

Dataset

This report utilizes the NTSB Census of US Civil Aviation Accidents. This dataset contains extensive information on the vast majority of aviation accidents and incidents since 2008. This includes everything from the date and location of the accident to the flight hours of each crew member and much, much more. This report uses data from April 23rd, 2024.

Data Collection

Data was downloaded from <http://data.nts.gov/avdata>. The NTSB distributes the data in an MDB file that must be loaded into Microsoft Access. I then exported the tables that I would be using as Excel files, which I finally converted into CSV files that were loaded into a MySQL server. Once in the MySQL server, I dropped columns of entirely null values and columns that consisted of unimportant data for this report such as the NTSB employee that input the data.

Background Information

This project is inspired by the work of Xiaoge Zhang, Prabhakar Srinivasan, and Sankaran Mahadevan in their paper "Sequential deep learning from NTSB reports for aviation safety prognosis" (<https://doi.org/10.1016/j.ssci.2021.105390>). During my reading of their report, I found their problem statement had the order of operations in reverse. In their paper, they began with the textual "Sequence of Events" and predicted the damage done to the aircraft, if fatalities occurred, and whether the event was categorized as an accident or incident. I believe that is backwards, as in an investigation, the sequence of events would be written last after all facts are known, and the predicted information would be some of the first information to be known.

Problem Statement

Accordingly, their research inspired two main questions that will be addressed in this report:

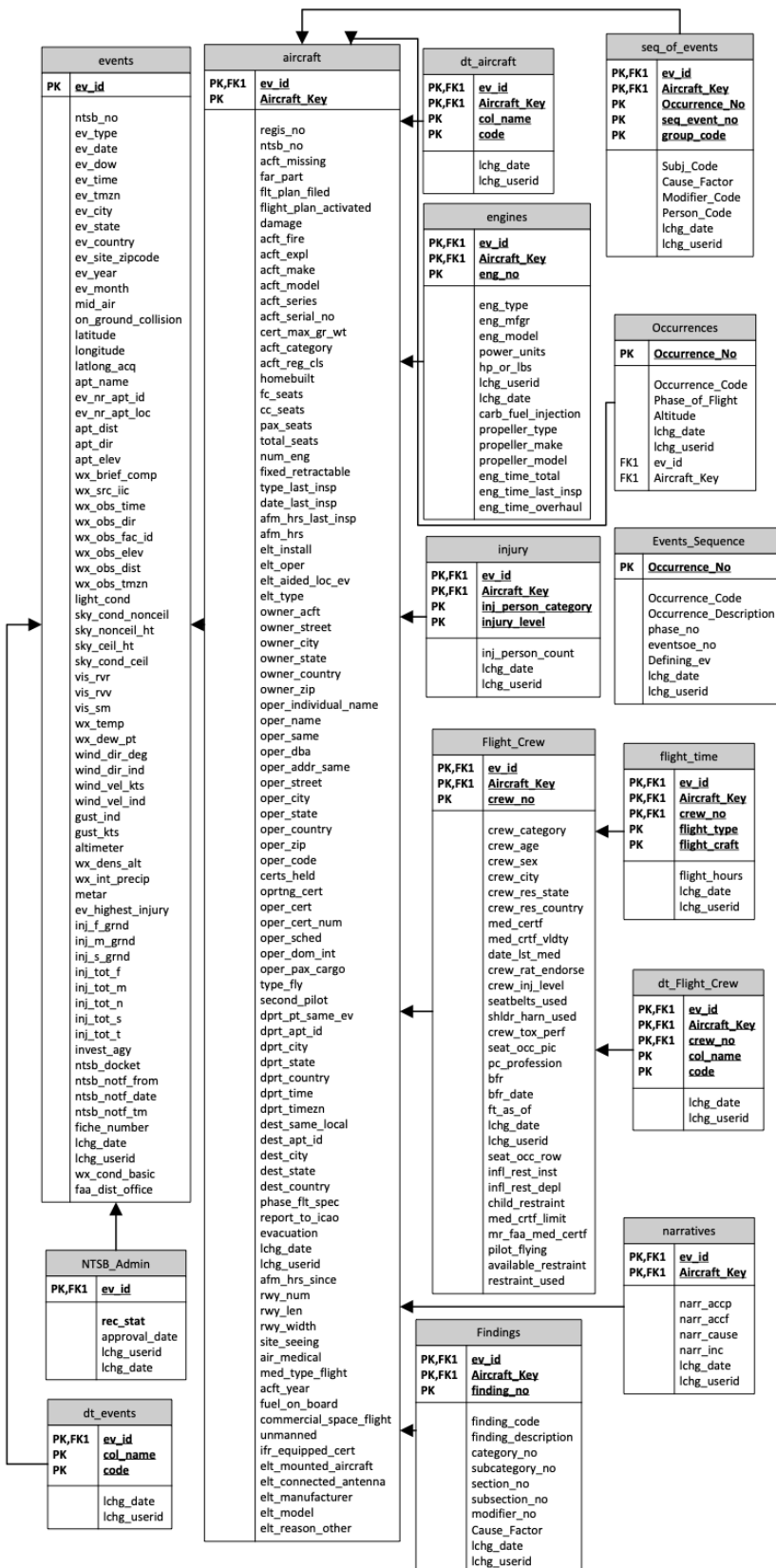
1. Can surface-level features (eg. data that can be immediately gathered at the beginning of an investigation) be used to predict the cause of the accident?
2. What does this tell us about the causes of accidents?

I hope that answering these two questions will add to the wealth of information surrounding the causality of aviation accidents and will help to reduce the amount of avoidable aviation accidents.

SQL Interactions

After my cursory cleaning in MySQL, I then imported the data into this Jupyter Notebook for further preparation and analysis.

-Data Schema



Cleaning

In this cleaning stage, I started with general filters on the events table, then cleaned each column of the table. I then merged the child tables onto the cleaned events table to filter away data on the events that were dropped from the events table. I then went through each column of the child tables cleaning each. Finally, I began early EDA by showing details about each column. Exact steps are shown in code.

Findings

As part of the cleaning process for the findings table, I dropped any of the findings beyond the 5th finding. This was to reduce the repetition in the training data to a manageable amount while preserving the quality of multiple findings.

EDA

*Beyond the analysis shown during the cleaning stage

For my EDA, I merged all the tables into one. While the above visuals show the distribution for each column, I needed to know the correlation and association of the variables. To resolve this, I performed a Chi² Test for Association between all the categorical variables, a Correlation Test between all of the continuous variables, and a One-Way-ANOVA Test between the categorical and continuous variables.

Encoding Categoricals

In my research, I found conflicting and confusing information on the many different strategies for encoding high cardinality nominal categorical features. Three, however, were regularly recommended: One-Hot-Encoding, Binary Encoding, and Feature Hash Encoding. I chose to test all three and pick the strategy that gave the best results. For the features that only had two options (Y/N | 1/0 | T/F | etc.), I binary encoded them in all three instances.

Main resources:

<https://www.kdnuggets.com/2021/05/deal-with-categorical-data-machine-learning.html>

<https://kantschants.com/complete-guide-to-encoding-categorical-features#heading-nominal-categorical-features>

Modelling (using Random Forest Classifier)

For modelling, I chose an 85%, 15% train-test-split to maintain a large training set, while still having plenty of records for validation. With the 85% test set, I then performed 10-fold cross validation to train a sklearn Random Forest Classifier for 17-class classification. Finally, the top performing model -the Binary Encoded model- predicted the top 4 causes for the accident.

Results

The results of the model validation against the test set are below.

Accuracy: 82.94%

Recall: 71.05%

Precision (Low-Bound): 46.43%

Precision (High-Bound): 71.05%

Category Performance

The model was most accurate at predicting classes 12, 16, 24, 32, 33.

The model had it's highest recall predicting classes 16, 24, 26, 32, 33.

Finally, the model was most precise on classes 16, 24, 26, 32, 33.

- 12 - Aircraft-Aircraft Systems
- 16 - Aircraft-Aircraft Oper/Perf/Capability
- 24 - Personnel-Action/Decision
- 26 - Personnel-Task Performance
- 32 - Environmental-Physical Environment
- 33 - Environmental-Conditions/Weather/Phenomena

Interpretation

To address my two main questions:

1. Can surface-level features be used to predict the cause of the accident?

In general, yes. The accuracy above 80% shows that the surface level features do hold information that indicates a potential cause for an accident. In addition, the continuous features seem to show more about the particular causes, so the model may be more accurate with features that I had to drop due to data cleanliness, like air pressure.

2. What does this tell us about the causes of accidents?

There are interactions between temperature, age of crew, wind velocity, total non-injuries, number of passenger seats, number of flight crew seats, day of week of the flight, and month of the flight that influence the cause of the accident. More research must be done on the particulars of the accident, but it is not unreasonable to assume that a summer month, weekend day of week, high age, and low number of flight crew will result in a higher chance of an accident caused by human error. Conversely, high temperatures reduce air density and lower overall aircraft performance, and high wind velocities can overspeed the aircraft. Therefore, those may point to environmental factors as the cause of the accident. Adding strength to these hypotheses, the model has the highest scores in predicting the causes that would be most closely associated with each.

Model Issues

The model was best at predicting the most common causes of the aviation accidents. While, the model still may be useful, this suggests it may have trouble identifying the more obscure causes for accidents. Possible future solutions could include upsampling the underrepresented causes, downsampling the most common causes, or using a model that allows for class weighting.

Acknowledgements

Finally, I would like to thank Dr. Shashi Jha (College of Charleston) for his guidance on this report; Jess Thomas (NTSB) for sharing the causes corresponding to the finding codes; and Xiaoge Zhang, Prabhakar Srinivasan, and Sankaran Mahadevan for inspiring this report.