

DATA 601 | Project 2| NY-DC/MD/VA Flights

Chiranjib Dutta, Asad Khan

University of Maryland Baltimore County

1000 Hilltop Circle, Baltimore, MD 21250

(Dated: December 19, 2021)

INTRODUCTION

In today's age flight delays have been a significant problem for aviation companies. This issue has cost billions of dollars to aviation companies and even deterred some customers because of losing time from delays. Some people might miss time sensitive meetings or may be unable to make it to important celebrations and events with their families due to these delays. Tackling this problem is important. We have a collection of flight data that was granted to us to produce reports to help analyze patterns and make inferences from the results. In this report, we will compare the data with the cause of cancellation/delays of flight, understand and predict the trend based on flight data from the year 2013 using data analysis and machine learning tools to help give insight as to what may cause flight delays.

1.1 Calculate the total number of seats for all the planned flights for each destination separately?

We begin with calculating the total number of seats for each planned flight that departs from their origin (EWR,JFK and LGA) to their destination airports (BWI,DCA,IAD) and we group the data by the origin to summarize flights from where they depart from. Figure 1 illustrates this step in the picture. We then list the seats for each of the flights that we have data on. Upon observation, we come to know that the flights that took off from LGA to destination (BWI,DCA,IAD) had the most seats in their respective groups. The result is illustrated in Figure 1.

		Total_seats
origin	dest	
EWR	BWI	46450
	DCA	94790
	IAD	67850
JFK	BWI	47540
	DCA	96460
	IAD	121530
LGA	BWI	21300
	DCA	714970

Figure: 1

1.2 What was the day of the year with the highest number of flights?

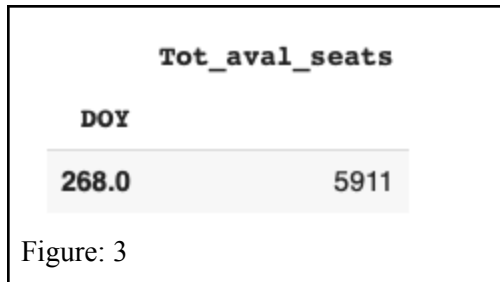
For this information, we started by gathering data from our main data frame which has all the information we need regarding flights. We grouped by the “day” column and applied an aggregation method on flights to count the total number of flights that took off on each day. This allowed us to sort the total flights in ascending order. By doing this we were able to find out that on day 86 the total number of flights were 56 which had the highest number of flights in comparison to the other days. Figure 2 illustrates the data we were working with sorted in descending order. The result is illustrated in Figure 2.

Total_flights	
DOY	
86.0	56
10.0	55
73.0	55
38.0	54
17.0	54
...	...
333.0	18
358.0	17
359.0	16
365.0	13
366.0	5

Figure: 2

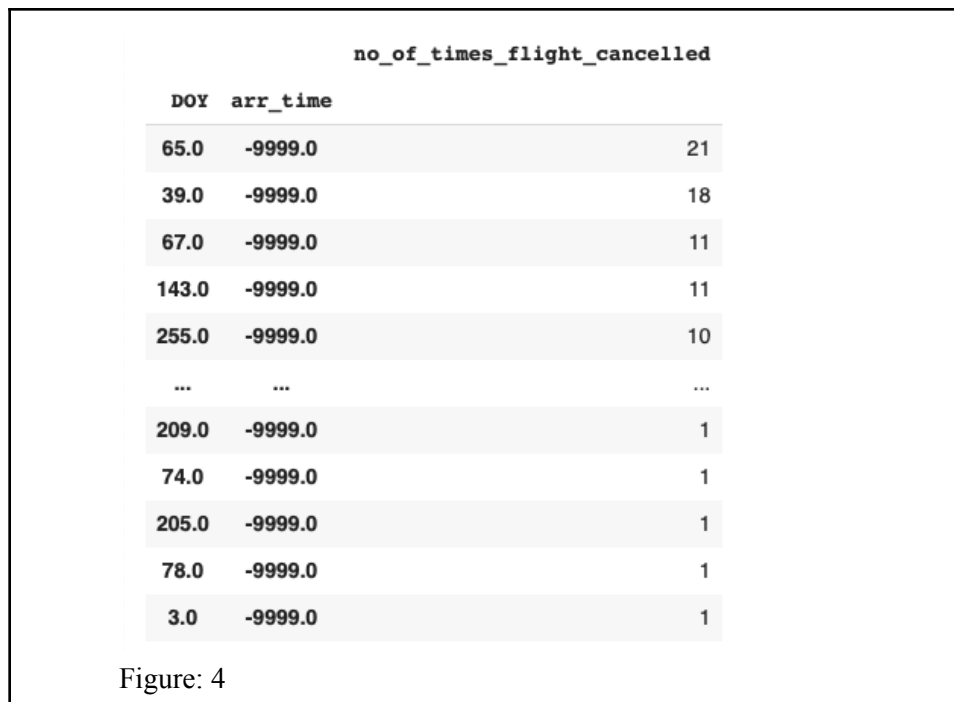
1.3 What was the day of the year with the highest number of seats?

In order for us to find out the day of the year with the highest number of seats, we again started by grouped by the “day” column and applied an aggregation method on seats to find the sum of the total number of seats that took off on each day. Next, we sorted the total seats in ascending order in order for us to be able to select and isolate the target number we were searching for. We were able to confirm that on day 268, the total number of seats were 5911 which was the highest number in comparison to the other days. The result is illustrated in Figure 3.



2.1 What day of the year do most cancellations happen?

For this day, we first decided to use the data in the “arrival time” field. The reason is because we know if the arrival time is (“Empty”, “None” or 0) this will mean that the flight never arrived or it got canceled. After that we start changing all the (“Empty”, “None” or 0) into an arbitrarily selected integer (such as -9999) to indicate that this will represent a cancellation. Since the question asked us to find the most cancellations in a single day, our next step was to use this data on arrival times where the value represents a cancellation and then group it by the “day” field. We performed a count values and sorting method to find the highest number of cancellations in a single day. By following this method, we observed that the most canceling flights in a day were on the sixty five day and that there were a total of 21 flights that were canceled on day sixty five. The result is illustrated in Figure 4.



2.2) Is there any relationship between the weather datasets and cancellations? If yes, describe it and justify with some numbers.

We are convinced that there is a relationship between the weather datasets and cancellations. We used two different data sets to investigate if there was an observable trend between the cancellation flights and daily weather in New York. We decided to compare canceled flights with each column of New York daily weather and we noticed that “Precipitation” likely had the most effect on the cancellation so we took the results from two different data frames containing weather and cancellation data and then compared them with one another. Since all the data in precipitation was less than 1 and greater than 0, we decided to change all the (string “T”) values to zeros which represented empty rows for us. After this step we used a date time module to separate the number of the day from the full date and then load the day value in a new column we labeled day. Our next step was to use precipitation data which was grouped by “day.” We performed a sum and sorting method to add all the precipitation and then compare it with the cancellation flights as it is shown below in the picture. Upon observation, we noticed that the number of times a flight gets canceled can be associated with weather conditions. It appears that snowfall and precipitation seem to be higher on the same day, or right next to the day a flight is canceled and that when there is no snowfall or precipitation, cancellations are not as high. It is within reason to conclude that there is a relationship between weather and cancellations based on the data. The result is illustrated in Figure 5.

				no_of_times_flight_cancelled
DOY	Precipitation	Snowfall	arr_time	
65.0	1.00	1.0	-9999.0	21
39.0	0.78	2.9	-9999.0	18
67.0	0.27	3.0	-9999.0	11
143.0	0.03	0.0	-9999.0	11
255.0	0.29	0.0	-9999.0	10
...
209.0	0.07	0.0	-9999.0	1
74.0	0.00	0.0	-9999.0	1
205.0	0.00	0.0	-9999.0	1
78.0	0.46	0.0	-9999.0	1
3.0	0.00	0.0	-9999.0	1

Figure: 5

2.3 Is there any relationship between the Federal Holiday Schedule and cancellations? If yes, describe it and justify with some numbers.

We are convinced that there is a relationship between Federal Holidays and cancellations. In order for us to make a relationship between canceled flights and federal holidays more apparent, we used two different data sets such as (df_with_fed_hol and df_cancelled_flights) and merged it based on the day of the year column “DOY.” We used a group by columns “DOY” and “Fel_hol” and counted the total number of canceled flights (using the count() aggregate function). We came to the conclusion that only Labor day and veterans day had an effect on the cancellation flight, which was surprising to us because our expectations on the holidays to have the greatest impact on cancellations was Thanksgiving and

Christmas. Nonetheless, there is a noticeable impact that Federal Holidays have on cancellations and the data suggests that there is a relationship between the Federal Holiday Schedule and cancellations. The result is illustrated in Figure 6.

no_of_times_flight_cancelled			
DOY	Fed_hol	arr_time	
245.0	Labor Day	-9999.0	2
315.0	Veterans Day	-9999.0	1

Figure:6

2.4 What is the total number of seats for the canceled flights? If we assume the average flight price of \$50, what is the total economic loss?

For this, we created a new data frame from the original data frame which only had the information of all the canceled flights such as (seats, tail number, carrier etc). After that we took a sum of all the seats in (canceled flights/new data frame) and multiplied it by fifty in order for us to find total economic loss in dollars. We found out that aviation companies are losing millions of dollars every year. The result is illustrated in Figure 7.

Total cancelled seats	27690
Total loss \$	1384500

Figure: 7

2.5 Determine the ratio of canceled flights/planned flights for each airline company, list it, and determine the most and least reliable airline company (most reliable = the one that has the smallest ratio of canceled/planned)

carrier	
ratio	
0.002901	US
0.004493	9E
0.011852	B6
0.033816	WN
0.062193	EV
0.106109	YV
0.178571	MQ

Figure:8

For this approach, we created a new dataframe that had only information of canceled flights. We used a grouped by and count method to get the total number of canceled flights per airline carrier from that data frame and then add this information to our original data frame under the column name 'sp_airline_cancelled_tot' using lambda function.

In order for us to find out the total number of planned flights that reached their destinations. We used a grouped by and count method to get the total number of planned flights for each carrier from the original data frame and then add this information to our original data frame under the column name 'planned_airline_tot_flights' using lambda function. We created a new column called 'ratio' found the ratio between $((sp_airline_canceled_tot) / (planned_airline_tot_flights))$ and sorted by most reliable carrier

We came to the conclusion that the top 3 reliable airlines (US is the 1st, 2nd is 9E and 3rd is B6) and last 3 reliable airlines are (MQ is in the bottom, then YV and then EV). The result is illustrated in Figure 8.

3.1 Calculate the average arrival delay for all the flights that took place on the same day and plot it (x = 1:365, y = daily average delay). On the same plot, please mark the Federal Holidays from the federal-holidays-2013.xlsx dataset.

For this, we created a new data frame called (df_New) from the old data frame and dropped all rows where arrival time equals -9999. As we mentioned in previous questions, we used -9999 to fill the empty rows in the df["arr_time"] column. Since the new data frame had a 'day' column for the day of year number, the team decided to create a new column called (day of year "DOY"). Next, we used a group by method to find an average delay per day. We plotted this information on a bar graph by setting the x value to the total number of days in a year. Since we needed to plot the holiday on the same graph as on "daily average delay". So my team decided that we use the "dayofyear method" to take the day from the date column and put it under "DOY". Doing this made it easier for us to plot and compare different data sets on the same graph as it is shown in the picture below.

We came to the conclusion that the months (May, July, Oct, Nov and Dec) had the most delays of flights due to federal holidays compared to the daily average delay. The result is illustrated in Figure 9.

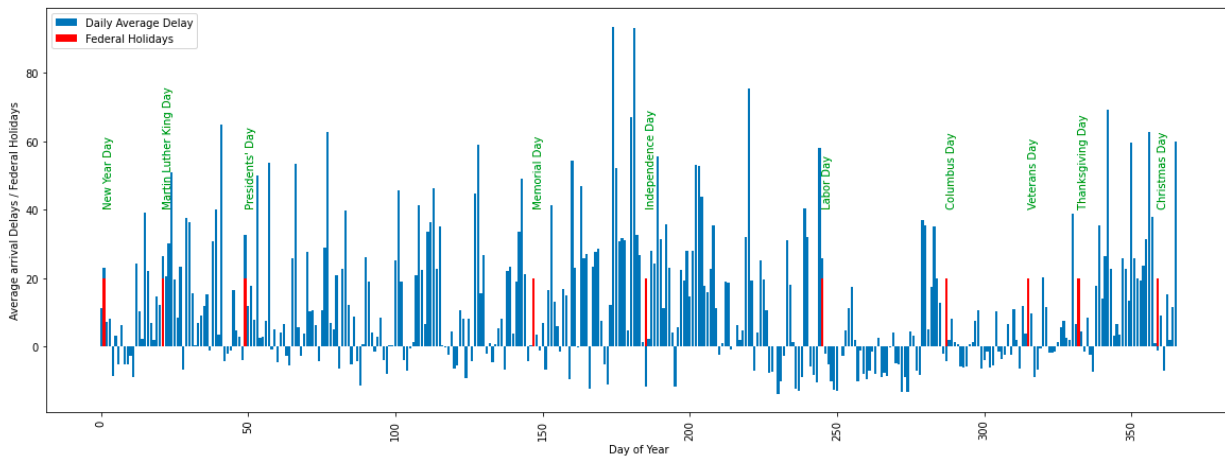
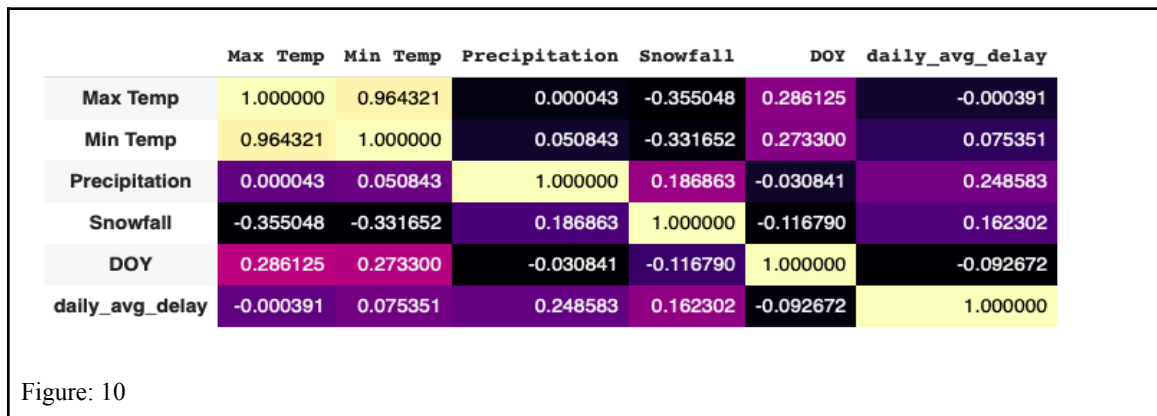


Figure: 9

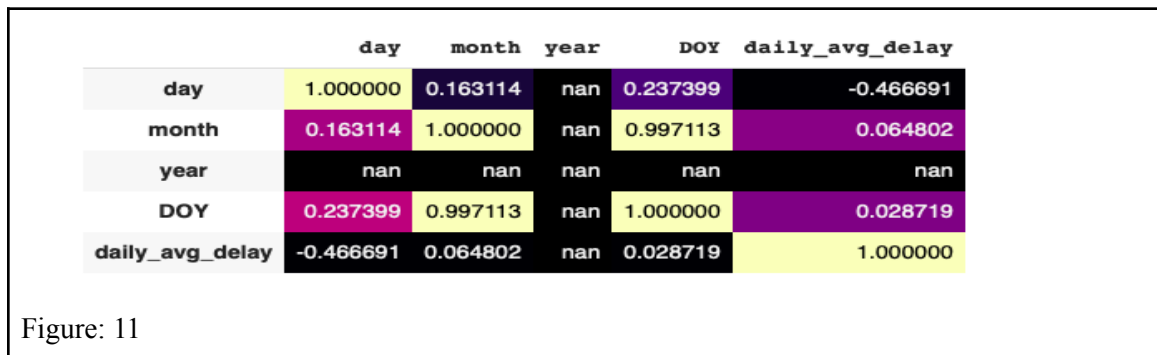
3.2 Is there a correlation between the weather datasets and daily average arrival delay?

We are convinced that there is a correlation between the weather and daily average arrival delay. We are starting with a data frame containing information on the daily weather of New York, “df_wthNYdaily” and set a new data frame called df_change = df_wthNYdaily while changing the index to day of year “DOY.” In our next step, we added the column called daily_avg_delay to the df_change. In order for us to find the relationship between 2 datasets we used a correlation method and compared daily average delay with every type of weather in the weather dataset. We came to the conclusion that most of the delayed flights were caused by precipitation and snowfall. For example, Precipitation had 25% and snowfall had 16% correlation with the daily flight average delay compared to other weather. The correlation map is illustrated in Figure 10.



3.3 Is there a correlation between the Federal Holiday Schedule and daily average arrival delay?

We are convinced that there is a correlation between Federal Holidays and daily average arrival delay. We started with the Federal Holiday Schedule and set a new data frame called df_change = Fed_hol while changing the index to day of year “DOY.” We then added a column called daily_avg_delay to the dataframe df_change to store the information we were looking for. In order for us to find the relationship between two datasets we used a correlation method and compared daily average delay with the Federal Holiday Schedule. We came to the conclusion that the daily average arrival delay is around 46% which is related to the federal holidays . The correlation map is illustrated in Figure 11.



3.4 Calculate the average arrival delay for all the flights for each arrival airport (e.g. IAD, DCA, and BWI) and determine most and least reliables (most reliable = the one that has the shortest average delay)

Using data that we mentioned earlier in 3.1 that we have in the dataframe called “df_New,” we chose to drop all the (empty “-9999 ”) rows in the arrival time for this task. We used the data frame to count the average arrival delay of each carrier using aggregate method and then grouped by the destination. We came to the conclusion that on the destination BWI, the most reliable carrier was WN with an average delay of 4.91 minutes and the least reliable carrier was EV which had an average delay of about 20 minutes. For the destination DCA, the most reliable carrier was DL which was 8 minutes earlier than the daily average and the least reliable carrier was MQ which had an average of about 28 minutes of average delay. For IDA, the most reliable carrier was UA which was 24 minutes earlier than the daily average and the least reliable YV carrier which had a delay average of 18.91 minutes. The result is illustrated in. Figure 12.

		aver_delay
dest	carrier	
BWI	9E	8.731288
	EV	20.056047
	WN	4.915000
DCA	9E	0.072206
	DL	-8.000000
	EV	21.340575
	MQ	28.195652
	UA	0.500000
	US	5.829000
IAD	9E	2.642623
	B6	12.805097
	EV	15.489133
	OO	3.000000
	UA	-24.000000
	YV	18.917266

Figure: 12

3.5 Calculate the average arrival delay for all the airlines and determine most and least reliables (most reliable = the one that has the shortest average delay)

In order for us to find arrival delay for all the airlines we used a carrier and arrival delay from df_New data frame. We then performed an aggregate method to find the average of delay and then grouped it by carriers and sorted by arrival delay. We came to the conclusion that the most reliable airline was DL which was 8 minutes earlier than the arrival time and the least reliable airline was MQ which was 28 minutes later than its arrival time. The result is illustrated in. Figure 13.

aver_delay	
carrier	
DL	-8.000
UA	-7.667
OO	3.000
9E	3.613
WN	4.915
US	5.829
B6	12.805
EV	17.360
YV	18.917
MQ	28.196

Figure: 13

3.6 What day of the week did we have the highest average delay?

weekday_delay	
week_day	
Wednesday	11.794330
Tuesday	11.176183
Thursday	13.106149
Sunday	10.094053
Saturday	11.124157
Monday	9.009183
Friday	10.903725

Figure: 14

To determine this, we used the datetime module to find a “day_name” from the date column and then set that value under the week_day column in df_New. We then perform an aggregate method to find the average delay and set the value under the new column called “weekday_avg” and then group it by week_day and sort by arrival delay. We came to the conclusion that most delays happened on wednesdays. The result is illustrated in. Figure 14.

3.7 Which one had a higher average delay: flights that took off in the morning (6 am to 10 am), noon (11 am to 2 pm), afternoon (3 pm to 5pm), or evening (6 pm to 10 pm)?

For this, we created a function called “day_session” and passed the “dep_time” column to find a higher average of delay for those flights that took off each session of the day. In order for us to return the correct part of the day. After that we created a new column in df_New called “session” and used the apply method to assign function return values under the new column. Furthermore, we perform an aggregate method to find the higher average departure delay. and then group it by session, sort by departure delay and then set these 2 columns to new dataframe called “dfNewSession”. We came to the conclusion that the flights that took off in the afternoon had the highest amount of delay which was 22.59 minutes of delay. The confusion matrix is illustrated in. Figure 15.

aver_delay	
day_session	
Afternoon (3pm-5pm)	22.599612
Evening (6pm-10pm)	21.839815
Noon (11am-2pm)	9.709850
Morning (6am-10am)	0.126671

Figure: 15

3.8 Determine the number of airplanes used in these flights manufactured by BOEING, EMBRAER, and AIRBUS separately.

manufacturer	
AIRBUS	4
AIRBUS INDUSTRIE	3929
BOEING	201
BOMBARDIER INC	3762
CANADAIR	1491
CANADAIR LTD	37
CESSNA	8
EMBRAER	4399
GULFSTREAM AEROSPACE	1
dtype: int64	

Figure: 16

To determine the number of airplanes used in these flights, we performed a group by method on the manufacturer column that we took from df_New and applied a size method to find the total number of planes by each manufacturer and then set those values to a new table called carrier_total_flights. We came to the conclusion that the Airbus industry had the most planes compared to other manufacturers. The result is illustrated in. Figure 16.

4. Build a linear regression model to estimate the arrival delay of the flights given in "flights_test_data.xlsx". Note that you have the full authority to decide what columns, what datasets (among the given datasets) to work with. In your report, please explain how you build the LR model and elaborate on its accuracy.

We have created a new data frame called "df_FTD" from flights_test_data.xlsx and inserted a column (day of year "DOY") to take the day from datetime and change it to day of year and then set that value under the DOY column in df_FTD. In order for us to prepare the data for testing we merged df_FTD and df_New on DOY and set this merging equals to df_test. In order for us to perform linear regression on arrival delay we will need to find out the column that has the strongest correlation with arrival delay column and we used a correlation matrix to find those columns such as (departure time, departure delay, DOY, arrival delay). We set all the values equal to X-Values that had the strongest correlation with arrival delay and set Y-Value equal to arrival delay. We split 80% of the data to train and another 20%, random state 1 and use it for testing while choosing a linear regression model. The conclusion is that we had a higher RMSE value of 16.21% and R2 was around 88%. We believe we could have got much more accurate results if we had larger amounts of data than we had to work with. The correlation map is illustrated in Figure 17. The RMSE and R2 are illustrated in. Figure 18.

Figure: 17

