

HOMEWORK 3 — L^AT_EX

Problem 1. XOR

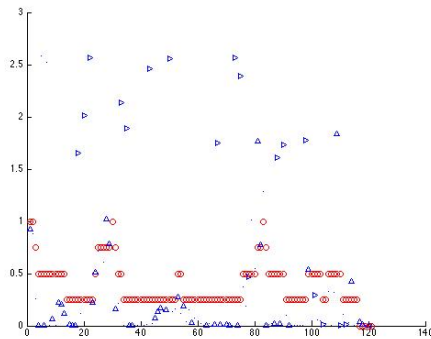
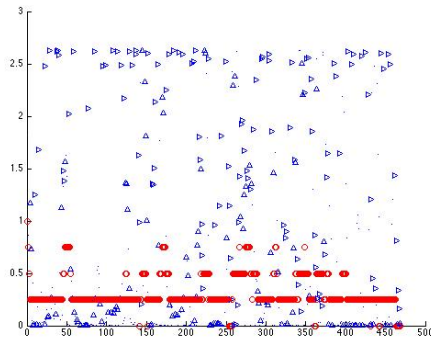
(a) The update rule for each layer is:

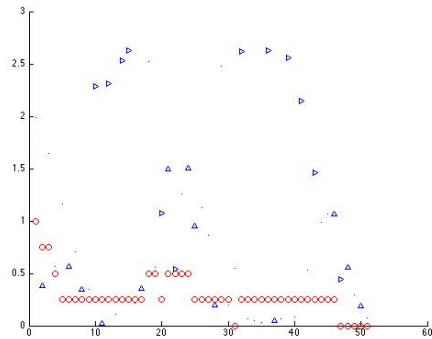
$$\Delta w_{i,j}(t) := \gamma \delta_j(t) a_i(t) + \alpha \Delta w_{i,j}(t-1)$$

$$w_{i,j} = w_{i,j} + \Delta w_{i,j}(t)$$

(b)

The red circles are the error rates. The blue arrows are the errors at each output; the direction of the arrows are used to differentiate between the four possible inputs that could have produced that output.





(d) Mean: 2.0602e+03

Standard deviation: 1.4958e+03

Variance: 2.2375e+06

Epochs: 1341, 4720, 835, 964, 777 , 2239, 1230, 2388, 1412, 4696

Second iteration (used different learning rates for each layer):

Epochs: 362, 1275, 438, 540, 1323, 560, 932, 1216, 1129, 1161

Mean: 893.6000

Variance: 1.4319e+05

Standard Deviation: 378.3993

After I fixed my broken tanh (thank you Akshat!): Epochs: 861, 50, 38, 44, 46, 33, 321, 56, 463, 1848

Mean: 376

Standard Deviation: 584.2655

Variance: 3.4137e+05

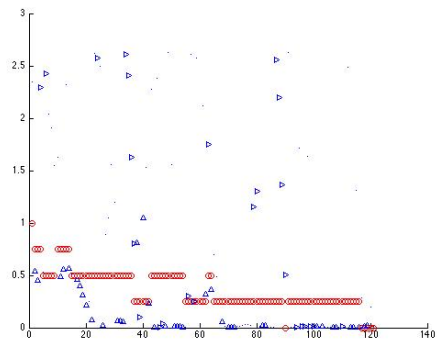
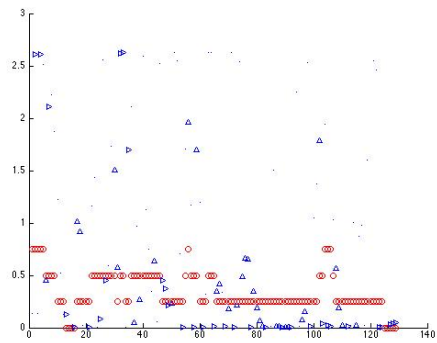
(e) In general, the number of epochs decreases as you increase the number of hidden units.

For 3 hidden units:

Epochs: 42, 58, 377, 62, 49, 118, 121, 129, 104, 58

Mean: 111.8000

Std: 98.8240

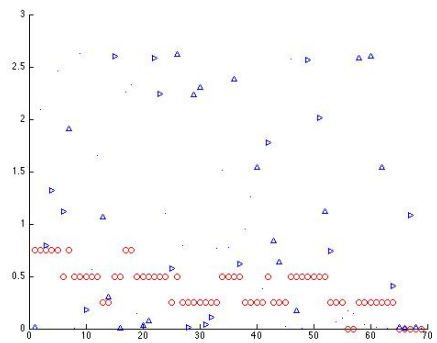


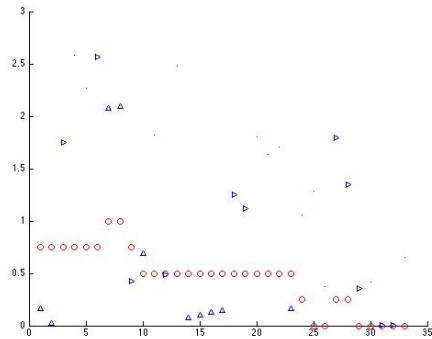
For 4 hidden units:

Epochs: 67, 94, 80, 85, 63, 73, 72, 67, 33, 69

Mean: 70.3000

Std: 16.1593





Problem 2. I converted the Y label vector to a 60,000 x 10 matrix so that I could get errors for each digit. I've included the average accuracy for each table.

I used two different tests to compare the output of the neural network to the label: First, a threshold value (if the probability was greater than 0.2, then the output becomes 1, otherwise it is a 0). Second, I chose the argument with the highest probability. These are the accuracies for each digit based on these techniques:

Accuracy Table: 100 hidden units, Sample 100 subsets at a time.

Theshold	Argmax
0.8279	0.6591
0.9334	0.7771
0.7914	0.8599
0.7593	0.8562
0.8066	0.8881
0.7957	0.9108
0.8727	0.9042
0.8912	0.9057
0.7643	0.9026
0.7956	0.8991

Accuracy Table: 50 hidden units, Sample 50 subsets at a time. Average accuracy: 0.8203, 0.8666.

Theshold	Argmax
0.8437	0.6798
0.9330	0.8111
0.7885	0.8502
0.8285	0.8817
0.8001	0.8951
0.7843	0.9107
0.8148	0.9042
0.8767	0.9311
0.7444	0.9026
0.7894	0.8991

Three examples of misclassified test instances and their corresponding true labels:

This table shows the counts of misclassifications: The rows indicate what the incorrect guesses were, the columns show what the incorrect guesses should have been.

	0	1	2	3	4	5	6	7	8	9
0	0	5	459	482	195	786	459	208	335	225
1	1	0	371	148	241	67	59	166	437	249
2	1	2	0	171	51	12	424	83	117	20
3	1	0	0	0	23	22	1	156	55	134
4	0	0	7	3	0	4	7	20	27	380
5	0	1	0	0	1	0	7	0	0	0
6	0	0	0	0	0	0	0	0	0	0
7	0	0	1	1	0	0	0	0	3	1
8	0	0	0	0	0	0	0	0	0	0
9	0	0	0	0	0	0	0	0	0	0

From the error rates and the table, it is clear that zero is the digit with the most errors. Zero is most commonly misclassified as six, one is most commonly misclassified as nine, and two is most commonly misclassified as seven.

This first example is a 5 misclassified as a 0:



This second example is a 7 misclassified as a 1:



The final example is a 9 misclassified as a 0:



Problem 3. Appendix

Listing 1: HW 3 Main file

```
1 % HW3 Main Script
2
```

```

3 % Problem 1: XOR
4 % Create a multi-layered neural net to implement XOR.
5 % User one hidden layer containing 2 units and an output layer containing
6 % 1 unit. Include a bias to each unit.
7 % Use backpropagation to train the weights.
8 % Your program should have two stopping criteria: it should stop when the
9 % error is below some threshold OR when a maximum number of epocks is
10 % reached, whatever comes first.
11
12
13 % PARAMETERS
14 N_inp = 2;
15 N_hid = 1;
16 N_out = 1;
17 N_bias = 1;
18 N_epochs = 5000;
19 numSuccess = 5;
20 alpha = 1;
21 momentum = 0.9;
22 K = 1; % Subset of training examples
23
24 % Create all combinations of inputs (plus bias)
25 X = ones(N_inp^2, N_inp + N_bias);
26 X(1:ceil(end/2), 1) = -1;
27 X(1:2:end, 2) = -1;
28 X(:,3) = 1;
29
30 Y = ones(1, 4);
31 Y([1, end]) = -1;
32
33 X_norm = zeros(size(X));
34 % Normalize the inputs and outputs
35 for j = 1:2 % exclude the bias
36     x_std = std(X(:,j));
37     x_mu = mean(X(:,j));
38     X_norm(:,j) = 1/x_std * (X(:,j) - x_mu);
39 end
40 X_norm(:,3) = 1;
41 disp('Starting_NN');
42
43
44 funny_tanh = @(x) 1.7159*tanh((2*x)/3) ;%+ twisting_slope*x;
45 funny_tanh_prime = @(x) 1.14393*(1-tanh((2*x)/3).^2) ;%+ twisting_slope*x;
46
47 SSE = @(t, x, func) sum((t-double(func(x))).^2,2)/2;
48 SSE_prime = @(t,x, func, func_prime) sum((t-double(func(x)))*double(func_prime(x))*x,

```

```

49
50 [W_l1, W_l2, ep] = backprop_train(N_inp, N_hid, N_out, N_epochs, X_norm, Y, ...
51     funny_tanh, funny_tanh_prime, numSuccess, alpha, momentum, K);
52
53 disp(ep);
54
55 % for p = 1:N_inp^2
56 %     A_h(:,p) = double(subs(funny_tanh, W_l1*X_norm(p,:)'));
57 %     A_o(:,p) = double(subs(funny_tanh, W_l2*[A_h(:,p); 1]));
58 %
59 % end
60
61 % tmp = [X(1,:) ', X(1,:) ', [A_h(:,1); 1]]';
62 % grad = ((Y(1) - A_o(1))*double(subs(diff(funny_tanh), W_l2*[A_h(:,p); 1])))*tmp;
63 % numgrad = (Y(1) - A_o(1)) * (Y(1) - A_o(1)) * (tmp - .0001);
64 %
65 % d = norm(grad-numgrad)/norm(grad+numgrad)
66
67 % grad_l1 = SSE_prime([Y; Y], W_l1*X(p,:) ', funny_tanh, diff(funny_tanh));
68 % numgrad_l1 = computeNumericalGradient(@SSE, [Y; Y], A_h, 1.0e-4)
69 % diff_l1 = norm(grad_l1-numgrad_l1)/norm(grad_l1+numgrad_l1)
70 %
71 % grad_l2 = double(diff(SSE(Y, A_o)))
72 % numgrad_l2 = computeNumericalGradient(SSE, Y, A_o, 1.0e-4)
73 % diff_l2 = norm(grad_l2-numgrad_l2)/norm(grad_l2+numgrad_l2)
74 %
75
76
77 % Problem 2: MNIST
78
79 % % Import dataset
80 load(' ./MNSIT_mats/training_images.mat ');
81 load(' ./MNSIT_mats/training_labels.mat ');
82 load(' ./MNSIT_mats/test_images.mat ');
83 load(' ./MNSIT_mats/test_labels.mat ');
84
85
86 % PARAMETERS
87 N_inp = 784;
88 N_hid = 50;
89 N_out = 10;
90 N_bias = 1;
91 N_epochs = 5000;
92 numSuccess = 5;
93 alpha = 1; % Global learning rate
94 momentum = 0.9;

```

```

95 twisting_slope = 0.00;
96 K = 100; % Subset of training examples
97 subset_size = 500;
98
99
100 funny_tanh = @(x) 1.7159*tanh((2*x)/3) ;%+ twisting_slope*x;
101 funny_tanh_prime = @(x) 1.14393*(1-tanh((2*x)/3).^2) ;%+ twisting_slope*x;
102
103
104 X_norm = [training_images , ones(size(training_images,1),1)];
105 X_subset = X_norm(1:subset_size,:);
106
107 Y = zeros(size(training_labels,1), 10);
108 for j = 1:size(training_labels,1)
109     Y(j,training_labels(j)+1) = 1;
110 end
111 Y_subset = Y(1:subset_size,:);
112
113
114 disp('Starting NN');
115 % [W_l1, W_l2, ep] = backprop_train_MNIST(N_inp, N_hid, N_out, N_epochs, X_norm, Y, ..
116 %     funny_tanh, funny_tanh_prime, numSuccess, alpha, momentum, K, true, true);
117
118 % TEST MNIST
119 X = [test_images , ones(size(test_images,1),1)];
120
121 Y = zeros(size(test_labels,1), 10);
122 for j = 1:size(test_labels,1)
123     Y(j,test_labels(j)+1) = 1;
124 end
125
126 thresh = 0.2;
127 testOut = zeros(N_out,size(test_images,1));
128 testOut2 = zeros(N_out,size(test_images,1));
129
130 for p = 1:size(test_images,1)
131     A_h = funny_tanh(W_l1*X(p,:));
132     A_o(:,p) = funny_tanh(W_l2*[A_h; 1]);
133     A_o(:,p) = exp(A_o(:,p))/sum(exp(A_o(:,p)));
134
135     for i = 1:size(A_o,1)
136         if(A_o(i,p) >= thresh)
137             testOut(i,p) = 1;
138         else
139             testOut(i,p) = 0;
140         end

```



```

141     end
142
143     testOut2(:,p) = 0;
144     [~, currArgmax] = max(A_o(:,p));
145     testOut2(currArgmax, p) = 1;
146
147
148
149 end
150 Yt = Y';
151 logical = testOut2 ~= Yt;
152 logical = logical - 2*Yt;
153
154 isVsShouldBe = zeros(10,10);
155
156 for i = 1:size(logical,2)
157     shouldBeInd = find(logical(:,i) == -1);
158     isInd = find(logical(:,i) == 1);
159     isVsShouldBe(isInd, shouldBeInd) = isVsShouldBe(isInd, shouldBeInd) + 1;
160 end
161
162 max(isVsShouldBe)
163
164 err_rate = sum(testOut ~= Y', 2)/size(Y,1);
165 err_rate2 = sum(testOut2 ~= Y', 2)/size(Y,1);
166 acc1 = sum(testOut == Y', 2)/size(Y,1);
167 acc2 = sum(testOut2 == Y', 2)/size(Y,1);
168
169
170 [err_rate, err_rate2]
171 [acc1, acc2]
172 mean([acc1, acc2], 1)

```

Listing 2: Backprop train

```

1 function [W_l1, W_l2, ep] = backprop_train(N_inp, N_hid, N_out, N_epochs, ...
2     X, Y, f, fprime, numSuccess, lrate, momentum, K)
3 %backprop_train Enter the size of the NN, the input data, error threshold,
4 % etc., this function will return the weights of the layers between
5 %the hidden layer and output layer.
6
7
8 % Create weight matrices
9 fan_in = 1/sqrt(N_inp+1);
10 W_l1 = fan_in * (-fan_in*rand(N_hid, N_inp+1)) + -fan_in;
11 W_l2 = fan_in * (-fan_in*rand(N_out, N_hid+1)) + -fan_in;
12

```

```

13 P = size(X,1);
14
15 % Create empty activation vectors
16 A_h = zeros(N_hid); %  $a^l = f(w^l * a^{l-1} + b^l)$ 
17 A_o = zeros(N_out, P);
18
19
20 err_new = 0;
21 ep = 0;
22
23 dW_l2_new = zeros(size(W_l2));
24 dW_l1_new = zeros(size(W_l1));
25
26
27 lrate_l2 = .01 * lrate; %0.01 * sqrt(N_hid+1); %sqrt(N_hid+1) * 0.5 * lrate
28 lrate_l1 = .5 * lrate; %0.1 * sqrt(N_hid+1); %sqrt(N_inp+1) * 0.99 * lrate
29 testOut = zeros(1,P);
30 figure; hold on
31 goodCount = 0;
32 while( goodCount < numSuccess && ep < N_epochs)
33
34     err_subtotal = 0;
35
36     for p = randsample(randperm(P), K)
37
38         % Forward pass
39         A_h = f(W_l1*X(p,:)');
40         A_o(:,p) = f(W_l2*[A_h; 1]);
41
42
43         testOut = sign(A_o);
44
45
46         % Calculate the output error with teaching signal
47         fpnet = fprime(W_l2*[A_h; 1]);
48         D_o = (Y(p) - A_o(:,p)).*fpnet;
49
50         err_subtotal = err_subtotal + sum(D_o.^2,1);
51
52
53         % Update weights in layer 2 with momentum
54         dW_l2_old = dW_l2_new;
55         for j = 1:N_out
56             dW_l2_new(j,:) = [lrate_l2*D_o(j).*[A_h; 1]]' + momentum*dW_l2_old(j,:);
57         end
58         W_l2 = W_l2 + dW_l2_new;

```

```

59
60     % Calculate hidden layer error
61     fpnet2 = [fprime(W_l1*X(p,:) ')] ' ;
62     D_h = (D_o'*W_l2(:,1:end-1)).*fpnet2;
63
64
65     % Change hidden layer weights
66     dW_l1_old = dW_l1_new;
67     for j = 1:N_hid
68         dW_l1_new(j,:) = lrate_l1*D_h(j)*X(p,:) + momentum*dW_l1_old(j,:);
69     end
70     W_l1 = W_l1 + dW_l1_new;
71
72 end
73
74 A_o;
75 err_rate = sum(testOut ~= Y)/size(Y,2)
76
77 if(err_rate == 0)
78     goodCount = goodCount + 1;
79 elseif(goodCount > 0 & err_rate ~= 0)
80     goodCount = goodCount - 1;
81 end
82
83 err_old = err_new;
84 err_new = err_subtotal;
85
86 ep = ep + 1
87
88 styles = ['b^', 'b>', 'bv', 'b<'];
89 plot(ep, err_rate, 'ro', ep, err_new, styles(mod(p,size(styles,2))+1));
90
91 end
92
93 if(goodCount >= numSuccess)
94     ep = ep - numSuccess;
95 end
96
97
98 end

```

Listing 3: Backprop train

```

1 function [W_l1, W_l2, ep] = backprop_train_MNIST(N_inp, N_hid, N_out, N_epochs, ...
2     X, Y, f, fprime, numSuccess, lrate, momentum, K, CE, SM)
3 %backprop_train Enter the size of the NN, the input data, error threshold,
4 % etc., this function will return the weights of the layers between

```

```

5  %the hidden layer and output layer.
6
7
8  % Create weight matrices
9  fan_in = 1/sqrt(N_inp+1);
10 W_l1 = fan_in - (-fan_in*rand(N_hid, N_inp+1)) + -fan_in;
11 W_l2 = fan_in - (-fan_in*rand(N_out, N_hid+1)) + -fan_in;
12
13 P = size(X,1);
14
15 % Create empty activation vectors
16 A_h = zeros(N_hid); %  $a^l = f(w^l * a^{l-1} + b^l)$ 
17 A_o = zeros(N_out, P);
18
19
20 err_new = 0;
21 ep = 0;
22
23 dW_l2_new = zeros(size(W_l2));
24 dW_l1_new = zeros(size(W_l1));
25
26 thresh = 0.1;
27
28 lrate_l2 = .01 * lrate; %0.01 * sqrt(N_hid+1); %sqrt(N_hid+1) * 0.5 * lrate
29 lrate_l1 = .1 * lrate; %0.1 * sqrt(N_hid+1); %sqrt(N_inp+1) * 0.99 * lrate
30 testOut = zeros(N_out,P);
31 testOut2 = zeros(N_out,P);
32 figure; hold on
33 goodCount = 0;
34 while( goodCount < numSuccess && ep < N_epochs)
35
36     err_subtotal = 0;
37
38     for p = randsample(randperm(P), K)
39
40         % Forward pass
41         A_h = f(W_l1*X(p,:) ');
42         A_o(:,p) = f(W_l2*[A_h; 1]);
43
44         if(SM == true)
45             A_o(:,p) = exp(A_o(:,p))/sum(exp(A_o(:,p)));
46         end
47
48         for i = 1:size(A_o,1)
49             if(A_o(i,p) >= thresh)
50                 testOut(i,p) = 1;

```

```

51         else
52             testOut(i,p) = 0;
53         end
54     end
55
56     testOut2(:,p) = 0;
57     [~, currArgmax] = max(A_o(:,p));
58     testOut2(currArgmax, p) = 1;
59
60
61     % Calculate the output error with teaching signal
62     if(CE == true)
63         fpnet = 1;
64     else
65         fpnet = fprime(W_l2*[A_h; 1]);
66     end
67
68     D_o = (Y(p,:) - A_o(:,p)).*fpnet;
69
70
71     err_subtotal = err_subtotal + sum(D_o,1);
72
73
74     % Update weights in layer 2 with momentum
75     dW_l2_old = dW_l2_new;
76     for j = 1:N_out
77         dW_l2_new(j,:) = [lrate_l2*D_o(j).*[A_h; 1]]' + momentum*dW_l2_old(j,:);
78     end
79
80     W_l2 = W_l2 + dW_l2_new;
81
82     % Calculate hidden layer error
83     if(CE == true)
84         fpnet2 = 1;
85     else
86         fpnet2 = [fprime(W_l1*X(p,:)')]';
87     end
88     D_h = (D_o'*W_l2(:,1:end-1)).*fpnet2;
89
90
91     % Change hidden layer weights
92     dW_l1_old = dW_l1_new;
93
94     for j = 1:N_hid
95         dW_l1_new(j,:) = lrate_l1*D_h(j)*X(p,:) + momentum*dW_l1_old(j,:);
96     end

```

```

97
98         W_l1 = W_l1 + dW_l1_new;
99
100     end
101
102     A_o;
103     err_rate = sum(testOut ~= Y',2)/size(Y,1);
104     err_rate2 = sum(testOut2 ~= Y',2)/size(Y,1);
105     [err_rate, err_rate2]
106
107     if(err_rate == 0)
108         goodCount = goodCount + 1;
109     elseif(goodCount > 0 & err_rate ~= 0)
110         goodCount = goodCount - 1;
111     end
112
113     err_old = err_new;
114     err_new = err_subtotal;
115
116     ep = ep + 1
117
118     %     styles = ['b^', 'b>', 'bv', 'b<'];
119     plot(ep, err_rate, 'r-', ep, err_rate2, 'b-');
120
121 end
122
123 if(goodCount >= numSuccess)
124     ep = ep - numSuccess;
125 end
126
127
128 end

```

Submitted by Alex Rosengarten on Feb 26, 2015.