

National University of Singapore
College of Design & Engineering - ECE

EE4400 Data Engineering and Deep Learning
Tutorial 3 - LSTM Solutions

Q1. LSTMs solve the problem using a unique additive gradient structure that includes direct access to the forget gate's activations, enabling the network to encourage desired behaviour from the error gradient using frequent gates updates on every time step of the learning process.

Q2. The key idea is that the network can learn what to store in the long-term state \mathbf{c}_t , what to throw away, and what to read from it:

- use forget gate to drop, then add what input gate selected
- the transformed long term state is combined with the output to form short term state \mathbf{h}_t (equal to output \mathbf{y}_t).

Q3. Refer to the Jupyter notebook with results and explanations.

```
In [242]: 1 # Tutorial 3 Question 3
2
3 # import Libraries
4
5 import tensorflow as tf
6 from tensorflow import keras
7 from tensorflow.keras.models import Sequential
8 from tensorflow.keras.layers import SimpleRNN
9 from tensorflow.keras.layers import LSTM
10 from tensorflow.keras.layers import BatchNormalization
11 from tensorflow.keras.layers import Dense
12
13 import numpy as np
14 from numpy import array
15 import os
16 import pandas as pd
17 import matplotlib.pyplot as plt
```

```
In [243]: 1 # define input sequence
2
3 raw_seq = [10, 20, 30, 40, 50, 60, 70, 80, 90]
4
5 # choose a number of time steps
6 seq_len = 3
7
8 # choose number of future time steps to predict
9 n_steps = 2
```

```
In [244]: 1 def split_sequence(sequence, seq_len, n_steps):
2     X, y = list(), list()
3     for i in range(len(sequence)):
4         # find the end of this pattern
5         end_ix = i + seq_len
6         # check if we are beyond the sequence
7         if end_ix + n_steps > len(sequence):
8             break
9         # gather input and output parts of the pattern
10        seq_x, seq_y = sequence[i:end_ix], sequence[end_ix:end_ix+n_steps]
11        X.append(seq_x)
12        y.append(seq_y)
13    return array(X), array(y)
```

```
In [245]: 1 # split into sequences
2
3 seq_X, seq_Y = split_sequence(raw_seq, seq_len, n_steps)
4
5 # summarize the data
6 for i in range(len(seq_X)):
7     print(seq_X[i], seq_Y[i])
```

```
[10 20 30] [40 50]
[20 30 40] [50 60]
[30 40 50] [60 70]
[40 50 60] [70 80]
[50 60 70] [80 90]
```

```
In [246]: 1 # split the sequence generated into final training X and Y sequence to be fed to model
2
3 X = seq_X.reshape((seq_X.shape[0], seq_X.shape[1], 1))
4 Y = seq_Y
5
6 print(X.shape, Y.shape)
```

```
(5, 3, 1) (5, 2)
```

```
In [247]: 1 #RNN
2 model = Sequential()
3 model.add(SimpleRNN(50, activation='sigmoid', input_shape=(seq_len, 1)))
4 model.add(Dense(n_steps))
5 model.compile(optimizer='adam', loss='mse')
6
7 print(model.summary())
8
9 # model training
10
11 history = model.fit(X, Y, epochs= 1000, verbose=0)
```

Model: "sequential_2"

Layer (type)	Output Shape	Param #
=====		
simple_rnn_1 (SimpleRNN)	(None, 50)	2600
dense_1 (Dense)	(None, 2)	102
=====		
Total params: 2,702		
Trainable params: 2,702		
Non-trainable params: 0		
None		

```
In [248]: 1 # Test Prediction
2
3 X_test = array([70, 80, 90])
4 X_test = X_test.reshape((1, seq_len, 1))
5 Y_test = model.predict(X_test)
6 print(Y_test)

1/1 [=====] - 0s 131ms/step
[[39.35645 43.52418]]
```

```
In [249]: 1 # Batch Normalization
2
3 model = Sequential()
4 model.add(BatchNormalization(input_shape=(seq_len, 1)))
5 model.add(SimpleRNN(50, activation='sigmoid'))
6 # note: using ReLU instead of sigmoid improves RNN even without BN
7 model.add(BatchNormalization())
8 model.add(Dense(n_steps))
9 model.compile(optimizer='adam', loss='mse')
10
11 print(model.summary())
```

Model: "sequential_3"

Layer (type)	Output Shape	Param #
=====		
batch_normalization_1 (Batch Normalization)	(None, 3, 1)	4
simple_rnn_2 (SimpleRNN)	(None, 50)	2600
batch_normalization_2 (Batch Normalization)	(None, 50)	200
dense_2 (Dense)	(None, 2)	102
=====		
Total params: 2,906		
Trainable params: 2,804		
Non-trainable params: 102		
None		

For Simple RNN layer:

Number of parameters is given by the formula:

(num_features + num_units)* num_units + num_units

num_units = equals the number of units in the RNN

num_features = equals the number features of your input

Now you have two things happening in your RNN.

First you have the recurrent loop, where the state is fed recurrently into the model to generate the next step. Weights for the recurrent step are:

recurrent_weights = num_units*num_units

The secondly you have new input of your sequence at each step.

input_weights = num_features*num_units

So now we have the weights, what's missing are the biases - for every unit one bias:

biases = num_units*1

So finally we have the formula:

$(\text{num_features} + \text{num_units}) * \text{num_units} + \text{num_units} = (50 + 1) * 50 + 50 = 2600$

For Batch Normalization Layer:

Batch Normalization computes 4 parameters (beta, gamma, mean, std deviation) for each feature. So no. parameters = $4 * 50 = 200$

Note that mean and std deviation are not trainable parameters as they are calculated.

(Work out the first BN layer yourself.)

For Output layer:

$**\text{output_channel_number} * (\text{input_channel_number} + 1) = 2 * (50 + 1) = 102$

```
In [250]: 1 # model training
          2
          3 history1 = model.fit(X, Y, epochs= 1000, verbose=0)
```

```
In [251]: 1 # Test Prediction
          2
          3 X_test = array([70, 80, 90])
          4 X_test = X_test.reshape((1, seq_len, 1))
          5 Y_test = model.predict(X_test)
          6 print(Y_test)
```

```
1/1 [=====] - 0s 161ms/step
[[ 99.656166 104.80522 ]]
```

```
In [252]: 1 #LSTM
          2 model = Sequential()
          3 model.add(LSTM(50, activation='relu', input_shape=(seq_len, 1)))
          4 model.add(Dense(n_steps))
          5 model.compile(optimizer='adam', loss='mse')
          6
          7 print(model.summary())
```

Model: "sequential_4"

Layer (type)	Output Shape	Param #
=====		
lstm (LSTM)	(None, 50)	10400
dense_3 (Dense)	(None, 2)	102
=====		
Total params: 10,502		
Trainable params: 10,502		
Non-trainable params: 0		
=====		
None		

For LSTM layer:

Number of parameters is given by the formula:

$4 * ((\text{num_features} + \text{num_units}) * \text{num_units} + \text{num_units})$

4 - The factor 4 is multiplied for LSTM because LSTM computes 4 types of weights per layer (W_{forget} , W_{input} , W_{output} , W_{cell})

num_units = equals the number of units in the LSTM

num_features = equals the number features of your input

recurrent_weights = $\text{num_units} * \text{num_units}$

input_weights = $\text{num_features} * \text{num_units}$

So now we have the weights, whats missing are the biases - for every unit one bias:

biases = $\text{num_units} * 1$

So finally we have the formula:

$4 * ((\text{num_features} + \text{num_units}) * \text{num_units} + \text{num_units}) = 4 * ((50 + 1) * 50 + 50) = 10400$

For Output layer:

$2 * (\text{input_channel_number} + 1) = 2 * (50 + 1) = 102$

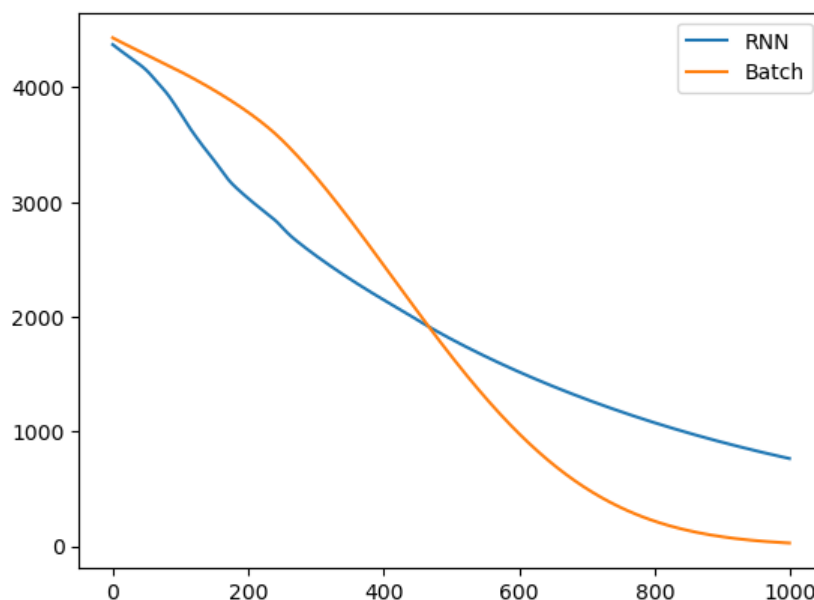
```
In [253]: 1 # model training
          2
          3 history2 = model.fit(X, Y, epochs= 1000, verbose=0)
```

```
In [254]: 1 # Test Prediction
          2
          3 X_test = array([70, 80, 90])
          4 X_test = X_test.reshape((1, seq_len, 1))
          5 Y_test = model.predict(X_test)
          6 print(Y_test)

1/1 [=====] - 0s 158ms/step
[[101.01683 111.027054]]
```

```
In [255]: 1 # plot training loss
          2 plt.plot(history.history["loss"])
          3 plt.plot(history1.history["loss"])
          4 plt.legend(("RNN", "Batch"))
```

Out[255]: <matplotlib.legend.Legend at 0x18ad701dbb0>



From the plot above we observe the loss function of the model without and with batch normalization.

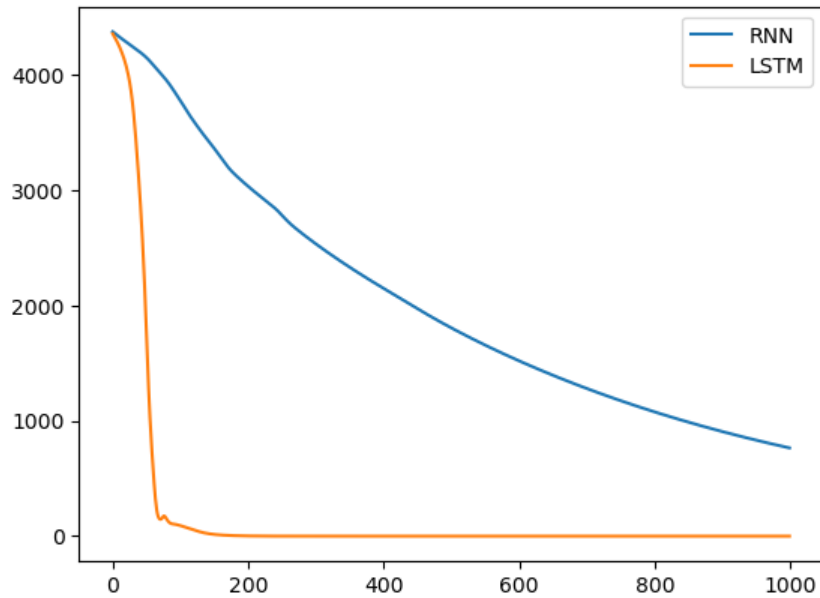
Batch normalization (BN) does not work well with RNN across time steps because the operation changes the characteristics of the sequential data presented in an RNN. The repeated application of the gain and offset in BN during the recurrent procedure may also lead to exploding or vanishing gradients.

However, batch normalization between RNN layers seems to be helpful when the activation function in the RNN is sigmoid.

Layer normalization can work well in an RNN as it captures statistics across inputs at a single time step.

```
In [256]: 1 plt.plot(history.history["loss"])
          2 plt.plot(history2.history["loss"])
          3 plt.legend(("RNN", "LSTM"))
```

Out[256]: <matplotlib.legend.Legend at 0x18ac5af2ca0>



From the plot above we can note that:

The performance of LSTM is better than RNN. To truly observe the benefit of LSTM, we can predict much further into the future. In our example, we only predict 2 time steps into the future. LSTM models are capable of predicting much further into the future than RNN with reasonable accuracy.

```
In [257]: 1 keras.backend.clear_session()
```