

Analysis of Flight Disruption Data to Determine Problem Areas Utilizing Database Applications

Fadjimata Issoufou Anaroua*, Jordan Sanders, Sean Jennings, Tyler Martin

Embry-Riddle Aeronautical University, Daytona Beach, FL, 32114, USA.

*Corresponding Author

Fadjimata Issoufou Anaroua, Embry-Riddle Aeronautical University, Daytona Beach, FL, 32114, USA.

Submitted: 06 Feb 2023; Accepted: 13 Feb 2023; Published: 20 Feb 2023

Citation: Anaroua, F. I. (2023). Analysis of Flight Disruption Data to Determine Problem Areas Utilizing Database Applications. *J Math Techniques Comput Math*, 2(2), 110-116.

Abstract

Flight delays represent a significant issue to airline profits and passenger satisfaction. Many factors can lead to a flight being delayed and/or cancelled. To compound the issue, COVID restrictions have created a shortage of skilled, willing labor and a disturbance in the supply chain. All aspects of air travel have been affected. For example, aircraft maintenance and new plane constructions have increased lead times due to parts unavailability. Mass layoffs or early retirements as a knee-jerk response to reduced travel in 2020 has exacerbated the problem. Pilot strikes protesting mandates led to mass flight cancellations in 2021. The study will evaluate flight delays, cancellations, and incident data with the goal of visualizing which airports, airlines, cities, or states are experiencing the highest number of flight disruptions relative to others. In the era of new technological development, database applications are commonly used for data analysis to allow pattern recognition and large data distribution and organization. This study will mainly serve the purpose of flight data retrieval, the compilation of data and database design and finally output data visualizations. The outcome of the research will be presented in a user-friendly interface where the user can easily generate visualizations which can be used to analyze flight delay information.

Index Terms: Aviation, Data, Database, Computers, Information, Technology, Delays, Cancellations, Visualization, Website, Sql, Mysql, Html, Css, Python

Introduction

In 2007, the U.S government had endured 31–40-billion-dollar downsides due to flight delays in 2017, 76% of the flights arrived on time” [1]. Air travel has led to a global interconnection of society. The desire to maximize the amount of time an aircraft is being flown, and therefore making a profit, drives airlines to decrease the amount of time in between flights. This greatly increases the chances of flight delays as there is a smaller time window in between arrivals and departures [2]. Any disturbance or cyber-attack can rapidly create global effects by causing monetary and reputational harm. Today, flight delay data has never been more relevant due to an increased focus on efficiency and improving passenger satisfaction. Our team believes it is urgent to address flight cancellation data and build robust systems capable of addressing and managing flight delay information on a global scale. “An accurate estimation of flight delay is critical for airlines because the results can be applied to increase customer satisfaction and incomes of airline agencies.

Datasets

The group has created a system to visualize flight delay data which can be used to recognize a pattern. Some of the data

points the group has collected include origin, origin city name, origin state abbreviation, destination, destination city name, destination state abbreviation, quarter, month, day of week, day of month, flight date, marketing unique carrier, departure time, departure delay, cancellations, and causes of delay [3, 4]. The data was retrieved from the Bureau of Transportation Statistics and uploaded to the database application using MySQL. MySQL was queried in the back-end to create the data visualization. A web interface was implemented using programming tools such as Python and HTML5/CSS framework. The website allows users to select an airport, airline, city, or state. This input connects to the MySQL database through Python and uses data visualization techniques to show the user delay trends relative to that input in 2021.

Literature Review

People who have traveled by plane are familiar with one of the most inconvenient aspects of flying: delays. The plane may arrive late, there may be only one line for takeoff or landing, or severe weather may impose multiple hour delays (sometimes resulting in flight cancellation); regardless of the reason, flight delays are a major inconvenience for air travel passengers. As a result, with flight data from more than 300 thousand U.S. flights annually, we can acquire significant insights from this data to better understand flight delays and the associated causes.

Furthermore, because of the availability of large data sets for visualization, functional database systems may be used to help visualize these flight delays. This might be extremely useful for both travelers and corporations. “In 2016, research for a post-flight data analysis using databases implied that in conventional systems, aircraft data was typically recorded in the form of files on a storage medium and handled by each flight sortie separately. As a result, as the number of flight sorties grows, so does the time spent searching for data for analysis. Instead of file-based flight data maintained by flight sortie, they proposed database-based data integration and management. This solution allows people to simply store and manage all of the real-time data in a database so that people may use it for visualization and in a variety of application programs.” [5]. Recent research has found it difficult to explain the principal reasons behind flight delays due to multiple factors being in play at a time.

Flight schedules can be subject to change. Because airline resources are closely connected, delays could quickly spiral out of control if suitable recovery measures are not enacted. Despite the complexities, certain patterns of flight delays are consistent with the airline's schedule performance. The case study yielded some interesting outcomes [6].

“Air Carrier: The cause of the cancellation or delay was due to circumstances within the airline’s control (e.g., maintenance or crew problems, aircraft cleaning, baggage loading, fueling, etc.) Extreme Weather: Significant meteorological conditions (actual or forecasted) that, in the judgment of the carrier, delays or prevents the operation of a flight such as tornado, blizzard or hurricane.

National Aerospace System (NAS): Delays and cancellations attributable to NAS that refer to a broad set of conditions, such as non-extreme weather conditions, airport operations, heavy traffic volume and air traffic control.

Late-Arriving Aircraft: A previous flight with the same aircraft arrived late, causing the present flight to depart late.

Security: Delays or cancellations caused by evacuation of a terminal or concourse, re-boarding of the aircraft because of a security breach, inoperative screening equipment and/or long lines more than 29 minutes at screening areas.” [6].

Our research uses the Relational Database Management System to give flight delay statistical information to the end user via a website. Moreover, this website allows the user to choose an airline, airport, or a geographical location. The website also shows the user a pie chart on the distribution of delays and a graph representing the delay-based days of the week or airline.

A previous study by Wesonga, Nabugoomu, and Jehopio did an analysis of flight delays. They used a logistics model with twelve attributes and determined that “the number of freighter movements and non-commercial flights per day significantly influence both arrival and departure delays” [7].

Elmasri and Navathe stated that relational database management systems “provide flexibility to develop new queries quickly and to reorganize the database as requirements change.” Additionally, relational database management systems are “the dominant type of database system for traditional database applications” [8].

System Design

The image below shows the initial system design concept as well as the programming languages that were used to implement the system. The database is accessed when the user interacts with the web site by choosing the type of data he or she would like to interact with. The website was created with HTML5 and CSS. Once the user selects, the website requests information from the database server and then displays the data processed by the web application. Then, the web application uses a batch process to analyze the data from the database. Both the batch process and the web application used Python.

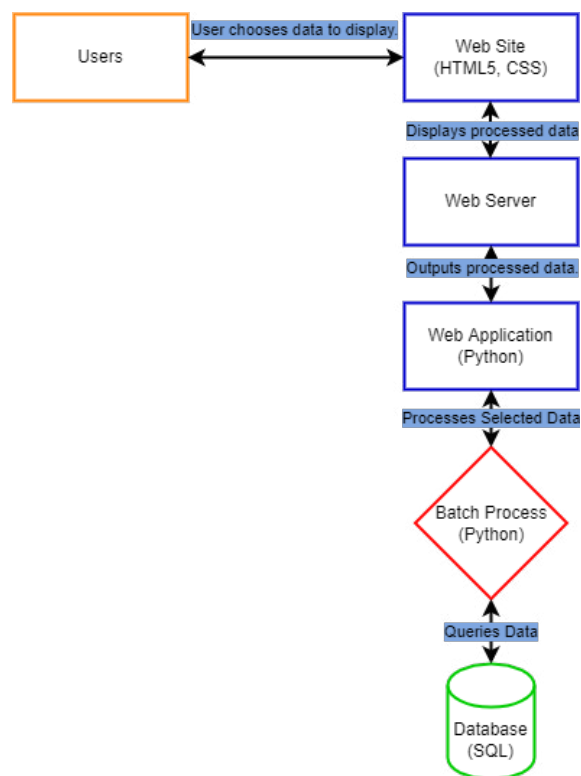


Figure 1: System Design Concept

Finally, the database is accessed by the batch process to query the data. The database system language used for this project is SQL.

Entities & Relationships

An entity is a real-world structure or a specific data set that we want to try to mimic in our database. They are frequently identified as the system's primary nouns. For our study, we have attempted to design and establish a direct link between an entity and multiple tables of data; as in a relational database, data for a single entity may be kept in multiple tables. Our entity relationship diagram is:

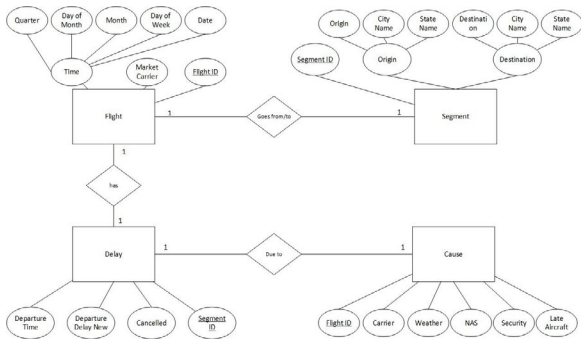


Figure 2: Entity Relationship Schema

As can be seen above, each the Flight and Segment entities are related to each other in a one-to-one relationship called Goes to/from. The Flight entity has a primary key called Flight_id. Furthermore, the Segment entity has a primary key of Segment_id which is a foreign key to the Flight_id of Flight. The primary keys of Delay and Cause also are foreign keys to the Flight_id of Flight. The Flight entity also has a one-to-one relationship with the Delay entity called Has. Additionally, the Delay entity has a one-to-one relationship with the Cause entity called Due to. The following image shows the UML (Unified Modeling Language) Diagram of the relational database.

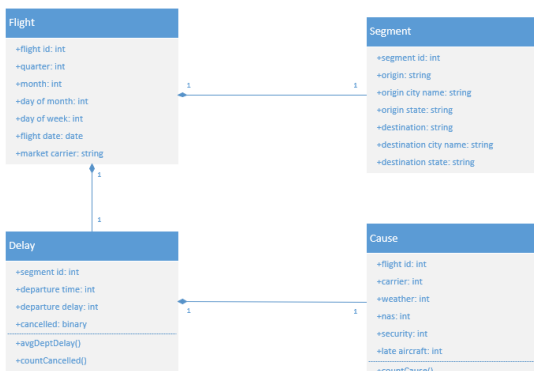


Figure 3: UML Diagram

The image below displays the use case diagram which illustrates how people interact with the system. The end user on the left side of the diagram selects his or her desired set of data on the delay prediction website from the front-end of the system. On the other side, the database management system (DBMS) interacts with the database server on the back-end of the system (i.e., creating the database).

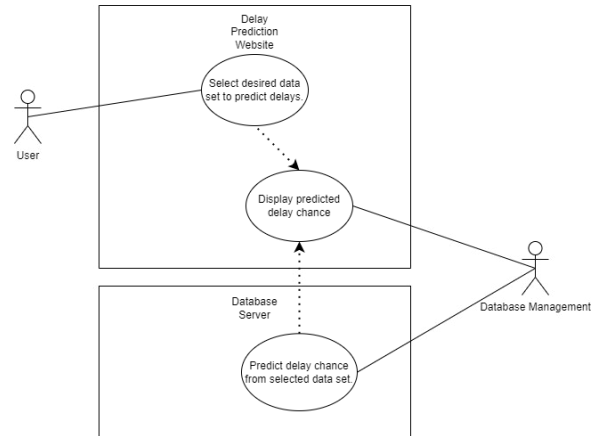


Figure 4: Use Case Diagram

The end user's desired set of data is analyzed to determine the specific 2021 delay information which is displayed to him or her. Once the data set has been selected, the DBMS accesses information stored in the database server to complete the task.

Other considerations

Once a pattern has been established, the website will present its findings to the user. This web interface provides the user with data visualizations of 2021's delays, allowing the user to make their own predictions about whether a flight is worth taking. For example, a user could insert an airport and see that Fridays had a large average delay compared to the other days of the week. The user could then decide that it might be better to fly out the day before or the day after.

Database Design

Database design, at its most basic level, entails defining entities to represent various types of data and designing relationships between those entities. By "entities," we mean how the data sets are related to each other. The tables below illustrate this project's four entities: Flights, Segments, Delays, and Causes. Inside each entity, there are attributes and primary keys which were discussed in more detail in the Entities Relationship section

Flights						
flight_id	Quarter	month	day_of_month	day_of_week	fl_date	Mkt_unk_carrier
1	1	1	23	6	01/23/2021	WN
2	1	1	7	4	01/07/2021	WN
3	1	1	7	4	01/07/2021	WN
4	1	1	7	4	01/07/2021	WN
5	1	1	8	5	01/08/2021	WN

segments						
segment_id	origin	origin_city_name	origin_state_abr	dest	dest_city_name	dest_city_abr
1	DAL	Dallas, TX	TX	SAT	San Antonio, TX	TX
2	MCO	Orlando, FL	FL	MDW	Chicago, IL	IL
3	LAX	Los Angeles, CA	CA	STL	St. Louis, MO	MO
4	LAS	Las Vegas, NV	NV	MCI	Kansas City, MO	MO
5	BNA	Nashville, TN	TN	ATL	Atlanta, GA	GA

Delays			
segment_id	dep_time	dep_delay_new	cancelled
1	928	0	0
2	923	23	0
3	816	16	0
4	1451	1	0
5	1436	113	0

Causes					
flight_id	carrier_delay	weather_delay	nas_delay	security_delay	late_aircraft_delay
1	44	null	0	null	null
2	54	0	null	null	0
3	11	0	7	0	null
4	null	null	null	0	0
5	null	0	0	0	15

Figure 5: Database Design

Results and Discussions

The relational database management system (RDBMS) was created via MySQL Workbench for the offline version of the program. The RDBMS consists of four tables: flight, cause, delay, and segment. These table structures can be seen in the Database Design section of this paper. The primary key of the flight table is flight_id. The other three tables have primary keys which are foreign keys to the flight table. Therefore, the flight table has the original unique id for each flight and each other table's id refers to the one created in the flight table. The code used to create this can be seen below.

```
1 CREATE TABLE FLIGHT
2
3 ( FLIGHT_ID INT NOT NULL AUTO_INCREMENT,
4   QTRER INT NOT NULL,
5   MNTH INT NOT NULL,
6   DAY_OF_MONTH INT NOT NULL,
7   DAY_OF_WEEK INT NOT NULL,
8   FL_DATE DATE,
9   MKT_UNIQUE_CARRIER CHAR(15),
10  PRIMARY KEY(FLIGHT_ID)
11 );
12
13 CREATE TABLE SEGMENT
14 ( SEGMENT_ID INT NOT NULL AUTO_INCREMENT,
15   ORIGIN CHAR(55) NOT NULL,
16   ORIGIN_STATE_ABR CHAR(55) NOT NULL,
17   DEST CHAR(55) NOT NULL,
18   DEST_CITY_NAME CHAR(55) NOT NULL,
19   DEST_STATE_ABR CHAR(55) NOT NULL,
20   PRIMARY KEY (SEGMENT_ID),
21   FOREIGN KEY (SEGMENT_ID) REFERENCES FLIGHT(FLIGHT_ID)
22 );
23
24 CREATE TABLE DELAY
25 ( DELAY_ID INT NOT NULL AUTO_INCREMENT,
26   DEP_TIME INT NOT NULL,
27   DEP_DELAY INT NOT NULL,
28   CANCELLED INT NOT NULL,
29   PRIMARY KEY (DELAY_ID),
30   FOREIGN KEY (DELAY_ID) REFERENCES FLIGHT(FLIGHT_ID)
31 );
32
33 CREATE TABLE CAUSE
34 ( CAUSE_ID INT NOT NULL AUTO_INCREMENT,
35   CARRIER_DELAY INT NOT NULL,
36   WEATHER_DELAY INT NOT NULL,
37   NAS_DELAY INT NOT NULL,
38   SECURITY_DELAY INT NOT NULL,
39   LATE_AIRCRAFT_DELAY INT NOT NULL,
40   PRIMARY KEY (CAUSE_ID)
```

Figure 6: Coding Practice

To load the data into these tables, the “import records from an external file” option was used. The csv files corresponding to each table were loaded using this option. The process of uploading the flight table data can be seen below

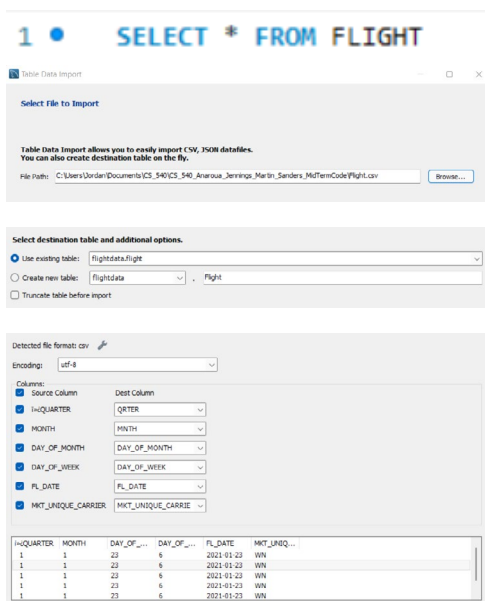


Figure 7: Process for uploading data

A website was created using HTML. It asked the user to search for an airport, city, state, or airline. The website then displayed data visualizations about delays associated with that particular airport, city, state, or airline. A static image of the website can be seen below

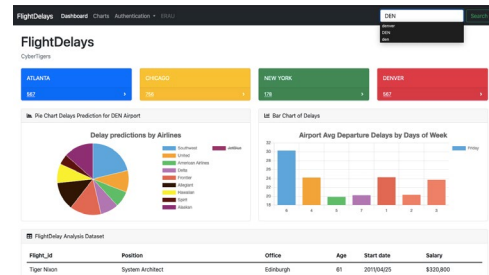


Figure 8: Website overview (Python Flask)

To connect the HTML website to the MYSQL RDBMS, the flask and flask_sqlalchemy Python libraries were used.

Testing Process

To test the queries needed and the data analysis to be shown to the user, some SQL queries were used. For example, delay by flight carriers. The delay_time column in the delay table is an integer value with a zero-value indicating no delay. Using this, we can calculate the average time of delay in minutes by averaging all values that are greater than zero with a SQL statement such as: `SELECT avg(delay.DEP_DELAY) FROM (delay INNER JOIN flight ON delay.DELAY_ID = flight.FLIGHT_ID) WHERE flight.MKT_UNIQUE_CARRIER = 'WN' AND delay.DEP_DELAY > 0`.

The group had the idea to create an internet accessible version of the website that would pull data from an SQL database hosted online on a web server. The web server was to interface with the user via PHPMyAdmin which is a program that was installed onto the server via Softaculous. PHPMyAdmin was chosen for the web server to make running SQL queries on the database more user-friendly through an intuitive interface. The csv files where be converted to SQL statements and queries to be entered in via the Poudel52_flightdata “SQL” tab of PHPMyAdmin. Ultimately the web server was not used for the final version of the project.

To test graphing and plotting techniques with the data, a Jupyter Notebook file was created. This file used the Matplotlib.pyplot library to implement the graphs. The Pandas library was also utilized to manipulate the data into the necessary slices to create the desired graphs. The data was sliced by Airport using an input to test one of the user inputs options on the website. After the user inputted the Airport Code, Pandas sliced the data in the different tables by that Airport Code. It was then further reduced by creating seven slices for each day of the week. The average delays for each of those days were plotted in a bar graph via Matplotlib. Pyplot as shown below.

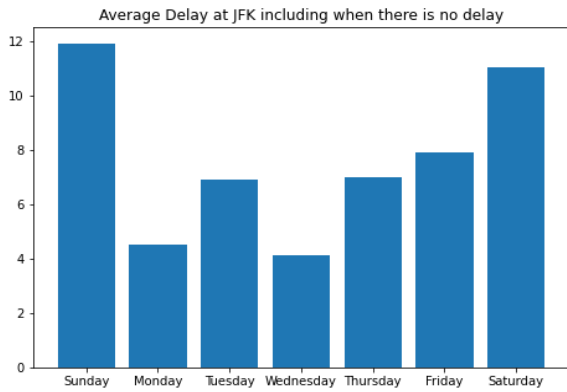


Figure 9: Average Delay bar chart (no delay)

This data was further sliced to only include instances that had a delay. Another bar graph was created to show the average delay of a flight at the user selected airport for each day of the week when there was a delay. An example of this graph is shown below for when the inputted airport is 'JFK'.

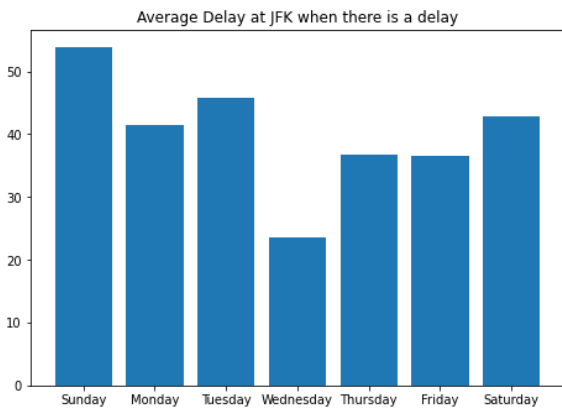


Figure 10: Average delay bar chart (with delay)

Furthermore, the Cause data set or table was sliced to only include instances from the inputted airport. Then, the total minutes delayed by all causes was found. The total minutes for each individual cause were calculated. Using the total minutes overall and the individual totals, percentages of each cause were found. This was then plotted via Matplotlib.Pyplot using the same input value of 'JFK'.

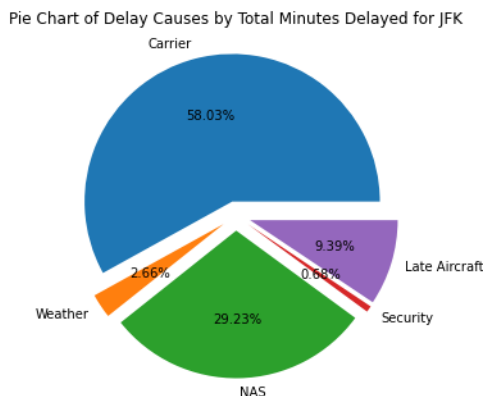


Figure 11: Pie Chart for causes of delay

At one point our group hosted our website, <https://flightdelays.poudel.tech/> using a web hosting service called HostGator <https://www.hostgator.com/> which is no longer in service. Files, such as index.html, were modified through cPanel. cPanel (<https://cpanel.net/>) is software installed under the server's operating system that provides a graphical user interface for interacting with a web server and modifying its files as an alternative to FTP (file transfer protocol).

A connection between MySQL and Python was created. This connection queried the flightDelay database and was utilized to plot the MySQL queries. The resulting charts displayed the averages and percentages of delays by choice of city, state, airlines, and airport within website search field.

The mysql.connector module was used to establish the Python connection to the MySQL database locally. Matplotlib, Pandas and NumPy were used in conjunction with mysql.connector to plot results of various delays and averages of delays within our data set in a form of python bar and pie chart.

The MySQL Database connection to Python allowed us to query average departure delays, percentages of delays by cause, and overall data set delays from our flightDelay database.

The selected queries were run through MySQL, and the result achieved allowed for displaying and graphically displaying various delays and problem areas within airports, airlines, cities, or states. This data set is comprised of ten (10) major airlines operating in the US from the BTS website. We included over 300 cities ranging from smaller to big cities with a higher amount of traffic, thus the chance of delay occurrence may be different dependent upon the location. Depending on the city, state, or specific airport, an overarching delay pattern was not determined for this study. Various results may appear with gaps, certainly due to the type of delay, the number of carriers' scheduled flights, and availability of data within the set. For the delay types from the overall data set, it was noticed that most delays were caused by carrier/airlines, NAS delay, and late aircraft delay. Furthermore, there were very few or no cancellations in most airports. Security delay and weather delay data appear to be non-significant. These delays have less chance of occurrence as they are often unpredictable which makes them less likely to cause huge delays within airports, cities, or states.

Some examples of queries and results are shown below.

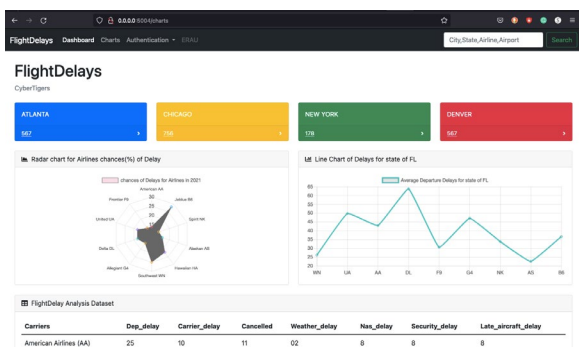


Figure 15: Average delay for Florida state

Lastly, the airline input results in two data visualizations: causes of departure delays as a line chart and chance of delays as a pie chart. An example of the airline input is shown below using Southwest (WN) as an example.

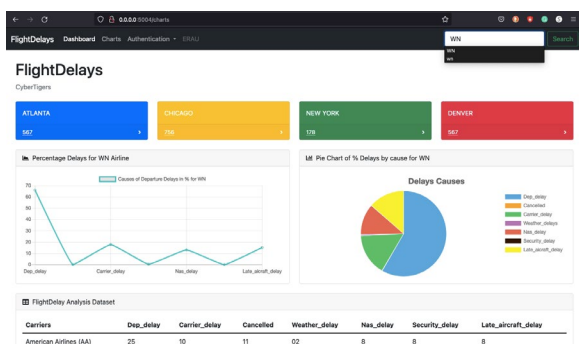


Figure 16: Delay prediction for Airlines

Conclusion

There are few modern-day inconveniences less anticipated than being stranded at the airport. Being able to make predictions and prudent choices based on historical flight delay information can save a significant amount of time, trouble, and tears. The group has created a system that directly interfaces with historical flight delay data in a database to visualize flight delay trends. An accessible, easy-to-understand graphical display of flight delay patterns from the past year could save travelers an untold amount of headache, stress, and trouble. Possible future applications of this concept could be delineated in the form of a phone app, an applet as part of an airline's website (for high-performing airlines), an informational service as part of a booking website, or simply a private independent resource for travelers and interested parties. Accessible information in a digestible format can be a powerful decision-making tool for a wide-reaching audience [9-12].

Future Work

This study could be furthered if more data is implemented into the DBMS. If there is more than just one year's worth of data, then certain machine learning algorithms could be applied. For example, logistics regression could be implemented to predict the likelihood of delay at a specified time and date. This would require attributes like time of day, time of week, month, and more. Another potential machine learning algorithm that could be implemented would be a form of regression (potentially linear or Gaussian-process regression) to predict the delay in minutes of a flight, which would require similar attributes. For any ma-

chine learning algorithm, more research would need to be done. It is important to note that the aviation industry is very volatile, meaning the slightest change in the economy or the general societal concern for safety can greatly affect the airline industry. For example, in 2001 and 2002, many people were scared to fly due to the events of 9/11. In 2008 and 2009, people preferred not to travel by plane due to the tight budget from the 2008 Financial Crisis. Lastly, the SARS-2-CoV-2 (COVID-19) pandemic greatly affected the air travel during 2020 and 2021. This means that machine learning for this application may not be as accurate due to the special circumstances surrounding the data set. Due to this fact (along with the current amount of data), this group decided it would be more beneficial to give the users the means to decide on their own what impact last year's delays would make on their flight choices today.

References

1. Yazdi, M. F., Kamel, S. R., Chabok, S. J. M., & Kheirabadi, M. (2020). Flight delay prediction based on deep learning and Levenberg-Marquardt algorithm. *Journal of Big Data*, 7, 1-28.
2. Rebollo, J. J., & Balakrishnan, H. (2014). Characterization and prediction of air traffic delays. *Transportation research part C: Emerging technologies*, 44, 231-241.
3. Bureau of Transportation Statistics (2021). Air Carriers: T-100 Domestic Segment (U.S. Carriers) (51276413) [Data set]. Bureau of Transportation Statistics.
4. Bureau of Transportation Statistics (2021). On-Time: Marketing Carrier On-Time Performance (Beginning January 2018) (51276413) [Data set]. Bureau of Transportation Statistics.
5. Shim, J. I., & Jo, G. S. Flight Data Visualization and Post-test Flight Data Analysis System by Using Database.
6. Anderson, A. B. A., Kumar, A. S., & Christopher, A. A. (2019). Analysis of flight delays in aviation system using different classification algorithms and feature selection methods. *The Aeronautical Journal*, 123(1267), 1415-1436.
7. Wesonga, R., Nabugoomu, F., & Jehopio, P. (2012). Parameterized framework for the analysis of probabilities of aircraft delay at an airport. *Journal of Air Transport Management*, 23, 1-4.
8. Elmasri, R., Navathe, S. B., Elmasri, R., & Navathe, S. B. (2000). *Fundamentals of Database Systems*. Addison-Wesley/publisher.
9. AhmadBeygi, S., Cohn, A., Guan, Y., & Belobaba, P. (2008). Analysis of the potential for delay propagation in passenger airline networks. *Journal of air transport management*, 14(5), 221-236.
10. Flightconnections. (2022). "all flights worldwide on a map!" (n.d.).
11. Wang, Y., Cao, Y., Zhu, C., Wu, F., Hu, M., Duong, V., ... & Stanley, H. E. (2020). Universal patterns in passenger flight departure delays. *Scientific reports*, 10(1), 6890.
12. Hassan, L. K., Santos, B. F., & Vink, J. (2021). Airline disruption management: A literature review and practical challenges. *Computers & Operations Research*, 127, 105137.

Copyright: ©2023 Fadjimata Issoufou Anaroua. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.