

CAPSTONE DATA SCIENCE PROJECT

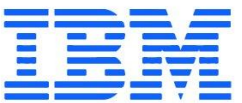
Michael Komer 06/2023



INTRODUCTION

My Name is Mike Komer, I will briefly share my background and the purpose of this report

- I currently work as an analyst and tableau data developer on the Enterprise Applications and Data Governance Team
- My goal was to complete this IBM Data Science Certification and to both further develop my toolset and methodologies, as well as show create a way to showcase them
- In this Capstone Course I will be sharing various course materials from prior sessions I was apart of to both showcase the newly added skills and tools from this certificate and on this report



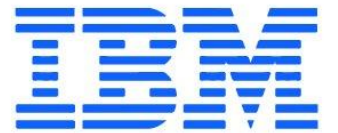


EXECUTIVE SUMMARY

**This Summary Report is for IBM Data Science Professional –
Certificate – Data Science Capstone Course**

Some of the things covered included:

- Data Wrangling, best practices and tools
- Exploratory Data Analysis SQL
- Data Visualization, various tools and formats
- Python Data Analysis
- Machine Learning Models and Visualizations using various Python libraries: Pandas, Seaborn, Scikit Learn, Matplotlib
- Predictive Analysis methodologies and tools available





DATA WRANGLING

Data Collection and Data Wrangling Process, Key Take Aways:

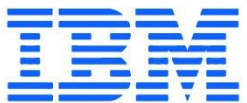
[Data Wrangling Project Github](#)

Data Wrangling is taking raw data and transforming it into a useable format.

Performed Exploratory Analysis on SpaceX Launch Data to find data patterns, anomalies and determine the label for training supervised models

Imported Data Sets, Prepped Data large datasets with Python in Jupyter Notebook

Imported Various Libraries Python Libraries to build visualizations during the exploratory phase; seaborn, Scikit Learn, matplotlib, pandas



INTERACTIVE VISUAL ANALYTICS

Interactive Visualizations:

Interactive Visual Analytics

Some of my favorite visuals made during certificate program displayed on the right.

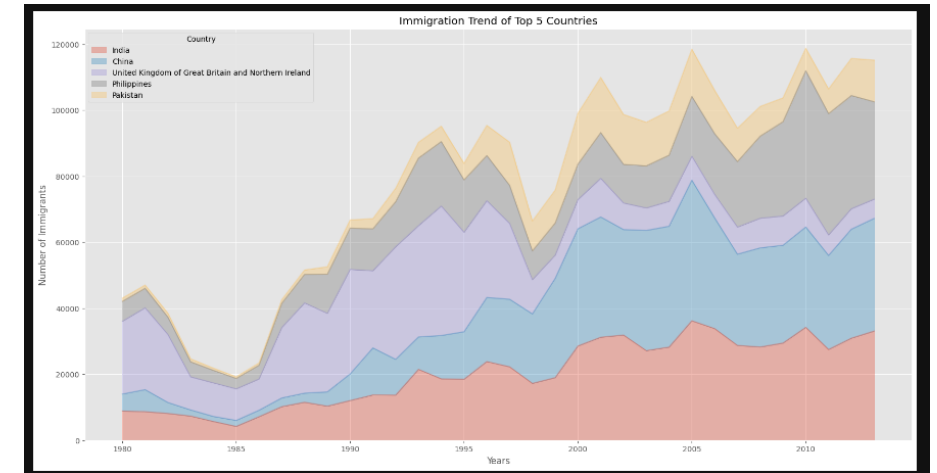
Performed Exploratory Data Analysis on SpaceX Launch Data to find data patterns, anomalies and determine the label for training supervised models

Explored different visualization types and use cases with SpaceX dataset

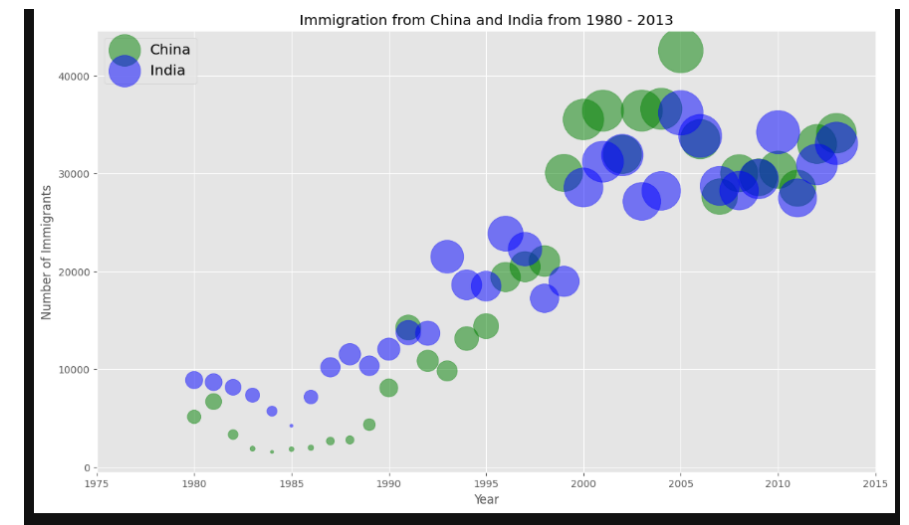
Began Classification Modeling and Testing strategy

Imported Various Libraries Python Libraries: seaborn, Scikit Learn, matplotlib, pandas

Stacked Area Plot Made During Certificate Program



Bubble Plot Time Series Data

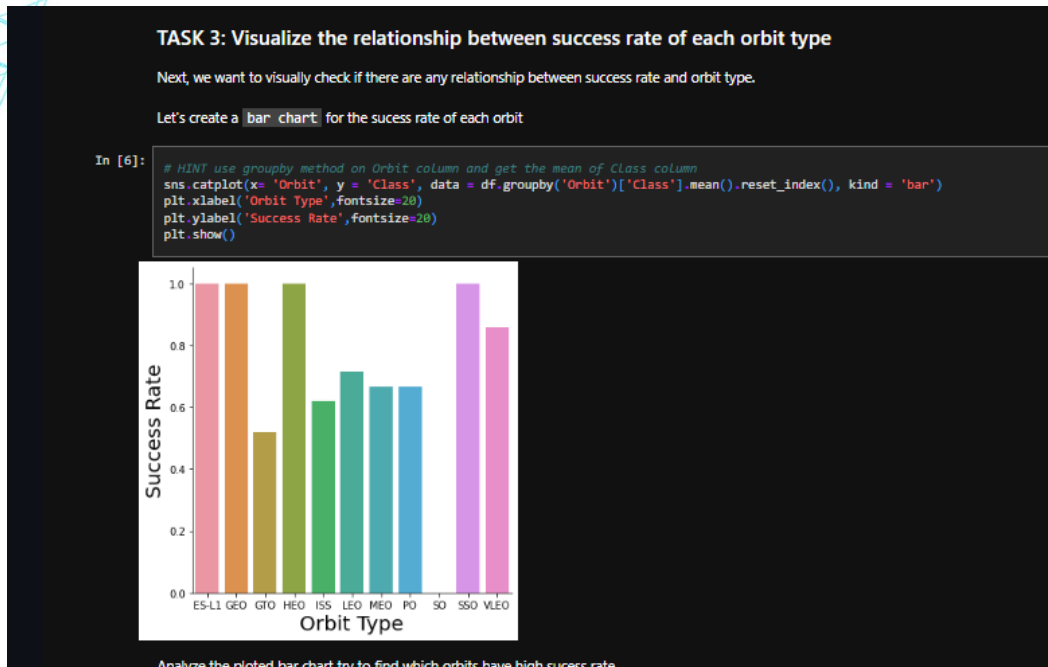


PREDICTIVE ANALYSIS

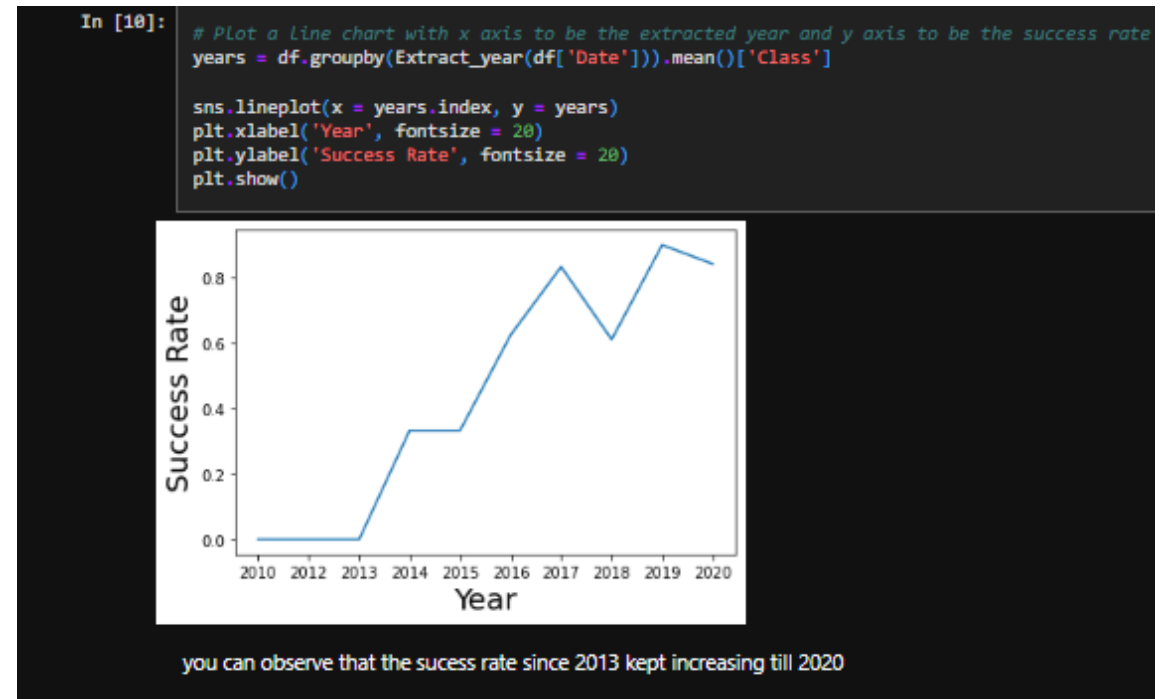
See the attached “Exploratory Data Analysis..” github project work
[Exploratory Data Analysis and Advanced Data Visualization](#)

Real Case Study Analysis on SpaceX data demonstrating how different data viz methods can be applied towards predictive analysis below:

LEFT: Orbit Type Relationship Success



RIGHT: Launch Success Trends





EDA WITH VISUALIZATION

See the attached "Predictive Analysis and Machine Learning Lab github project work

[Predictive Analysis and Machine Learning Lab](#)

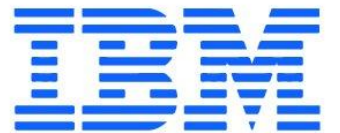
Exploratory Data Analysis with Viz lab seen above.

Real Case Study Analysis on SpaceX data and Machine Learning Prediction Models

Standardized the dataset, split data into groups, and began various data modeling tests and visual analysis on dataset

Tested various Machine Learning Models against dataset and applied different visualization tools and analysis strategies

Determined the machine learning method best fit for current data across Logistic Regression, Classification Trees, and Support Vector Management Models

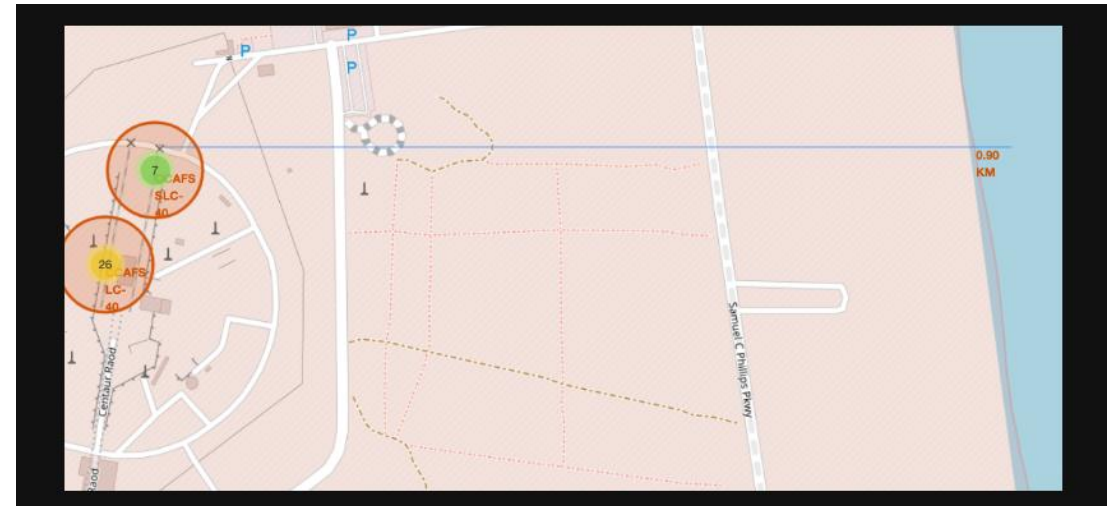
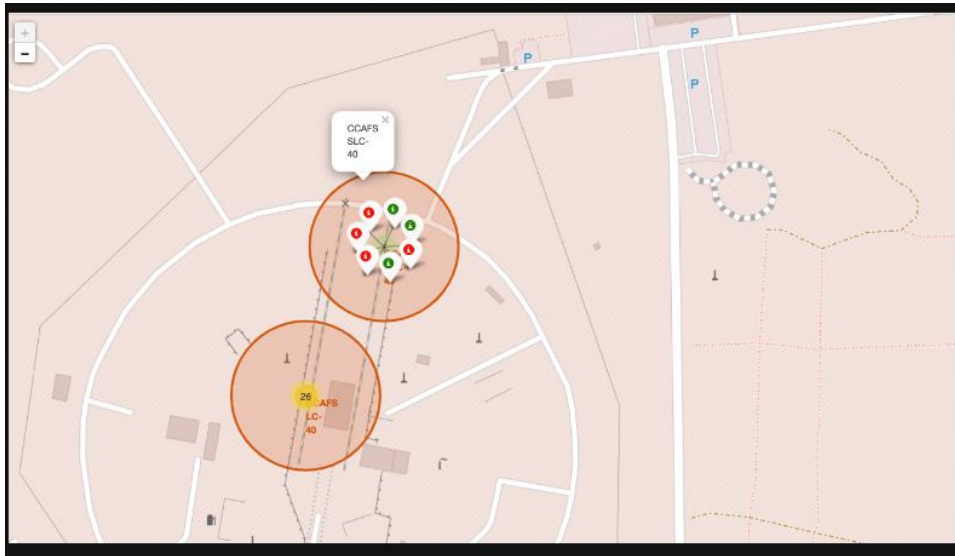


FOLIUM INTERACTIVE MAP

See below github project link

[Interactive Visual Analytics](#)

Focused on the Map Data Launch Analysis, added elements to the maps within Python that displayed in depth details of launch sites, activity levels, trends in graph visualization marks



PLOTLY VIZ DASHBOARD RESULTS

Plotly Dashboard:

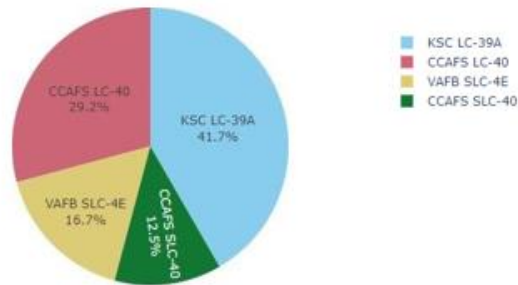
[Plotly Viz Link](#)

Continuation of SpaceX Case Study,

Dashboard is made up of pie chart and scatter plots

Pie chart displays distribution of successful landing and to visualize launch site success rates where the Scatter plot helps show the varying factors to success across launch sites

Successful Launches Across Launch Sites



Payload Mass vs. Success vs. Booster Version Category



PREDICTIVE ANALYSIS CLASSIFICATION

Predictive Analysis(Classification):

Predictive Analysis Classification

Continuation of SpaceX Case Study

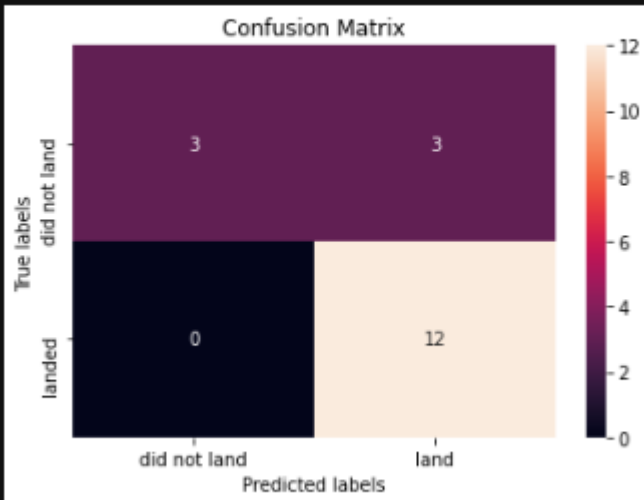
Best Performing Methodology:

Out[31]:	LogReg	SVM	Tree	KNN
Jaccard_Score	0.833333	0.845070	0.882353	0.819444
F1_Score	0.909091	0.916031	0.937500	0.900763
Accuracy	0.866667	0.877778	0.911111	0.855556

Tree Model
considered most
popular and
shown here yields
best accuracy

Confusion Matrix Summary

```
In [29]: knn_yhat = knn_cv.predict(X_test)
plot_confusion_matrix(Y_test, knn_yhat)
```



INNOVATIVE INSIGHTS(Confusion Matrix):

- Since all models performed the same for the test set, the confusion matrix is the same across all models.
- The models depict 12 successful landings when the true label was successful landing.
- The models predicted 3 unsuccessful landings when the true label was unsuccessful landing.
- The models predicted 3 successful landings when the true label was unsuccessful landings (false positives).
- The models appear to over predict successful landings, possibly in part due to sample size.

CONCLUSION

FINAL THOUGHTS(UPDATED POST COMPLETEION):

Overall satisfied with IBM Data Science Professional Certificate Program

Was able to test out new tools and experiment with different visualization models and strategies in a semi structured self paced environment

Expanded on my Jupyter Notebook, GitLab, Python, Python Libraries, Python Visualization, and Machine Learning fundamentals

Below are some of my badges acquired:

