

Homework 2 - Introduction to machine learning and artificial neural networks

Ali Nehrani

Resumen

Implementing simple neural network for learning given data. Modifying the gradient decent and plotting decicion boundary. There are some parameters that we can tune them in order to make the algorithm work for the given data. Beside implementing decicion boundary and Cross Entropy cost, I provided different version of the gradient decent and controlling.

Contents

1	Introduction	1
2	Parameters and My codes	1
2.1	Cross Entropy cost function	2
2.2	Decision Boundary code	2
2.3	Accelerated gradient control	2
2.4	Armijo linesearch method	2
3	Text Outputs	3
4	Plottings	3
5	Conclusion	4
6	Final Codes	6
6.1	MSE efficient code	6
6.2	Cross entropy efficient code	8

1. Introduction

In this home work I prepared 2 codes; one for MSE cost function and the other for cross entropy objective function. I designed a neural network with "5" neuron with MSE cost that can learn the data using the Nestrov's accelerated gradient decent and for cross entropy function my neural network can learn with 6 hidden layers. I described the parameters and corresponding results in the following.

2. Parameters and My codes

There are some parameters that we were allowed to modify and I modified them. I listed the my modified parameters in the following:

nV: number of neurons

eta: the learning rate in the gradient decent.

I used accelerated gradient decent which requires another

parameter *mu* and I set the values as follows for the MSE case:

```
1 % Parameters for MSE cost function code
2 nV = 5; %TUNE (number of hidden units)
3 eta = 0.87; %TUNE (learning rate)
4 mu = 0.087 % Accelerated gradient parameter
```

The average iterations for converging in this case is: The value of the parameters are almost same in the Cross Entropy case

```
1 % Parameters for Cross Entropy cost function code
2 nV = 6; %TUNE (number of hidden units)
3 eta = 0.47; %TUNE (learning rate)
4 mu = 0.077 % Accelerated gradient parameter
```

The average iterations for converging in this case is:

I used Nestrov's accelerated gradient decent (NAG) to achieve faster convergence. My NAG code is

```
1 % Parameters for MSE and Cross Entropy cost function
  code
2 if it==1
3     vect_p = -eta*dw;
4 end
5 % Accelerated Gradient
6 vect = mu*vect_p - eta*dw ; %
7 w = w - mu * vect_p + (1+mu)*vect; % + beta*
  randn(nV,3)
8 vect_p = vect;
```

The NAG that I used in both cases are almost same.

Because Starting point is random, sometimes the algorithm fall into local minima and doesnt converge. To solve this problem I check the convergency after several iterations and if the average error is above some value, then I reset the starting point. Beside I modify the learning rate to help it to learn faster.

```
1 % Parameters for MSE and Cross Entropy cost function
  code
2 % Trying to reset the starting point
3 if (mod(it,5000) == 0)
4     if E >= 1e-3
5         eta = eta + .001;
6         mu = mu + .001;
7         W = (rand(1,nV)-0.5); % hidden->output
  weights
```

```

8      w = (rand(nV,nX)-0.5); % input->hidden
      weights
9      elseif (eta > .1 && mu > .1)
10         eta = eta - 0.01;
11         mu = mu - .01;
12     end

```

2.1. Cross Entropy cost function

I changed the MSE code into cross entropy code straight forward

```

1 % Cross Entropy Cost function
2 E = -T.*log(O) - (1-T).*log(1-O);

```

The gradient of the cross entropy cost function is coded as

```

1 % Gradient of Cross Entropy Cost function
2
3 %% Gradient of the input-hidden layer
4 dW = -(T*sigder(Oh)/O).*V + ((1-T)*sigder(Oh)/(1-O)).*V;
5 %% Gradient of the hidden-output layer
6 dw = -(T*sigder(Oh)/O).*W.*sigder(Vh)'*x + ((1-T)*sigder(Oh)/(1-O)).*W.*sigder(Vh)'*x;

```

2.2. Decision Boundary code

I coded the decision boundary straight forward by using for loops and sigmoid function and trained neuron weights. It consists of two parts; at the first part I make the decision background:

```

1 xmin = -2; %
2 xmax = 2; %
3 ymin = -2; %
4 ymax = 2; %
5 figure;
6 hold on
7 dx=0.1;
8 dy=0.1;
9
10 for x1=xmin:dx:xmax
11     for y1=ymin:dy:ymax
12         activation = forwardPass([x1 y1 1],w,W);
13         if activation > 0.5
14             plot(x1,y1,'c','markersize',5)
15         else
16             plot(x1,y1,'m','markersize',5)
17         end
18     end
19 end
20 legend('Class 1', 'Class 0');

```

At the second part I draw the output of the neural net on the decision areas

```

1 nn_output = zeros(size(tr_x,1),1);
2 for i=1:size(tr_x,1)
3     nn_output(i) = forwardPass(tr_x(i,:),w,W);
4 end
5 %
6 class1 = tr_x(find(nn_output>=0.5),1:2);%
7 class0 = tr_x(find(nn_output<0.5),1:2);%
8 plot(class1(:,1),class1(:,2),'b>');
9 plot(class0(:,1),class0(:,2),'go');
10 % include legend
11 legend('Class 1', 'Class 0');
12 %legend();
13 % label the axes.

```

```

14 xlabel('x');
15 ylabel('T');

```

2.3. Accelerated gradient control

I implemented accelerated gradient and as the input is random, sometimes the algorithm sometimes falls into local minima. In order to avoid this I used some control in which the information of the error function is used. I sampled the error at some iterations and make decision based on the mean or max and min of the sampled data. The code for different controlling strategy is provided in the following:

```

1 % controlling section
2 % For analysing the state of convergence I take
3 % samples and based on
4 % the average of the samples
5 if (mod(it,500) == 0)
6     MeanE = [MeanE E];
7     MeanE = MeanE_old + .9*(E-MeanE_old);%
8     MeanE_old = MeanE
9 end
10
11 % dw_old = dw;
12 if (mod(it,5000) == 0)
13     if mean(MeanE) >= 1e-1
14         vect_p = -eta*dw;
15         eta = eta + .00001;
16         mu = mu + .001;
17 % if not converged then start from point
18 % w = .005*(rand(nV,nX)-0.5); % input->
19 % hidden weights
20 W = (rand(1,nV)-0.5)*0.05;
21 else
22     eta = eta - 0.0001;
23     mu = mu - 0.001;
24 if max(MeanE) < 1e-5
25     beta = beta*.0001;
26 elseif min(MeanE) > 1e-4
27     beta = beta*1.2;
28 end
29 end
30 MeanE=[];
31 end

```

2.4. Armijo linesearch method

I also implemented the Armijo line search method in order to choose the optimal stepsize. But according to the stochastic nature of the starting point, this doesn't work in this case.

```

1 % Armijo line search
2 E_old = E;
3 j = 1;
4 while (j>0)
5     w1 = w + eta*dw;
6     [O1, Oh1, V1, Vh1] = forwardPass(x,w1,W);
7     E1 = 0.5*(T-O)^2;
8     dw1 = -(T-O1)*sigder(Oh1)*(W.*sigder(Vh1))
9     '*x;
10
11     if (E1 <= E_old+gamma*eta*(dw1'*dw1))
12         j = 0;
13         alpha_armijo = eta;
14         mu = eta*delta;
15     else

```

```

15         eta = eta*delta
16     end
17     E_old = E1;
18 end
19 if j == 0
20     clear O1 Oh1 V1 Vh1 E_old w1
21 end

```

3. Text Outputs

With these parameters in both cases, my modified code was converged. Output of 5-layer perceptron after $10e5$ iterations:

```

1 Desired Output 1.00 == 1.00 (Network Ouput)
2 Desired Output 1.00 == 1.00 (Network Ouput)
3 Desired Output 1.00 == 1.00 (Network Ouput)
4 Desired Output 1.00 == 1.00 (Network Ouput)
5 Desired Output 0.00 == 0.00 (Network Ouput)
6 Desired Output 0.00 == 0.00 (Network Ouput)
7 Desired Output 0.00 == 0.00 (Network Ouput)
8 Desired Output 0.00 == 0.00 (Network Ouput)
9 Desired Output 1.00 == 1.00 (Network Ouput)
10 Desired Output 1.00 == 1.00 (Network Ouput)
11 Desired Output 1.00 == 1.00 (Network Ouput)
12 Desired Output 1.00 == 1.00 (Network Ouput)
13 Desired Output 0.00 == 0.00 (Network Ouput)
14 Desired Output 0.00 == 0.00 (Network Ouput)
15 RMS error over data points:0.00188

```

For more iterations I got even better results in both cases. In the following I provided output of 6-layer perceptron after $10e5$ iterations for cross entropy cost:

```

1 Desired Output 1.00 == 1.00 (Network Ouput)
2 Desired Output 1.00 == 1.00 (Network Ouput)
3 Desired Output 1.00 == 1.00 (Network Ouput)
4 Desired Output 1.00 == 1.00 (Network Ouput)
5 Desired Output 0.00 == 0.00 (Network Ouput)
6 Desired Output 0.00 == 0.00 (Network Ouput)
7 Desired Output 0.00 == 0.00 (Network Ouput)
8 Desired Output 0.00 == 0.00 (Network Ouput)
9 Desired Output 1.00 == 1.00 (Network Ouput)
10 Desired Output 1.00 == 1.00 (Network Ouput)
11 Desired Output 1.00 == 1.00 (Network Ouput)
12 Desired Output 1.00 == 1.00 (Network Ouput)
13 Desired Output 0.00 == 0.00 (Network Ouput)
14 Desired Output 0.00 == 0.00 (Network Ouput)
15 RMS error over data points:0.00002

```

Based on these results I can conclude that the Cross entropy cost function works better than MSE cost function in terms of accuracy with constant equal iterations.

4. Plottings

I provided different plots from different running and I provided captions for them. As we have two classes, I denoted them with different colors. The training points are also denoted in the class areas that they are belonging (with different shapes).

Figure 1 represents sharp convergence with cross entropy cost function and Figure 2 represents the corresponding decision boundary. I set number of iterations into $5 \times 10e5$. For comparing Cross Entropy case with MSE I run

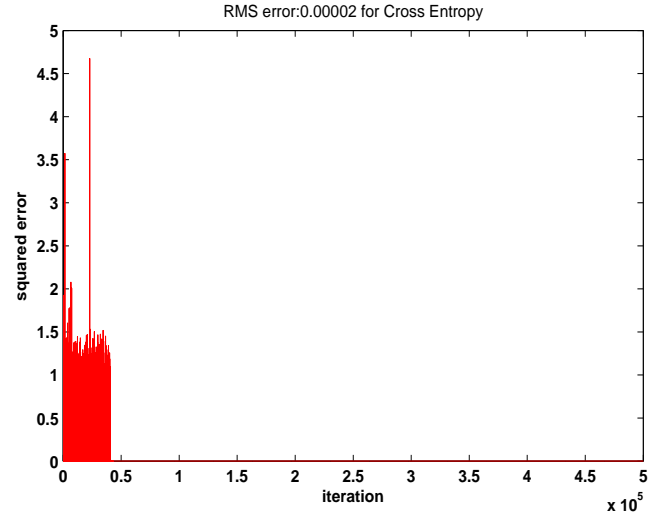


Figure 1: RMS error for 6 layer net with cross entropy case - sharp convergence

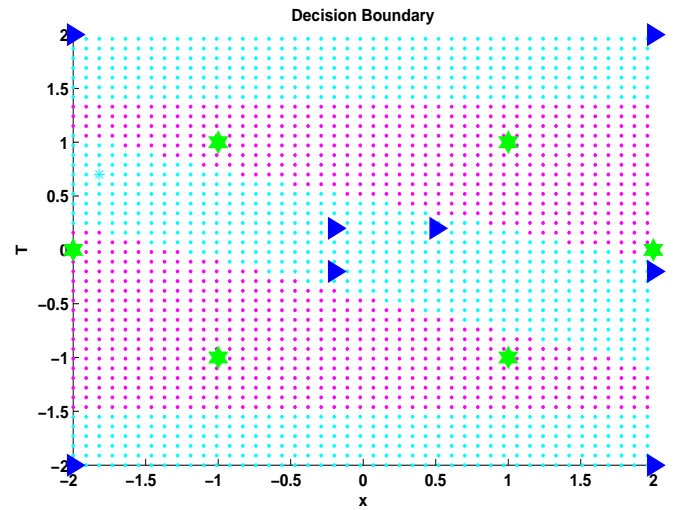


Figure 2: Decision boundary for 6-layer perceptron- after less than $10e5$ iterations converged-cross entropy case. Blue triangle are corresponding to class 1 and green stars with class 0 (cyan 1-area and magnet 0-area)

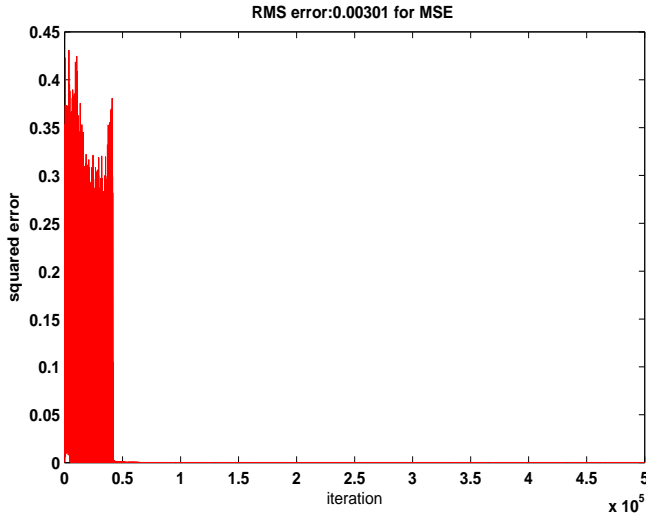


Figure 3: MSE error for 6-layer perceptron- after less than 5×10^5 iterations converged-MSE case.

the code for MSE cost with the same parameters and set number of layers into 6 .

Figure 3 represents the error vs iteration report with MSE cost function and Figure 4 is the corresponding decision boundary.

I also find the proper parameters and also proper changes in the code so that it can solve the 5-layer network. I provided the code at the end of the report. Figure 5 represents the error vs iteration and Figure 6 represents the corresponding decision boundary.

In order to complete the report I run the codes with smooth gradient (the codes are noted in the MATLAB file and also in the report) and I provided the results in the report. I tested different layers with these softer gradients and as we can see smoother gradient reduction in the error plot. Figure 7 represents running for 5-layer perceptron, and Figure 8 represents the corresponding decision boundary. Figure 9 represents the error report over 6-layer network.

5. Conclusion

I noticed that the patterns that created over running the code with two cost functions; MSE and Cross Entropy are different. This difference are sometimes significant and sometimes the corresponding patterns are similar. I successfully designed 5-layer network with MSE but doing this task with Cross Entropy was difficult and in the best case I successfully modified parameters for 6-layer network although I tried different steepest decent algorithms (Accelerated, Armijo and other controls). With MSE I the network learned slower comparing to the Cross Entropy and accuracy is also worse w.r.t Cross Ent. In conclusion for this problem MSE converges better but slower and with slightly higher error.

When I set parameters, Cross Entropy case works better

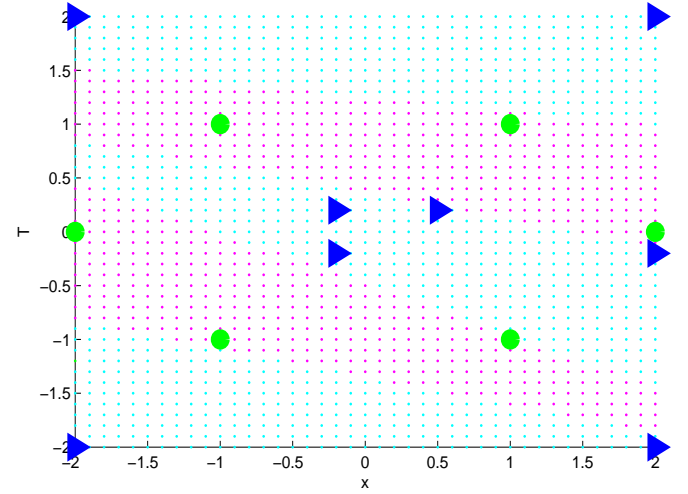


Figure 4: Decision boundary for 6-layer perceptron- after less than 10000 iterations converged-MSE case. Blue triangle are corresponding to class 1 and green stars with class 0 (cyan 1-area and magnet 0-area)

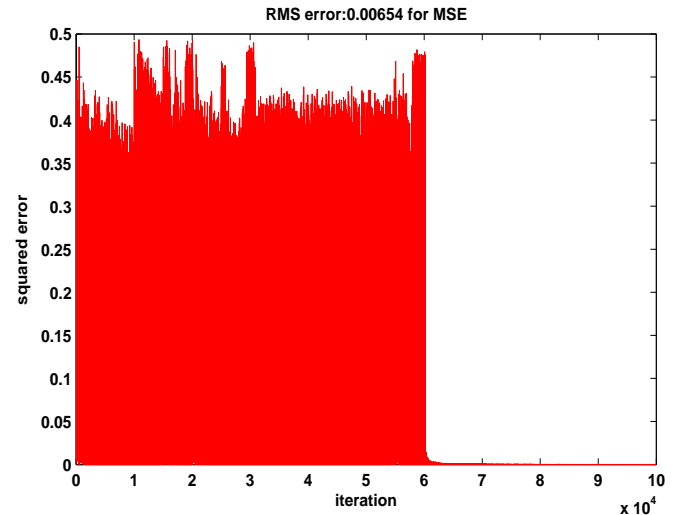


Figure 5: MSE error for 5-layer perceptron- after less than 10×10^4 iterations with MSE cost case

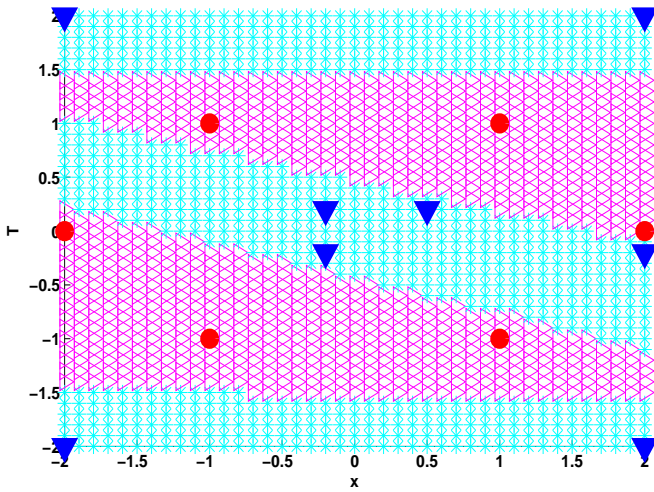


Figure 6: Decision boundary for 5-layer perceptron- after less than 10^5 iterations with MSE cost case. Blue circles belonging to class 0 denoted with magnet area, blue triangles belong to class 1 denoted with cyan area

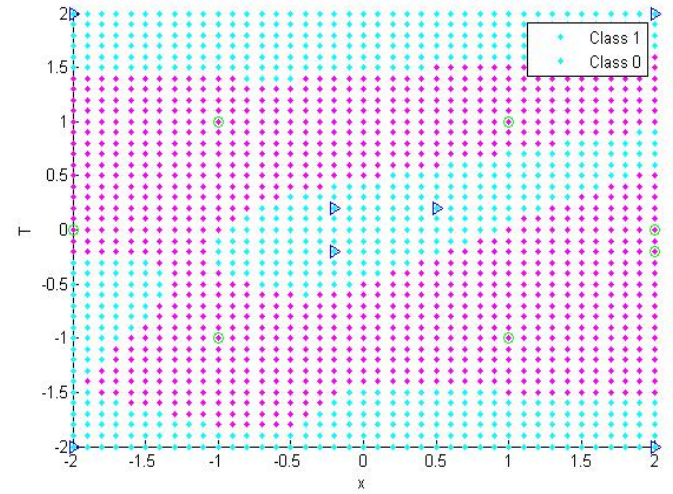


Figure 8: Decision boundary for 7-layer perceptron- after less than 10000 iterations not converged-cross entropy case

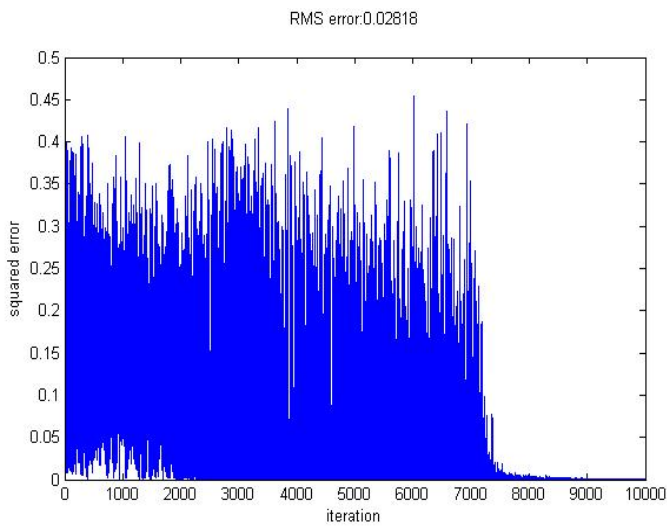


Figure 7: RMS error - 5 layer net, after less than 8000 iterations converged- MSE case

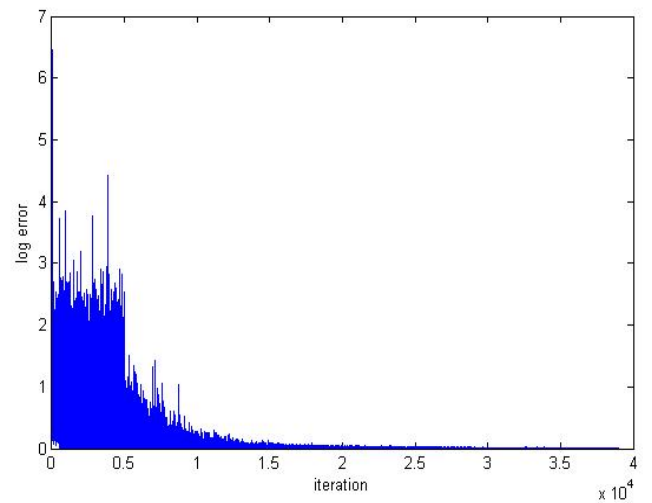


Figure 9: RMS error - 6 layer net, still fast converged- cross entropy case

than MSE in finding better local minimas. This fact is shown with the sum square error which is much less in the Cross Entropy case.

6. Final Codes

In the following I provide the MATLAB code that I used for 5-layer perceptron with MSE cost function.

6.1. MSE efficient code

I provided the whole source code, the main part and the functions. In this code I implemented Nestrov's accelerated gradient which responds well in this case.

```

1 function simpleBP2_hw
2 % ENTER YOUR NAME HERE:
3 % for 5 hidden layer
4 %%
5 clear all
6 clc
7 nX = 3; % input dimension (do not change)
8 nV = 5; % TUNE (number of hidden units)
9 eta = .87; % TUNE (learning rate)
10
11 x = zeros(1,nX); % input layer
12 V = zeros(1,nV); % hidden layer
13 O = 0; % output layer
14 %%
15 W = (rand(1,nV)-0.5)*0.05; % hidden->output
    weights
16 w = (rand(nV,nX)-0.5)*0.01; % input->hidden weights
17
18 % Do not change the training data
19 tr_data = [ -2 -2 1 1;
20             2 -2 1 1;
21             -2 2 1 1;
22             2 2 1 1;
23             -1 -1 1 0;
24             1 -1 1 0;
25             -1 1 1 0;
26             1 1 1 0;
27             -0.2 -0.2 1 1;
28             2 -0.2 1 1;
29             -0.2 0.2 1 1;
30             0.5 0.2 1 1;
31             2 0.0 1 0;
32             -2 0.0 1 0];
33 tr_x = tr_data(:,1:nX);
34 tr_T = tr_data(:,end);
35
36 [tr_m, tr_n] = size(tr_data);
37 figure(1); clf;
38
39 MAXIT = 1000000; % TUNE
40 % MAXIT = 50000; % TUNE
41 err = zeros(MAXIT,1);
42 %%
43 dw_old = (rand(nV,3)-0.5)
44 mu = .087
45 %%
46 for it = 1:MAXIT,
47
48     k = floor(rand*tr_m+1);
49     x = tr_x(k,:);
50     T = tr_T(k);
51     [O, Oh, V, Vh] = forwardPass(x,w,W);

```

```

53     E = 0.5*(T-O)^2; % SUM SQUARED ERROR
54 % E = -T.*log10(O) - (1-T).*log10(1-O); % SUM
    CROSS ENTROPY
55     err(it+1)=E;
56
57     if (mod(it,1000)==0) % show learnig progress
58         cla;
59         plot(0:it, err(1:it+1));
60         xlabel('iteration');
61         ylabel('squared error');
62         drawnow;
63     end;
64 % -----MSE
65     dW = -(T-O)*sigder(Oh)*V; % SUM SQUARED ERROR
    BASED dW
66 % -----Cross Ent
67 % dW = -(T*sigder(Oh)/O).*V + ((1-T)*sigder(Oh)
    /(1-O)).*V; % my Tuning - SUM CROSS ENTROPY
68 % ----- MSE
69     dw = -(T-O)*sigder(Oh)*(W.*sigder(Vh))'*x; % SUM
    SQUARED ERROR BASED dw
70 % -----Cross Ent
71 % dw = -(T*sigder(Oh)/O).(W.*sigder(Vh))'*x +
    ((1-T)*sigder(Oh)/(1-O)).*(W.*sigder(Vh))'*x; %
    My tuning
72     W = W - eta*dW; % SUM CROSS ENTROPY
73
74     W = W - eta*dW;
75 % w = w - eta* dw+0.005*(rand(nV,3)-0.5); %
    TUNE
76 % grw = beta*dw_old + dw
77 % dw1 = alpha*dw - mtm*dw
78 %
79 %%
80     if it==1
81         vect_p = -eta*dw;
82     end
83
84     %w = sqmul*grw; % 0.01*dw_old; %0.0005*(rand(nV
    ,3)-0.5); % TUNE
85
86 % Accelerated Gradient
87     vect = mu*vect_p - eta*dw; %
88     w = w - mu * vect_p + (1+mu)*vect; %
89     vect_p = vect;
90
91 % dw_old = dw;
92     if (mod(it,5000) == 0)
93         if E >= 1e-3
94             vect_p = -eta*dw;
95             eta = eta + .001;
96             mu = mu + .001;
97             W = (rand(1,nV)-0.5); % hidden->output
    weights
98             w = (rand(nV,nX)-0.5); % input->hidden
    weights
99         elseif (eta > .1 && mu > .1)
100             eta = eta - 0.01;
101             mu = mu - .01;
102         end
103     end
104
105 end
106 %%
107 plot((0:length(err)-1),err);
108 xlabel('iteration');
109 ylabel('squared error');
110 drawnow;
111
112 ms_err = 0;

```



```

114 for k=1:tr_m,
115     x = tr_x(k,:);
116     T = tr_T(k);
117     [O, Oh, V, Vh] = forwardPass(x,w,W);
118     ms_err = ms_err + (T-O)^2;
119
120     fprintf('Desired Output %2.2f == %2.2f (Network
121           Ouput)\n',T,O);
122 end
123 ms_err = ms_err/tr_m;
124 rms_err = ms_err^0.5;
125 fprintf('RMS error over data points:%3.5f\n',rms_err
126 );
127 title(sprintf('RMS error:%3.5f\n',rms_err));
128
129 %% ENIER YOUR code to show decision boundary here
130 % -----
131 xmin = -2; %min(tr_x(:,1));
132 xmax = 2;%max(tr_x(:,1));
133 ymin = -2; %min(tr_x(:,2));
134 ymax = 2; %max(tr_x(:,2));
135 figure;
136 hold on
137 dx=0.1;
138 dy=0.1;
139
140 for x1=xmin:dx:xmax
141     for y1=ymin:dy:ymax
142         activation = forwardPass([x1 y1 1],w,W);
143         if activation > 0.5
144             plot(x1,y1,'*c', 'markersize', 9)
145         else
146             plot(x1,y1,'>m', 'markersize', 9)
147         end
148     end
149 end
150 legend('Class 1', 'Class 0');
151
152 nn_outpt = zeros(size(tr_x,1),1);
153 for i=1:size(tr_x,1)
154     nn_outpt(i) = forwardPass(tr_x(i,:),w,W);
155 end
156
157 %class = zeros(n_pairs, size(x,2)+size(O,2))
158 class1 = tr_x(find(nn_outpt>=0.5),1:2);%(1:7,1:2)
159 class0 = tr_x(find(nn_outpt<0.5),1:2);%(8:14,1:2)
160 plot(class1(:,1),class1(:,2), 'b>');
161 plot(class0(:,1),class0(:,2), 'go');
162 % include legend
163 legend('Class 1', 'Class 0');
164 %legend();
165 % label the axes.
166 xlabel('x');
167 ylabel('T');
168 % -----
169
170 end % main
171 %%
172
173
174 function [O,Oh,V,Vh] = forwardPass(x,w,W)
175 Vh = x*w';
176 V = sigmoid(Vh);
177 V(end) = 1;
178 %V(end-1) = 0;
179 Oh = V*W';
180 O = sigmoid(Oh);
181 end
182

```

```

183 function ret = sigmoid(x)
184 B = 2;
185 ret = 1./(1+exp(-B*x));
186 end
187
188 function ret = sigder(x)
189 B = 2;
190 ret = B*sigmoid(x).*(1-sigmoid(x));
191 end

```

6.2. Cross entropy efficient code

In the following code works well for the cross entropy cost function. All parts are same unless the main part.

```

1 nX = 3;           % input dimension (do not change)
2 nV = 6;           % TUNE (number of hidden units)
3 eta = .271;       % TUNE (learning rate)
4
5 x = zeros(1,nX); % input layer
6 V = zeros(1,nV); % hidden layer
7 O = 0;           % output layer
8
9 W = (rand(1,nV)-0.5)*0.05; % hidden->output
    weights
10 w = (rand(nV,nX)-0.5)*0.01; % input->hidden weights
11 %%
12 % Do not change the training data
13 tr_data = [ -2    -2    1    1;
14             2     -2    1    1;
15             -2     2    1    1;
16             2     2    1    1;
17             -1    -1    1    0;
18             1     -1    1    0;
19             -1     1    1    0;
20             1     1    1    0;
21             -0.2  -0.2  1    1;
22             2     -0.2  1    1;
23             -0.2  0.2   1    1;
24             0.5   0.2   1    1;
25             2     0.0  1    0;
26             -2    0.0  1    0];
27 tr_x = tr_data(:,1:nX);
28 tr_T = tr_data(:,end);
29
30 [tr_m, tr_n] = size(tr_data);
31 figure(1); clf;
32
33 MAXIT = 500000;    % TUNE
34 % MAXIT = 50000;   % TUNE
35 err = zeros(MAXIT,1);
36 %%
37 dw_old = (rand(nV,3)-0.5);
38 MeanE = [];
39 beta = 0.0005;
40 mu = .00217;
41
42 % armijo line search parameters
43 delta = 0.5;
44 gamma = 1e-4;
45 alpha = eta;
46 %%
47 for it = 1:MAXIT
48
49     k = floor(rand*tr_m+1);
50     x = tr_x(k,:);
51     T = tr_T(k);
52     [O, Oh, V, Vh] = forwardPass(x,w,W);
53

```

```

54
55 E = 0.5*(T-O)^2;           % SUM SQUARED ERROR
56
57 err(it+1)=E;
58
59 if (mod(it,1000)==0) % show learnig progress
60     cla;
61     plot(0:it, err(1:it+1));
62     xlabel('iteration');
63     ylabel('squared error');
64     drawnow;
65 end;
66
67 dW = -(T-O)*sigder(Oh)*V; % SUM SQUARED ERROR
    BASED dW
68
69
70 dw = -(T-O)*sigder(Oh)*(W.*sigder(Vh))'*x; % SUM
    SQUARED ERROR BASED dw
71
72 W = W - eta*dW; % SUM CROSS ENITROPY
73
74 W = W - eta*dW;
75 %%
76 if it==1
77     vect_p = -eta*dw;
78 end
79 %w - sqmul*grw; % - 0.01*dw_old; %0.0005*(rand(nV
    ,3)-0.5); % TUNE
80 % Accelerated Gradient
81 vect = mu*vect_p - eta*dw ; %
82
83 w = w - mu * vect_p + (1+mu)*vect + beta *(rand(
    nV,3)-0.5) ; % + beta*randn(nV,3)
84
85 vect_p = vect;
86 % For analysing the state of convergence I take
    samples and based on
87 % tne average of the samples
88 if (mod(it,500) == 0)
89     MeanE = [MeanE E];
90 %     MeanE = MeanE_old + .9*(E-MeanE_old);%
91 %     MeanE_old = MeanE
92 end
93
94 % dw_old = dw;
95 if (mod(it,5000) == 0)
96     if mean(MeanE) >= 1e-1
97 %
98 if max(MeanE) < 1e-5
99     beta = beta*.0001;
100 elseif min(MeanE) > 1e-4
101     beta = beta*1.2;
102 end
103 end
104 MeanE=[];
105 end

```