

EEL6935
Programming Assignment 3
Abhinav Reddy Podduturi 1599-2633

GitHub Link: <https://github.com/Abhinav-Reddy/BigDataAssignment3>

Google Doc Link:

https://docs.google.com/document/d/1r1lf5l1Cz3kktEDQF49AB7DxSJpONDY_zpqQEVwFa3w/edit?usp=sharing

Architecture Overview:

Architecture consists of three main layers embedding layer, LSTM and sigmoid layer. DNA sequence is converted into a embeddings and fed to LSTM. Output from LSTM is fed to sigmoid layer which outputs the probability of the input sequence being a TFBS. LSTM was chosen mainly for their good performance on sequences and their ability to remember long term dependencies.

Preprocessing:

Each nucleotide was assigned a unique number between 1 and 4. One hot vector representation of each word was used as embeddings in the embeddings layer.

LSTM:

Embeddings from the embedding layer are given as input to the LSTM layer. This layer outputs a vector of length 8.

Sigmoid Layer:

Output vector from LSTM is fed to sigmoid layer which outputs the final probability.

Training:

Adadelta optimizer was used with a batch size of 8 and gradient clipping norm of 1.25. Model was trained for 10 epochs with a training set of 1900.

Results:

Secured 8th position and a score of 84.615 on Kaggle.

Tried following enhancements:

Different things were tried to improve the performance of the system.

- 1) Different optimizers such as SGD, adanorm were used to check their impact on the performance on the system. No significant improvements were observed.
- 2) A probability of 0.4 was used as a cutoff to classify the sequences. This improved the F1 score by 1 percent.

- 3) Deep LSTM's with different number of layers from 2 to 4 were tried. None of them could offer any significant improvements.
- 4) Instead of considering one letter in the sequence as an individual event to LSTM, two adjacent letters were clubbed together and given as input to LSTM. Correspondingly one hot vector size was increased to 16. This gave an improvement of 2 percent. Also clubbing 3 or 4 letters didn't provide any improvement.

Screenshot:

```
[12 13 9 ..., 13 13 13]]
2018-04-09 02:17:55.022473: I tensorflow/core/common_runtime/gpu/gpu_device.cc:1120] Creating TensorFlow device (/device:GPU:0) -> (device
id: 4fa7:00:00.0, compute capability: 3.7)
Started training
Train on 1900 samples, validate on 100 samples
Epoch 1/10
1900/1900 [=====] - 10s 5ms/step - loss: 0.1839 - acc: 0.7558 - val_loss: 0.1262 - val_acc: 0.8100
Epoch 2/10
1900/1900 [=====] - 10s 5ms/step - loss: 0.1083 - acc: 0.8505 - val_loss: 0.1170 - val_acc: 0.8300
Epoch 3/10
1900/1900 [=====] - 10s 5ms/step - loss: 0.1011 - acc: 0.8621 - val_loss: 0.1096 - val_acc: 0.8400
Epoch 4/10
1900/1900 [=====] - 10s 5ms/step - loss: 0.0975 - acc: 0.8632 - val_loss: 0.1093 - val_acc: 0.8400
Epoch 5/10
1900/1900 [=====] - 10s 5ms/step - loss: 0.0944 - acc: 0.8674 - val_loss: 0.1058 - val_acc: 0.8400
Epoch 6/10
1900/1900 [=====] - 10s 5ms/step - loss: 0.0937 - acc: 0.8695 - val_loss: 0.1081 - val_acc: 0.8600
Epoch 7/10
1900/1900 [=====] - 10s 5ms/step - loss: 0.0920 - acc: 0.8747 - val_loss: 0.1098 - val_acc: 0.8500
Epoch 8/10
1900/1900 [=====] - 10s 5ms/step - loss: 0.0910 - acc: 0.8784 - val_loss: 0.1090 - val_acc: 0.8600
Epoch 9/10
1900/1900 [=====] - 10s 5ms/step - loss: 0.0901 - acc: 0.8795 - val_loss: 0.1063 - val_acc: 0.8600
Epoch 10/10
1900/1900 [=====] - 10s 5ms/step - loss: 0.0892 - acc: 0.8853 - val_loss: 0.1077 - val_acc: 0.8400
Training time finished.
3 epochs in 0:01:38.151121
abhinav@BigData:~/BigDataAssignment3/src$
```