

MINI PROJECT

1 Question 1

In this question we are expected to implement an auto-encoder neural network architecture such that the architecture will have a single hidden layer to extract features from natural images in an unsupervised way. In this process the cost function below will be used for optimization.

$$J_{ae} = \frac{1}{2N} \sum_{i=1}^N \|d(m) - o(m)\|^2 + \frac{\lambda}{2} \left[\sum_{b=1}^{L_{hid}} \sum_{a=1}^{L_{in}} (W_{a,b}^{(1)})^2 + \sum_{c=1}^{L_{out}} \sum_{b=1}^{L_{hid}} (W_{b,c}^{(2)})^2 \right] + \beta \sum_{b=1}^{L_{hid}} KL(\rho | \tilde{\rho}_b)$$

In the equation above the first term is the average squared-error(MSE) between the network output and wanted result across training samples. The second term is a regularization term which is Tykhonov regularization with parameter λ . Finally the third term is Kullback-Leibler divergence (KL divergence). This term help to make the hidden unit activation sparse. Weight of KL divergence is controlled with β while the sparsity level is tuned with ρ in a bernoulli distribution with mean ρ and another with mean $\tilde{\rho}_b$, which is the mean activation of b.

1.1 Part a

In this part preprocessing of the data has been completed. In this step firstly 16x16 RGB images are read and they converted to grayscale using luminosity model, which is given below.

$$Y = 0.2126 * R + 0.7152 * G + 0.0722 * B$$

Then we are asked to normalize the data. To realize that firstly means of each image are subtracted from themselves. Then standard deviation is calculated and the data are clipped at range $[-3 * std, 3 * std]$. Then normalization formula is applied and after that the data is mapped into $[0.1, 0.9]$ range, $(0.8 * \text{Normalized} + 0.1)$. The normalization formula is given below. Also you can see the random 200 images before pre-processing and after pre-processing in fig.1 and fig.2.

$$\text{Normalized} = \frac{x - \min(x)}{\max(x) - \min(x)}$$

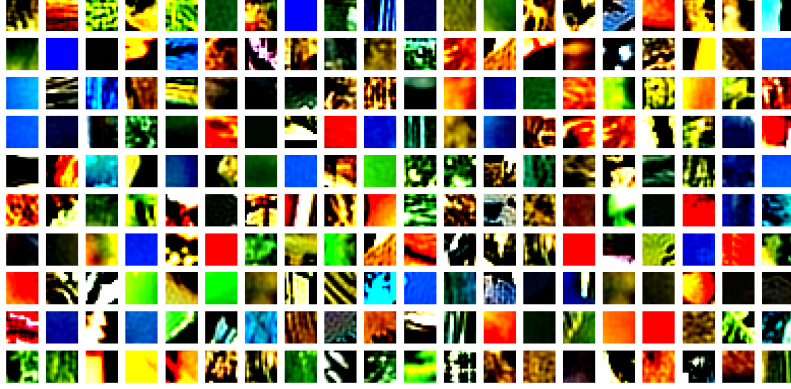


Figure 1: 200 random images before processing



Figure 2: 200 random images after processing

When we examine the figure 1 and figure 2, it can be clearly seen that most of the gray scale images are quite similar to their original versions and they are just gray versions of those images. I mean the color differences can be understood from the different shades. However, some of the grayscale images are slightly different than their original version. The reason behind this discrepancy is that we have clipped the data. However, in any case we are now ready to extract the features.

1.2 Part b

In this part we are firstly going to initialize weights and bias in the interval $[-\omega_o, \omega_o]$ where $\omega_o = \sqrt{\frac{6}{L_{pre} + L_{post}}}$, then a cost function which is asCost will be written. This function calculates the cost and the partial derivatives of weights and biases. The outputs of aeCost will be J and J_{grad} . Those terms will be given to gradient descent solver to minimize the cost that we have defined before after define and calculate parameters this solver is just an classical neural network architecture, with activation function sigmoid ($\sigma(x) = \frac{e^x}{1+e^x}$). The all equations and partial derivatives are given below for the cost function.

$$W_{all} = \begin{bmatrix} W \\ \beta \end{bmatrix}$$

1.2.1 MSE

$$MSE = \frac{1}{2N} \sum_{i=1}^N \|X - \tilde{X}\|^2$$

$$Y = \sigma(\sigma(\tilde{X} * W_{all_{hidden}}))$$

$$\delta_{out} = -\sigma'(V_{out}) \odot (X - \tilde{X})$$

$$\delta_{hid} = \delta'(V_{hid}) \odot W_{out}^T \delta_{out}$$

Partial derivatives are

$$\frac{\partial}{\partial w_{out}} MSE = W_{hid}^T \delta_{out}$$

$$\frac{\partial}{\partial w_{hid}} MSE = W_{out}^T \delta_{hid}$$

$$\frac{\partial}{\partial b_{out}} MSE = \frac{1}{N} \sum \delta_{out}$$

$$\frac{\partial}{\partial b_{hid}} MSE = \frac{1}{N} \sum \delta_{hid}$$

1.2.2 Tykhonov Regularization

Tykhonov Regularization=TRL

$$TRL = \frac{\lambda}{2} \left[\sum_{b=1}^{L_{hid}} \sum_{a=1}^{L_{in}} (W_{a,b}^{(1)})^2 + \sum_{c=1}^{L_{out}} \sum_{b=1}^{L_{hid}} (W_{b,c}^{(2)})^2 \right]$$

Partial derivatives are

$$\frac{\partial}{\partial w_{out}} TRL = \lambda W^{(2)}$$

$$\frac{\partial}{\partial w_{hid}} TRL = \lambda W^{(1)}$$

$$\frac{\partial}{\partial b_{out}} TRL = 0$$

$$\frac{\partial}{\partial b_{hid}} TRL = 0$$

1.2.3 Kullback-Leibler

Kullback-Leibler=KL

$$\beta \sum_{b=1}^{L_{hid}} KL(\rho | \tilde{\rho}_b)$$

Given that the distribution is bernoulli with parameter ρ given $\tilde{\rho}$, so the corresponding equation is given below.

$$\tilde{\rho} = \frac{1}{N} \sum \sigma(X \odot W_{hid} + b_{hid})$$

$$KL = \beta \sum_{b=1}^{L_{hid}} \rho \log\left(\frac{\rho}{\tilde{\rho}}\right) + (1 - \rho) \log\left(\frac{1 - \rho}{1 - \tilde{\rho}}\right)$$

Partial derivatives are

$$\frac{\partial}{\partial w_{out}} KL = 0$$

$$\frac{\partial}{\partial w_{hid}} KL = \beta \left[\frac{1 - \rho}{1 - \tilde{\rho}} - \frac{\rho}{\tilde{\rho}} \right]$$

$$\frac{\partial}{\partial b_{out}} KL = 0$$

$$\frac{\partial}{\partial b_{hid}} KL = 0$$

Notice that $\frac{\partial}{\partial w_{hid}} KL$ is a column vector and the operations are done according to this.

1.2.4 Summation

Known that $J_{ae} = MSE + KL + TRL$; therefore we also have the equation below.

$$\frac{\partial}{\partial x} J_{ae} = \frac{\partial}{\partial x} MSE + \frac{\partial}{\partial x} KL + \frac{\partial}{\partial x} TRL$$

This equation gives:

$$\frac{\partial}{\partial w_{out}} J_{ae} = W_{hid}^T \delta_{out} + \lambda W^{(2)}$$

$$\frac{\partial}{\partial w_{hid}} J_{ae} = W_{out}^T \delta_{hid} + \lambda W^{(1)} + \beta \left[\frac{1 - \rho}{1 - \tilde{\rho}} - \frac{\rho}{\tilde{\rho}} \right]$$

$$\frac{\partial}{\partial b_{out}} J_{ae} = \frac{1}{N} \sum \delta_{out}$$

$$\frac{\partial}{\partial b_{hid}} J_{ae} = \frac{1}{N} \sum \delta_{hid}$$

Those equations are used in aeCost function, you can check the function from appendix.

1.3 Part c

In this part we are asked to obtain the optimized weights, and display the first layer of connection weights for each hidden neuron as separate images. Those images can be seen below.

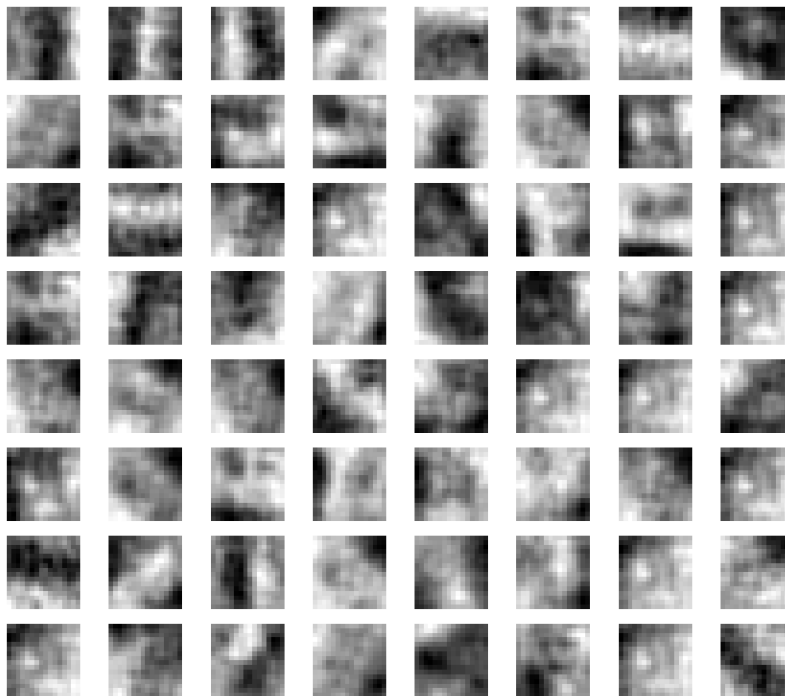


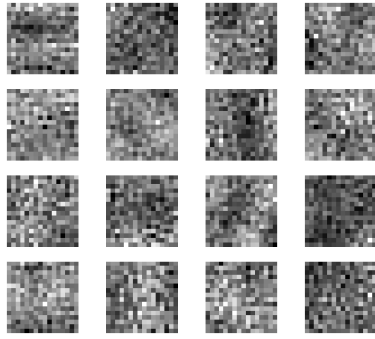
Figure 3: Visualization of First Layer of Connection Weights

When the figure above is examined it can be said that those images look like the small parts of big images. I mean there are lots of line-like or curve-like with different orientations can be observed from those weights. Therefore, it can be said that those images are not the variations of the original data, they are the parts that create the original images. Namely those are the features of the data.

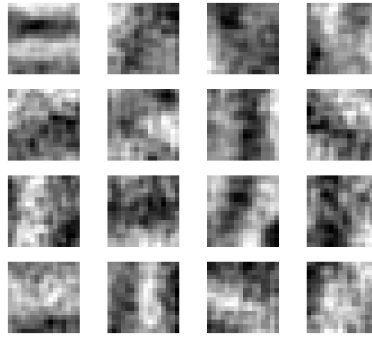
1.4 Part d

In this section we are asked to retrain the network for different values. In this training we are needed to choose 3 different $L_{hid} \in [10, 100]$ and 3 different $\lambda \in [0, 10^{-3}]$. Therefore, in total we are going to train this network for 9 more times for those parameters. The results of those trains can be seen below.

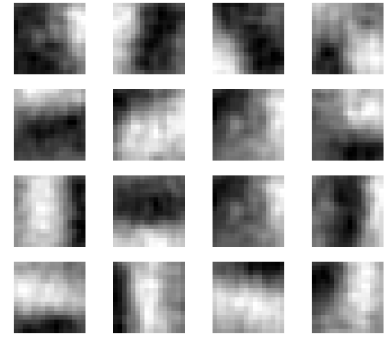
In figure 4, 5, 6 we can observe those 9 different cases. As you can see L_{hid} is the number of images. Also it can be said that the number of hidden layers affects the overall learning process. This is because, when it is increased, we can observe more different and complex plots in those figures. Actually number of those neurons are related to network capacity. However, the increase of network capacity does not mean better results every time, it can lead to overfitting in some cases since it starts to learn complex details. On the other hand,



(a) $L_{hid} = 16, \lambda = 0$



(b) $L_{hid} = 16, \lambda = 2 * 10^{-4}$



(c) $L_{hid} = 16, \lambda = 10^{-3}$

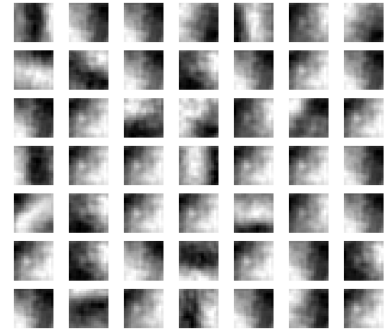
Figure 4: Images for $L_{hid} = 16$



(a) $L_{hid} = 49, \lambda = 0$



(b) $L_{hid} = 49, \lambda = 2 * 10^{-4}$



(c) $L_{hid} = 49, \lambda = 10^{-3}$

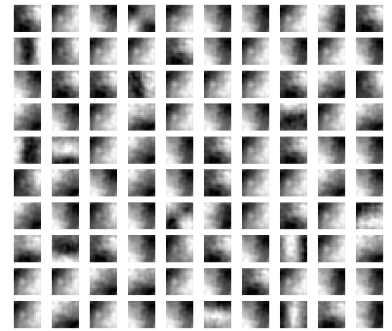
Figure 5: Images for $L_{hid} = 49$



(a) $L_{hid} = 100, \lambda = 0$



(b) $L_{hid} = 100, \lambda = 2 * 10^{-4}$



(c) $L_{hid} = 100, \lambda = 10^{-3}$

Figure 6: Images for $L_{hid} = 100$

small number of neurons may not be enough to learn all necessary features. Therefore, medium number can be said to be best.

Additionally in figure 4,5,6 we can examine the images for different λ values. As you can see when λ is higher the images are smoother. It can be said that λ helps the network in learning process by preventing it from overfitting. In this way the networks learn better rather than memorize the original data. However, as we can see from the equation, λ is in the regularization term. Therefore, we need to arrange the value of λ according to this terms, since if this term is too large then the cost will be bigger and it will corrupt the learning. Therefore, we should pick a medium λ .

2 Question 2

In the second question, we considered a model for examining sequences of words. The task was to get a prediction of the fourth word from the given 3 preceding words (Trigram). The architecture of this network is shown below.

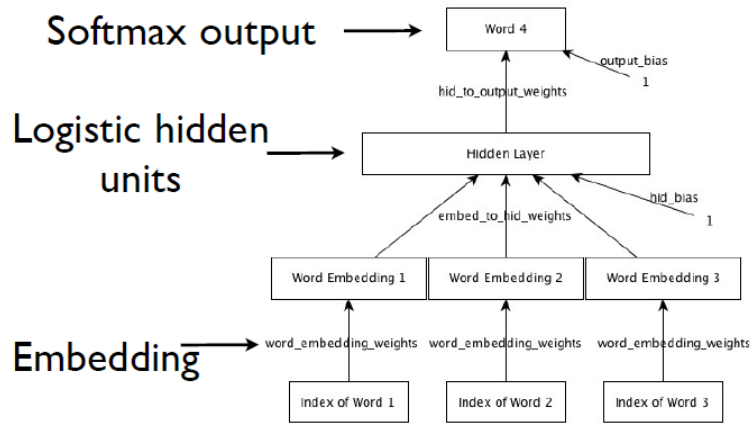


Figure 7: Network Architecture

As seen from the architecture we need to start to task from embedding layer. In order to give the inputs to this layer, indexes of the words in the given datasets has been converted to vector representations using one hot encoding method. The length of those vectors are the given dictionary size. Additionally in the question, it is indicated that the activation function of the hidden layer is sigmoid function while the output activation function is softmax function. ($Softmax : o_i = \frac{e^{z_i}}{\sum_{j=1}^{250} e^{z_j}}$, (250 neurons)), ($Sigmoid : \sigma(x) = \frac{e^x}{1+e^x}$)

2.1 Part a

In this part of the question we are required to implement stochastic gradient algorithm using the following parameters.

- Mini Batch Size=200

-
- Learning Rate (η)=0.15
 - Momentum Rate= α =0.85
 - Max Epoch=50
 - Gaussian variables weights and biases with std= 0.01
 - (D,P)=(32,256),(16,128),(8,64)

In order to implement the gradient descent I firstly needed to give the one hot encoded indices to embedding layer. In order to continue after this process I thought the embedding layer with an activation. To avoid any error the activation is $g(x)=x$, which is linear activation. Therefore, in total I have three different activation function. Therefore, I have added act_f function to select proper activation. Additionally, I had to add a function to choose the neuron, and a function to get the derivatives of activations. The derivatives are given below.

$$\frac{\partial}{\partial x}\sigma(x) = \frac{\partial}{\partial x}o_i(x) = x(1 - x)$$

The forward propagations of the network is quite standard, it gives the outputs of the activation functions to the next layer. The backpropagation is the most complex part of the question. In the backpropagation part δ is calculated like question 1, then update parameter is calculated. Using this update variable and momentum new weights are calculated.

$$\delta^{(t)} = (act)' * ((W^{(t+1)})^T \odot \delta^{(t+1)})$$

Where $(act)'$ represents the derivative of last activation function.

$$UP^{(t)} = \eta * \delta \odot V$$

$$W_{all} = W_{all} + \frac{UP}{batchsize} + \alpha * \frac{UP^{(t-1)}}{batchsize}$$

As a loss function I have used cross entropy error, in the train function. The train function is a classical train function of the neural network codes.

$$L(X, \tilde{X}) = \frac{-1}{N} \sum_{j=1}^N \sum_{i=1}^M x_{j,i} \log(x_{pred,j,i}) = \frac{-1}{N} \sum_{j=1}^N \tilde{X} * \log(X^T)$$

While the loss function is Cross Entropy loss the delta is given below.

$$\delta_{out} = X - \tilde{X}$$

The summary of the logic of the code is over, all the codes can be found in appendix part.

2.1.1 Optimization

I tried the codes with different parameters for several times and finally I have choose the parameters below.

- Learning Rate (η)=0.25
- Momentum Rate= α =0.8
- Gaussian variables weights and biases with std= 0.25

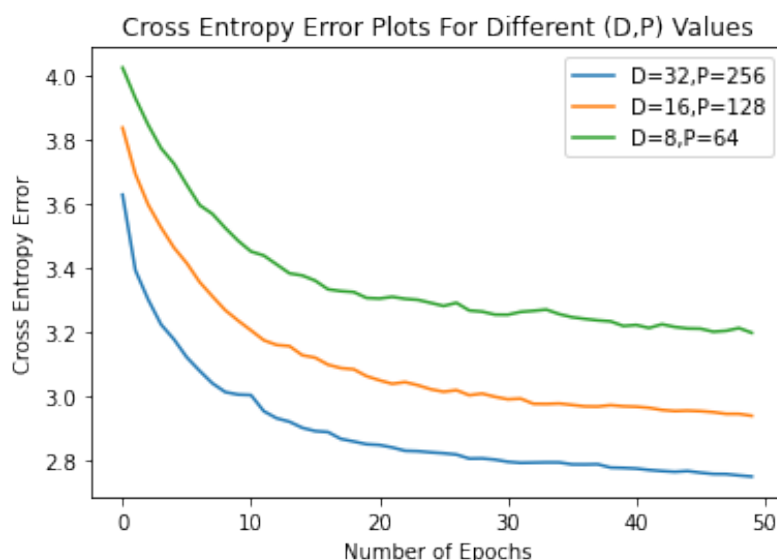


Figure 8: Cross Entropy Error curves for 3 different (D,P)

As seen from the figure, as expected loss is decreasing with epochs. However, after some points, losses converge. Additionally, with the decrease of (D,P) the error becomes bigger. The reason for this increase in error is that when D decreases mapping of words become more lossy. Also, when the hidden neuron number decreases learning capacity of the model decreases which is not a desired case as we know from previous question.

Additionally for this report I have not used early stopping based on cross entropy but in the code I have uploaded there is an early stopping. If the cross entropy loss is not changed more than 0.5% then it stops.

2.2 Part B

In this part of the question we are asked to predict 10 most possible words of the 4th word of 5 different Trigram. For this purpose I have written a predict function which is quite straightforward. It calculate forward propagation then sort and takes the highest 10 probability outputs.

I have used the predict function to create predictions for 3 of (D,P). You can see the results below.

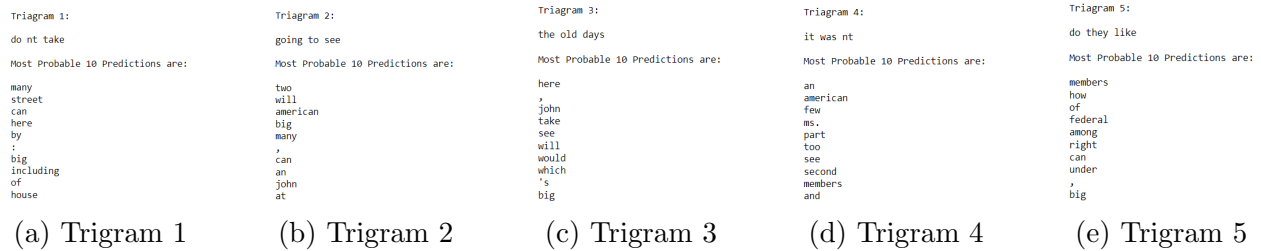


Figure 9: Most probable 10 predictions for 5 Triagram when (D,P)=(32,256)

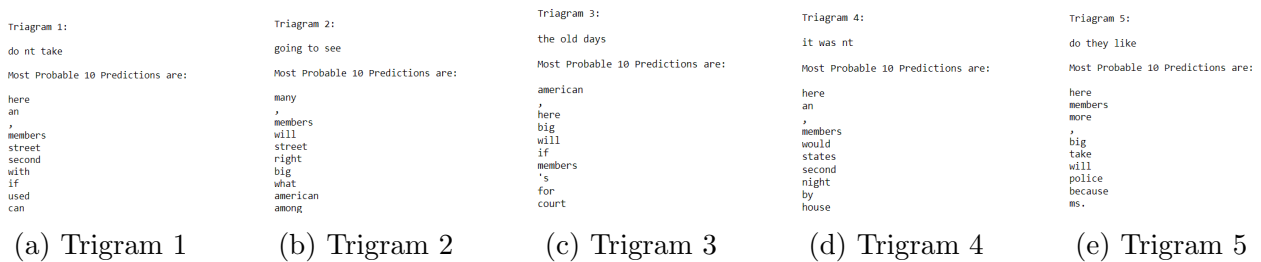


Figure 10: Most probable 10 predictions for 5 Triagram when (D,P)=(16,128)

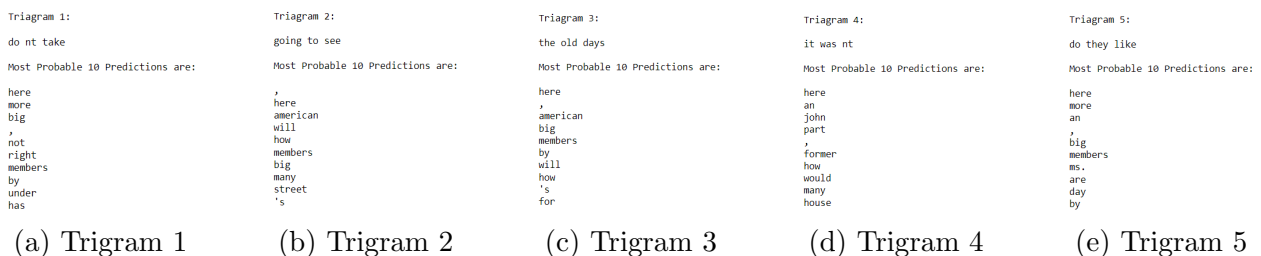


Figure 11: Most probable 10 predictions for 5 Triagram when (D,P)=(8,64)

As seen from the figure 9,10,11 we get 10 predictions for each 5 trigrams. While (D,P)=(32,256) has more meaningful predictions like "going to see two ..." (D,P)=(8,64) predictions are mostly have no meanings. It was expected from previous part and by figure 8.

3 Question 3

In this question we are going to create a network that will classify the human activities (downstairs=1, jogging=2, sitting=3, standing=4, upstairs=5, walking=6) according to movement signals, those signals are obtained from three different sensors for 150 time units. We have two different data-sets, which are training set with 3000 samples and test set with 600 samples.

3.1 Part A

Firstly, we will create a single layered RNN with hidden 128 neurons ,using the back propagation through time. In this process hyperbolic tangent (tanh) activation function. Then it will be followed by multi-layer perceptron network with softmax activation function will be used for classification. Throughout this processes I will use a stochastic gradient descent algorithm with learning rate ($\eta = 0.1$),momentum rate($\alpha = 0.85$) and maximum 50 epochs while the distribution of the weights and biases are Xavier Uniform distribution, which is given below.

$$W = [-a, a]$$

$$a = \sqrt{\frac{6}{L_{prev} + L_{next}}}$$

Then layer class initialization is similar to question 2, but the distribution for weights are different. For *act_f* function tanh and its derivative are added. Also a new function (*rec_{act}*) is added to lyr class, because now we are dealing with time series problem and we need to use previous sample as well.

After layer class, I have created RNN class, which contains forward propagation, back propagation, train, etc. The main idea and equation of the RNN network forward propagation is given below.

$$h(t) = f_H(W_{IH}x(t) + W_{HH}h(t-1))$$

$$y(t) = f_o(W_{HO}h(t))$$

The forward propagation function is divided to two parts which are that one for the recurrent layer and the other one is for the MLP layers. For the recurrent layer, it starts at time unit 0 and continue till the end time. At each iteration, the *rec_{act}* function is used and the current value is updated and its value is assigned for next prevupdate. The forward propagation of the MLP layers is the classical one like the previous questions, but differently

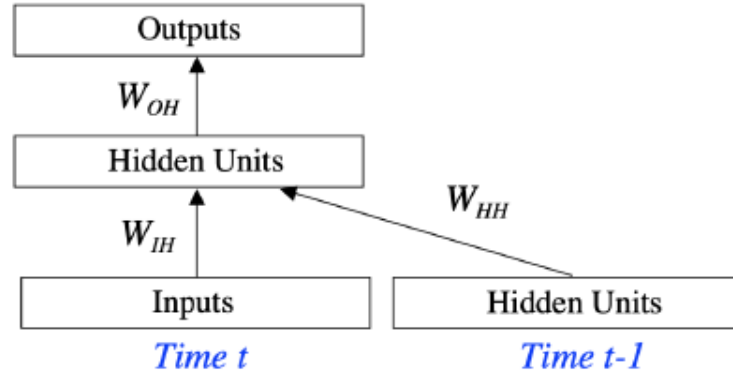


Figure 12: RNN Schematic

they use ReLU activation function.

For back propagation, both MLP back propagation and Backpropagation Through Time (BPTT) are used. Then the combination of those are used for update of weights.

$$W = W - (\alpha * W_{moment} + \eta * W_{up})$$

$$W_{up} = \sum_t^T \frac{\partial E(t)}{\partial W}$$

Similarly W_{moment} is updated at each epoch. The full backpropagation can be seen in back function from RNN class.

The other elements of the RNN class are classical functions like before. However, there are two different functions which are predict and conf. Predict function brings the accuracy of the prediction, conf function creates the confusion matrix of the predictions over the testset. The results of the network are given below in figure 13 and 14.

As it can be seen from figure 13 and 14 the results are not good. Firstly it is obvious that the loss plot is not stable. Namely, we cannot see any converge behaviour in the loss. Also from confusion matrices it can be seen that the network is prone to predict a certain behaviour. This instability is come from the gradients of the weights since they become very high or low during learning process. This can stop the learning of algorithm or it can cause the program to get warnings like Nan outputs. The reason for gradient explosion is probable due to the cumulative effect of Backpropagation Through Time (BPTT). This algorithm is used to train recurrent neural networks. It works by propagating the error backwards through time. This can lead to the gradients to accumulate to extremely high values, which can lead to instability. There are several potential solutions for this problem. Reducing the learning rate can help. It will slow the accumulation of weights, and make the network less vulnerable. Using clipping of gradients can help. It will limit the size of the gradients, and prevent them to become too large. Using normalization of layers can help. These layers can normalize the gradients and map them to a specific range. However, those solutions can be

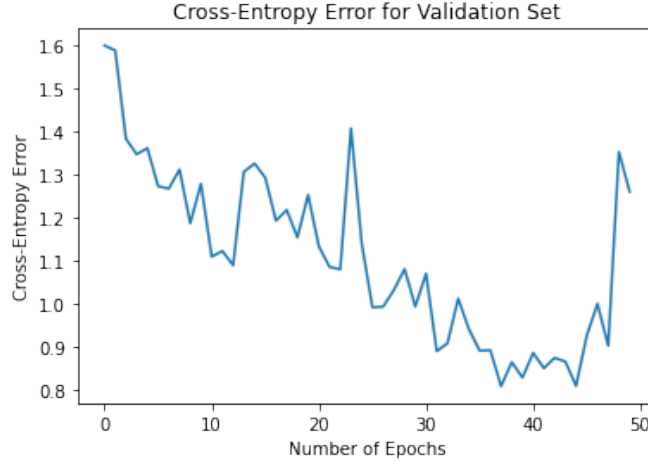
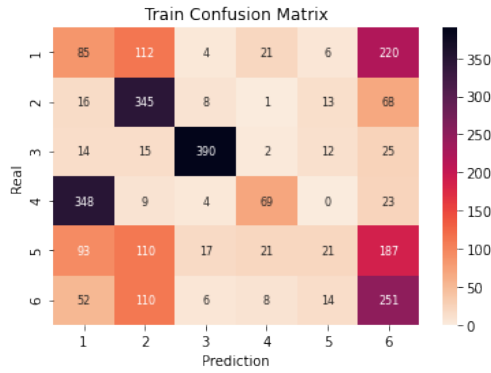
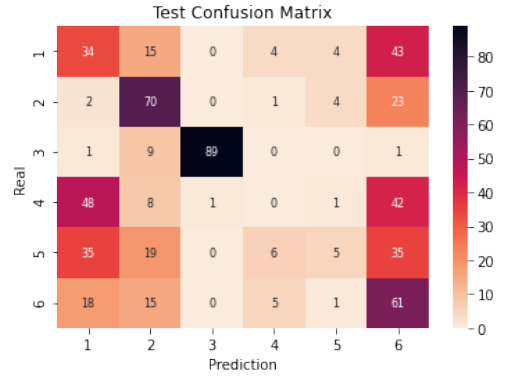


Figure 13: Cross Entropy Loss per Epoch for 50 epochs for RNN



(a) Confusion Matrix for Train, Train Accuracy= 43.0%



(b) Confusion Matrix for Test, Test Accuracy= 43.164%

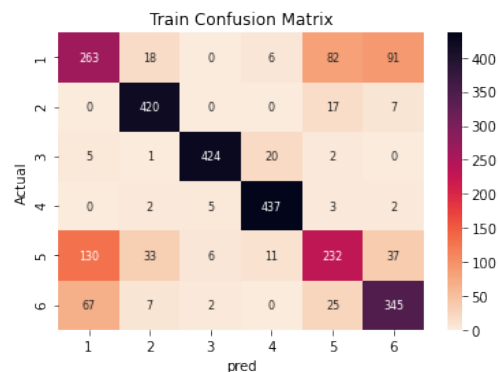
Figure 14: Confusion Matrices

used together, or also there is possibility that none of them works.

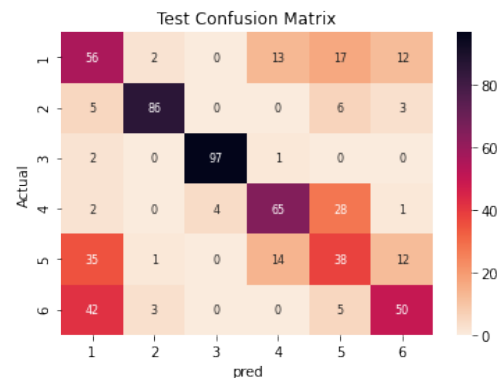
3.2 Part B

In this part we are going to repeat the previous part but rather than RNN we are going to use Long Short Term Memory(LSTM). This algorithm is better to handle the recurrent layers problems like extremely high or low gradients. In figure 15 you can see an LSTM cell. As seen from figure there are several gates in this structure and the equations for those gates are given below.

$$\begin{aligned}
 f_t &= \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) \\
 i_t &= \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) \\
 c_t &= \tanh(W_v \cdot [h_{t-1}, x_t] + b_v)
 \end{aligned}$$



(a) Confusion Matrix for Train,Train
Accuracy: 78.555%



(b) Confusion Matrix for Test,Test Ac-
curacy: 65.33%

Figure 17: Confusion Matrices

seems more prone to converge, the accuracy becomes better. I got an 78.55% train accuracy and 65.33% test accuracy. Even if those are not perfect, they are still huge improvements because they were respectively 43% and 43.164% in RNN results. Also from confusion matrices it can be seen that there is no single prediction problem anymore. Therefore, it can be said that the network can learn better.

4 Part C

In this part we are going to repeat the previous part but rather than RNN or LSTM we are going to use Gated Recurrent Unit, which is an cell-based network like LSTM. Gru has less gates than LSTM. In figure 18 you can see an LSTM cell. As seen from figure there are several gates in this structure and the equations for those gates are given below.

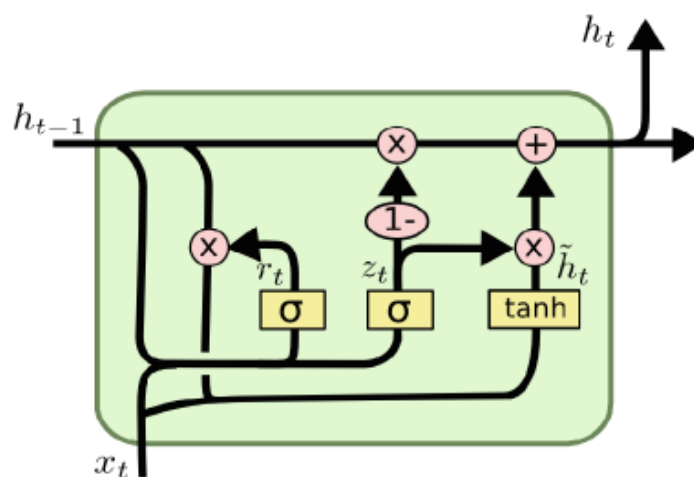


Figure 18: Schematic of GRU cell

$$\begin{aligned}
z_t &= \sigma(x_t.W_z + U_z.ht - 1 + b_z) \\
r_t &= \sigma(x_t.W_r + U_r.ht - 1 + b_r) \\
\tilde{h}(t) &= \tanh(x_t.W_h + U_h.ht - 1 * r_t + b_h) \\
h_t &= (1 - z_t) * h_{t-1} + z_t * \tilde{h}(t)
\end{aligned}$$

The backpropagation of the GRU is similar to BPTT of LSTM. The full implementation of the BPTT of GRU can be seen from back function in the GRU class. The rest of the GRU class is also almost same with the LSTM class.

After implementing the GRU network, I got the results below.

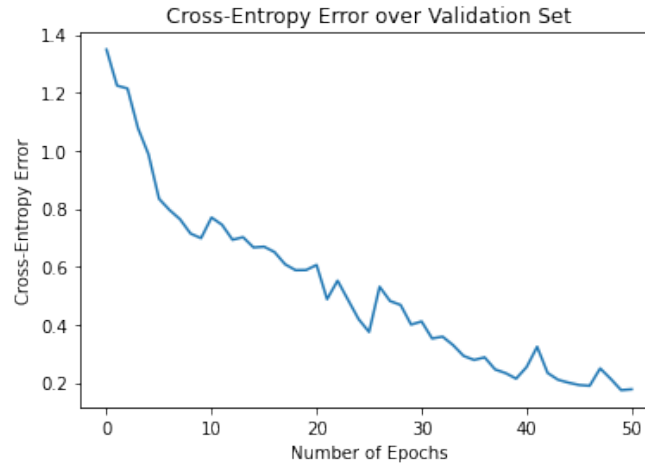
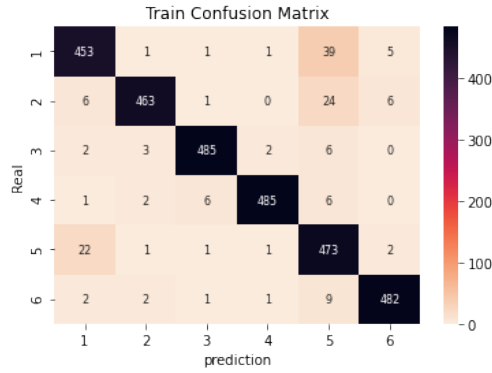


Figure 19: Cross Entropy Loss per Epoch for 50 epochs for GRU

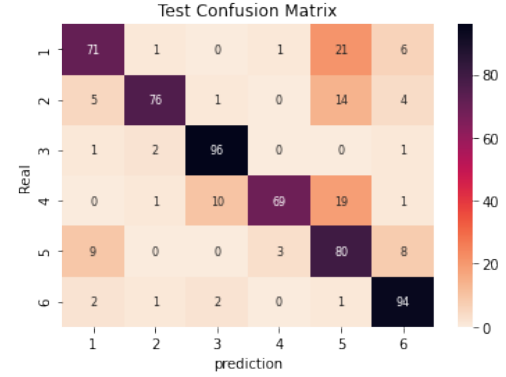
As we can see from figure 19 and 20 gru brought a big improvement to LSTM. When we observe the figures, we can easily see that the losses seems more prone to converge, the accuracy becomes better. I got an 94.85% train accuracy and 81% test accuracy. Even if those are not perfect, they are still huge improvements because they were respectively 78.55% and 65.33% in LSTM results. Also from confusion matrices it can be seen that the orientations of model to predict some behaviours changed and its much more better now. Also it is observed that GRU is faster than LSTM.

Even if GRU bring better results LSTM is more stable than GRU. In addition to that there is a high error possibility when trying to work with high learning rates. Therefore, I made the learning rate 0.3 for GRU.

When comparing all the results it can be said that RNN is the worst among those 3 models. If we want the training to be time efficient or if there are low time units GRU can be chosen, if there are high time units LSTM can be used. Since it has less vulnerable to overflows and can be more accurate with higher learning rates and high time units.



(a) Confusion Matrix for Train,Train
Accuracy: 94.85%



(b) Confusion Matrix for Test,Test Ac-
curacy: 81%

Figure 20: Confusion Matrices

5 Appendix- Codes

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 import seaborn as sn
4 import h5py
5 import os
6 import sys
7 import warnings
8 warnings.filterwarnings("ignore")
9
10 def initialize_parameters(Lin, Lhid):
11
12     np.random.seed(443)
13     Lout=Lin
14     w_0 =np.sqrt(6/(Lin+Lhid))
15     w_1=np.sqrt(6/(Lhid+Lout))
16
17     W1 = np.random.uniform(-w_0,w_0,(Lin,Lhid))
18     b1 = np.random.uniform(-w_0,w_0,(1,Lhid))
19
20
21     W2 = W1.T # at the same time W2=W1.T
22     b2 = np.random.uniform(-w_1,w_1,(1,Lout))
23     return (W1, W2, b1, b2), (0,0,0,0)
24
25 sigmoid = lambda x:np.exp(x)/ (1+np.exp(x))
26 derSigmoid= lambda x: sigmoid(x)*(1-sigmoid(x))
27
28 def forward(W_e, data):

```

```

29     W1, W2, b1, b2 = W_e
30
31     W1_ = data.dot(W1)+b1
32     z = sigmoid(W1_)
33     dz= derSigmoid(W1_)
34
35
36     W2_ = z.dot(W2)+b2
37     z2 = sigmoid(W2_)
38     dz2= derSigmoid(W2_)
39
40     cac=(data,dz,dz2)
41
42     return z, z2,cac
43
44 def aeCost(W_e, data, params):
45     (Lin, Lhid, lmb, beta, rho) = params
46     W1, W2, b1, b2 = W_e
47     N= len(data)
48
49     hid, out, (_,d_hid,d_out) = forward(W_e, data)
50     hid_mean = hid.mean(axis=0,keepdims=True)
51
52     ASE = (1/(2*N)) * np.sum(np.power((data - out),2)) #mse
53
54     TYK = (lmb/2) * (np.sum(W1**2) + np.sum(W2**2)) # Tykhonov
55     KL = rho*np.log(rho/hid_mean) + (1-rho)*np.log((1-rho)/(1-hid_mean))
56     KL = beta * KL.sum() # Kullback-Leibler
57
58     J = ASE + TYK + KL ## J_ae
59
60     d_ASE=-(data-out)/N
61     d_TYK=W1*lmb , W2*lmb
62     d_KL=beta*(- rho/hid_mean + (1-rho)/(1 - hid_mean))/N
63
64     d1 = d_ASE * d_out
65     dW2 = hid.T.dot(d1) + d_TYK[1]
66     db2 = d1.sum(axis=0, keepdims=True)
67
68     d2 = d_hid * (d1 .dot( W2.T) + d_KL)
69
70     dW1 = data.T .dot(d2) + d_TYK[0]
71     db1 = d2.sum(axis=0, keepdims=True)
72
73
74     J_grad=(dW1,dW2,db1,db2)
75     return J, J_grad

```

```

76
77 def update_parameters(We, mW, dW, m, l_rate):
78     mW = m * np.array(dW, dtype=object) + l_rate * np.array(mW, dtype=object)
79
80     We = np.array(We, dtype=object) - mW
81
82     We = tuple(We)
83     mW = tuple(mW)
84
85     return We, mW
86
87 def solver( data, params, eta, alpha, epoch, batch):
88
89     loss_list = []
90     if batch is None:
91         batch = len(data)
92
93     Lin = params[0]
94     Lhid = params[1]
95
96     We, mWe = initialize_parameters(Lin, Lhid)
97
98     iter_ep = int(len(data) / batch)
99     # j_prev=0 ##early stopping but not used
100    for k in range(epoch):
101        J_sum = 0
102
103        temp_start=0
104        temp_end = batch
105
106
107
108        for j in range(iter_ep):
109
110            data_temp = data[temp_start:temp_end]
111
112            J, Jgr = aeCost(We, data_temp, params)
113
114
115
116            We, mWe = update_parameters(We, mWe, Jgr, eta, alpha)
117
118            temp_start = temp_end
119            temp_end += batch
120            J_sum += J
121
122        J_sum = J_sum/iter_ep

```

```

123     #earl=abs(J_sum-j_prev)/abs(J_sum) ## early stopping but not used
124     #j_prev=J_sum ## early stopping but not used
125
126     #if earl<0.0001:
127     #     break
128
129     print("Loss: {:.2f} [Epoch {} of {}]".format(J_sum, k+1, epoch))
130     loss_list.append(J_sum)
131
132     return We, loss_list
133
134 def plot_weights(we):
135     (w1,w2,b1,b2) = we
136     fig=plt.figure(figsize=(18, 16))
137     plot_shape = int(np.sqrt(w1.shape[1]))
138     for i in range(w1.shape[1]):
139         plt.subplot(plot_shape,plot_shape,i+1)
140         plt.imshow(np.reshape(w1[:,i],(16,16)), cmap='gray')
141         plt.axis('off')
142     plt.show()
143
144 def label_1h(y, size):
145     output = np.zeros((len(y), size))
146     for q in range(len(y)):
147         temp = np.zeros(size)
148         temp[y[q]-1] = 1
149         output[q,:] = temp
150     return output
151
152 def data_1h(x, size):
153     output = np.zeros((len(x), x.shape[1], size))
154     for q in range(len(x)):
155         for p in range(x.shape[1]):
156             temp = np.zeros(size)
157             temp[x[q,p]-1] = 1
158             output[q,p,:] = temp
159     return output
160
161 class Layer:
162     def __init__(self,dim_in,neur_N,avg,std, act):
163         self.dim_in = dim_in
164         self.neur_N = neur_N
165         self.act = act
166         self.prevU = 0
167         self.Last_Act=None
168         self.err_layer=None
169         self.delta_layer=None

```

```

170     if self.act == 'softmax' or self.act == 'sigmoid':
171         self.all_W= np.random.normal(avg,std, (neur_N,dim_in+1))
172         self.weight=self.all_W[:, :-1]
173         self.bias=self.all_W[:, -1:]
174     else:
175         self.D = dim_in
176         self.S_dict = neur_N
177         self.weight = np.random.normal(avg, std, (self.S_dict,self.D))
178
179     def act_f(self, x):
180
181         if(self.act == 'softmax'):
182             e_x = np.exp(x - np.max(x))
183             return e_x/np.sum(e_x, axis=0)
184
185         elif(self.act == 'sigmoid'):
186             return np.exp(2*x)/(1+np.exp(2*x))
187         else:
188             return x
189
190     def act_N(self,x):
191         if self.act == 'sigmoid' or self.act == 'softmax':
192             samp_N = x.shape[1]
193             inp_temp = np.r_[x, [np.ones(samp_N)*-1]]
194             self.Last_Act = self.act_f(self.all_W.dot(inp_temp))
195
196         else:
197             Embed = np.zeros((x.shape[0],x.shape[1], self.D))
198             for m in range(Embed.shape[0]):
199                 Embed[m,:,:] = self.act_f(x[m,:,:].dot(self.weight))
200             Embed = Embed.reshape((Embed.shape[0], Embed.shape[1] *
201                                     ↪ Embed.shape[2]))
202             self.Last_Act = Embed.T
203         return self.Last_Act
204
205     def act_der(self, x):
206         if(self.act == 'sigmoid' or self.act == 'softmax'):
207             return x*(1-x)
208         else:
209             return np.ones(x.shape)
210
211
212
213 class Neural:
214     def __init__(self):
215         self.lyrs=[]

```

```

216
217 def layer_add(self, layer):
218     self.lyrs.append(layer)
219
220 def Forward(self, data_train):
221     IN=data_train
222     for layer in self.lyrs:
223         IN=layer.act_N(IN)
224     return IN
225
226 def Back(self, l_rate, size_batch, data_train, label_train, momentum):
227     out_forwa = self.Forward(data_train)
228     for i in reversed(range(len(self.lyrs))):
229         lyr = self.lyrs[i]
230
231         if(lyr == self.lyrs[-1]):
232             lyr.delta_layer=label_train.T-out_forwa
233         else:
234             layer_next = self.lyrs[i+1]
235             lyr.err_layer = np.matmul(layer_next.weight.T,
236                                     ↪ layer_next.delta_layer)
237             der=lyr.act_der(lyr.Last_Act)
238             lyr.delta_layer=der*lyr.err_layer
239
240     for i in range(len(self.lyrs)):
241         lyr = self.lyrs[i]
242         if(i == 0):
243             inp_temp = data_train
244         else:
245             samp_N = self.lyrs[i - 1].Last_Act.shape[1]
246             inp_temp = np.r_[self.lyrs[i - 1].Last_Act, [np.ones(samp_N)*-1]]
247
248         if(lyr.act == 'sigmoid' or lyr.act == 'softmax'):
249             update = l_rate*np.matmul(lyr.delta_layer, inp_temp.T)
250             lyr.all_W+= update/size_batch + (momentum*lyr.prevU)
251         else:
252             emb_delta = lyr.delta_layer.reshape((3,size_batch,lyr.D))
253             inp_temp = np.transpose(inp_temp, (1,0,2))
254             update = np.zeros((inp_temp.shape[2], emb_delta.shape[2]))
255             for i in range(emb_delta.shape[0]):
256                 update += l_rate * np.matmul(inp_temp[i,:,:].T,
257                                     ↪ emb_delta[i,:,:])
258             update = update
259             lyr.weight += update/size_batch + (momentum*lyr.prevU)
260             lyr.prevU = update/size_batch

```

```

261 def Train(self,l_rate,size_batch,data_train,label_train, data_test, lbl_test,
    ↪ num_ep,momentum,S_dict):
262     losses = []
263     temp_loss=0
264     for ep in range(num_ep):
265
266         print("Ep:",ep)
267         indexing=np.random.permutation(len(data_train))
268         data_train=data_train[indexing,:]
269         label_train=label_train[indexing]
270         batch_N = int(np.floor(len(data_train)/size_batch))
271         for j in range(batch_N):
272             data_one_hot =
    ↪ data_1h(data_train[j*size_batch:size_batch*(j+1)], S_dict)
273             label_one_hot =
    ↪ label_1h(label_train[j*size_batch:size_batch*(j+1)], S_dict)
274             self.Back(l_rate,size_batch,data_one_hot,label_one_hot,momentum)
275
276             out_val = self.Forward(data_test)
277             los_cros = - np.sum(np.log(out_val) * lbl_test.T)/out_val.shape[1]
278             print('C-E Error ', los_cros)
279             losses.append(los_cros)
280             print(abs(los_cros-temp_loss)/abs(los_cros))
281             if abs(los_cros-temp_loss)/abs(los_cros)<0.005:
282                 print("stopped based on cross-entropy")
283                 break
284             temp_loss=los_cros
285
286         return losses
287
288
289 def predict(self, inputIMG, k):
290     out = self.Forward(inputIMG)
291     return np.argsort(out, axis=0)[:,:k]
292
293
294 class lyr: #lyr Class for RNN
295     def __init__(self,dim_in,neur_N,act,beta):
296         self.dim_in = dim_in
297         self.neur_N = neur_N
298         self.act = act
299         self.beta=beta
300         self.prevU = 0
301         self.prevU_RNN = 0
302         self.Last_Act=None
303         self.err_lyr=None
304         self.delta_lyr=None

```

```

305
306     self.XUD=np.sqrt(6/(dim_in+neur_N))
307     self.all_W =np.random.uniform(-self.XUD,self.XUD,(dim_in+1,neur_N))
308     self.W2=np.random.uniform(-self.XUD,self.XUD,(neur_N,neur_N))
309
310
311
312
313     def act_f(self, x):
314
315         if(self.act == 'softmax'):
316             e_x = np.exp(x - np.max(x))
317             return e_x/np.sum(e_x, axis=1, keepdims=True)
318         elif(self.act == 'hyperbolic'):
319             return np.tanh(x*self.beta)
320         elif(self.act=="sigmoid"):
321             return np.exp(2*x)/(1+np.exp(2*x))
322         elif(self.act=="relu"):
323             return np.maximum(0,x)
324         else:
325             return x
326
327     def act_N(self,x):
328         samp_N = len(x)
329         inp_temp = np.concatenate((x, -np.ones((samp_N, 1))), axis=1)
330         self.Last_Act = self.act_f(inp_temp.dot(self.all_W))
331         return self.Last_Act
332
333     def rec_Act(self,x,hid):
334         samp_N = len(x)
335         inp_temp = np.concatenate((x, -np.ones((samp_N, 1))), axis=1)
336         last= hid.dot(self.W2)+ inp_temp.dot(self.all_W)
337         self.Last_Act = self.act_f(last)
338         return self.Last_Act
339
340     def act_derv(self, x):
341         if(self.act=="sigmoid" or self.act == 'softmax'):
342             return x*(1-x)
343         elif (self.act=="relu"):
344             return (x>0)*1
345         elif(self.act == 'hyperbolic'):
346             return self.beta*(1-x*x)
347         else:
348             return np.ones(x.shape)
349
350 class RNN_:
351     def __init__(self,data_train):

```



```

352     self.samp_T= data_train.shape[1]
353     self.del_rec=np.empty((32, self.samp_T, 128))
354     self.loss_rec=np.empty((32, self.samp_T, 128))
355     self.p_hid=np.zeros((32, 128))
356     self.lyrs=[]
357
358     def lyr_add(self,lyr):
359         self.lyrs.append(lyr)
360
361     def Forward(self,data_train):#Foward Prop
362         samp_N, samp_T, D = data_train.shape
363         inp=np.empty((samp_N, samp_T, 128))
364         self.p_hid=np.zeros((samp_N, 128))
365         for t in range(samp_T): #all time units are needed
366             x = data_train[:, t]
367             inp[:, t]=self.lyrs[0].rec_Act(x,self.p_hid)
368             self.p_hid=inp[:, t]
369         output = inp[:, -1]
370
371         for lyr in self.lyrs[1:len(self.lyrs)]: # MLP
372             output=lyr.act_N(output) #last sample contain memory
373         return inp,output
374
375     def back(self,l_rate,size_batch,data_train,lab_train,momentum):
376         inp,output = self.Forward(data_train)
377         out_forw = output
378         for i in reversed(range(len(self.lyrs))):# backprop til recurrent
379             lyr = self.lyrs[i]
380             #out layer
381             if(lyr == self.lyrs[-1]):
382                 lyr.delta_lyr=lab_train-out_forw
383             else :
384                 layer_next = self.lyrs[i+1]
385                 lyr.err_lyr =
386                     ↪ layer_next.delta_lyr.dot(layer_next.all_W[:len(layer_next.all_W)-1].T)
387                 derv=lyr.act_derv(lyr.Last_Act)
388                 lyr.delta_lyr=derv*lyr.err_lyr
389                 if (lyr == self.lyrs[0]):
390                     self.loss_rec[:, -1]=lyr.err_lyr
391                     self.del_rec[:, -1]=lyr.delta_lyr
392
393         d_all_weight=0
394         d_hid_weight=0
395         samp_N, samp_T, D = data_train.shape
396
397         for t in reversed(range(samp_T)):
398             lyr=self.lyrs[0]

```

```

398         if t > 0:
399             u = inp[:, t-1]
400         else:
401             u = np.zeros((samp_N, 128))
402
403         derv=lyr.act_derv(u)
404         d_hid_weight+=u.T.dot(self.del_rec[:,t])
405         d_all_weight+=np.concatenate((data_train[:,t-1],
406             ↪ -np.ones((len(data_train[:,t-1]), 1))),
407             ↪ axis=1).T.dot(self.del_rec[:,t])
408
409         # Recc dlt updt
410         self.loss_rec[:,t-1]=self.del_rec[:,t].dot(lyr.W2.T)
411         self.del_rec[:,t-1]=self.loss_rec[:,t-1]*derv
412
413     for i in range(len(self.lyrs)): #update all weights
414         lyr = self.lyrs[i]
415         if(i == 0):
416             lyr.prevU = d_all_weight*l_rate/(150*size_batch)
417             lyr.prevU_RNN = d_hid_weight*l_rate/(150*size_batch)
418             lyr.all_W+= (momentum*lyr.prevU) +lyr.prevU
419             lyr.W2+= (momentum*lyr.prevU_RNN) +lyr.prevU_RNN
420
421         else:
422             samp_N = len(self.lyrs[i - 1].Last_Act)
423             inp_temp=np.concatenate((self.lyrs[i - 1].Last_Act,
424                 ↪ -1*np.ones((samp_N, 1))), axis=1)
425             lyr.prevU = inp_temp.T.dot(lyr.delta_lyr)*l_rate/size_batch
426             lyr.all_W+= lyr.prevU*momentum +lyr.prevU
427
428     def train(self,l_rate,size_batch,data_train,lab_train, inp_test, lab_test,
429         ↪ N_ep,momentum):
430         cr_los_list = []
431         train_res=[]
432         for ep in range(N_ep):
433             print("ep:",ep)
434             ind=np.random.permutation(len(data_train))
435             data_train=data_train[ind]
436             lab_train=lab_train[ind]
437             batch_N = int(len(data_train)/size_batch)
438             for j in range(batch_N):
439                 train_data = data_train[j*size_batch:(j+1)*size_batch]
440                 train_labels = lab_train[j*size_batch:(j+1)*size_batch]
441                 self.back(l_rate,size_batch,train_data,train_labels,momentum)
442             _, out_val = self.Forward(inp_test)
443             _, out_train = self.Forward(data_train)

```

```

441         cr_los = np.sum(-lab_test*np.log(out_val))/len(out_val)
442         cr_los1 = np.sum(-lab_train*np.log(out_train))/len(out_train)
443         print('C-E Error of Validation', cr_los)
444         print('C-E Error Error of Train', cr_los1)
445         cr_los_list.append(cr_los)
446         train_res.append(cr_los1)
447     return cr_los_list, train_res
448
449     def Predict(self, inps, realout):
450         _, out = self.Forward(inps)
451         out = out.argmax(axis=1)
452         realout = realout.argmax(axis=1)
453         return ((out == realout).mean()*100)
454
455     def conf(self, inp, outp):
456         _, pred = self.Forward(inp)
457         pred = pred.argmax(axis=1)
458         outp = outp.argmax(axis=1)
459         q = len(np.unique(outp))
460         conf = np.zeros((q, q))
461         for p in range(len(outp)):
462             conf[outp[p]][pred[p]] += 1
463         return conf
464
465     class LSTM_lyr: #LSTM lyr Class
466     def __init__(self, dim_in, neur_N, beta):
467         self.dim_in = dim_in
468         self.neur_N = neur_N
469         self.beta = beta
470         self.Last_Act = None
471         self.err_lyr = None
472         self.delta_lyr = None
473         self.prevU_f, self.prevU_i, self.prevU_c, self.prevU_o = 0, 0, 0, 0
474
475         self.XUD = np.sqrt(6/(dim_in+neur_N))
476
477
478         self.Wf = np.random.uniform(-self.XUD, self.XUD, (dim_in+1, neur_N)) #
479         ↪ forget gate
480         self.Wi = np.random.uniform(-self.XUD, self.XUD, (dim_in+1, neur_N)) # input
481         ↪ gate
482         self.Wc = np.random.uniform(-self.XUD, self.XUD, (dim_in+1, neur_N)) # cell
483         ↪ gate
484         self.Wo = np.random.uniform(-self.XUD, self.XUD, (dim_in+1, neur_N)) #
485         ↪ output gate

```

```

484 def act_f(self, x,act):
485     if(act == 'softmax'):
486         e_x = np.exp(x - np.max(x))
487         return e_x/np.sum(e_x, axis=1, keepdims=True)
488     elif(act == 'tanh'):
489         return np.tanh(x*self.beta)
490     elif(act=="sigmoid"):
491         return np.exp(2*x)/(1+np.exp(2*x))
492     elif(act=="relu"):
493         return np.maximum(0,x)
494     else:
495         return x
496
497
498 def act_N(self,x, w, act):
499     samp_N = len(x)
500     inp_temp = np.concatenate((x, -np.ones((samp_N, 1))), axis=1)
501     self.Last_Act = self.act_f(inp_temp.dot(w),act)
502     return self.Last_Act
503
504 def rec_Act(self,x,hid, act):
505     samp_N = len(x)
506     inp_temp = np.concatenate((x, -np.ones((samp_N, 1))), axis=1)
507     last=hid.dot(self.W2) + inp_temp.dot(self.all_W)
508     self.Last_Act = self.act_f(last,act)
509     return self.Last_Act
510
511 def act_derv(self, x,act):
512     if(act=="sigmoid" or act == 'softmax'):
513         return x*(1-x)
514     elif (act=="relu"):
515         return (x>0)*1
516     elif(act == 'tanh'):
517         return self.beta*(1-x*x)
518     else:
519         return np.ones(x.shape)
520
521 class LSTM_:
522     def __init__(self,data_train):
523
524         self.lyrs=[]
525
526     def lyr_add(self,lyr):
527         self.lyrs.append(lyr)
528
529     def Forward(self,data_train):
530

```

```

531     sampN, sampT, sampD = data_train.shape
532     sampH=128
533     lyr = self.lyrs[0]
534
535     mem = np.empty((sampN, sampT, sampH))
536     i_t = mem
537     f_t = mem
538     c_t = mem
539     o_t = mem
540     tanhc = mem
541
542
543     h_t_1=np.zeros((sampN, sampH))
544     c_prv=h_t_1
545
546
547     z = np.empty((sampN, sampT, sampD + sampH))
548
549
550
551     #Apply functions
552     for i in range(sampT):
553         z[:, i] = np.concatenate((h_t_1, data_train[:, i]),axis=1)
554         zt = z[:, i]
555
556         i_t[:, i] = lyr.act_N(zt, lyr.Wi, "sigmoid")
557         f_t[:, i] = lyr.act_N(zt, lyr.Wf, "sigmoid")
558         c_t[:, i] = lyr.act_N(zt, lyr.Wc, "tanh")
559         o_t[:, i] = lyr.act_N(zt, lyr.Wo, "sigmoid")
560
561         mem[:, i] = c_t[:, i]*i_t[:, i] + c_prv*f_t[:, i]
562         c_prv=mem[:, i]
563
564         tanhc[:, i] = lyr.act_f(c_prv, "tanh")
565         h_t_1 = o_t[:, i] * tanhc[:, i]
566
567
568         cac = {"z_summ": z, #Summation of h_t-1 and x_t
569               "memory": mem, #Memory
570               "tanhc": (tanhc), # tanh memor
571               "f_t": f_t, # f_t out
572               "i_t": (i_t), # i_t out
573               "c_t": (c_t),# c_t out
574               "o_t": (o_t)}# o_t out
575
576
577     for lyr in self.lyrs[1:len(self.lyrs)]:

```

```

578         #For MLP lyrs
579         h_t_1=lyr.act_N(h_t_1)
580     OUT= h_t_1
581     return cac,OUT
582
583 def back(self,l_rate,size_batch,data_train,lab_train,momentum):
584     cac,output = self.Forward(data_train)
585     out_forw = output
586     z = cac["z_summ"]
587     c=cac["memory"]
588     tanhc=cac["tanhc"]
589     f_t=cac["f_t"]
590     i_t=cac["i_t"]
591     c_t=cac["c_t"]
592     o_t=cac["o_t"]
593
594     for q in reversed(range(len(self.lyrs))):# backprop til LSTM part
595         lyr = self.lyrs[q]
596
597         if(lyr == self.lyrs[-1]): #output
598             lyr.delta_lyr=lab_train-out_forw
599         elif(lyr==self.lyrs[0]):
600             layer_next = self.lyrs[q+1]
601             lyr.err_lyr =
602             ↪ layer_next.delta_lyr.dot(layer_next.all_W[0:len(layer_next.all_W)-1].T)
603             lyr.delta_lyr=lyr.err_lyr
604
605         else:
606             layer_next = self.lyrs[q+1]
607             lyr.err_lyr =
608             ↪ layer_next.delta_lyr.dot(layer_next.all_W[0:len(layer_next.all_W)-1].T)
609             derv=lyr.act_derv(lyr.Last_Act)
610             lyr.delta_lyr=derv*lyr.err_lyr
611
612     dWf, dWi, dWc, dWo = 0,0,0,0 #init grads to zero
613     H=128
614     T = z.shape[1]
615     samp_N = len(data_train)
616
617     init_lyr=self.lyrs[0]
618     delta=init_lyr.delta_lyr
619
620     for t in reversed(range(T)):#BPTT OF LSTM
621         u = z[:, t]
622
623         if t > 0:

```

```

623         c_prv = c[:, t - 1]
624     else:
625         c_prv = 0
626
627
628     dc = delta * o_t[:, t] * init_lyr.act_derv(tanhc[:, t], "tanh")
629
630     dc_t = dc * i_t[:, t] * init_lyr.act_derv(c_t[:, t], "tanh")
631     di_t = dc * c_t[:, t] * init_lyr.act_derv(i_t[:, t], "sigmoid")
632     df_t = dc * c_prv * init_lyr.act_derv(f_t[:, t], "sigmoid")
633     do_t = delta * tanhc[:, t] * init_lyr.act_derv(o_t[:, t], "sigmoid")
634
635
636     dWc += np.concatenate((u, -np.ones((samp_N, 1))), axis=1).T.dot(dc_t)
637     dWi += np.concatenate((u, -np.ones((samp_N, 1))), axis=1).T.dot(di_t)
638     dWf += np.concatenate((u, -np.ones((samp_N, 1))), axis=1).T.dot(df_t)
639     dWo += np.concatenate((u, -np.ones((samp_N, 1))), axis=1).T.dot(do_t)
640
641     # upd gradients
642     duc = dc_t.dot(init_lyr.Wc.T[:, :H])
643     dui = di_t.dot(init_lyr.Wi.T[:, :H])
644     duf = df_t.dot(init_lyr.Wf.T[:, :H])
645     duo = do_t.dot(init_lyr.Wo.T[:, :H])
646
647     delta = duc+dui+duf+duo
648
649
650     for i in range(len(self.lyrs)): #update the Weights
651         lyr = self.lyrs[i]
652         if(i == 0):
653
654             up_f, up_i, up_c, up_o = np.array([dWf, dWi, dWc, dWo]) * l_rate / size_batch
655
656             lyr.Wf += up_f + lyr.prevU_f * momentum
657             lyr.Wi += up_i + lyr.prevU_i * momentum
658             lyr.Wc += up_c + lyr.prevU_c * momentum
659             lyr.Wo += up_o + lyr.prevU_o * momentum
660
661             lyr.prevU_f
662             ↪ , lyr.prevU_i, lyr.prevU_c, lyr.prevU_o = np.array([up_f, up_i, up_c, up_o])
663
664         else:
665             samp_N = len(self.lyrs[i - 1].Last_Act)
666             inp_temp = np.concatenate((self.lyrs[i - 1].Last_Act,
667                                     ↪ -np.ones((samp_N, 1))), axis=1)
668             upd = (inp_temp.T.dot(lyr.delta_lyr)) * l_rate / size_batch
669             lyr.all_W += upd + lyr.prevU * momentum

```

```

668         lyr.prevU = upd
669
670     def train(self,l_rate,size_batch,data_train,lab_train, inp_test, lab_test,
        ↪ N_ep,momentum):
671         cr_los_list = []
672         train_res=[]
673         for ep in range(N_ep):
674             print("ep:",ep)
675             ind=np.random.permutation(len(data_train))
676
677             data_train=data_train[ind]
678             lab_train=lab_train[ind]
679             batch_N = int(len(data_train)/size_batch)
680             for j in range(batch_N):
681                 train_data = data_train[j*size_batch:size_batch*(j+1)]
682                 train_labels = lab_train[j*size_batch:size_batch*(j+1)]
683                 self.back(l_rate,size_batch,train_data,train_labels,momentum)
684             _, out_val = self.Forward(inp_test)
685             _, out_train = self.Forward(data_train)
686             cr_los = np.sum(-np.log(out_val) * lab_test)/len(out_val)
687             cr_los1 = np.sum(-np.log(out_train) * lab_train)/len(out_train)
688             print('C-E Error of Validation', cr_los)
689             print('C-E Error of Train', cr_los1)
690             cr_los_list.append(cr_los)
691             train_res.append(cr_los1)
692         return cr_los_list, train_res
693
694     def Predict(self,inps,realout):
695         _,out = self.Forward(inps)
696         out = out.argmax(axis=1)
697         realout = realout.argmax(axis=1)
698         return ((out == realout).mean()*100)
699
700
701     def conf(self,inp,outp):
702         _,pred= self.Forward(inp)
703         pred = pred.argmax(axis=1)
704         outp = outp.argmax(axis=1)
705         q = len(np.unique(outp))
706         conf=np.zeros((q,q))
707         for p in range(len(outp)):
708             conf[outp[p]][pred[p]] += 1
709         return conf
710
711 class GRU_lyr: #GRU lyr Class
712     def __init__(self,dim_in,neur_N,beta):
713         self.dim_in = dim_in

```



```

714     self.neur_N = neur_N
715     self.beta = beta
716
717     self.Last_Act=None
718     self.err_lyr=None
719     self.delta_lyr=None
720     self.prevU_Wz,self.prevU_Wr,self.prevU_Wh=0,0,0
721     self.prevU_Uz,self.prevU_Ur,self.prevU_Uh=0,0,0
722
723     self.XUD=np.sqrt(6/(dim_in+neur_N))
724     self.w1=np.sqrt(6/(neur_N+neur_N))
725
726     self.Uz = np.random.uniform(-self.w1, self.w1, size=(neur_N, neur_N))
727     self.Wz = np.random.uniform(-self.XUD,self.XUD,(dim_in+1,neur_N))
728
729     self.Ur = np.random.uniform(-self.w1, self.w1, size=(neur_N, neur_N))
730     self.Wr = np.random.uniform(-self.XUD,self.XUD,(dim_in+1,neur_N))
731
732     self.Uh = np.random.uniform(-self.w1, self.w1, size=(neur_N, neur_N))
733     self.Wh = np.random.uniform(-self.XUD,self.XUD,(dim_in+1,neur_N))
734
735
736 def act_f(self, x,act):
737     if(act == 'softmax'):
738         e_x = np.exp(x - np.max(x))
739         return e_x/np.sum(e_x, axis=1, keepdims=True)
740     elif(act == 'tanh'):
741         return np.tanh(x*self.beta)
742     elif(act=="sigmoid"):
743         return np.exp(x)/(1+np.exp(x))
744     elif(act=="relu"):
745         return np.maximum(0,x)
746     else:
747         return x
748
749 def act_N(self,x, w, h, u, act):
750     samp_N = len(x)
751     inp_temp = np.concatenate((x, -np.ones((samp_N, 1))), axis=1)
752     self.Last_Act = self.act_f(inp_temp.dot(w)+h.dot(u),act)
753     return self.Last_Act
754
755 def act_derv(self, x,act):
756     if(act=="sigmoid" or act == 'softmax'):
757         return x*(1-x)
758     elif (act=="relu"):
759         return (x>0)*1
760     elif(act == 'tanh'):

```

```

761         return self.beta*(1-x*x)
762     else:
763         return np.ones(x.shape)
764
765
766 class GRU_:
767     def __init__(self,data_train):
768
769
770         self.lyrs=[]
771     def lyr_add(self,lyr):
772         self.lyrs.append(lyr)
773
774     def Forward(self,data_train):
775
776         lyr = self.lyrs[0]      #GRU First lyr
777         sampN, sampT, _ = data_train.shape
778         sampH=128
779         h_t_1=np.zeros((sampN, sampH))
780         h_t = np.empty((sampN, sampT, sampH))
781         h_t_t = np.empty((sampN, sampT, sampH))
782
783         z_t = np.empty((sampN, sampT, sampH))
784         r_t = np.empty((sampN, sampT, sampH))
785
786         #apply funcs
787         for t in range(sampT):
788             x = data_train[:, t]
789             z_t[:, t] = lyr.act_N(x, lyr.Wz, h_t_1, lyr.Uz , "sigmoid")
790             r_t[:, t] = lyr.act_N(x, lyr.Wr, h_t_1, lyr.Ur , "sigmoid")
791             h_t_t[:, t] = lyr.act_N(x, lyr.Wh, (r_t[:, t] * h_t_1), lyr.Uh,
792                                     ↪ "tanh")
793             h_t[:, t] = (1 - z_t[:, t]) * h_t_1 + z_t[:, t] * h_t_t[:, t]
794             h_t_1 = h_t[:, t]
795
796             cac = {"z_t": z_t,
797                   "r_t": r_t,
798                   "h_t_t": (h_t_t),
799                   "h_t": h_t}
800
801         for ly in self.lyrs[1:len(self.lyrs)]: # MLP
802             h_t_1=ly.act_N(h_t_1)
803
804         oup= h_t_1
805         return cac,oup
806

```

```

807 def back(self,l_rate,size_batch,data_train,lab_train,momentum):
808     cac,OUT = self.Forward(data_train)
809     out_forw = OUT
810     z_t = cac["z_t"]
811     r_t=cac["r_t"]
812     h_t_t=cac["h_t_t"]
813     h_t=cac["h_t"]
814
815
816     for q in reversed(range(len(self.lyrs))):# backprop til GRU
817         lyr = self.lyrs[q]
818         #outputlyr
819         if(lyr == self.lyrs[-1]):
820             lyr.delta_lyr=lab_train-out_forw
821         elif(lyr==self.lyrs[0]):
822             layer_next = self.lyrs[q+1]
823             lyr.err_lyr =
824                 ⇨ layer_next.delta_lyr.dot(layer_next.all_W[0:len(layer_next.all_W)-1].T)
825             lyr.delta_lyr=lyr.err_lyr
826
827         else:
828             layer_next = self.lyrs[q+1]
829             lyr.err_lyr =
830                 ⇨ layer_next.delta_lyr.dot(layer_next.all_W[0:len(layer_next.all_W)-1].T)
831             derv=lyr.act_derv(lyr.Last_Act)
832             lyr.delta_lyr=derv*lyr.err_lyr
833
834         # initialize gradients to zero
835         dWz,dWr,dWh = 0,0,0
836         dUz,dUr,dUh = 0,0,0
837
838
839
840     sampH=128
841     samp_N, sampT, sampD = data_train.shape
842
843
844     init_lyr=self.lyrs[0]
845     delta=init_lyr.delta_lyr
846
847     for t in reversed(range(sampT)):
848         x = data_train[:, t]
849         if t > 0:
850             h_t_1 = h_t[:, t - 1]
851         else:
852             h_t_1 = np.zeros((samp_N, sampH))

```

```

851     dz = (h_t_t[:, t] - h_t_1) * init_lyr.act_derv(z_t[:, t], "sigmoid")
852     ↪ *delta
853     dh_t_t = z_t[:, t] * init_lyr.act_derv(h_t_t[:, t], "tanh")*delta
854     dr = dh_t_t.dot(init_lyr.Uh.T) * init_lyr.act_derv(r_t[:,
855     ↪ t], "sigmoid")* h_t_1
856
857     dWz += np.concatenate((x, -np.ones((samp_N, 1))), axis=1).T.dot(dz)
858     dWh += np.concatenate((x, -np.ones((samp_N, 1))),
859     ↪ axis=1).T.dot(dh_t_t)
860     dWr +=np.concatenate((x, -np.ones((samp_N, 1))), axis=1).T.dot(dr)
861
862     dUz += h_t_1.T.dot(dz)
863     dUh += h_t_1.T.dot(dh_t_t)
864     dUr += h_t_1.T.dot(dr)
865
866     # update the gradients
867
868     d9 = (1 - z_t[:, t])*delta
869     d11 = dz.dot(init_lyr.Uz.T)
870     d13 = dh_t_t.dot(init_lyr.Uh.T) * (r_t[:, t] + h_t_1 *
871     ↪ init_lyr.act_derv(r_t[:, t], "sigmoid").dot(init_lyr.Ur.T))
872
873     delta = d9+d11+d13
874
875     for i in range(len(self.lyrs)):
876         lyr = self.lyrs[i]
877         if(i == 0):
878
879             up_Wz,up_Wr,up_Wh=np.array([dWz,dWr,dWh])*l_rate/size_batch
880             up_Uz,up_Ur,up_Uh=np.array([dUz,dUr,dUh])*l_rate/size_batch
881
882             lyr.Wz+= up_Wz + lyr.prevU_Wz*momentum
883             lyr.Uz+= up_Uz + lyr.prevU_Uz*momentum
884             lyr.Wr+= up_Wr + lyr.prevU_Wr*momentum
885             lyr.Ur+= up_Ur + lyr.prevU_Ur*momentum
886             lyr.Wh+= up_Wh + lyr.prevU_Wh*momentum
887             lyr.Uh+= up_Uh + lyr.prevU_Uh*momentum
888
889             lyr.prevU_Wz
890             ↪ ,lyr.prevU_Wr,lyr.prevU_Wh=np.array([up_Wz,up_Wr,up_Wh])
891             lyr.prevU_Uz
892             ↪ ,lyr.prevU_Ur,lyr.prevU_Uh=np.array([up_Uz,up_Ur,up_Uh])
893
894     else:

```

```

892         samp_N = len(self.lyrs[i - 1].Last_Act)
893         inp_temp=np.concatenate((self.lyrs[i - 1].Last_Act,
894         ↪      -np.ones((samp_N, 1))), axis=1)
894         upd = (inp_temp.T.dot(lyr.delta_lyr))*l_rate/size_batch
895         lyr.all_W+= upd + lyr.prevU*momentum
896         lyr.prevU = upd
897
898     def train(self,l_rate,size_batch,data_train,lab_train, inp_test, lab_test,
899     ↪      N_ep,momentum):
900         cr_los_list = []
901         train_res=[]
902         for ep in range(N_ep):
903             print("ep:",ep)
904             ind=np.random.permutation(len(data_train))
905
906             data_train=data_train[ind]
907             lab_train=lab_train[ind]
908             batch_N = int(len(data_train)/size_batch)
909             for j in range(batch_N):
910                 train_data = data_train[j*size_batch:size_batch*(j+1)]
911                 train_labels = lab_train[j*size_batch:size_batch*(j+1)]
912                 self.back(l_rate,size_batch,train_data,train_labels,momentum)
913             _, out_val = self.Forward(inp_test)
914             _, out_train = self.Forward(data_train)
915             cr_los = np.sum(-np.log(out_val) * lab_test)/len(out_val)
916             cr_los1 = np.sum(-np.log(out_train) * lab_train)/len(out_train)
917             print('C-E Error of Validation', cr_los)
918             print('C-E Error of Train', cr_los1)
919             cr_los_list.append(cr_los)
920             train_res.append(cr_los1)
921         return cr_los_list, train_res
922
923     def Predict(self,inps,realout):
924         _,out = self.Forward(inps)
925         out = out.argmax(axis=1)
926         realout = realout.argmax(axis=1)
927         return ((out == realout).mean()*100)
928
929     def conf(self,inp,outp):
930         _,pred= self.Forward(inp)
931         pred = pred.argmax(axis=1)
932         outp = outp.argmax(axis=1)
933         q = len(np.unique(outp))
934         conf=np.zeros((q,q))
935         for p in range(len(outp)):
936             conf[outp[p]][pred[p]] += 1

```

```

937     return conf
938
939
940
941 def q1():
942     filename = 'data1.h5'
943     data=h5py.File(filename, 'r')['data'][()]
944     data_gray = 0.2126*data[:,0,:,:] + 0.7152*data[:,1,:,:] +
945     ↪ 0.0722*data[:,2,:,:]
946     mean_data = np.mean(data_gray, axis=(1,2))
947
948     for k in range(len(mean_data)):
949         data_gray[k,:,:] -= mean_data[k]
950
951     std_data = np.std(data_gray)
952
953     normalized = np.clip(data_gray, std_data*(-3),std_data*3)
954
955     normalization = lambda x:(x-x.min()) / (x.max()-x.min())
956
957     data_nor= normalization(normalized)*0.8 + 0.1
958
959     data_T=np.transpose(data,(0,2,3,1))
960     random_sample=np.random.randint(0,len(data_nor),size=(200))
961
962     row = 10
963     col = 20
964     fig=plt.figure(figsize=(20, 10), dpi= 100)
965     for j in range(row*col):
966         plt.subplot(row,col,j+1)
967         rand_pic=random_sample[j]
968         plt.imshow(data_T[rand_pic,:,:,:])
969         plt.axis('off')
970     plt.show()
971
972     fig=plt.figure(figsize=(20, 10), dpi= 100)
973     for j in range(row*col):
974         plt.subplot(row,col,j+1)
975         rand_pic=random_sample[j]
976         plt.imshow(data_nor[rand_pic,:,:],cmap='gray')
977         plt.axis("off")
978     plt.show()
979
980     num_pixel=data_nor.shape[1]
981     Lin = Lout = num_pixel**2
982     Lhid = 64

```

```

983     batch_size = 32
984     epoch = 200
985     lmb = 5e-4
986     alpha = 0.85
987     rho = 0.025
988     eta = 0.075
989     beta = 2
990
991     params= (Lin, Lhid, lmb, beta, rho)
992
993     data_solve=data_nor.reshape(data_nor.shape[0],data_nor.shape[1]**2)
994     w,j=solver(data_solve, params, eta, alpha, epoch, batch_size)
995
996     plot_weights(w)
997     hid_list=[16,49,100]
998     lmb_list=[0,2e-4,1e-3]
999     j_list=[j]
1000    w_list=[w]
1001    for i in range(len(hid_list)):
1002        for j in range(len(lmb_list)):
1003            params= (Lin, hid_list[i], lmb_list[j], beta, rho)
1004            print("parameters: Lin={}, Lhid={}, lmb={}, beta={}, rho={}
                  ↪ ".format(Lin, hid_list[i], lmb_list[j], beta, rho))
1005            w,j=solver(data_solve, params, eta, alpha, epoch, batch_size)
1006            j_list.append(j)
1007            w_list.append(w)
1008    for i in w_list[1:]:
1009        plot_weights(i)
1010
1011
1012 def q2():
1013     filename="data2.h5"
1014     testx = h5py.File(filename, 'r')['testx'][()]
1015     traind = h5py.File(filename, 'r')['traind'][()]
1016     trainx = h5py.File(filename, 'r')['trainx'][()]
1017     vald = h5py.File(filename, 'r')['vald'][()]
1018     valx = h5py.File(filename, 'r')['valx'][()]
1019     words = h5py.File(filename, 'r')['words'][()]
1020     v_size = 250
1021     lr = 0.15
1022     moment = 0.85
1023     batch = 200
1024     epoch = 50
1025
1026
1027
1028     valx1h = data_1h(valx, v_size)

```

```

1029     vald1h = label_1h(vald, v_size)
1030
1031     P_1 = 256
1032     D_1 = 32
1033
1034     model1 = Neural()
1035     model1.layer_add(Layer(D_1, v_size, 0, 0.25, 'emb'))
1036     model1.layer_add(Layer(3*D_1, P_1, 0, 0.25, 'sigmoid'))
1037     model1.layer_add(Layer(P_1, v_size, 0, 0.25, 'softmax'))
1038
1039     loss1 =
1040         ↪ model1.Train(lr, batch, trainx, traind, valx1h, vald1h, epoch, moment, v_size)
1041
1042     P_2 = 128
1043     D_2 = 16
1044
1045     model2 = Neural()
1046     model2.layer_add(Layer(D_2, v_size, 0, 0.25, 'emb'))
1047     model2.layer_add(Layer(3*D_2, P_2, 0, 0.25, 'sigmoid'))
1048     model2.layer_add(Layer(P_2, v_size, 0, 0.25, 'softmax'))
1049
1050     loss2 =
1051         ↪ model2.Train(lr, batch, trainx, traind, valx1h, vald1h, epoch, moment, v_size)
1052
1053     P_3 = 64
1054     D_3 = 8
1055
1056     model3 = Neural()
1057     model3.layer_add(Layer(D_3, v_size, 0, 0.25, 'emb'))
1058     model3.layer_add(Layer(3*D_3, P_3, 0, 0.25, 'sigmoid'))
1059     model3.layer_add(Layer(P_3, v_size, 0, 0.25, 'softmax'))
1060
1061     loss3 =
1062         ↪ model3.Train(lr, batch, trainx, traind, valx1h, vald1h, epoch, moment, v_size)
1063
1064     plt.plot(loss1, label="D=32,P=256")
1065     plt.plot(loss2, label="D=16,P=128")
1066     plt.plot(loss3, label="D=8,P=64")
1067     plt.legend()
1068     plt.title('Cross Entropy Error Plots For Different (D,P) Values')
1069     plt.xlabel('Number of Epochs')
1070     plt.ylabel('Cross Entropy Error')
1071     plt.show()
1072
1073     rnd = np.random.permutation(len(testx))[0:5]

```



```
1073     samp = testx[rnd,:]
1074
1075
1076
1077     samp1h = data_1h(samp, 250)
1078
1079     preds1 = model1.predict(samp1h, 10)
1080
1081     for i in range(5):
1082         print('Triagram {}: \n'.format(i+1))
1083         st=""
1084         for q in range(3):
1085             st+=str(words[samp[i,q]-1].decode("utf-8"))+" "
1086
1087         print(st+"\n")
1088         print('Most Probable 10 Predictions are:\n ')
1089         w_list=[]
1090         for j in range(10):
1091             w_list.append(str(words[preds1[j,i]-1].decode("utf-8")))
1092         for j in w_list:
1093             print(j)
1094         print("\n")
1095
1096
1097     preds2 = model2.predict(samp1h, 10)
1098
1099     for i in range(5):
1100         print('Triagram {}: \n'.format(i+1))
1101         st=""
1102         for q in range(3):
1103             st+=str(words[samp[i,q]-1].decode("utf-8"))+" "
1104
1105         print(st+"\n")
1106         print('Most Probable 10 Predictions are:\n ')
1107         w_list=[]
1108         for j in range(10):
1109             w_list.append(str(words[preds2[j,i]-1].decode("utf-8")))
1110         for j in w_list:
1111             print(j)
1112         print("\n")
1113
1114
1115     preds3 = model3.predict(samp1h, 10)
1116
1117     for i in range(5):
1118         print('Triagram {}: \n'.format(i+1))
1119         st=""
```

```

1120     for q in range(3):
1121         st+=str(words[samp[i,q]-1].decode("utf-8"))+" "
1122
1123     print(st+"\n")
1124     print('Most Probable 10 Predictions are:\n ')
1125     w_list=[]
1126     for j in range(10):
1127         w_list.append(str(words[preds3[j,i]-1].decode("utf-8")))
1128     for j in w_list:
1129         print(j)
1130     print("\n")
1131 def q3(t):
1132     file = h5py.File("data3.h5", "r") # get data
1133
1134     dat_tra = np.array(file[list(file.keys())[0]])
1135     lab_tra = np.array(file[list(file.keys())[1]])
1136     dat_test = np.array(file[list(file.keys())[2]])
1137     lab_test = np.array(file[list(file.keys())[3]])
1138
1139
1140
1141     neuron_n=128
1142     batch = 32
1143     epoch = 30
1144     alpha = 0.85 # momentum
1145     eta = 0.03 #learning rate
1146     ind=np.random.permutation(len(dat_tra))
1147     dat_tra=dat_tra[ind]
1148     lab_tra=lab_tra[ind]
1149
1150     size_v = int(len(dat_tra) / 10)
1151     dat_v=dat_tra[:size_v]
1152     lab_v=lab_tra[:size_v]
1153     dat_tra1=dat_tra[size_v:]
1154     lab_tra1=lab_tra[size_v:]
1155
1156     if t==1:
1157         model_rnn= RNN_(dat_tra1)
1158         model_rnn.lyr_add(lyr(3, neuron_n,'hyperbolic',1)) #3 val sensor
1159         model_rnn.lyr_add(lyr(neuron_n,70,'relu',1))
1160         model_rnn.lyr_add(lyr(70,30,'relu',1))
1161         model_rnn.lyr_add(lyr(30,6,'softmax',1))
1162         creL_rnn, lis_tra_rnn =
            ↪ model_rnn.train(eta,batch,dat_tra1,lab_tra1,dat_v,lab_v,epoch,alpha)
1163
1164     plt.plot(creL_rnn)
1165

```

```

1166     plt.xlabel('Number of Epochs')
1167     plt.ylabel('C-E Error ')
1168     plt.title('C-E Error of of Validation for RNN')
1169     plt.show()
1170
1171     acc_test_rnn=model_rnn.Predict(dat_test,lab_test)
1172     acc_train_rnn=model_rnn.Predict(dat_tra1,lab_tra1)
1173     print("Test set Accuracy: "+str(acc_test_rnn)+"%")
1174     print("Train set Accuracy: "+str(acc_train_rnn)+"%")
1175
1176
1177
1178
1179     conf_test_rnn=model_rnn.conf(dat_test,lab_test)
1180     sn.heatmap(conf_test_rnn, yticklabels=[1, 2, 3, 4, 5, 6], xticklabels=[1,
1181     ↪ 2, 3, 4, 5, 6], cmap=sn.cm.rocket_r, fmt='g')
1182     plt.xlabel("Prediction")
1183     plt.title("Confusion Matrix of test set for rnn")
1184     plt.ylabel("Real")
1185     plt.show()
1186
1187     conf_train_rnn=model_rnn.ConfusionMatrix(dat_tra1,lab_tra1)
1188
1189     sn.heatmap(conf_train_rnn, yticklabels=[1, 2, 3, 4, 5, 6],
1190     ↪ xticklabels=[1, 2, 3, 4, 5, 6], cmap=sn.cm.rocket_r, fmt='g')
1191     plt.xlabel("Prediction")
1192     plt.title("Confusion Matrix of train set for rnn")
1193     plt.ylabel("Real")
1194     plt.show()
1195
1196     elif t==2:
1197         model_lstm= LSTM_(dat_tra1)
1198         model_lstm.lyr_add(LSTM_lyr(131, neuron_n,1)) #3 val from sensor 128 prev
1199         model_lstm.lyr_add(lyr(neuron_n,70,'relu',1))
1200         model_lstm.lyr_add(lyr(70,30,'relu',1))
1201         model_lstm.lyr_add(lyr(30,6,'softmax',1))
1202         creL_lstm, lis_tra_lstm =
1203         ↪ model_lstm.train(eta,batch,dat_tra1,lab_tra1,dat_v,lab_v,epoch,alpha)
1204
1205     plt.plot(creL_lstm)
1206
1207
1208     plt.xlabel('Number of Epochs')
1209     plt.ylabel('C-E Error ')
1210     plt.title('C-E Error of Validation for LSTM')
1211     plt.show()
1212
1213     acc_test_lstm=model_lstm.Predict(dat_test,lab_test)

```

```

1210     acc_train_lstm=model_lstm.Predict(dat_tra1,lab_tra1)
1211     print("Test set Accuracy: "+str(acc_test_lstm)+"%")
1212     print("Train set Accuracy: "+str(acc_train_lstm)+"%")
1213
1214
1215
1216
1217     conf_test_lstm=model_lstm.conf(dat_test,lab_test)
1218     sn.heatmap(conf_test_lstm, yticklabels=[1, 2, 3, 4, 5, 6],
1219         ↪ xticklabels=[1, 2, 3, 4, 5, 6], cmap=sn.cm.rocket_r, fmt='g')
1219     plt.xlabel("Prediction")
1220     plt.title("Confusion Matrix of test set for lstm")
1221     plt.ylabel("Real")
1222     plt.show()
1223
1224     conf_train_lstm=model_lstm.ConfusionMatrix(dat_tra1,lab_tra1)
1225     sn.heatmap(conf_train_lstm, yticklabels=[1, 2, 3, 4, 5, 6],
1226         ↪ xticklabels=[1, 2, 3, 4, 5, 6], cmap=sn.cm.rocket_r, fmt='g')
1226     plt.xlabel("Prediction")
1227     plt.title("Confusion Matrix of train set for lstm")
1228     plt.ylabel("Real")
1229     plt.show()
1230     elif t==3:
1231         model_gru = GRU_(dat_tra1)
1232         model_gru.lyr_add(GRU_lyr(3, neuron_n,1))
1233         model_gru.lyr_add(lyr(neuron_n,70,'relu',1))
1234         model_gru.lyr_add(lyr(70,30,'relu',1))
1235         model_gru.lyr_add(lyr(30,6,'softmax',1))
1236         creL_gru, lis_tra_gru =
1237         ↪ model_gru.train(eta,batch,dat_tra1,lab_tra1,dat_v,lab_v,epoch,alpha)
1237
1238     plt.plot(creL_gru)
1239
1240     plt.xlabel('Number of Epochs')
1241     plt.ylabel('C-E Error ')
1242     plt.title('C-E Error ofValidation for GRU')
1243     plt.show()
1244
1245     acc_test_gru=model_gru.Predict(dat_test,lab_test)
1246     acc_train_gru=model_gru.Predict(dat_tra1,lab_tra1)
1247     print("Test set Accuracy: "+str(acc_test_gru)+"%")
1248     print("Train set Accuracy: "+str(acc_train_gru)+"%")
1249
1250
1251
1252
1253     conf_test_gru=model_gru.conf(dat_test,lab_test)

```

```
1254     sn.heatmap(conf_test_gru, yticklabels=[1, 2, 3, 4, 5, 6], xticklabels=[1,
    ↪ 2, 3, 4, 5, 6], cmap=sn.cm.rocket_r, fmt='g')
1255     plt.xlabel("Prediction")
1256     plt.title("Confusion Matrix of test set for gru")
1257     plt.ylabel("Real")
1258     plt.show()
1259
1260     conf_train_gru=model_gru.ConfusionMatrix(dat_tra1,lab_tra1)
1261
1262     sn.heatmap(conf_train_gru, yticklabels=[1, 2, 3, 4, 5, 6],
    ↪  xticklabels=[1, 2, 3, 4, 5, 6], cmap=sn.cm.rocket_r, fmt='g')
1263     plt.xlabel("Prediction")
1264     plt.title("Confusion Matrix of train set for gru")
1265     plt.ylabel("Real")
1266     plt.show()
1267
1268     ## in order to run the first question please use q1() function
1269     ## in order to run the second question please use q2() function
1270     ##in order to run the third question please use the q3(t) function, such that t
    ↪ can only take 1,2 or 3
1271     ## those t values represent the part of the third question
1272     ## Thus, in order to run RNN code use q3(1)
1273     ## in order to run LSTM code use q3(2)
1274     ## in order to run GRU code use q3(3)
```