

1. Overview

The networks used this time are all Recurrent Neural Networks.

LSTM → The model consists of a single multi layer Long Short Term Memory with **2** layers. The input size is **160**. The sequence length is **100**. The hidden size is **128**. The batch size is **768**. The classes to classify are **10** for the 10 instrument families. The input and output channels are as follows, LSTM → (160, 128, 2) Fc1 → (128, 10)

BLSTM→ The model consists of a single multi layer Bidirectional Long Short Term Memory with **2** layers. The input size is **160**. The sequence length is **100**. The hidden size is **128**. The batch size is **768**. The classes to classify are **10** for the 10 instrument families. The input and output channels are as follows, BLSTM → (160, 128, 2) Fc1 → (128, 10)

Both networks use the same input size and sequence length as specified in the project statement. The size of the input data is **16,000** since we are looking at only the 1st second of the waveform. This size allows the network to be trained in a reasonable amount of time while looking at the features of one second of the waveform across families.

Findings → RNNs are best suited to data which occurs in a sequence, such as the data we have at hand which is sound data.

→ LSTMs and BLSTMs train faster than CNNs in general.

→ LSTMs train faster relative to BLSTMs as they are not bidirectional, hence a two way computational overhead is eliminated.

→ The accuracy of the LSTM is higher than BLSTM. (More in discussion)

→ Batch size affects the overall accuracy because the smaller the mini batches the more the weight updates. Very large batch size (2048) helped to train the model faster but at the cost of accuracy. An optimal trade-off between accuracy and batch size was to select 768 as the batch size.

→ Hidden size also affects training time and accuracy, as a very large number of hidden units slows down training (with respect to time) and also the model might learn as well due to issue of vanishing gradients.

2. Training

For both the networks, a batch size of 768 was used as this size utilized a significant portion of the memory and also sped up the training times more.

LSTM → The **optimizer** used was **Adam Optimizer**. The criterion is **Cross Entropy Loss (CLE)**.

To train the LSTM, the learning rate was set to **0.003**, this learning rate performed well with the Advanced Net from project 1 and hence was a suitable candidate for this implementation as well. This learning rate performed well and the model decreased **training loss** per epoch as can be seen in the loss curve below.

BLSTM → The **optimizer** used is **Adam Optimizer**. The criterion is **Cross Entropy Loss (CLE)**.

To train the BLSTM, the learning rate was set to **0.003**, same as the LSTM.

The main benefits (while training) with LSTMs and BLSTMs is that they can be trained faster as compared to CNNs. This was definitely seen in the training time per epochs. For our advanced net, we saw per epoch timings of 6:26 vs 3:02 for the BLSTM.

With LSTM the training time per epoch was even less with timings of 2:03 for a batch size of 768.

With 283704 examples in the training dataset and a batch size of 768, the number of iterations per epoch were $283704/768 = 369.40625 \rightarrow 369$ (since last batch is dropped).

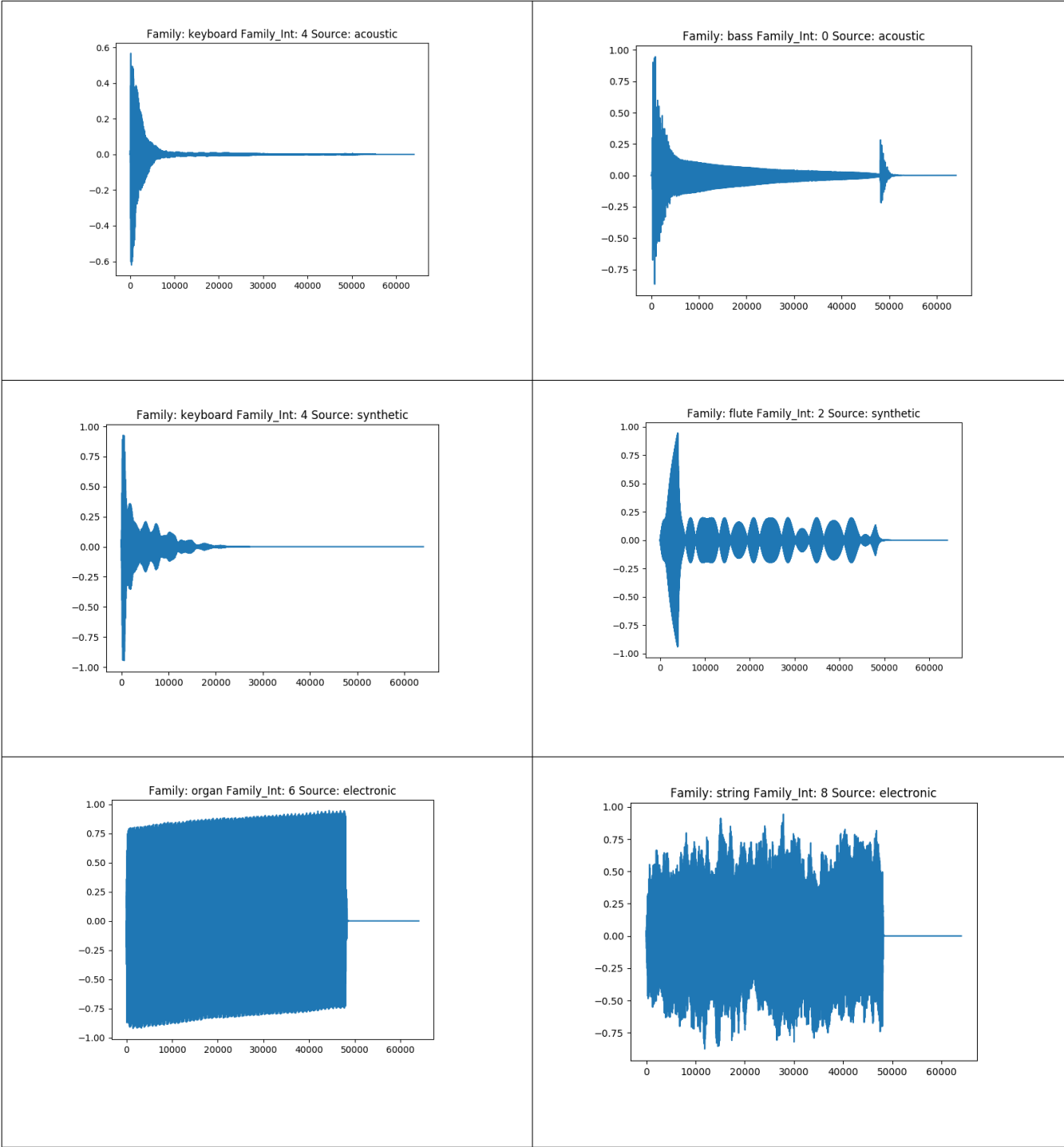
The waveforms presented below are for frame size of 160. Looking at the waveforms for size 160 rather than 16,000 helps to see if the patterns at this size help the model learn better or not.

```
Epoch [3/30], Step [100/369], Loss: 1.5675
Training Accuracy : 41 %
Epoch [3/30], Step [200/369], Loss: 1.5286
Training Accuracy : 43 %
Epoch [3/30], Step [300/369], Loss: 1.4258
Training Accuracy : 44 %
Epoch 3, validation_loss: 1.4782496690750122
Validation Accuracy : 44 %
```

Monitoring progress in an epoch

3. Results

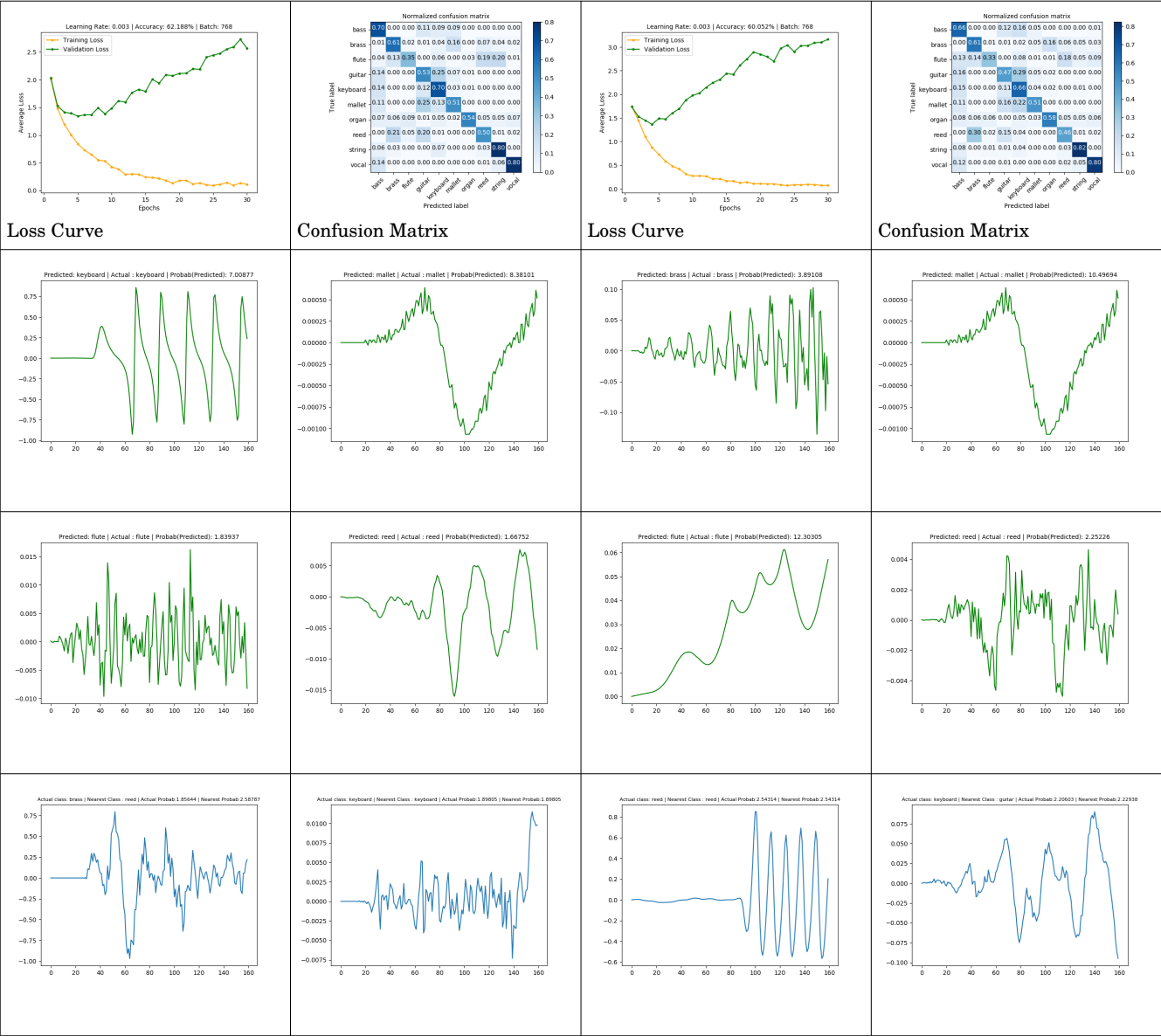
Waveform examples across 3 sources and any 2 instrument families -
(Waveforms shown without transforms picked directly from train dataset)



Accuracy results included in the Loss curves.

LSTM

BLSTM



Waveform for samples where correct class probability is very high. (160 input size)

Waveform for samples near decision boundary. (160 input size)

LSTM

BLSTM

Please zoom to view

```
Epoch [30/30], Step [100/369], Loss: 0.0533
Training Accuracy : 56 %
Epoch [30/30], Step [200/369], Loss: 0.0831
Training Accuracy : 97 %
Epoch [30/30], Step [300/369], Loss: 0.0573
Training Accuracy : 97 %
Epoch 30, Validation Loss: 2.584319372177124
Validation Accuracy : 62 %

Finished Training
Accuracy of the network on the 4096 sound samples: 62 %

/usr/local/lib/python3.6/dist-packages/matplotlib/pyplot.py:522: RuntimeWarning: More than 20 figures have been opened. Figures created through the pyplot interface ('matplotlib.pyplot.figure') are retained until explicitly closed and may consume too much memory. (To control this warning, see the rcParam 'figure.max_open_warning').
max_open_warning, RuntimeWarning)

Confusion Matrix: [[551  0  1  88  70  72  0  0  0  0]
 [ 2 156  0  3 11  42  1 19  9  0]
 [ 0 22 59  0 11  0  5 32 34  1]
 [ 87  2  3 323 152  41  4  1  0  0]
 [103  0  1  83 582  22  5  2  3  0]
 [21  0  0  47 24 96  0  0  0  0]
 [34 27 42  3 24  9 232 22 21 31]
 [ 0 48 11 44  2  1  0 112  2  4]
 [18  0  0  1 20  0  0  1 10 232  0]
 [18  0  0  0  0  0  0  1  8 105]]

Normalized confusion matrix
[[0.70460338 0. 0.00127877 0.11253197 0.08951407 0.09287161
 0. 0. 0. 0. ]
 [0.00783514 0.01704771 0.01355291 0.01172477 0.00319725 0.04705088
```

Overall Accuracy

```
0.0070522 0.02830875 0.04504555 0.89301702]]
Accuracy of 0 : 67 %
Accuracy of 1 : 68 %
Accuracy of 2 : 45 %
Accuracy of 3 : 48 %
Accuracy of 4 : 65 %
Accuracy of 5 : 47 %
Accuracy of 6 : 48 %
Accuracy of 7 : 58 %
Accuracy of 8 : 89 %
Accuracy of 9 : 66 %
```

Per Class Accuracy

```
Epoch [30/30], Step [100/369], Loss: 0.0565
Training Accuracy : 98 %
Epoch [30/30], Step [200/369], Loss: 0.0491
Training Accuracy : 98 %
Epoch [30/30], Step [300/369], Loss: 0.0888
Training Accuracy : 97 %
Epoch 30, validation_loss: 3.16330316543579
Validation Accuracy : 60 %

Finished Training
Accuracy of the network on the 4096 sound samples: 60 %

/usr/local/lib/python3.6/dist-packages/matplotlib/pyplot.py:522: RuntimeWarning: More than 20 figures have been opened. Figures created through the pyplot interface ('matplotlib.pyplot.figure') are retained until explicitly closed and may consume too much memory. (To control this warning, see the rcParam 'figure.max_open_warning').
max_open_warning, RuntimeWarning)

Confusion Matrix: [[514  0  0  92 129 36  1  0  1  9]
 [ 2 155  2  3 12  42 15 12  7]
 [22 23 56  0 13  1  1 30  9 15]
 [ 36  4  3 221 27 32 11  0  1 13]
 [105  0  0 79 479 26 10  3 10 11]
 [ 20  0  0 31 42 95  0  0  0  0]
 [ 36 26 26  1 21 15 269 21 21 29]
 [ 0 68  4 33  9  0  1 103  2  4]
 [ 23  2  4  3 11  0  0  9 229  0]
 [ 10  0  0  1  0  0  0  3  7 105]]

Normalized confusion matrix
[[0.657289 0. 0. 0.11764786 0.16496164 0.04603581
 0.00127877 0. 0.00127877 0.01308093]
 [0.00392157 0.60784314 0.00784314 0.01176471 0.02352941 0.04785882
 0.1647588 0.0582313 0.04785882 0.02745098]
 [0.12941376 0.13579432 0.32941376 0. 0.07647059 0.00588235
 0.00588235 0.17647059 0.05294118 0.00823529]
 [0.2568686 0.0083132 0.0083132 0.4472452 0.28874388 0.05220228
 0.01794454 0. 0.00161132 0.00161132]
 [0.14561107 0. 0. 0.10957884 0.06435506 0.0306103
 0.05495321 0.00416089 0.01388963 0.00130696]
 [0.10638298 0. 0. 0.10489362 0.22340428 0.50531915
 0. 0. 0. 0. ]]
```

Overall Accuracy

```
Normalized confusion matrix
[[0.74377224 0. 0. 0.06998814 0.10083037 0.07829181
 0.00237248 0. 0.00237248 0.00237248]
 [0.16791045 0.64552239 0.0119403 0.01865672 0.0119403 0.02238886
 0. 0. 0.05597015 0.04850746 0.01865672]
 [0.16111111 0.15555556 0.20666667 0.01111111 0.01111111 0.01666667
 0. 0.13888889 0.15555556 0.08333333]
 [0.19354839 0.0030722 0.0015361 0.48694316 0.25192012 0.05376344
 0.00460829 0. 0.0030722 0.0015361]
 [0.17385621 0. 0.00392157 0.11895425 0.64183807 0.04313725
 0.00653595 0.00522876 0.00053595 0. ]
 [0.0990099 0. 0. 0.21782178 0.15841584 0.52475248
 0. 0. 0. 0. ]
 [0.0876494 0.03784861 0.0438247 0.00398466 0.04780876 0.02191235
 0.56772908 0.06374502 0.05178486 0.07171315]
 [0. 0.33333333 0.04273504 0.15384615 0.02136752 0.01282051
 0. 0.35897436 0.04700855 0.02991453]
 [0.07189542 0.00326797 0.00980392 0.0138719 0.0620915 0.
 0. 0.01633987 0.01372549 0.00980392]
 [0.20557376 0. 0. 0. 0. 0. 0. 0. 0. 0. ]
 [0.0078922 0.0078922 0.06382979 0.71631286]]
Accuracy of 0 : 70 %
Accuracy of 1 : 84 %
Accuracy of 2 : 27 %
Accuracy of 3 : 37 %
Accuracy of 4 : 63 %
Accuracy of 5 : 52 %
Accuracy of 6 : 55 %
Accuracy of 7 : 41 %
Accuracy of 8 : 84 %
Accuracy of 9 : 66 %
```

Per Class Accuracy

4. Discussion

Using the Adam optimizer [2] which varies the learning rate with the changes in parameters the curve is more gradual and reflective of expected training behavior with the training loss decreasing after each epoch. However the validation loss decreases at the beginning but then increases. This suggests that the model is overfitting to the train data.

From the confusion matrix, it can be seen that both the LSTM and BLSTM are getting confused at similar classes. For example in the row for bass both the LSTM and BLSTM are confused it for class guitar and keyboard.

In terms of speed and accuracy, as discussed above in the training section, both these models training much faster and also improved upon the accuracy from project1. For project1 the reported accuracy was 57% and with

the these 2 models we are seeing an accuracy of 62% which is a gain of 5%. The per epoch timings were about half of what we saw for advanced net from project 1.

Weight initialization used is random weight initialization. As with our advanced model, on subsequent runs our LSTM and BLSTM model wasn't being thrown off too much in terms of our accuracy, whereas in our advanced model there was a variability of about 5% on subsequent runs. Hence a Xavier weight initialization as with our advanced model didn't seem necessary. [4]

For a run of Epochs → 30 | Learning Rate → 0.003 | Batch Size → 768 | Hidden Units → 128 | Input Size → 100 | Sequence length → 160

The final accuracies were

LSTM → 62 %

BLSTM → 60 %

Finally, dropout layers with a p of 0.5 (drops 50% of the nodes), were added to perform regularization to prevent overfitting and penalizes the loss function but these caused the model not to learn at all and was getting stuck in a local minima (judged by the loss curves). [1] Hence, dropout was later removed and the accuracies presented here are for runs without dropout.

In short both LSTM and BLSTM performed similarly in terms of accuracies, per class accuracies and decision boundaries. This can be seen in the confusion matrices, loss curves and waveform visualizations seen above.

Again it can be seen that the “memory” aspect of RNN helps them to learn the sequence data better. However, the accuracies are still low and this raises questions towards the quality of the dataset which also affected the accuracies in project 1. [5]

Data → Analyzing a few waveforms from the different families(refer per family per source waveforms above (more were analyzed)), a few waveforms look similar to the waveforms from other families.[3] This might be a reason as to why the model confuses one family with another as can be seen with the case of guitar in the confusion matrix. It is difficult to say exactly how many instances of such waveforms are present in the dataset without going through it entirely. The distribution of data is also uneven which dips the per class accuracy.

With the data being “noisy” with waveforms similar across families, makes the models job difficult to differentiate between instrument families and hence the overall accuracy is relatively low.

Since softmax was not used for the fully connected layer, the “probabilities” are at a different metric and not between 0 and 1 per class. But the logic to derive the highest probability remains the same as project 1.

6. References

- [1] Galeone, P. (2017). *Analysis of Dropout*. [online] P. Galeone's blog. Available at: <https://pgaleone.eu/deep-learning/regularization/2017/01/10/anaysis-of-dropout/>.
- [2] Kingma, D. and Ba, J. (2014). *Adam: A Method for Stochastic Optimization*. [online] Arxiv.org. Available at: <https://arxiv.org/abs/1412.6980>.
- [3] Shuvaev, S., Koulakov, A. and Giaffar, H. (2017). *Representations of Sound in Deep Learning of Audio Features from Music*. [online] Arxiv.org. Available at: <https://arxiv.org/pdf/1712.02898.pdf>.
- [4] Hochreiter Sepp, Koulakov, A.(1997). *Long Short Term Memory*. [online] Available at: <http://www.bioinf.jku.at/publications/older/2604.pdf>
- [5] Yin, W., Kann, K., Yu, M. and Schütze, H. (2017). *Comparative Study of CNN and RNN for Natural Language Processing*. [online] Arxiv.org. Available at: <https://arxiv.org/abs/1702.01923>