_horizon + 1:
560            counter = 0;
561            for p in range(0,early_stopping_horizon):
562                if accuracy[i-p] <= accuracy[i-p-1]:
563                    counter+=1
564                else:
565                    break

567            if(counter == early_stopping_horizon):
568                return True
569            else:
570                return False

572    def calculate_error(weights, data, label):
573        error = 0.0;
574        for j in range(0,len(data)):
575            prediction = sigmoid(numpy.dot(weights, data[j]))
576            if(prediction >0.5 and label[j]!=1):
577                error += 1;
578            elif (prediction <=0.5 and label[j]!=0):
579                error += 1;
580        return error;

582    def dropFirstColumn(weights):
583        return numpy.array(weights)[0][1:]

585    def fit(training_data, training_label, test_data, test_label, validation_data,
586            validation_label, digits, learning_rate=0.0001, iteration=200, batch_size=0,
587            T=5000, lamda_vals=[0], norm=2):

589        t = len(training_data)
590        accuracy_plot_array = []
591        log_likelihood_array = []
592        weight_vector_length_array = []
593        accuracy_plot_training_array = []
594        weights_for_all_lamda = []
595        test_error_array = []
596        org_learning_rate = learning_rate

```

```
598        for lamda in lamda_vals:
599            weights = numpy.matrix(numpy.zeros(len(training_data[0])))
600            weights_array = []
601            weight_vector_length = []
602
603            accuracy_plot_training = []
604            accuracy_plot_validation = []
605            accuracy_plot_testing = []
606
607            log_likelihood_training = []
608            log_likelihood_validation = []
609            log_likelihood_testing = []
610            min_error_index = 0
611            learning_rate = org_learning_rate
612
613            for i in range(0, iteration):
614                # initialise Gradient
615                gradient = numpy.matrix(numpy.zeros(len(training_data[0])))
616                cnt = 0
617
618                norm_term = []
619                if (norm == 2):
620                    norm_term = l2_norm(lamda, weights)
621                else:
622                    norm_term = l1_norm(lamda, weights)
623
624                # calculate gradient over all the samples
625                for j in range(0,t):
626                    # update gradient
627                    gradient += ((training_label[j]) - sigmoid(numpy.dot(weights,
628                    training_data[j])))*numpy.array(training_data[j])
629                    cnt += 1
630                    if (batch_size != 0 and cnt == (int)(t/batch_size)):
631                        cnt = 0
632                        weight_vector_length.append(numpy.linalg.norm(weights))
633                        weights = weights + learning_rate * (gradient - norm_term)
634                        # re-initialise Gradient
635                        gradient = numpy.matrix(numpy.zeros(len(training_data[0])))
636                        # calculating log likelhood of training, validation and test
637                        dataset
638                        log_likelihood_training.append(calculate_log_liikelihood(training_data,
639                         training_label, weights, t))
640                        log_likelihood_validation.append(calculate_log_liikelihood(validation_data
641                        validation_label, weights, len(validation_data)))
642                        log_likelihood_testing.append(calculate_log_liikelihood(test_data,
643                        test_label, weights, len(test_data)))
644
645                if (batch_size == 0):
646                    weight_vector_length.append(numpy.linalg.norm(weights))
647                    # update weights vector according to the update rule of Gradient descent metho
648                    weights = weights + learning_rate * (gradient - norm_term)
649
650                    # calculating log likelhood of training, validation and test dataset
651                    log_likelihood_training.append(calculate_log_liikelihood(training_data,
652                    training_label, weights, t))
653                    log_likelihood_validation.append(calculate_log_liikelihood(validation_data,
654                    validation_label, weights, len(validation_data)))
655                    log_likelihood_testing.append(calculate_log_liikelihood(test_data, test_label,
656                    weights, len(test_data)))
657
658                # anneling of learning rate
659                learning_rate = learning_rate/(1+i/T)
660
661                # calculating error percentage on train, test and validation data
662                accuracy_plot_training.append((len(training_data) - calculate_error(weights,
```

```
663              training_data, training_label))*100/len(training_data))
664              accuracy_plot_validation.append((len(validation_data) − calculate_error(weights,
665              validation_data, validation_label))*100/len(validation_data))
666              accuracy_plot_testing.append((len(test_data) − calculate_error(weights, test_data,
667              test_label))*100/len(test_data))
668
669              weights_array.append(weights)
670
671              # check for early stopping
672              early_stopping_horizon = 15
673              min_error_index = i
674              if(early_stopping(early_stopping_horizon, accuracy_plot_validation, i) and i >
675              early_stopping_horizon):
676                  min_error_index = i−early_stopping_horizon;
677                  weights = weights_array[min_error_index]
678                  break
679
680          weight_vector_length_array.append(weight_vector_length)
681          weights_for_all_lamda.append(weights)
682
683          log_likelihood_array.append(log_likelihood_training);
684          log_likelihood_array.append(log_likelihood_validation);
685          log_likelihood_array.append(log_likelihood_testing);
686          accuracy_plot_array.append(accuracy_plot_training)
687          accuracy_plot_array.append(accuracy_plot_validation)
688          accuracy_plot_array.append(accuracy_plot_testing)
689          accuracy_plot_training_array.append(accuracy_plot_training)
690
691          test_error = (calculate_error(weights, test_data, test_label)*100)/len(test_data)
692          validation_error = (calculate_error(weights, validation_data, validation_label)*100)
693          /len(validation_data)
694          print('Error on validation dataset : ' + str(validation_error) + '%');
695          print('Error on test dataset : ' + str(test_error) + '%');
696          test_error_array.append(test_error)
697
698      #For single lambda value
699      if (len(lamda_vals) == 1):
700          plotlyGraphs(log_likelihood_array, ['Training Set','Validation Set','Test Set'],
701          'Log Likelihood Plot')
702          plotlyGraphs(accuracy_plot_array, ['Training Set','Validation Set','Test Set'],
703          'Percent Correct Plot')
704      #else:
705          #For multiple lambda values
706          #plotlyGraphsRegularization(accuracy_plot_training_array, lamda_vals, "Accuracy vs
707          Epoch")
708          #plotlyGraphsRegularization(weight_vector_length_array, lamda_vals, "Weight Vector
709          Length vs Epoch")
710          #plotlyErrorVsLamda(test_error_array, lamda_vals)
711
712      # printing error on training and testing dataset
713      return weights_for_all_lamda
714
715  def plot_weights(weights):
716      lamda_vals = [1000, 100, 10, 1, 0.1, 0.001, 0.0001]
717      i = 0
718      for w in weights:
719          # Plot weights as image after removing bias terms. The rest of columns are pixels
720          print ('lambda = ' + str(lamda_vals[i]))
721          i += 1
722          pixels1 = dropFirstColumn(w)
723          pixels = numpy.reshape(pixels1, (28, 28))
724          plt.imshow(pixels, cmap='gray')
725          plt.show()
726
727  def l2_norm(lamda, weights):
```

```
728          return 2*lamda*weights;
729
730  def l1_norm(lamda, weights):
731      w = numpy.ones(len(weights))
732      for i in range(len(weights)):
733          if (weights[i] < 0):
734              w[i] = -1
735      return lamda*w;
736
737  # 5.d - Final test error v/s Lambda values graph
738  # Generates bar graphs
739  def plotlyErrorVsLamda(test_error_array, lamda_vals):
740      lamda_vals1 = [math.log(lamda) for lamda in lamda_vals]
741      pyl.sign_in('chetang', 'vil7vTAuCSWt2lEZvaH9')
742      trace = []
743      colors = ['rgb(104,224,204)', 'rgb(204,204,204)', 'rgb(49,130,189)', 'rgb(41,180,129)']
744      for i in range(0, len(test_error_array)):
745          yl = test_error_array[i]
746          x1 = lamda_vals1[i]
747
748          trace1 = go.Bar(
749              x=x1,
750              y=y1,
751              name='Lambda = log(' + (str)(lamda_vals[i]) + ')',
752              marker=dict(
753              color=colors[i]
754              )
755          )
756          trace.append(trace1)
757
758      layout = go.Layout(
759          xaxis=dict(tickangle=-45),
760          barmode='group',
761      )
762      fig = go.Figure(data=trace, layout=layout)
763      pyl.iplot(fig, filename='Final Error vs Lambda')
764
765  #————————————————————————————Main function————————————————————————————
766
767  if __name__ == "__main__":
768      numpy.random.seed(0)
769
770      N = 20000
771      N_test = 2000
772      iteration = 200
773      batch_size = 0
774      T = 2000
775
776      #lamda_vals = [0]
777      lamda_vals = [1000, 100, 10, 1, 0.1, 0.001, 0.0001] # regularization
778      weightage parameter.
779       Put [0] for no regularization
780      norm = 2 # 2 for l-2 norm, 1 for l-1 norm
781
782      full_training_data, full_training_label, full_test_data, full_test_label,
783      full_validation_data, full_validation_label = get_data(N, N_test)
784
785      # Parameters for training data on 2's and 3's
786      digits = {2:1, 3:0}
787      training_data, training_label = extract_digit_specific_data(digits,
788      full_training_data,
789      full_training_label)
790      validation_data, validation_label = extract_digit_specific_data(digits,
791      full_validation_data,
792       full_validation_label)
```

```
793        test_data , test_label = extract_digit_specific_data ( digits , full_test_data ,
794        full_test_label )
795        learning_rate = 0.0001
796
797        weights23 = fit ( training_data , training_label , test_data , test_label ,
798                        validation_data , validation_label , learning_rate , iteration ,
799                        batch_size , digits , T, lamda_vals , norm )
800        plot_weights ( weights23 )
801
802        # Parameters for training data on 2's and 8's
803        digits = {2:1 , 8:0}
804        training_data , training_label = extract_digit_specific_data ( digits ,
805        full_training_data ,
806        full_training_label )
807        validation_data , validation_label = extract_digit_specific_data ( digits ,
808        full_validation_data ,
809        full_validation_label )
810        test_data , test_label = extract_digit_specific_data ( digits , full_test_data ,
811        full_test_label )
812        learning_rate = 0.1
813        lamda_vals = [0] #Regularization not asked for with 2/8 case
814
815        weights28 = fit ( training_data , training_label , test_data , test_label ,
816                        validation_data , validation_label , digits , learning_rate ,
817                        iteration ,
818                        batch_size , T, lamda_vals , norm )
819        plot_weights ( weights28 )
820
821        weights = weights28 [0] − weights23 [0]
822
823        plot_weights ( weights )
824
825  ──────────────────────────────────────────────────────────────
826
827  Softmax_Regression . py
828
829  """
830  Softmax Regression on MNIST Data set to perform 10−way classification
831  """
832  import numpy
833  import math
834  import plotly . plotly as py1
835  import plotly . graph_objs as go
836
837  from LoadMNIST import load_mnist
838
839  #─────────────────────────────── Utility functions ────────────────────────────
840  def get_data (N, N_test ) :
841      #load MNIST data using libraries available
842      training_data , training_labels = load_mnist ( ' training ' )
843      test_data , test_labels = load_mnist ( ' testing ' )
844
845      training_data = flatArray (N, 784 , training_data ) #training_data is N x
846      784 matrix
847      training_labels = training_labels [:N]
848      test_data = flatArray ( N_test , 784 , test_data )
849      test_labels = test_labels [: N_test ]
850
851      # adding column of 1s for bias
852      training_data = addOnesColAtStart ( training_data )
853      test_data = addOnesColAtStart ( test_data )
854
855      # Last 10% of training data size will be considered as the validation
856       set
857      N_validation = int (N / 10)
```

```
858         validation_data = training_data[N-N_validation:N]
859         validation_labels = training_labels[N-N_validation:N]
860     N=N-N_validation
861     #update training data to remove validation data
862         training_data = training_data[:N]
863         training_labels = training_labels[:N]
864
865     return training_data, training_labels, test_data, test_labels,
866     validation_data,
867     validation_labels
868
869 def flatArray(rows, cols, twoDArr):
870     flattened_arr = numpy.zeros(shape=(rows, cols))
871     for row in range(0, rows):
872         i=0
873         for element in twoDArr[row]:
874             for el1 in element:
875                 flattened_arr[row][i] = el1
876                 i = i+1
877     return flattened_arr
878
879 def addOnesColAtStart(matrix):
880     Ones = numpy.ones(len(matrix))
881     newMatrix = numpy.c_[Ones, matrix]
882     return newMatrix
883
884 # custom sigmoid function; if -x is too large, return value 0
885 def exp(x):
886     if(x < 709):
887         return math.exp(x)
888     else:
889         return math.exp(709)
890
891 def addRowToMatrix(matrix, row):
892     newMatrix = numpy.zeros(shape = ((len(matrix) + 1), len(matrix[0])))
893     for i in range(len(matrix)):
894         newMatrix[i] = matrix[i]
895     newMatrix[i+1] = row
896     return newMatrix
897
898 def l2_norm(lamda, weights):
899     return 2*lamda*weights;
900
901 def l1_norm(lamda, weights):
902     w = numpy.ones((len(numpy.array(weights)),
903     len(numpy.array(weights)[0])))
904     for i in range(len(w)):
905         for j in range(len(w[0])):
906             if (weights[i,j] < 0):
907                 w[i][j] = -1
908     return lamda*numpy.matrix(w);
909
910 # custom sigmoid function; if -x is too large, return value 0
911 def sigmoid(x):
912     if(-x < 709):
913         return 1 / (1 + math.exp(-x))
914     else:
915         return 1 / (1 + math.exp(708))
916
917 def calculate_log_liikelihood(data, label, weights, t):
918     log_likelihood = 0.0
919     for j in range(0,t):
920         log_likelihood += (label[j]*numpy.log(sigmoid(numpy.dot(weights,
921         data[j])))) +
922         ((1-label[j])*numpy.log(sigmoid(-1*numpy.dot(weights, data[j]))))
```

```
923
924          return −1∗log_likelihood/t
925
926  def plotlyGraphs(error_plot_array, labels, name):
927          pyl.sign_in('cgandotr', '3c9fho4498')
928          trace = []
929          for i in range(len(labels)):
930              yl = error_plot_array[i]
931              yl = [k for k in yl]
932              #yl = yl[:min_index[i]]
933              xl = [(j+1) for j in range(len(yl))]
934
935              trace1 = go.Scatter(
936                  x=xl,
937                  y=yl,
938                  name = str(labels[i]), # Style name/legend entry with html tags
939                  connectgaps=True
940              )
941
942              trace.append(trace1)
943          data = trace
944          fig = dict(data=data)
945          pyl.iplot(fig, filename=name)
946
947  def early_stopping(early_stopping_horizon, accuracy, i):
948          if i > early_stopping_horizon:
949              counter = 0;
950              for p in range(0,early_stopping_horizon):
951                  if accuracy[i−p] <= accuracy[i−p−1]:
952                      counter+=1
953                  else:
954                      break
955
956              if(counter == early_stopping_horizon):
957                  return True
958              else:
959                  return False
960
961  def calculate_error(weights, data, label, k = 10):
962          error = 0.0;
963          for j in range(0,len(data)):
964              softmax_denom = 0.0
965              softmax_num = numpy.zeros(k);
966              for x in range(0,k):
967                  softmax_num[x] = exp(numpy.dot(numpy.transpose(weights[:,x]),
968                   data[j]))
969                  softmax_denom += softmax_num[x]
970
971              prediction = numpy.argmax(softmax_num/softmax_denom)
972              if(prediction != label[j]):
973                  error += 1;
974          return error;
975
976  def softmax_loss(weights, labels, data, c):
977          loss = 0.0
978          for i in range(len(data)):
979              softmax_denom = 0.0
980              softmax_num = numpy.zeros(c);
981              for x in range(0,c):
982                  softmax_num[x] = exp(numpy.dot(numpy.transpose(weights[:,x]),
983                   data[i]))
984                  softmax_denom += softmax_num[x]
985              softmax = softmax_num[labels[i]]/softmax_denom
986              if (softmax > 0):
987                  log_softmax = math.log(softmax)
```

31

```
988              else :
989                  log_softmax = 0
990           loss += ( log_softmax )
991       return −1∗loss /( len ( data ))

993  def getSoftmax(k, weights, training_data, j):
994       softmax_denom = 0.0
995       softmax_num = numpy.zeros(k);
996       for x in range (0,k):
997           softmax_num [x] = exp(numpy.dot(numpy.transpose(weights[:,x]) ,
998           training_data [j]))
999           softmax_denom += softmax_num [x]
1000      return softmax_num/softmax_denom

1002 def fit(training_data, training_label, test_data, test_label, validation_data,
1003          validation_label, iteration = 1000, T=2000, lamda=0.001,
1004          learning_rate = 0.0001, norm = 2):
1005      k = len (numpy.unique(training_label))
1006      t = len(training_data)
1007      weights = numpy.matrix(numpy.zeros((len(training_data[0]), k)))
1008      weights_array = []

1010      early_stopping_horizon = 15
1011      error_plot = numpy.zeros(iteration)
1012      min_error_index = 0
1013      loss_array = []
1014      loss_training = []
1015      loss_validation = []
1016      loss_testing = []

1018      accuracy_plot_array = []
1019      accuracy_plot_training = []
1020      accuracy_plot_validation = []
1021      accuracy_plot_testing = []
1022      test_error = 0.0

1024      for i in range(0, iteration):
1025          # initialise Gradient
1026          gradient = numpy.matrix(numpy.zeros((len(training_data[0]), k)))

1028          # calculate gradient over all the samples
1029          for j in range(0,t):
1030              modified_label = numpy.zeros(k);
1031              modified_label[training_label[j]] = 1
1032              softmax = getSoftmax(k, weights, training_data, j)
1033              gradient += numpy.transpose(numpy.transpose(numpy.matrix(
1034              modified_label − softmax)
1035              ) ∗ numpy.matrix(training_data[j]))

1037          norm_term = 0.0
1038          if (norm == 2):
1039              norm_term = l2_norm(lamda, weights)
1040          else :
1041              norm_term = l1_norm(lamda, weights)
1042          # update weights vector according to the update rule of Gradient
1043          descent method
1044          weights = weights + learning_rate ∗ (gradient − norm_term)

1046          loss_training.append(softmax_loss(weights, training_label,
1047          training_data, k))
1048          loss_validation.append(softmax_loss(weights, validation_label,
1049          validation_data, k))
1050          loss_testing.append(softmax_loss(weights, test_label, test_data, k))

1052          # calculating error percentage on train, test and validation data
```

```python
1053            training_error = calculate_error(weights, training_data,
1054            training_label)
1055            validation_error = calculate_error(weights, validation_data,
1056             validation_label)
1057            test_error = calculate_error(weights, test_data, test_label)
1058
1059            accuracy_plot_training.append((len(training_data) - training_error)
1060            *100/len(training_data))
1061            accuracy_plot_validation.append((len(validation_data) -
1062            validation_error)*100/len(validation_data))
1063            accuracy_plot_testing.append((len(test_data) - test_error)*100/
1064            len(test_data))
1065
1066            learning_rate = learning_rate/(1+i/T)
1067
1068            error_plot[i] = validation_error*100/len(validation_data);
1069            weights_array.append(weights)
1070
1071            # check for early stopping
1072            if (early_stopping(early_stopping_horizon, accuracy_plot_validation,
1073             i)):
1074                min_error_index = i-early_stopping_horizon;
1075                weights = weights_array[min_error_index]
1076                break
1077            min_error_index = i
1078
1079        loss_array.append(loss_training)
1080        loss_array.append(loss_validation)
1081        loss_array.append(loss_testing)
1082
1083        accuracy_plot_array.append(accuracy_plot_training)
1084        accuracy_plot_array.append(accuracy_plot_validation)
1085        accuracy_plot_array.append(accuracy_plot_testing)
1086
1087        # printing error on training and testing dataset
1088        print('Error on validation dataset : ' + str(error_plot[min_error_index])
1089         + '%');
1090        print('Error on test dataset : ' + str(test_error*100/len(test_data)) +
1091        '%');
1092
1093        return weights, loss_array, accuracy_plot_array
1094
1095    #───────────────────────────────Main function─────────────────────────────
1096
1097    if __name__ == "__main__":
1098        numpy.random.seed(0)
1099        learning_rate = 0.0001
1100        N = 20000
1101        N_test = 2000
1102        lamda = 0.001          # regularization weightage parameter
1103        T = 2000
1104        iteration = 1000
1105        training_data, training_label, test_data, test_label, validation_data, validation_label
1106        = get_data(N, N_test)
1107
1108        weights, loss_array, accuracy_plot_array = fit(training_data, training_label,
1109                                                 test_data, test_label, validation_data,
1110                                                 validation_label)
1111
1112        plotlyGraphs(loss_array, ['Training Set','Validation Set','Test Set'], "Loss Function and
1113         Iterations")
1114        plotlyGraphs(accuracy_plot_array, ['Training Set','Validation Set','Test Set'], "Accuracy
1115         and Iterations")
```