

Implementation of a Feedforward Neural Network from Scratch for Image Classification

Course: Machine Learning (INF267)

Chryssa Nampouri (t8150096)
Department of Management Science and Technology
Athens University of Economics and Business
Athens, Greece

Supervisor: Prof. Prodromos Malakasiotis

1 Project Description

The purpose of this project is the implementation of Stochastic Gradient Ascent, i.e. the process of maximizing, instead of minimizing, a loss function, in order to train a Feedforward Neural Network with one hidden layer. Then, the implemented algorithm will be used for the classification of images of the Mnist and Cifar-10 data sets.

```
1 import re
2 import pickle
3 import numpy as np
4 import pandas as pd
5 import matplotlib.pyplot as plt
6 import warnings
7 warnings.simplefilter('error', RuntimeWarning)
```

2 Data Sets

2.1 Mnist data set

In the data folder there is the data set of mnist. Mnists consists of 28x28 grayscale images. In total there are 10 training files train0.txt, train1.txt, ..., train9.txt where each rows of traink.txt corresponds to an example that belongs to the class k .

The testing data follows the same format.

In total we have $6 * 10^5$ training examples and 10^3 testing examples.

```
1 def load_mnist_data( dataset ):
2
3     """ Load the dataset. Reads the training and testing files and creates
4         matrices.
5
6     :param dataset: The data set folder
7     :return:
8         train_data: the matrix with the training data
9         test_data: the matrix with the data that will be used for testing
10        y_train: the matrix consisting of one
11                hot vectors on each row (ground truth for training)
12        y_test: the matrix consisting of one
13                hot vectors on each row (ground truth for testing)
14
15     """
16
17     # Load the train files
18     df = None
19
20     y_train = []
21
22     for i in range( 10 ):
23         tmp = pd.read_csv( dataset + 'train%d.txt' % i, header=None, sep=" " )
24         # Build labels - one hot vector
25         hot_vector = [ 1 if j == i else 0 for j in range(0,10) ]
26
27         for j in range( tmp.shape[0] ):
28             y_train.append( hot_vector )
29         # Concatenate dataframes by rows
30         if i == 0:
31             df = tmp
32         else:
33             df = pd.concat( [df, tmp] )
34
35     train_data = df.as_matrix()
36     y_train = np.array( y_train )
37
38     # Load test files
39     df = None
40
41     y_test = []
42
43     for i in range( 10 ):
44         tmp = pd.read_csv( dataset + 'test%d.txt' % i, header=None, sep=" " )
45         # Build labels - one hot vector
46         hot_vector = [ 1 if j == i else 0 for j in range(0,10) ]
```

```

46
47     for j in range( tmp.shape[0] ):
48         y_test.append( hot_vector )
49
50     # Concatenate dataframes by rows
51     if i == 0:
52         df = tmp
53     else:
54         df = pd.concat( [df, tmp] )
55
56     test_data = df.as_matrix()
57     y_test = np.array( y_test )
58
59     return train_data , test_data , y_train , y_test

```

2.2 CIFAR-10 data set

In the data folder there is, also the data set of **Cifar-10**. The archive contains the files `data_batch_1`, `data_batch_2`, ..., `data_batch_5`, as well as `test_batch`. Each of these files is a Python "pickled" object produced with `cPickle` and contains a dictionary with the following elements:

- **data**: a 10000x3072 numpy array. Each row of the array stores a 32x32 colour image. The first 1024 entries contain the red channel values, the next 1024 the green, and the final 1024 the blue. The image is stored in row-major order, so that the first 32 entries of the array are the red channel values of the first row of the image.
- **labels**: a list of 10000 numbers in the range 0-9, indicating the category each example belongs to.

In total we have $5 * 10^4$ training examples and 10^4 testing examples.

Unpickle data set:

```

1 def unpickle( file ):
2
3     """ Routine which will open pickled objects and return a dictionary.
4
5     :param file: Each batch of the data set
6     :return: A dictionary
7
8     """
9
10    with open( file , 'rb' ) as fo:
11        dict = pickle.load( fo , encoding='bytes' )
12    return dict

```

Load Cifar-10 data set

```
1 def load_cifar_data():
2
3     """ Load cifar data set and make necessary transformations.
4
5     :return:
6         X_train: the matrix with the training data
7         X_test: the matrix with the data that will be used for testing
8         y_train: the matrix consisting of one
9             hot vectors on each row (ground truth for training)
10        y_test: the matrix consisting of one
11            hot vectors on each row (ground truth for testing)
12
13     """
14
15     X_train = []
16     y_train = []
17     X_test = []
18     y_test = []
19
20     train_tmp = []
21     test_tmp = []
22
23     for i in range(1,6):
24         batch_dict = unpickle("../datasets/cifar-10-batches-py/data_batch_%d" %
25 i)
26         X_train.extend(batch_dict[b"data"])
27         train_tmp.extend(batch_dict[b"labels"])
28
29     for x in range(len(train_tmp)):
30         one = [ 1 if train_tmp[x] == j else 0 for j in range(0,10) ]
31         y_train.append(one)
32
33     batch_dict = unpickle("../datasets/cifar-10-batches-py/test_batch")
34     X_test.extend(batch_dict[b"data"])
35     test_tmp.extend(batch_dict[b"labels"])
36
37     for x in range(len(test_tmp)):
38         one = [ 1 if test_tmp[x] == j else 0 for j in range(0,10) ]
39         y_test.append(one)
40
41     X_train = np.asarray(X_train)
42     X_test = np.asarray(X_test)
43     y_train = np.asarray(y_train)
44     y_test = np.asarray(y_test)
45
46
47     return X_train, X_test, y_train, y_test
```

3 Load the desired data set

Select fom standard input if you want to train the Mnist (option 1) or the Cifar-10 (option 2) data set.

```
1 while True:
2     try:
3         # Read integer from stdin
4         data_set = int(input("Select 1 of the following numbers based on the
5         desired data set:\n \
6                               \n 1: Mnist data set \n 2: Cifar10 data set\n"))
7         print("")
8         if(data_set in range(1,3)):
9             break;
10        else:
11            raise ValueError('Invalid input. Please select again an integer
12            between 1-2!')
13        except ValueError:
14            print("")
15            print("Invalid input. Please select again an integer between 1-2!")
16            print("")
17
18 # Case1: input = 1
19 if(data_set == 1):
20     print("You selected Mnist data set!")
21     X_train, X_test, y_train, y_test = load_mnist_data("../datasets/mnistdata/"
22     )
23
24 # Case2: input = 2
25 elif (data_set == 2):
26     print("You selected Cifar-10 data set!")
27     X_train, X_test, y_train, y_test = load_cifar_data()
```

4 Plot Mnist data set

```
1 def plot_mnist():
2
3     """ Plot 100 random images from the mnist training set. """
4
5     n = 100
6     sqrt_n = int( n**0.5 )
7     samples = np.random.randint(X_train.shape[0], size=n)
8
9     plt.figure( figsize=(11,11) )
10
11     cnt = 0
12     for i in samples:
13         cnt += 1
14         plt.subplot( sqrt_n, sqrt_n, cnt )
15         plt.subplot( sqrt_n, sqrt_n, cnt ).axis('off')
16         plt.imshow( X_train[i].reshape(28,28), cmap='gray' )
17
18     plt.show()
```

5 Plot Cifar-10 data set

```
1 def plot_cifar():
2
3     """ Plot 100 random images from the cifar training set. """
4
5     n = 100
6     sqrt_n = int( n**0.5 )
7     samples = np.random.randint(X_train.shape[0], size=n)
8
9     plt.figure( figsize=(11,11) )
10
11     cnt = 0
12     for i in samples:
13         arr = X_train[i]
14         R = arr[0:1024].reshape(32,32)/255.0
15         G = arr[1024:2048].reshape(32,32)/255.0
16         B = arr[2048:].reshape(32,32)/255.0
17
18         img = np.dstack((R,G,B))
19         cnt += 1
20         plt.subplot( sqrt_n, sqrt_n, cnt )
21         plt.subplot( sqrt_n, sqrt_n, cnt ).axis('off')
22         plt.imshow(img, interpolation='bicubic')
23
24     plt.show()
```

6 View of the selected data set

```
1 if(data_set == 1):  
2     plot_mnist()  
3 elif(data_set == 2):  
4     plot_cifar()
```

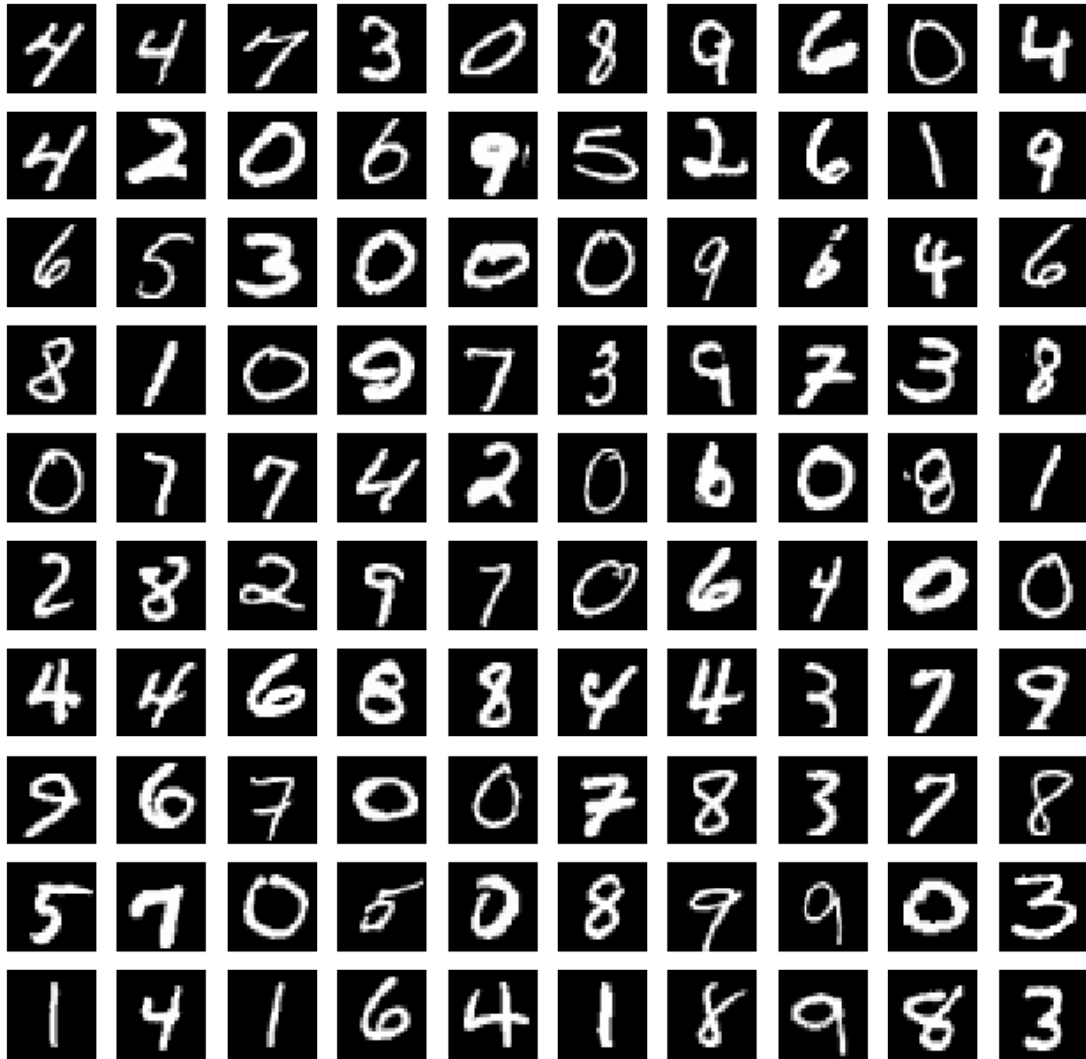


Figure 1: MNIST Data Set

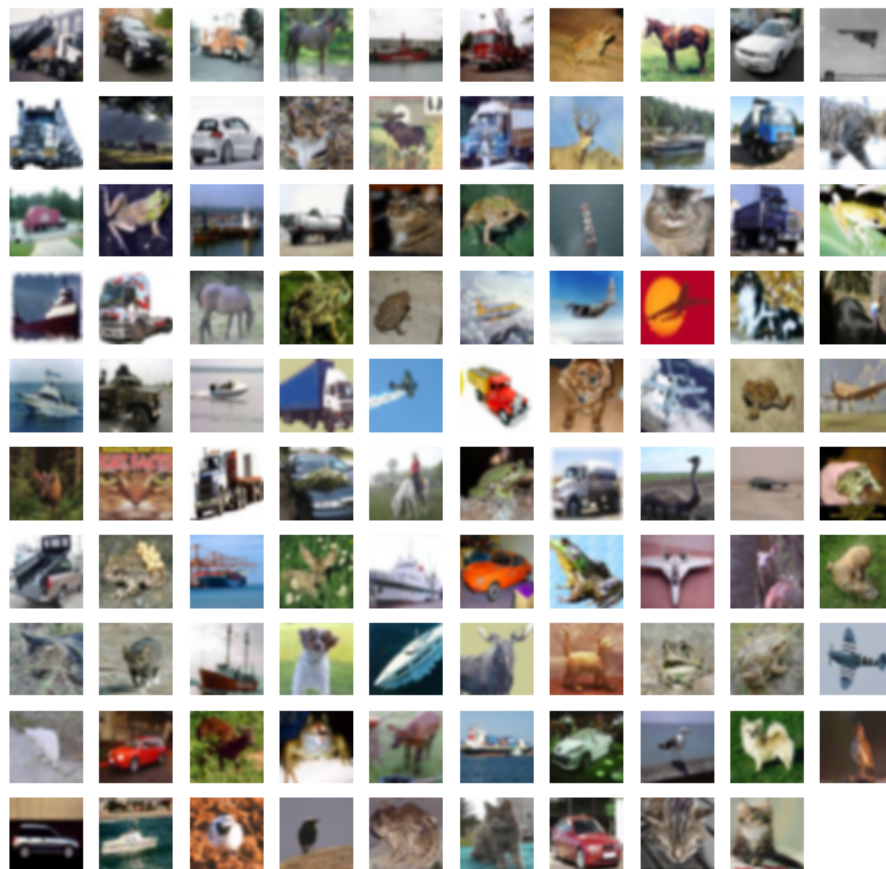


Figure 2: CIFAR-10 Data Set

7 Normalize the data set

Pixel values are integers that range from 0 (black) to 255 (white). So, we divide each feature by the maximum value, in order to normalize our data set in the range $[0, 1]$.

```
1 # Normalize the data set
2 X_train = X_train.astype(float)/255
3 X_test = X_test.astype(float)/255
```

8 Add bias parameter in the data set

Insert a column of 1's as the first entry in the feature vector — this is a little trick that allows us to treat the bias as a trainable parameter within the weight matrix rather than an entirely separate variable.

```
1 # Add bias in train and test set
2 X_train = np.hstack( (np.ones((X_train.shape[0],1) ), X_train) )
3 X_test = np.hstack( (np.ones((X_test.shape[0],1) ), X_test) )
```

9 Activation Functions

Using non-linear Activations we are able to generate non-linear mappings from inputs to output and learn something more complex and complicated from data. The below function implements the ***Logarithm Activation Function***, the ***Tangent Activation Function*** and the ***Cosine Activation Function*** and convert linear output of the current hidden layer of the neural network into non-linear. Then, the non-linear output will be used as input to the next layer.

Activation Function	Equation	Derivatives
logarithm	$Z(a) = \log(1 + e^a)$	$\frac{\partial Z}{\partial a} = \frac{e^a}{1+e^a}$
tangent	$Z(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$	$\frac{\partial Z}{\partial a} = 1 - Z^2(a)$
cosine	$Z(a) = \cos a$	$\frac{\partial Z}{\partial a} = -\sin a$

```

1 def activations(activation, a, ax=1):
2
3     """ Calculates the non-linear output of the current layer and the
4     derivative of the activation function used.
5
6     :param activation: The chosen non-linear activation function
7     :param a: The N x (M+1) matrix with the linear output of the current layer
8     :return:
9         z: N x (M+1) matrix with the non-linear output of the current layer and
10        grad_z: N x (M+1) matrix with the derivatives of the activation
11        function
12
13     """
14
15     # Case1: logarithm activation function
16     if(activation == "log"):
17         z = np.logaddexp(0.0, a)
18         grad_z = np.exp(a)/(1 + np.exp(a))
19     # Case2: tangent activation function
20     elif(activation == "tan"):
21         z = (np.exp(a) - np.exp(-a))/(np.exp(a) + np.exp(-a))
22         grad_z = np.ones(z.shape) - (z**2)
23     # Case3: cosine activation function
24     elif(activation == "cos"):
25         z = np.cos(a)
26         grad_z = -np.sin(a)
27
28     return z, grad_z

```

10 Softmax Fuction

The softmax function is defined as:

$$S_{nk} = \frac{e^{y_{nk}-m}}{\sum_{k=1}^K e^{y_{nk}-m}} \quad (1)$$

```
1 def softmax(y, ax=1):
2
3     """ Implements Softmax function and turns output numbers from logits layer
4     into probabilities that sum to one.
5
6     :param y: The N x K matrix with the linear output of the last hidden layer
7     :param ax=1: use by default, when the array id 2D
8     :return:
9         s: The N x K matrix with the probabilities of each train example n
10        belong to each category
11
12        """
13
14        # Find maximum elemeny per row
15        m = np.max(y, axis=ax, keepdims=True)
16        # Implement softmax function
17        # Subtract the maximum element so as to avoid overflow
18        p = np.exp(y - m)
19        s = (p / np.sum(p, axis=ax, keepdims=True))
20
21    return s
```

11 The Model

A Neural Network with one hidden layer, which classifies each example in one out of ten categories.

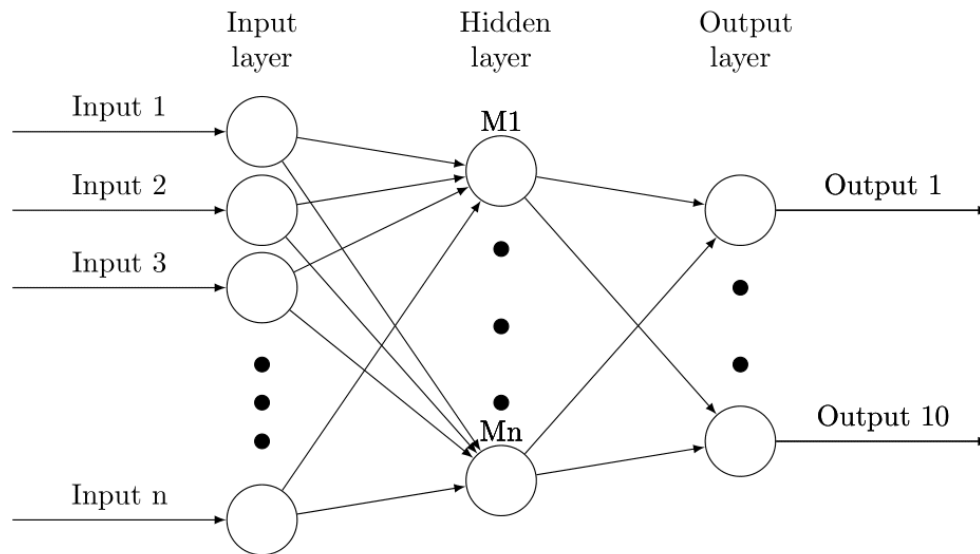


Figure 1: Neural Network with one hidden layer

12 Feed Forward - Cost Function

The cost function (logLikelihood plus regularization term) we want to maximize for the problem of classifying N number of data in K categories/classes is:

$$E(W) = \sum_{n=1}^N \sum_{k=1}^K t_{nk} \log s_{nk} - \frac{\lambda}{2} \left[\left(\sum_{k=1}^K \|\mathbf{w}_k^{(2)}\|^2 \right) + \left(\sum_{j=1}^M \|\mathbf{w}_j^{(1)}\|^2 \right) \right],$$

where s_{nk} is the softmax function defined as:

$$s_{nk} = \frac{e^{y_{nk}}}{\sum_{j=1}^K e^{y_{nj}}},$$

where y_{nk} is the linear combination of the parameters in the hidden layer defined as:

$$y_{nk} = \mathbf{z}_n (\mathbf{w}_k^{(2)})^T$$

where z_n is the output of the selected activation function in the input layer defined as:

$$\mathbf{z}_n(a), \quad a = \mathbf{x}_n (\mathbf{w}_j^{(1)})^T,$$

$\mathbf{W}^{(2)}$ is a $K \times (M + 1)$ matrix, where each line represents the vector $\mathbf{w}_k^{(2)}$,

$\mathbf{W}^{(1)}$ is a $(M + 1) \times (D + 1)$ matrix, where each line represents the vector $\mathbf{w}_j^{(1)}$.

The cost function can be simplified in the following form:

$$E(W) = \sum_{n=1}^N \left[\left(\sum_{k=1}^K t_{nk} (\mathbf{z}_n (\mathbf{w}_k^{(2)})^T) \right) - \log \left(\sum_{j=1}^K e^{\mathbf{z}_n (\mathbf{w}_j^{(2)})^T} \right) \right] - \frac{\lambda}{2} \left[\left(\sum_{k=1}^K \|\mathbf{w}_k^{(2)}\|^2 \right) + \left(\sum_{j=1}^M \|\mathbf{w}_j^{(1)}\|^2 \right) \right],$$

In the above formula we have used the fact that $\sum_{k=1}^K t_{nk} = 1$.

We use the logsumexp trick, where m is the maximum element:

$$\log \sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{z}_n} = \log \left(\sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{z}_n + m - m} \right) \quad (2)$$

$$= \log \left(\sum_{j=1}^K e^m e^{\mathbf{w}_j^T \mathbf{z}_n - m} \right) \quad (3)$$

$$= \log \left(e^m \sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{z}_n - m} \right) \quad (4)$$

$$= \log e^m + \log \left(\sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{z}_n - m} \right) \quad (5)$$

$$= m + \log \left(\sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{z}_n - m} \right) \quad (6)$$

13 Partial Derivatives of $\mathbf{w}^{(1)}$ & $\mathbf{w}^{(2)}$ Values

The $\mathbf{w}^{(2)}$ and $\mathbf{w}^{(1)}$ values arise from the following variables of the Cost Function:

$\mathbf{w}^{(2)}$: $s_{nk} \Rightarrow y_{nk} \Rightarrow w_k^{(2)}$ and the Regularization Term &

$\mathbf{w}^{(1)}$: $s_{nk} \Rightarrow y_{nk} \Rightarrow z_n \Rightarrow a \Rightarrow w_j^{(1)}$ and the corresponding Regularization Term

So, the partial derivatives of the values $\mathbf{W}^{(2)}$ of the Cost Function are given by the following equation:

$$\frac{\partial E}{\partial \mathbf{W}^{(2)}} = \frac{\partial E}{\partial S} \frac{\partial S}{\partial Y} \frac{\partial Y}{\partial \mathbf{W}^{(2)}}, \quad (7)$$

where

$$\frac{\partial E}{\partial Y} = (T - S)^T, \text{ a } K \times N \text{ matrix} \quad (8)$$

&

$$\frac{\partial Y}{\partial \mathbf{W}^{(2)}} = Z, \text{ a } N \times (M + 1) \text{ matrix} \quad (9)$$

\hookrightarrow **So, the final result for $\mathbf{W}^{(2)}$ is a $K \times (M+1)$ matrix as follows:**

$$(\mathbf{T} - \mathbf{S})^T \times \mathbf{Z}$$

where T is a $N \times K$ matrix with the truth values of the training data, such that $[T]_{nk} = t_{nk}$, S is the corresponding $N \times K$ matrix that holds the softmax probabilities and Z is the $N \times (M + 1)$ matrix that holds the output of the activation function in the input layer.

As for $W^{(1)}$, the partial derrivatives of these values are given by the following equation:

$$\frac{\partial E}{\partial W^{(1)}} = \frac{\partial E}{\partial S} \frac{\partial S}{\partial Y} \frac{\partial Y}{\partial Z} * \frac{\partial Z}{\partial A} \frac{\partial A}{\partial W^{(1)}}, \quad (10)$$

where

$$\frac{\partial E}{\partial Y} = (T - S), \text{ a } N \times K \text{ matrix}, \quad (11)$$

$$\frac{\partial Y}{\partial Z} = W^{(2)}, \text{ a } K \times (M + 1) \text{ matrix}, \quad (12)$$

$$\frac{\partial Z}{\partial A} = Z'(A), \text{ a } N \times (M + 1) \text{ matrix}, \quad (13)$$

$$\frac{\partial A}{\partial W^{(1)}} = X, \text{ a } N \times (D + 1) \text{ matrix} \quad (14)$$

\hookrightarrow So, the final result for $W^{(1)}$ is a $(M+1) \times (D+1)$ matrix as follows:

$$((T - S) W^{(2)} * Z'(A))^T X$$

$(*)$: *element – wise product*

where T is the matrix with the truth values of the training data, such that $[T]_{nk} = t_{nk}$, S is the corresponding matrix that holds the softmax probabilities, $W^{(2)}$ is the matrix with the values of weights between the hidden layer and the output layer, $Z'(A)$ is the matrix with the derivative of the selected activation function and X is the matrix of the input data.

```

1 def cost_grad_softmax(w1, w2, batchX, activation, batchY, lamda):
2
3     """ Compute the cost function and the partial derivatives of the weights.
4
5     :param w1: The (M+1) x (D+1) matrix with the values of weights between the
6     input layer and the hidden layer
7     :param w2: The K x (M+1) matrix with the values of weights between the
8     hidden layer and the output layer
9     :param batchX: The Nb x (D+1) matrix with the current mini batch of data
10    :param activation: The chosen activation function
11    :param batchY: The Nb x K matrix with the binary labels of the data
12    :param lamda: The positive regularization parameter
13
14    :return:
15        E(w): the cost of the current mini batch,
16        grad_w1: (M+1) x (D+1) matrix with the partial derivatives of the
17        weights w1 and
18        grad_w2: K x (M+1) matrix with the partial derivatives of the weights w2
19
20    """
21
22    a = batchX.dot(w1.T)
23
24    z, grad_z = activations(activation, a)
25    y = z.dot(w2.T)
26    s = softmax(y)
27
28    max_error = np.max(y, axis=1)
29
30    # Compute the cost function to check convergence
31    # Using the logsumexp trick for numerical stability
32    Ew = np.sum(batchY * y) - np.sum(max_error) - \
33        np.sum(np.log(np.sum(np.exp(y - np.array([max_error, ] * y.shape[1]).T
34        ), 1))) - \
35        (0.5 * lamda) * (np.sum(np.square(w1)) + np.sum(np.square(w2)))
36
37    # Calculate gradient for w2
38    grad_w2 = (batchY - s).T.dot(z) - lamda * w2
39    # Calculate gradient for w1
40    grad_a = (batchY - s).dot(w2)
41    grad_w1 = np.multiply(grad_a, grad_z).T.dot(batchX) - lamda * w1
42
43    return Ew, grad_w1, grad_w2

```


14 Generate batches

In case of Big Data we can apply ***Stochastic Gradient Ascent*** i.e. we assume that data are coming in small batches each time, instead of one sample and we make an update of the parameters by using a mini-batch of the data set. The below function generates a mini-batch every time it is called.

```
1 def next_batch(X, y, batchSize):
2
3     """ Generate mini batches of the training examples.
4
5     :param X: The N x (D+1) input matrix with the training examples
6     :param y: The N x K matrix with binary labels of the examples in X
7     indicating the 10 categories
8     :param batchSize: The chosen size of the batches
9     :return:
10         X: the produced batch containing some of the examples of X
11         y: the respectively labels of the examples in the batch
12
13     """
14
15     # Loop over our data set 'X' in mini-batches of size 'batchSize'
16     for i in np.arange(0, X.shape[0], batchSize):
17         # Yield a tuple of the current batched data and labels
18         yield (X[i:i + batchSize], y[i:i + batchSize])
```

15 Runner Method - Back Propagation

The below function runs the whole procedure. At each iteration it generates stochastic mini batches, that are fed to the algorithm in the above function `*cost_grad_softmax()*`. Through the training of a mini-batch, we calculate the cost of our predictions and the partial derivatives of the weights. Then, we update the weights using the derivatives and continue to the next mini-batch of the same iteration. When all mini-batches are trained in one iteration, we keep the last cost and move on to the next iteration using different mini-batches (randomly chosen).

```
1 def ml_softmax_train(t, X, lamda, w1_init, w2_init, options, activation):
2
3     """ Back Propagation: Call cost_grad_softmax function and update the values
4         of weights.
5
6         :param t: The N x K matrix with binary labels of the examples in X
7                   indicating the 10 categories
8         :param X: The N x (D+1) input data matrix with ones already added in the
9                   first column
10        :param lamda: The positive regularizarion parameter
11        :param w1_init: The (M+1) x (D+1) matrix with the initial values of the
12                       parameters w1
13        :param w2_init: The K x (M+1) matrix with the initial values of the
14                       parameters w2
15        :param options: options(1) is the maximum number of iterations
16                       options(2) is the tolerance
17                       options(3) is the learning rate eta
18                       options(4) is the size of batches
19        :param activation: The chosen activation function
20        :return:
21                w1: the trained (M+1) x (D+1) matrix of the parameters w1
22                w2: the trained K x (M+1) matrix of the parameters w2
23                costs: a list containing all the results by cost function
24
25        """
26
27        # Generate the initial weights w1 & w2 randomly using Xavier initialization
28        # method
29        w1 = np.random.rand(*w1_init.shape) * np.sqrt(1/w1_init.shape[1])
30        w2 = np.random.rand(*w2_init.shape) * np.sqrt(1/w2_init.shape[1])
31
32        # Maximum number of iteration of gradient ascend
33        _iter = options[0]
34
35        # Tolerance
36        tol = options[1]
```

```

33 # Learning rate
34 eta = options[2]
35 # Size of batches
36 batchSize = options[3]
37
38 Ewold = -np.inf
39 costs = []
40
41 for i in range(1, _iter+1 ):
42
43     # Shuffle randomly the training examples and respectively their labels
44     # on each iteration ,
45     # in order to implement stochastic gradient ascent
46     permutation = np.random.permutation(len(X))
47     X = X[permutation,:]
48     t = t[permutation,:]
49
50     # Generate mini-batches after shuffling the data set
51     for (batchX, batchY) in next_batch(X, t, batchSize):
52
53         # Calculate cost and partial derivatives of the parameters w1 & w2
54         # for each mini-batch and iteration
55         Ew, grad_w1, grad_w2 = cost_grad_softmax(w1, w2, batchX, activation
56         , batchY, lamda)
57
58         # Update parameters based on gradient ascent
59         w1 = w1 + eta * grad_w1
60         w2 = w2 + eta * grad_w2
61
62         # Save cost produced by the last mini batch
63         costs.append(Ew)
64
65         # Show the current cost function on screen
66         if i % 50 == 0:
67             print('Iteration : %d, Cost function :%f ' % (i, Ew))
68
69         # Break if you achieve the desired accuracy in the cost function
70         if np.abs(Ew - Ewold) < tol:
71             break
72
73     Ewold = Ew
74
75 return w1, w2, costs

```

16 Select Activation Function

Through the three options described above (logarithm, tangent, cosine), user has to select the desired one from the standard input.

```
1 def selectActivation():
2
3     """ Selects the number from standard input, which indicates the desired
4     activation function.
5
6     :return: The selected activation function
7
8     """
9
10    # While you do not select an integer between 1–3 continue
11    while True:
12        try:
13            # Read integer from stdin
14            act = int(input("Select 1 of the following numbers based on the
15            desired Activation Function:\n \
16            \n 1: logarithm function \n 2: tangent function \n
17            3: cosine function\n"))
18            print("")
19            if act in range(1,4):
20                break;
21            else:
22                raise ValueError('Invalid input. Please select again an integer
23                between 1–3!')
24        except ValueError:
25            print("")
26            print("Invalid input. Please select again an integer between 1–3!")
27            print("")
28
29    # Case1: input = 1
30    if act == 1:
31        activation = "log"
32        print("You selected Logarithm Activation Function!")
33    # Case2: input = 2
34    elif act == 2:
35        activation = "tan"
36        print("You selected Tangent Activation Function!")
37    # Case3: input = 3
38    elif act == 3:
39        activation = "cos"
40        print("You selected Cosine Activation Function!")
41
42    return activation
```

17 Gradcheck

When implementing gradient-based methods, it is suggested to include numerical gradient check (gradcheck).

Numerical approximation of the partial derivatives of $W^{(1)}$ & $W^{(2)}$:

$$\frac{\partial E}{\partial W^{(1)}} \approx \frac{E(W^{(1)} + \varepsilon) - E(W^{(1)} - \varepsilon)}{2\varepsilon} \quad (15)$$

$$\frac{\partial E}{\partial W^{(2)}} \approx \frac{E(W^{(2)} + \varepsilon) - E(W^{(2)} - \varepsilon)}{2\varepsilon} \quad (16)$$

($\varepsilon = 10^{-6}$)

```
1 def gradcheck_softmax(w1_init, w2_init, X, t, lamda, activation):
2
3     """ Check if the equations of the partial derivatives are correct.
4
5     :param w1_init: The (M+1) x (D+1) matrix with the initial values of the
6     parameters w1
7     :param w2_init: The K x (M+1) matrix with the initial values of the
8     parameters w2
9     :param X: The N x (D+1) input data matrix with ones already added in the
10    first column
11    :param t: The N x K matrix with binary labels of the examples in X
12    indicating the 10 categories
13    :param lamda: The positive regularization parameter
14    :param activation: The chosen activation function
15    :return:
16        grad_w1: The computed partial derivatives of the weights w1
17        numericalGrad1: The approximate values of the derivatives of the
18        weights w1
19        grad_w2: The computed partial derivatives of the weights w2
20        numericalGrad2: The approximate values of the derivatives of the
21        weights w2
22
23    """
24
25    # Generate the initial weights w1 & w2 randomly using Xavier initialization
26    method
27    w1 = np.random.rand(*w1_init.shape) * np.sqrt(1/w1_init.shape[1])
28    w2 = np.random.rand(*w2_init.shape) * np.sqrt(1/w2_init.shape[1])
29
30    epsilon = 1e-6
31
32    _list = np.random.randint(X.shape[0], size=5)
33    x_sample = np.array(X[_list, :])
```

```

27     t_sample = np.array(t[_list, :])
28
29     Ew, grad_w1, grad_w2 = cost_grad_softmax(w1, w2, x_sample, activation,
30     t_sample, lamda)
31
32     numericalGrad1 = np.zeros(grad_w1.shape)
33     numericalGrad2 = np.zeros(grad_w2.shape)
34
35     # Compute all numerical gradient estimates for w1 and store them in
36     # the matrix numericalGrad1
37     for k in range(numericalGrad1.shape[0]):
38         for d in range(numericalGrad1.shape[1]):
39
40             # Add epsilon to the w[k,d]
41             w_tmp1 = np.copy(w1)
42             w_tmp1[k, d] += epsilon
43             e_plus1, _, _ = cost_grad_softmax(w_tmp1, w2, x_sample, activation
44             , t_sample, lamda)
45
46             # Subtract epsilon to the w[k,d]
47             w_tmp1 = np.copy(w1)
48             w_tmp1[k, d] -= epsilon
49             e_minus1, _, _ = cost_grad_softmax(w_tmp1, w2, x_sample,
50             activation, t_sample, lamda)
51
52             # Approximate gradient ( E[ w[k,d] + theta ] - E[ w[k,d] - theta ]
53             ) / 2*e
54             numericalGrad1[k, d] = (e_plus1 - e_minus1) / (2 * epsilon)
55
56     # Compute all numerical gradient estimates for w2 and store them in
57     # the matrix numericalGrad2
58     for m in range(numericalGrad2.shape[0]):
59         for n in range(numericalGrad2.shape[1]):
60
61             # Add epsilon to the w[k,d]
62             w_tmp = np.copy(w2)
63             w_tmp[m, n] += epsilon
64             e_plus2, _, _ = cost_grad_softmax(w1, w_tmp, x_sample, activation,
65             t_sample, lamda)
66
67             # Subtract epsilon to the w[k,d]
68             w_tmp = np.copy(w2)
69             w_tmp[m, n] -= epsilon
70             e_minus2, _, _ = cost_grad_softmax(w1, w_tmp, x_sample, activation
71             , t_sample, lamda)
72
73             # Approximate gradient ( E[ w[k,d] + theta ] - E[ w[k,d] - theta ]
74             ) / 2*e
75             numericalGrad2[m, n] = (e_plus2 - e_minus2) / (2 * epsilon)
76
77     return (grad_w1, numericalGrad1, grad_w2, numericalGrad2)

```

18 Run Gradcheck Function

The below function runs the Gradcheck function and finds the estimated difference between the computed partial derivatives and the approximate ones. We expect insignificant difference between them, in order to ensure about the correctness of our computations, otherwise there is an error in the equations of the derivatives.

```
1 # N: number of training data
2 # D: number of feautres , plus the one of bias
3 N, D = X_train.shape
4
5 # Number of categories
6 K = 10
7 # Number of hidden units
8 M = 100
9
10 # Initialize w1 & w2 weights for gradient ascent , such it has same number of
    columns as our input features
11 w1_init = np.zeros((M, D))
12 w2_init = np.zeros((K, M))
13
14 # Regularization parameter
15 lamda = 0.1
16
17 # Activation function
18 activation = selectActivation()
19
20 # Calculate the partial derivatives and the approximate derivatives of the
    weights w1 & w2
21 grad_w1, numericalGrad1, grad_w2, numericalGrad2 = gradcheck_softmax(w1_init,
    w2_init, X_train, y_train, lamda, activation)
22
23 # Compare partial derivatives with the approximate ones and print the estimated
    difference
24 print( "The difference estimate for gradient of w1 is : ", np.max(np.abs(
    grad_w1 - numericalGrad1)) )
25 print( "The difference estimate for gradient of w2 is : ", np.max(np.abs(
    grad_w2 - numericalGrad2)) )
```

19 Training

```
1 # N: number of training data
2 # D: number of feautures , plus the one of bias
3 N, D = X_train.shape
4
5 # Number of categories
6 K = 10
7 # Number of hidden units
8 M = 300
9
10 # Initialize w1 & w2 weights for gradient ascent
11 w1_init = np.zeros((M, D))
12 w2_init = np.zeros((K, M))
13
14 # Regularization parameter
15 lamda = 0.01
16
17 # Check if activation variable is defined , otherwise select activation function
18 try:
19     activation
20 except NameError:
21     activation = selectActivation()
22
23 # options for gradient ascent
24 # i.e. maximum number of iterations , tolerance , learning rate , size of batches
25 options = [300, 1e-6, 0.001, 256]
26
27 # Train the model
28 w1, w2, costs = ml_softmax_train(y_train, X_train, lamda, w1_init, w2_init,
    options, activation)
```


20 Plot Cost Function

```
1 # Plot cost versus number of iterations
2 plt.plot(np.squeeze(costs))
3 plt.ylabel('cost')
4 plt.xlabel('iterations (per tens)')
5 plt.title("Learning rate =" + str(format(options[2], 'f')))
6 plt.show()
```

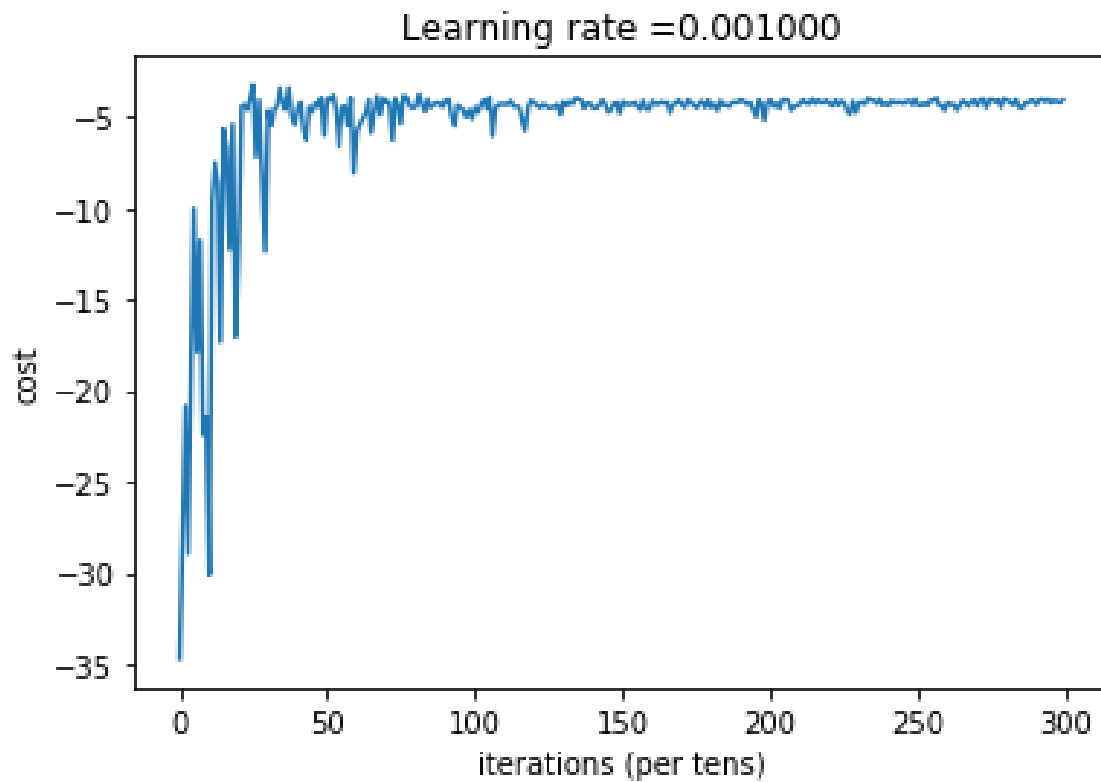


Figure 3: Cost Function for MNIST Data Set

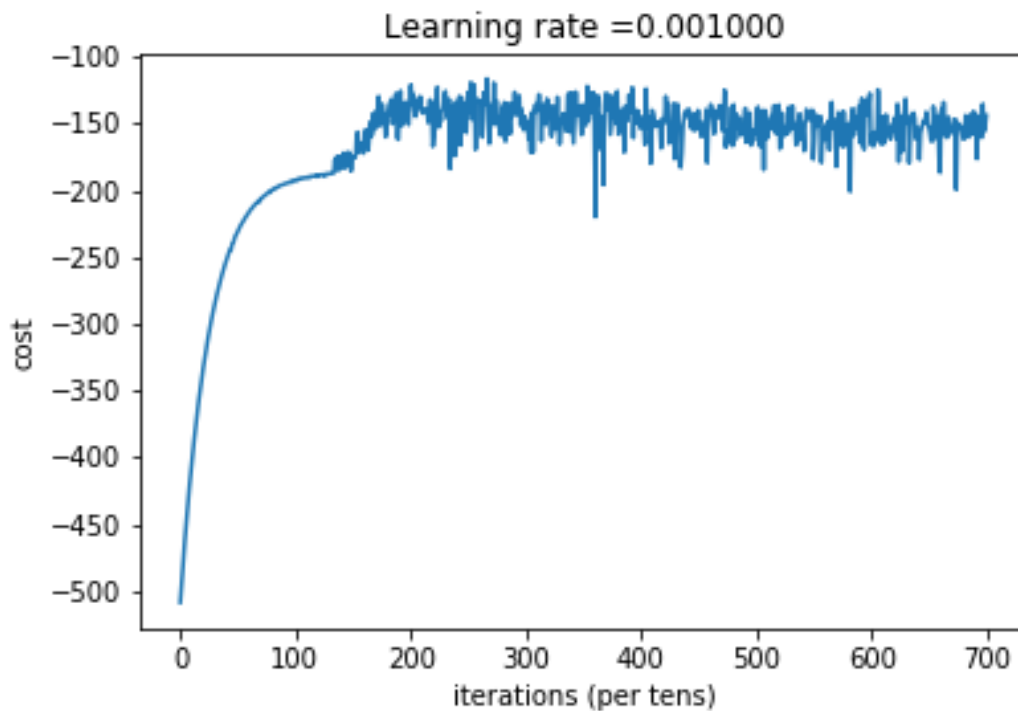


Figure 4: Cost Function for Cifar Data Set

21 Testing

After training our model, we have to test it using the test data set without the labels.

```

1 def ml_softmax_test(w1, w2, X_test, activation):
2
3     """ Finds the category in which each example belong to.
4
5     :param w1: The trained (M+1) x (D+1) matrix of the parameters w1
6     :param w2: The trained K x (M+1) matrix of the parameters w2
7     :param X_test: The matrix with the data that will be used for testing
8     :param activation: The chosen activation function
9     :return:
10         ttest: The (N x 1) matrix that contains the predict category of each
11         example
12
13     """
14     a = X_test.dot(w1.T)

```

```

15     z, grad_z = activations(activation , a)
16     y = z.dot(w2.T)
17     ytest = softmax(y)
18     # Keep the position with the biggest probability , as the category a test
    example belongs to
19     ttest = np.argmax(ytest , 1)
20
21     return ttest

```

21.1 Call Test Function for Train Data Set

```

1 pred = ml_softmax_test(w1, w2, X_train , activation)

```

Train Accuracy

```

1 # Compare our predictions with the real values and compute the train accuracy
    of the model
2 np.mean( pred == np.argmax(y_train,1) )

```

21.2 Call Test Function for Test Data Set

```

1 pred = ml_softmax_test(w1, w2, X_test , activation)

```

Test Accuracy

```

1 # Compare our predictions with the real values and compute the test accuracy of
    the model
2 np.mean( pred == np.argmax(y_test,1) )

```

22 Misclassified Test Data

22.1 Plot Mnist's Faults

```

1 def plot_mnists_faults():
2
3     """ Plot 25 random misclassified images from the Mnist training set. """
4
5
6     faults = np.where(np.not_equal(np.argmax(y_test,1),pred))[0]
7
8     # plot 25 misclassified examples from the test set
9     n = 25
10    samples = np.random.choice(faults , n)
11    sqrt_n = int( n ** 0.5 )
12
13    plt.figure( figsize=(11,13) )

```

```

14
15     cnt = 0
16     for i in samples:
17         cnt += 1
18         plt.subplot( sqrt_n, sqrt_n, cnt )
19         plt.subplot( sqrt_n, sqrt_n, cnt ).axis( 'off' )
20         plt.imshow( X_test[i,1:].reshape(28,28)*255, cmap='gray' )
21         plt.title( "True: "+str(np.argmax(y_test,1)[i])+ "\n Predicted: "+ str(
pred[i]))
22
23     plt.show()

```

22.2 Plot Cifar-10 Faults

```

1 def plot_cifar_faults():
2
3     """ Plot 25 random misclassified images from the Cifar-10 training set. """
4
5     faults = np.where(np.not_equal(np.argmax(y_test,1),pred))[0]
6
7     n = 25
8     sqrt_n = int( n**0.5 )
9     samples = np.random.choice( faults , n)
10
11     plt.figure( figsize=(15,15) )
12
13     cnt = 0
14     for i in samples:
15         arr = X_test[i, 1:] * 255
16         R = arr[0:1024].reshape(32,32)/255.0
17         G = arr[1024:2048].reshape(32,32)/255.0
18         B = arr[2048:].reshape(32,32)/255.0
19
20         img = np.dstack((R,G,B))
21         cnt += 1
22         plt.subplot( sqrt_n, sqrt_n, cnt )
23         plt.subplot( sqrt_n, sqrt_n, cnt ).axis( 'off' )
24         plt.imshow(img,interpolation='bicubic')
25         plt.title( "True: "+str(np.argmax(y_test,1)[i])+ "\n Predicted: "+ str(
pred[i]))
26
27     plt.show()

```

22.3 View Misclassified Data

```

1 if(data_set == 1):
2     plot_mnists_faults()
3 elif(data_set == 2):
4     plot_cifar_faults()

```

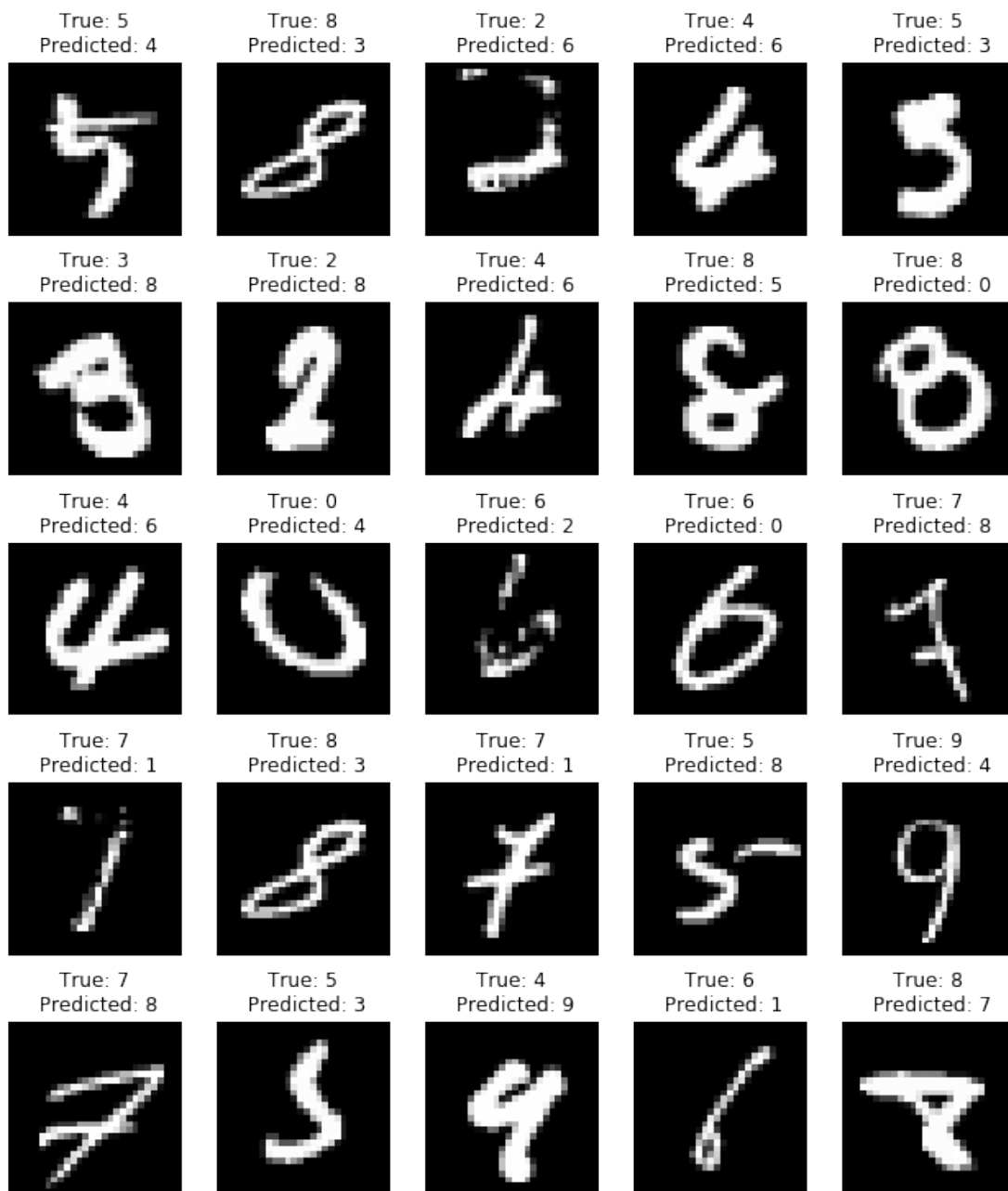


Figure 5: Mnist Misclassified Data

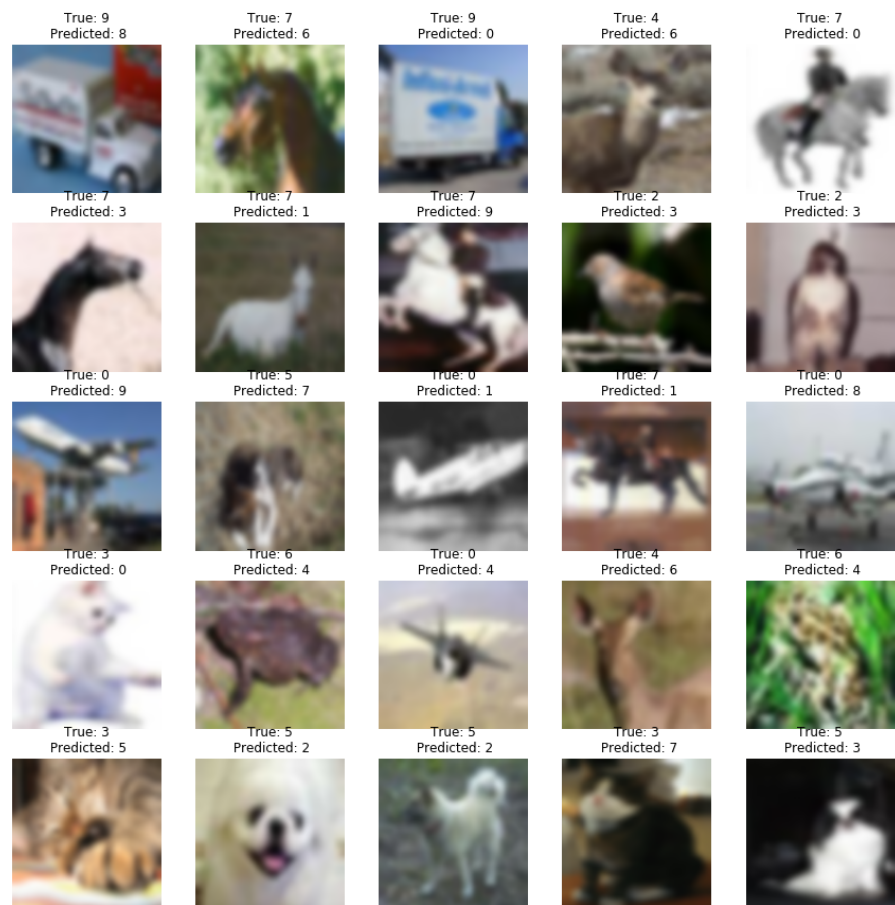


Figure 6: CIFAR-10 Misclassified Data

23 Fine-Tuning & Results

23.1 MNIST Data Set

Activation Function	Hidden Units	Iterations	Learning Rate	lamda	Accuracy
<i>tan</i>	100	500	10^{-3}	10^{-1}	0.9797
<i>tan</i>	200	500	10^{-3}	10^{-1}	0.9816
<i>tan</i>	300	700	10^{-3}	10^{-1}	0.9802
<i>tan</i>	300	500	10^{-4}	10^{-1}	0.9769
<i>tan</i>	300	1000	10^{-4}	10^{-1}	0.9065
<i>tan</i>	200	500	10^{-2}	10^{-1}	0.9313
<i>tan</i>	200	500	10^{-3}	10^{-2}	0.9827
<i>tan</i>	300	300	10^{-3}	10^{-2}	0.9834
<i>tan</i>	200	300	10^{-3}	10^{-2}	0.9794
<i>log</i>	100	500	10^{-3}	10^{-1}	0.9749
<i>log</i>	200	500	10^{-3}	10^{-1}	0.9746
<i>log</i>	300	500	10^{-3}	10^{-1}	0.9734
<i>log</i>	300	700	10^{-3}	10^{-1}	0.9777
<i>log</i>	300	500	10^{-4}	10^{-1}	0.9724
<i>log</i>	200	350	10^{-2}	10^{-1}	0.7865
<i>log</i>	200	500	10^{-2}	10^{-1}	0.9597
<i>log</i>	200	500	10^{-3}	10^{-2}	0.9508
<i>log</i>	300	300	10^{-3}	10^{-2}	0.9421
<i>log</i>	200	300	10^{-3}	10^{-2}	0.9288
<i>cos</i>	100	500	10^{-3}	10^{-1}	0.9828
<i>cos</i>	200	500	10^{-3}	10^{-1}	0.9604
<i>cos</i>	300	500	10^{-3}	10^{-1}	0.9763
<i>cos</i>	300	700	10^{-3}	10^{-1}	0.9813
<i>cos</i>	300	500	10^{-4}	10^{-1}	0.9809
<i>cos</i>	300	1000	10^{-4}	10^{-1}	0.9821
<i>cos</i>	200	500	10^{-2}	10^{-1}	0.4604
<i>cos</i>	200	500	10^{-3}	10^{-2}	0.9817
<i>cos</i>	300	300	10^{-3}	10^{-2}	0.9822
<i>cos</i>	200	300	10^{-3}	10^{-2}	0.9798

Figure 7: MNIST Results

In Mnist data set our model achieved an accuracy over 97% in most cases.

23.2 CIFAR-10 Data Set

Activation Function	Hidden Units	Iterations	Learning Rate	lamda	Accuracy
<i>tan</i>	100	500	10^{-3}	10^{-1}	0.487
<i>tan</i>	200	500	10^{-3}	10^{-1}	0.4843
<i>tan</i>	300	700	10^{-3}	10^{-1}	0.4788
<i>tan</i>	300	1000	10^{-3}	10^{-1}	0.4861
<i>tan</i>	200	500	10^{-3}	10^{-2}	0.4644
<i>tan</i>	200	300	10^{-3}	10^{-2}	0.4761
<i>log</i>	100	500	10^{-3}	10^{-1}	0.478
<i>log</i>	200	500	10^{-3}	10^{-1}	0.4965
<i>log</i>	300	700	10^{-3}	10^{-1}	0.5141
<i>log</i>	300	500	10^{-4}	10^{-1}	0.4758
<i>log</i>	200	500	10^{-2}	10^{-1}	0.1435
<i>log</i>	200	500	10^{-3}	10^{-2}	0.1104
<i>log</i>	300	300	10^{-3}	10^{-2}	0.1242
<i>log</i>	200	300	10^{-3}	10^{-2}	0.1
<i>cos</i>	100	500	10^{-3}	10^{-1}	0.1152
<i>cos</i>	200	500	10^{-3}	10^{-1}	0.1221
<i>cos</i>	300	500	10^{-3}	10^{-1}	0.1352
<i>cos</i>	300	700	10^{-3}	10^{-1}	0.1267
<i>cos</i>	300	1000	10^{-3}	10^{-1}	0.1207
<i>cos</i>	200	500	10^{-2}	10^{-1}	0.10156
<i>cos</i>	200	500	10^{-3}	10^{-2}	0.1

Figure 8: CIFAR-10 Results

The accuracy of the model on Cifar-10 is extremely low peaking at 51% accuracy due to the complexity of the data set and especially in case of changing the learning rate parameter. In contrast to Mnist data set, Cifar-10 has fewer training examples and much more feautres (colour images), so their score difference is completely justified. In order to enhance the performance of the model, we could add a second hidden layer. The cosine activation function is the least effective one and it is not suggested.