

Airline Delay Analysis

Big Data Spring 2022

Introduction

Information

- Analyze a series of data sets related to airlines and airports
- Create metrics to visualize what may cause delays, and their trend
- Provide a series of visualizations that show these trends
- This will be helpful as we all travel this season

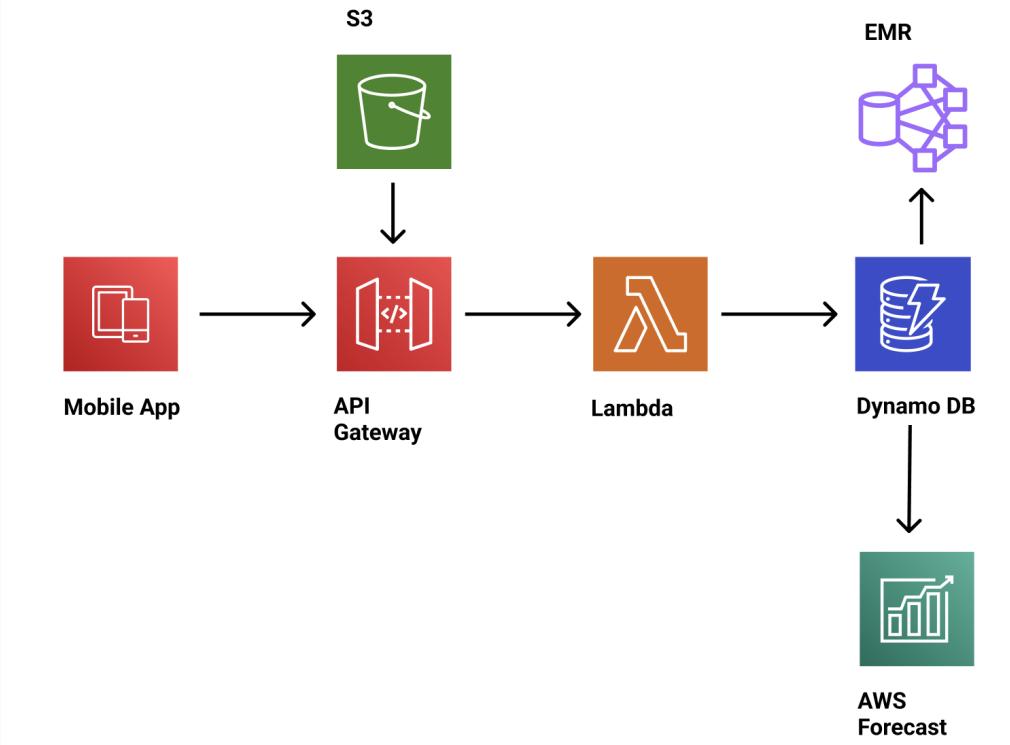
Architecture

- Used Kaggle, Amstat
- Spark Implementation, NumPy, Pandas
- Processing 10gb of data requires a distributed system
- Scalable programming that can be implemented on future data streams

Big Data

- Large Scale Data Set
- Attempted running our calculations locally
- Errors occurred and calculations were taking way to long
- Sparks distributed nature solved our problem
- Data cleansing/extraction was practically impossible locally

Mobile App/ Prediction



Preliminary Analysis

specifications and considerations

- representation of time
 - *points vs durations*
 - float data type
- standardization
 - centralized record-keeping
 - nomenclature consistency: IATA airport codes
 - objectivity of time data
 - time zone recording
- comprehensiveness of data fields
 - events/stages of flight process
 - *planned and actual times*

Preliminary Analysis (cont'd)

tradeoffs and implications

- choice of local time zones over single standard
 - difficult/unintuitive time comparisons
 - real-world practicality
- inclusion of seemingly-redundant fields
 - data volume
 - elimination of complicated calculations

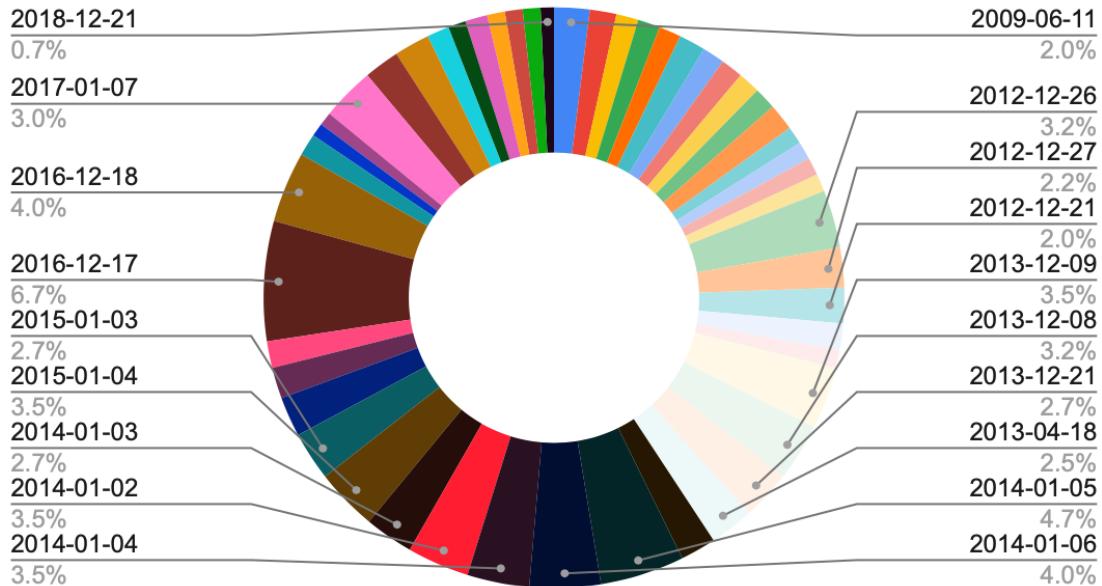
Preliminary Analysis (cont'd)

recognizing missing/
incomplete data

- “expected”/necessary fields
- cancelled flights
- exceptions

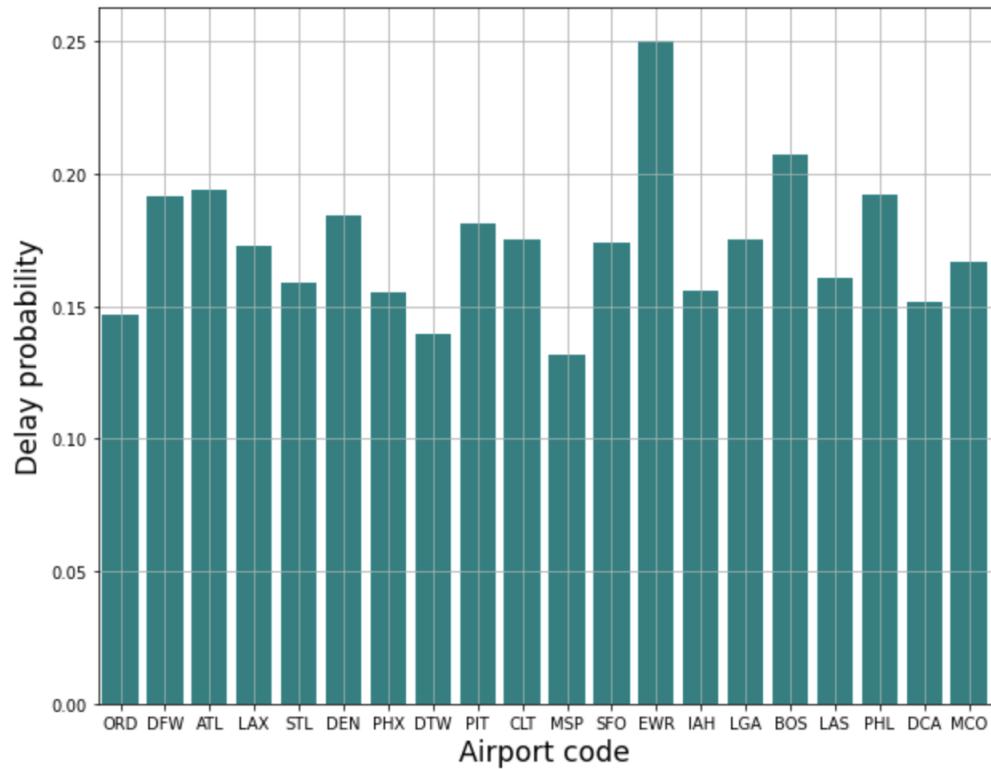
Least Reliable Day to Fly

Delay vs. Date



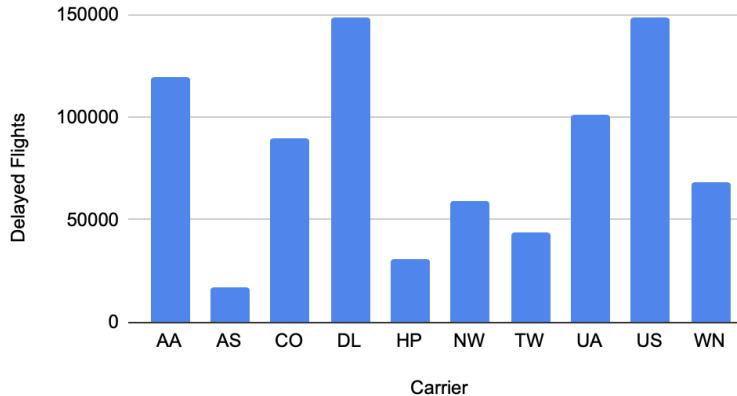
1994

Delay probability per top 20 busiest airports

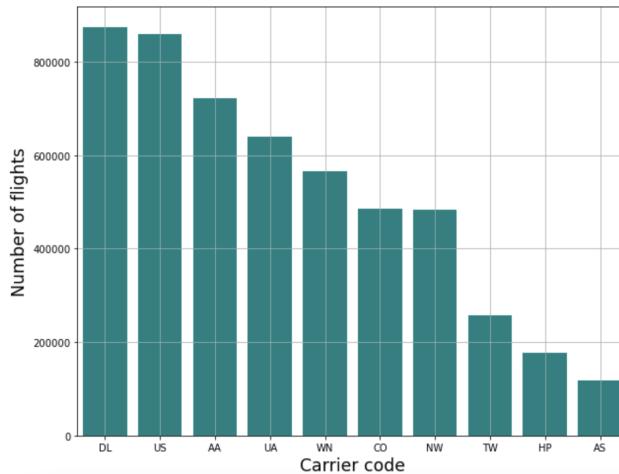


1994

Delayed Flights vs. Carriers

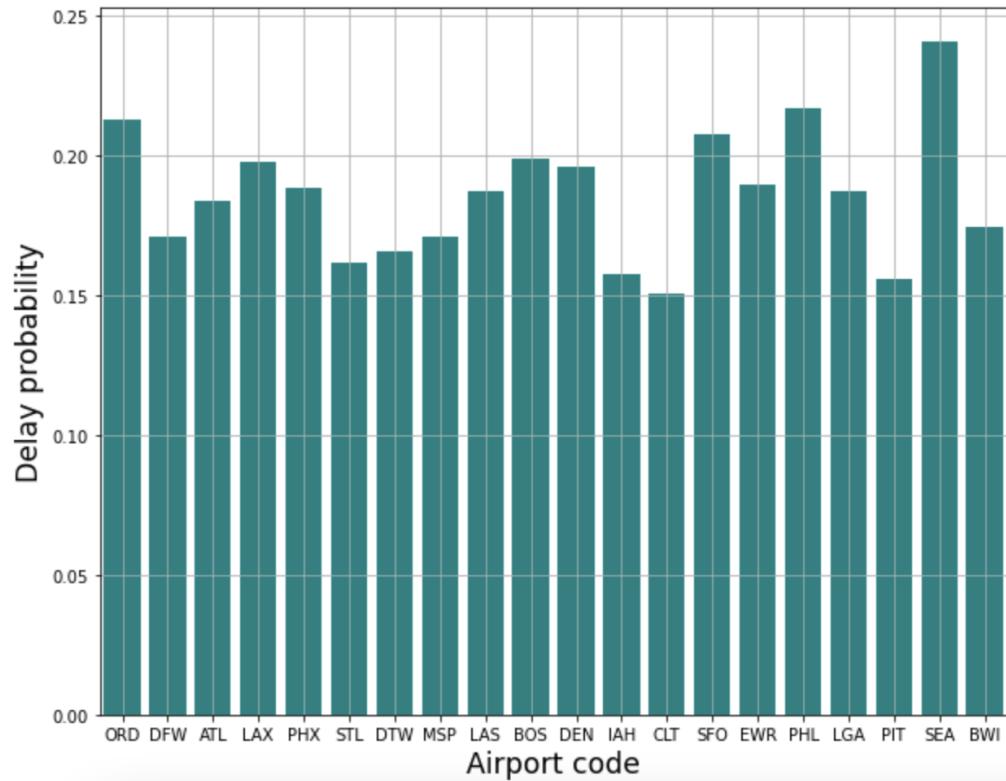


Top 20 busiest carriers



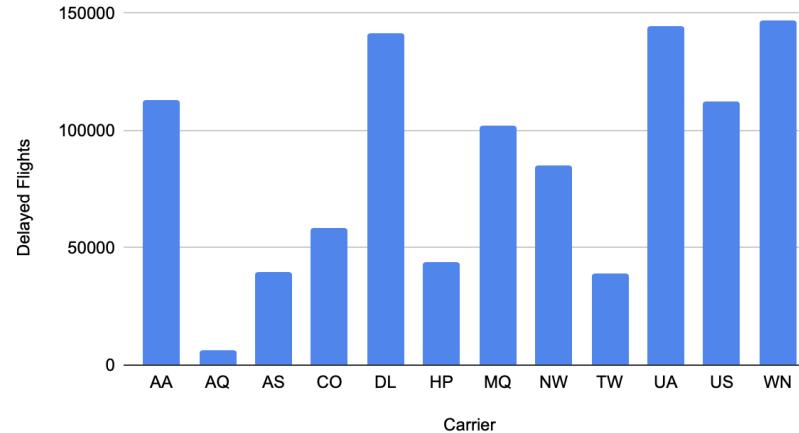
2001

Delay probability per top 20 busiest airports

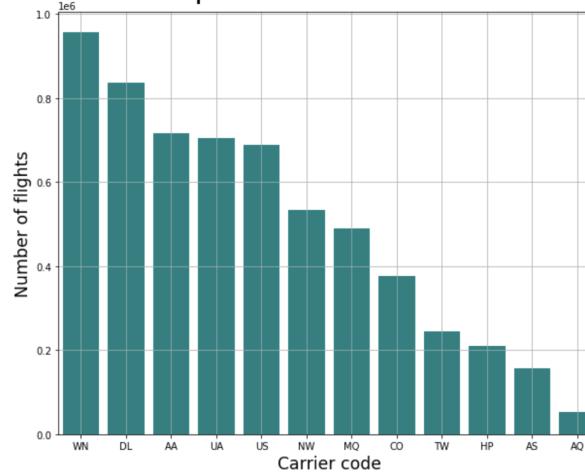


2001

Delayed Flights vs. Carrier

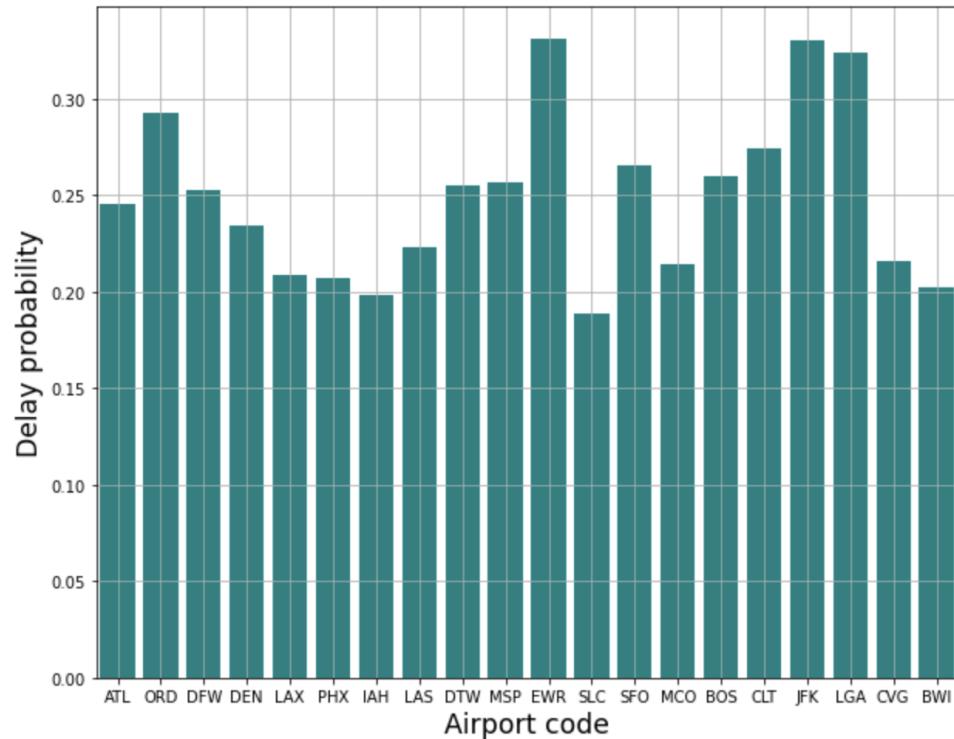


Top 20 busiest carriers



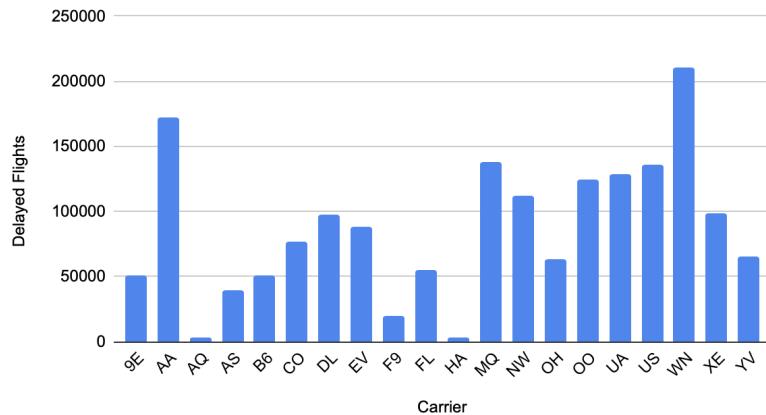
2007

Delay probability per top 20 busiest airports

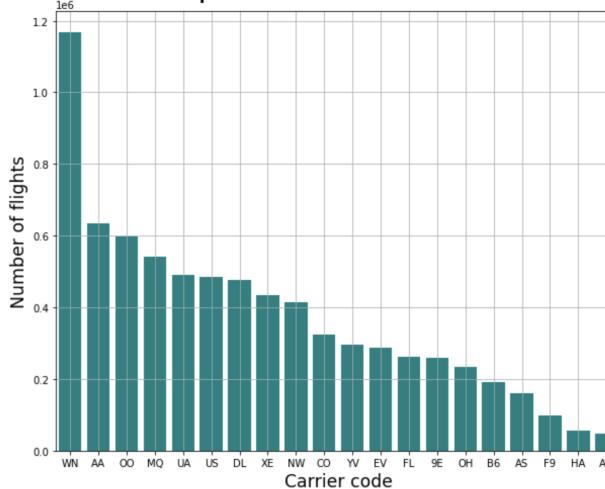


2007

Delayed Flights vs. Carrier

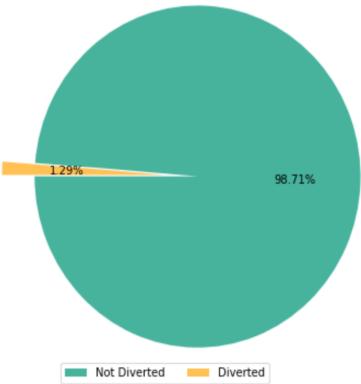


Top 20 busiest carriers

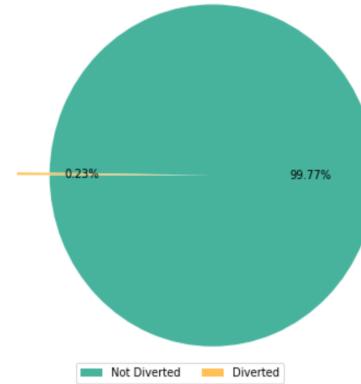


1994

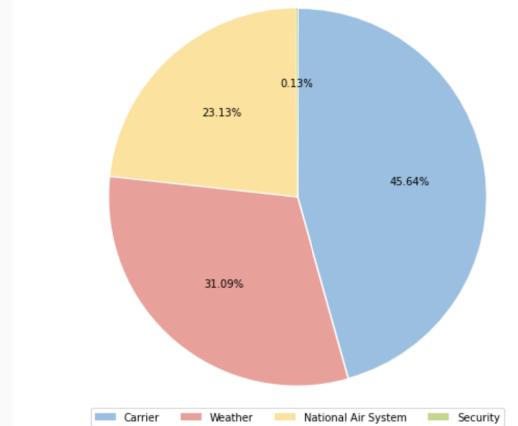
Cancelled Flights Ratio



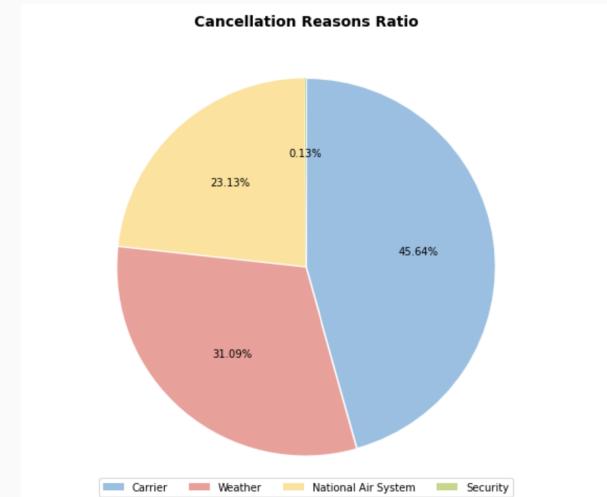
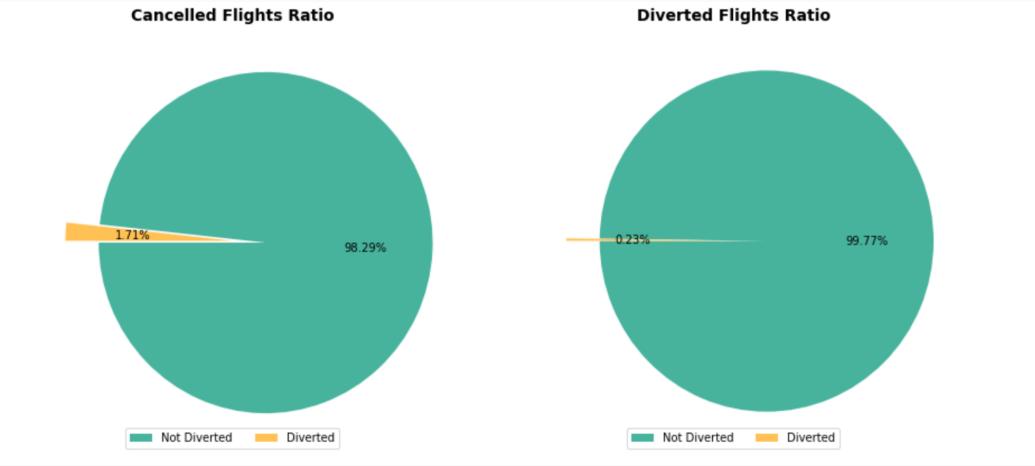
Diverted Flights Ratio



Cancellation Reasons Ratio

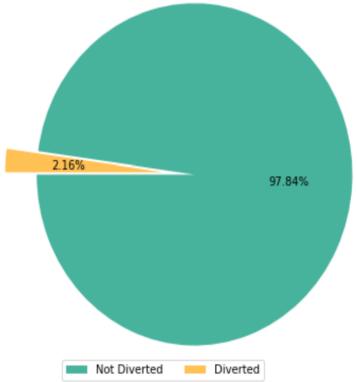


2006

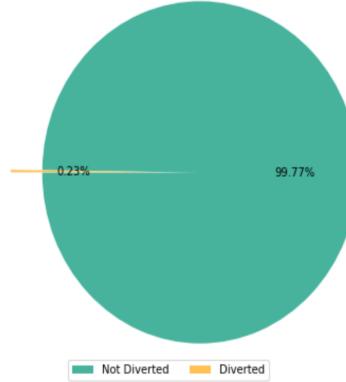


2007

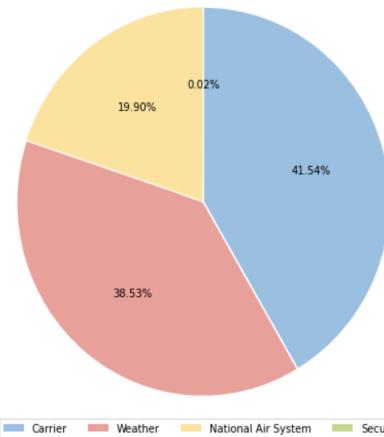
Cancelled Flights Ratio



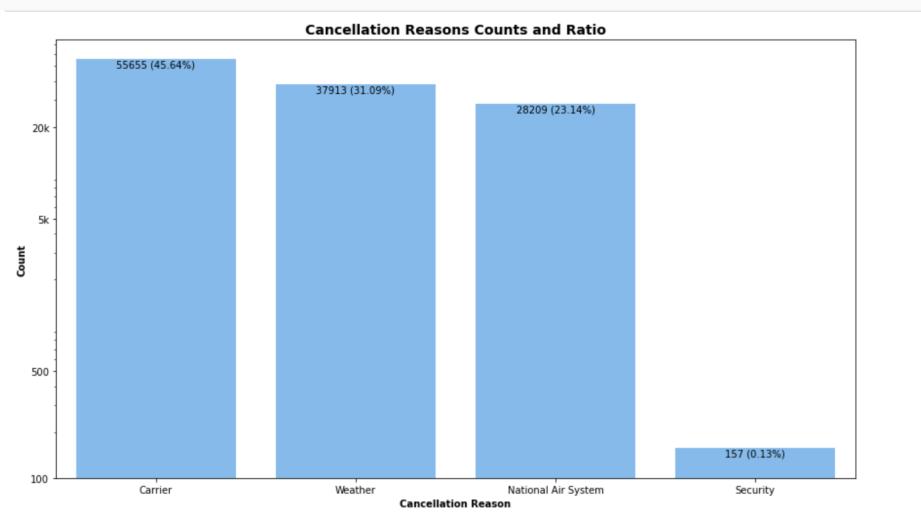
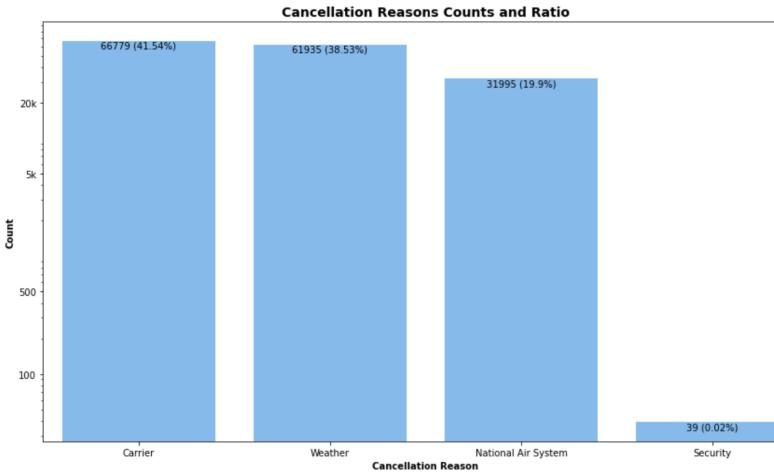
Diverted Flights Ratio



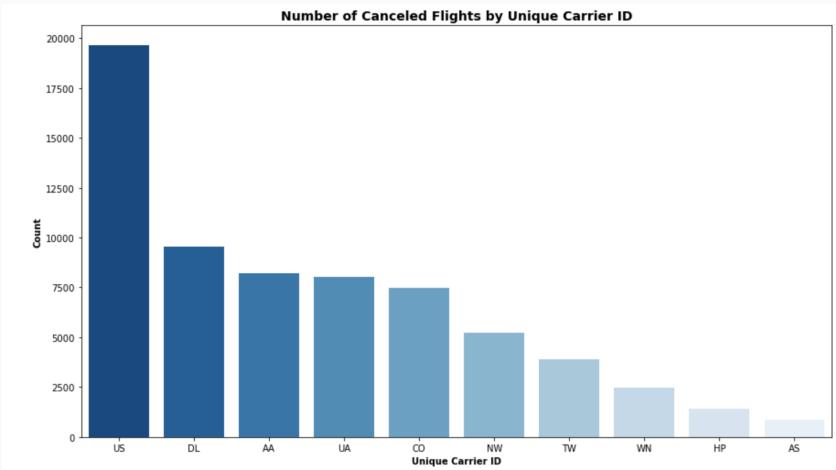
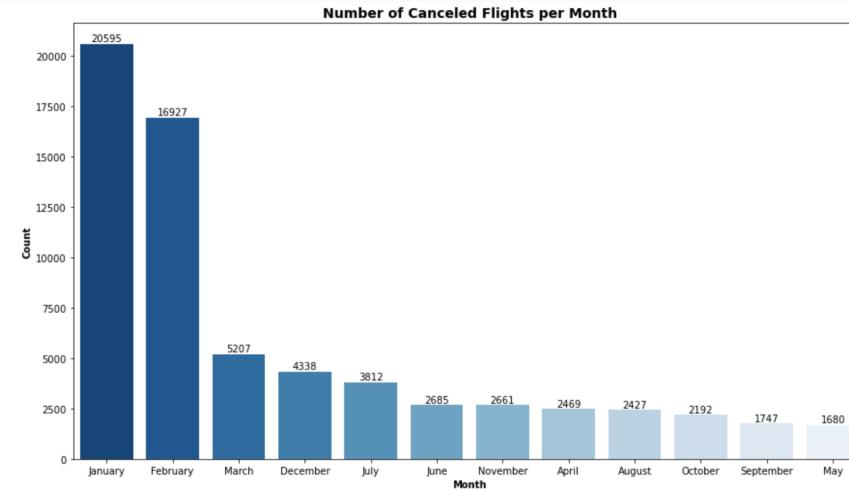
Cancellation Reasons Ratio



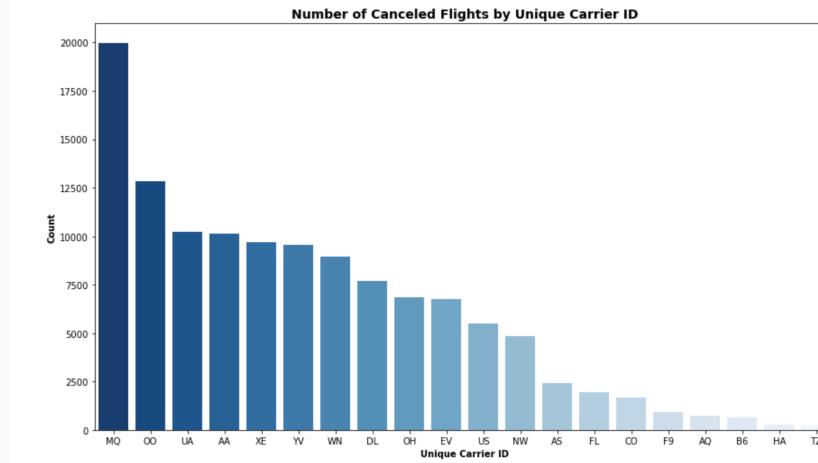
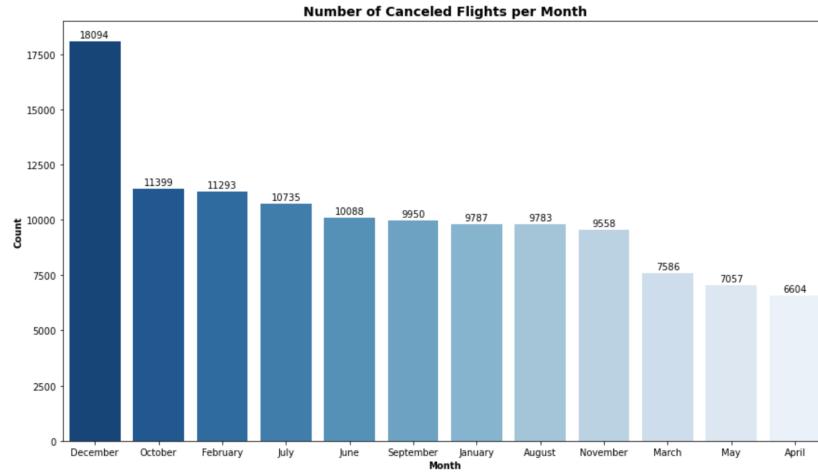
2006 v 2007



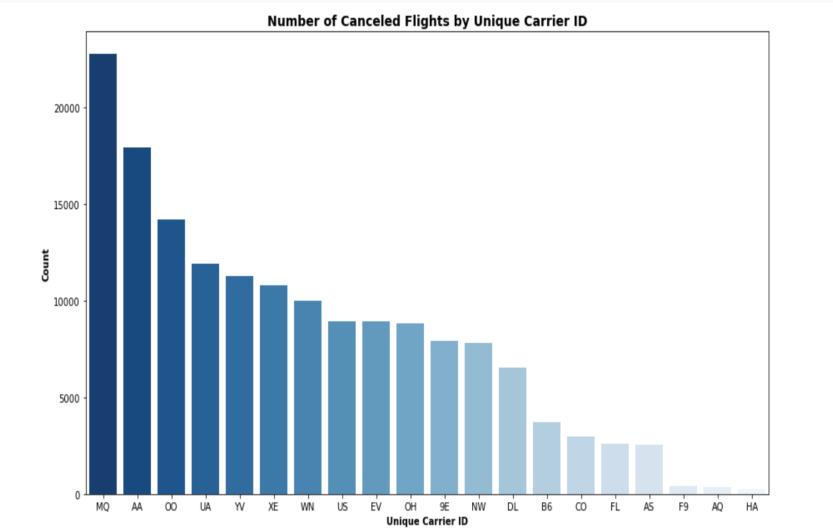
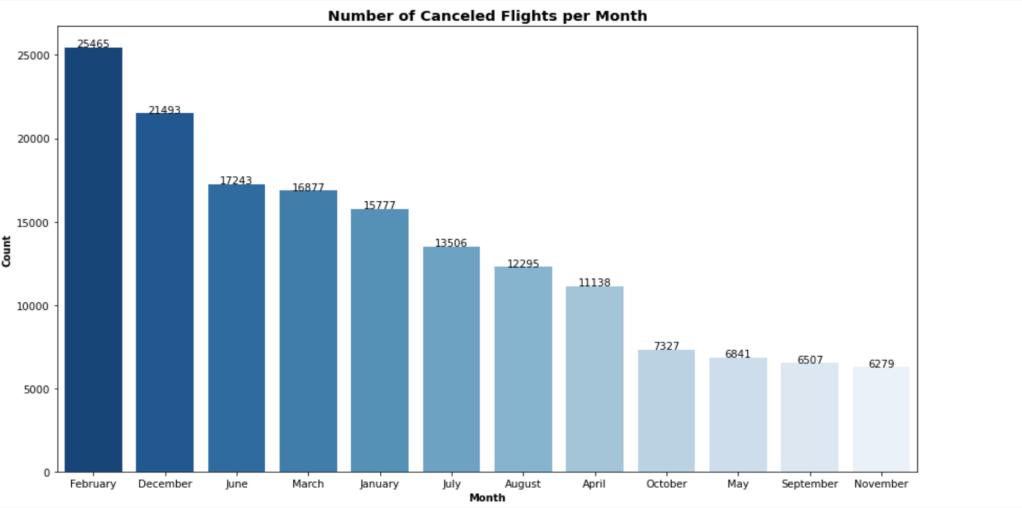
1994



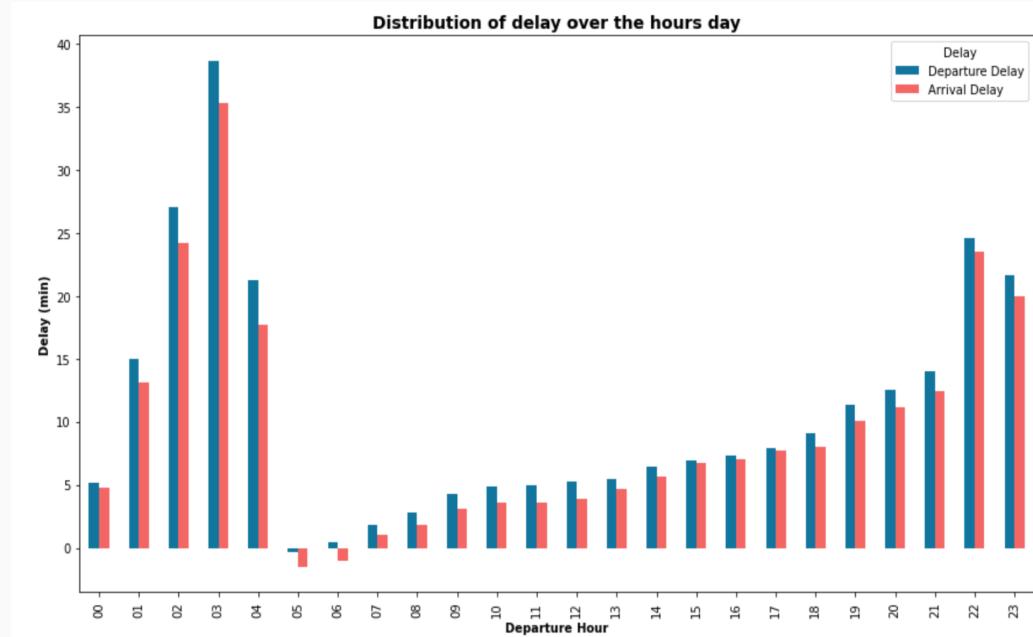
2006



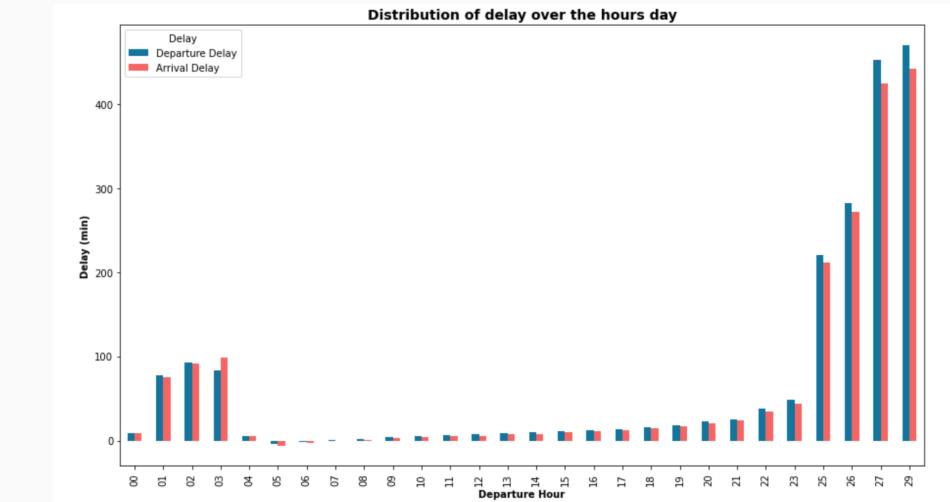
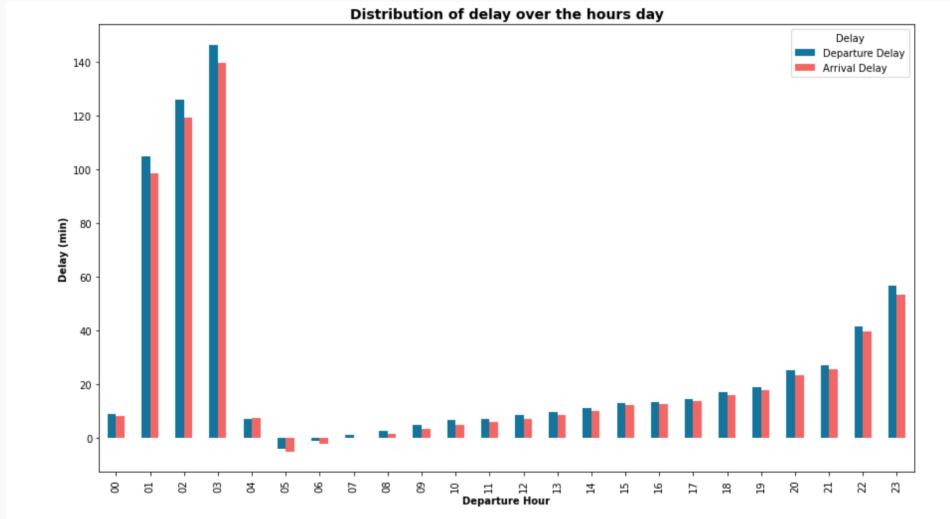
2007



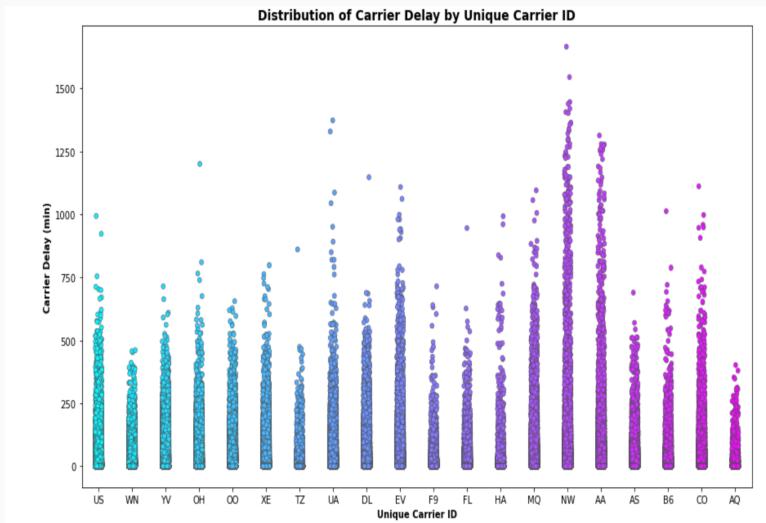
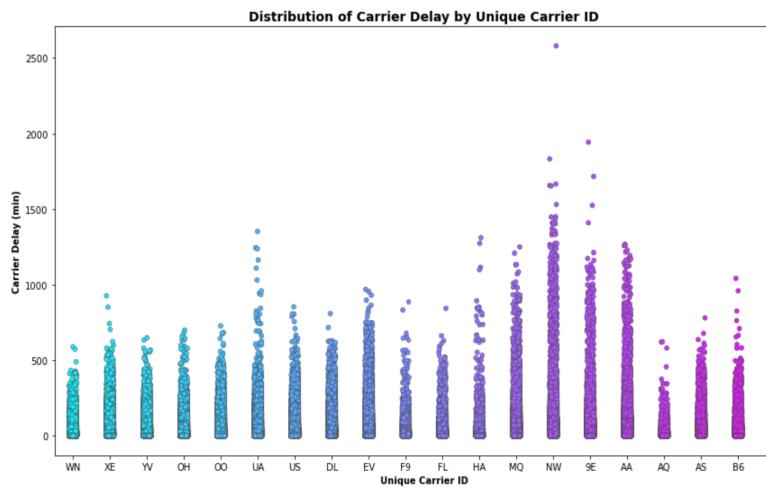
1994



2006 v 2007

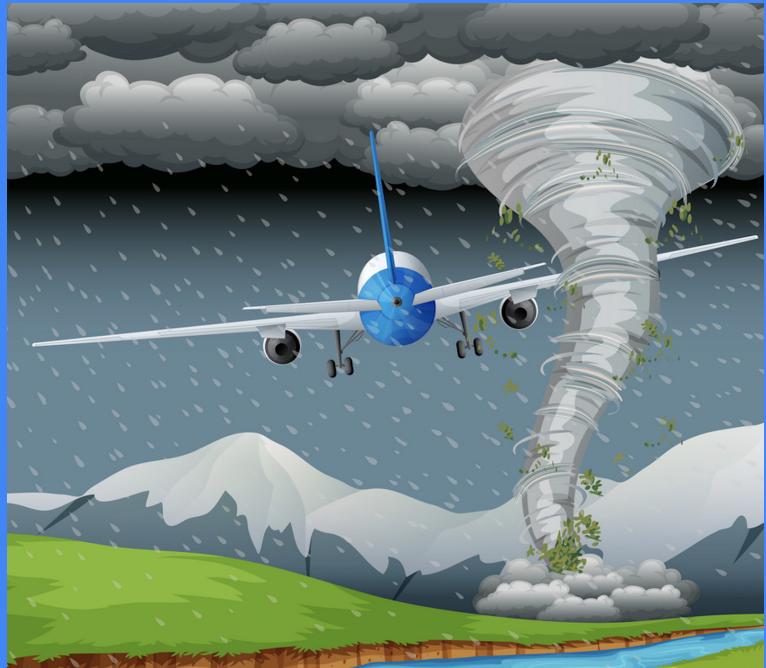


2006 v 2007



Conclusion

- Climate And Carrier



Thanks!

