

Flight Delay Prediction

Vishrut Raj¹, Viran Raj², Satyam Singh³, Adityanath Mishra⁴
^{1,2,3,4}Rajkumar Goel Institute of Technology/AKTU

Abstract - light delay prediction is fundamental to establish the more efficient airline business. The development of accurate prediction models for flight delays became cumbersome due to the complexity of air transportation system, the number of methods for prediction, and the deluge of flight data. In this context, this paper presents a thorough literature review of approaches used to build flight delay prediction model. Airlines delays make immense loss for business field as well as in budget loss for a country. Flight delays hurt airlines, airports, and passengers. We are proposing machine learning algorithms like Linear regression Techniques. The aim of this research work is to predict Flight Delay, Which is highest economy producing field for many countries and among many transportation this one is fastest and comfort, so to identify and reduce flight delays, can dramatically reduce the flight delays to saves huge amount of turnovers, using machine-learning algorithms. Flight delays could always be annoying, especially in the case when the period of delay was so long that there was even a danger to miss the next flight. However, if there was a way to predict whether there would be a delay or even better – how long the delay could be, then people could make earlier preparation to reschedule following flights in an earlier manner.

Index Terms - Data prediction, Machine Learning, Linear Regression Techniques.

I.INTRODUCTION

Flight delays could always be annoying, especially in the case when the period of delay was so long that there was even a danger to miss the next flight. However, if there was a way to predict whether there would be a delay or even better – how long the delay could be, then people could make earlier preparation to reschedule following flights in an earlier manner. For that consideration, we adopted a dataset containing airline delayed time and other air liner information provided by Kaggle to building a model, mainly aiming to solve the following questions. Whether there would be a delay with certain publicly reachable resources; and 2. How long delayed time

one could expect with the same information given. We deployed python s k learn and pandas library to build our model, and evaluate our model based on R-Square for linear regression and accuracy rate for logistic regression. As a brief result of our project, we found, it would be helpful to use the following factors in evaluating our model: week, month, airline carrier reference, planned elapsed time (in air time), distance between two departure and destinations, flight planned departure time, departure airport code, and taxi-in and taxi out time.(1)

The continuous increase of storage capacities and computational power is currently pulling the development of data analytics. Indeed, companies (and especially IT-intensive ones) are collecting massive volume of data (often referred as Big Data), such as web logs, customer information, production and sales tracking, etc. Analyzing these datasets, with data mining algorithms for example, allows the extraction of information that can help a company to gain knowledge (for example on customs' behaviors) or to use the information as a basis for new products or services.

Dataset Used

We chose the “Airlines Delay” data from www.kaggle.com/datasets, which was actually provided by the U.S. Department of Transportation's (DOT) Bureau of Transportation Statistics (BTS), that tracks the on time performance of domestic flights operated by large air carriers. Summary information on the number of on-time, delayed, cancelled and diverted flights appears in the data. We initially intend to use the whole dataset covering 10 years to build our model, however, due to our limited computation resources, we had to cut the dataset of the latest dataset of the year or 2018. In the dataset, 7213446 rows included, and 27 variables involved.

II.LITERATURE SURVEY

A.LOGISTIC REGRESSION MODEL

No data mining projects could be finished without thoroughly understand the data first. So, in order to better understand data, we start our project by exploring the data first. We found the original dataset includes 28 attributes/columns and while most of the data were in float format, some of them were object types). In addition, as shown, there were also many null values in the original datasets. So, we need to first clean the columns with null values and change data types of objects into suitable types (mostly integers) for the convenience of machine learning.

B.LINEAR REGRESSION MODEL

We renamed the original data column names and validated the nulls, however with a little different approach. We first plotted a density plot for chosen attributes. After plotting the density plot with columns with “nan” values, we found none of the columns strictly follows normal distribution, and most of them were largely skewed and concentrated to only few values. Replacing methods, we tried included applying fillna () method to replace “nan” and replacing missing “nan” values with the mean of corresponding columns. However, none of the methods enable us to develop model with desirable results. So instead of replacing “nan” with normal distribution, we decided to use merely replace “nan” with extreme values that without the original data range.

C.INITIAL DATA EXPLORING

After data cleaning we start the first process of exploring our data if there were any patterns within the independent variables.

The above graph shows the no of delays airline wise. On the left side you can see there is a false value, which means instances when an airline has not been delayed. On the right-side true values suggests that there is a delay. We can see and conclude that maximum delay is caused by Southwest Airlines. Also, in the next graph we can see that the maximum number of flights are from Southwest Airlines, which compel us to think that one of the reasons for the delay is the operational process of airline. And these delays are known as career delay, we can reduce this delay with effective planning strategies.

D.DIMENSIONALITY REDUCTION

So as there were 28 columns and we wanted our model to be very precise, before going ahead we wanted to be sure there should not be any kind of correlation between the predictor variables, otherwise our model will be overfitting. So, we used correlation matrix and the criteria was, if 2 variables have correlation greater than 0.4 or less than -0.4, we will drop one of those variables.

The most common correlated variables are actual departure time, actual arrival time, planned arrival time planned departure time among others, this makes sense because these factors are directly impacting the delay so there is no point of adding those variables. For example, the website FlightCaster exploit several sources of information (airports, airlines, weather and possibly historical data) to provide probabilities of being on-time, less than one hour late or more than one hour late, to travelers. However, this website is using the same estimations for all the flights when no short-term information is available.

III.METHODOLOGY

Due to the huge data (7213446 lines of data included), it would both be impossible and impractical for us to manually explore and find patterns between flight delay information and related influencers. Therefore, we decide to first manually clean the data, and then adopt machine learning with Python sklearn library. For data cleaning, we firstly change data-types of certain columns with “object” type, and replace ‘null’ values with certain values, to make the data suitable for machine learning. Afterwards, used pandas, seaborn and matplotlib to make initial exploration in order to find some intuitive relationship between variables. Finally, we deploy machine learning method to dig out factors and their correlation with flight delays – to be specific, we used linear regression to predict the expected delay time of flights and used logistic regression to predict whether a flight might or might not be delayed.(2)

A.USING LINEAR REGRESSION

Since logistic regression is appropriate for categorical values, and we expect to predict the potential delayed time, which is a numerical valuable, it makes more sense to apply Linear Regression for our model. Therefore, we applied sklearn linear model, and used r2_score to evaluate our model. We set

Delay_Departure_Time, which is a set of both positive and negative figures to imply how long exactly a plane departure delayed/early². We then also include week, month, airline carrier reference, planned elapsed time (in airtime), distance between two departure and destinations, flight planned departure time, departure airport code, and taxi-in and taxi-out³ time.

We split the test-train sets into 2:8 ratio, and got a R Square score of 0.806, which was acceptable. (4)

B.USING LOGISTIC REGRESSION

Actual Arrival time - Expected Arrival Time + (Actual Departure time - Expected Departure Time). But this independent variable could be any random number so we created another column in which if total delay > 0 then value will be 1 else it will be 0. This made us think that we should run logistic regression on this model and predict the factors responsible for delay.

C.FACTORS AFFECTING FLIGHT DELAY

- Airline_Carrier - (AA, UA, ...)
- Month - Which Month (Jan, Feb...)
- Week - Which date of the week (Monday, Tuesday, etc.)
- Planned_Elapsed_Time - Planed in-airtime
- Tax In/Tax Out
- DISTANCE - that’s straightforward
- Planned_Departure_Time
- Airport_Departure_Code - Which airport the plane is going to start

IV.HELPFUL HINTS

As shown in (Figure 1), there were also many null values in the original datasets. So, we need to first clean the columns with null values and change data types of objects into suitable types (mostly integers) for the convenience of machine learning. As shown in (Figure 2), there were also many null values in the original datasets. So, we need to first clean the columns with null values and change data types of objects into suitable types (mostly integers) for the convenience of machine learning. We renamed the original data column names and validated the nulls, however with a little different approach. We first plotted a density plot for chosen attributes. After plotting the density plot with columns with “nan” values, we found none of the columns strictly follows normal distribution, and most of them were largely

skewed and concentrated to only few values (see figure 3). As we can see the above graph (Figure 4). Taxi in And Taxi out time are almost similar for most of the airlines but there for endeavour airlines & republic airways, the taxi out time is much higher than taxi in time. We were not able to find out the exact cause for this, so we consider it an anomaly. (Figure5) We then also include week, month, airline carrier reference, planned elapsed time (in air time), distance between two departure and destinations, flight planned departure time, departure airport code, and taxi-in and taxi-out³ time.

We split the test-train sets into 2:8 ratio, and got a R Square score of 0.806, which was acceptable. For testing we use 25 % of the dataset and we used 75 % of dataset for training our model. After running our model, we found out that our model is 82 % accurate and below is the conclusion matrix (Figure6).

```
In [20]: a = airline.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 7213446 entries, 0 to 7213445
Data columns (total 28 columns):
FL_DATE      object
OP_CARRIER  object
OP_CARRIER_FL_NUM  int64
ORIGIN       object
DEST         object
CRS_DEP_TIME int64
DEP_TIME     float64
DEP_DELAY    float64
TAXI_OUT     float64
WHEELS_OFF   float64
WHEELS_ON   float64
TAXI_IN      float64
CRS_ARR_TIME int64
ARR_TIME     float64
ARR_DELAY    float64
CANCELLED    float64
CANCELLATION_CODE  object
DIVERTED     float64
CRS_ELAPSED_TIME float64
ACTUAL_ELAPSED_TIME float64
AIR_TIME     float64
DISTANCE     float64
CARRIER_DELAY float64
WEATHER_DELAY float64
NAS_DELAY    float64
SECURITY_DELAY float64
LATE_AIRCRAFT_DELAY float64
Unnamed: 27 float64
```

Figure 1. Original Data Types

```
In [20]: for col in airline.columns:
print("There are: %d null values in column %s" % (airline[col].isnull().sum(),col))
There are: 0 null values in column FL_DATE
There are: 0 null values in column OP_CARRIER
There are: 0 null values in column OP_CARRIER_FL_NUM
There are: 0 null values in column ORIGIN
There are: 0 null values in column DEST
There are: 0 null values in column CRS_DEP_TIME
There are: 112317 null values in column DEP_TIME
There are: 117234 null values in column DEP_DELAY
There are: 115830 null values in column TAXI_OUT
There are: 115829 null values in column WHEELS_OFF
There are: 119346 null values in column WHEELS_ON
There are: 119346 null values in column TAXI_IN
There are: 0 null values in column CRS_ARR_TIME
There are: 119345 null values in column ARR_TIME
There are: 137848 null values in column ARR_DELAY
There are: 0 null values in column CANCELLED
There are: 7896862 null values in column CANCELLATION_CODE
There are: 0 null values in column DIVERTED
There are: 18 null values in column CRS_ELAPSED_TIME
There are: 134442 null values in column ACTUAL_ELAPSED_TIME
There are: 134442 null values in column AIR_TIME
There are: 0 null values in column DISTANCE
There are: 5869736 null values in column CARRIER_DELAY
There are: 5869736 null values in column WEATHER_DELAY
There are: 5869736 null values in column NAS_DELAY
There are: 5869736 null values in column SECURITY_DELAY
There are: 5869736 null values in column LATE_AIRCRAFT_DELAY
There are: 7213446 null values in column Unnamed: 27
```

Figure 2.Columns with Null Values

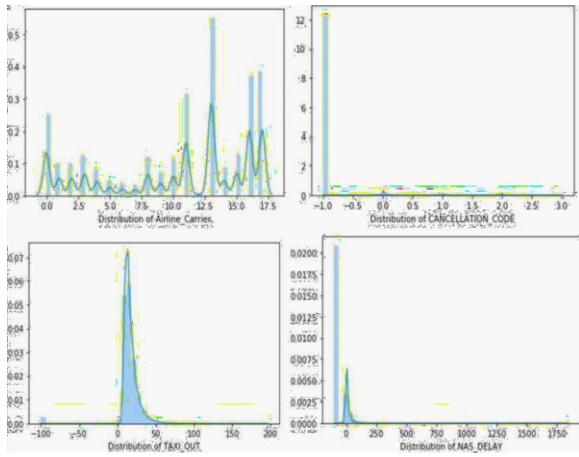


Figure 3. Plotting the Density Plot Graph

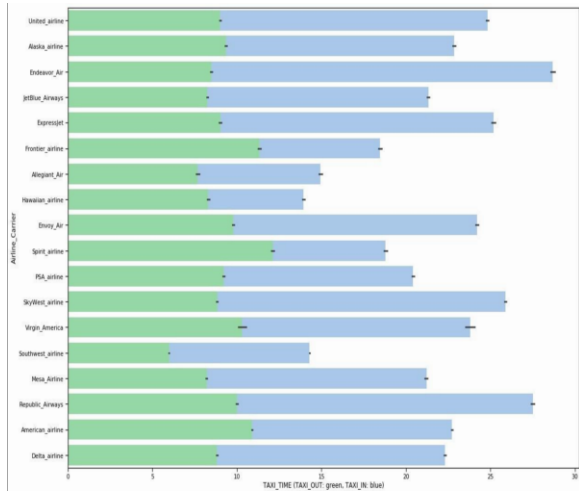


Figure 4. Graph showing Taxi In and Taxi Out

Linear model accuracy (with the test set): 0.8862686700094936

Figure 5. Linear model accuracy

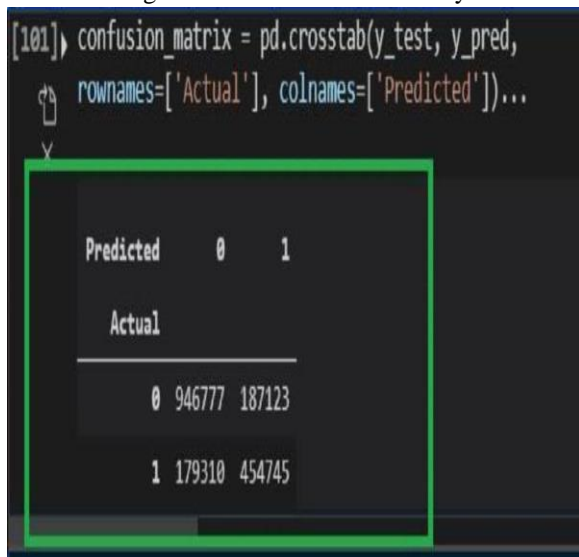


Figure 6.Conclusion Matrix

V.CONCLUSION

After applying both the models for predicting whether a flight should be delayed, as well as how much one would expect a flight should be delayed, we found the following factors to be important: week, month, airline carrier reference, planned elapsed time (in airtime), distance between two departure and destinations, flight planned departure time, departure airport code, and taxi-in and taxi-out4 time. By applying our model, on the data collected, one could be able to predict whether a flight might be delayed, and more importantly, how long delayed time she/he would expect.

However, there is some limitation in our model, first, our model only included one-year data due to our computation capability, as more years of data included, the prediction could be easier. In addition, some other related information such as airplane type, e.g., detailed weather data specific to airport were not included. Therefore, researchers could try to collect more related data and deploy better computational powers to build a better model. This paper presented a methodology for predicting aggregate flight departure delays in airports by exploring supervised learning methods. This way, we may be able to predict the delays of a new flight, without needing several months of data to build a prediction model. Another step forward would be to generalize the model to flights of the entire world, or at least to exploit more data sources, to build more complete predictions. Finally, the most interesting step would be to integrate such a model into a flight booking tool, to provide the delay prediction to future passengers, even this would require a strong confidence in the information provided, considering the possible impact in terms of reservations.

REFERENCES

- [1] We chose the “Airlines Delay” data from www.kaggle.com/datasets.
- [2] 2018 U.S. Airlines Delay Analysis- <http://www.milantomin.com/2018-u-sairlinesdelay-analysis/>.
- [3] Taxi in/out time means the time when a flight wheel was on/off to the time the flight gate in/on time. See Aviation System Performance Metrics: https://aspmhelp.faa.gov/index.php/ASPM_Taxi_Times:_Definitions_of_Variables.

- [4] statisticsbyjim.com/regression/choosing-regression-analysis/.
- [5] towardsdatascience.com/the-poisson-distribution-and-poisson-process-explained-4e2cb17d459.
- [6] Fox news reporter Rick Seaney article “Do flights ever leave early? And 4 other common travel questions”, <https://www.foxnews.com/travel/do-flights-everleaveearly-and-4-other-common-travelquestions>.
- [7] <https://developers.amadeus.com/flight-delay-predictionmachine>.
- [8] developers.amadeus.com/flight-delay-prediction-machinelearning.
- [9] developers.amadeus.com/flight-delay-prediction-machinelearning.
- [10] developers.amadeus.com/flight-delay-prediction-machinlearning.
- [11] towardsdatascience.com.
- [12] blogs.mprnews.org/newscut/2014/12/bigdata-and-the-delayed-flight.
- [13] www.icao.int/annual-report-2018/Pages/the-world-of-airtransport-in-2018.aspx.
- [14] www.icao.int/annual-report2018/Pages/the-world-of-air-transport-in-2018.aspx.
- [15] developers.amadeus.com/flight-delayprediction-machine-learning.
- [16] www.foxnews.com/travel/do-flights-ever-leave-early-and-4othercommon-travelquestions