# Machine Learning in the Clouds

Jeremy Beard

MSDS 696

Capstone

December 2023

> **"TRAVEL IS NOT REWARD FOR WORKING, IT'S EDUCATION FOR LIVING "**

Anthony Bourdain

# AGENDA

- Problem Statement
- Methods and Operations
- Notable EDA Plots
- ML Overview
- ML Results
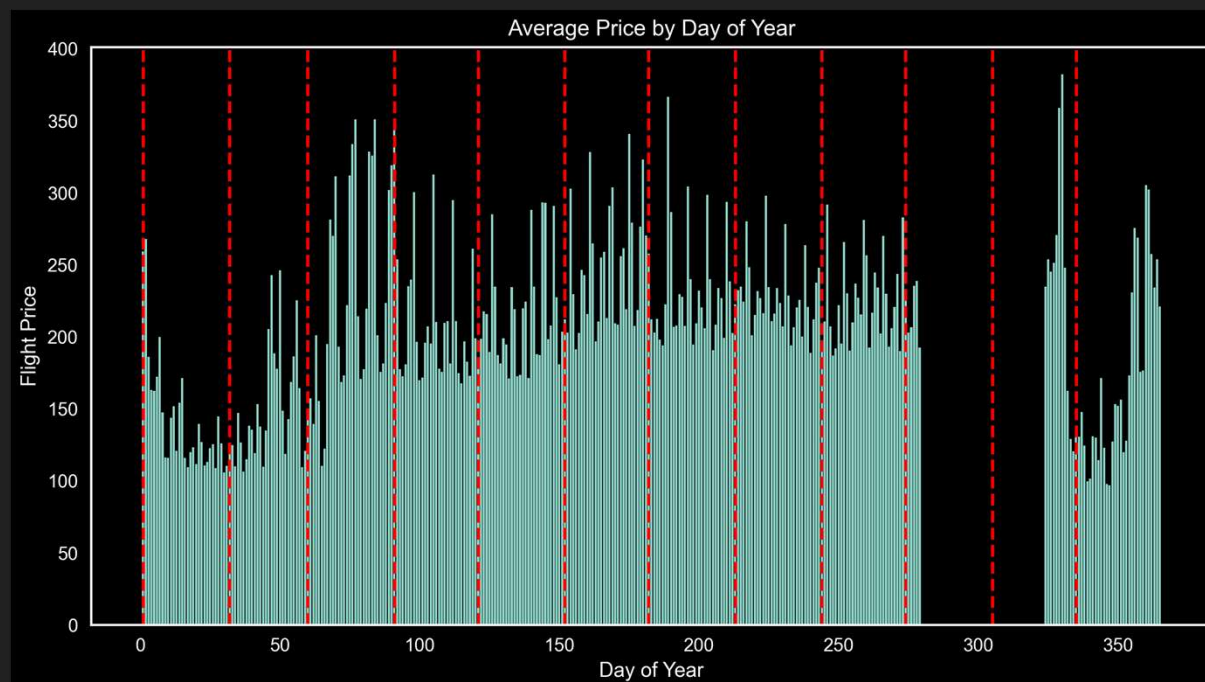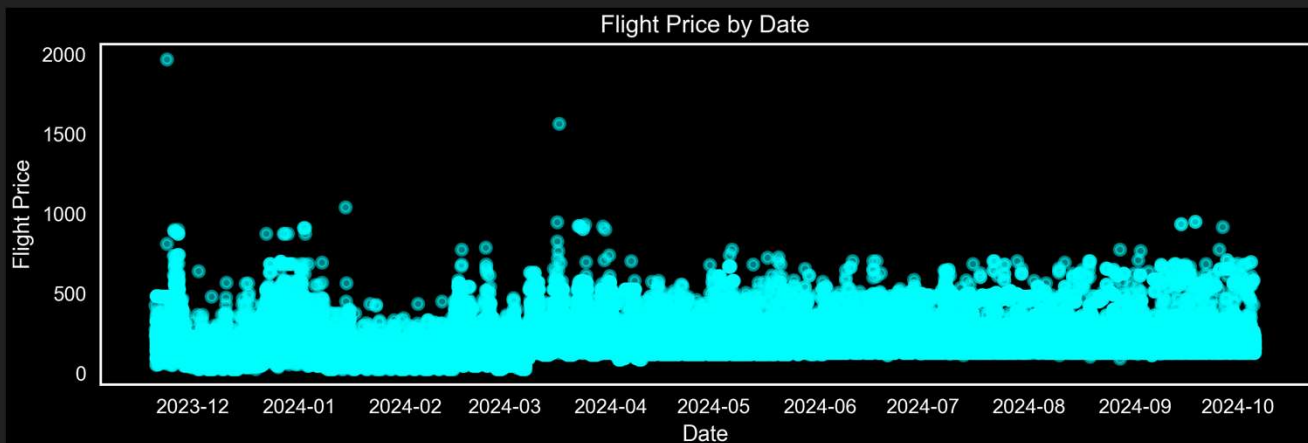- Future Work
- Conclusions
- References

# Problem Statement

"How can we reliably predict trends and prices for air travel based on common real-life inputs to leverage the same ML tactics corporations are utilizing, but on consumer-side?"
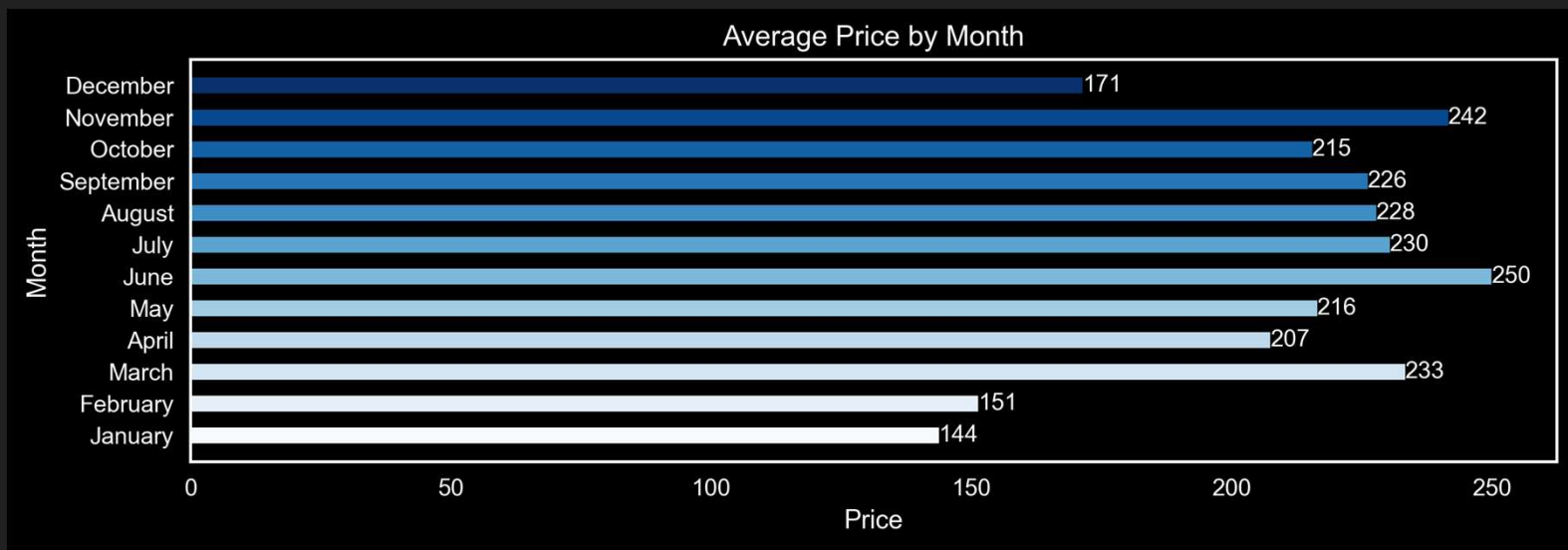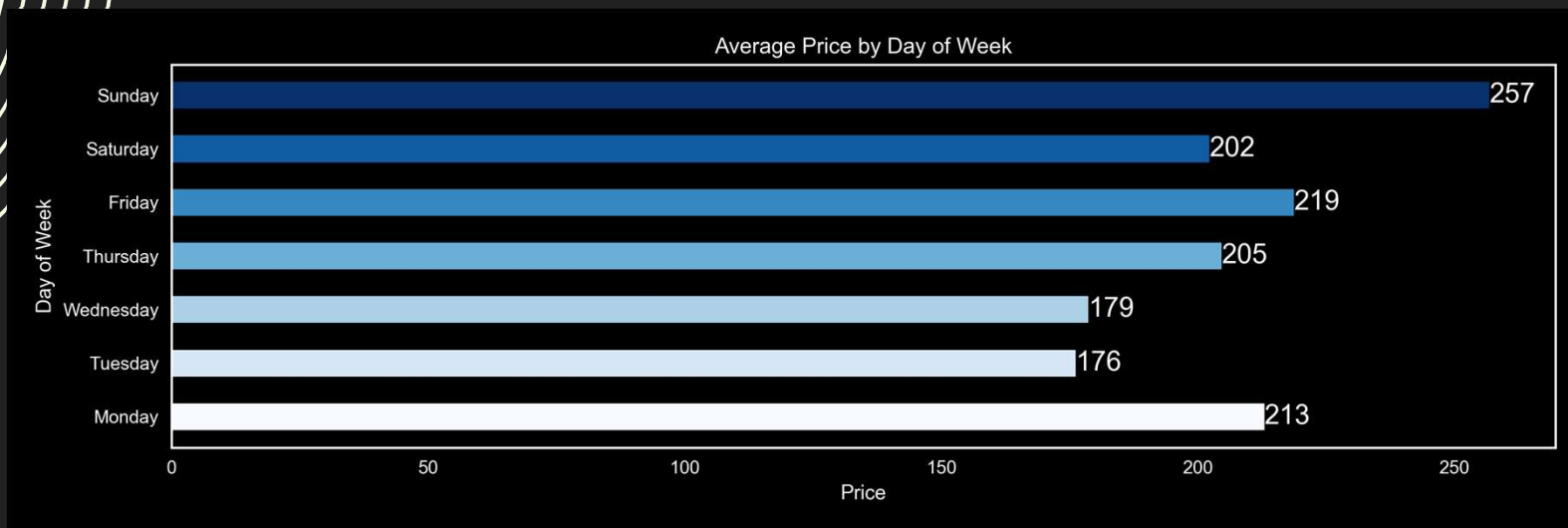
- The challenge is to implement a web-scraping solution that extracts relevant data, including departure details, flight duration, layovers, class types, and current/historical pricing information, from a relevant online source (or sources).

- Solution will enable the creation of a precise predictive model and empower stakeholders with valuable insights for informed decision-making in the airline industry.
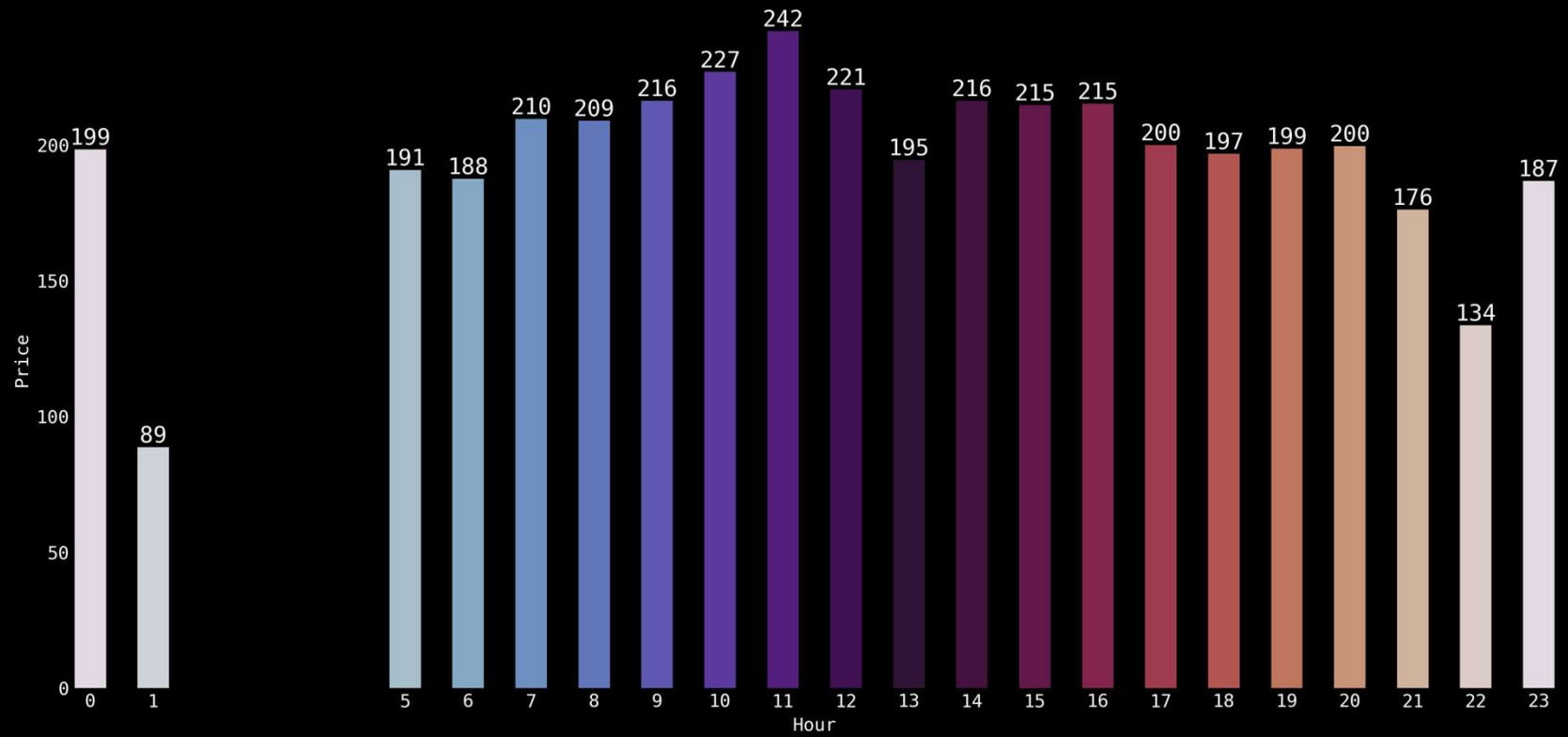
4

# Methods, Operations & Processes

- Web-scraping (via API)

- EDA

- Tableau Dashboarding

- Model Creation

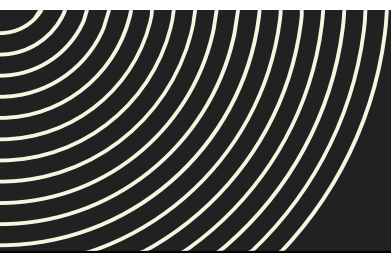- Model Comparison

- Final Model Selection

- Summary

# Notable EDA

**Flight Price by Date**

**Average Price by Day of Year**

Average Price by Day of Week

| Day of Week | Price |
|-------------|-------|
| Sunday | 257 |
| Saturday | 202 |
| Friday | 219 |
| Thursday | 205 |
| Wednesday | 179 |
| Tuesday | 176 |
| Monday | 213 |

Average Price by Month

| Month | Price |
|-------|-------|
| December | 171 |
| November | 242 |
| October | 215 |
| September | 226 |
| August | 228 |
| July | 230 |
| June | 250 |
| May | 216 |
| April | 207 |
| March | 233 |
| February | 151 |
| January | 144 |

Average Price by Hour of Day

# Tableau

FLIGHTS DASHBOARD

Frontier $130.1 | United $205.1 | American $207.1 | Spirit $148.4

Created by Jeremy Beard - MSDS 696 - Regis University

Earliest Flight Data 11/20/2023
Latest Flight Data 10/5/2024
Total Flights 44,367

# Machine Learning Overview

# ML Overview

- Utilized **pycaret** Python package for rapidly iterating on multiple ML models

- Compared multiple dimensions of model parameters
  - Focus: timeseries handling
  - Compared 3-5 different methods of handling datetimes

- Created **1** ML model baseline and **8** pycaret iterations through 8-10 ML models each

# ML Model Results

# ML Model Results

- **Mean Absolute Error:** Average of the absolute differences between the predicted and actual values



Average Error per Model

| Model | Mean Absolute Error |
|---|---|
| Ordinal Date, KFolds gen | 22.05 |
| Parsed Date, KFolds gen | 22.23 |
| Ordinal Date, Timeseries gen | 25.64 |
| Parsed Date, Timeseries gen | 25.96 |
| mkTime Date, KFolds gen | 27.26 |
| mkTime Date, Timeseries gen | 30.07 |
| Sklearn Linear Regression | 54.68 |
| fullDateTime, Timeseries gen | 55.73 |
| dateTime, Timeseries gen | 55.83 |

# ML Model Results

- Overall good performance
  - MAE: $22.05
  - $R^2$: 0.83
- Pycaret timeseries generator did not perform as well as expected
- Random Forest and Extra Trees

| Model # | Date Handling | Pycaret Generator | Top Model | MAE |
|---------|---------------|-------------------|-----------|-----|
| 1 | Parsed Datetime | N/A | Linear Regression | 54.68 |
| 2 | Parsed Datetime | KFold | Extra Trees | 22.23 |
| 3 | Ordinate + Hour | KFold | Random Forest | 22.05 |
| 4 | mktime | KFold | Extra Trees | 27.26 |
| 5 | deptDatetime | Timeseries | Random Forest | 55.73 |
| 6 | deptDate + Hour | Timeseries | LightGBM | 55.83 |
| 7 | Parsed Datetime | Timeseries | Extra Trees Regressor | 25.96 |
| 8 | Parsed Datetime | Timeseries | Extra Trees Regressor | 25.64 |
| 9 | Parsed Datetime | Timeseries | Random Forest | 30.07 |

| | Model | MAE | R2 |
|---|-------|-----|-----|
| rf | Random Forest Regressor | 22.0532 | 0.8309 |

# Future Work

- More data!
  - More airports, more airlines, continually update the sqlite3 database
  - 5-sec latency API queries made 45,000 datapoints take all weekend
  - Research other API's
  - Implement selenium?
- Implement more feature engineering
- Experiment with more Tableau dashboarding

# Conclusions

- Using the Booking.com API, we were able to create a model with $22 Mean Absolute Error prediction

- Data is biased according to when it was scraped

- Best Models: Random Forest, Extra Trees

  - Using Ordinal datetime

- Timeseries still a question in my mind, it should perform better

# References

- https://github.com/jeremyabeard5/MSDS696
- https://booking-com13.p.rapidapi.com/flights/one-way
- https://stackoverflow.com/questions/40217369/python-linear-regression-predict-by-date
- https://matplotlib.org/stable/gallery/style_sheets/style_sheets_reference.html
- https://pandas.pydata.org/docs/reference/api/pandas.DatetimeIndex.dayofweek.html
- https://pandas.pydata.org/docs/reference/api/pandas.Series.dt.month.html
- https://stackoverflow.com/questions/71419004/how-to-plot-vertical-lines-at-specific-dates-in-matplotlib
- https://medium.com/analytics-vidhya/calendar-heatmaps-a-perfect-way-to-display-your-time-series-quantitative-data-ad36bf81a3ed
- https://calplot.readthedocs.io/en/latest/index.html
- https://matplotlib.org/stable/users/explain/colors/colormaps.html
- https://medium.com/analytics-vidhya/calendar-heatmaps-a-perfect-way-to-display-your-time-series-quantitative-data-ad36bf81a3ed
- https://pypi.org/project/july/
- https://datascience.stackexchange.com/questions/2368/machine-learning-features-engineering-from-date-time-data
- https://datascience.stackexchange.com/questions/112357/feature-engineering-for-datetime-column
- https://www.reddit.com/r/learnpython/comments/chunas/correlation_with_day_of_week/
- https://mikulskibartosz.name/time-in-machine-learning
- https://www.pycaret.org/tutorials/html/REG101.html
- https://pycaret.readthedocs.io/en/latest/api/regression.html

# THANK YOU

Jeremy Beard

jeremyab5@gmail.com

https://www.linkedin.com/in/jeremyab5