

# PROJECT 1: HAND WRITTEN DIGITS RECOGNITION PROGRAM

ALI NEHRANI

**ABSTRACT.** In this project, I implemented a neural network technique in order to recognize hand written digits. I used simple multiple layer neural network to address this problem and tested it with min square error and cross entropy objective functions. Data hand written digits data is very important for training the network so I created more data with supported hand written digit creation function. In network I used around 20 to 30 hidden layers. Results are reported through figures.

## 1. PREPROCESSING

There are some data provided for this project with special structure  $1 \times 200$ . This amount of data is not enough to get 90% accuracy. This amount of accuracy requires more proper data. For generating more data provided Matlab function `getUserTraj.m` is applied. I write a function `DatasetMake.m` which request user to draw a digit at the given area (with mouse) then it added that to the dataset. Using this function I created 4300 data. In the following I will discuss the structure and the coverage of the created data.

Figure 1 presents the process of preparing training data. The training data is saved as a Matlab data structure `Data09.mat` and provided along the programs.

The structure of the `Data09` is a matrix  $4300 \times 201$  where rows are training data for all digits 0 – 9 and first 100 columns are the training datum's  $x$ -axis and columns 101 – 200 are its  $y$ -axis and column 201 is the target. The following simple Matlab function creates training data set:

```
1 dt=0.01
2 fig = figure(1)
3 if ~ exist('DataN');
4     DataN = [];
5 end
6 target = input('enter the target you want to collect data: ');
7 for i=1:189
8     [x,y] = getUserTraj(dt,fig);
9     DataN = [DataN; [x y] target];
10    y = input('Do you want to continue? ');
11    if y == 'n'
12        break
13    end
14 end
```

---

*Date:* Dec. 2017.

Thanks for your time to read my report.

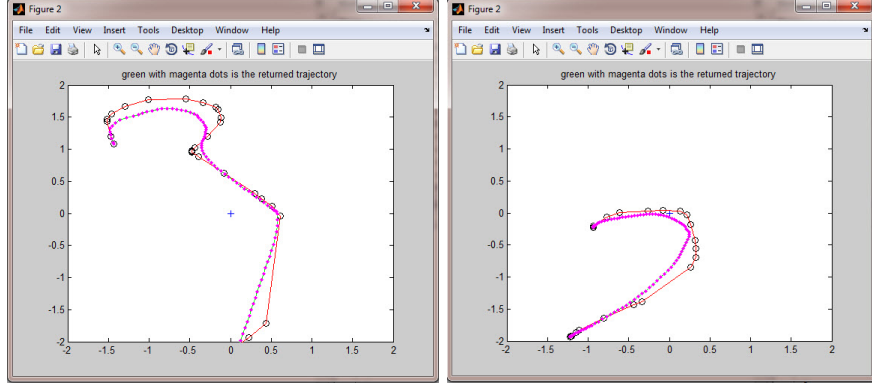


FIGURE 1. Building training data for hand written digits; Training digit "3" left, Training digit "7" right

FIGURE 2. Testing hand written data for scale adaptation - digit "6"

**1.1. Training Data.** Data is very important in increasing the accuracy of the method. The network should be scale and rotation invariant as well as ability to recognize digits that not recognizable by humans somehow (up to 90%). Theoretically there is a solution for scaling by making all data standard. This standardization is also supported by providing scaled data (scaled inputs) as it is shown in Figure 2. The following Matlab function is implemented for making data standard:

```

1 function [SData, mean_X, std_X] = standardize(varargin)
2 switch nargin
3     case 1
4         mean_X = mean(varargin{1});
5         std_X = std(varargin{1});
6
7         SData = varargin{1} - repmat(mean_X, [size(varargin{1}, 1) 1]);
8
9         for i = 1:size(SData, 2)
10             SData(:, i) = SData(:, i) / std(SData(:, i));
11         end
12     case 3
13         mean_X = varargin{2};
14         std_X = varargin{3};
15         SData = varargin{1} - repmat(mean_X, [size(varargin{1}, 1) 1]);
16         for i = 1:size(SData, 2)
17             SData(:, i) = SData(:, i) / std_X(:, i);
18         end
19 end

```

For Rotation, I used two tricks; first I applied mathematical rotation technique and rotated the data randomly and added those rotated data to the existing data (by adding these data the total amount of data became twice and also more). The Matlab function `rand_rotation.m` does the rotation

```

1 function r_data = rand_rotation(Data, rtheta)
2     Data_x = Data(:,1:100);
3     Data_y = Data(:,101:200);
4
5     sinus = sin(rtheta*pi/180);
6     cosinus = cos(rtheta*pi/180);
7     r_data = [Data_x*cosinus + Data_y*sinus - Data_x*sinus + Data_y*cosinus
8               ...
9               Data_x * cosinus + Data_y*sinus - Data_x*sinus + Data_y *
               cosinus];
9 end

```

The other technique that I applied to improve the ability of recognizing rotation is to providing rotated hand-written digits. This approach is also useful in training network to recognize bad written digits. In Figure 3 it was shown that what kind of data is feeded to the network.

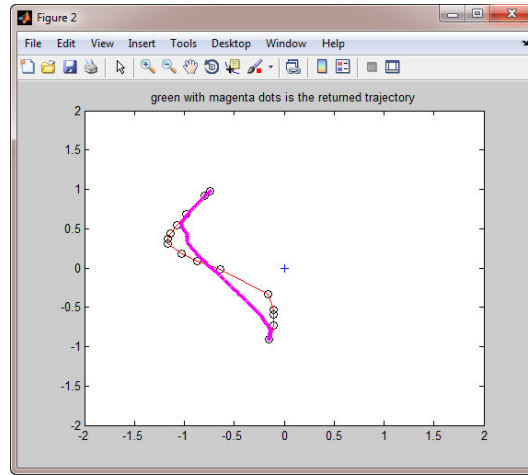


FIGURE 3. Training hand written data for scale adaptation - digit "7"

**1.2. Testing Data.** For creating test data I prepared a Matlab function `TestData.m` which gets the number of testing data and for that number; first it asks for the target (the number that use wants to write), second loads the `getUserTraj.m` for creating handwritten digit by the user, then adds this data to the testing data set and predicts the handwritten digit, compares it with the given target and returns the accuracy percentage. Figure 4 shows a process of creating test data and the output of the program.

I created 4300 data in order to cover any kind of rotation and bad writing cases. The codes are provided in the following:

```

1 acprt=0;
2 if ~ exist('Tdata')
3     Tdata = [];
4 end
5 for j = 1:n_test
6     y_true = input('Enter true y: ')

```

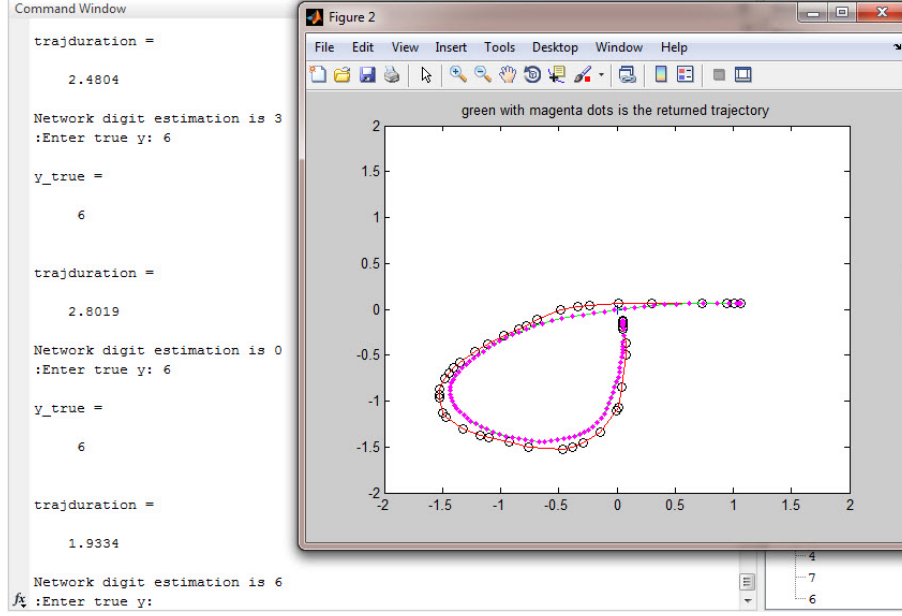


FIGURE 4. Testing hand written data and the prediction result of the program- digit "6"

```

7 [x,y] = getUserTraj(dt, fig);
8 test_data = [x y];
9 Tdata=[Tdata; x y y_true];
10 [O,Oh,V,Vh] = forwardPass(test_data,w,W);
11 [o_m,O_i]=max(O);
12 y_pred = O_i-1;
13 fprintf('Network digit estimation is %d\n:', y_pred);
14
15 if y_pred == y_true
16     acprt = acprt + 1;
17 end
18 end
19 fprintf('Testing accuracy equals to: %f\n',acprt*100/j );

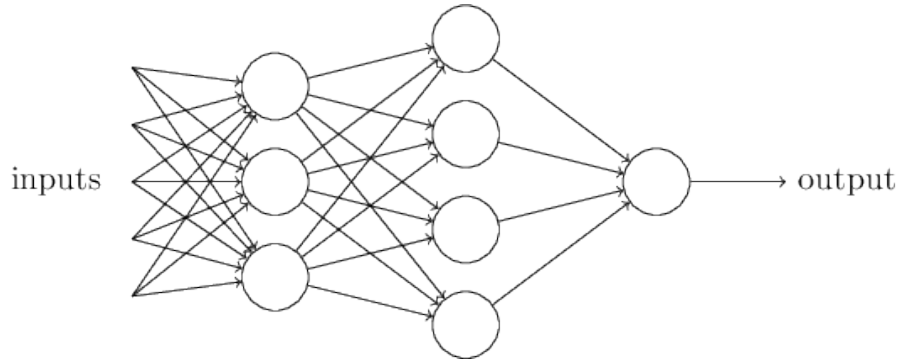
```

In applying Cross-validation I separate 10% of the data for testing process. the cross validation results are provided in Table 1.

## 2. NEURAL NETWORK

The structure Neural Network was the same as the provided multi layer perceptron as shown in the Figure 5. It consists of three layers (like all other networks) and the hidden layer itself has two layer. For training the hand written digits I applied 15 unites for the hidden layers. The input has 200 units and output has 10 units. For this there is two weight matrices; input to hidden layer  $w$  with size  $200 \times$  and hidden to output layer weight  $W$  with size  $30 \times 10$ .

In the final part **Softmax** layer is included to improve the recognition. The softmax

FIGURE 5. Structure of the Applied Network:  $200 \times 15 \times 15 \times 10$ 

is formulated as:

$$O_i = \frac{\exp l_i}{\sum_i \exp l_i}$$

where  $l_i$  is the output of the *hidden-out* layer.

**2.1. Hyper Parameters of the Neural Network.** Neural networks have some hyperparameters which should be set before training. There are strategies for some of them and some other are set based on experience. In applied neural network there are some hyper parameters in learning, perceptrons and also type of cost function. I listed the hyperparameters with their value in the following:

Number of hidden layers:	25
Maximum iteration:	2000000
learning rate eta ( $\eta$ ):	0.000013
slack learning rate nu ( $\nu$ ):	0.0000107
Sigmoid parameter B:	1.8
Leaky RELU parameter $\alpha$ :	0.15
Random noise coefficient beta ( $\beta$ ):	$10^{-10}$

I applied Nesterov version of the accelerated gradient descent method which itself has two hyperparameters for the learning process which has two hyper parameters  $\eta$  and  $\nu$ .  $\eta$  weights the gradient effect on the updating of the new weight and  $\nu$  weights the guiding vector. I studied the effect of some parameters in learning.

**2.2. Activation Functions.** One of the most important part of the neural networks is activation functions. I applied three different activation functions in this project. Sigmoid activation function which is represented in Figure ?? with its derivative. Sigmoid function is represented as

$$y = \frac{1}{1 + \exp(-z)}, \quad y' = y(1 - y)$$

The Matlab implementation of these function is provided in the following

```

1 % activation function
2 function ret = sigmoid(x)
3 B = .8;
4 ret = 1./(1+exp(-B*x));
5 end

```

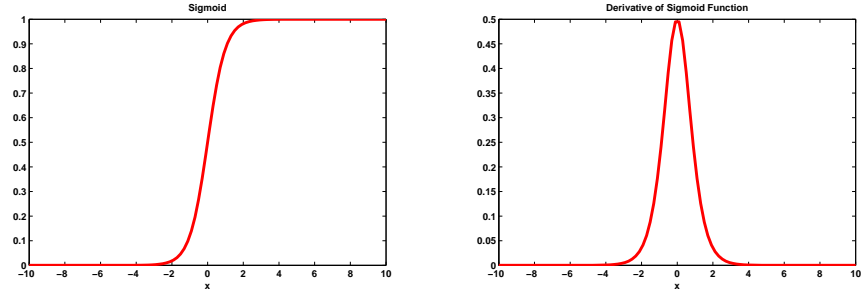


FIGURE 6. Sigmoid activation function (left) and it's derivative (right)

```

6 % derivative of sigmoid function
7 function ret = sigder(x)
8 B = .8;
9 ret = B*sigmoid(x).*(1-sigmoid(x));
10 end

```

The next activation function that I applied is Leaky RELU which is represented with it's derivative in Figure ???. The mathematical formula of this activation function is

$$y = x \text{ if } x \geq 0, \quad y = \alpha x \text{ if } x < 0,$$

The Matlab implementation of these function is provided in the following

```

1 % leaky relu
2 function ret =lrel(x)
3 ret = x;
4 alpha = 0.19;
5 ret(find(x<0)) = alpha .* x(find(x<0));
6 end
7 % leaky relu der
8 function ret = lrelder(x)
9 ret = ones(size(x));
10 alpha = 0.19;
11 ret(find(x<0)) = alpha;
12 end

```

The next activation function that I applied is x-Sigmoid function which is represented in Figure ???. This activation function is obtained by multiplying  $x$  to the sigmoid function to prevent ints dying effects. x-sigmoid is similar to the Leaky RELU. The Matlab implementation of these function is provided in the following

```

1 %x sigmoid
2 function ret = sigmoid(x)
3 B = 1.8;
4 ret = x./(1+exp(-B*x));
5 end
6 % x sig derivative
7 function ret = xsigder(x)
8 B = 1.8;
9 ret = x.*(B*sigmoid(x).*(1-sigmoid(x))) + 1./(1+exp(-B*x));
10 end

```

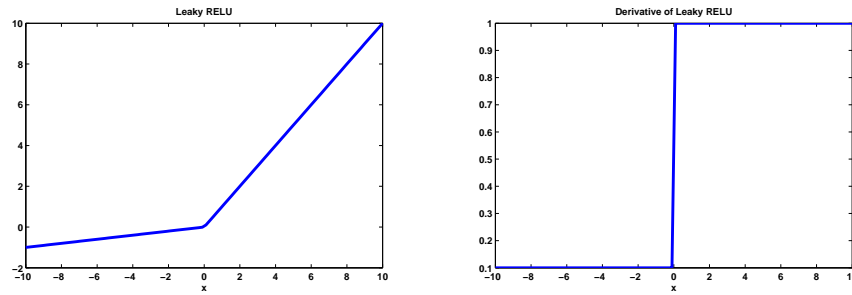


FIGURE 7. Leaky RELU activation function (left) and it's derivative (right) with parameter  $\alpha = 0.2$

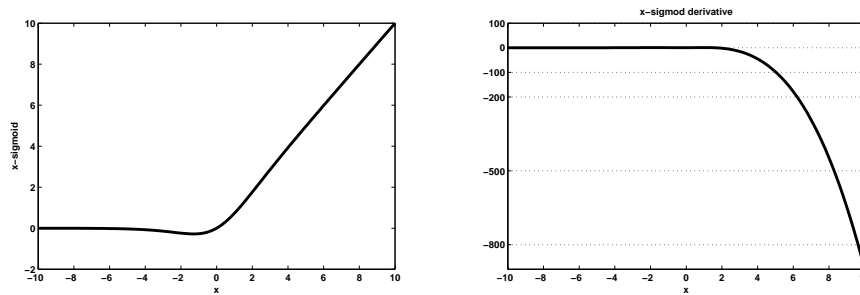


FIGURE 8. x-Sigmoid activation function (left) and it's derivative (right) with parameter  $\alpha = 0.2$

### 3. LEARNING WITH GRADIENT DECENT

For learning, we should consider to select cost function and optimization method. I introduced two type of cost functions: Cross entropy objective function and mean square error. Each of these cost function has propertiese and effective, for some problems cross entropy works better and for others mean square error. Mathematical formula of MSE is

$$MSE = \frac{1}{2} \sum_{i=1}^N (y_i - \hat{y}_i)^2 = \frac{1}{2} \|y - \hat{y}\|_2^2$$

where  $y$  is target and  $\hat{y}$  is prediction. In Figure 9 I plotted the responce of MSE for two random vectors. Matlab implementation of MSE and its derivative is in the following code box

```
1 if strcmp(cost, 'mse')
2     E = 0.5*norm(T-O,2)^2;
3     dW = -((T-O).*sigder(Oh))*V;
4     dw = -((((T-O).*sigder(Oh))'*W).*sigder(Vh))*x;
```

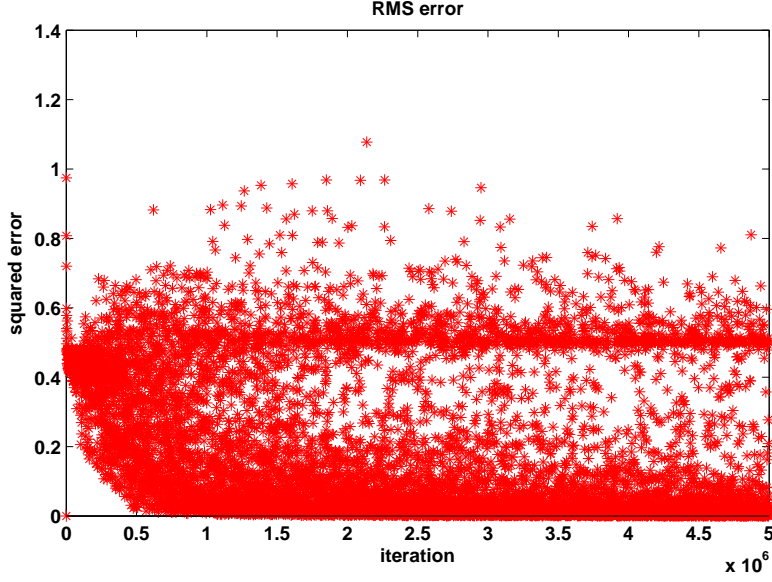


FIGURE 9. MSE function behaviour for random vectors uniformly distributed in  $[0\ 1]$

Cross Entropy cost function is also one good measure that alternatively is used in data processing and machine learning. Mathematical formula of cross entropy is

$$CE = - \sum_{i=1}^N y_i \log(\hat{y}_i)$$

where  $N$  here is number of classes. The behaviour of cross entropy cost function for two random vectors and two classes is presented in Figure 10. Matlab implementation of cross entropy for two classes case and its derivative is in the following code box

```

1 % Cross validation objective function
2 E = -sum(T.*log10(O) + (1-T).*log10(1-O));
3 dW = -(1/(log2(10)))*((T.*sigder(Oh)./O)*V + ((1-T).*sigder(Oh)./(1-O))*V);
   % my Tuning
4 dw = -(1/(log2(10)))*(((T.*sigder(Oh)./O)'*W).*sigder(Vh))'*x + (1/(log2
   (10)))*(((1-T).*sigder(Oh)./(1-O))'*W).*sigder(Vh))'*x;
```

**3.1. Optimization and Backward propagation.** My optimization algorithm which Nestrov's accelerated gradient decent is implemented as:

```

1 if it==1
2     vect_p = -eta*dw;
3 end;
4 % Accelerated Gradient
5 vect = mu*vect_p - eta*dw + beta*rand(nV, tr_n-1); %
6 w = w - mu * vect_p + (1+mu)*vect;
7 vect_p = vect;
```



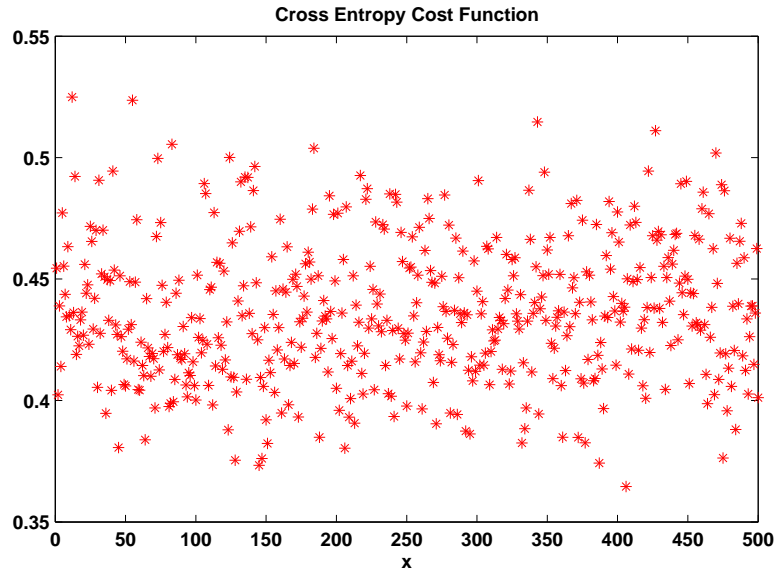


FIGURE 10. Cross entropy cost function behaviour for random vectors uniformly distributed in  $[0\ 1]$

In this implementation, `eta` and `nu` are hyperparameters, `vect` and `vect_p` are guiding vectors and `nV` is number of hidden layers. As it is well known, Gradient decent based methods most of the times fall into local minima and do not converge for large data. Without convergence we can not reach to proper accuracy (especially the accuracy that is requested). To resolve this problem I used rerunning trick which worked well.

In this trick after for example  $10^6$  iteration which no convergency is achieved, I saved the weights so that in the next run I used those trained weights. I did this many times until I reached very small errors.

**3.2. Feedforward.** Feedforward process which is the same as lms code. In the feedforward we apply computed weights to predict the target.

```

1 % forward pass of neural network
2 function [O,Oh,V,Vh] = forwardPass(x,w,W)
3 Vh = x*w';
4 V = sigmoid(Vh);
5     V(end) = 1;
6 Oh = W*V';
7 O = sigmoid(Oh);
8 end

```

#### 4. EVALUATING WITH CROSS VALIDATION

In order to evaluate the network properly, we should apply machine learning tools, at least the techniques that we learned from the course. Cross validation helps to evaluate the performance of the network. I applied 10 fold Cross validation.

- (1) Data is divided into 10 parts (randomly saved in different mat files)

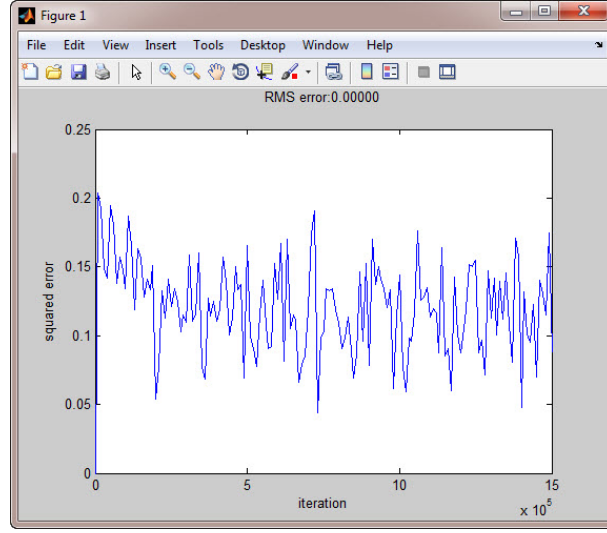


FIGURE 11. Error for multiple training- several times of training with 4300 data

(2) Each time 9 parts are applied for training and 1 is used for testing

In the following paper I provided accuracy results for each of 10 times train-test process:

Table 1: Cross-validation results with optimal number of layers and many iterations

Training	1	2	3	4	5	6	7	8	9	10
Achieved accuracy	95%	92%	85%	80%	91%	96%	89%	80%	87%	94%

How ever the program is written so that manuly I can change number of folds. Results of Table 1 the results are shown for more that  $100 \times 10^6$  and using 25 hidden layers (which is optimal number of hidden layers according to cross validation). The out put of the program is

```

1 Error in fold 1
2 is equal to 95.220141 percent
3 cost =mse
4 Error in fold 2
5 is equal to 92.580796 percent
6 cost =mse
7 Error in fold 3
8 is equal to 84.840749 percent
9 cost =mse
10 Error in fold 4
11 is equal to 80.435597 percent
12 cost =mse
13 Error in fold 5
14 is equal to 91.346604 percent
15 cost = mse
16 Error in fold 6
17 is equal to 96.409836 percent
18 cost =mse
19 Error in fold 7

```

```

20 is equal to 89.960187 percent
21 cost = mse
22 Error in fold 8
23 is equal to 80.814988 percent
24 cost = mse
25 Error in fold 9
26 is equal to 87.112412 percent
27 cost =mse
28 Error in fold 10
29 is equal to 94.346604 percent
30 Average Error in 10 folds is: 89.3068

```

In the report I reduced number of iterations into 10000 in order to reduce running time. In this amount of iterations I can not reach acceptable accuracy, but in order to demonstrate how cross validation works for evaluation the model I applied this low amount of iterations. I applied cross validation in order to obtain optimal amount of hidden layers. The Results of cross validation with small number of iterations (10000) is shown in Figure 12. According to this figure optimal number of hidden layers are about 25.

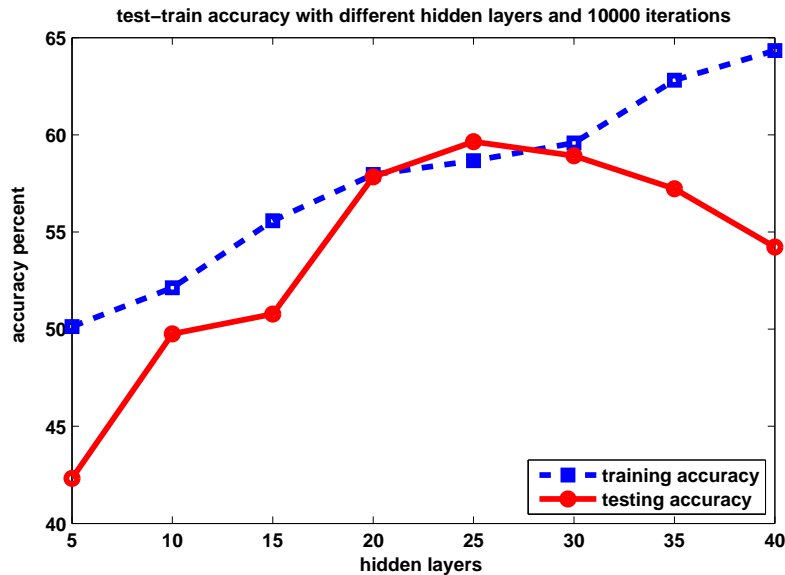


FIGURE 12. Cross entropy results for obtaining optimal number of hidden layers

## 5. CONCLUSION

This project was challenging to get more than 90% accuracy with small amount of data 6000 and it requires to apply many techniques. Beside If we increase number of data, convergence become problematic and using gradient decent type methods most of the times fall in local minimas. So we have some technical challenges from this type for which I applied tricks to resolve them.

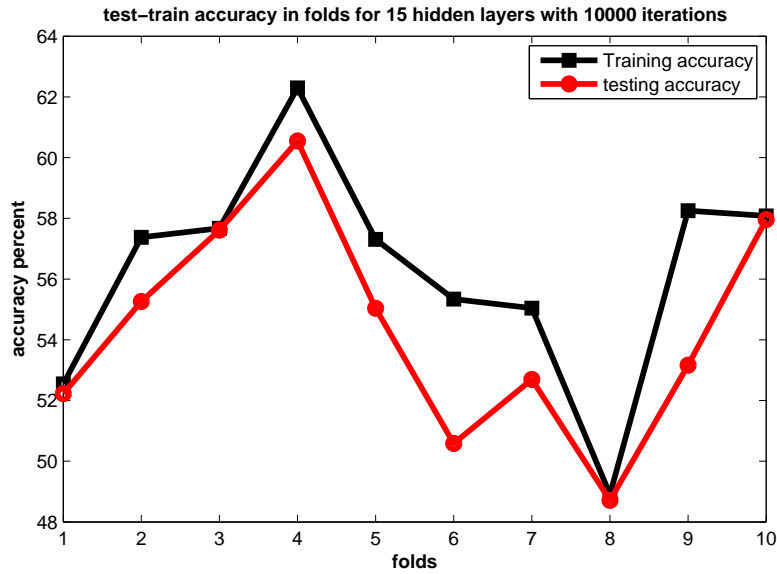


FIGURE 13. Results of accuracy in different folds with 15 hidden layers

Another important issue that we should notice is creating data carefully to cover many possible cases. Putting all together helped me to achieve this amount of accuracy.

There are other technical methods, including application of SVM method or other data processing techniques also can help to improve the accuracy even more, but with this amount of data I don't think that they can improve significantly and for my case results are obtained after many times training.

## 6. APPENDIX: MATLAB IMPLEMENTATIONS

### Main Part of the Code

```

1 %
2 % main This program do cross validation
3 %
4 clear all
5 close all
6 clc
7
8 %% Define parameters
9 MAXIT = 200000; % maximum iteration of training
10 nV = 25; % number of hidden units
11 batch_size = 10; % batch size
12 beta = 1e-10; % randomizer
13 mu = 0.0001; % optimization parameter
14 eta = 0.0001; % learning rate
15 fold = 10; % number of folds
16 %% Cross Validation data splitting

```

```

17 %try
18 %    data = [struct2cell(load('Data09.mat')){:}]; % IMPORTANT please
    Uncomment this when using OCTAVE
19 %catch
20     tr_tg = load('Data09.mat'); % IMPORTANT please Uncomment
    this when using MAILAB
21     data = struct2array(tr_tg); clear tr_tg % IMPORTANT please Uncomment
    this when using MAILAB
22 %end
23 %
24 %
25 %load weight1.mat % In the case of testing ( This loads Trained weights)
26 %%
27 [m_data, n_data] = size(data);
28 % _____ using 10-fold cross validation
    _____
29 total = linspace(1,m_data,m_data);
30 test_sel = total;
31
32 Train_data_ad = {}; Test_data_ad = {};
33 % number of folds
34 test_num = floor(m_data/fold);
35 % Cross validation splitting data – the address of data stored only
36 for n_fold = 1:fold
37     T_sample = randi(size(test_sel,2),test_num,1); % selecting 10% of the
    data randomly
38     Test_data_ad{n_fold} = T_sample;
39     Train_data_ad{n_fold} = setdiff(total,T_sample);
40     test_sel = setdiff(test_sel,T_sample);
41 end
42 clear T_sample test_num
43 %% Iteration on the folds
44 E_fold = [];
45 for n_fold = 1:fold
46     %% training process
47     % select type of the cost function:
48     %cost = 'cross'
49     cost = 'mse'
50     %cost = input('which kind of cost function do you prefer: enter "cross"
    for cross entropy and "mse" for mean square error: ')
51 [W,w] = NN_train(data(Train_data_ad{n_fold},:), cost, nV, eta, batch_size,
    beta, mu, MAXIT);
52
53 %% Testing Process
54 n_true = 0;
55 test_data = data(Test_data_ad{n_fold},:); % data is loaded from test
    fold
56 for i=1:size(test_data,1)
57     [O,Oh,V,Vh] = forwardPass(test_data(i,1:end-1),w,W);
58     [o_m,O_i]=max(O);
59     y_pred = O_i-1;
60     if y_pred == test_data(i,end)
61         n_true = n_true + 1;

```

```

62     end
63 end
64 E_fold = [E_fold n_true/size(test_data,1)*100];
65 fprintf('Test accuracy in fold %d is equal to: %f percent \n', n_fold,
        n_true/size(test_data,1)*100);
66
67 % clculating training accuracy
68 train_data = data(Train_data_ad{n_fold},:);
69 E_train = [];
70 n_true = 0;
71 for i=1:size(train_data,1)
72     [O,Oh,V,Vh] = forwardPass(train_data(i,1:end-1),w,W);
73     [o_m,O_i]=max(O);
74     y_pred = O_i-1;
75     if y_pred == train_data(i,end)
76         n_true = n_true + 1;
77     end
78 end
79 E_train = [E_train n_true/size(train_data,1)*100];
80 fprintf('Train accuracy in fold %d is equal to: %f percent \n', n_fold
        , n_true/size(train_data,1)*100);
81 clear W w test_data
82 end
83 fprintf('Average Test accuracy in %d folds is: %f \n', fold, mean(E_fold))
84 ;
85 fprintf('Average Train accuracy in %d folds is: %f \n', fold, mean(E_train
86 ));

```

### Training Neural Network

```

1 %Handwritten Recognition Project
2 % I did the simple one just for zero, next time for other digits too
3 % text time testing to see the performance if it was qualified
4 % thinking about possible improvements
5 % thinking about the demo
6 %
7 %
8 function [W, w] = NN_train(tr_data, cost, nV, eta, batch_size, beta, mu,
    MAXIT)
9 % ENTER YOUR NAME HERE:
10 %% preparing the data
11 % preprocessing the data
12 % determine which data to load
13
14 [tr_m, tr_n] = size(tr_data);
15 tr_data = [standardize(tr_data(:,1:end-1),0.5*ones(1,tr_n-1),0.5*ones(1,
    tr_n-1)) tr_data(:,end)];
16
17 %% initial weights
18 .....
19 if ~ (exist('W') && exist('w'))
20     W = (rand(10,nV)-0.5)*0.05; % hidden->output weights
21     w = (rand(nV,tr_n-1)-0.5)*0.01; % input->hidden weights
22 end;

```

```

22
23 %% initial parameters
24     .....
25     x = zeros(1, tr_n-1); % input layer
26     V = zeros(nV, tr_n-1); % hidden layer
27     O = 0; % output layer
28 %% setting the parameters
29     figure(1); clf;
30     Mean_E = 0;
31     err = zeros(MAXIT,1);
32     for it = 1:MAXIT,
33         tr_T = randi(tr_m, batch_size, 1);
34         x = tr_data(tr_T, 1:tr_n-1);
35         T = zeros(10, size(tr_T, 1));
36         % design target
37         for i=1:batch_size
38             T(tr_data(tr_T(i), tr_n)+1, i)=1; % training target, temporary
39         end
40         [O, Oh, V, Vh] = forwardPass(x, w, W);
41         %
42         % reshaping
43         Oh = reshape(Oh, 10*batch_size, 1);
44         O = reshape(O, 10*batch_size, 1);
45         T = reshape(T, 10*batch_size, 1);
46         if strcmp(cost, 'mse')
47             E = 0.5*norm(T-O, 2)^2; % SUM SQUARED ERROR
48             % MSE derivative case
49             dW = -reshape((T-O).*relder(Oh), 10, batch_size)*V; % SUM
50             SQUARED ERROR BASED dW
51             dw = -((reshape(((T-O).*relder(Oh)), batch_size, 10)*W).*relder
52             (Vh))'*x; % 0.04*w; % SUM SQUARED ERROR BASED dw
53         elseif strcmp(cost, 'cross')
54             % Cross validation objective function
55             E = -sum(T.*log10(O) + (1-T).*log10(1-O));
56             dW = -(1/(log2(10)))*((T.*sigder(Oh)./O)*V + ((1-T).*sigder(
57             Oh)./(1-O))*V); % my Tuning
58             dw = -(1/(log2(10)))*(((T.*sigder(Oh)./O)*W).*sigder(Vh))
59             '*x + (1/(log2(10)))*(((1-T).*sigder(Oh)./(1-O))*W).*sigder(Vh))*x
60             ; % My tuning
61         end
62         Mean_E = Mean_E + ((sqrt(E)/batch_size - Mean_E)/MAXIT);
63         err(it+1) = sqrt(Mean_E)/batch_size;
64         if (mod(it, 10000)==0) % show learnig progress
65             cla;
66             plot(0:10000:it, err(1:10000:it+1));
67             xlabel('iteration');
68             ylabel('squared error');
69             drawnow;
70         end;
71     end
72 %
73 if it==1

```

```

68     vect_p1 = -eta*dW;
69 end;
70     vect1 = mu*vect_p1 - eta*dW ; %
71     W = W - mu * vect_p1 + (1+mu)*vect1;
72     vect_p1 = vect1;
73 %
74 if it==1
75     vect_p = -eta*dw;
76 end;
77 % Accelerated Gradient
78     vect = mu*vect_p - eta*dw + beta*rand(nV, tr_n-1); %
79     w = w - mu * vect_p + (1+mu)*vect;
80     vect_p = vect;
81
82 end %iteration
83 if (it > 9000)
84     plot((0:10000:length(err)-1), err(1:10000:end));
85     xlabel('iteration');
86     ylabel('squared error');
87     drawnow;
88 end
89     ms_err = 0;
90     ms_err = ms_err/tr_m;
91     rms_err = ms_err^0.5;
92     if mod(it, 5000)==0
93         fprintf('RMS error over data points:%3.5f\n', rms_err);
94         title(sprintf('RMS error:%3.5f\n', rms_err));
95
96     end
97 %%
98 end % end for trainer
99 %end % main

```