

DD2424 - Assignment 2

Pierre Rudin

May 24, 2017

1 Introduction

The aim of this assignment is to train a two layer network using a mini-batch gradient descent with momentum, so that the network can classify images from the CIFAR-10 data set. This will be done by implementing the calculations and algorithms given by the instructions for this assignment in Matlab.

1.1 Data set

CIFAR-10 is a data set of 60000 labeled 32 x 32 pixel coloured images separated in 5 training batches and 1 test batch, each batch holding 10000 images. In this assignment 2 training batches and the test batch i used.

1.2 Algorithm

This algorithm make use of some pretty simple linear algebraic equations. Since the equations are given by the instructions this section will contain my implementation step by step.

1. The data sets are loaded and the network initialized. The data set loads to \mathbf{X} (size $d \times N$), \mathbf{Y} ($K \times N$) and \mathbf{y} ($1 \times N$). The network is initialized by creating \mathbf{W}_1 ($nodesInHidden \times d$), \mathbf{W}_2 ($K, nodesInHidden$), \mathbf{b}_1 ($nodesInHidden, 1$) and \mathbf{b}_2 ($K \times 1$) which will contain randomly generated numbers. This is also where we create the mini-batches ($\mathbf{XBatches}$ and $\mathbf{YBatches}$).
2. Take the j :th batches from $\mathbf{XBatches}$ and $\mathbf{YBatches}$ and make label predictions on each image in $\mathbf{XBatches}$.
3. Compute gradients on the predictions and batches from previous step.
4. Add momentum
5. Update the network.
6. Repeat from step 2 until end of epochs or other condition is satisfied.

2 Results

2.1 Analytical gradient checking

The analytically computed gradients were compared to numerically computed gradients, to compute the numerical gradients the code given with the assignment was used. The gradients are checked on a two layer network

	Sum	Mean	Min	Max
b1	6.191410e-06	1.238282e-07	6.644981e-13	6.188640e-06
b2	1.374946e-05	2.749892e-08	2.733445e-08	2.767633e-08
W1	4.240115e-03	2.760491e-08	2.723168e-08	1.127854e-06
W2	4.500134e-06	4.500134e-07	4.499698e-07	4.501177e-07

Table 1: Absolute difference between numerically and analytically computed gradients

Since the biggest difference between the numerically and analytically computed gradients are smaller than $1e-5$, it's safe to say that the analytically computed gradients are accurate.

2.2 Momentum

To test the effects of momentum, a small sub-set of the training data is selected to perform many training epochs in a reasonable time. It is then possible to compare the output graphs to determine the effects of the added momentum.

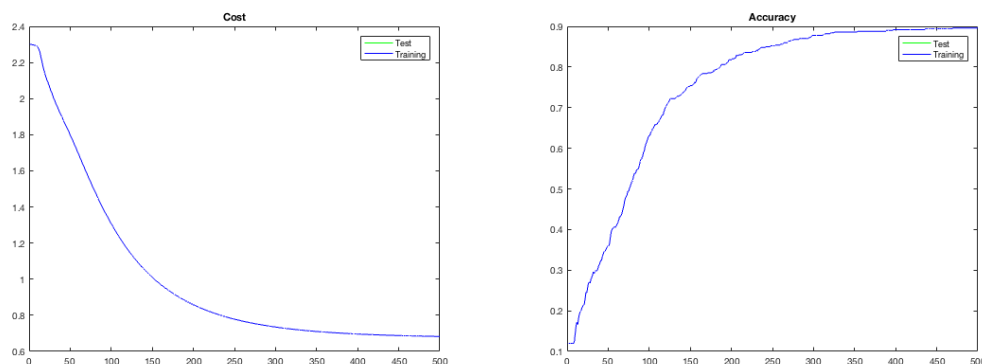


Figure 1: Graphs displaying cost and accuracy development for each epoch, when $\rho = .2$.

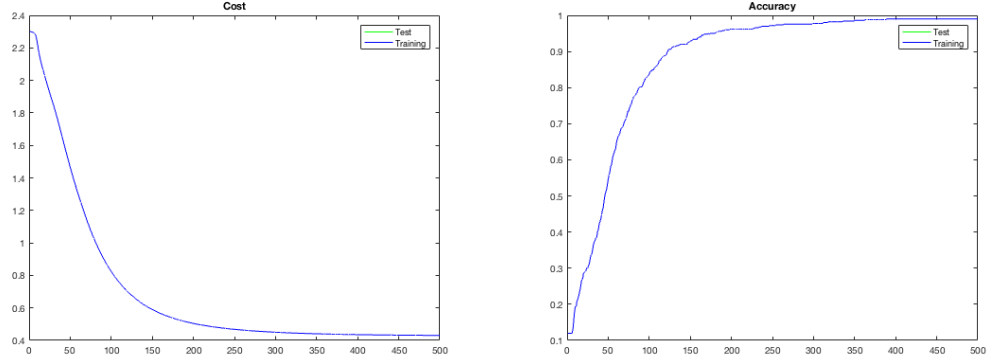


Figure 2: Graphs displaying cost and accuracy development for each epoch, when $\rho = .45$.

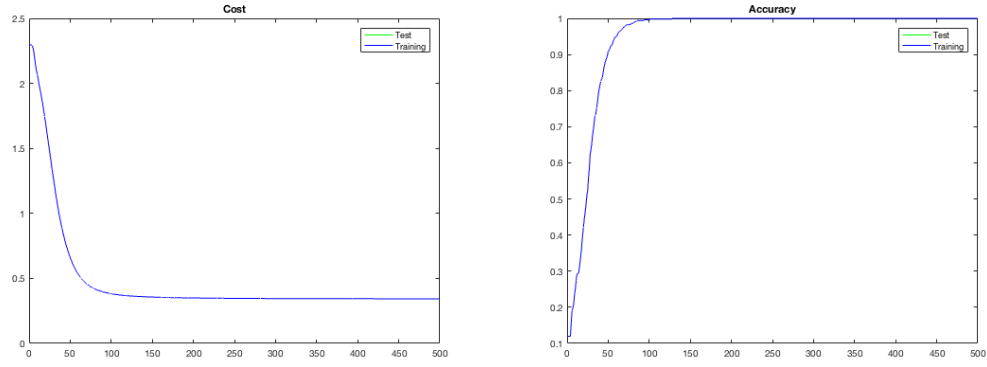


Figure 3: Graphs displaying cost and accuracy development for each epoch, when $\rho = .7$.

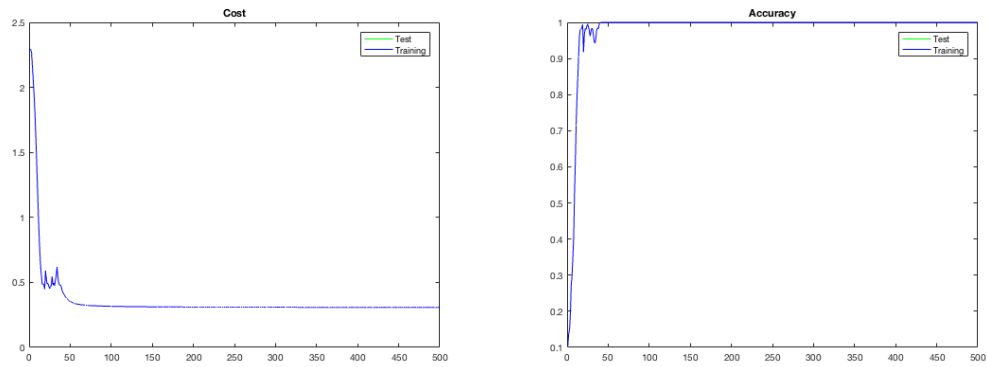


Figure 4: Graphs displaying cost and accuracy development for each epoch, when $\rho = .95$.

As shown in the graphs the accuracy curve grows faster and the cost curve descends faster when the momentum increases.

2.3 Coarse search

The initial search for our parameters starts with these values;

$$\eta = [.01, .025, .05, .1, .25, .5, .75]$$

$$\lambda = [0, .001, .0025, .005, .01, .025, .05, .1]$$

Which combined adds up to 56 different combination. Each combination is used to train the network for 5 epochs.

2.3.1 #1

Parameters:

$$\eta = .025$$

$$\lambda = .001$$

Accuracy test: .4243 Accuracy train: .4884

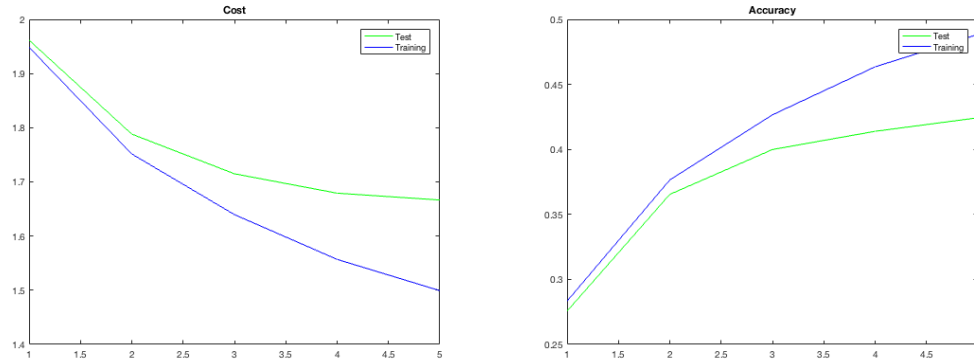


Figure 5: Graphs displaying cost and accuracy development for each epoch, when $\eta = .025$ and $\lambda = .001$.

2.3.2 #2

Parameters:

$$\eta = .025$$

$$\lambda = 0$$

Accuracy test: .4243 Accuracy train: .4938

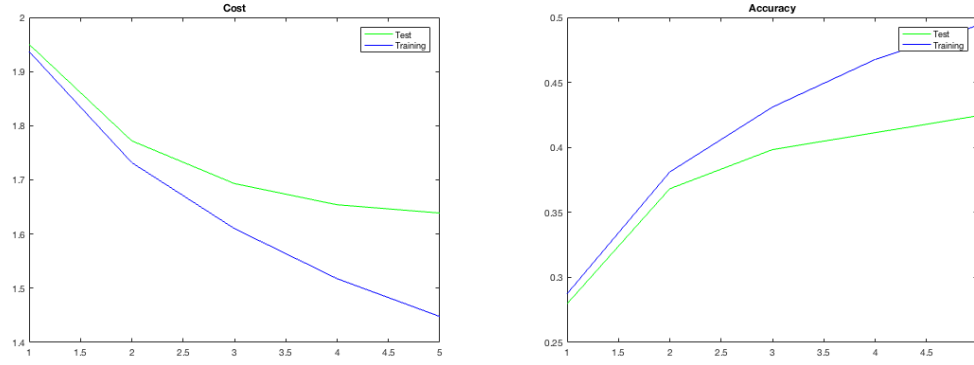


Figure 6: Graphs displaying cost and accuracy development for each epoch, when $\eta = .025$ and $\lambda = 0$.

2.3.3 #3

Parameters:

$$\eta = .025$$

$$\lambda = .0025$$

Accuracy test: .4146 Accuracy train: .4764

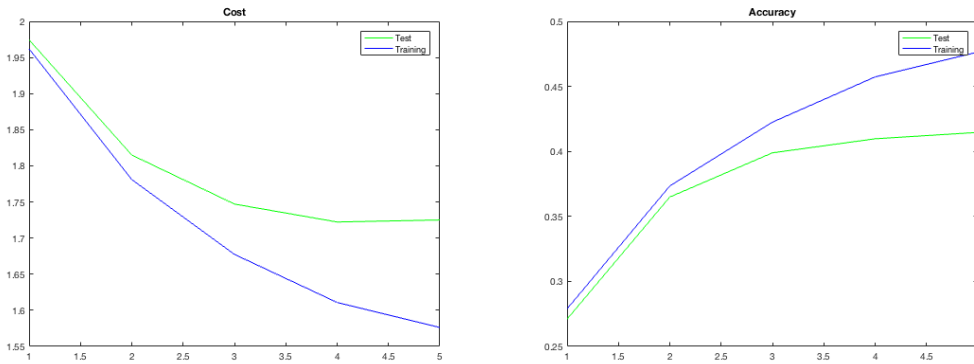


Figure 7: Graphs displaying cost and accuracy development for each epoch, when $\eta = .025$ and $\lambda = .0025$.

2.4 Fine search

With the result from the Coarse search in mind, the following parameters were chosen for the Fine search;

$$\eta = (.04-.01) .* \text{rand}(1,8) + .01$$

$$\lambda = (.005) .* \text{rand}(1,8)$$

Which combined adds up to 64 different combinations. Each combination is used to train the network for 10 epochs.

2.4.1 #1

Parameters:

$$\eta = .016228$$

$$\lambda = .002828$$

Accuracy test: .4413, Accuracy train: .5242

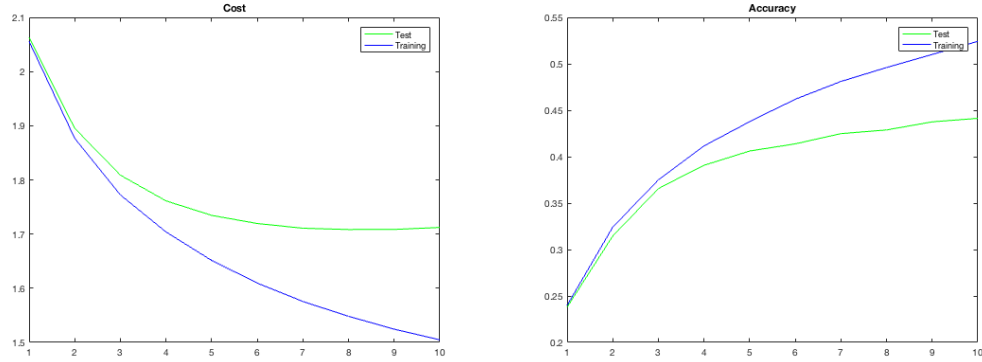


Figure 8: Graphs displaying cost and accuracy development for each epoch, when $\eta = .016228$ and $\lambda = .002828$.

2.4.2 #2

Parameters:

$$\eta = .022175$$

$$\lambda = .001514$$

Accuracy test: .4376, Accuracy train: .5532

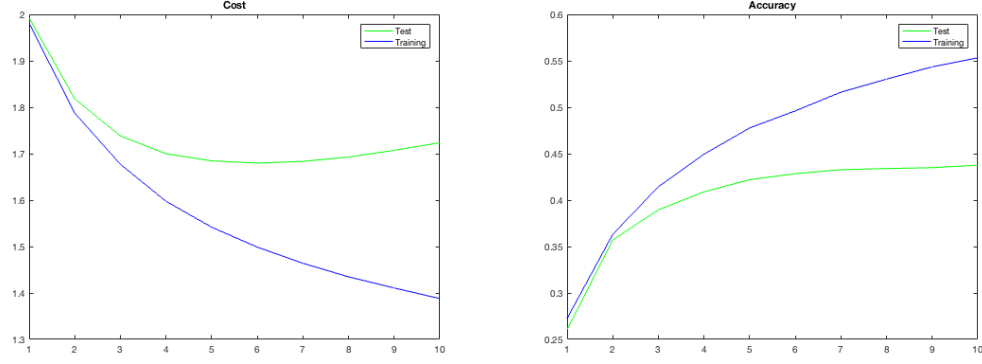


Figure 9: Graphs displaying cost and accuracy development for each epoch, when $\eta = .022175$ and $\lambda = .001514$.

2.4.3 #3

Parameters:

$$\eta = .022175$$

$$\lambda = .003633$$

Accuracy test: .4354, Accuracy train: .5253

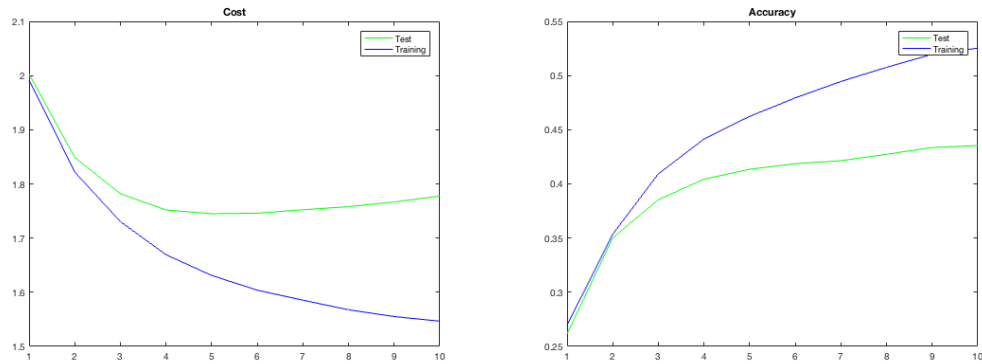


Figure 10: Graphs displaying cost and accuracy development for each epoch, when $\eta = .022175$ and $\lambda = .003633$.

2.5 Best parameters

When the best parameters are found, a further run is done where the network is trained for 30 epochs on the full training set except for the last 1000 images that are used for validation. Parameters:

$$\eta = .016228$$

$$\lambda = .002828$$

Accuracy test: .742, Accuracy train: .648222

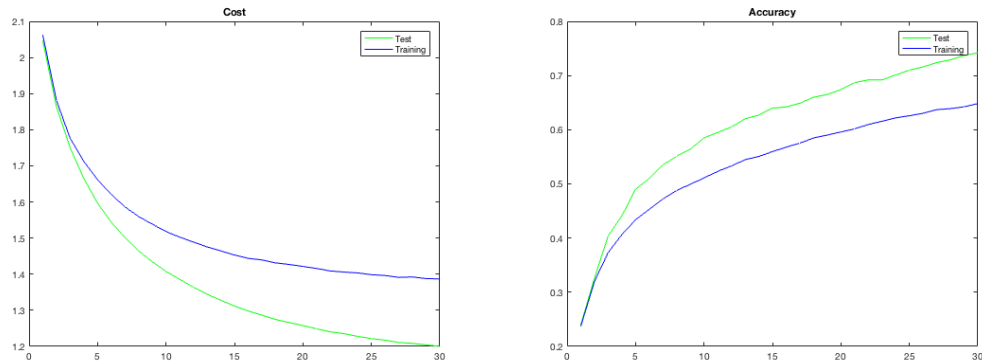


Figure 11: Graphs displaying cost and accuracy development for each epoch, when $\eta = .016228$ and $\lambda = .002828$.

3 Code

3.1 Main

```
1 clear all;
2 clc;
3 close all;
4 format longEng
5
6 addpath Datasets/cifar-10-batches-mat;
7
8 % Parameters
9 n_batch = 50; %
10 n_epochs = 5; %
11 h = 1e-5; %
12 nodes_in_hidden_layer = 100; % Number of nodes in hidden layer
13
```



```

14 % Hyper parameters
15 eta = .044576; % learning rate
16 lambda = .00998; % regularization
17 decay_rate = .998; % decay in learning rate
18 rho = .8; % Momentum
19
20 % Data setup
21 [X,Y,y,mean_X] = LoadBatch('data_batch_1.mat');
22 [XValid, YValid, yValid] = LoadData('data_batch_2.mat', mean_X);
23 [XTest, YTest, yTest] = LoadData('test_batch.mat', mean_X);
24 [XBatches, YBatches] = GetMiniBatches(X, Y, n_batch);
25
26 [W, b] = InitModel(X,nodes_in_hidden_layer);
27
28 [W,b, costs_train , costs_test , accs_train , accs_test] = trainingLoop(
    XBatches, YBatches, W, b, n_epochs, eta, lambda, rho, decay_rate,
    nodes_in_hidden_layer, X, Y, y, XTest, YTest, yTest);
29
30 PTest = EvaluateClassifier(XTest, W, b);
31 acc_test = ComputeAccuracy(PTest, yTest);
32
33 fprintf('Accuracy test\t%f\tCost test\t', acc_test);
34 for n = 1:length(costs_test)
35     fprintf('%f\t', costs_test(n));
36 end;
37
38 P = EvaluateClassifier(X, W, b);
39 acc_train = ComputeAccuracy(P, y);
40
41 fprintf('Accuracy train\t%f\tCost train\t', acc_train);
42 for n = 1:length(costs_train)
43     fprintf('%f\t', costs_train(n));
44 end;
45
46
47 % Plots evolution of the cost
48 figure(2);
49 x = 1:1:n_epochs;
50 plot(x, costs_test, 'g', x, costs_train, 'b');
51 title('Cost')
52 legend('Test', 'Training')
53

```

```

54 % Plots evolution of the accuracy
55 figure(3);
56 x = 1:1:n_epochs;
57 plot(x, accs_test, 'g', x, accs_train, 'b');
58 title('Accuracy')
59 legend('Test', 'Training')

```

3.2 LoadBatch()

```

1 function [X, Y, y, mean_X] = LoadBatch(filename)
2     A = load(filename);
3     X = double(A.data') / 255;
4     y = A.labels';
5     y = y + uint8(ones(1, length(y))); % Add one to simplify
        indexing
6
7     mean_X = mean(X, 2);
8     X = X - repmat(mean_X, [1, size(X, 2)]);
9
10    % Create image label matrix
11    Y = zeros(10, length(X));
12    for i = 1:length(Y)
13        Y(y(i), i) = 1;
14    end;

```

3.3 LoadData()

```

1 function [X, Y, y] = LoadData(filename, mean_X)
2     A = load(filename);
3     X = double(A.data') / 255;
4     y = A.labels';
5     y = y + uint8(ones(1, length(y))); % Add one to simplify
        indexing
6
7     X = X - repmat(mean_X, [1, size(X, 2)]);
8
9     % Create image label matrix
10    Y = zeros(10, length(X));
11    for i = 1:length(Y)
12        Y(y(i), i) = 1;
13    end;

```

3.4 InitModel()

```

1 function [W, b] = InitModel(X)
2     [d, ~] = size(X);
3     W = normrnd(0, .01, 10, d);
4     b = normrnd(0, .01, 10, 1);

```

3.5 GetMiniBatches()

```

1 function [XBatches, YBatches] = GetMiniBatches(Xtrain, Ytrain,
    n_batch)
2     [d, N] = size(Xtrain);
3     [K, ~] = size(Ytrain);
4
5     XBatches = zeros(N/n_batch, d, n_batch);
6     YBatches = zeros(N/n_batch, K, n_batch);
7
8     for j=1:N/n_batch
9         j_start = (j-1)*n_batch + 1;
10        j_end = j*n_batch;
11        Xbatch = Xtrain(:, j_start:j_end);
12        Ybatch = Ytrain(:, j_start:j_end);
13
14        XBatches(j, :, :) = Xbatch;
15        YBatches(j, :, :) = Ybatch;
16    end
17
18    % Permute to simplify picking out image representations from
    % the
19    % matrices.
20    XBatches = permute(XBatches, [3 2 1]);
21    YBatches = permute(YBatches, [3 2 1]);

```

3.6 trainingLoop()

```

1
2 function [W, b, costs_train, costs_test, accs_train, accs_test] =
    trainingLoop(XBatches, YBatches, W, b, n_epochs, eta, lambda, rho,
    decay_rate, nodes_in_hidden_layer, X, Y, y, XTest, YTest, yTest)
3     accs_train = double.empty();
4     accs_test = double.empty();
5     costs_train = double.empty();
6     costs_test = double.empty();
7
8     %Momentum
9     moment_W = {zeros(size(W{1})), zeros(size(W{2}))};

```

```

10 moment_b = {zeros(size(b{1})), zeros(size(b{2}))};
11
12 [~,~,l] = size(XBatches);
13
14 for i = 1:n_epochs
15     for j = 1:l
16         % Get j:th batch
17         XBatch = XBatches(:, :, j)';
18         YBatch = YBatches(:, :, j)';
19
20         %hidden_layer = MakeHiddenLayer(XBatch, W{1}, b{1});
21         P = EvaluateClassifier(XBatch, W, b);
22
23         [grad_W, grad_b] = ComputeGradients(XBatch, YBatch, P
24             , W, b, lambda, nodes_in_hidden_layer);
25
26         % Calculate momentum
27         moment_W{1} = eta * grad_W{1} + rho * moment_W{1};
28         moment_b{1} = eta * grad_b{1} + rho * moment_b{1};
29         moment_W{2} = eta * grad_W{2} + rho * moment_W{2};
30         moment_b{2} = eta * grad_b{2} + rho * moment_b{2};
31
32         % Update W's and b's
33         W{1} = W{1} - moment_W{1};
34         b{1} = b{1} - moment_b{1};
35         W{2} = W{2} - moment_W{2};
36         b{2} = b{2} - moment_b{2};
37
38     end;
39     eta = decay_rate * eta;
40
41     %hidden_layer = MakeHiddenLayer(X, W{1}, b{1});
42     P = EvaluateClassifier(X, W, b);
43     acc_train = ComputeAccuracy(P, y);
44     cost_train = ComputeCost(X, Y, W, b, lambda);
45
46     %hidden_layer = MakeHiddenLayer(XTest, W{1}, b{1});
47     PTest = EvaluateClassifier(XTest, W, b);
48     acc_test = ComputeAccuracy(PTest, yTest);
49     cost_test = ComputeCost(XTest, YTest, W, b, lambda);
50
51     accs_train = [accs_train acc_train];

```

```

51         accs_test = [accs_test acc_test ];
52
53         costs_train = [costs_train cost_train];
54         costs_test = [costs_test cost_test ];
55
56         % Displays the cost and accuracy of each epoch
57         fprintf('Cost train: %f\n', cost_train);
58         fprintf('Accuracy train: %f\n', acc_train);
59         fprintf('Cost test: %f\n', cost_test);
60         fprintf('Accuracy test: %f\n\n', acc_test);
61
62         fprintf('Epoch %f is done.\n', i);
63     end;

```

3.7 EvaluateClassifier()

```

1 function P = EvaluateClassifier(X, W, b)
2     hiddenLayer = MakeHiddenLayer(X, W{1}, b{1});
3     s2 = W{2}*hiddenLayer+b{2};
4     P = exp(s2) ./ sum(exp(s2));

```

3.8 MakeHiddenLayer()

```

1 function hiddenLayer = MakeHiddenLayer(X, W, b)
2     s = W * X + b;
3     hiddenLayer = max(0,s);

```

3.9 ComputeGradients()

```

1 function [grad_W, grad_b] = ComputeGradients(X, Y, P, W, b,
2     lambda, nodes_in_hidden_layer)
3     [hY, ~] = size(Y);
4     [hX, lX] = size(X);
5
6     grad_W = {zeros(nodes_in_hidden_layer, hX), zeros(hY,
7         nodes_in_hidden_layer)};
8     grad_b = {zeros(nodes_in_hidden_layer, 1), zeros(hY, 1)};
9
10    hidden_layer = MakeHiddenLayer(X, W{1}, b{1});
11
12    for i = 1:lX
13        yy = Y(:, i);
14        xx = X(:, i);

```

```

14     p = P(:, i);
15     [hp, ~] = size(p);
16     diag_p = eye(hp) .* p;
17     g = -((yy' / (yy' * p)) * (diag_p - p * p'));
18
19     grad_b{2} = grad_b{2} + g';
20     grad_W{2} = grad_W{2} + g' * hidden_layer(:, i)';
21
22     s1 = W{1} * xx + b{1};
23     [hs1, ~] = size(s1);
24
25     g = g * W{2};
26     diag_s1 = eye(hs1) .* (s1 > 0);
27     g = g * diag_s1;
28
29
30
31     grad_b{1} = grad_b{1} + g';
32     grad_W{1} = grad_W{1} + g' * xx';
33 end;
34
35     grad_b{1} = grad_b{1} ./ lX;
36     grad_W{1} = grad_W{1} ./ lX + 2 * lambda .* W{1};
37
38     grad_b{2} = grad_b{2} ./ lX;
39     grad_W{2} = grad_W{2} ./ lX + 2 * lambda .* W{2};

```

3.10 ComputeCost()

```

1 function J = ComputeCost(X, Y, W, b, lambda)
2     P = EvaluateClassifier(X, W, b);
3     c = 0;
4     [~, x] = size(X);
5     for i = 1:x
6         y = Y(:, i);
7         p = P(:, i);
8         c = c + -log(y' * p);
9     end;
10    J = c/x + lambda * sum(sum(W{1}.^2)) + lambda * sum(sum(W
        {2}.^2));

```

3.11 ComputeAccuracy()

```

1 function acc = ComputeAccuracy(P, y)

```

```
2     [~, prediction] = max(P);  
3     acc = sum(prediction==y) / length(P);
```