# KTH Royal Institute of Technology

## DD2424
### Deep Learning in Data Science

### Year 2017 - 2018

# Assignment 2

Picó Oristrell Sandra                    940531−2308    sandrapo@kth.se

*Professors :*
  Sullivan Josephine

August 5, 2018

# Contents

The main goal of this second assignment is to train and test a two-layer network with multiple outputs to classify images from the CIFAR-10 dataset. As in the first assignment, the network will be trained using mini-batch gradient descent with a cost function that computes the cross-entropy loss and with an L$_2$ regularization term on the weight matrix. The main structure of the code will be similar to assignment 1, but adding more parameters as now we are implementing one extra layer.

In order to explain and show the results obtained in a proper way I will divide the report into different subsections.

At the same time, in the appendix section, I will append all the implemented functions to be able to check the code ,if needed.

# 1 Background

The first section of this assignment will be focused on explaining the background needed for the assignment implementation.

## 1.1 Background 1: Mathematical background.

Following the same structure than in the first assignment, the network that we are going to implement will be a classifier where the output will correspond to a vector of probabilities. Then, the main structure of the network will be as follows:



Figure 1: Main structure of the implemented network.

Then, the predicted class will correspond to the output label with higher probability, i.e the *arg max* of the probability's vector.

In order to compute the output vector we will implement the following two-layer neural network:

$$s_1 = W_1 x + b_1 \tag{1}$$

$$h = max(0, s_1) \tag{2}$$

$$s = W_2 h + b_2 \tag{3}$$

$$\mathbf{p} = softmax(s) = \frac{exp(s)}{1^T exp(s)} \tag{4}$$

This means that during the training phase we will learn the parameters $W_1$, $W_2$, $b_1$ and $b_2$.

Then, defining the model as: $\Theta = [W_1, W_2, b_1, b_2]$, we will solve the following optimization problem to be able to set all the mentioned trainable parameters:

$$\Theta^* = argmin_\Theta J(D, \lambda, \Theta) \tag{5}$$

Where *J* is a cost function that minimize the cross-entropy plus a regularization term on *W*.

$$J(D, \lambda, \Theta) = \frac{1}{|D|} \sum_{x,y \in D} \log(y^T p) + \lambda \sum_{l=1}^{2} \sum_{i,j} W_{l_{i,j}}^2 \tag{6}$$

The problem will be solve via mini-batch gradient descent with momentum.

For the mini-batch implementation we will begin with sensible random initialization for W and b, and then estimate the parameters for k = 1,2:

$$W_k^{(t+1)} = W_k^{(t)} - \eta \frac{\partial J((t+1), \lambda, \Theta)}{\partial W_k} \tag{7}$$

$$b_k^{(t+1)} = b_k^{(t)} - \eta \frac{\partial J((t+1), \lambda, \Theta)}{\partial b_k} \tag{8}$$

## 1.2 Background 2: Speeding up training, add momentum.

The vanilla version of the mini-batch gradient descent is really slow when using a sensible learning rate and with the size of data that we are going to use for this implementation. For this reason, we will need to speed up the training.

One of the possible ways to do it is adding a momentum term in the update step.

To achieve it, you first need to initialize a momentum vector (matrix) $v_0$ for each parameter of the network. The vector will be initialized to zero and will have the same dimension as the correspondence parameter. Once initialized, per each time step t:

$$v_t = p * v_{t-1} + \eta * \frac{\partial J}{\partial \theta} \tag{9}$$

$$\theta_t = \theta_{t-1} - v_t \tag{10}$$

where $\eta$ is the learning rate, p $\in$ [0, 1] and $\theta$ is a generic placeholder to represent one of the parameters of the model.

In this assignment we will not implement an adaptive learning rate but we will decay the learning rate by some factor after each epoch.

## 1.3 Data structure

As in the previous assignment, the data that we are going to use is from the CIFAR-10 dataset. The dataset has a size of 10000 x 3072 following the next structure:
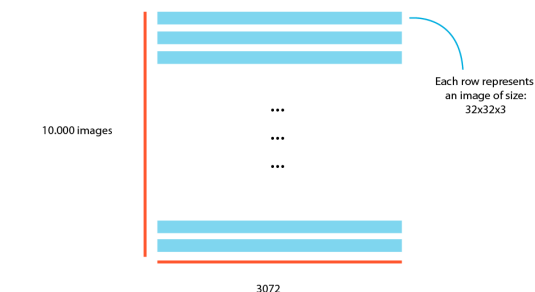


Figure 2: Data structure of CIFAR-10 dataset

# 2 Exercises

## 2.1 Exercise 1: Read in the data and initialize the parameters of the network.

To begin with this assignment we will split the data as follows:

- **Training data:** data_batch_1.mat

- **Validation data:** data_batch_2.mat

- **Test data:** test_batch.mat

To achieve it, we are going to use the function LoadBatch, already implemented for Assignment 1. Nevertheless, we will also need to pre-process the input data to have zero mean. We need to compute the mean for training data and then subtract this mean to the validation and test sets.

$$mean\_X = mean(X, 2); \tag{11}$$

$$X = X - repmat(mean\_X, [1, size(X, 2)]); \tag{12}$$

Next step will be initialize and set up all the parameters of the network. In this 2 layer neural network we will have m=50 nodes in the hidden layer. As the two weight matrices and bias vectors will have different dimensions, we will use cell arrays to store them.

The weight matrices will be initialized with a Gaussian distribution with 0 mean and 0.001 of standard deviation. On the other hand, bias vectors will be initialized to zero.

Then, the first assignment will be just the initialization of the network. The implementation can be found in the appendix section.

## 2.2 Exercise 2: Compute the gradients for the network parameters

After initialize and define all the network parameters, we will start writing the functions to compute the gradients. To do so, we will re-use the functions already implemented in Assignment 1.

Specifically, we are re-writing the function: ComputeGradients and transforming it to be able to compute a two-layer network using cell-arrays. EvaluateClassifier and ComputeCost are also modified.

Once the two mentioned functions are already implemented, it is important to check if the gradients are calculated properly. To be able to compare it, we will use the same methodology as in the first assignment; compare the gradients achieved with the numerical ones.

The achieved values are as follows:

|  | Values |
|---|---|
| $g\_W_1$ | 1.9894e-05 |
| $g\_b_1$ | 7.1271e-04 |
| $g\_W_2$ | 5.0976e-08 |
| $g\_b_2$ | 6.2052e-10 |

As the values achieved are small enough, we will consider that the gradients are properly implemented.

The last experiment that we are going to implement to be sure that gradient computations and mini-batch gradient descent algorithm are okay is to train the network and check if it is possible to over-fit to training data and get a small loss after enough number of epochs.

Specifically, we are going to use the following configurations:

- Small amount of training data: 100 examples.

- Regularization term turned off: lambda = 0.

- Enough number of epochs: 300.

- Reasonable learning rate: 0.01

As can be seen in the following plot, when we train the two layer network using the previous parameters,we are able to over-fit into the training data:
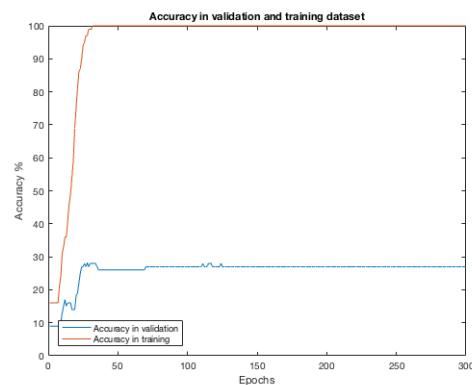


Figure 3: Training and testing accuracy with 100 examples.

We can conclude that the training is working because we are able to over-fit with only 29 epochs.

## 2.3    Exercise 3: Add momentum to your update step

During this third exercise, we have implemented the momentum vectors into the MiniBatch function to speed up the learning.

   We will check its behaviour using different values of rho = [0.5,0.9,0.99] and fixing the learning rate into 0.001. The number of epochs will be set up to 10.

   At the same time, to fully check the importance of it and the speed up achieved using momentum, I will do the following experiments using all the training data of the data_batch_1.mat dataset.

   Here we can see the comparison between using momentum or without using it.

- **With momentum:**


   - **rho = 0.5**



Figure 4: Cost with 10 epochs, rho = 0.5
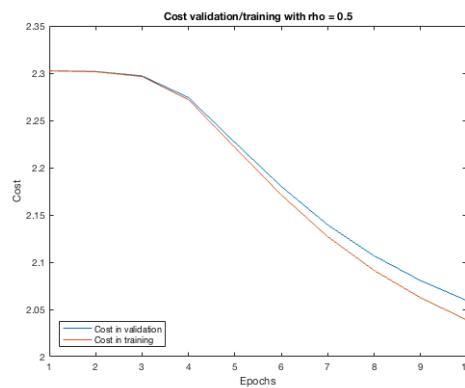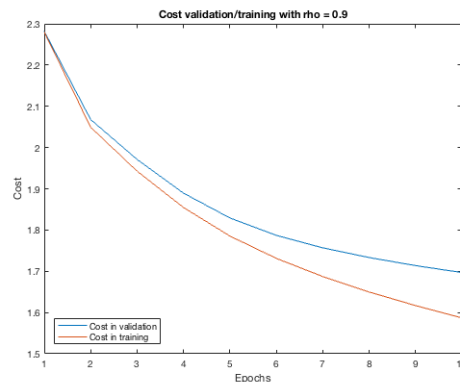

   - **rho = 0.9**



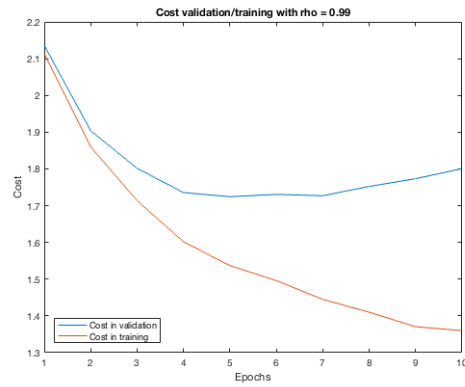Figure 5: Cost with 10 epochs, rho = 0.9

– **rho = 0.99**



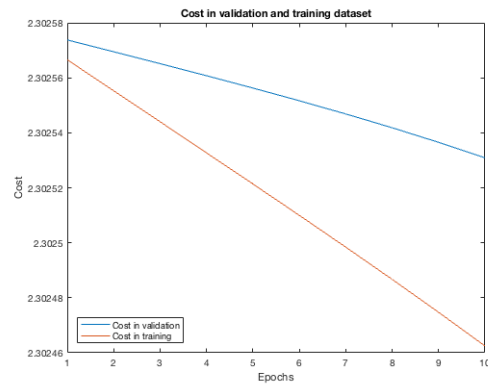Figure 6: Cost with 10 epochs, rho = 0.99

- **Without momentum:**



Figure 7: Cost with 10 epochs with no momentum.

As can be concluded from the previous images, using momentum helps to converge faster than not using it. Also, with momentum and rho = 0.90 we achieve the best performance.

7

## 2.4 Exercise 4: Training your network

After doing the necessary experiments to verify that the momentum vectors and the gradients were well implemented, we will proceed to train the network correctly.

The first part of this assignment will be focused on finding a reasonable range of values for the learning rate parameter.

- **Find reasonable range of values for learning rate.**

  During this exercise, we are going to use all the samples from the training data; i.e not only using 100 examples of the data.

  As specified in the assignment description, we will use $rho = 0.9$, a small value for the regularization term and also only 10 epochs.

  To obtain a reasonable value of the learning rate, we will test all the next values:

  eta_values $= [0.1, 0.01, 0.001, 0.0001, 1e-05, 1e-06, 1e-7]$;

  The obtained results can be seen in the next figure:



Figure 8: Training loss for different learning rate values.

As can be seen in the figure previously attached, the value of eta greatly affects the loss obtained while training in the network. In the mentioned plot, we can conclude that when the value of the learning rate is too small; as in the cases of 1e-5,1e-6 and 1e-7, we practically do not learn; the loss is exactly the same during the first 10 epochs of the network.

Therefore, for the next experiment, we will consider that the first range of values of learning rate must be from 0.001 to 0.1.

- **Coarse random search.**

  In this step we are going to find the best combinations for lambda and eta values, fixing rho = 0.9.

  The main intention will be performing training and then checking the learned network's best performance on the validation set; via the accuracy on it.

  The first search will be focused over the feasible learning-rates that we identified in the previous experiment with a search over a very broad range of values for lambda.

  Specifically, the range of values that we used for this first search are the next ones:

  | Network parameters | |
  |---|---|
  | rho | 0.9 |
  | epochs | 10 |
  | eta range values | 0.001 to 0.1 |
  | lambda range values | 1e-7 to 0.1 |
  | number of pairs | 70 |

  Using the mentioned configurations, the 3 best results obtained using this first search are:

  | Eta | Lambda | Validation accuracy |
  |---|---|---|
  | 0.011778 | 0.000323 | 43.54% |
  | 0.013877 | 0.000165 | 43.28% |
  | 0.018709 | 0.007597 | 43.20% |

- **Coarse-to-fine random search to set lambda and eta.**

  During this search, we used a small range of values for lambda, to achieve better and accurate results. Specifically, we have been based on the previous best results. Then, during the coarse-to-fine random search, the parameters and the results obtained are as follows:

  | Network parameters | |
  |---|---|
  | rho | 0.9 |
  | epochs | 10 |
  | eta range values | 0.01 to 0.07 |
  | lambda range values | 1e-6 to 0.005 |
  | number of pairs | 70 |

  The results obtained in the coarse-to-fine random search are the next ones:

  | Eta | Lambda | Validation accuracy |
  |---|---|---|
  | 0.014001 | 0.003347 | 44.00% |
  | 0.019811 | 0.003421 | 44.00% |
  | 0.018525 | 0.000001 | 43.82% |

  We can state that during this second search, we are improving the achieved results.

- **Best hyper-parameters**

  Finally, using the best hyper-parameters combination found previously, we trained the model for longer and we checked the performance obtained on the test set.

  During the training, we also break out early once the training cost is ever > 3 * original training cost.

Also, we are using only three of the training sets as training data instead of using all of them since exceeds Matlab maximum array dimensions.

Using the best hyper-parameters selection:

| rho | lambda | eta |
|-----|--------|-----|
| 0.9 | 0.003347 | 0.014001 |

As a result of training the network during 30 epochs and with the previous specified hyper-parameters selection, we obtain 49.46% of accuracy in testing and 55.64% of accuracy in training data.

On the other hand,the validation and training loss achieved per each epoch, can be state in the figure:



Figure 9: Final training and validation loss for 30 epochs.

# 3  Conclusions

During this assignment I learned about the implementation of a two-layer neural network, but above all the importance of correctly selecting the network hyper-parameters to be able to train and get better performance.

It should also be noted that without momentum, the experiments would have been much slower.

Therefore, with all the results obtained during the different experiments, we can conclude that the implementation of the 2-Layer neural network is properly implemented and verified.

# Appendices

## A  Assignment 2 implementation:

```matlab
1  %Author: Sandra Pic  Oristrell
2  %Data: 13 of August 2018.
3  %Assignment 2: 2 layer neural network.
4
5  %% Exercise 1: Read in the data and initialize the parameters of the network.
6  %Load the data using the function implemented in assignment 1.
7
8  [X_training, Y_training, y_training] = LoadBatch('./data_batch_1.mat');
9  [X_validation, Y_validation, y_validation] = LoadBatch('./data_batch_2.mat');
10 [X_test, Y_test, y_test] = LoadBatch('./test_batch.mat');
11
12 %Apply more pre-processing to the raw input data.
13 %Transform it to have zero mean.
14 mean_X = mean(X_training, 2);
15 X_training = X_training - repmat(mean_X, [1, size(X_training, 2)]);
16 X_validation  = X_validation  - repmat(mean_X, [1, size(X_validation,  2)]);
17 X_test  = X_test  - repmat(mean_X, [1, size(X_test,  2)]);
18
19 %Variables
20 d = size(X_training,1);
21 N = size(X_training,2);
22 K = size(Y_training,1);
23
24 %Data structure for the parameters of the network and initialize the
25 %values.
26
27 %50 nodes in the hidden layer
28 m = 50;
29 lambda = 0;
30 %Initialize the network:
31 [W,b] = Initialize_Network(m,d,K);
32
33 %% Exercise 2: Compute the gradients for the network parameters.
34
35 %Re-write or update gradients functions from assignment 1.
36 %% 2.1 Check that the gradients are implemented properly.
37
38 batch_size = 100;
39
40 [P,h] = EvaluateClassifier(X_training(:,1:batch_size),W,b);
41 [grad_W, grad_b] = ...
     ComputeGradients(X_training(:,1:batch_size),Y_training(:,1:batch_size),P,h,W,b,lambda);
42 [ngrad_W, ngrad_b] = ...
     ComputeGradsNumSlow(X_training(:,1:batch_size),Y_training(:,1:batch_size),W,b,lambda,1e-5);
43
44 %Check the comparison between the numerical and the computed gradients:
45 eps = 1e-10;
46 %Layer 1:
47 gradient_b1_comparison = sum(abs(ngrad_b{1} - grad_b{1})/max(eps, sum(abs(ngrad_b{1}) ...
     + abs(grad_b{1})))));
48 gradient_W1_comparison = sum(sum(abs(ngrad_W{1} - grad_W{1})/max(eps, ...
     sum(sum(abs(ngrad_W{1}) + abs(grad_W{1}))))));
49 %Layer 2:
50 gradient_b2_comparison = sum(abs(ngrad_b{2} - grad_b{2})/max(eps, sum(abs(ngrad_b{2}) ...
     + abs(grad_b{2})))));
51 gradient_W2_comparison = sum(sum(abs(ngrad_W{2} - grad_W{2})/max(eps, ...
     sum(sum(abs(ngrad_W{2}) + abs(grad_W{2}))))));
```

```matlab
52
53   %Check that the gradients have small number.
54   fprintf("Results for the calculated gradients:");
55   fprintf("Layer 1:");
56   fprintf("W1: %f",gradient_W1_comparison);
57   fprintf("b1: %f",gradient_b1_comparison);
58   fprintf("Layer 2:");
59   fprintf("W2: %f",gradient_W2_comparison);
60   fprintf("b2: %f",gradient_b2_comparison);
61
62   %% 2.2 Try if you can overfit in training data with 200 epochs and reasonable learning ...
         rate.
63
64   GDparams.eta = 0.01;
65   GDparams.n_batch = 1;
66   GDparams.n_epochs = 300;
67   lambda = 0;
68
69   %Save all the information per each epoch.
70   cost_training_list = zeros(1, GDparams.n_epochs);
71   cost_validation_list = zeros(1, GDparams.n_epochs);
72   accuracy_training_list = zeros(1, GDparams.n_epochs);
73   accuracy_validation_list = zeros(1, GDparams.n_epochs);
74   epochs_list = zeros(1, GDparams.n_epochs);
75
76   %Try with only 100 examples.
77   examples = 100;
78   X_train = X_training(:,1:examples);
79   Y_train = Y_training(:,1:examples);
80   y_train = y_training(:,1:examples);
81
82   X_val = X_validation(:,1:examples);
83   Y_val = Y_validation(:,1:examples);
84   y_val = y_validation(:,1:examples);
85
86   for i = 1:GDparams.n_epochs
87       [W,b] = MiniBatchGD(X_train, Y_train,y_train,GDparams,W,b,lambda);
88       fprintf("Epoch: %d\n",i);
89       epochs_list(i) = i;
90       cost_training_list(i) = ComputeCost(X_train, Y_train, W, b, lambda);
91       accuracy_training_list(i) = ComputeAccuracy(X_train,y_train,W,b);
92       accuracy_validation_list(i) = ComputeAccuracy(X_val,y_val,W,b);
93       cost_validation_list(i) = ComputeCost(X_val, Y_val, W, b,lambda);
94   end
95
96   %Plot cost in training and validation per epoch.
97   title_text = "Cost in validation and training dataset";
98   PlotCost(epochs_list, cost_validation_list, cost_training_list,title_text);
99
100  %Plot accuracy in training and validation per epoch.
101  title_text = "Accuracy in validation and training dataset";
102  PlotAccuracy(epochs_list,accuracy_validation_list, accuracy_training_list,title_text);
103
104  %Final accuracy in validation dataset and training dataset:
105  accuracy_final_training = ComputeAccuracy(X_train,y_train,W,b);
106  fprintf("Final training accuracy: %f %", accuracy_final_training);
107
108  accuracy_final_validation = ComputeAccuracy(X_val,y_val,W,b);
109  fprintf("Final validation accuracy: %f %",accuracy_final_validation);
110
111  %%  Exercise 3: Add momentum to your update step.
112  %To help speed up training times you should add momentum terms into your mini-batch ...
         update steps
113
114  GDparams.eta = 0.001;
```

```matlab
115  GDparams.n_batch = 100;
116  GDparams.n_epochs = 10;
117  lambda = 0;
118
119  eta_decay_rate = 0.95;
120
121  %Compare rho values = [0.5,0.9,0.99];
122  rho_values = [0.50,0.9,0.99];
123
124  for rho = rho_values
125      [W,b] = Initialize_Network(m,d,K);
126      GDparams.eta = 0.01;
127      cost_training_list = zeros(1, GDparams.n_epochs);
128      cost_validation_list = zeros(1, GDparams.n_epochs);
129      accuracy_training_list = zeros(1, GDparams.n_epochs);
130      accuracy_validation_list = zeros(1, GDparams.n_epochs);
131      epochs_list = zeros(1, GDparams.n_epochs);
132      for i = 1:GDparams.n_epochs
133          [W,b] = MiniBatchGD_withMomentum(X_training, ...
                   Y_training,y_training,GDparams,W,b,lambda,rho);
134          fprintf("Epoch: %d\n",i);
135          epochs_list(i) = i;
136          cost_training_list(i) = ComputeCost(X_training, Y_training, W, b, lambda);
137          accuracy_training_list(i) = ComputeAccuracy(X_training,y_training,W,b);
138          accuracy_validation_list(i) = ComputeAccuracy(X_validation,y_validation,W,b);
139          cost_validation_list(i) = ComputeCost(X_validation, Y_validation, W, b,lambda);
140          GDparams.eta  = GDparams.eta * eta_decay_rate;
141      end
142      %Plot cost in training and validation per epoch.
143      title_text = "Cost validation/training with rho = " + num2str(rho);
144      PlotCost(epochs_list, cost_validation_list, cost_training_list,title_text);
145      %Plot accuracy in training and validation per epoch.
146      title_text = "Accuracy validation/training with rho = " + num2str(rho);
147      PlotAccuracy(epochs_list,accuracy_validation_list, accuracy_training_list,title_text);
148      %Final accuracy in validation dataset and training dataset:
149      accuracy_final_training = ComputeAccuracy(X_training,y_training,W,b);
150      fprintf("Final training accuracy: %f %", accuracy_final_training);
151      accuracy_final_validation = ComputeAccuracy(X_validation,y_validation,W,b);
152      fprintf("Final validation accuracy: %f %",accuracy_final_validation);
153  end
154
155  %% Exercise 4: Training your network
156  % All the experiments will be using all the examples.
157
158  %%  4.1 Find reasonable values for the learning rate:
159
160
161  %Regularization term to small value:
162  lambda = 0.000001;
163  rho = 0.9;
164  eta_values = [0.1,0.01,0.001,0.0001,0.00001,0.000001,0.0000001];
165  GDparams.n_batch = 100;
166  GDparams.n_epochs = 10;
167
168  final_cost_training_list = {};
169  for eta = eta_values
170      GDparams.eta = eta;
171      cost_training_list = zeros(1, GDparams.n_epochs);
172      [W,b] = Initialize_Network(m,d,K);
173      fprintf("Eta value: %f\n",eta);
174      for i = 1:GDparams.n_epochs
175          fprintf("Epoch: %d\n",i);
176          [W,b] = MiniBatchGD_withMomentum(X_training, ...
                   Y_training,y_training,GDparams,W,b,lambda,rho);
177          cost_training_list(i) = ComputeCost(X_training, Y_training,W,b,lambda);
```

```matlab
178        end
179        final_cost_training_list{end+1} = cost_training_list;
180    end
181    epochs_list = zeros(1,GDparams.n_epochs);
182    for i = 1: GDparams.n_epochs
183        epochs_list(i) = i;
184    end
185    %Plot the cost per each case of eta.
186    PlotCost_Eta_values(epochs_list, eta_values, final_cost_training_list);
187
188    %% 4.2 Coarse-search random.
189
190    GDparams.n_batch = 100;
191    GDparams.n_epochs = 10;
192    decay_rate= 0.95;
193    rho=0.9;
194    n_pairs= 70;
195    %Learning rate range:
196    eta_max = 0.1;
197    eta_min = 0.001;
198    %Lambda range:
199    lambda_max = 0.1;
200    lambda_min = 1e-7;
201
202    validation_accuracy_list = zeros(1,n_pairs);
203    eta_values_list = zeros(1,n_pairs);
204    lambda_values_list = zeros(1,n_pairs);
205
206    for j = 1:n_pairs
207        fprintf("\n Pair number: %d \n ", j);
208        eta_exp = log10(eta_min) + (log10(eta_max) - log10(eta_min))*rand(1, 1);
209        eta = 10^eta_exp;
210        lambda_exp = log10(lambda_min) + (log10(lambda_max) - log10(lambda_min))*rand(1, 1);
211        lambda = 10^lambda_exp;
212
213        GDparams.eta = eta;
214        [W,b] = Initialize_Network(m,d,K);
215
216        for i=1: GDparams.n_epochs
217            fprintf("Epoch: %d\n", i);
218            [W, b] = MiniBatchGD_withMomentum(X_training, Y_training,y_training, GDparams, ...
                    W, b, lambda, rho);
219            GDparams.eta = GDparams.eta*decay_rate;
220        end
221        validation_accuracy_list(j) = ComputeAccuracy(X_validation,y_validation,W,b);
222        eta_values_list(j) = GDparams.eta;
223        lambda_values_list(j) = lambda;
224    end
225
226    %Sort the array:
227    [validation_accuracy_list,validation_indexs] = sort(validation_accuracy_list,'descend');
228    eta_values_list = eta_values_list(validation_indexs);
229    lambda_values_list = lambda_values_list(validation_indexs);
230
231    %Write file:
232    filename = 'exercise_4.2_coarse_search.txt';
233    fid = fopen(filename,'wt');
234    fprintf(fid,"\nAccuracy values:\n");
235    fprintf(fid,'%f\t',validation_accuracy_list);
236    fprintf(fid,"\nEta values:\n");
237    fprintf(fid,'%f\t',eta_values_list);
238    fprintf(fid,"\nLambda values:\n");
239    fprintf(fid,'%f\t',lambda_values_list);
240    fprintf(fid,"\n");
241    fclose(fid);
```

```matlab
242
243  %% 4.2 Coarse-fine search random.
244  %Better values for lambda search.
245
246  GDparams.n_batch = 100;
247  GDparams.n_epochs = 10;
248  decay_rate= 0.95;
249  rho=0.9;
250  n_pairs= 70;
251  %Learning rate range:
252  eta_max = 0.07;
253  eta_min = 0.01;
254  %Lambda range:
255  lambda_max = 0.005;
256  lambda_min = 1e-6;
257
258  validation_accuracy_list = zeros(1,n_pairs);
259  eta_values_list = zeros(1,n_pairs);
260  lambda_values_list = zeros(1,n_pairs);
261
262  for j = 1:n_pairs
263      fprintf("\n Pair number: %d \n ", j);
264      eta_exp = log10(eta_min) + (log10(eta_max) - log10(eta_min))*rand(1, 1);
265      eta = 10^eta_exp;
266      lambda_exp = log10(lambda_min) + (log10(lambda_max) - log10(lambda_min))*rand(1, 1);
267      lambda = 10^lambda_exp;
268
269      GDparams.eta = eta;
270      [W,b] = Initialize_Network(m,d,K);
271
272      for i=1: GDparams.n_epochs
273          fprintf("Epoch: %d\n", i);
274          [W, b] = MiniBatchGD_withMomentum(X_training, Y_training,y_training, GDparams, ...
275              W, b, lambda, rho);
275          GDparams.eta = GDparams.eta*decay_rate;
276      end
277      validation_accuracy_list(j) = ComputeAccuracy(X_validation,y_validation,W,b);
278      eta_values_list(j) = GDparams.eta;
279      lambda_values_list(j) = lambda;
280  end
281
282  %Sort the array:
283  [validation_accuracy_list,validation_indexs] = sort(validation_accuracy_list,'descend');
284  eta_values_list = eta_values_list(validation_indexs);
285  lambda_values_list = lambda_values_list(validation_indexs);
286
287  %Write file:
288  filename = 'exercise_4.2_coarse_fine_search.txt';
289  fid = fopen(filename,'wt');
290  fprintf(fid,"\nAccuracy values:\n");
291  fprintf(fid,'%f\t',validation_accuracy_list);
292  fprintf(fid,"\nEta values:\n");
293  fprintf(fid,'%f\t',eta_values_list);
294  fprintf(fid,"\nLambda values:\n");
295  fprintf(fid,'%f\t',lambda_values_list);
296  fprintf(fid,"\n");
297  fclose(fid);
298
299  %% 4.3 Best hyper-parameters for training the network.
300  %Best combination
301
302  m= 50;
303
304  %Train with all the training data:
305  [X_training1, Y_training1, y_training1] = LoadBatch('./data_batch_1.mat');
```

```matlab
306  [X_training2, Y_training2, y_training2] = LoadBatch('./data_batch_2.mat');
307  [X_training3, Y_training3, y_training3] = LoadBatch('./data_batch_3.mat');
308  [X_training4, Y_training4, y_training4] = LoadBatch('./data_batch_4.mat');
309  [X_training5, Y_training5, y_training5] = LoadBatch('./data_batch_5.mat');
310
311  %Test dataset
312  [X_test, Y_test, y_test] = LoadBatch('./test_batch.mat');
313
314  %except 1000 samples for validation: using the validation dataset
315  %(data_batch_2.mat)
316  X_validation = X_training2(:, 1:1000);
317  Y_validation = Y_training2(:, 1:1000);
318  y_validation = y_training2(:, 1:1000);
319  X_training2 = X_training2(:, 1001:10000);
320  Y_training2 = Y_training2(:, 1001:10000);
321  y_training2 = y_training2(:, 1001:10000);
322
323  %All training data together:
324  X_training = [X_training1, X_training2, X_training3];
325  Y_training = [Y_training1, Y_training2, Y_training3];
326  y_training = [y_training1, y_training2, y_training3];
327
328  %Pre-processing:
329  mean_X = mean(X_training, 2);
330  X_training = X_training - repmat(mean_X, [1, size(X_training, 2)]);
331  X_validation  = X_validation  - repmat(mean_X, [1, size(X_validation,  2)]);
332  X_test  = X_test  - repmat(mean_X, [1, size(X_test,  2)]);
333
334  %Variables
335  d = size(X_training,1);
336  N = size(X_training,2);
337  K = size(Y_training,1);
338
339  GDparams.n_batch = 100;
340  GDparams.n_epochs = 30;
341  decay_rate= 0.95;
342  rho= 0.9;
343  %Best combination:
344  GDparams.eta = 0.014001;
345  lambda = 0.003347;
346
347  %Plot the training and validation cost after each epoch of training
348  cost_training_list = zeros(1, GDparams.n_epochs);
349  cost_validation_list = zeros(1, GDparams.n_epochs);
350  [W,b] = Initialize_Network(m,d,K);
351  %Original training cost
352  original_training_cost = ComputeCost(X_training, Y_training,W,b,lambda);
353  for i=1:GDparams.n_epochs
354      [W, b] = MiniBatchGD_withMomentum(X_training, Y_training,y_training, GDparams, W, ...
             b, lambda, rho);
355      cost_training_list(i) = ComputeCost(X_training, Y_training,W,b,lambda);
356      cost_validation_list(i) = ComputeCost(X_validation, Y_validation,W,b,lambda);
357      GDparams.eta = GDparams.eta*decay_rate;
358      if cost_training_list(i) > 3*original_training_cost
359          fprintf("Cost_training(i) > 3*original_training_cost");
360          i = GDparams.n_epochs;
361      end
362  end
363
364  epochs_list = zeros(1, GDparams.n_epochs);
365  for i=1:GDparams.n_epochs
366      epochs_list(i) = i;
367  end
368
369  %Plot the loss per epochs in validation and in training:
```

16

```matlab
370  title_text = "Validation/Training cost per epoch";
371  PlotCost(epochs_list, cost_validation_list, cost_training_list,title_text)
372
373  %Final accuracy in testing and training:
374  final_accuracy_test = ComputeAccuracy(X_test,y_test,W,b);
375  fprintf("Final accuracy in test: %f % \n", final_accuracy_test);
376  final_accuracy_training = ComputeAccuracy(X_training,y_training,W,b);
377  fprintf("Final accuracy in training: %f % \n", final_accuracy_training);
378
379  %% Functions implementation
380
381
382  function PlotCost_Eta_values(epochs_list, eta_values, final_cost_training_list)
383      figure;
384      for i = 1:length(eta_values)
385          text_label = "eta = " + num2str(eta_values(i));
386          plot(epochs_list,final_cost_training_list{i},'DisplayName',text_label);
387          hold on;
388      end
389      hold off;
390      xlabel('Epochs');
391      ylabel('Cost');
392      legend;
393  end
394
395  %Plot accuracy for validation/test - training dataset per epoch.
396  function PlotAccuracy(epochs, accuracy_validation, accuracy_training,title_text)
397      figure;
398      plot(epochs,accuracy_validation,epochs,accuracy_training);
399      title(title_text);
400      xlabel('Epochs');
401      ylabel('Accuracy %');
402      legend('Accuracy in validation', 'Accuracy in training','Location','southwest');
403  end
404
405  %Plot the loss/cost for validation/test -training dataset per epoch.
406  function PlotCost(epochs, cost_validation, cost_training,title_text)
407      figure;
408      plot(epochs,cost_validation,epochs,cost_training);
409      title(title_text);
410      xlabel('Epochs');
411      ylabel('Cost');
412      legend('Cost in validation', 'Cost in training','Location','southwest');
413  end
414
415  %Function to initialize the values of the Weight matrices.
416  function [W,b] = Initialize_Network(m,d,K)
417      W1 = randn(m,d)*0.001;
418      W2 = randn(K,m)*0.001;
419      b1 = zeros(m, 1);
420      b2 = zeros(K, 1);
421      W = {W1,W2};
422      b = {b1, b2};
423  end
424
425  function [X,Y,y] = LoadBatch(filename)
426      A = load(filename);
427      X = double(A.data');
428      X= X/255;
429      y = double(A.labels') + 1;
430      vec = ind2vec(y);
431      Y = full(vec);        %One-hot representation
432  end
433
434  %Function to compute the cost of the 2-layer network.
```

```matlab
435   function J = ComputeCost(X,Y,W,b,lambda)
436       [P,h] = EvaluateClassifier(X,W,b);
437       crossentropy_term = sum(diag(-log(double(Y)'*P)));
438       reg_term = sum(W{1}(:).^2)+sum(W{2}(:).^2);
439       J = (1/(size(X,2))*crossentropy_term)+(lambda*reg_term);
440   end
441
442   function [P,h] = EvaluateClassifier(X,W,b)
443       s1 = W{1}*X + b{1};
444       h = max(0, s1);
445       s = W{2} * h + b{2};
446       P = exp(s)./sum(exp(s));
447   end
448
449   %Function to compute the accuracy in a 2-layer neural network:
450   function acc = ComputeAccuracy(X,y,W,b)
451       [P,h] = EvaluateClassifier(X, W, b);
452       [¬,pred_idx] = max(P);
453       acc=((sum(pred_idx==y)/size(X,2)*100));
454   end
455
456   function [grad_W, grad_b] = ComputeGradients(X, Y, P, h, W,b,lambda)
457       k_x = size(X,2);
458       W1 = cell2mat(W(1));
459       b1 = cell2mat(b(1));
460       W2 = cell2mat(W(2));
461       b2 = cell2mat(b(2));
462       grad_W1 = zeros(size(W1));
463       grad_W2 = zeros(size(W2));
464       grad_b1 = zeros(size(b1));
465       grad_b2 = zeros(size(b2));
466
467       for i= 1:k_x
468           P_i = P(:, i);
469           h_i = h(:, i);
470           Y_i = Y(:, i);
471           X_i = X(:, i);
472           g = -(Y_i-P_i)';
473           grad_b2 = grad_b2 + g';
474           grad_W2 = grad_W2 + g'*h_i';
475
476           h_i(find(h_i > 0)) = 1;
477           g = g*W2*diag(h_i);
478
479           grad_b1 = grad_b1 + g';
480           grad_W1 = grad_W1 + g'*X_i';
481       end
482       grad_b1 = grad_b1/k_x;
483       grad_W1 = grad_W1/k_x + 2*lambda*W1;
484       grad_b2 = grad_b2/k_x;
485       grad_W2 = grad_W2/k_x + 2*lambda*W2;
486
487       grad_b = {grad_b1, grad_b2};
488       grad_W = {grad_W1, grad_W2};
489   end
490
491   function [grad_W, grad_b] = ComputeGradsNumSlow(X, Y, W, b, lambda, h)
492       grad_W = cell(numel(W), 1);
493       grad_b = cell(numel(b), 1);
494       for j=1:length(b)
495           grad_b{j} = zeros(size(b{j}));
496           for i=1:length(b{j})
497               b_try = b;
498               b_try{j}(i) = b_try{j}(i) - h;
499               c1 = ComputeCost(X, Y, W, b_try, lambda);
```

```matlab
500            b_try = b;
501            b_try{j}(i) = b_try{j}(i) + h;
502            c2 = ComputeCost(X, Y, W, b_try, lambda);
503            grad_b{j}(i) = (c2-c1) / (2*h);
504        end
505    end
506    for j=1:length(W)
507        grad_W{j} = zeros(size(W{j}));
508        for i=1:numel(W{j})
509            W_try = W;
510            W_try{j}(i) = W_try{j}(i) - h;
511            c1 = ComputeCost(X, Y, W_try, b, lambda);
512            W_try = W;
513            W_try{j}(i) = W_try{j}(i) + h;
514            c2 = ComputeCost(X, Y, W_try, b, lambda);
515            grad_W{j}(i) = (c2-c1) / (2*h);
516        end
517    end
518 end
519
520
521 %Mini batch function implementation without momentum and learning rate
522 %decay.
523 function [Wstar, bstar] = MiniBatchGD(X,Y,y,GDparams,W,b,lambda)
524    N = size(X,2);
525    eta = GDparams.eta;
526    n_batch = GDparams.n_batch;
527
528    for j=1:N/n_batch
529        j_start = (j-1)*n_batch + 1;
530        j_end = j*n_batch;
531        indx = j_start:j_end;
532        Xbatch = X(:,indx);
533        Ybatch = Y(:,indx);
534        [P,h] = EvaluateClassifier(Xbatch,W,b);
535        [grad_W, grad_b] = ComputeGradients(Xbatch,Ybatch,P,h, W,b,lambda);
536        W{1} = W{1} - eta*grad_W{1};
537        W{2} = W{2} - eta*grad_W{2};
538        b{1} = b{1} - eta*grad_b{1};
539        b{2} = b{2} - eta*grad_b{2};
540    end
541    bstar = b;
542    Wstar = W;
543 end
544
545 function [v_W, v_b] = InitializeMomentum(W,b)
546    v_b = {zeros(size(b{1})), zeros(size(b{2}))};
547    v_W = {zeros(size(W{1})), zeros(size(W{2}))};
548 end
549
550 function [Wstar, bstar] = MiniBatchGD_withMomentum(X, Y,y,GDparams,W,b,lambda,rho)
551    N = size(X,2);
552    eta = GDparams.eta;
553    n_batch = GDparams.n_batch;
554    [v_W,v_b] = InitializeMomentum(W,b);
555
556    for j=1:N/n_batch
557        j_start = (j-1)*n_batch + 1;
558        j_end = j*n_batch;
559        indx = j_start:j_end;
560        Xbatch = X(:,indx);
561        Ybatch = Y(:,indx);
562        [P,h] = EvaluateClassifier(Xbatch,W,b);
563        [grad_W, grad_b] = ComputeGradients(Xbatch,Ybatch,P,h, W,b,lambda);
564        %Update with momentum:
```

```matlab
565         v_W{1} = rho*v_W{1} + eta*grad_W{1};
566         v_b{1} = rho*v_b{1} + eta*grad_b{1};
567         v_W{2} = rho*v_W{2} + eta*grad_W{2};
568         v_b{2} = rho*v_b{2} + eta*grad_b{2};
569         W{1} = W{1} - v_W{1};
570         b{1}  = b{1} - v_b{1};
571         W{2} = W{2} - v_W{2};
572         b{2}  = b{2} -  v_b{2};
573     end
574     bstar = b;
575     Wstar = W;
576 end
```