

Machine Learning II
Final Project
Professor Amir Jafari

Sequence Prediction with Recurrent Neural Networks: Predicting the Next Airport a Plane Will Visit

Group 1

*Gayathri Chandrasekaran,
Jonathan Giguere,
Patrick Maus*

Table of Contents

Introduction	4
<i>Description of the Data Set</i>	4
Summary	5
Exploratory Data Analysis	5
Analysis of Autocorrelation and Stationarity	6
My Contribution to the Project	8
Percent of Code from the Internet	8
<i>Description of the Deep Learning Networks and Training Algorithms</i>	8
Summary	8
Recurrent Neural Networks (RNN)	9
Long Short Term Memory (LSTM)	9
Gated Recurrent unit (GRU)	10
Sequence2Sequence	11
<i>Experimental Setup</i>	12
Data Preprocessing	12
Framework	13
Model Implementation	13
Single layer LSTM Network:	13
Stacked LSTM Network:	13
Ensemble Method:	14
Sequence2sequence	14
Model Evaluation Metrics	14
<i>Results</i>	16
Random Sequences Output:	16
Single Layer Lstm Model:	16
Stacked Lstm Network Results:	17
Stacked Ensemble Methods:	18
Sequence2Sequence Method:	19
Summary	21
Summary and Conclusions	22
Evaluating Stationarity and Autocorrelation	22
Repeated Sequences	22
Sequence Length	22
Looking at Other Airlines	22
Best Model	22
	2

<i>References</i>	23
<i>Appendix A: Code for Project</i>	24
Appendix B: Performance for All Standard RNN Models	25
Appendix C: Performance for All Sequence2sequence Models	29

Introduction

Our group was interested in exploring Recurrent Neural Networks (RNNs), so we chose a project focused on predicting the next value in a given sequence. Specifically, we analyzed 3.6 million US airline flights over a period of 6 months. Each record in our dataset includes the airline name, a unique identifier for each aircraft, the origin and destination airports, and the time of the flight. We defined a flight as one unique aircraft flying from one of about 350 airports to another. Our research question is “Given a sequence of N airports visited by a unique plane, can we predict the next airport ($N+1$)?”

We believe this approach is applicable to numerous different problems beyond predicting the next airport. Our world is awash in devices that record their time and location, leading to an explosion of geospatial data that can inform everything from advertising to pandemic responses. Often, the data is too dense for traditional geospatial analysis methods and a common approach is to represent a dataset as a network or sequence of known locations visited.¹ Applying a similar methodology as the one proposed will likely lead to additional insights in multiple fields and business cases.

For our project, we first conducted extensive exploratory data analysis to better understand the distribution of our data across airlines, airports, and airplanes. We ultimately decided to segment our data to subsets grouped by airlines as it appeared aircraft within different airlines followed their own patterns. Next, we examined the data for stationarity and autocorrelation to ensure there was a signal embedded within the sequences that our models could identify and then learn for prediction. We found strong evidence that our data was stationary and had “colored” signals rather than white noise.

Next, we conducted preprocessing on the data to create equal-length tokenized sequences which could be used as inputs for an RNN. Since we were unsure of the ideal sequence length for this problem, we needed to test multiple different airline subsets with multiple different sequence lengths as a data input for each of our models. To facilitate efficient model development, training, and evaluation in an environment with multiple different input datasets, we developed a data ingest and preprocessing pipeline using a series of Python scripts which employed different “model_name” and “run_name” variables to track activity across different RNN models (LSTM, GRU, seq2seq, etc) and airline/sequence length “runs”.

Our team developed, trained, and evaluated models including LSTM, GRU, 1dCONV, bi-directional LSTM, sequence2sequence, etc.

To evaluate our models, we built another script to plot the training history and visualize the precision, recall and F1 measure of our models’ predictions for the target ($N+1$) airport.

¹ Wang, Senzhang, J. Cao and Philip S. Yu. “Deep Learning for Spatio-Temporal Data Mining: A Survey.” *ArXiv abs/1906.04928* (2019): n. pag.

Description of the Data Set

Summary

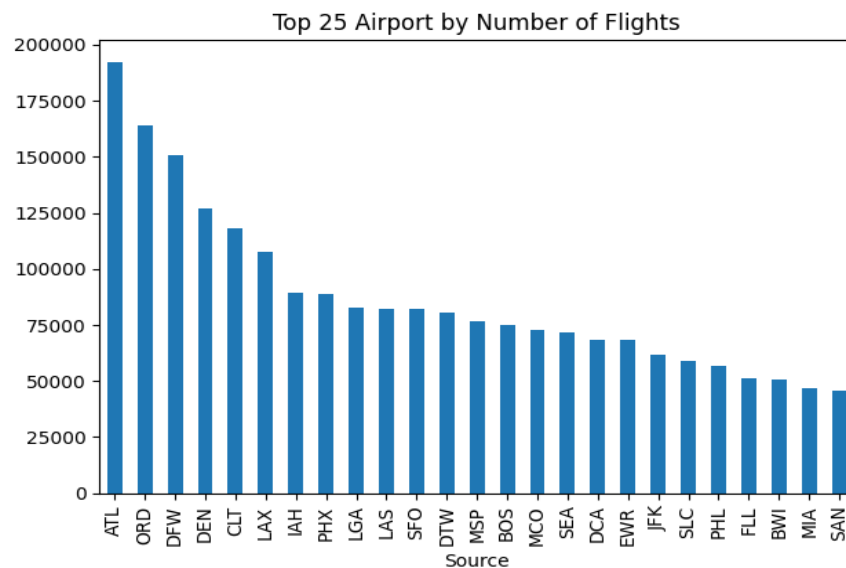
The data set was downloaded from the Bureau of Transportation Statistics, a US Department of Transportation agency charged to collect, analyze, and store transportation statistics.² Using six months of data from September 2019 to February 2020, we reviewed about 3.7 million flights from 5,716 unique aircraft. We choose to only use data up to February 2020 to avoid impacts from the Coronavirus pandemic to affect the sequence patterns. Within our data, there are 19 different airlines, but the top five represents about 2/3 of the entire data. Several sample rows are included below.

Table 1. Data sample from February 1, 2020.

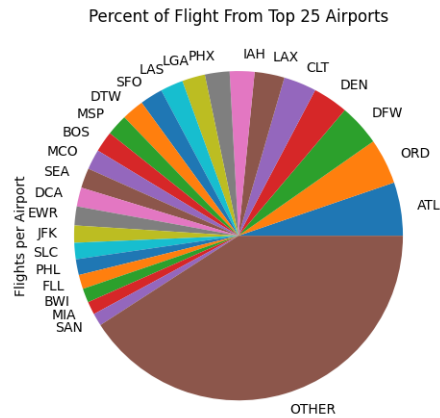
FL_DATE	OP_UNIQUE_CARRIER	TAIL_NUM	OP_CARRIER_FL_NUM	ORIGIN	DEST	DEP_TIME	ARR_TIME
2020-02-01	MQ	N269NN	3825	ORD	TUL	1646	1820
2020-02-01	MQ	N908AE	3829	JFK	BNA	1336	1458
2020-02-01	MQ	N663AR	3831	GNV	MIA	844	1020
2020-02-01	MQ	N618AE	3833	DFW	SJT	852	955
2020-02-01	MQ	N618AE	3833	SJT	DFW	1024	1132

Exploratory Data Analysis

Next we reviewed the distribution of flight activity across the airports in the data. Our analysis found that the data distribution of number of flights per airport was left-tailed, with a few airports much busier than others.



² <https://www.bts.gov>



The centralization was greater in some airports than others. For example, Delta Airlines has a major hub in Atlanta, which accounts for about 25% of all flight activity for Delta flights. As we built, trained, and evaluated a model that predicts the next sequence, we had to be aware of this centralization because our models could 1) generate artificially higher accuracy scores by more frequently predicting the busiest airport or 2) learn to predict flight activity out of the largest airports but poorly generalize to smaller airports.



Figure 1. A weighted network plot of Delta's flight activity in 2018 snapped to a map of the continental US. Darker, thicker lines represent greater numbers of flights and clearly show Delta's employment of a "hub and spoke" model.

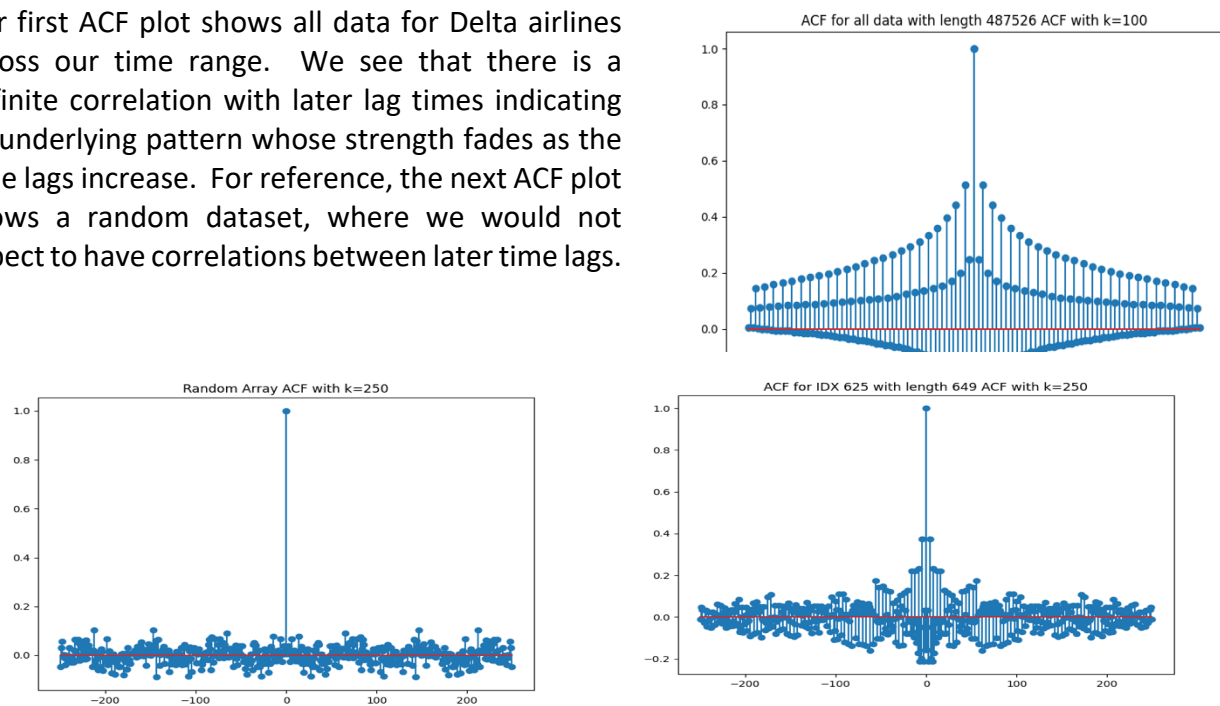
Analysis of Autocorrelation and Stationarity

Once we completed our exploratory data analysis, we next focused on ensuring our data was fit for use in sequence prediction. Any machine learning model depends on an actual "signal" to be present in the data, which by definition has a pattern that can be mapped from the domain (or training data) to the range (or test data). In our case, if aircraft do not follow an observable pattern when flying from airport to airport, we would be trying to predict randomness, or "white noise."

Our first check was to plot the autocorrelation of the entire dataset. An autocorrelation plot shows how well a signal correlates with itself at equal time lags. If the plot shows an impulse of 1 at time lag 0 but then very little or known correlation at subsequent time lags, the data is likely

to be random white noise. On the other hand, an autocorrelation plot with strong correlation for multiple lags indicates that previous events influence the current event and suggest an underlying pattern.

Our first ACF plot shows all data for Delta airlines across our time range. We see that there is a definite correlation with later lag times indicating an underlying pattern whose strength fades as the time lags increase. For reference, the next ACF plot shows a random dataset, where we would not expect to have correlations between later time lags.



We also surveyed a selection of random aircraft to see if individual aircraft ACF plots showed good autocorrelation. We found multiple examples similar to the above figure for IDX 625, which corresponds to aircraft tail number 'N849DN'. This shows strong evidence of autocorrelation with several time lags that can decrease in intensity as the time lags increase.

Next, we examined stationarity. A stationary time series is a dataset that lacks a trend or seasonality. Most models, including our RNNs, assume a dataset to be stationary so we should ensure our data is most likely stationary. Since our sequences are drawn from a dataset that likely has both trend and seasonality *in the volume of activity*, we need to be particularly cautious. Fortunately, our selection of sequence order somewhat insulates our data from underlying trend or seasonality. For example, Aircraft A could make a flight between three airports every week, Aircraft B could make multiple flights to multiple airports in the first month before going into maintenance, and Aircraft C could be dormant for the first 5 months and only active the last month of our data. Because we only care about the ordered aspects of airports visited for each flight, our data is somewhat protected from trend and seasonality effects on the volume of activity.

To test for stationarity, we used a similar methodology to the ACF plots. We first used an Augmented Dickey-Fuller (ADF) test against the whole dataset, subsets for different airlines, and finally on individual aircraft. Our results were also similar to the ACF. We found that across the entire data and at the airline levels, we obtained very low test statistics and small p-values which

suggest the data is stationary. For example in the Delta airlines data, we obtained a test statistic of -30.293, far below our cutoff of -2.862 for a .05 p-value. This suggests that we can reject the null hypothesis the data is not stationary and treat it as stationary.

When we examined individual aircraft, we found similar results to the ACF plots. Most aircraft yielded low ADF test statistics, suggesting we can treat the data as stationary. There were several ADF tests which yielded high test statistics and p-values, suggesting some aircraft trips were non-stationary. However, this was a very small number of aircraft. In the Delta Airlines data, 94% of all aircraft yielded a p-value less than 0.05 during the ADF Test. Since this was a small number of aircraft, we decided that we could treat the data overall as stationary.

My Contribution to the Project

In this project, I focused on setting up the data pipeline, conducting EDA of the data, exploring the data for autocorrelation and stationarity, building an initial baseline model, analyzing the results, and producing evaluation metrics. I built the pipeline to ingest the raw data files from the Bureau of Transportation Statistics, including the methodology to parse a flat file of individual flights into an ordered sequence of airports visited. Some of this code was reused from an exploration of ship traffic during a Directed Studies course over the Summer 2020 semester.

I also conducted the analysis into autocorrelation and stationarity. All of the code was written by me, but some of it was drawn from previous work in my Multivariable Modeling class from Spring of 2020. In the preprocessing script, I did use this stack overflow answer to generate a broadcast function (<https://stackoverflow.com/questions/40084931/taking-subarrays-from-numpy-array-with-given-stride-stepsize>). I also used the Keras documentation ([Tokenizer](#)) to build the tokenized sequences.

For the modeling script, I leverage previous work in this class and used some sources from Professor Jafari's GitHub (<https://github.com/amir-jafari/Deep-Learning>). For the modeling_eval script, I used this guide (<https://www.machinecurve.com/index.php/2019/10/22/how-to-use-binary-categorical-crossentropy-with-keras/>) as a source for building the training loss/validation plot.

Percent of Code from the Internet

Most of the code was created by me for this project or in previous classes. The total number of lines of code I wrote across all the scripts is 518 lines. I used about 30 lines of code from the internet, mostly focused on using the Keras tokenizer to generate sequences. I modified about another 50 lines of code. That means about 15% was from the internet or modified while 85% was generated by me for this project.

Description of the Deep Learning Networks and Training Algorithms

Summary

As mentioned previously, to forecast the next airport a plane will visit is dependent on the previous travel history, time component of the data series plays a vital role and this needs to be captured. Hence, we decided to go for recurrent neural network models where prior information

from previous states helps for sequence prediction. We implemented three major recurrent network models for our data. These models include: LSTM, GRU, and sequence2sequence. This section will briefly describe how each of these models works at a high level.

Recurrent Neural Networks (RNN)

RNN is one of the powerful neural networks that is capable of capturing the temporal dependencies present in the data. The output of the network not only depends on the current input but also on the previous output. This network carries a delay component (previous outputs) as a memory unit to the hidden layers. This helps the information to flow from one step to the next.

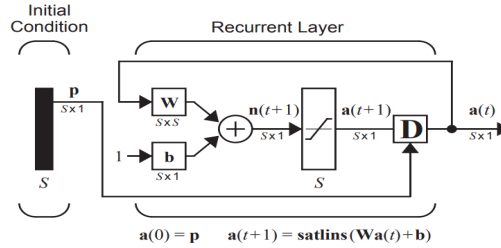
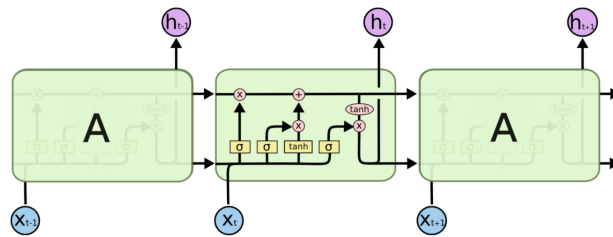


Figure 2. Recurrent Network

Long Short Term Memory (LSTM)

When the sequences are long enough, Recurrent Neural networks (RNN) suffer from the short term memory, where there are chances of information loss in the process. Vanishing gradient is one of the major problems for loss of information present in the initial layers. While updating the network weights in the back propagation step, gradients shrink as it back propagates in time. Because of this, initial layers won't contribute much to the learning process. In order to overcome this, an internal mechanism called gates/cell states are introduced in the networks that regulates the flow of information. This gate helps us to understand the sequence in the information and determine which data to keep or throw away.



The repeating module in an LSTM contains four interacting layers.

Figure 3. Long Short Term Memory Network

In the LSTM network, each neural layer consists of four stages that helps to remember and pick the information from the input sequences.

First stage decides what information to keep/throw away. This is done by a sigmoid layer which outputs 0 to forget the current input cell and 1 to keep the value. This stage is called the forget gate layer (f_t).

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

Next stage consists of two parts, where the sigmoid layer (input layer) decides which new information we have to store in the network, while the tanh layer creates a new vector to store this block.

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

These informations are combined in the next stage and summed up to the output of the ft (forget layer) to force the network to forget the data that we decided earlier in the first stage and add this new information to the memory

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

In the output layer stage, the sigmoid layer decides which parts of the cell is going to output and this is combined with the tanh layer output of the previous stage results to produce the filtered part that we decided to show.

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

Gated Recurrent unit (GRU)

GRU is similar to the LSTM with few changes in the network structure. GRU is simpler than the LSTM model. Unlike four gates in LSTM, this unit has two gates only. The forget gate and input gates are combined into a single layer called update gate. Reset gate is used to decide how much past information to forget. This architecture contains only the hidden state to transfer information not the cell state. Due to this simpler version, they are speedier than LSTM and also have less tensor operation than standard LSTM.

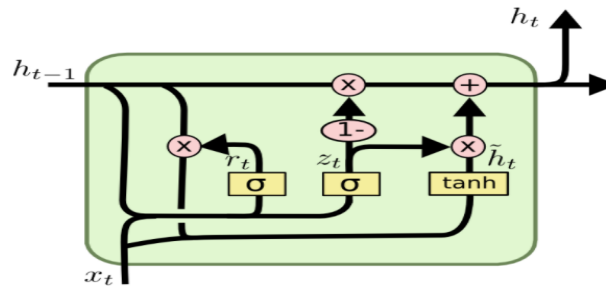


Figure 4. Gated Recurrent Network

$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

Sequence2Sequence

This deep learning method is unique in that it uses two separate RNN models. One is called an encoder which learns states based on an input sequence. The other is called a decoder and uses the states from the encoder to predict an output sequence one item at a time. This strategy is often used for machine translation.

There are two different “modes” for sequence2sequence models; training and inference. Using French to English machine translation as an example, during training, an English phrase is given to the encoder character by character, its states are then passed to the decoder with a special character that indicates the beginning of a French phrase. From here, the decoder makes a prediction and error is calculated by comparing the predicted character against the ground truth next character. Normally, the decoder uses its predictions recursively to generate character after character.

To prevent the decoder from making erroneous predictions in the sequence that would get worse as the French sentence goes on, we use something called teacher forcing during training. In this translation example, teacher forcing would be implemented by feeding the decoder the ground truth next French character no matter what the decoder predicts. This allows the decoder to properly learn the French sequences of characters. To exemplify this, we will translate the English word *go* to its French translation *va*. *Go* is fed to the encoder character by character. The encoder states and a special character ('\t' for example) is passed to the decoder to make a prediction. If the decoder predicts the letter 'r' when the ground truth character is 'v', on the next iteration the decoder is fed 'v' instead of the 'r'.

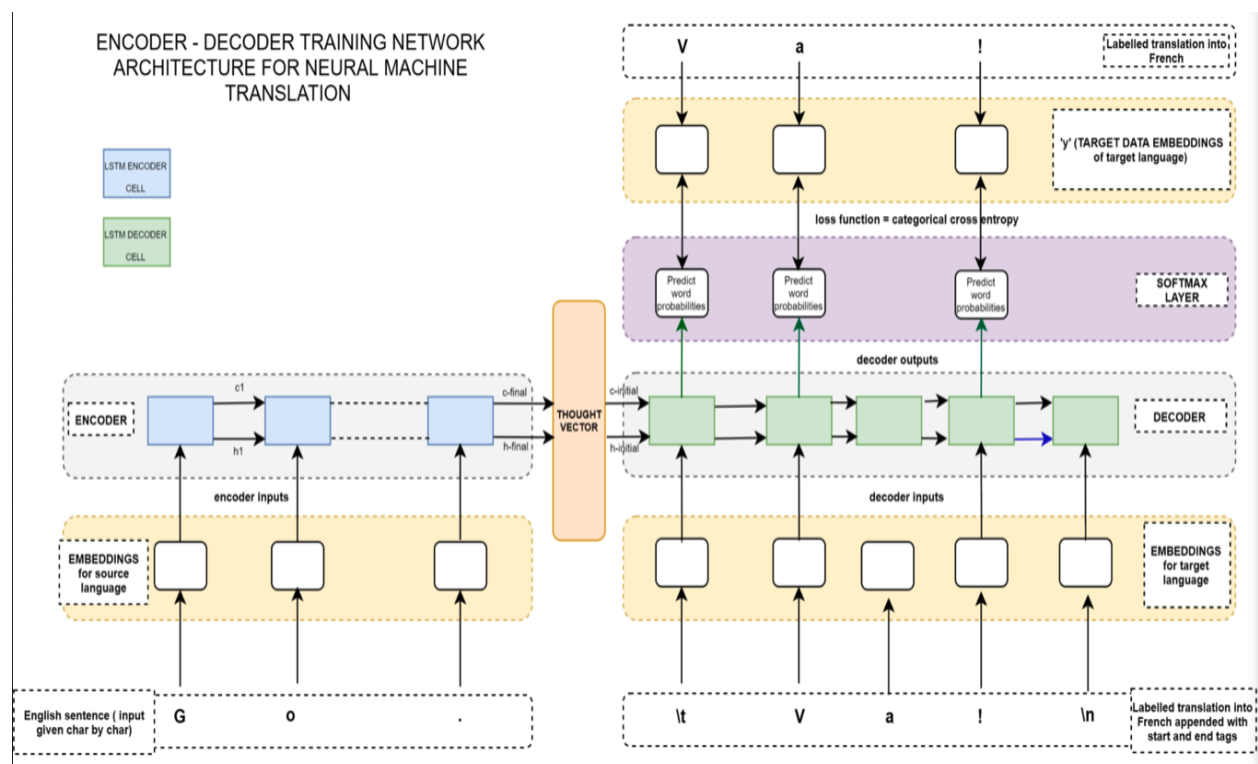


Figure 5. Sequence2sequence model in training mode using teacher forcing.

At inference time, this teacher forcing is not used and the model is free to recursively predict the next element of a sequence one at a time using the previously predicted character and the updated states from the decoder. In the next section we will show how sequence2sequence was applied to our problem domain.

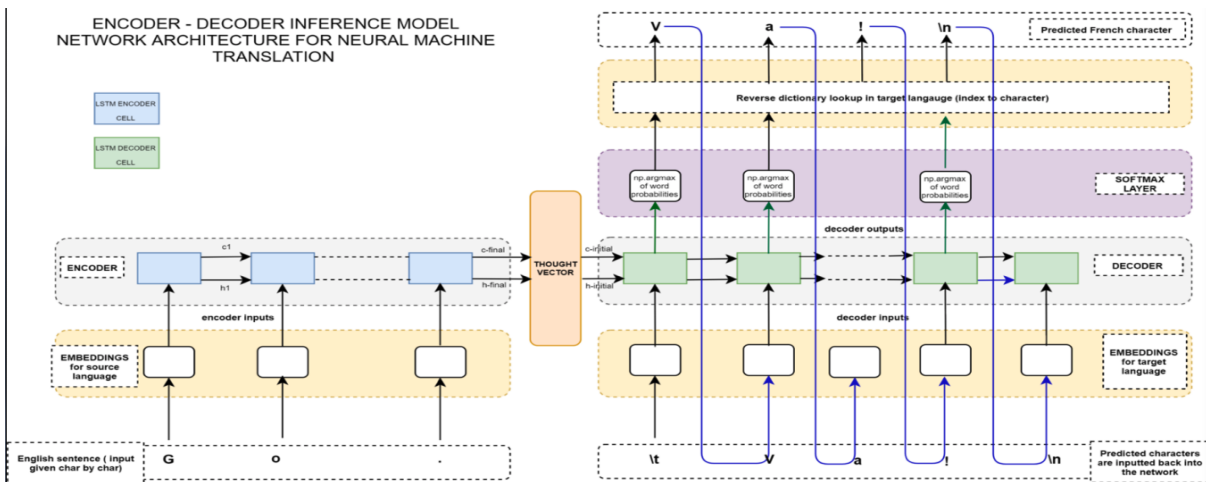


Figure 6. Sequence2sequence model in inference mode not using teacher forcing.

Experimental Setup

We built, trained, and evaluated a wide-range of models including LSTM with one or two layers, GRU-based models, bi-directional LSTMs, 1D-convolutional models, and Sequence-to-Sequence models. Out of which LSTM & sequence2sequence produces best results for our dataset. Let's discuss in brief about the experimental setup of these models.

Data Preprocessing

Our RNN models work best with sequences of equal length, so a key part of our data preprocessing was splitting the airports visited into equal-length arrays. To accomplish this, we read in the DataFrame from the EDA_and_cleaning.py script, where each row includes the string of all airports visited by each aircraft. This str is then split into a Numpy array, and then broadcast through a function that takes the Numpy array, a window length, and a stride to generate an array of arrays each the length of the window. For example, if there were originally 7 elements [0,1,2,3,4,5,6] and the window is 3 and stride is 1, we would get

```
[[0,1,2],
 [1,2,3],
 [2,3,4],
 [3,4,5],
 [4,5,6]]
```

as the resulting array. We performed this processing for every unique aircraft and concatenated all the resulting sequences together. This results in an array with a shape of n (number of series) by the window length.

Next, we used the Keras tokenizer function to map the unique string elements representing airports into integers. After tokenization, we generated y (the target variable) by slicing off the last column of the resulting series of sequences. The target variable is also one-hot encoded to a binary representation. This data array will be our target for one-step ahead prediction. The

remaining data is now in the shape of n (number of series) by the window length - 1 and will be our feature set. Finally, we split the data into a train and test set and save the resulting arrays for use in modeling.

Framework

Keras was used to implement all our RNN models. Keras is a user-friendly and simple to understand deep learning framework that the members of our team have experience using. Additionally, it allows for easy implementation of layers with GRU and LSTM units.

Model Implementation

All models were built using Keras. We compiled our models with categorical cross entropy as the loss function and accuracy as the metric. We then fit the model, determined the accuracy, and then saved the model for further evaluation in the model_evaluate.py script.

Single layer LSTM Network:

Input sequence represents airport location information, the sequence is encoded into numeric values. The First layer of the network is the embedding layer where input sequences are embedded into the patterns, thereby preparing the data to feed into the lstm network.

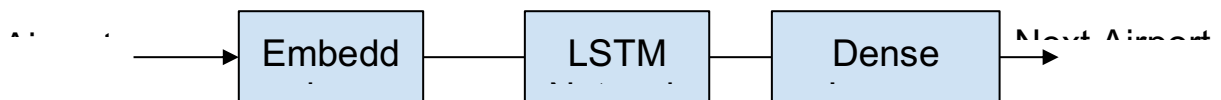


Figure 6. Experiment Set up - Single layer LSTM network

Parameters Used:

Neurons - 400 Epoch - 10 Batch Size - 256 Embedding Size - 200 Learning rate - 1e-3
 Airport Sequence length - 25,50 Optimizer - Adam

Stacked LSTM Network:

Multiple hidden LSTM layers are stacked on top of each other to form a stacked LSTM network. Three layers of LSTM are stacked up together and sent to a dense layer output with softmax activation to get the output.

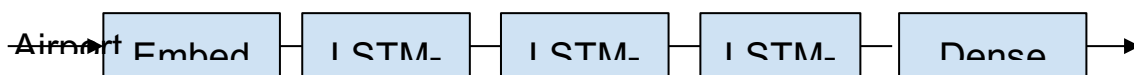


Figure 7. Experiment Set up - Stacked LSTM network

Parameters Used:

Neurons - 400-300-200 | Epoch - 10 | Batch Size - 256 | Embedding Size - 200 | Learning rate - 1e-3
 Airport Sequence length - 25,50 Optimizer - Adam

Ensemble Method:

In this method, I have stacked up above network results as input to the logistic regression model and value with maximum probability is chosen to be the next in sequence value.

Base learners : single layer lstm network output & stacked layer lstm network output

Meta learners : Logistic Regression

Loss: Cross entropy

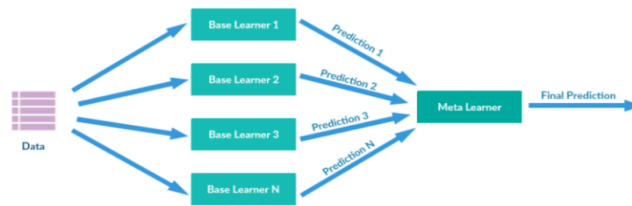


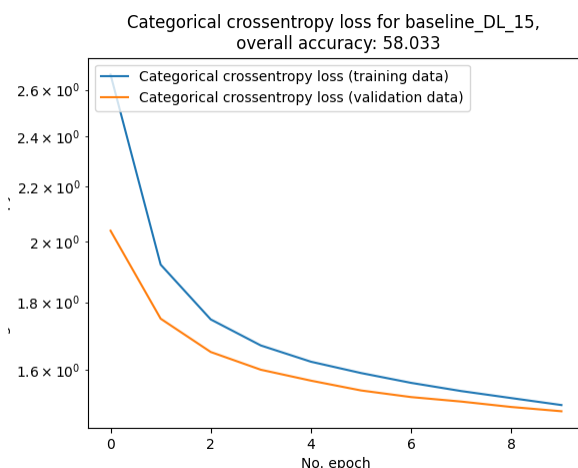
Figure 8. Experiment Set up - Stacked Ensemble Method

Sequence2sequence

For our sequence2sequence implementation, we also used categorical cross entropy as the loss function and accuracy as the metric. The encoder and decoder both use a single LSTM layer with 256 neurons. The model used sequences of flights in lengths 5, 10, 15, 25, and 50. One consideration that had to be made was how to split the sequences into beginning and ending sequences. For example, if using sequences of length 25, you could feed the encoder 13 airports as the beginning sequence and 12 to the decoder as the ending sequence. Ultimately, we got the best performance from sequence2sequence by feeding the encoder all but one of the airports in the sequence and having the decoder only predict the final destination in the sequence. It is important to mention that the ending sequences were padded with starting number '1111' and ending number '9999'. This allowed us to feed the decoder a special beginning character similar to the machine translation example described earlier in this report. All sequence2sequence models were trained for 10 epochs with a batch size of 64. The Adam optimizer was used with a 0.001 learning rate. The sequence2sequence implementation can be found in modeling_seq2seq.py and it includes training, testing, and evaluation all in one script.

Model Evaluation Metrics

To evaluate the performance of our model, we employed several different metrics. All models were saved as hdf5 files with filenames corresponding to the model name and run type (airline and fixed sequence length, eg. "DL_15" uses data only from Delta airlines with a sequence length of 15). Additionally, the training history of each model was saved as a Pandas DataFrame for future retrieval and visualization.



All model evaluations occur in the model_evaluation.py script (except for sequence2sequence). In this script, required x and y data, models, corresponding history files and tokenizers are loaded from disk. Each model's training

history is then plotted to compare the training and validation loss to ascertain if the model overfit or underfit. The example on the right shows the training history for our baseline model on the Delta data with a sequence length of 15 (“baseline_DL_15”). In this plot, we can see that both losses drop and converge near the tenth epoch. Additional training may improve performance on the training data but would likely overfit the data and not improve the validation data. This supports our use of 10 epochs in training the baseline model.

Next, we determined the overall accuracy of each model. Simply put, this metric determines for each predicted next airport, how many times did the model get the correct answer divided by the total number of predictions. We also calculated the Cohen’s Kappa score, which is scaled from -1 (indicating no agreement) to 1 (complete agreement). Cohen’s Kappa is usually used to compare the results from two different manual annotations but can also be used to evaluate models in cases where we are concerned about the likelihood of randomly guessing affecting our accuracy. Because of its additional inclusion of uncertainty, Cohen’s Kappa score is a complimentary metric to overall accuracy.

In addition to these metrics, we generated a confusion matrix and a classification report to evaluate where our models did well and where they did poorly. As mentioned earlier in our EDA section, there is a risk that our models may only learn patterns associated with larger, busier airports. Since the airports have more flight activity, correctly predicting activity at these airports would inflate the accuracy of the model across all the data. To analyze if this was occurring, our team developed a scatter plot to show the precision and recall of a model for each airport. By sizing the scatterplot mark for each airport to the number of flights at that airport, we can create a visualization to help us understand how well the model performed across airports of different sizes.

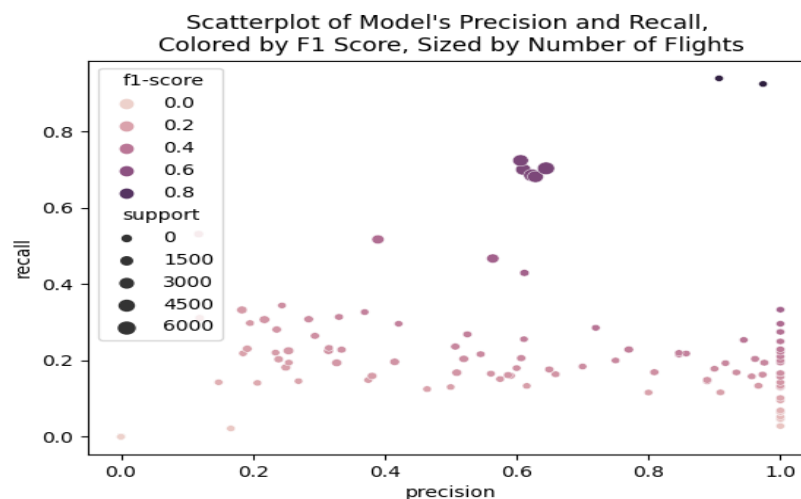


Figure 9. In this scatterplot from the baseline model using United Airlines data, we can see a tight grouping of the largest airports at around 0.6 precision and 0.7 recall. Smaller airports struggle to break 0.4 recall but many earn high precision scores.

Finally, our scripts recorded the overall accuracy and Cohen’s Kappa score for each model against each subset of data in a log. This log allowed for our team to review how well a model performed and determine the best performing model.

Results

In order to evaluate the performance of the model, we generated dummy random sequences and fed this to our experimental set up and observed the results. This forms the base of the experiment to evaluate the other model performance.

Random Sequences Output:

```
Final accuracy on validations set: 24.448
Cohen kappa score 0.0
Precision: 0.059769239279687005
Recall: 0.24447748215262485
F1: 0.09605515589772128
```

From the above result, we could see that the accuracy of the random sequences is close to 24%. The Cohen kappa score suggests that output is completely random. If any of our experimental networks produces accuracy less than this base accuracy, then the model is underfitting or performing bad.

Single Layer Lstm Model:

Airline: Delta

Sequence Length: 25

```
Modal Evalution Results

Accuracy: 62.422
Cohen kappa score 0.5919228700612527
Precision: 0.6325509588477811
Recall: 0.6242243679737651
F1: 0.622350160687898
```

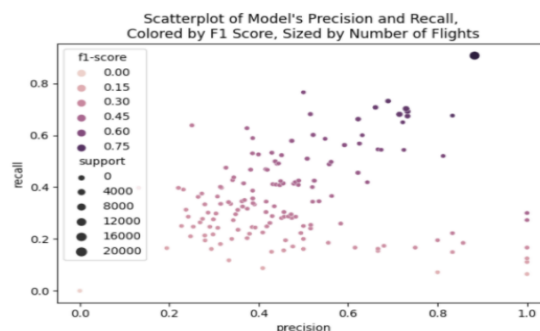


Figure 10. Scatterplot Precision vs Recall - LSTM network on Test set

Airline: Delta

Sequence Length: 50

Modal Evalution Results

Accuracy: 62.646

Cohen kappa score 0.5933283419217596

Precision: 0.6360756942665542

Recall: 0.626464252053319

F1: 0.623502422167004

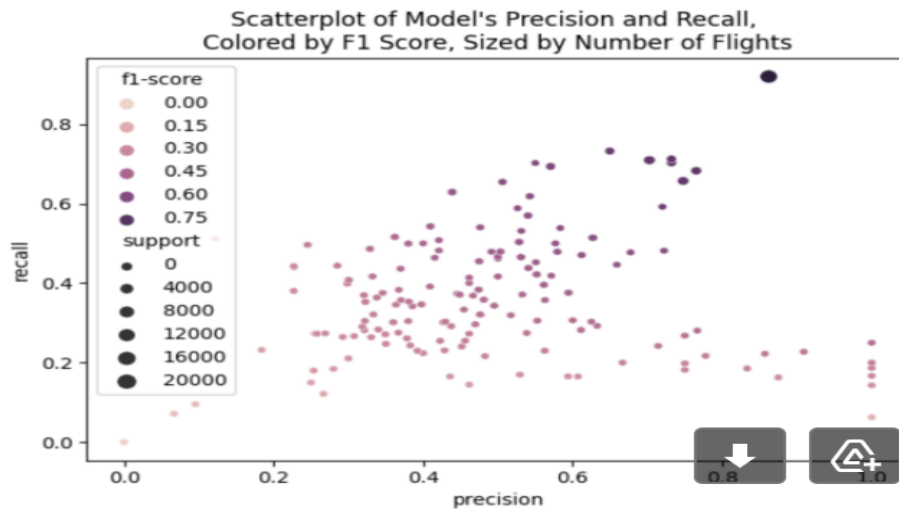


Figure 11.Scatterplot Precision vs Recall - LSTM network on Test set

Stacked Lstm Network Results:

Airline: Delta

Sequence Length: 25

Modal Evalution Results

Accuracy: 61.911

Cohen kappa score 0.5850593497951324

Precision: 0.6364303404942205

Recall: 0.619112430474435

F1: 0.6149982165744742

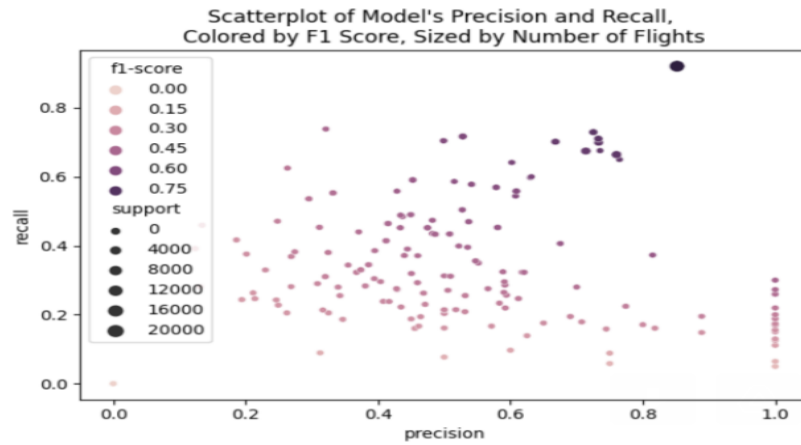


Figure 12. Scatterplot Precision vs Recall - Stacked LSTM network on Test set

Airline: Delta
Sequence Length: 50

```

Modal Evaluation Results

Accuracy: 62.227
Cohen kappa score 0.5895667642018765
Precision: 0.636641565032489
Recall: 0.622267851532696
F1: 0.6209776412113642
  
```

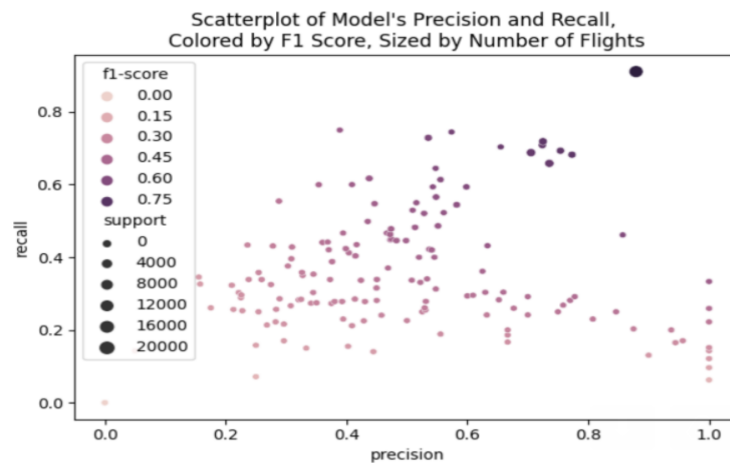


Figure 13. Scatterplot Precision vs Recall - Stacked LSTM network on Test set

Stacked Ensemble Methods:

Airline: Delta
Sequence Length: 25

```
Final Best Model Results
Final accuracy on validations set: 63.429820707097775 %
Cohen Kappa 0.6016462538775594
F1 score 0.37488921840421413
```

Airline: Delta

Sequence Length: 50

```
Final Best Model Results
Final accuracy on validations set: 63.58444414523585 %
Cohen Kappa 0.6033147613362919
F1 score 0.3825873847459288
```

Sequence2Sequence Method:

Airline: Delta

Sequence Length: 25

```
Accuracy: 0.6262
Cohens Kappa: 0.5920248164166245
Precision: 0.6318722285441556
Recall: 0.6262
F1: 0.6221704140365036
```

Scatterplot of seq2seq_DL_25 Precision and Recall,
Colored by F1 Score, Sized by Number of Flights

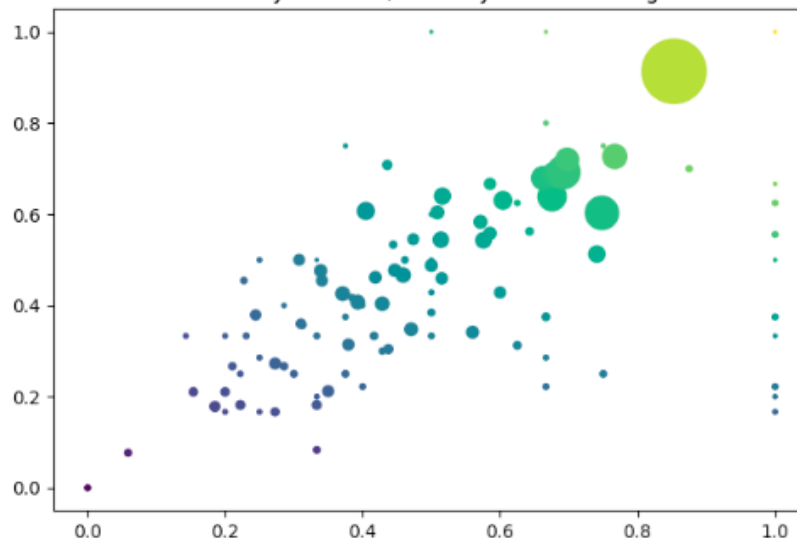


Figure 14. Scatterplot Precision vs Recall - Seq2seq network on Test set

Airline: Delta

Sequence Length: 50

Accuracy: 0.613
Cohens Kappa: 0.5793893933873405
Precision: 0.6229796477216917
Recall: 0.613
F1: 0.611046099536135

Scatterplot of seq2seq_DL_50 Precision and Recall,
Colored by F1 Score, Sized by Number of Flights

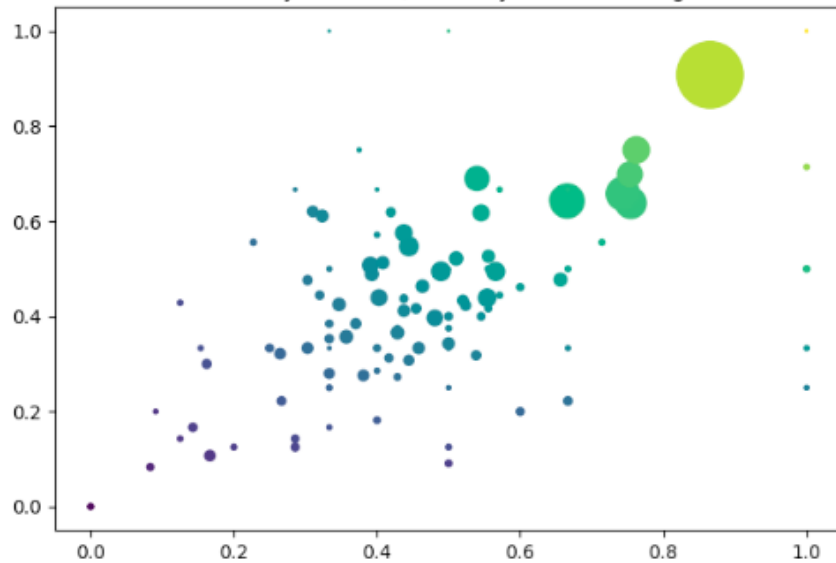


Figure 15. Scatterplot Precision vs Recall - Seq2seq network on Test set

Airline: Delta
Sequence Length: 15

Accuracy: 0.597
Cohens Kappa: 0.5641579221749138
Precision: 0.5996325790310026
Recall: 0.597
F1: 0.5930424510046131

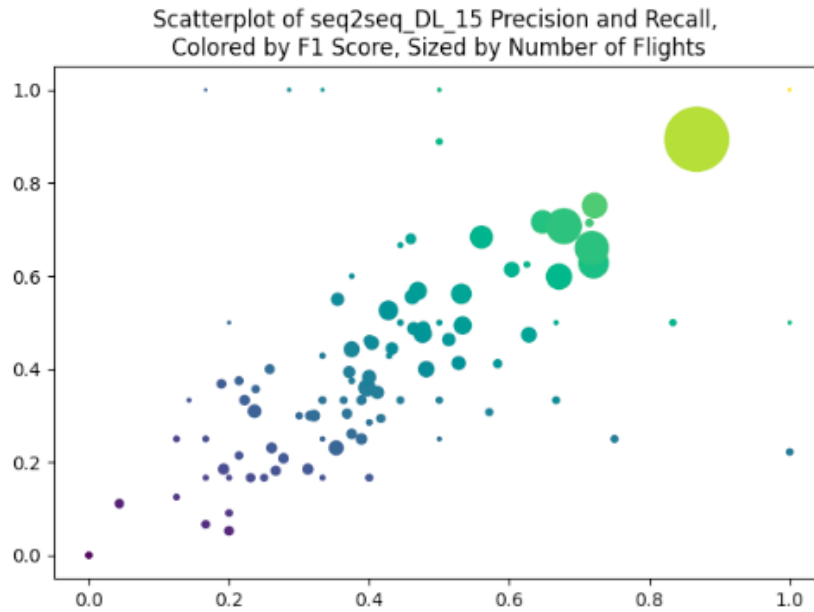


Figure 15. Scatterplot Precision vs Recall - Seq2seq network on Test set

Although the ensemble method improves the accuracy of the model and kappa score, we could see F1 score is reduced. Hence, from other network model performance, LSTM network and sequence2sequence method performs moderately better than other networks.

Summary

Our team used Delta Airlines sequences in lengths 5, 10, 15, 25, and 50 for all of our models in order to perform a comparative analysis with the evaluation criteria described above. Below is a table showing the performance of our top five models on a test set.

Table 2. Results from top 5 models compared with Baseline performance.

Model	Accuracy	Cohen's Kappa	F1	Sequence length
LSTM network	0.6264	0.5933	0.6235	50
seq2seq	0.6262	0.5920	0.6222	25
LSTM network	0.6242	0.5919	0.6223	25
Stacked LSTM	0.6222	0.5895	0.6209	50
seq2seq	0.6130	0.5794	0.6110	50
Baseline (pure random)	0.2444	0.0	0.09	50

Summary and Conclusions

Evaluating Stationarity and Autocorrelation

Our group took multiple steps to verify that our data was a relatively stationary dataset with an embedded signal we could learn through our models. In addition to the described work with ACF plots and stationarity tests with the Augmented Dickey-Fuller test, we also applied a truly randomly generated dataset against a trained model and received far lower scores than on the actual data. This suggests there is an underlying pattern.

Repeated Sequences

During our analysis, our team noticed repeated elements in a sequence, eg an aircraft departing one airport to visit that same airport. In our research, this was usually accompanied by a time gap of several days. This activity could be related to flight maintenance, erroneous data, or non-commercial flights missing from our data. After discussion with our team, we decided it would be best to keep the data as is with these anomalies rather than correct them and introduce additional, unknown biases.

Sequence Length

In Table 2, we see that sequences of length 25 or 50 yielded the best results from our models. Our group theorized that sequences of more than 50 airports might perform better but would be challenging to use in a production setting. The longer the sequence, the less robust our model would be in picking up on drastic changes to the airport patterns.

Looking at Other Airlines

To limit the scope of the project, we only focused on Delta Airlines flights. As seen in Figure 1, Delta has a central hub in Atlanta which accounts for over 25% of all of its activity. This led to our models performing very well on predicting when a flight would go to Atlanta or other busy airports, but the models struggled to perform well on smaller airports. We tried to model United and American Airlines data and found that these less centralized datasets did worse than the Delta data with our models.

Best Model

Our group expected sequence2sequence to perform the best on this task but it was outperformed by a simple single layer LSTM. As explained previously, different sequence lengths were tested with sequence2sequence but all hyperparameters stayed the same. Tuning the hyperparameters after finding the best sequence length would probably improve seq2seq's performance.

References

https://scikit-learn.org/stable/modules/generated/sklearn.metrics.accuracy_score.html?highlight=accuracy#sklearn.metrics.accuracy_score

https://scikit-learn.org/stable/modules/generated/sklearn.metrics.cohen_kappa_score.html

<https://towardsdatascience.com/neural-machine-translation-using-seq2seq-with-keras-c23540453c74>

[Understanding LSTM Networks -- colah's blog](#)

[Illustrated Guide to LSTM's and GRU's: A step by step explanation | by Michael Phi | Towards Data Science](#)

Appendix A: Code for Project

Link to all code on GitHub: <https://github.com/redhairedcelt/Final-Project-Group1>

Appendix B: Performance for All Standard RNN Models

Model	Airline	Sequence length	Embedding size	Number of neurons	optimizer	Best Accuracy - Validation Set
Single Layer LSTM without Dropout	DL	5	300	100	Adam	59.15(50 epoch)
Single Layer LSTM without Dropout	DL	5	300	200	Adam	59.40(20 epoch)
Single Layer LSTM without Dropout	DL	5	300	400	Adam	59.58(20 epoch)
Single Layer LSTM without Dropout	DL	5	150	400	Adam	59.70(20 epoch)
Single Layer LSTM without Dropout	DL	5	125	400	Adam	59.57(20 epoch)
Stacked Layer LSTM(two) without Dropout	DL	5	125	400 - 300	Adam	59.60(20 epoch)
Stacked Layer LSTM(Three) without Dropout	DL	5	125	400 - 300 - 200	Adam	59.50(20 epoch)
Stacked Layer LSTM(Three) without Dropout	DL	5	125	400 - 400 - 400	Adam	59.45(10 epoch)
Stacked Layer LSTM(four) without Dropout	DL	5	125	400 - 300 - 200 - 100	Adam	59.37(20 epoch)
Bidirectional LSTM (Single Layer) without Dropout	DL	5	125	400	Adam	59.54(20 epoch)
Bidirectional LSTM (two Layer) without Dropout	DL	5	125	400 - 400	Adam	59.28(20 Epoch)
Bidirectional LSTM (two Layer) without Dropout	DL	5	125	400 - 300	Adam	59.40(20 Epoch)
Bidirectional LSTM (three Layer) without Dropout	DL	5	125	400 - 300 - 200	Adam	59.40(20 Epoch)
Conv1d-Max-Conv1d-Global Max	DL	5	300	400(kernel - 1 & Pooling - 2)	Adam	38.70(20 Epoch)
Conv1d-Max-Conv1d-Global Max	DL	5	300	200(kernel - 2 & Pooling - 2)	Adam	56.12(20 epoch)
Conv1d-Max-Conv1d-Global Max	DL	5	200	200(kernel - 2 & Pooling - 2)	Adam	56.62(20 epoch)
GRU - Dense without Dropout	DL	5	150		Adam	35.20(20 epoch)
GRU - Dense with Dropout -0.2	DL	5	200		Adam	34.04(10 epoch)
Stacked GRU (2 Layer) - Dense with dropout	DL	5	200		Adam	34.5(20 Epoch)
conv1d-GRU	DL	5	200	200(kernel - 2 & Pooling - 2)		49.30(20 epoch)

Single Layer LSTM without Dropout	DL	10	300	200	Adam	60.86(20 epoch)
Single Layer LSTM without Dropout	DL	10	300	300	Adam	61.43(10 epoch)
Single Layer LSTM without Dropout	DL	10	300	400	Adam	61.33(10 epoch)
Single Layer LSTM without Dropout	DL	10	300	500	Adam	61.49(7 epoch)
Stacked Layer LSTM(two) without Dropout	DL	10	200	500-400	Adam	61.45(20 epoch)
Stacked Layer LSTM(four) without Dropout	DL	10	200	500-400-300-200	Adam	61.52(10 epoch)
Bidirectional LSTM (Single Layer) without Dropout	DL	10	125	400	Adam	60.54(10 epoch)
Bidirectional LSTM (two Layer) without Dropout	DL	10	125	400 - 300	Adam	60.87(10 epoch)
Bidirectional LSTM (two Layer) without Dropout	DL	10	125	200 - 100	Adam	60.32(10 epoch)
Conv1d-Max-Conv1d-Global Max	DL	10	300	200(kernel - 1 & Pooling - 2)	Adam	53.43(10 epoch)
GRU - Dense without Dropout	DL	10	200		Adam	43.05(15 epoch)
GRU - Dense with Dropout -0.2	DL	10	200		Adam	40.12(15 epoch)
Stacked GRU (2 Layer) - Dense with dropout	DL	10	200		Adam	44.25(20 Epoch)
conv1d-GRU	DL	10	200	200(kernel - 2 & Pooling - 2)	Adam	56.82(50 Epoch)
Single Layer LSTM without Dropout	DL	25	200	200	Adam	62.11(20 epoch)
Single Layer LSTM without Dropout	DL	25	200	300	Adam	62.39(20 epoch)
Single Layer LSTM without Dropout	DL	25	200	400	Adam	62.59(10 epoch)
Single Layer LSTM without Dropout	DL	25	200	500	Adam	62.55(10 epoch)
Stacked Layer LSTM(two) without Dropout	DL	25	200	200 - 100	Adam	61.61(22 Epoch)
Stacked Layer LSTM(two) without Dropout	DL	25	200	400-300	Adam	61.50(10 Epoch)
Stacked Layer LSTM(three) without Dropout	DL	25	200	400-300-200	Adam	62.28(10 Epoch)

Bidirectional LSTM (Single Layer) without Dropout	DL	25	200	400	Adam	61.18(8 epoch)
Bidirectional LSTM (two Layer) without Dropout	DL	25	200	400 - 300	Adam	62.20(8 epoch)
Bidirectional LSTM (two Layer) without Dropout	DL	25	200	200 - 100	Adam	61.65(10 epoch)
Conv1d-Max-Conv1d-Global Max	DL	25	200	200(kernel - 8 & Pooling - 2)	Adam	45.01(10 epoch)
Conv1d-Max-Conv1d-Global Max	DL	25	200	200(kernel - 4 & Pooling - 2)	Adam	39.53(12 Epoch)
GRU - Dense without Dropout	DL	25	200		Adam	53.15(50 epoch)
GRU - Dense with Dropout -0.2	DL	25	200		Adam	53.15(30 epoch)
Stacked GRU (2 Layer) - Dense with dropout	DL	25	200		Adam	53.17(20 Epoch)
conv1d-GRU	DL	25	200	200(kernel - 2 & Pooling - 2)	Adam	44.98(50 Epoch)
Single Layer LSTM without Dropout	DL	50	200	200	Adam	62.09(20 epoch)
Single Layer LSTM without Dropout	DL	50	200	300	Adam	62.96(14 epoch)
Single Layer LSTM without Dropout	DL	50	200	400	Adam	62.88(14 epoch)
Single Layer LSTM without Dropout	DL	50	200	500	Adam	63.17(10 epoch)
Stacked Layer LSTM(two) without Dropout	DL	50	200	200 - 100	Adam	61.32(22 Epoch)
Stacked Layer LSTM(two) without Dropout	DL	50	200	400-300	Adam	62.67(10 Epoch)
Stacked Layer LSTM(three) without Dropout	DL	50	200	400-300-200	Adam	62.36(10 Epoch)
Bidirectional LSTM (Single Layer) without Dropout	DL	50	200	400	Adam	61.44(10 epoch)
Conv1d-Max-Conv1d-Global Max	DL	50	200	400(kernel - 8 & Pooling - 5)	Adam	38.39(10 epoch)
GRU - Dense without Dropout	DL	50	200		Adam	57.55(50 epoch)
GRU - Dense with Dropout -0.2	DL	50	200		Adam	55.57(30 epoch)
Stacked GRU (2 Layer) - Dense with dropout	DL	50	200		Adam	57.15(20 Epoch)

conv1d-GRU	DL	50	200	200(kernel - 2 & Pooling - 2)	Adam	38.77(20 Epoch)
------------	----	----	-----	-----------------------------------	------	-----------------

Appendix C: Performance for All Sequence2sequence Models

Model	Accuracy	Cohen's Kappa	F1	Sequence length
seq2seq	0.6262	0.5920	0.6222	25
seq2seq	0.6130	0.5794	0.6110	50
seq2seq	0.5970	0.5642	0.5930	15
seq2seq	0.5930	0.5607	0.5930	10
seq2seq	0.4968	0.4526	0.4900	5