# EEE 543: Neural Networks -Final Project Report

Oytun Gunes, ID: 20901170, Mucahit Gumus ID: 20900664

July 20, 2019

## Abstract

In this report 10 class classification task is done on commonly used CIFAR-10 data. A data analysis is performed before the design, cross correlation of the images are examined. Two different methods are implemented: conjugate gradient cost function and line search and single layer convolutional neural network (CNN). Many hyperparameters are optimized to obtain accuracies. In the results part cross entropy cost, validation accuracy and confusion matrices are illustrated. 30 percent test accuracy is obtained with the first method. CNN increased the test accuracy up to 40 percent.

## 1  Introduction

In the last five years CNN (Convolutional Neural Networks) has achieved outperforming results compared to standard classification algorithms.This has been started with the breaktrough of ImageNet classification results [1]. Soon after, Donahue et al. [2] have display the results of this network on several image dataset. They have found that CNN is an effective feature extractor. kaggle.com, a popular website in machine learning, has organized a competition on CIFAR-10 data. It is very interesting to work on a popular image data and perform 10 class classification. Our expectation on the task is to obtain a high accuracy of course. However, most people in the literature have constructed a complex and at least 3-5 layers on libraries on CNN structure to obtain a higher accuracy, whereas in our task the challenge is to code a complex architecture without using libraries and toolboxes and optimize parameters on a big data.

## 2  Methods

### 2.1  Analysis on Data

In this project the aim is to classify objects and animals correctly given training data. The data given is a popular data in the literature known as CIFAR-10. It consists of 50000 training and 10000 test data. There are 10 classes in the data. The sample images from each class is shown below.
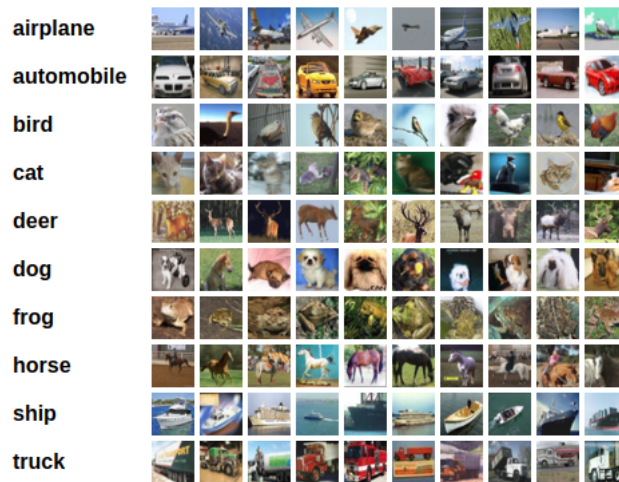


Figure 1: Sample Images from the Data

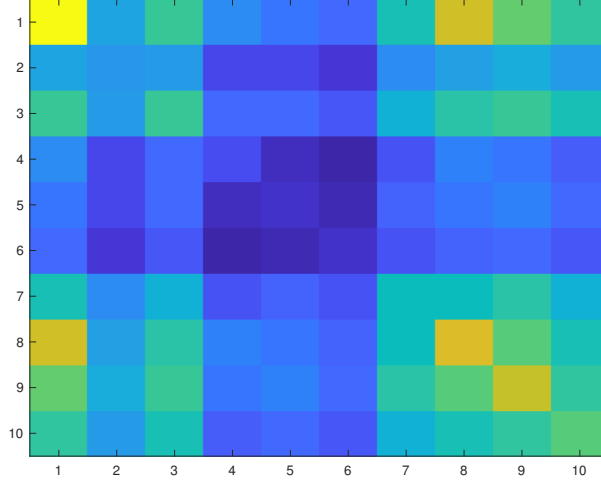The correlation between the image is shown below.

Figure 2: Visualization of Correlation Matrix of different classes

The correlation matrix above explains the cross correlation between the images. In the diagonals it explains the autocorrelation of the images which are high as expected. Off diagonal entries of the correlation matrix is also high which means there is high correlation among different classes. Therefore, classification over the dataset becomes challenging.

## 2.2 Training with conjugate gradient cost function and line search

We firstly tried a fully connected multi-layer neural network structure in order for classifying the data of 10 kind of objects. In this multi-layer structure, we have 40 neurons in the first hidden layer, 20 neurons in the second hidden layer and 10 neurons are used in the output layer in order to classify 10 objects. At the output, we implement softmax decision and calculate the cross entropy as an error function and using backpropogation we found the update rules for the weights and biases. The error is defined as:

$$E = -\sum_{i=1}^{10} (d_i) log(o_i)$$

where d is the desired vector, the correct object's index is 1 and others' are 0. The derivative of softmax function is found and as $(o_i - d_i)$. The gradient of the output layer is found as:

$$\frac{\partial E}{\partial W_{ji}} = \frac{\partial E}{\partial z_i} \frac{\partial z_i}{\partial W_{ji}}$$

$$\frac{\partial E}{\partial W_{ji}} = (o_i - d_i)h_j$$

where $h$ is the output of the hidden layer.

The gradient of the second hidden layer is found again by the chain rule as:

$$\frac{\partial E}{\tilde{W}_{aj}} = z_a h_j (1 - h_j)\lambda \sum_{i=1}^{10} (o_i - d_i)W_{ji}$$

Continuing in the same manner, gradient for the first layer is also found.

Then, those gradient values and the current value of the error is used in the fmincg function to calculate the updates. In fmincg, The Polack- Ribiere flavour of conjugate gradients is used to compute search directions, and a line search using quadratic and cubic polynomial approximations and the Wolfe-Powell stopping criteria is used together with the slope ratio method for guessing initial step sizes.

Before training, we firstly preprocess data in order to protect from saturation and decreasing complexity. Hence, we modify the pictures to gray scale and subtract the mean. Lastly, we zip the data to [-3std, 3 std] interval.

2

## 2.3    Convolutional Neural Networks (CNN) Explanation [3]

Convolutional Network as mentioned above is a highly effective method in classification. It solves some of the problems in feed-forward neural network approach. One of them is fully connected structure has an increasing number of parameters since each node in layer L is connected to a node in layer L-1. It becomes unmanagable to work on with a huge number of parameters. Secondly, computing the linear activations of the hidden units would be computationally costly. Local connectivity of CNN solves these problems.

**Local Connectivity:**

Each hidden unit is connected only to a subregion(patch) of the image. CNN process patches independently.
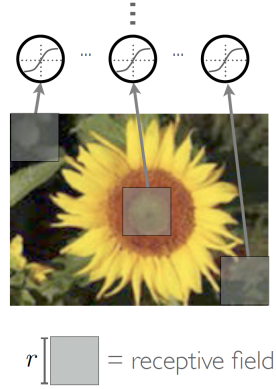


Figure 3: Patches in the Image

In addition, units are connected to all channels. 1 channel if grayscale image, 3 channels (R, G, B) if color image as shown below.
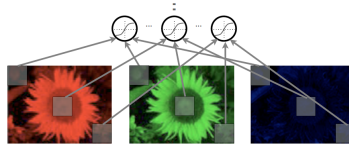


Figure 4: Channels

**Parameter Sharing:**

In CNN it shares a matrix of parameters across certain units. Units organized into the same feature map share parameters. Hidden units within a feature map cover different positions in the image. This reduces the number of parameters and it will help to extract features at every position. Feature maps of the image is shown below.
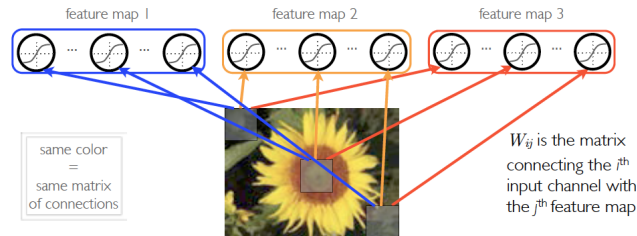


Figure 5: Feature Maps in the Image

**Convolution Layer:**

Feature maps described above is computed with a discrete convolution ( *) of a kernel matrix $k_{ij}$ which is the hidden weights matrix $W_{ij}$ with its rows and columns flipped.
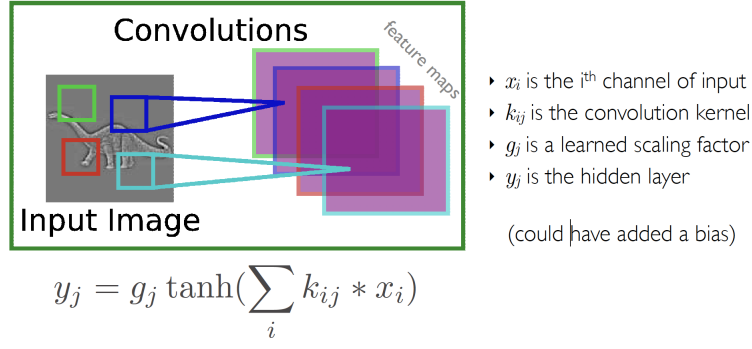
$$y_j = g_j \tanh\left(\sum_i k_{ij} * x_i\right)$$

- $x_i$ is the i<sup>th</sup> channel of input
- $k_{ij}$ is the convolution kernel
- $g_j$ is a learned scaling factor
- $y_j$ is the hidden layer

(could have added a bias)

Figure 6: Obtaining Feature maps

The convolution operation of an image $x$ with a kernel $k$ is as below:



Figure 7: Convolution Operation on the Image

Convolution operation in the image can be considered as the weighted sum between two signals. In image processing the convolution of the location (x,y) is done by extracting a patch of kxk dimension small image. Then filter this patch with a filter which has dimensions also kxk by multiplying elementwise and adding the values and obtain the result. Similarly this operation should be done for all kxk pairs in the image by sliding one by one. After the convolution layer it passes through an activation function such as sigmoid or commonly Relu.

**Max Pooling Layer:** Max pooling layer is done after the convolution layer in order to reduce further the number of parameters and introduces invariance to local translations. As an illustration 2x2 max pooling is done on the image below. The maximum element is selected.



Figure 8: Max pooling Illustration

**Output Layer:** Output layer is a regular, fully connected layer with softmax non-linearity output provides an estimate of the conditional probability of each class. The network is trained by stochastic gradient descent.

**Gradient of Convolution Layer:**

We have used an error function of cross-entropy. Using the backpropogation we have updated the kernels and output weights. Let l be the loss function. For convolution operation $y_j = x_i * k_{ij}$ the gradient for $x_i$ is :

$$\nabla_{x_i} l = \sum_j (\nabla_{y_j} l) * (W_{ij})$$

where * is the convolution with zero padding and $x_i$ is the row/column flipped version of $x_i$ .

4

Gradient for $W_{ij}$:

$$\nabla_{W_{ij}} l = (\nabla_y l) * \tilde{x}_i$$

**Gradient of Pooling Layer:**
Let $l$ be the loss function:
For max pooling operation $y_{ijk} = \max_{p,q} x_{i,j+p,k+q}$ the gradient for $x_{ijk}$.

$$\nabla_{x_{ijk}} l = 0 \text{ except for } \nabla_{x_{i,j+p',k+q'}} l = \nabla_{y_{ijk}} l$$

The overall structure of CNN is to alternate between the convolution and pooling layers as shown below:



Figure 9: Overall CNN architecture

**Simulation Setup:**
Our simulation setup for CNN is as follows:
Single convolutional layer with 16 kernels with a size of 5x5.
Sigmoid layer with $\lambda = 1$
Max pooling layer with 2x2
Fully connected layer 3136x10
Output Layer 10x1
Softmax Layer
**Parameter Selection:**
Parameters are selected using the toolbox of MATLAB since it is computationally efficient.

# 3 Results

## 3.1 Results of conjugate gradient cost function and line search

In figure 10, we have visualized the first hidden layer weights as images. We have observed that the weights are distributed differently, however we could not observe specific features. That is because training with the fully connected structure and line seaerch methedology is not successful enough to extract hidden features of the objects although these neural structure achieves other classification problems such as letter recognition well.

Figure 10: Visualization of hidden weights after training.

In figure 11 and 12, we have visualized confusion matrix for train and test data. We have observed that confusion matrices are similar for both of the dataset which shows that there is no overfitting with the learning. It also demonstrates that decision of cat is always less than other classes hence algorithm decided on dog mostly for cat images. Accuracy in classifying dog and ship images is higher than other 8 classes.



Figure 11: Confusion Matrix after training.



Figure 12: Confusion Matrix after testing.

As it is also obvious from the confusion matrix, the classification with this method does not achieve a good rate hence accuracy is around 30 %.

## 3.2 Results of CNN

In this section we present different types of results to demonstrate training and validation performance of CNN algorithm described in methods section. Training 5 epochs lasted around 50 minutes without using GPU.

Firstly, cross entropy cost over 5 epochs are shown below. It can be observed from the figure that training provides with decrease in cross entropy cost.



Figure 13: Cross Entropy Cost over epochs

Secondly, validation accuracy over epochs are shown below. It can be seen that validation accuracy increases up to 40%. It is obvious from figure 13 and **??** that validation accuracy and training cost are positively related.



Figure 14: Validation Accuracy over epochs

Thirdly the kernels of the convolution layers after training are shown below. Each kernel has extracted a different feature as, there are bright points at different pixels in the figure.



Figure 15: Final kernels of the convolution layer

Finally confusion matrix of the test data is shown below. As expected diagonals are higher than the off-diagonals which means that training provides to decide on true class for all of the 10 classes.



Figure 16: Confusion Matrix of test data

Implemented CNN structure provides 41 percent test accuracy.

# 4  Discussion

CIFAR-10 dataset is a challenging data as we have seen from the literature since it has objects in the image. Some people have obtained accuracies above 90 with a very complex structure of CNN 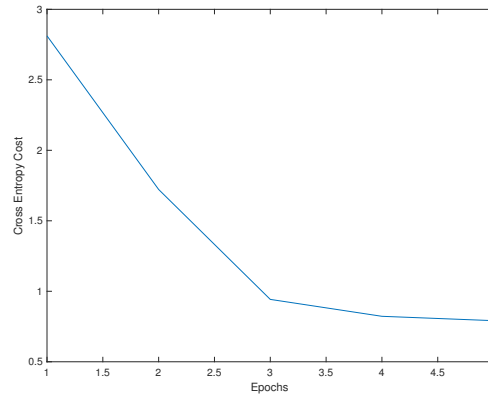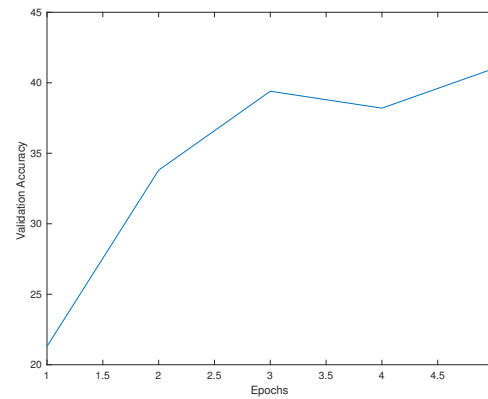with the help of phyton libraries. They have mostly used 3 layer convolution with dropout and batch normalization layers. Our design of single layer CNN reaches up to 40 percent which is comparably reasonable result for a single layer CNN. Optimizing hyper-parameters was a really tough challenge since we had around total of 10 parameters such as initializing weights, kernel size, max pooling size, learning rate, $\lambda$, momentum constant $alpha$ and activation function type. 50 minutes of training time did not allow us to try many parameters to obtain a higher accuracy. We have also tried 2 layer CNN since it takes around 1.5 hours to train all the data with CPU we could not optimize the parameters enough to obtain a higher accuracy compared to single layer CNN.

The first method we have used has a fully connected neural network structure and use the given fmincg function before, gave around 30 percent accuracy. Training time of this design is around 30 minutes for 5 epochs.

Of course, there are many ways to increase the test accuracy. One of them is to add data augmentation before the input layer and add more data to training data. Another one is to add rectification layer which

performs well in object recognition which makes the training translation invariant.

# References

[1] Deng, Jia, et al. "Imagenet: A large-scale hierarchical image database." Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on. IEEE, 2009.

[2] Donahue, Jeff, et al. "Decaf: A deep convolutional activation feature for generic visual recognition." International conference on machine learning. 2014.

[3] Lecture Notes of Neural Network Bilkent University 2017-2018 Fall.

# Appendix: MATLAB CODES

```matlab
1  clear;
2  close all;
3  clc;
4  %loading data and normalization
5  tic;
6  load('final_project_dataset.mat');
7  pixelNum=32;
8  for i=1:size(testdat,1)
9      aa=testdat(i,:);
10     testdatShaped(:,:,:,i)=reshape(aa, [pixelNum pixelNum 3]);
11 end
12 for i=1:size(traindat,1)
13     aa=traindat(i,:);
14     traindatShaped(:,:,:,i)=reshape(aa, [pixelNum pixelNum 3]);
15 end
16 % preprocess test data for w-b and normalize
17 nbWTrainData= zeros(pixelNum,pixelNum,size(traindatShaped,4));
18 for i=1:size(traindatShaped,4)
19     nbWTrainData(:,:,i)=rgb2gray(traindatShaped(:,:,:,i));
20 end
21 % preprocess train data for w-b and normalize
22 nbWTestData= zeros(pixelNum,pixelNum,size(testdatShaped,4));
23 for i=1:size(testdatShaped,4)
24     nbWTestData(:,:,i)=rgb2gray(testdatShaped(:,:,:,i));
25 end
26 nbWTrainData=double(nbWTrainData);
27 nbWTestData=double(nbWTestData);
28 %modify train and test labels by increasing 1
29 trainlbl = trainlbl +1;
30 testlbl = testlbl + 1;
31
32 %initialize weights and kernel
33 ker1L=5;
34 ker1Num=16;
35 pl1=2;
36 numOfHid= ker1Num*((pixelNum-ker1L+1)/pl1)^2;
37 numOfOut=10;
38 kerCons1= sqrt(6/(ker1L*ker1L+900));   % how to get 4 is ker1L^2*1/pl1^2
39 ker1= -1*kerCons1+ (kerCons1*2)*rand(ker1L,ker1L,ker1Num);
40 ker1=0.05*randn(ker1L,ker1L,ker1Num);
41 biasKer1=-1*kerCons1+ (kerCons1*2)*rand(pixelNum-ker1L+1,pixelNum-ker1L+1,
       ker1Num);
42 biasKer1=0.05*randn(pixelNum-ker1L+1,pixelNum-ker1L+1,ker1Num);
43 weightConsOut= sqrt(6/(numOfHid+900)); %how to get 200
```

```matlab
44  weightOut= −1*weightConsOut+ (weightConsOut*2)*rand(numOfHid,numOfOut);
45  biasOut=−1*weightConsOut+ (weightConsOut*2)*rand(1,numOfOut);
46  lambda=1;
47  eta=0.0001;
48  alpha=0;
49  toc;
50  epochNum=5;
51  miniBatchNum=128;
52
53
54  confusiontest=zeros(numOfOut,numOfOut);
55
56  accEpoch=zeros(1,epochNum);
57  CostEpoch=zeros(1,epochNum);
58
59  for epochIter =1:epochNum
60      eta=0.0001*(0.5^(epochIter−1));% decrease learning rate over epochs
61      tic;
62      picSequence = randperm(length(trainlbl));
63      prevWODel= zeros(numOfHid,numOfOut);
64      prevBODel=zeros(1,numOfOut);
65      prevK1Del=zeros(ker1L,ker1L,ker1Num);
66      prevBiasKer1Del=zeros(pixelNum−ker1L+1,pixelNum−ker1L+1,ker1Num);
67      % mini−batch
68      for batchIter =1:floor(length(picSequence)/miniBatchNum)
69          storeWODel= zeros(numOfHid,numOfOut);
70          storeBODel=zeros(1,numOfOut);
71          storeK1Del=zeros(ker1L,ker1L,ker1Num);
72          storeBiasKer1Del=zeros(pixelNum−ker1L+1,pixelNum−ker1L+1,ker1Num);
73          for mBIter =1:miniBatchNum
74              X=nbWTrainData(:,:,picSequence((batchIter −1)*miniBatchNum+mBIter));
75
76              %convolutional layer
77              conv1Out=zeros(pixelNum−ker1L+1,pixelNum−ker1L+1,ker1Num);
78              for ker1Iter =1:ker1Num
79                  conv1Out(:,:,ker1Iter) = conv2(X,ker1(:,:,ker1Iter),'valid');
80              end
81              conv1OutAct=1./(1+exp(−lambda*(conv1Out−biasKer1)));
82
83              % max pooling layer
84              pl1Out= zeros((pixelNum−ker1L+1)/pl1,(pixelNum−ker1L+1)/pl1,ker1Num)
                    ;
85              pl1OutIndex=zeros(2,(pixelNum−ker1L+1)/pl1,(pixelNum−ker1L+1)/pl1,
                    ker1Num);
86              for ker1Iter =1:ker1Num
87                  for i =1:(pixelNum−ker1L+1)/pl1
88                      for k=1:(pixelNum−ker1L+1)/pl1
89                          dumy= conv1OutAct((i −1)*pl1+1:i*pl1  ,(k−1)*pl1+1:k*pl1  ,
                                ker1Iter);
90                          [maxVal, maxInd] = max(dumy(:));
91                          [Irow, Icol] = ind2sub(size(dumy),maxInd);
92                          pl1Out(i,k,ker1Iter) = maxVal;
93                          pl1OutIndex(1,i,k,ker1Iter)=Irow;
94                          pl1OutIndex(2,i,k,ker1Iter)=Icol;
95                      end
96                  end
97              end
98              %fully connected layer
99              hiddenOut=pl1Out(:);
100             neuralOut=transpose(hiddenOut)*weightOut−biasOut;
101             % output layer
```

```matlab
102                sumOut=sum(exp(neuralOut));
103                neuralOutAct= exp(neuralOut)*1/(sumOut);
104            %backpropogation
105            d=zeros(1,numOfOut);
106            d(1,trainlbl(picSequence((batchIter-1)*miniBatchNum+mBIter)))=1;
107            gradOut= neuralOutAct - d;
108            gradHidden= weightOut*(transpose(gradOut));
109            gradpl1=reshape(gradHidden, [(pixelNum-ker1L+1)/pl1  (pixelNum-ker1L
                   +1)/pl1  ker1Num]);

110
111            gradconv1=zeros(pixelNum-ker1L+1,pixelNum-ker1L+1,ker1Num);
112            for ker1Iter=1:ker1Num
113                for i=1:(pixelNum-ker1L+1)/pl1
114                    for k=1:(pixelNum-ker1L+1)/pl1
115                        indRow=(i-1)*pl1+pl1OutIndex(1,i,k,ker1Iter);
116                        indCol=(k-1)*pl1+pl1OutIndex(2,i,k,ker1Iter);
117                        gradconv1(indRow,indCol,ker1Iter)=gradpl1(i,k,ker1Iter)*
                            conv1OutAct(indRow,indCol,ker1Iter)*(1-conv1OutAct(
                            indRow,indCol,ker1Iter));
118                    end
119                end
120            end
121            storeWODel=storeWODel -1*eta*(hiddenOut)*gradOut;
122            weightOut =weightOut+  -1*eta*(hiddenOut)*gradOut + alpha*prevWODel;
123            prevWODel=alpha*prevWODel+-1*eta*(hiddenOut)*gradOut;
124            storeBODel=storeBODel -1*eta*(-1)*gradOut;
125            biasOut=biasOut+ -1*eta*(-1)*gradOut;
126            prevBODel=alpha*prevBODel -1*eta*(-1)*gradOut;
127
128            ker1Update=zeros(ker1L,ker1L,ker1Num);
129            for ker1Iter=1:ker1Num
130                for i=1:ker1L
131                    for k=1: ker1L
132                        ind1=ker1L-i+1;
133                        ind2=pixelNum-i+1;
134                        ind3=k;
135                        ind4=k+size(conv1Out,1)-1;
136                        ker1Update(:,:,ker1Iter)=ker1Update(:,:,ker1Iter)-1*(eta
                            )*sum(sum(X(ind1:ind2,ind3:ind4).*gradconv1(:,:,
                            ker1Iter)));  %*1/((pixelNum-ker1L+1)^2)
137                    end
138                end
139            end
140            storeK1Del=storeK1Del+ker1Update;
141            storeBiasKer1Del=storeBiasKer1Del -1*eta*(-1)*gradconv1;
142            aaaaa=0;
143        end
144        % update the kernels and weights
145        ker1=ker1+(1/miniBatchNum)*storeK1Del+alpha*prevK1Del;
146        prevK1Del=alpha*prevK1Del+(1/miniBatchNum)*storeK1Del;
147        biasKer1=biasKer1+(1/miniBatchNum)*storeBiasKer1Del + alpha*
               prevBiasKer1Del;
148        prevBiasKer1Del=prevBiasKer1Del++(1/miniBatchNum)*storeBiasKer1Del;
149        weightOut =weightOut+  (1/miniBatchNum)*storeWODel + alpha*prevWODel;
150        prevWODel=alpha*prevWODel+(1/miniBatchNum)*storeWODel;
151        biasOut=biasOut+ (1/miniBatchNum)*storeBODel+ alpha*prevBODel;
152        prevBODel=alpha*prevBODel + (1/miniBatchNum)*storeBODel;
153
154    end
155    toc;
156
```

```matlab
157        %test validation data
158        tic;
159        testError=0;
160        testShuf=randperm(size(nbWTrainData,3));
161        for testIter=1:1000
162            X=nbWTrainData(:,:,testShuf(testIter));
163            conv1Out=zeros(pixelNum-ker1L+1,pixelNum-ker1L+1,ker1Num);
164            for ker1Iter=1:ker1Num
165                conv1Out(:,:,ker1Iter) = conv2(X,ker1(:,:,ker1Iter),'valid');
166            end
167            conv1OutAct=1./(1+exp(-lambda*(conv1Out-biasKer1)));
168            pl1Out= zeros((pixelNum-ker1L+1)/pl1,(pixelNum-ker1L+1)/pl1,ker1Num);
169            pl1OutIndex=zeros(2,(pixelNum-ker1L+1)/pl1,(pixelNum-ker1L+1)/pl1,
                   ker1Num);
170            for ker1Iter=1:ker1Num
171                for i=1:(pixelNum-ker1L+1)/pl1
172                    for k=1:(pixelNum-ker1L+1)/pl1
173                        dumy= conv1OutAct((i-1)*pl1+1:i*pl1  ,(k-1)*pl1+1:k*pl1  ,
                           ker1Iter);
174                        [maxVal, maxInd] = max(dumy(:));
175                        [Irow, Icol] = ind2sub(size(dumy),maxInd);
176                        pl1Out(i,k,ker1Iter) = maxVal;
177                        pl1OutIndex(1,i,k,ker1Iter)=Irow;
178                        pl1OutIndex(2,i,k,ker1Iter)=Icol;
179                    end
180                end
181            end
182            hiddenOut=pl1Out(:);
183            neuralOut=transpose(hiddenOut)*weightOut-biasOut;
184            sumOut=sum(exp(neuralOut));
185            neuralOutAct= exp(neuralOut)*1/(sumOut);
186            d=trainlbl(testShuf(testIter));
187            [mV, mI]= max(neuralOutAct);
188            if (mI ~= d)
189                testError=testError+1;
190            end
191            dvec=zeros(1,numOfOut);
192            dvec(1,d)=1;
193            CostEpoch(1,epochIter)=-dvec*transpose(log(neuralOutAct));
194
195        end
196        testError*(100/1000)
197        accEpoch(1,epochIter)=testError*(100/1000);
198
199        toc;
200    end
201
202 % Calculating test accuracy
203
204 testError=0;
205 for testIter=1:size(nbWTestData,3)
206     X=nbWTestData(:,:,testIter);
207     conv1Out=zeros(pixelNum-ker1L+1,pixelNum-ker1L+1,ker1Num);
208     for ker1Iter=1:ker1Num
209         conv1Out(:,:,ker1Iter) = conv2(X,ker1(:,:,ker1Iter),'valid');
210     end
211     conv1OutAct=1./(1+exp(-lambda*(conv1Out-biasKer1)));
212     pl1Out= zeros((pixelNum-ker1L+1)/pl1,(pixelNum-ker1L+1)/pl1,ker1Num);
213     pl1OutIndex=zeros(2,(pixelNum-ker1L+1)/pl1,(pixelNum-ker1L+1)/pl1,ker1Num);
214     for ker1Iter=1:ker1Num
215         for i=1:(pixelNum-ker1L+1)/pl1
```

```
216              for k=1:(pixelNum-ker1L+1)/pl1
217                   dumy= conv1OutAct((i-1)*pl1+1:i*pl1  ,(k-1)*pl1+1:k*pl1  ,ker1Iter
                        );
218                   [maxVal, maxInd] = max(dumy(:));
219                   [Irow, Icol] = ind2sub(size(dumy),maxInd);
220                   pl1Out(i,k,ker1Iter) = maxVal;
221                   pl1OutIndex(1,i,k,ker1Iter)=Irow;
222                   pl1OutIndex(2,i,k,ker1Iter)=Icol;
223              end
224          end
225      end
226      hiddenOut=pl1Out(:);
227      neuralOut=transpose(hiddenOut)*weightOut-biasOut;
228      sumOut=sum(exp(neuralOut));
229      neuralOutAct= exp(neuralOut)*1/(sumOut);
230      d=testlbl(testIter);
231      [mV, mI]= max(neuralOutAct);
232      %confusion matrix
233      confusiontest(d,mI)=confusiontest(d,mI)+1;
234      if (mI ~= d)
235          testError=testError+1;
236      end
237
238  end
239  testerrorfinal = testError*(100/size(nbWTestData,3));
240  % plotting   curves
241  plot(CostEpoch)
242  xlabel 'Epochs';
243  ylabel 'Cross Entropy Cost';
244
245  figure
246  accuracyepochs=100-accEpoch;
247  plot(accuracyepochs)
248  xlabel 'Epochs';
249  ylabel 'Validation Accuracy ';
250
251  figure
252  for i=1:size(ker1,3)
253      subplot(4,4,i)
254      imagesc(ker1(:,:,i));
255  end
256
257
258
259  imagesc(confusiontest)

 1  clear;
 2  close all;
 3  clc;
 4  %loading data and normalization
 5  load('final_project_dataset.mat');
 6  pixelNum=32;
 7  for i=1:size(testdat,1)
 8      aa=testdat(i,:);
 9      testdatShaped(:,:,:,i)=reshape(aa, [pixelNum pixelNum 3]);
10  end
11  for i=1:size(traindat,1)
12      aa=traindat(i,:);
13      traindatShaped(:,:,:,i)=reshape(aa, [pixelNum pixelNum 3]);
14  end
15  % preprocess test data for w-b and normalize
```

```matlab
16  bWTestData= zeros(pixelNum,pixelNum,size(testdatShaped,4));
17  nbWTestData= zeros(pixelNum,pixelNum,size(testdatShaped,4));
18  for i=1:size(testdatShaped,4)
19      bWTestData(:,:,i)= 0.2126*testdatShaped(:,:,1,i) + 0.7152*testdatShaped
            (:,:,2,i) + 0.0722*testdatShaped(:,:,3,i);
20      nbWTestData(:,:,i)=bWTestData(:,:,i)-ones(pixelNum,pixelNum)*sum(sum(
            bWTestData(:,:,i)))*1/(pixelNum*pixelNum);
21  end
22  stdData=std(nbWTestData(:));
23  nbWTestData((nbWTestData > 3*stdData))=3*stdData;
24  nbWTestData((nbWTestData < -3*stdData))=-3*stdData;
25  nbWTestData=nbWTestData*(4/(30*stdData))+0.5;
26  % preprocess train data for w-b and normalize
27  bWTrainData= zeros(pixelNum,pixelNum,size(traindatShaped,4));
28  nbWTrainData= zeros(pixelNum,pixelNum,size(traindatShaped,4));
29  for i=1:size(traindatShaped,4)
30      bWTrainData(:,:,i)= 0.2126*traindatShaped(:,:,1,i) + 0.7152*traindatShaped
            (:,:,2,i) + 0.0722*traindatShaped(:,:,3,i);
31      nbWTrainData(:,:,i)=bWTrainData(:,:,i)-ones(pixelNum,pixelNum)*sum(sum(
            bWTrainData(:,:,i)))*1/(pixelNum*pixelNum);
32  end
33  stdData2=std(nbWTrainData(:));
34  nbWTrainData((nbWTrainData > 3*stdData2))=3*stdData2;
35  nbWTrainData((nbWTrainData < -3*stdData2))=-3*stdData2;
36  nbWTrainData=nbWTrainData*(4/(30*stdData2))+0.5;
37  %modify train and test labels by increasing 1
38  trainlbl = trainlbl +1;
39  testlbl = testlbl + 1;
40  %save processed data to call in CEcost
41  save('dataProcessed.mat','nbWTestData','nbWTrainData', 'testlbl','trainlbl');
42  %initialization of some parameters
43  numOfHid=40;
44  numOfOut=10;
45  lamda=1;
46  %uniform random initial of weights
47  weightConsHid= sqrt(6/(pixelNum*pixelNum+numOfOut));
48  weightHid = -1*weightConsHid+ (weightConsHid*2)*rand(pixelNum,pixelNum,numOfHid)
        ;
49  biasHid=-1*weightConsHid+ (weightConsHid*2)*rand(numOfHid,1);
50  weightConsOut= sqrt(6/(numOfHid));
51  weightOut= -1*weightConsOut+ (weightConsOut*2)*rand(numOfOut,numOfHid);
52  biasOut=-1*weightConsOut+ (weightConsOut*2)*rand(numOfOut,1); %reshaping all
        weights for function
53  W=[reshape(weightHid,[pixelNum*pixelNum*numOfHid 1]);biasHid;reshape(weightOut,[
        numOfHid*numOfOut 1]); biasOut];
54  conjIter=1;
55  for i=1:conjIter
56      %call of fmingrad
57      [W, fX, iterations] = fmincg('CEcost',W,optimset('MaxIter',50));
58  end
59  %reshaping weights
60  weightHid= reshape(W(1:pixelNum*pixelNum*numOfHid), [pixelNum pixelNum numOfHid
        ]);
61  biasHid=W(pixelNum*pixelNum*numOfHid+1:pixelNum*pixelNum*numOfHid+numOfHid);
62  weightOut= reshape(W(pixelNum*pixelNum*numOfHid+numOfHid+1:pixelNum*pixelNum*
        numOfHid +numOfHid +numOfHid*numOfOut), [numOfOut numOfHid]);
63  biasOut=W(pixelNum*pixelNum*numOfHid +numOfHid +numOfHid*numOfOut+1:length(W));
64  %ploting weights as image for all hidden neurons
65  figure(1);
66  for hiddenIter=1:numOfHid
67      subplot(8,5,hiddenIter);
```

```matlab
68        imagesc(weightHid(:,:,hiddenIter));
69    end
70  %cost and confusion matrix calculation for train and test data
71  Jtrain=0;
72  Ctrain=zeros(10,10);
73  Ctest=zeros(10,10);
74  for trainIter=1:length(trainlbl)
75      X=nbWTrainData(:,:,trainIter);
76      hiddenOut=zeros(numOfHid,1);
77      for hiddenIter=1:numOfHid
78          hiddenOut(hiddenIter,1)=sum(sum(X.*weightHid(:,:,hiddenIter)))-biasHid(
                hiddenIter,1);
79      end
80      hiddenOutAct=1./(1+exp(-1*hiddenOut));
81      neuralOut=weightOut*hiddenOutAct-biasOut;
82      neuralOutAct=1./(1+exp(-1*neuralOut));
83      [m ind] = max(neuralOutAct);
84      Ctrain(trainlbl(trainIter),ind)=Ctrain(trainlbl(trainIter),ind)+1;
85      d=zeros(numOfOut,1);
86      d(trainlbl(trainIter))=1;
87      Jtrain=Jtrain+transpose(-1*d)*log(neuralOutAct)-transpose(1-d)*(log(1-
            neuralOutAct));
88  end
89  Jtest=0;
90  for testIter=1:length(testlbl)
91      X=nbWTestData(:,:,testIter);
92      hiddenOut=zeros(numOfHid,1);
93      for hiddenIter=1:numOfHid
94          hiddenOut(hiddenIter,1)=sum(sum(X.*weightHid(:,:,hiddenIter)))-biasHid(
                hiddenIter,1);
95      end
96      hiddenOutAct=1./(1+exp(-1*hiddenOut));
97      neuralOut=weightOut*hiddenOutAct-biasOut;
98      neuralOutAct=1./(1+exp(-1*neuralOut));
99      [m ind] = max(neuralOutAct);
100     Ctest(testlbl(testIter),ind)=Ctest(testlbl(testIter),ind)+1;
101     d=zeros(numOfOut,1);
102     d(testlbl(testIter))=1;
103     Jtest=Jtest+transpose(-1*d)*log(neuralOutAct)-transpose(1-d)*(log(1-
            neuralOutAct));
104 end
105 %plotting confusion matrix and normalizing cost results
106 Jtrain=Jtrain/length(trainlbl);
107 Jtest=Jtest/length(testlbl);
108 figure(2);
109 imagesc(Ctest, [0 max(max(Ctest))]);
110 figure(3);
111 imagesc(Ctrain, [0 max(max(Ctrain))]);
112 disp ('Cross entropy cost for test data: '); disp(num2str(Jtest));
113 disp ('Cross entropy cost for train data: '); disp(num2str(Jtrain));

1   function [JOut,JGradOut] = CEcost(Weights)
2   %initialize some params
3       aa=Weights;
4       pixelNum=32;
5       numOfHid=40;
6       numOfOut=10;
7       lamda=1;
8       load('dataProcessed.mat');
9       %taking weights for use in iteration
10      weightHid= reshape(Weights(1:pixelNum*pixelNum*numOfHid), [pixelNum pixelNum
```

```matlab
            numOfHid]);
11      biasHid=Weights(pixelNum*pixelNum*numOfHid+1:pixelNum*pixelNum*numOfHid+
            numOfHid);
12      weightOut= reshape(Weights(pixelNum*pixelNum*numOfHid+numOfHid+1:pixelNum*
            pixelNum*numOfHid +numOfHid +numOfHid*numOfOut), [numOfOut numOfHid]);
13      biasOut=Weights(pixelNum*pixelNum*numOfHid +numOfHid +numOfHid*numOfOut+1:
            length(Weights));
14      picSequence = randperm(length(trainlbl));
15      J=0; %initial cost before iteration
16      storeHiddenDer= zeros(pixelNum,pixelNum,numOfHid);
17      storeHidBiasDer=zeros(numOfHid,1);
18      storeOutDer=zeros(numOfOut,numOfHid);
19      storeOutBiasDer=zeros(numOfOut,1);
20      for batchIter=1:length(picSequence) %batch iter
21          X=nbWTrainData(:,:,picSequence(batchIter)); %shuffling data
22          %output and gradient calculations
23          hiddenOut=zeros(numOfHid,1);
24          for hiddenIter=1:numOfHid
25              hiddenOut(hiddenIter,1)=sum(sum(X.*weightHid(:,:,hiddenIter)))-
                    biasHid(hiddenIter,1);
26          end
27          hiddenOutAct=1./(1+exp(-1*hiddenOut));
28          neuralOut=weightOut*hiddenOutAct-biasOut;
29          neuralOutAct=1./(1+exp(-1*neuralOut));
30          d=zeros(numOfOut,1);
31          d(trainlbl(picSequence(batchIter)))=1;
32
33          derivativeFOut = neuralOutAct.*(1-neuralOutAct);
34          gradOut=((-1*d).*(1./neuralOutAct)+(1-d).*(1./(1-neuralOutAct))).*
                derivativeFOut;
35          %gradOut= neuralOutAct - d;
36          derivativeHidden= hiddenOutAct.*(1-hiddenOutAct);
37          gradHidden= derivativeHidden.*transpose((transpose(gradOut)*weightOut));
38          %summing cost and derivatives
39          J=J+transpose(-1*d)*log(neuralOutAct)-transpose(1-d)*(log(1-neuralOutAct
                ));
40          %J=J+transpose(-1*d)*log(neuralOutAct);
41          storeOutDer= storeOutDer+ gradOut*transpose(hiddenOutAct);
42          for hiddenIter=1:numOfHid
43              storeHiddenDer(:,:,hiddenIter)=storeHiddenDer(:,:,hiddenIter)+
                    gradHidden(hiddenIter,1)*X;
44          end
45          storeHidBiasDer= storeHidBiasDer -1*gradHidden;
46          storeOutBiasDer=storeOutBiasDer+-1*gradOut;
47      end
48      %normalizing
49      storeOutDer=storeOutDer*(1/length(picSequence));
50      storeHiddenDer=storeHiddenDer*(1/length(picSequence));
51      storeHidBiasDer=storeHidBiasDer*(1/length(picSequence));
52      storeOutBiasDer=storeOutBiasDer*(1/length(picSequence));
53      JOut=J*(1/length(picSequence));
54      JGradOut=[reshape(storeHiddenDer,[pixelNum*pixelNum*numOfHid 1]);
            storeHidBiasDer;reshape(storeOutDer,[numOfHid*numOfOut 1]);
            storeOutBiasDer];
55  end


1  function [X, fX, i] = fmincg(f, X, options, P1, P2, P3, P4, P5)
2  % Minimize a continuous differentialble multivariate function.
3  % Starting point is given by "X" (D by 1), and the function named in the string
       "f" must
4  % return a function value and a vector of partial derivatives. The Polack-
```

```matlab
 5  % Ribiere flavour of conjugate gradients is used to compute search directions,
 6  % and a line search using quadratic and cubic polynomial approximations and the
 7  % Wolfe-Powell stopping criteria is used together with the slope ratio method
 8  % for guessing initial step sizes.
 9  %
10  % If the function terminates within a few iterations, it could be an indication
        that the function value
11  % and derivatives are not consistent (ie, there may be a bug in the
12  % implementation of your "f" function). The function returns the found
13  % solution "X", a vector of function values "fX" indicating the progress made
14  % and "i" the number of iterations (line searches or function evaluations,
15  % depending on the sign of "length") used.
16  %
17  % Usage: [X, fX, i] = fmincg(f, X, options, P1, P2, P3, P4, P5)
18  %
19  % (C) Copyright 1999, 2000 & 2001, Carl Edward Rasmussen
20  %
21  % Permission is granted for anyone to copy, use, or modify these
22  % programs and accompanying documents for purposes of research or
23  % education, provided this copyright notice is retained, and note is
24  % made of any changes that have been made.

26  % Read options
27  if exist('options', 'var') && ~isempty(options) && isfield(options, 'MaxIter')
28      length = options.MaxIter;
29  else
30      length = 100;
31  end


34  RHO = 0.01;                            % a bunch of constants for line searches
35  SIG = 0.5;       % RHO and SIG are the constants in the Wolfe-Powell conditions
36  INT = 0.1;    % don't reevaluate within 0.1 of the limit of the current bracket
37  EXT = 3.0;                    % extrapolate maximum 3 times the current bracket
38  MAX = 20;                         % max 20 function evaluations per line search
39  RATIO = 100;                                   % maximum allowed slope ratio

41  argstr = ['feval(f, X'];                       % compose string used to call
        function
42  for i = 1:(nargin - 3)
43    argstr = [argstr, ',P', int2str(i)];
44  end
45  argstr = [argstr, ')'];

47  if max(size(length)) == 2, red=length(2); length=length(1); else red=1; end
48  S=['Iteration '];

50  i = 0;                                        % zero the run length counter
51  ls_failed = 0;                             % no previous line search has failed
52  fX = [];
53  [f1 df1] = eval(argstr);                      % get function value and gradient
54  i = i + (length<0);                                        % count epochs?!
55  s = -df1;                                     % search direction is steepest
56  d1 = -s'*s;                                            % this is the slope
57  z1 = red/(1-d1);                           % initial step is red/(|s|+1)

59  while i < abs(length)                                      % while not finished
60    i = i + (length>0);                                    % count iterations?!

62    X0 = X; f0 = f1; df0 = df1;               % make a copy of current values
63    X = X + z1*s;                                           % begin line search
```

17

```matlab
64      [f2 df2] = eval(argstr);
65      i = i + (length<0);                                      % count epochs?!
66      d2 = df2'*s;
67      f3 = f1; d3 = d1; z3 = -z1;              % initialize point 3 equal to point 1
68      if length>0, M = MAX; else M = min(MAX, -length-i); end
69      success = 0; limit = -1;                      % initialize quanteties
70      while 1
71        while ((f2 > f1+z1*RHO*d1) | (d2 > -SIG*d1)) & (M > 0)
72          limit = z1;                                         % tighten the bracket
73          if f2 > f1
74            z2 = z3 - (0.5*d3*z3*z3)/(d3*z3+f2-f3);            % quadratic fit
75          else
76            A = 6*(f2-f3)/z3+3*(d2+d3);                        % cubic fit
77            B = 3*(f3-f2)-z3*(d3+2*d2);
78            z2 = (sqrt(B*B-A*d2*z3*z3)-B)/A;        % numerical error possible - ok!
79          end
80          if isnan(z2) | isinf(z2)
81            z2 = z3/2;                        % if we had a numerical problem then bisect
82          end
83          z2 = max(min(z2, INT*z3),(1-INT)*z3);  % don't accept too close to limits
84          z1 = z1 + z2;                                       % update the step
85          X = X + z2*s;
86          [f2 df2] = eval(argstr);
87          M = M - 1; i = i + (length<0);                      % count epochs?!
88          d2 = df2'*s;
89          z3 = z3-z2;                       % z3 is now relative to the location of z2
90        end
91        if f2 > f1+z1*RHO*d1 | d2 > -SIG*d1
92          break;                                              % this is a failure
93        elseif d2 > SIG*d1
94          success = 1; break;                                 % success
95        elseif M == 0
96          break;                                              % failure
97        end
98        A = 6*(f2-f3)/z3+3*(d2+d3);                        % make cubic extrapolation
99        B = 3*(f3-f2)-z3*(d3+2*d2);
100       z2 = -d2*z3*z3/(B+sqrt(B*B-A*d2*z3*z3));        % num. error possible - ok!
101       if ~isreal(z2) | isnan(z2) | isinf(z2) | z2 < 0   % num prob or wrong sign?
102         if limit < -0.5                              % if we have no upper limit
103           z2 = z1 * (EXT-1);                   % the extrapolate the maximum amount
104         else
105           z2 = (limit-z1)/2;                                % otherwise bisect
106         end
107       elseif (limit > -0.5) & (z2+z1 > limit)        % extraplation beyond max?
108         z2 = (limit-z1)/2;                                  % bisect
109       elseif (limit < -0.5) & (z2+z1 > z1*EXT)       % extrapolation beyond limit
110         z2 = z1*(EXT-1.0);                             % set to extrapolation limit
111       elseif z2 < -z3*INT
112         z2 = -z3*INT;
113       elseif (limit > -0.5) & (z2 < (limit-z1)*(1.0-INT))   % too close to limit?
114         z2 = (limit-z1)*(1.0-INT);
115       end
116       f3 = f2; d3 = d2; z3 = -z2;              % set point 3 equal to point 2
117       z1 = z1 + z2; X = X + z2*s;                   % update current estimates
118       [f2 df2] = eval(argstr);
119       M = M - 1; i = i + (length<0);                        % count epochs?!
120       d2 = df2'*s;
121     end                                              % end of line search
122
123     if success                                        % if line search succeeded
124       f1 = f2; fX = [fX' f1]';
```

```matlab
125          fprintf('%s %4i | Cost: %4.6e\r', S, i, f1);
126          s = (df2'*df2-df1'*df2)/(df1'*df1)*s - df2;        % Polack-Ribiere direction
127          tmp = df1; df1 = df2; df2 = tmp;                           % swap derivatives
128          d2 = df1'*s;
129          if d2 > 0                                        % new slope must be negative
130            s = -df1;                                 % otherwise use steepest direction
131            d2 = -s'*s;
132          end
133          z1 = z1 * min(RATIO, d1/(d2-realmin));          % slope ratio but max RATIO
134          d1 = d2;
135          ls_failed = 0;                                   % this line search did not fail
136        else
137          X = X0; f1 = f0; df1 = df0;  % restore point from before failed line search
138          if ls_failed | i > abs(length)            % line search failed twice in a row
139            break;                                  % or we ran out of time, so we give up
140          end
141          tmp = df1; df1 = df2; df2 = tmp;                           % swap derivatives
142          s = -df1;                                                 % try steepest
143          d1 = -s'*s;
144          z1 = 1/(1-d1);
145          ls_failed = 1;                                   % this line search failed
146        end
147        if exist('OCTAVE_VERSION')
148          fflush(stdout);
149        end
150      end
151      fprintf('\n');
152    end
```