# Homework 4 - Neural Network vs LogReg

Soheil Sadeghi Eshkevari

07-06-2018

## 1   Introduction

In this report, results of the last homework is going to be presented. In this project, we trained our machine to classify MNIST dataset, which is a well-known dataset of digit handwritings. For this purpose, both Logistic Regression and Neural Network were tried. For each, three different solvers, Stochastic Gradient Decent (SGD), SGD with Momentum and ADAM were exploited. Since all these optimization techniques are stochastic, various batch sizes were tried and results compared. In the next sections, results will be presented.

## 2   Logistic Regression Results

The first classifier that was tried was the Logistic Regression. This method is basically a binary classifier, thus to be able to train our dataset using it, we turned our labels to binary. For this purpose, label 5 is assigned to be +1 and other labels to be -1. Results are shown below:
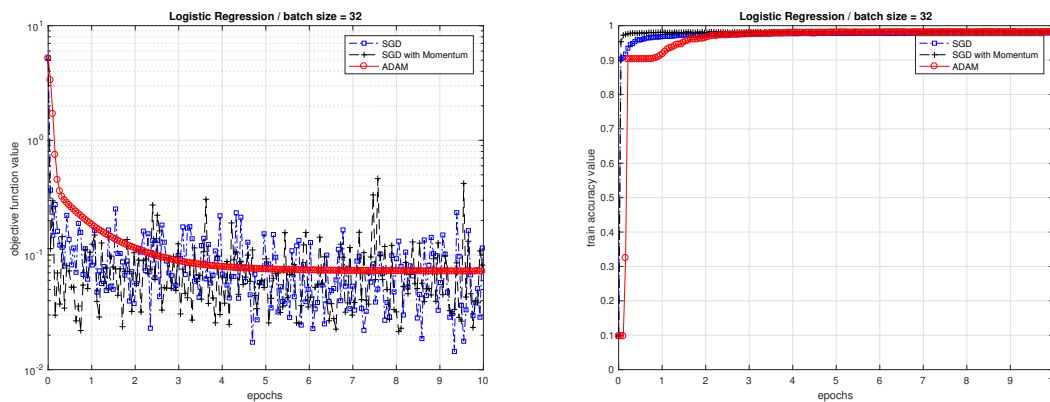


Figure 1: Performance plots vs. number of epochs with batch size 32

By comparing plots, following observations can be made. Less computational cost the problem needs when smaller batch sizes are set. In addition, with the setting I used for the training, ADAM is the most stable, while the least desirable one in terms of the fastest accuracy. Both SGD methods are very fast and suitable, however there show considerable fluctuations in their performance. This implies that the average of parameters can be a better alternative to get more stable results.
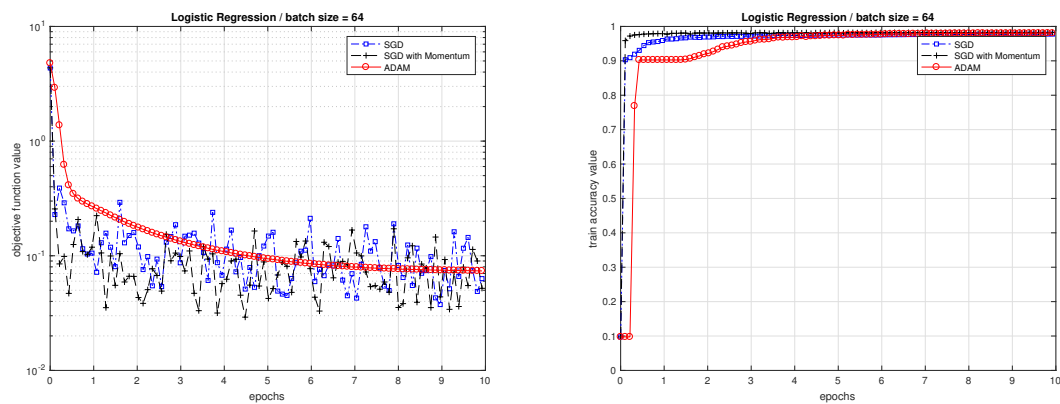
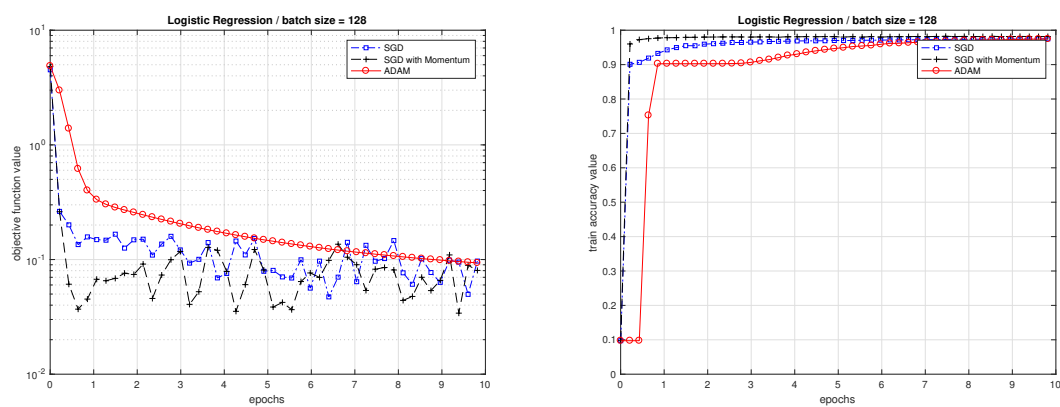Figure 2: Performance plots vs. number of epochs with batch size 64



Figure 3: Performance plots vs. number of epochs with batch size 128

# 3 Neural Network Results, Multi-Label Classification

The next set of results belongs to the multi-class classification of the MNIST data using a 2 layer Neural Network. Sigmoid activation functions were used to connect layers. At the output layer, cross entropy function turns values to probabilities for each label. The number of nodes in the hidden layer also is set to be 40, based on both my own tuning and the suggested number from Mohammad and Mertcan's course project. As before, three different batch sizes were analyzed and results are shown below:
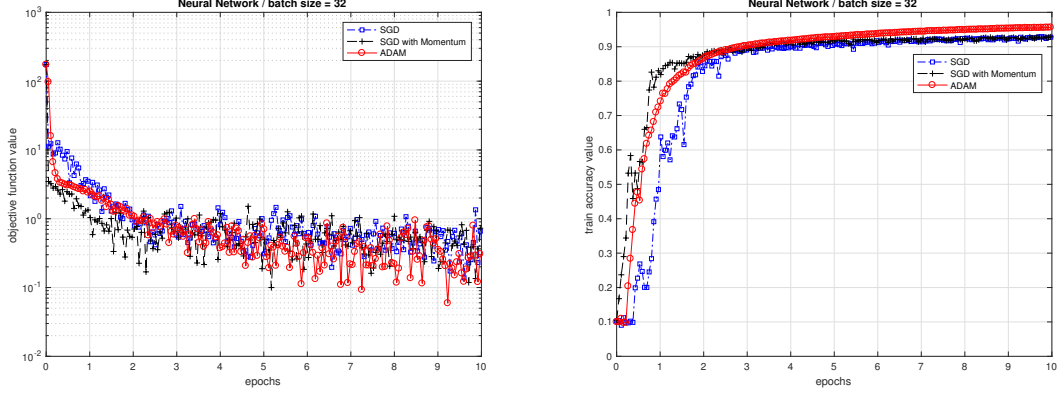


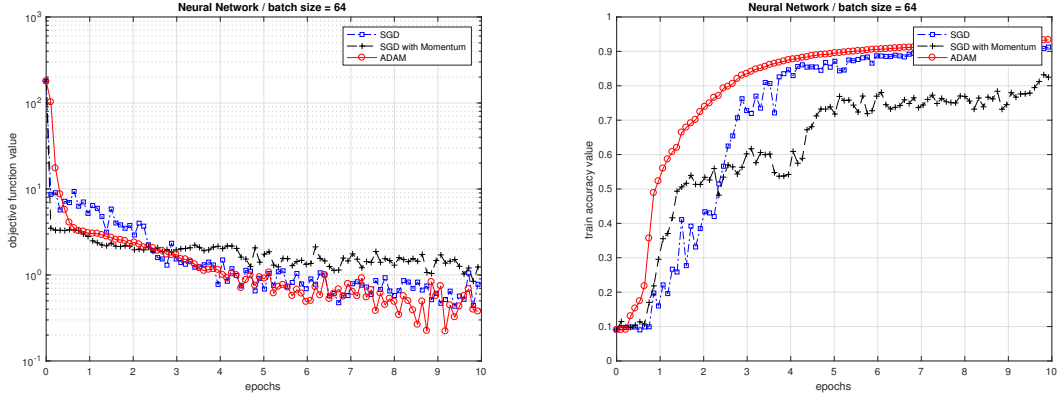Figure 4: Performance plots vs. number of epochs with batch size 32



Figure 5: Performance plots vs. number of epochs with batch size 64

The most challenging part for NN training by SGD and SGD with momentum was to tune step sizes in a way that they converge to accurate results. In order to find the best values, first a fixed step was given and the best value for that was found (around 0.8). Next, a decaying step size in the format of $\frac{a}{b+k}$ was tuned to vary within 0.9 to 0.7. The step size was successful in the majority of cases, however, as it is noticeable in Figure 6, SGD with momentum in case of batch size equal to 128 was not able to get converged and resulted inaccurate predictions. In addition, by comparing different batch sizes, it can be interpreted that ADAM is obviously the dominant solver for NN since it is both stable and accurate. It also has no sensitivity to the step size and default settings worked appropriately.
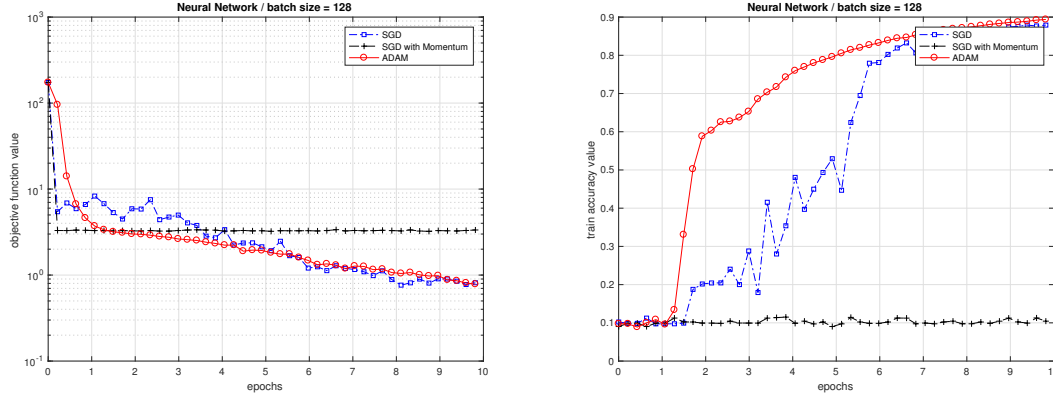
Figure 6: Performance plots vs. number of epochs with batch size 128

# 4    Fair comparison - Binary NN vs. Logistic Regression

As requested in the homework definition by Professor Scheinberg, in order to have a fair comparison between two methods, a binary classification problem with same logic explained for Logistic Regression was held using Neural Networks on MNIST dataset. The same analyses were done and results are presented hereafter:



Figure 7: Performance plots vs. number of epochs with batch size 32

By comparing results, it can be deduced that NN is also very successful for binary classification of MNIST. The best accuracy of NN is slightly higher than the best accuracy derived from Logistic Regression. However, as before, for higher batch sizes, less desired behavior can be observed by NN. The most desired setting between these two methods is to use NN with batch size equal to 32 and using ADAM or SGD with momentum as the solver.

4

Figure 8: Performance plots vs. number of epochs with batch size 64
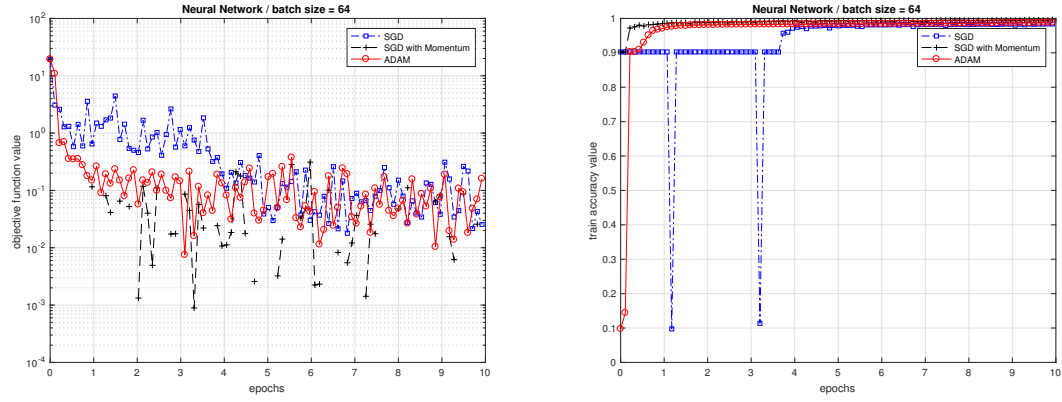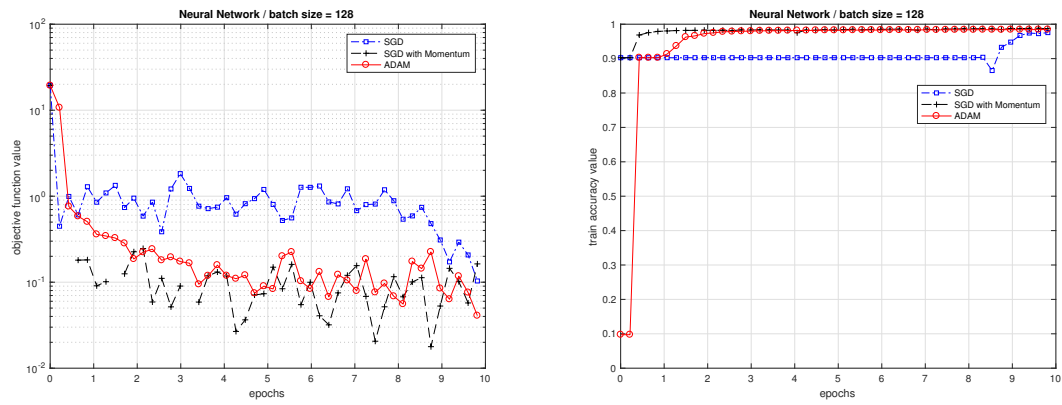


Figure 9: Performance plots vs. number of epochs with batch size 128

# 5 Code

In this section, the main body of my NN code for this homework is presented. However, subroutines and associated files will be sent in a package. Please note that a similar file exists for Logistic Regression and can be found in the package.

```matlab
1  % DATABASE 4: mnist.mat (always run this section in advance, for all for four algorithms)
2
3  clc;
4  clear;
5
6  rng shuffle
7
8  data_1 = load('new_data/mnist_mult.mat');
9
10 X_temp = data_1.TrainX;
11 Y = data_1.TrainY;
12 Xt_temp = data_1.TestX;
13 Y_t = data_1.TestY;
14
15 % fair comparison (Binary class NN) ————
16 Y_temp = zeros(size(Y,1),2);
17 Y_temp(:,1) = Y(:,5)==1;
18 Y_temp(:,2) = Y(:,5)==0;
19 Y = Y_temp;
20
21 Y_temp = zeros(size(Y_t,1),2);
22 Y_temp(:,1) = Y_t(:,5)==1;
23 Y_temp(:,2) = Y_t(:,5)==0;
24 Y_t = Y_temp;
25 Y_temp = [];
26 % ———————————————————————————
27
28 X = zeros(60000,28*28);
29 for i = 1:60000
30     temp = squeeze(X_temp(:,:,i));
31     X(i,:) = temp(:);
32 end
33
34 X_t = zeros(10000,28*28);
35 for i = 1:10000
36     temp = squeeze(Xt_temp(:,:,i));
37     X_t(i,:) = temp(:);
38 end
39
40 %% Scaling and parameters
41
42 for i = 1:size(X,1)
43     X(i,:) = X(i,:)/norm(X(i,:),2);
44 end
45 for i = 1:size(X_t,1)
46     X_t(i,:) = X_t(i,:)/norm(X_t(i,:),2);
47 end
48
49 input_layer_size  = 28*28;% 28x28 Input Images of Digits
50 hidden_layer_size = 40;   % 40 hidden units
51 num_labels = 2;           % 10 labels, from 1 to 10
52                           % (note that we have mapped "0" to label 10)
53 lambda = 0.0;
54
55 % w_1 = rand(hidden_layer_size, (input_layer_size + 1));
56 % w_2 = rand(num_labels, (hidden_layer_size + 1));
```

```matlab
57  w_1 = rand(hidden_layer_size, (input_layer_size));
58  w_2 = rand(num_labels, (hidden_layer_size));
59  nn_params = [w_1(:); w_2(:)];
60  % a_2 = 1./(1+X'*w_1);
61  % p = 1./(1+a_2*w_2);
62  % hs = -log(exp(p)./sum(exp(p),2));
63
64  batch_size = 128; % default = 50
65  epochs = 10;
66  iter = ceil(epochs*size(X,1)/batch_size);
67
68  %% solver (ADAM)
69
70  alfa = 0.002; % 0.002
71  beta_1 = 0.9;
72  beta_2 = 0.999;
73  eps = 1e-8;
74  m = zeros(numel(nn_params),1);
75  v = zeros(numel(nn_params),1);
76  t = 0;
77  thresh = 1e-2;
78
79  grad = 1;
80  J_tr_ADAM = zeros(iter,1);
81  J_ts_ADAM = zeros(iter,1);
82  nn_params_ADAM = zeros(iter+1,numel(nn_params));
83  nn_params_ADAM(1,:) = nn_params;
84
85  for j = 1:iter
86      t = t + 1;
87      Sk = randi(size(X,1),1,batch_size);
88      [J_tr_ADAM(j), grad] = nnCostFunction(nn_params_ADAM(j,:), ...
89                                    input_layer_size, ...
90                                    hidden_layer_size, ...
91                                    num_labels, ...
92                                    X(Sk,:), Y(Sk,:), lambda);
93      m = beta_1*m + (1-beta_1)*grad;
94      v = beta_2*v + (1-beta_2)*(grad.^2);
95      m_hat = m/(1-beta_1^t);
96      v_hat = v/(1-beta_2^t);
97      nn_params_ADAM(j+1,:) = nn_params_ADAM(j,:) - (alfa * m_hat./(sqrt(v_hat)+eps))';
98      J_ts_ADAM(j) = nnCost(nn_params_ADAM(j+1,:), ...
99                                    input_layer_size, ...
100                                   hidden_layer_size, ...
101                                   num_labels, ...
102                                   X_t, ...
103                                   Y_t, lambda);
104                              j
105 %      obj_fnc = J_tr_ADAM(j)
106 %      norm(grad,2)
107 end
108
109 % plot(1:20:iter,[J_tr_ADAM(1:20:end) J_ts_ADAM(1:20:end)]);
110
111 %% solver (SGD)
112
113 a = 28350;
114 b = 30500;
115 teta = 0.95;
116
117 k = 1:iter;
118 nu = a./(b+k);
119
120 J_tr_SGD = zeros(iter,1);
121 J_ts_SGD = zeros(iter,1);
```

```matlab
122  nn_params_SGD = zeros(iter+1,numel(nn_params));
123  nn_params_SGD(1,:) = nn_params;
124
125  v = zeros(length(nn_params),1);
126  v_prev = zeros(length(nn_params),1);
127
128  for j = 1:iter
129      Sk = randi(size(X,1),1,batch_size);
130      [J_tr_SGD(j), grad] = nnCostFunction(nn_params_SGD(j,:), ...
131                                    input_layer_size, ...
132                                    hidden_layer_size, ...
133                                    num_labels, ...
134                                    X(Sk,:), Y(Sk,:), lambda);
135      v = teta * v_prev + grad;
136      v_prev = v;
137      nn_params_SGD(j+1,:) = nn_params_SGD(j,:) - nu(j) * grad';
138  %      nn_params = nn_params - nu(j) * grad;
139  %      nn_params = nn_params - 0.1 * v;
140      J_ts_SGD(j) = nnCost(nn_params_SGD(j+1,:), ...
141                                    input_layer_size, ...
142                                    hidden_layer_size, ...
143                                    num_labels, ...
144                                    X_t, ...
145                                    Y_t, lambda);
146                               j
147  %      J_tr_SGD(j)
148  %      norm(grad,2)
149  end
150
151  % plot(J_ts_SGD(2:j)); hold on; plot(J_tr_SGD(2:j));
152
153  %% solver (SGD + Momentum)
154
155  a = 28350;
156  b = 30500;
157  teta = 0.95;
158
159  k = 1:iter;
160  nu = a./(b+k);
161
162  J_tr_SGDM = zeros(iter,1);
163  J_ts_SGDM = zeros(iter,1);
164  nn_params_SGDM = zeros(iter+1,numel(nn_params));
165  nn_params_SGDM(1,:) = nn_params;
166  v = zeros(length(nn_params),1);
167  v_prev = zeros(length(nn_params),1);
168
169  for j = 1:iter
170      Sk = randi(size(X,1),1,batch_size);
171      [J_tr_SGDM(j), grad] = nnCostFunction(nn_params_SGDM(j,:), ...
172                                    input_layer_size, ...
173                                    hidden_layer_size, ...
174                                    num_labels, ...
175                                    X(Sk,:), Y(Sk,:), lambda);
176      v = teta * v_prev + grad;
177      v_prev = v;
178  %      nn_params = nn_params - nu(j) * grad;
179  %      nn_params = nn_params - nu(j) * grad;
180      nn_params_SGDM(j+1,:) = nn_params_SGDM(j,:) - nu(j) * v';
181      J_ts_SGDM(j) = nnCost(nn_params_SGDM(j+1,:), ...
182                                    input_layer_size, ...
183                                    hidden_layer_size, ...
184                                    num_labels, ...
185                                    X_t, ...
186                                    Y_t, lambda);
```

```matlab
187                                     j
188 %        J_tr_SGDM(j)
189 %        norm(grad,2)
190 end
191
192 % plot(J_ts_SGDM(2:j)); hold on; plot(J_tr_SGDM(2:j));
193
194 %% plots
195 step = 100;
196 % function in three methods
197 semilogy((1:step:iter)*batch_size/size(X,1),J_tr_SGD(1:step:end), '-.bs'); hold on
198 semilogy((1:step:iter)*batch_size/size(X,1),J_tr_SGDM(1:step:end), '--k+'); hold on
199 semilogy((1:step:iter)*batch_size/size(X,1),J_tr_ADAM(1:step:end), '-ro');
200 grid on
201 legend('SGD','SGD with Momentum','ADAM');
202 xlabel('epochs');
203 ylabel('objective function value');
204 title(['Neural Network / batch size = ',num2str(batch_size)])
205
206 % accuracy in three methods
207 y_vec = zeros(size(Y,1),1);
208 for i = 1:size(Y,1)
209     y_vec(i) = find(Y(i,:)==1);
210 end
211
212 pred_SGD = zeros(1,length(1:step:iter));
213 pred_SGDM = zeros(1,length(1:step:iter));
214 pred_ADAM = zeros(1,length(1:step:iter));
215 for i = 1:step:iter
216     pred_SGD(i) = mean(double(predict(nn_params_SGD(i+1,:), ...
217         X,hidden_layer_size,input_layer_size,num_labels) == y_vec));
217     pred_SGDM(i) = mean(double(predict(nn_params_SGDM(i+1,:), ...
        X,hidden_layer_size,input_layer_size,num_labels) == y_vec));
218     pred_ADAM(i) = mean(double(predict(nn_params_ADAM(i+1,:), ...
        X,hidden_layer_size,input_layer_size,num_labels) == y_vec));
219 end
220 figure;
221 plot((1:step:iter)*batch_size/size(X,1),pred_SGD(1:step:end), '-.bs'); hold on
222 plot((1:step:iter)*batch_size/size(X,1),pred_SGDM(1:step:end), '--k+'); hold on
223 plot((1:step:iter)*batch_size/size(X,1),pred_ADAM(1:step:end), '-ro');
224 grid on
225 legend('SGD','SGD with Momentum','ADAM');
226 xlabel('epochs');
227 ylabel('train accuracy value');
228 title(['Neural Network / batch size = ',num2str(batch_size)])
```

# 6   Conclusions

In this homework, the main catch was to have a hands-on experience on coding some *off the shelf* methods and solvers and try to play with parameters to get a deeper understanding about these, rather than just superficially use them as blackbox toolsuites.