

DD2424 - Assignment 3

Pierre Rudin

June 7, 2017

1 Introduction

The aim of this assignment is to implement the methods used in previous assignments so that they can be used with any number of layers, and there after train a 3 layer network using a mini-batch gradient descent with momentum and batch normalization, 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 ($nodes_{hiddenLayer1} \times d$), $\mathbf{W}_{2,...,n-1}$ ($nodes_{thisLayer} \times nodes_{previousLayer}$) \mathbf{W}_n ($K, nodes_{hiddenLayern-1}$), \mathbf{b}_1 ($nodes_{hiddenLayer1}, 1$), $\mathbf{b}_{2,...,n-1}$ ($nodes_{thisLayer}, 1$) and \mathbf{b}_n ($K \times 1$) 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 with batch normalization added 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.

2.1.1 2-Layer Network

	Sum	Mean	Min	Max
b1	2.821483e-01	5.642965e-03	2.116609e-04	2.344918e-02
b2	4.500311e-06	4.500311e-07	4.481301e-07	4.540629e-07
W1	1.003890e+02	6.535743e-04	2.087295e-09	5.148941e-03
W2	1.062019e-04	2.124039e-07	1.658791e-07	2.622531e-07

Table 1: Absolute difference between numerically and analytically computed gradients on a 2-layer network

With batch normalization the difference between the analytically and numerically calculated gradients are bigger than in previous assignments. But since the mean difference still are in the range of $1e-3$ or smaller I would still consider the analytically calculated gradients to be quite accurate.

2.1.2 3-Layer Network

	Sum	Mean	Min	Max
b1	1.555098e-01	3.110196e-03	8.404835e-05	1.305131e-02
b2	4.115022e-01	1.371674e-02	3.729078e-04	4.355665e-02
b3	4.499778e-06	4.499778e-07	4.492763e-07	4.504875e-07
W1	1.126712e+02	7.335363e-04	8.897409e-09	4.812008e-03
W2	1.043999e+01	6.959991e-03	5.697405e-06	5.041149e-02
W3	3.832258e-05	1.277419e-07	7.448231e-08	1.879620e-07

Table 2: Absolute difference between numerically and analytically computed gradients on a 3-layer network

Adding a layer increases the difference between the methods of calculating the gradients. This however is expected, and it is concluded that the numbers for 3 layers look good.

2.1.3 4-Layer Network

	Sum	Mean	Min	Max
b1	2.708227e-01	2.708227e-03	3.645406e-05	9.005303e-03
b2	3.209452e-01	6.418905e-03	1.289368e-04	2.986614e-02
b3	3.332370e-01	1.110790e-02	4.138029e-06	3.213377e-02
b4	4.499423e-06	4.499423e-07	4.491167e-07	4.509904e-07
W1	2.191894e+02	7.135072e-04	3.616336e-09	6.332208e-03
W2	2.316350e+01	4.632700e-03	9.506812e-08	3.893993e-02
W3	7.959463e+00	5.306309e-03	8.669500e-06	2.761229e-02
W4	3.708043e-05	1.236014e-07	7.681426e-08	2.101081e-07

Table 3: Absolute difference between numerically and analytically computed gradients on a 4-layer network

Adding yet another layer increases the difference between the methods of calculating the gradients even further. This, again, however is expected, and it is concluded that the numbers for 4 layers looks good.

2.2 Batch Normalization

To test the effects of Batch Normalization, 3 test runs where made with Batch Normalization, and 3 with out Batch Normalization. Each run is then presented with a graph in the section bellow.

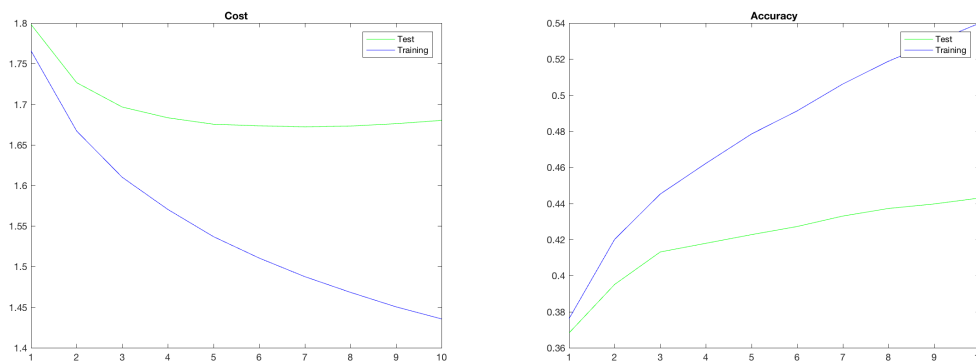


Figure 1: Graphs displaying cost and accuracy development for each epoch, with Batch Normalization and $\eta = .025$.

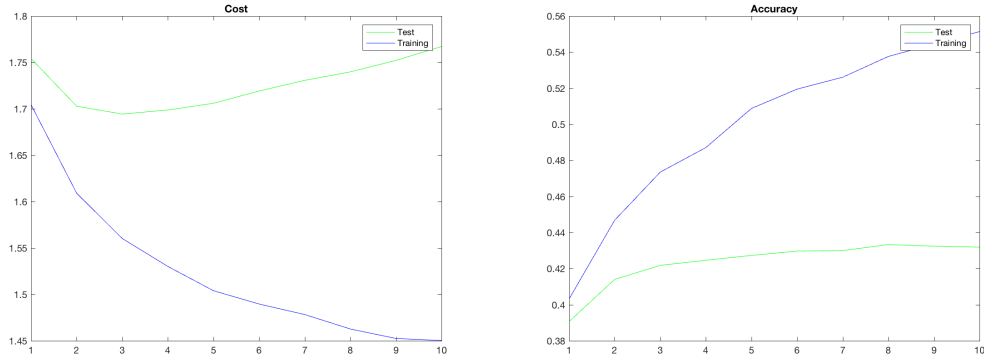


Figure 2: Graphs displaying cost and accuracy development for each epoch, with Batch Normalization and $\eta = .05$.

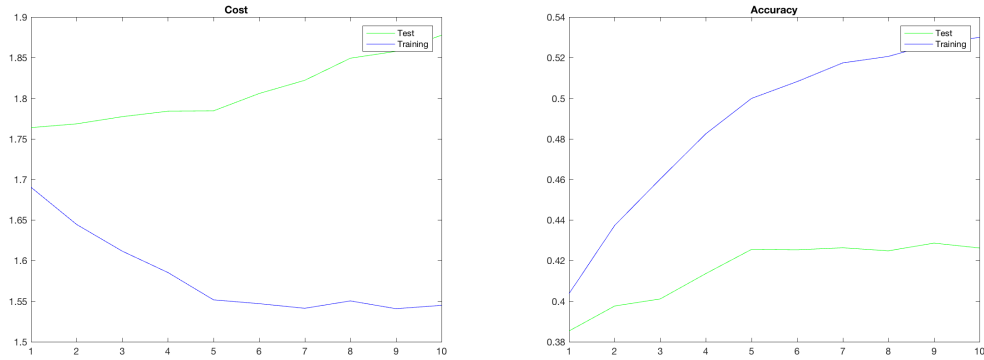


Figure 3: Graphs displaying cost and accuracy development for each epoch, with Batch Normalization and $\eta = .1$.

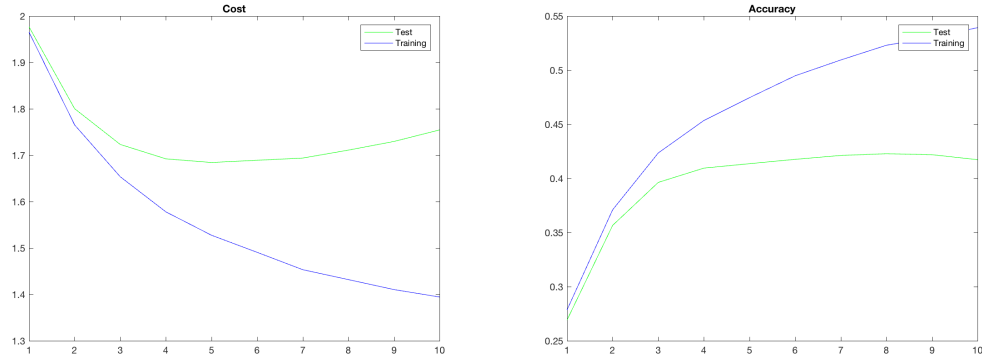


Figure 4: Graphs displaying cost and accuracy development for each epoch, without Batch Normalization and $\eta = .025$.

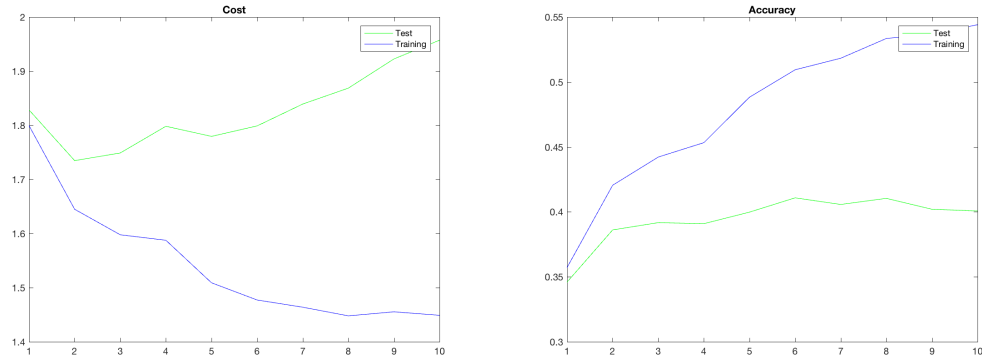


Figure 5: Graphs displaying cost and accuracy development for each epoch, without Batch Normalization and $\eta = .05$.

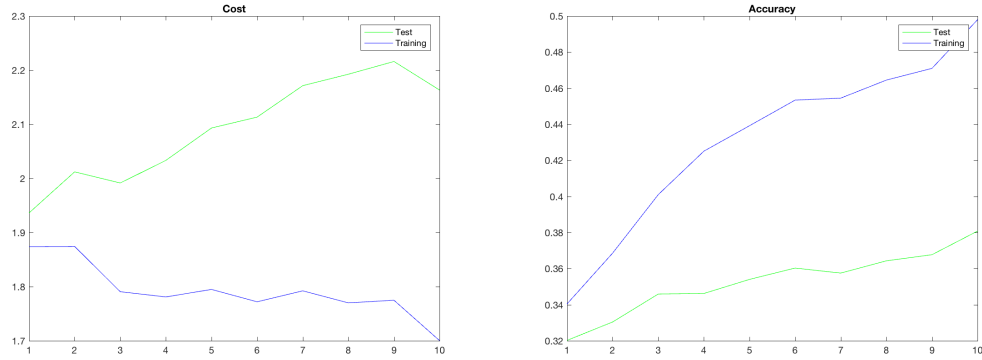


Figure 6: Graphs displaying cost and accuracy development for each epoch, without Batch Normalization and $\eta = .1$.

As we can see in the graphs, the gap between the cost on the training data and the test data is smaller when Batch Normalization is used. It's also noted that the overall accuracy on the test data is higher when Batch Normalization is used.

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:

$$\eta = .1$$

$$\lambda = 0$$

Accuracy test: .4181 Accuracy train: .4695

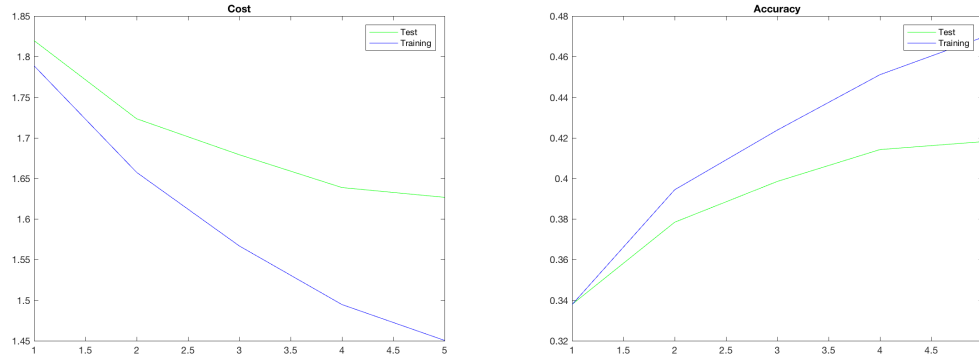


Figure 7: Graphs displaying cost and accuracy development for each epoch, when $\eta = .1$ and $\lambda = 0$.

2.3.2 #2

Parameters:

$$\eta = .1$$

$$\lambda = .0005$$

Accuracy test: .414 Accuracy train: .4665

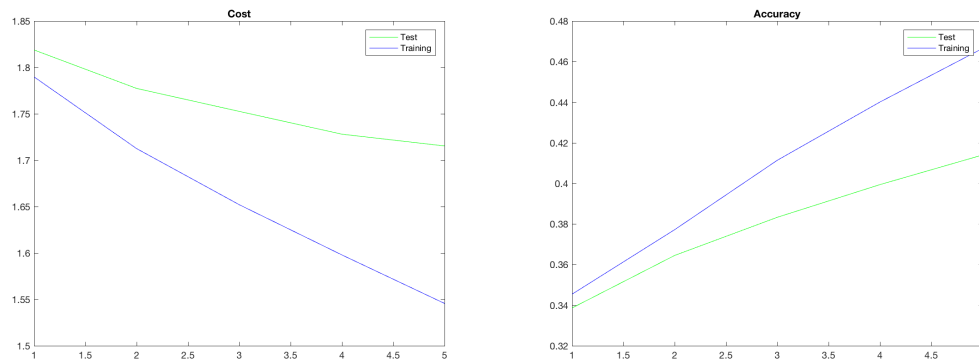


Figure 8: Graphs displaying cost and accuracy development for each epoch, when $\eta = .1$ and $\lambda = .0005$.

2.3.3 #3

Parameters:

$$\eta = .05$$

$$\lambda = 0$$

Accuracy test: .3852 Accuracy train: .4102

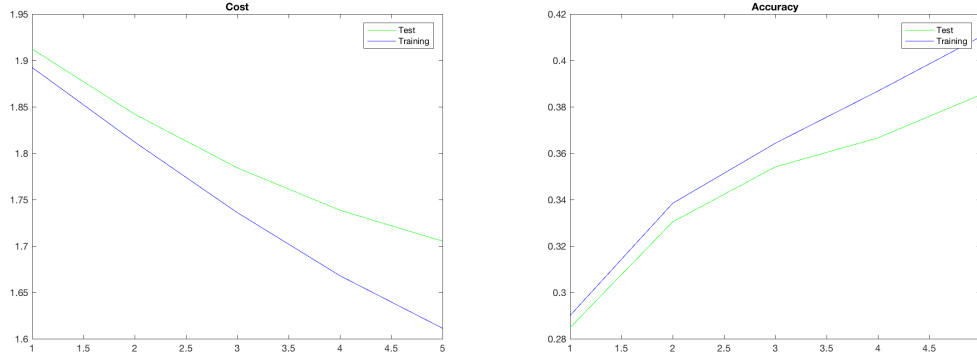


Figure 9: Graphs displaying cost and accuracy development for each epoch, when $\eta = .05$ and $\lambda = 0$

2.4 Fine search

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

$$\eta = (.15-.03) \cdot \text{rand}(1,8) + .03$$

$$\lambda = .001 \cdot \text{rand}(1,8)$$

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

2.4.1 #1

Parameters:

$$\eta = .045828$$

$$\lambda = .000662$$

Accuracy test: .4266, Accuracy train: .4744

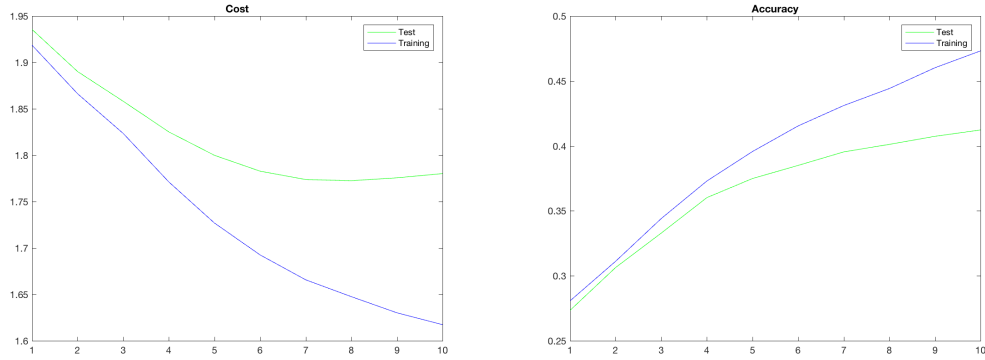


Figure 10: Graphs displaying cost and accuracy development for each epoch, when $\eta = .045828$ and $\lambda = .000662$.

2.4.2 #2

Parameters:

$$\eta = .045828$$

$$\lambda = .000378$$

Accuracy test: .42, Accuracy train: .5008

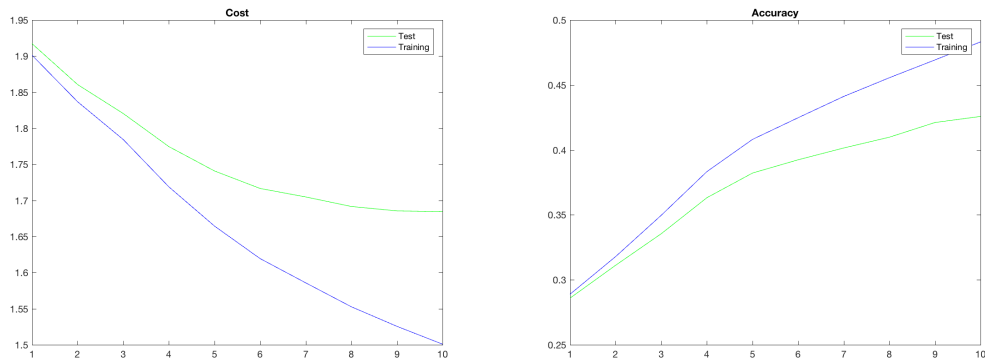


Figure 11: Graphs displaying cost and accuracy development for each epoch, when $\eta = .045828$ and $\lambda = .000378$.

2.4.3 #3

Parameters:

$$\eta = .045828$$

$$\lambda = .000261$$

Accuracy test: .4092, Accuracy train: .5033

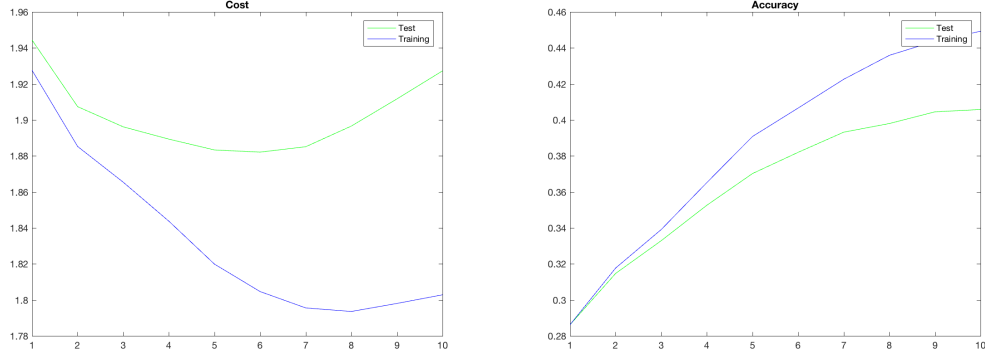


Figure 12: Graphs displaying cost and accuracy development for each epoch, when $\eta = .045828$ and $\lambda = .000261$.

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:

$$\eta = .045828$$

$$\lambda = .000662$$

Accuracy test: .432, Accuracy train: .619778

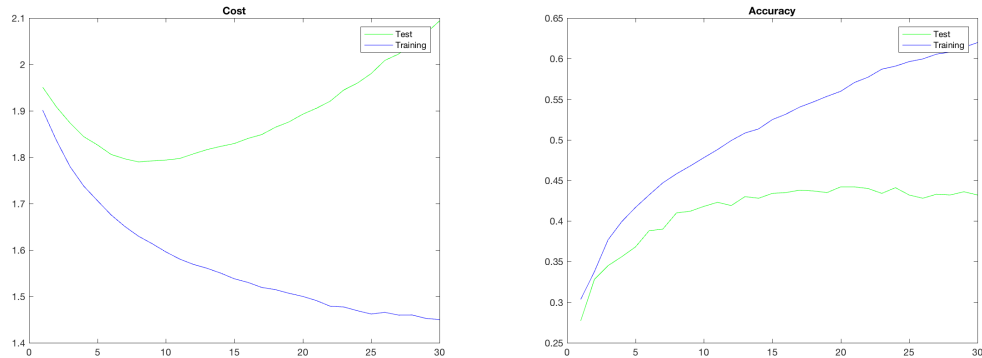


Figure 13: Graphs displaying cost and accuracy development for each epoch, when $\eta = .045828$ and $\lambda = .000662$.

3 Code

3.1 Main

```
1 clear all;
2 clc;
3 close all;
4
5 addpath Datasets/cifar-10-batches-mat;
6
7 % Parameters
8 n_batch = 50;
9 n_epochs = 10;
10 h = 1e-5; %
11 nodes_in_hidden_layers = [50, 30]; % Number of nodes in hidden
    layers
12
13 % Hyper parameters
14 eta = 0.01; % learning rate
15 lambda = 0.001; % regularization
16 decay_rate = .998; % decay in learning rate
17 rho = .8; % Momentum
18
19 % Data setup
20 [X,Y,y,mean_X] = LoadBatch('data_batch_1.mat');
21 [XValid, YValid, yValid] = LoadData('data_batch_2.mat', mean_X);
22 [XTest,YTest,yTest] = LoadData('test_batch.mat', mean_X);
23 [XBatches, YBatches] = GetMiniBatches(X, Y, n_batch);
24 [W, b] = InitModel(X,nodes_in_hidden_layers);
25
26 trainingLoop(XBatches,YBatches,W,b,n_epochs,eta,lambda,rho,
    decay_rate,nodes_in_hidden_layers,X,Y,y,XTest,YTest,yTest,h);
```

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, nodes_in_hidden_layers)
2     W = cell(length(nodes_in_hidden_layers)+1, 1);
3     b = cell(length(nodes_in_hidden_layers)+1, 1);
4     [d, ~] = size(X);
5
6     W{1} = normrnd(0, .001, nodes_in_hidden_layers(1), d);
7     b{1} = zeros(nodes_in_hidden_layers(1), 1);
8
9     for i = 2:length(nodes_in_hidden_layers)
10        W{i} = normrnd(0, .001, nodes_in_hidden_layers(i),
            nodes_in_hidden_layers(i-1));
11        b{i} = zeros(nodes_in_hidden_layers(i), 1);
12    end;
13
14    W{length(nodes_in_hidden_layers)+1} = normrnd(0, .001, 10,
        nodes_in_hidden_layers(length(nodes_in_hidden_layers)));
15    b{length(nodes_in_hidden_layers)+1} = zeros(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_layers,X,Y,y,XTest,YTest,yTest,h)
3     epsilon = 1e-5;
4
5     accs_train = double.empty(n_epochs,0);
6     accs_test = double.empty(n_epochs,0);
7     costs_train = double.empty(n_epochs,0);
8     costs_test = double.empty(n_epochs,0);
9
10    n = length(nodes_in_hidden_layers);
11
12    %Momentum
13    moment_W = cell(n+1,1);
```

```

14     moment_b = cell(n+1,1);
15     mu_exp = cell(n+1,1);
16     v_exp = cell(n+1,1);
17     for i = 1:n+1
18         moment_W{i} = zeros(size(W{i}));
19         moment_b{i} = zeros(size(b{i}));
20     end
21
22
23     [~,~,l] = size(XBatches);
24
25     for i = 1:n_epochs
26         for j = 1:l
27             % Get j:th batch
28             XBatch = XBatches(:, :, j)';
29             YBatch = YBatches(:, :, j)';
30
31             [P, mu_exp, v_exp] = EvaluateClassifier(XBatch, W, b,
32                 epsilon, 'train', mu_exp, v_exp);
33             [grad_W, grad_b, mu_exp, v_exp] = ComputeGradients(
34                 XBatch, YBatch, P, W, b, lambda,
35                 nodes_in_hidden_layers, epsilon, 'train', mu_exp,
36                 v_exp);
37
38             % Calculate momentum
39             for m = 1:n+1
40                 moment_W{m} = eta * grad_W{m} + rho * moment_W{m};
41                 moment_b{m} = eta * grad_b{m} + rho * moment_b{m};
42             end
43
44             % Update W's and b's
45             for m = 1:n+1
46                 W{m} = W{m} - moment_W{m};
47                 b{m} = b{m} - moment_b{m};
48             end
49         end;
50     end;
51     eta = decay_rate * eta;

```

```

50     [P, mu_exp, v_exp] = EvaluateClassifier(X, W, b, epsilon,
51         'train', mu_exp, v_exp);
52     accs_train(i) = ComputeAccuracy(P, y);
53     costs_train(i) = ComputeCost(X, Y, W, b, lambda, epsilon,
54         'train', mu_exp, v_exp);
55
56     [PTest, ~, ~] = EvaluateClassifier(XTest, W, b, epsilon,
57         'test', mu_exp, v_exp);
58     accs_test(i) = ComputeAccuracy(PTest, yTest);
59     costs_test(i) = ComputeCost(XTest, YTest, W, b, lambda,
60         epsilon, 'test', mu_exp, v_exp);
61
62     %Displays the cost and accuracy of each epoch
63     fprintf('Cost train: %f\n', costs_train(i));
64     fprintf('Accuracy train: %f\n', accs_train(i));
65     fprintf('Cost test: %f\n', costs_test(i));
66     fprintf('Accuracy test: %f\n\n', accs_test(i));
67
68     fprintf('Epoch %f is done.\n', i);
69 end;
70
71 %% Plots evolution of the cost
72 figure(2);
73 x = 1:1:n_epochs;
74 plot(x, costs_test, 'g', x, costs_train, 'b');
75 title('Cost')
76 legend('Test', 'Training')
77
78 % Plots evolution of the accuracy
79 figure(3);
80 x = 1:1:n_epochs;
81 plot(x, accs_test, 'g', x, accs_train, 'b');
82 title('Accuracy')
83 legend('Test', 'Training')
84
85 close all;

```

3.7 EvaluateClassifier()

```

1 function [P, mu_exp, v_exp] = EvaluateClassifier(X, W, b, epsilon
2     , mode, mu_exp, v_exp)
3     hiddenLayer = X;
4     for i = 1 : length(W) - 1

```

```

4         [hiddenLayer, ~, ~, ~, mu_exp{i}, v_exp{i}] =
            MakeHiddenLayer(hiddenLayer, W{i}, b{i}, epsilon, mode
                , mu_exp{i}, v_exp{i});
5     end;
6     s = W{length(W)}*hiddenLayer+b{length(W)};
7     P = exp(s) ./ sum(exp(s));

```

3.8 MakeHiddenLayer()

```

1 function [hiddenLayer, s0, mu, v, mu_exp, v_exp] =
    MakeHiddenLayer(X, W, b, epsilon, mode, mu_exp, v_exp)
2     s = W * X + b;
3     s0 = s;
4     [s, mu, v] = BatchNormalize(s, epsilon, mode, mu_exp, v_exp);
5     hiddenLayer = max(0,s);
6
7     if strcmp(mode, 'train')
8         alpha = .99;
9         mu_exp = alpha * mu + (1 - alpha) * mu;
10        v_exp = alpha * v + (1 - alpha) * v;
11    end

```

3.9 ComputeGradients()

```

1 function [grad_W, grad_b, mu_exp, v_exp] = ComputeGradients(X, Y,
    P, W, b, lambda, nodes_in_hidden_layers, epsilon, mode, mu_exp
    , v_exp)
2     [hY, ~] = size(Y);
3     [hX, lX] = size(X);
4
5     n = length(nodes_in_hidden_layers);
6     hidden_layers = cell(n,1);
7     mu = cell(n-1,1);
8     v = cell(n-1,1);
9     S = cell(n-1,1);
10    hidden_layers{1} = X;
11    for i = 2:n+1
12        [hidden_layers{i}, S{i-1}, mu{i-1}, v{i-1}, mu_exp{i-1},
            v_exp{i-1}] = MakeHiddenLayer(hidden_layers{i-1}, W{i-1},
            b{i-1}, epsilon, mode, mu_exp, v_exp);
13    end;
14
15    grad_W = cell(1, n + 1);
16    grad_b = cell(1, n + 1);

```



```

17
18 grad_W{1} = zeros(nodes_in_hidden_layers(1),hX);
19 grad_b{1} = zeros(nodes_in_hidden_layers(1),1);
20
21 for i = 2:n
22     grad_W{i} = zeros(nodes_in_hidden_layers(i),
23         nodes_in_hidden_layers(i-1));
24     grad_b{i} = zeros(nodes_in_hidden_layers(i),1);
25 end;
26
27 grad_W{length(nodes_in_hidden_layers) + 1} = zeros(hY,
28     nodes_in_hidden_layers(n));
29 grad_b{length(nodes_in_hidden_layers) + 1} = zeros(hY, 1);
30
31 gs = zeros(lX,nodes_in_hidden_layers(n));
32
33 % Calculate gradients for last layer.
34 for i = 1:lX
35     y = Y(:,i);
36     p = P(:,i);
37
38     [hp,~] = size(p);
39     diag_p = eye(hp) .* p;
40     g = -((y' / (y' * p)) * (diag_p - p * p'));
41
42     grad_b{n+1} = grad_b{n+1} + g';
43     grad_W{n+1} = grad_W{n+1} + g' * hidden_layers{n+1}(:,i)
44         ';
45
46     g = g * W{n+1};
47     s = S{n}(:,i);
48     [hs, ~] = size(s);
49     diag_s = eye(hs) .* (s>0);
50     g = g * diag_s;
51
52     gs(i,:) = g;
53 end;
54
55 grad_b{n+1} = grad_b{n+1} ./ lX;
56 grad_W{n+1} = grad_W{n+1} ./ lX + 2 * lambda .* W{n+1};
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99

```

```

56     for i = n:-1:1
57         x = hidden_layers{i};
58         gs = BatchNormBackPass(gs, S{i}, mu{i}, v{i}, epsilon);
59
60
61         grad_b{i} = sum(gs, 1)' ./ lX;
62         grad_W{i} = gs' * x' ./ lX + 2 * lambda .* W{i};
63
64         gs = gs * W{i};
65         if i > 1
66             for j = 1:lX
67                 s = S{i-1}(:, j);
68                 [hs, ~] = size(s);
69                 diag_s = eye(hs) .* (s>0);
70                 gs(j, :) = gs(j, :) * diag_s;
71             end;
72         end;
73     end;

```

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);

```

3.12 BatchNormalize()

```

1  function [s, mu, v] = BatchNormalize(s, epsilon, mode, mu, v)
2      if strcmp(mode, 'train')
3          mu = mean(s, 2);

```

```

4         v = var(s,0,2);
5     end
6
7     s2 = s - mu;
8     v2 = eye(length(v)) .* (v + epsilon);
9     v2 = (v2 ^ (-1/2));
10    s = v2 * s2;

```

3.13 BatchNormBackPass()

```

1 function g = BatchNormBackPass(g,s,mu,v,epsilon)
2
3     [n, m] = size(g);
4     for i = 1:n
5         diags = eye(m) .* (s(:,i) - mu)';
6         diagv = eye(m) .* (v + epsilon);
7
8         gradv = - sum(g(i,:) * (diagv ^ (-3 / 2)) * diags, 1) /
9                 2;
10        gradmu = - sum(g(i,:) * (diagv ^ (-1 / 2)), 1);
11
12        g(i,:) = g(i,:) * (diagv ^ (-1 / 2)) + 2 * (gradv * diags
13                + gradmu) / n;
14    end;

```