

Adaptive Intelligence Assignment 1:

Supervised learning of handwritten letters

Wenjie Zhang

The University of Sheffield, Sheffield, United Kingdom

Abstract—This report focuses on the implementation of a supervised learning algorithm – an artificial neural network (ANN) based multi-layer perceptron with Mean Square Error, which is used to classify images. Regularization is done by adding penalty term L1. The data set used in this report is EMNIST, but samples per class restrict to 1000 for training and 250 for testing. For validation, we further split the train set into train set and development set with portion 8:2.

Index Terms—EMNIST, MLP, ANN network, MSE, L1,

I. INTRODUCTION

Supervised learning is a task of learning a function that maps inputs into output based on previously seen examples. It deals with two types of problems – classification and regression problems. This report mainly focuses on the classification problem, our objective is to find the pattern in the data that correlate with desired output and hence predict correct labels for unseen input data. In order to do it, there are some key ingredients: desired output(label) need to be provided along with the input, there has to be some feedback for the predicted output (i.e. how well the results against desired output), offline analysis is used rather than real-time analysis.

The process of training a classifier using supervised learning usually contains 4 stages: 1. Input raw data with corresponding desired output(label). 2. Using the algorithm to characterize input and according to features to predict outputs. 3. Using provided desired output to measure how well the algorithm is and generate errors(feedback). 4. According to the errors(feedback), tune parameters in the algorithm to minimize the error and repeat the process to reach the optimum results.

II. METHODS

A. Data Normalisation

As most commonly in machine learning, the features of data are measured at different scales and do not contribute equally to model fitting, model learning might end up creating a bias and take long time to learn and convergence. Thus, normalization is an important step for machine learning, which is especially useful for MLP. As we are doing back-propagation in MLP, normalization offer more stability and faster speed during the process. In this report, we choose Min-Max normalization, comparing to Standard deviation it converges more quickly and the process is faster. The reason for it is Min-Max has better distribution, as the data is pixel, value from 0-255, most value is 0. The difference between max and min is large and mean usually small.

B. Weight Initialisation

Besides normalization, weight initialization also affects the process of learning to a large extent. The objective of weight initialization is to prevent gradient exploding or vanishing. Gradient exploding means changes in nth layer quicker than n-1th layer, which leads to greater weight and gradient exploding, usually appears in multi-layer network and large initial weight. Gradient vanishing means latter layer (i.e. input layer to first hidden layer) times lots of < 1 numbers, which leads the results to 0, thus weight does not update, usually appears in multi-layer network and unsuitable activation function, such as sigmoid. Both due to the process of backpropagation. Desired weight initialization should have a mean value equal to 0. Xavier initialization initializes weight from a standard normal distribution and then scale by square root of $1/n$. The most important factor here is square root $1/n$, which adjusts variance and leads to faster convergence. Comparing to other methods of weight initialization, such as uniform distribution (0,1) only considering variance of input layer, Xavier considering variance of input and output layer and try to keep the same for each layer. If we initialize weight to zero or fixed value, the model may fail (no convergence). Thus, in this report, we choose Xavier initialization. But He initialization may be a better way to initialize for this report since we are using RELU activation function, $x_i > 0$.

C. Activation Function

The activation function enables the network to solve non-linear problems, it's a crucial element for neuron networks. Here we take sigmoid, tanh, and RELU as examples of activation functions to comparing. Sigmoid function has advantages, outputs are sensitive to input, and parameters updates more stable. But it has drawbacks, 1. no matter the value of inputs, the derivative also less than 0.5, and for large inputs the derivatives approach to 0 which may lead to gradient vanish. 2. The parameters of the layer are all updated in the positive or negative direction, which leads to slow convergence. 3. Long runtime since there is an exponential function. Another activation function tanh, it's quite similar to Sigmoid but it avoids the second drawback of Sigmoid. The activation function used in this report is RELU function, since the calculation is much simpler, faster convergence. Also, for positive numbers, the derivative of RELU is always 1 which avoids gradient explode. However, it also has some drawbacks, such as the derivatives always 0 for negative numbers, which means it unable to pass gradient for negative values. The

neuron with negative values won't update and propagate to n-1th layer, which is known as the dead RELU problem. We can avoid this problem by decreasing learning rates so that fewer dead neurons, a small amount of dead neurons won't affect too much.

D. Brief description of algorithm

Classification problems can be solved with numerous amounts of algorithms, such as Linear Classifiers, Support Vector Machines, Decision Trees, K-Nearest Neighbour and Random Forest, etc, choices of algorithms depend on the data and situation. In this report, as we are going to classify images of handwritten letters, the algorithm chosen is an Artificial neural network (ANN) based multi-layer perceptron. To build a more complex ANN, we need to build a basic single-layer network (perceptron). Perceptron includes input layer (pre-synaptic neurons), output layer (post-synaptic neurons), activation function, and weight from each pre-synaptic neuron to each post-synaptic neuron. Each input layer neuron match one pixel in data (sample data shape 1×784 , each sample has 784 pixels), each output layer neurons match to one letter class (26 in total), weight between each input neuron, and each output neuron is initialized by Xavier initialization, activation function used is RELU. For training model, first feed-forward input (input x \rightarrow $h = (w, x)$ \rightarrow $y = \text{RELU}(h)$) and calculate the error between output and desired output by MSE. Then do backpropagation use the error to adjust our weight. During this process, we need to know how weight changes affect error, so dE/dW . Since error related to $\text{RELU}(h)$, h related to w and x , by chain rule we get $dE/dW = dE/dY * dY/dH * dH/dW$. Since we are using Gradient descent, $dW = -\text{learning rate} * \text{derivative of MSE} * \text{derivative of RELU} * x$, the negative sign absorbs in derivative of MSE, we can define derivatives of $\text{MSE} * \text{derivative of RELU}$ as our local gradient. Thus, dW simplifies to learning rate $* \text{local gradient} * \text{input}$, then uses dW to adjust our weight. Since dY/dH is derivative of RELU equal to x ; $x=0$ $x \neq 0$ $x=1$, but $x=0$ is undefined in RELU derivative. We define derivative $x=0$ as equal to 1. In this report, we implement 2 hidden layer networks, using matrix instead of for loop to access data and update weights so that run time is decreased significantly.

III. RESULTS AND DISCUSSION

A. Question 3,4

Single layer network model is our baseline for comparing performance between different number of hidden layers. As the Fig.1 shows, the weight update after 400 epochs is converged close to 0. In [1] paper, the author achieves 55.78% accuracy on the EMNIST Letter data-set by linear classifier; by contrast, our model shows accuracy around 63% \pm 10%. This may due to Min-Max is used and runs 400 epochs.

B. Question 5,6

Overfitting usually occurs when lots of noises in training data, during training the model will learn these noises to decrease the error. The reason for this is model is too complex

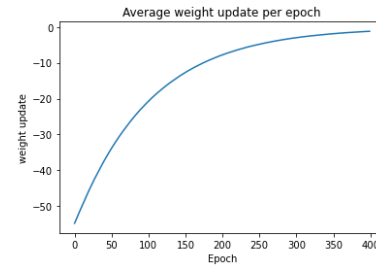


Fig. 1. weight update convergence.

with parameters more than the number of input patterns. This leads to decrease performance on unseen data, but our objective is a generalized model. Regularisation can prevent overfitting (see Fig.2) since it reduces variance problems which is essential in machine learning. Thus, regularisation methods need to be used.

There are plenty of methods of regularisation that stop overfitting, such as Early Stopping, Data Augmentation, Dropout, and L1, L2 which we used in this report. Dropout prevents overfitting by randomly ignore half of the hidden units during training, simply change the network during training. Data Augmentation preventing overfitting from more data samples, it changes original data to get more data, such as rotation, add noise, crop part of data (i.e. 784 pixels to 300 pixels). Early stopping is stopping the learning process before it reaches optimum. There is a point with the lowest error on training and unseen data, which is usually before the lowest error on training data. Early stopping is by stopping at that point to prevent overfitting. L1, L2 is used to reduce the complexity of our model by adding a penalty term to the MSE. Both are preventing overfitting by decreasing dW . There are differences between L1(Fig.3) and L2(Fig.4), we can see it though differentiate them. 1. dW decrease by a constant in L1 while dW decrease as proportional to W . 2. When there is large $|w|$ — $L1$ decrease smaller than $L2$, small $|w|$ — $L1$ decrease larger than $L2$. In this report, L1 is used according to the data and difference between L1 and L2 mentioned.

Since L1 can be seen as square root of square of w , the derivative is Fig.5. So, basically we are adding a $-\text{learning rate} * \lambda * \text{sign}(w)$ to dW , which is a constant. This penalty term does aid our learning decreasing validation error as Fig.6 and Fig.7. The optimum value of λ found is 0.00003

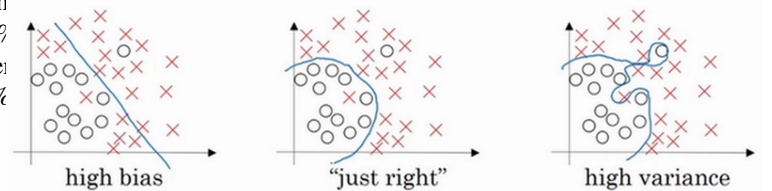


Fig. 2. variance effect.

$$L(x, y) \equiv \sum_{i=1}^n (y_i - h_{\theta}(x_i))^2 + \lambda \sum_{i=1}^n |\theta_i|$$

Fig. 3. adding L1.

$$L(x, y) \equiv \sum_{i=1}^n (y_i - h_{\theta}(x_i))^2 + \lambda \sum_{i=1}^n \theta_i^2$$

Fig. 4. adding L2.

$$|x|' = \frac{1}{2\sqrt{x^2}} \cdot 2x = \frac{x}{\sqrt{x^2}} = \frac{x}{|x|}$$

Fig. 5. adding L2.

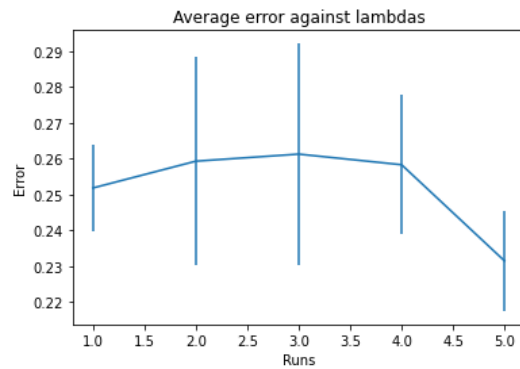


Fig. 6. without L1.

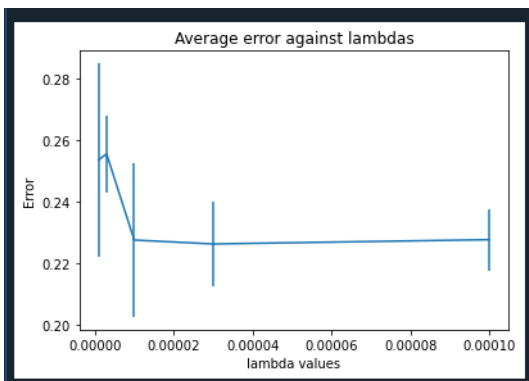


Fig. 7. with L1.

C. Question 7

In this report, we've tested single hidden layer neuron with number 50,100,200,400, the results is Fig.8. The error bar used is mean and standard deviation. This result may slightly different from running code, since i forgot to change data normalization method to Min-Max, n/255 is used in this result. Also, due to long runtime, epoches is 100 in these runs although matrix operation is performed to cancel inner for loop. Expect higher accuracy for longer time(converge).

As we can see from graph, the result is increased as number of neurons increasing in hidden layer. The reason for this might be, more neurons in hidden layer means it can remember more patterns between input and desired output In general, the result of adding a hidden layer is better than perceptron.

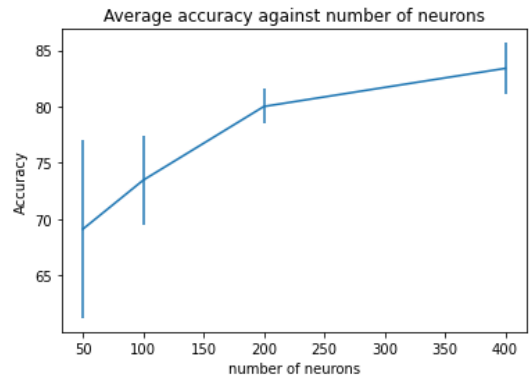


Fig. 8. question7.

D. Question 8

In this report, we've tested second hidden layer neuron with number 50,100,200,400(fixed first hidden layer neuron number = 100), the results is Fig.9. The error bar used is mean and standard deviation.

Again we can see a increasing trend as number of neurons increasing, same reason for question7. In general, more hidden layer and number of neurons both produce higher accuracy. But there is a upper bound, since there is noise in data. Also, hand written letters is a hard task to classify.

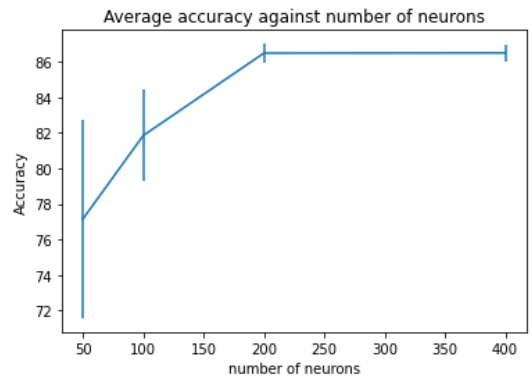


Fig. 9. question8.

ACKNOWLEDGMENT

This report used part of code from lab2.

There is 5 code file in total for different question

APPENDIX: EXAMPLE CODE FOR QUESTION8

```

normalize data by min-max Normalization  $x_{train}$  =
np.array([(samples - min(samples))/(max(samples) -
min(samples)) for samples in  $x_{train}$ ]) normalised $_train$  =
 $x_{train}$ / $x_{test}$  = np.array([(samples
min(samples))/(max(samples)
min(samples)) for samples in  $x_{test}$ ]) normalised $_test$  =
 $x_{test}$  randomised data with corresponding label randomised $_train$  =
shuf fle(normalised $_train$ , train_labels, random $_state$  =
0) randomised $_test$ , randomised $_test$  labels =
shuf fle(normalised $_test$ , test_labels, random $_state$  = 0)
split train data into dev set 8:2 train $_dev\_percent$  =
int(0.2 * randomised $_train$ .shape[0]) randomised $_dev$  =
randomised $_train$ [: train $_dev\_percent$ ] randomised $_dev$  labels =
randomised $_train$  labels[: train $_dev\_percent$ ] randomised $_train$  =
randomised $_train$ [train $_dev\_percent$ 
:] randomised $_train$  labels =
randomised $_train$  labels[train $_dev\_percent$  :]
 $n\_samples$ ,  $img\_size$  = randomised $_train$ .shape
The EMNIST contains A-Z so we will set the number of
labels as 26 n_labels = 26
transfer labels into array, index of element
equal to 1 is the label of that sample  $y_{train}$  =
np.zeros((randomised $_train$  labels.shape[0], n_labels))  $y_{test}$  =
np.zeros((randomised $_test$  labels.shape[0], n_labels))
for i in range(0, randomised $_train$  labels.shape[0]) :
 $y_{train}[i, randomised\_train\_labels[i].astype(int)] = 1$ 
for i in range(0, randomised $_test$  labels.shape[0]) :
 $y_{test}[i, randomised\_test\_labels[i].astype(int)] = 1$ 
The number of epochs is a hyperparameter that defines the
number times that the learning algorithm will work through
the entire training dataset.
The batch size is a hyperparameter that defines the number
of samples to work through before updating the internal model
parameters.
ref: https://machinelearningmastery.com/difference-between-a-batch-and-an-epoch/
 $n\_epoch$  = 250 batch $_size$  = 50 100 n_batches =
int(math.ceil( $n\_samples$ /batch $_size$ ))
define the size of each of the layers in the
network  $n\_input\_layer$  =  $img\_size$   $n\_hidden\_layer$  =
100 number of neuronsof the hidden layer. 0 delete this layer n_
n_labels
Add another hidden layer  $n\_hidden\_layer2$  =
0 100 number of neuronsof the hidden layer. 0 delete this layer
eta is the learning rate eta = 0.05
set lambda value lambda list = [0.000001, 0.000003,
0.00001, 0.00003, 0.0001] lamda = 0.00003 how many runs
different $_run$  = 5 record validation error validation $_error$  =
np.zeros((len(lambda list), different $_run$ ))

```

```

number of neuron num $_neuron$  =
[50, 100, 200, 400] record accuracy accuracy list =
np.zeros((len(num $_neuron$ ), different $_run$ ))
for question 6 for num $_neurons$  in range(0, len(num $_neuron$ )) :
for num $_tries$  in range(0, different $_run$ ) :
Initialize as simple network For W1 and W2 columns are the input and the
Number of columns(input) need to be equal to the number of features of
Number of columns(input) need to be equal to the number of neuronsof
 $n\_hidden\_layer2$  = num $_neuron$ [num $_neurons$ ] Xavier $_init$  =
True
if Xavier $_init$  : if  $n\_hidden\_layer$  > 0 : W1 =
np.random.randn( $n\_hidden\_layer$ ,  $n\_input\_layer$ ) *
np.sqrt(1/( $n\_input\_layer$ )) if  $n\_hidden\_layer2$  > 0 : W2 =
np.random.randn( $n\_hidden\_layer2$ ,  $n\_hidden\_layer$ ) *
np.sqrt(1/( $n\_hidden\_layer$ )) W3 =
np.random.randn( $n\_output\_layer$ ,  $n\_hidden\_layer2$ ) *
np.sqrt(1/( $n\_hidden\_layer2$ )) else : W2 =
np.random.randn( $n\_output\_layer$ ,  $n\_hidden\_layer$ ) *
np.sqrt(1/( $n\_hidden\_layer$ )) else : W1 =
np.random.randn( $n\_output\_layer$ ,  $n\_input\_layer$ ) *
np.sqrt(1/( $n\_input\_layer$ )) else :
if  $n\_hidden\_layer$  > 0 : W1 =
np.random.uniform(0, 1, ( $n\_hidden\_layer$ ,  $n\_input\_layer$ )) W2 =
np.random.uniform(0, 1, ( $n\_output\_layer$ ,  $n\_hidden\_layer$ ))
The following normalises the random weights
so that the sum of each row = 1 W1 =
np.divide(W1, np.matlib.repmat(np.sum(W1, 1)[:None], 1,  $n\_input\_layer$ ))) W2 =
np.divide(W2, np.matlib.repmat(np.sum(W2, 1)[:
, None], 1,  $n\_hidden\_layer$ ))
if  $n\_hidden\_layer2$  > 0 : W3 =
np.random.uniform(0, 1, ( $n\_output\_layer$ ,  $n\_hidden\_layer2$ )) W3 =
np.divide(W3, np.matlib.repmat(np.sum(W3, 1)[:
, None], 1,  $n\_hidden\_layer2$ ))
W2 = np.random.uniform(0, 1, ( $n\_hidden\_layer2$ ,  $n\_hidden\_layer$ )) W2 =
np.divide(W2, np.matlib.repmat(np.sum(W2, 1)[:
, None], 1,  $n\_hidden\_layer$ )) else : W1 =
np.random.uniform(0, 1, ( $n\_output\_layer$ ,  $n\_input\_layer$ )) W1 =
np.divide(W1, np.matlib.repmat(np.sum(W1, 1)[:
, None], 1,  $n\_input\_layer$ ))
Initialize the biases bias $_W1$  =
np.zeros(( $n\_output\_layer$ ,)) if  $n\_hidden\_layer$  > 0 :
bias $_W1$  = np.zeros(( $n\_hidden\_layer$ ,)) bias $_W2$  =
np.zeros(( $n\_output\_layer$ ,)) if  $n\_hidden\_layer2$  > 0 :
bias $_W3$  = np.zeros(( $n\_output\_layer$ ,)) bias $_W2$  =
np.zeros(( $n\_hidden\_layer2$ ,))
Train the network for i in range(0,  $n\_epoch$ ) :
Initialise the gradients for each batch dW1 =
np.zeros(W1.shape) if  $n\_hidden\_layer$  > 0 : dW2 =
np.zeros(W2.shape)
We will shuffle the order of the
samples each epoch shuffled $_dxs$  =
np.random.permutation( $n\_samples$ ) shuffled $_data$ ,  $s\_et$  =
randomised $_train$ [shuffled $_dxs$ ] shuffled $_label$ ,  $s\_et$  =
 $y_{train}$ [shuffled $_dxs$ ]
for batch in range(0, n_batches) :
Initialise the gradients for each batch dW1 =

```

```

np.zeros(W1.shape)dbias_W1 =
np.zeros(bias_W1.shape)if n_hidden_layer
0 : dW2 = np.zeros(W2.shape)dbias_W2 =
np.zeros(bias_W2.shape)if n_hidden_layer2 > 0 : dW3 =
np.zeros(W3.shape)dbias_W3 = np.zeros(bias_W3.shape)
shuffled_input_idx = shuffled_idx[batch * batch_size :
batch * batch_size + batch_size]
x0 = shuffled_data[batch * batch_size : batch * batch_size +
batch_size]
Form the desired output, the correct neuron should have 1
the rest 0
desired_output = shuffled_label[batch * batch_size :
batch * batch_size + batch_size]
h1 = np.matmul(W1,x0.T)+np.tile(bias_W1,(batch_size,1)).T
np.maximum(h1,0)
if n_hidden_layer > 0 : Neuralactivation :
hidden_layer- > outputlayerh2 = np.matmul(W2,x1) +
np.tile(bias_W2,(batch_size,1)).T
Apply the RELU function x2 = np.maximum(h2,0)
if n_hidden_layer2 > 0 : Neuralactivation :
hidden_layer1- > hidden_layer2h3 = np.dot(W3,x2) +
np.tile(bias_W3,(batch_size,1)).T
Apply the RELU function x3 = np.maximum(h3,0)
Compute the error signal e_n = desired_output - x3.T
Backpropagation: output layer -i hidden layer 2 delta3 =
e_n.T * np.where(x3 <= 0, 0, 1)
dW3 += np.matmul(delta3,x2.T) dbias_W3+ =
np.sum(delta3)
Backpropagation: hidden layer -i input layer delta2
= np.where(x2 != 0, 0, 1) * np.matmul(W3.T,
delta3) else: Compute the error signal e_n =
desired_output - x2.TBackpropagation : outputlayer- >
hidden_layerdelta2 = e_n.T * np.where(x2 <= 0, 0, 1)
dW2 += np.matmul(delta2, x1.T)
dbias_W2 = np.sum(delta2)
Backpropagation: hidden layer -i input layer delta1 =
np.where(x1 != 0, 0, 1) * np.dot(W2.T, delta2) else:
e_n = np.subtract(desired_output, x1.T)
delta1 = np.where(x1 != 0, 0, 1) * e_n
dW1 += np.matmul(delta1,x0)
dbias_W1+ = np.sum(delta1)
After each batch update the weights using accumulated
gradients as we can't let 0 undefined, so simply
set when x = 0, sign(x)=1 this part code is
for no L1 W1 += eta*dW1/batch_sizebias_W1+ =
eta * dbias_W1/batch_sizeif n_hidden_layer > 0 :
W2+ = eta * dW2/batch_sizebias_W2+ = eta *
dbias_W2/batch_sizeif n_hidden_layer2 > 0 : W3+ =
eta * dW3/batch_sizebias_W3+ = eta * dbias_W3/batch_size
W1 += eta*dW1/batch_size-eta*lamda*np.where(W1 <
0, -1, 1)bias_W1+ = eta * dbias_W1/batch_size
if n_hidden_layer > 0 : W2+ = eta * dW2/batch_size -
eta * lamda * np.where(W2 < 0, -1, 1)bias_W2+ =
eta * dbias_W2/batch_sizeif n_hidden_layer2 > 0 : W3+ =
eta * dW3/batch_size - eta * lamda * np.where(W3 <
0, -1, 1)bias_W3+ = eta * dbias_W3/batch_size

```

```

TODO: use the test set to compute the network's accuracy
n = randomised_test.shape[0]
p_ra = 0correct_value = np.zeros((n,))predicted_value =
np.zeros((n,))
for num in range(0, n): x0 = randomised_test[num]y =
y_test[num]
correct_value[num] = np.argmax(y)
h1 = np.dot(W1, x0) + bias_W1x1 =
np.maximum(h1,0)if n_hidden_layer > 0 :
h2 = np.dot(W2,x1) + bias_W2x2 =
np.maximum(h2,0)if n_hidden_layer2 > 0 : h3 =
np.dot(W3,x2) + bias_W3x3 = np.maximum(h3,0)
predicted_value[num] = np.argmax(x3)else :
predicted_value[num] = np.argmax(x2)else :
predicted_value[num] = np.argmax(x1)
if predicted_value[num] == correct_value[num] : p_ra =
(p_ra + 1)
accuracy = 100*p_ra/nprint("numberofneuron =
",num_neuron[num_neurons],"numberoftries =
",num_tries,"Accuracy =
",accuracy)accuracy_list[num_neurons][num_tries] =
accuracyyerr = np.zeros(len(num_neuron))mean_error =
np.zeros(len(num_neuron))for v in range(0, len(num_neuron)) :
mean_error[v] = np.mean(accuracy_list[v])yerr[v] =
np.std(accuracy_list[v])
plt.xlabel('number of neurons') plt.ylabel('Accuracy')
plt.errorbar(num_neuron, mean_error, yerr =
yerr)plt.title('Averageaccuracyagainstnumberofneurons')plt.show()

```

REFERENCES

- [1] G. Cohen, S. Afshar, J. Tapson, and A. van Schaik, "EMNIST: an extension of MNIST to handwritten letters." Retrieved from: <http://arxiv.org/abs/1702.05373>, (2017).
- [2] L. Prechelt, Early Stopping — But When?, pp. 53–67. Berlin, Heidelberg: Springer Berlin Heidelberg, 2012.

IEEE conference templates contain guidance text for composing and formatting conference papers. Please ensure that all template text is removed from your conference paper prior to submission to the conference. Failure to remove the template text from your paper may result in your paper not being published.