

Guidance Note - Use of LSTM Code

Issac Pang

April 12, 2023

This guidance note provides a brief summary of how to use the developed LSTM code for testing. In general, the LSTM code consists of two main modules: **Training** and **Utility Functions**. The **Training** module contains all LSTM algorithms (e.g., LSTM(Vanilla) and TCN). If you wish to use/develop the LSTM algorithms to train a new model, please go to the **Training** module. The **Utility Functions** module is mainly used to generate the 3D CSMIP data for training and to evaluate the performance of trained LSTM models.

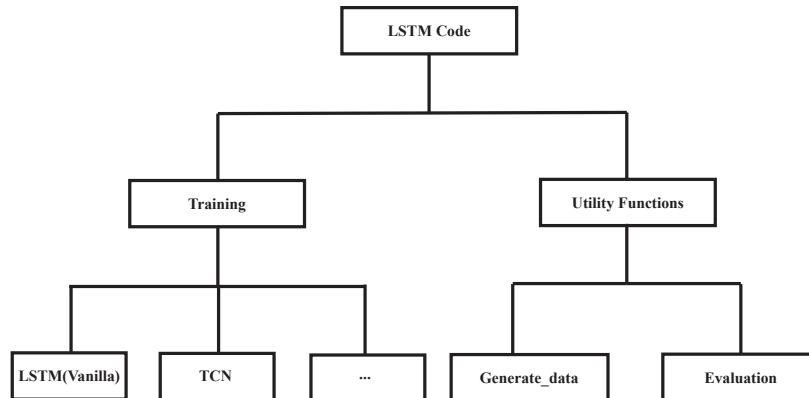


Figure 1: LSTM Code Organization

To use the code for testing or for other research purposes, please download all files in the **Training** and **Utility Functions** modules. Then, update the directory of **train_folder_path** (line 9) and **modelPath** (line 39). The **train_folder_path** is used to store all CSMIP training data, and the **modelPath** is used to save the trained LSTM model. Number of time steps may also need to be updated (line 10), depending on the CSMIP data.

Example: LSTM(vanilla).py

```
1  ...
2
3  if __name__ == "__main__":
4
5      #specify the train_folder_path
6      #specify time step (each seismic event has different timesteps, so need to define manually)
7      #specify output_response (e.g. Accel data in Channel 3, 5, 8 => then 3)
8      #specify window size (stack size, please refer to Zhang2019 for details)
9      train_folder_path = r'C:\Users\BRACE2\Desktop\CSMIP\data\training'
10     time_step = 7200
11     output_response = 3
12     window_size = 2
13
14     #CSMIP data is read into "train_file".
15     train_file = Read_CSMIP_data.Read_CSMIP_data(train_folder_path, time_step, output_response,
16     ↪ window_size)
17
18     #The required 3d array will be generated for training
19     datax, datay = train_file.generate_3d_array()
20
21     _, n_step, featurex = datax.shape
22     _, n_step, featurey = datay.shape
23
24     print(featurex)
25     print(featurey)
```

```

26 print(n_step)
27
28 n_unit = 30
29 n_epoch = 10000
30 min_lr = 1e-5 #1e-5
31
32 inputs = Input(shape=(n_step, featurex))
33 #print (inputs)
34 outputs = TimeDistributed(Dense(featurey))(LSTM(n_unit, return_sequences=True)(inputs))
35 print (outputs)
36
37 model1 = Model(inputs=inputs, outputs=outputs)
38
39 modelPath = r'C:\Users\BRACE2\Desktop\CSMIP\model\'
40 mcp_save = ModelCheckpoint(modelPath+ "vanilla_" +str(n_unit)+".hdf5", save_best_only=True,
    ↳ monitor='loss', mode='min')
41 reduce_lr_loss = ReduceLROnPlateau(monitor='loss', factor=0.9, patience=500, verbose=1,
    ↳ mode='min', min_lr=min_lr) #1e-5 50
42 early_stopping = EarlyStopping(monitor='loss', patience=20000, verbose=False,
    ↳ restore_best_weights=True) #2000
43
44 #adam = optimizers.Adam(lr=1e-5)
45 model1.compile(optimizer='adam', loss='mean_squared_error')
46 starttime = time.time()
47 # history = model1.fit(datax, datay, shuffle=True, epochs=n_epoch, verbose=1,
    ↳ callbacks=[mcp_save, reduce_lr_loss, early_stopping])
48
49 history = model1.fit(datax, datay, shuffle=True, epochs=n_epoch, verbose=1)
50
51 runtime = time.time()-starttime
52
53 dictionary = {}
54 dictionary['model'] = model1
55 dictionary['history'] = history
56 dictionary['runtime'] = runtime
57 save_model_dict(dictionary, "vanilla_"+str(n_unit)+"_addFC", modelPath)

```

After training the LSTM model, **Evaluation.py** can be used to evaluate the accuracy of the prediction. Please update the directory of **train_folder_path** (line 32), **test_folder_path** (line 33) and **modelPath** (line 48). The **train_folder_path** and **test_folder_path** are used to store all CSMIP training and testing data respectively, and the **modelPath** is used to load the trained LSTM model. Please also update the name of the trained LSTM model (e.g., "vanilla_30_addFC", no need to include "_model.h5"). The number of time steps may also need to be updated (line 10), depending on the CSMIP data. Lines 60-92 are used to plot the acceleration responses of the predicted results and the target results for both training and testing sets. The "3" in lines 62, 70, 79 and 87 indicates that the index of the interested CSMIP data (it should be the third item in the **train_folder_path** and **test_folder_path**). If required, please change the number to plot the interested CSMIP dataset. Lines 99-104 are used to calculate the correlation coefficient of the predicted and the target results. It should be noted that *np.corrcoef()* generates a 2×2 coefficient matrix, which shows 1.0 in both diagonal elements, indicating that each dataset is perfectly correlated with itself, and ≤ 1.0 in the off-diagonal elements, indicating that how the two arrays are correlated with each other. The code (lines 99-104) automatically extracts the off-diagonal element, showing the correlation between the predicted and the target results. For the details of calculation, please go to: <https://numpy.org/doc/stable/reference/generated/numpy.corrcoef.html>

Example: Evaluation.py

```

1 ...
2 def save_model_dict(dictionary, name, modelPath):
3     dictionary['model'].save(modelPath+r"model_response\\"+name+"_model.h5")
4     f = open(modelPath+r"model_response\\"+name+"_history.pkl", "wb")
5     pickle.dump(dictionary['history'].history, f)
6     f.close()
7     f = open("runtime.txt", "a")
8     f.write(name+" runtime: ")
9     f.write(str(dictionary['runtime']/60))

```

```

10     f.write("\n")
11     f.close()
12
13 def load_model_dict(name, modelPath):
14     dictionary = {}
15     try:
16         model = load_model(modelPath+r"model_response\\"+name+".hdf5")
17     except:
18         model = load_model(modelPath+r"model_response\\"+name+"_model.h5")
19     dictionary['model'] = model
20     f = open(modelPath+r"model_response\\"+name+"_history.pkl", 'rb')
21     history = pickle.load(f)
22     f.close()
23     dictionary['history'] = history
24     return dictionary
25
26 if __name__ == "__main__":
27
28     #specify the train_folder_path and the test_folder_path
29     #specify time step (each seismic event has different timesteps, so need to define manually)
30     #specify output_response (e.g. Accel data in Channel 3, 5, 8 => then 3)
31     #specify window size (stack size, please refer to Zhang2019 for details)
32     train_folder_path = r'C:\Users\BRACE2\Desktop\CSMIP\data\training'
33     test_folder_path = r'C:\Users\BRACE2\Desktop\CSMIP\data\testing'
34
35     time_step = 7200
36     output_response = 3
37     window_size = 2
38
39     #CSMIP data is read into "train_file" and "test_file".
40     train_file = Read_CSMIP_data.Read_CSMIP_data(train_folder_path, time_step, output_response,
41     ↪ window_size)
42     test_file = Read_CSMIP_data.Read_CSMIP_data(test_folder_path, time_step, output_response,
43     ↪ window_size)
44
45     #The required 3d array will be generated for training/testing
46     datax, datay = train_file.generate_3d_array()
47     testx, testy = test_file.generate_3d_array()
48
49     #load the trained model
50     modelPath = r'C:\Users\BRACE2\Desktop\CSMIP\model\'
51     model_dict = load_model_dict("vanilla_30_addFC", modelPath)
52     model = model_dict['model']
53
54     #perform prediction and load the results in datapredict and testpredict
55     datapredict = model.predict(datax)
56     testpredict = model.predict(testx)
57     print("train_loss:")
58     print(model.evaluate(datax, datay, verbose=0))
59     print("test_loss:")
60     print(model.evaluate(testx, testy, verbose=0))
61
62     #Sample Plot of training data
63     plt.figure()
64     plt.plot(model.predict(datax)[3,:,0], color='blue', lw=1.0)
65     plt.plot(datay[3,:,0], ':', color='red', alpha=0.8, lw=1.0)
66     plt.title('Training Set: 3rd Floor Acceleration (x-direction)')
67     plt.legend(["Predicted", "Real"])
68     plt.xlabel("Time Step")
69     plt.ylabel("Acceleration (cm/sec$^2$)")
70
71     plt.figure()
72     plt.plot(model.predict(datax)[3,:,1], color='blue', lw=1.0)

```

```

71 plt.plot(datay[3,:,1],':', color='red', alpha=0.8, lw=1.0)
72 plt.title('Training Set: Roof Acceleration (x-direction)')
73 plt.legend(["Predicted", "Real"])
74 plt.xlabel("Time Step")
75 plt.ylabel("Acceleration (cm/sec$^2$)")
76
77 #Sample Plot of testing data
78 plt.figure()
79 plt.plot(model.predict(testx)[3,:,0], color='blue', lw=1.0)
80 plt.plot(testy[3,:,1],':', color='red', alpha=0.8, lw=1.0)
81 plt.title('Testing Set: 3rd Floor Acceleration (x-direction)')
82 plt.legend(["Predicted", "Real"])
83 plt.xlabel("Time Step")
84 plt.ylabel("Acceleration (cm/sec$^2$)")
85
86 plt.figure()
87 plt.plot(model.predict(testx)[3,:,1], color='blue',lw=1.0)
88 plt.plot(testy[3,:,0],':', color='red', alpha=0.8, lw=1.0)
89 plt.title('Testing Set: Roof Acceleration (x-direction)')
90 plt.legend(["Predicted", "Real"])
91 plt.xlabel("Time Step")
92 plt.ylabel("Acceleration (cm/sec$^2$)")
93
94 # Correlation Coefficient
95 # Note: The resulting matrix from np.corrcoef shows this by having 1.0
96 # in both diagonal elements, indicating that each array is perfectly correlated
97 # with itself, and < 1.0 in the off-diagonal elements, indicating that how the two arrays
98 # are correlated with each other.
99 print("training corr")
100 train_corr = np.corrcoef(datapredict.flatten(), datay.flatten())[0,1]
101 print(train_corr)
102 print("testing corr")
103 test_corr = np.corrcoef(testpredict.flatten(), testy.flatten())[0,1]
104 print(test_corr)
105
106
107 # Error - evaluate the error between the predicted result and the real result
108 errors = np.array([])
109
110 x = (datapredict[:, :, 0] - datay[:, :, 0]) / np.max(np.abs(datay[:, :, 0]),
111 ↪ axis=1).reshape((-1,1))
112 hist = np.histogram(x.flatten(), np.arange(-0.2, 0.201, 0.001))[0]
113 errors = np.append(errors, hist)
114 x = (datapredict[:, :, 1] - datay[:, :, 1]) / np.max(np.abs(datay[:, :, 1]),
115 ↪ axis=1).reshape((-1,1))
116 hist = np.histogram(x.flatten(), np.arange(-0.2, 0.201, 0.001))[0]
117 errors = np.append(errors, hist)
118 x = (testpredict[:, :, 0] - testy[:, :, 0]) / np.max(np.abs(testy[:, :, 0]),
119 ↪ axis=1).reshape((-1,1))
120 hist = np.histogram(x.flatten(), np.arange(-0.2, 0.201, 0.001))[0]
121 errors = np.append(errors, hist)
122
123 errors = errors.reshape((-1, 4, 400))
124 np.save("errors_new.npy", errors)
125 error = np.load("errors_new.npy")
126
127 print(error.shape)
128
129 # Print the error graph, a better result will lead to an error curve centralized to 0.

```

```

130 plt.figure()
131 plt.plot(np.arange(-20, 20, 0.1), error[0][0] / (np.sum(error[0][0]) * 0.001))
132 plt.plot(np.arange(-20, 20, 0.1), error[0][1] / (np.sum(error[0][1]) * 0.001))
133 plt.legend(["Third floor", "Roof"])
134 plt.xlim(-20,20)
135 plt.xlabel("Normalized Error (%)")
136 plt.ylabel("PDF")
137 plt.title('Training Set')
138
139 plt.figure()
140 plt.plot(np.arange(-20, 20, 0.1), error[0][2] / (np.sum(error[0][2]) * 0.001))
141 plt.plot(np.arange(-20, 20, 0.1), error[0][3] / (np.sum(error[0][3]) * 0.001))
142 plt.legend(["Third floor", "Roof"])
143 plt.xlim(-20,20)
144 plt.xlabel("Normalized Error (%)")
145 plt.ylabel("PDF")
146 plt.title('Testing Set')

```

Lines 130-146 are used to generate the normalized error curves of different floors. The errors of all the response cases and time steps are collected and the distribution can be plotted. The more the distribution is centered around 0.0, the better the model.