



Winning Space Race with Data Science

Supawish Kanokpongsakorn
March 5th , 2022



Outline

- Executive Summary
- Introduction
- Methodology
- Results
- Conclusion
- Appendix

Executive Summary

- Summary of methodologies
 - Data Collection through API
 - Data Collection with Web Scraping
 - Data Wrangling
 - Exploratory Data Analysis with SQL
 - Exploratory Data Analysis with Data Visualization
 - Interactive Visual Analytics with Folium
 - Machine Learning Prediction
- Summary of all results
 - Exploratory Data Analysis result
 - Interactive analytics in screenshots
 - Predictive Analytics result

Introduction

- Project background and context

Space X advertises Falcon 9 rocket launches on its website with a cost of 62 million dollars; other providers cost upward of 165 million dollars each, much of the savings is because Space X can reuse the first stage. Therefore, if we can determine if the first stage will land, we can determine the cost of a launch. This information can be used if an alternate company wants to bid against space X for a rocket launch. This goal of the project is to create a machine learning pipeline to predict if the first stage will land successfully.

- Problems you want to find answers

- What factors determine if the rocket will land successfully?
- The interaction amongst various features that determine the success rate of a successful landing.
- What operating conditions needs to be in place to ensure a successful landing program.

Section 1

Methodology

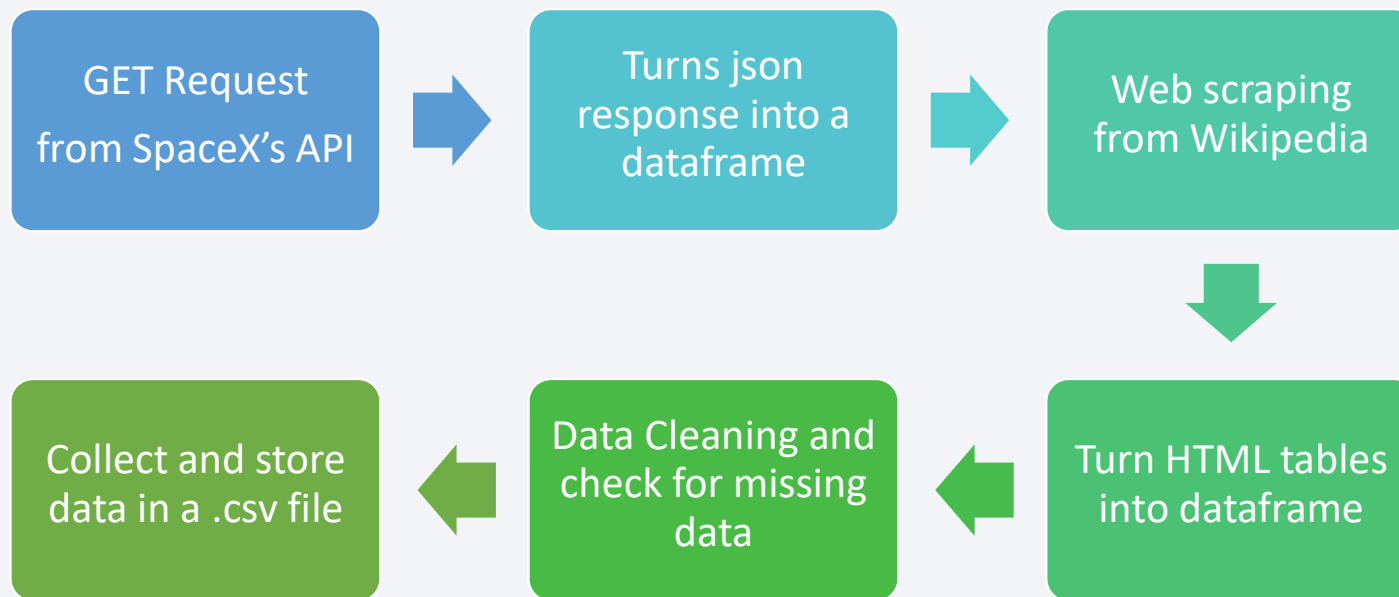
Methodology

Executive Summary

- Data collection methodology:
 - Data was collected using SpaceX API and web scraping from Wikipedia.
- Perform data wrangling
 - One-hot encoding was applied to categorical features
- Perform exploratory data analysis (EDA) using visualization and SQL
- Perform interactive visual analytics using Folium and Plotly Dash
- Perform predictive analysis using classification models
 - Build, tune, evaluate multiple classification models and use the one with highest accuracy score

Data Collection

- Describe how data sets were collected.



Data Collection – SpaceX API

- Present your data collection with SpaceX REST calls using key phrases and flowcharts
- [datasci-capstone-proj/Data Collection.ipynb at master · playerB/datasci-capstone-proj \(github.com\)](#)

Make a GET request to SpaceX's API and receive a response



Convert response to json



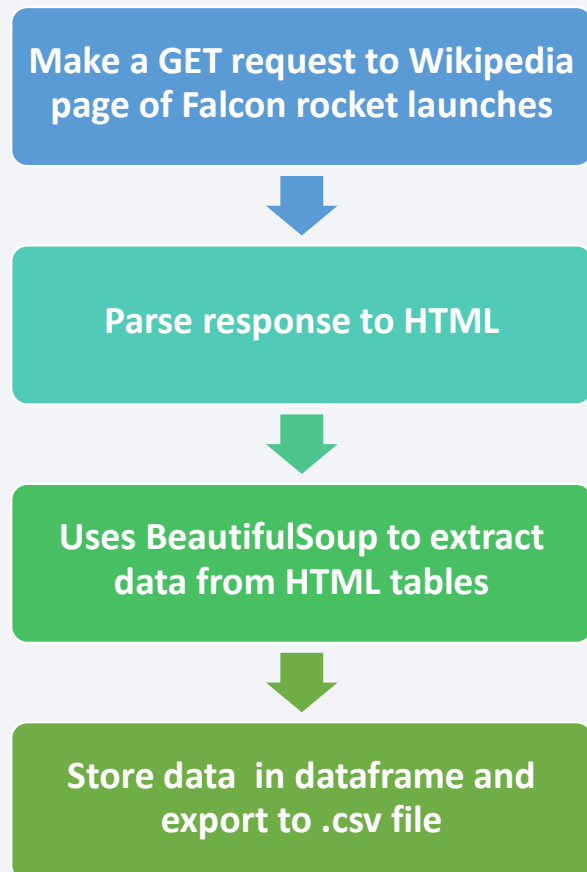
Convert json to Pandas dataframe



Imputing missing values in PayloadMass column

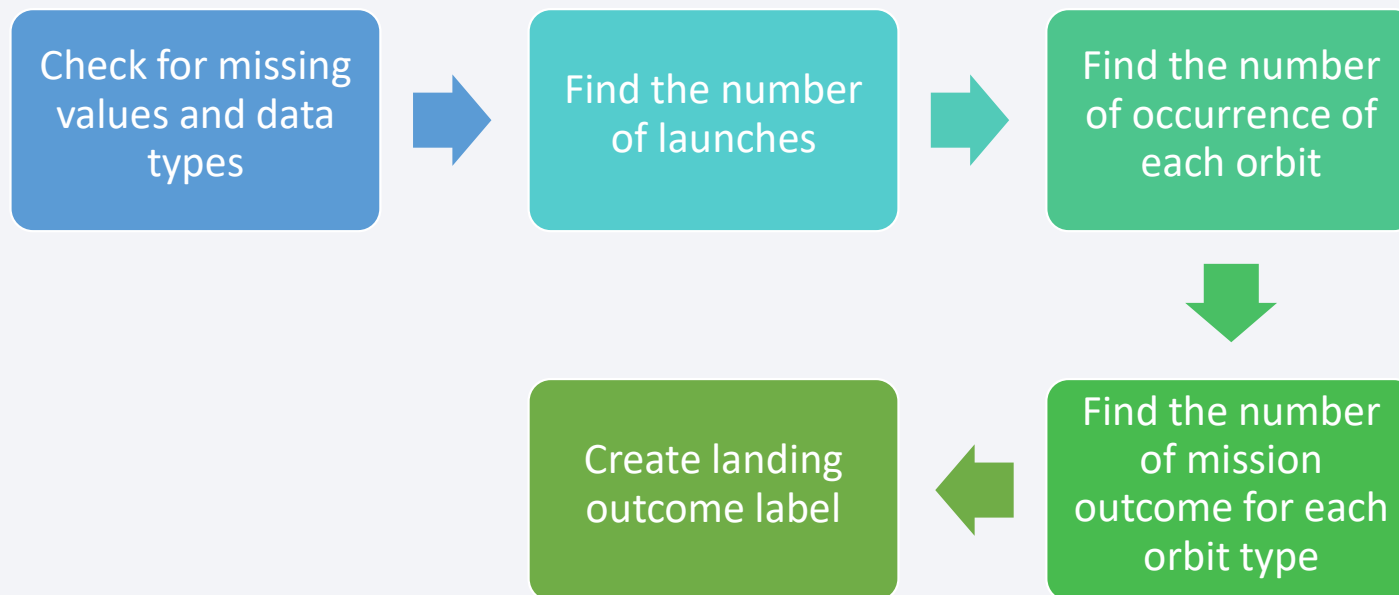
Data Collection - Scraping

- Present your web scraping process using key phrases and flowcharts
- [datasci-capstone-proj/Data Collection Scraping.ipynb at master · playerB/datasci-capstone-proj \(github.com\)](https://github.com/playerB/datasci-capstone-proj)

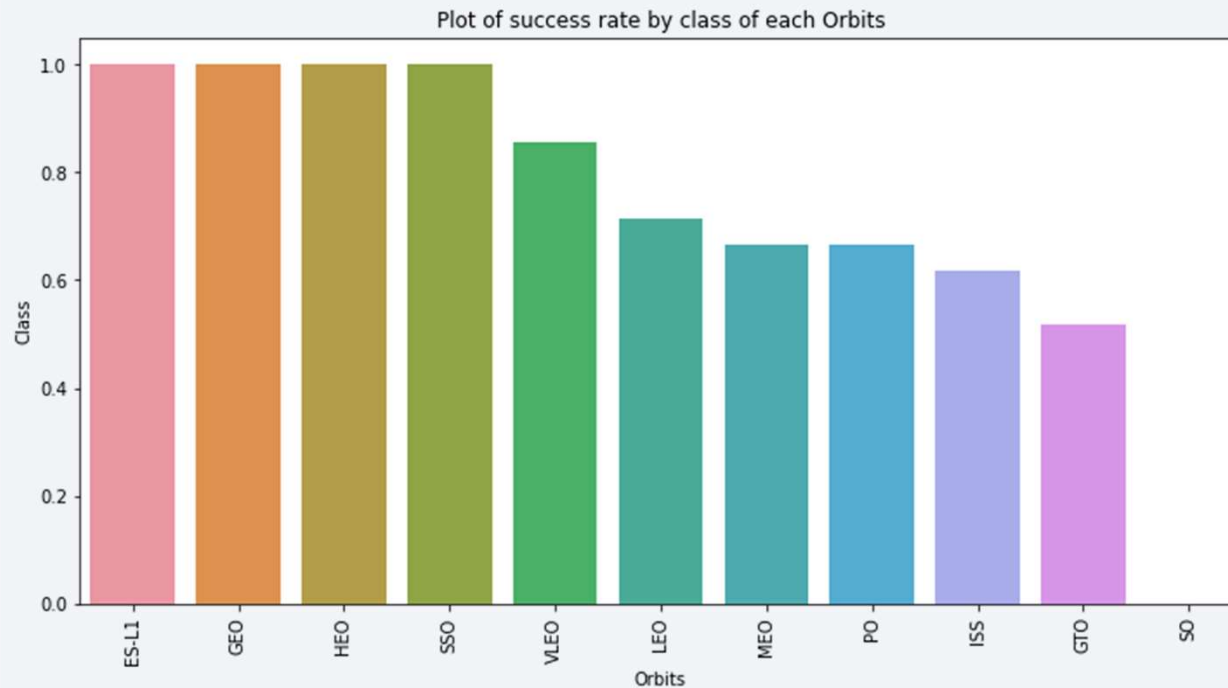


Data Wrangling

- Describe how data were processed
- [datasci-capstone-proj/Data Wrangling.ipynb at master · playerB/datasci-capstone-proj \(github.com\)](#)

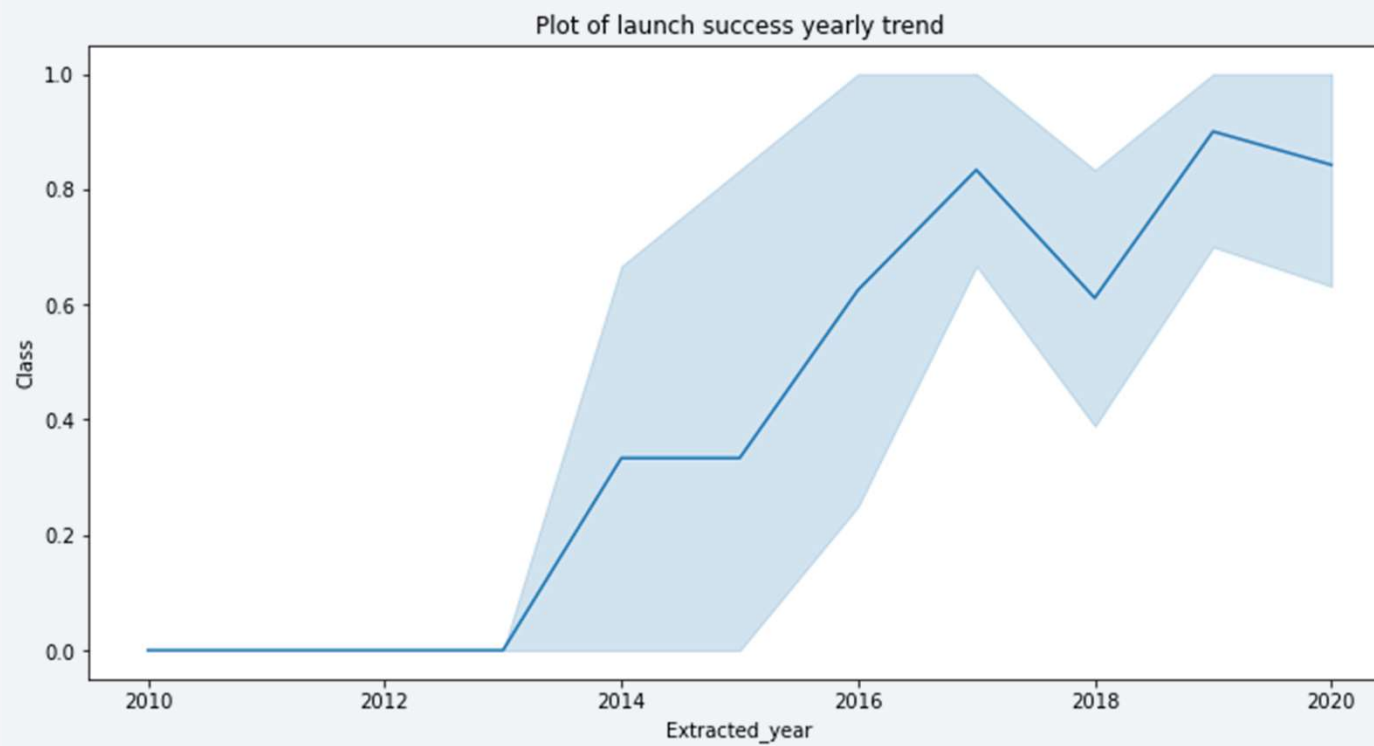


EDA with Data Visualization



- We explored the data by visualizing the relationship between flight number and launch Site, payload and launch site, success rate of each orbit type, flight number and orbit type, the launch success yearly trend.

EDA with Data Visualization



EDA with SQL

- We applied EDA with SQL to get insight from the data. We wrote queries to find out about interesting information such as:
 - The names of unique launch sites in the space mission.
 - The total payload mass carried by boosters launched by NASA (CRS)
 - The average payload mass carried by booster version F9 v1.1
 - The total number of successful and failure mission outcomes
 - The failed landing outcomes in drone ship, their booster version and launch site names.
- [datasci-capstone-proj/EDA with SQL.ipynb at master · playerB/datasci-capstone-proj \(github.com\)](#)

Build an Interactive Map with Folium

- We marked all launch sites, and added map objects such as markers, circles, lines to mark the success or failure of launches for each site on the folium map.
- We assigned the feature launch outcomes (failure or success) to class 0 and 1.i.e., 0 for failure, and 1 for success.
- Using the color-labeled marker clusters, we identified which launch sites have relatively high success rate.
- We calculated the distances between a launch site to its proximities. We answered some question for instance:
 - Are launch sites near railways, highways and coastlines.
 - Do launch sites keep certain distance away from cities.
- [datasci-capstone-proj/Interactive Visual Analytics with Folium lab.ipynb at master · playerB/datasci-capstone-proj \(github.com\)](#)

Build a Dashboard with Plotly Dash

- We built an interactive dashboard with Plotly dash
- We plotted pie charts showing the total launches by a certain sites
- We plotted scatter graph showing the relationship with Outcome and Payload Mass (Kg) for the different booster version.
- [datasci-capstone-proj/dash_interactivity.py at master · playerB/datasci-capstone-proj \(github.com\)](https://github.com/playerB/datasci-capstone-proj)

Predictive Analysis (Classification)

- We loaded the data using numpy and pandas, transformed the data, split our data into training and testing.
- We built different machine learning models and tune different hyperparameters using GridSearchCV.
- We used accuracy as the metric for our model, improved the model using feature engineering and algorithm tuning.
- We found the best performing classification model.
- [datasci-capstone-proj/Machine Learning Prediction.ipynb at master · playerB/datasci-capstone-proj \(github.com\)](#)

Results

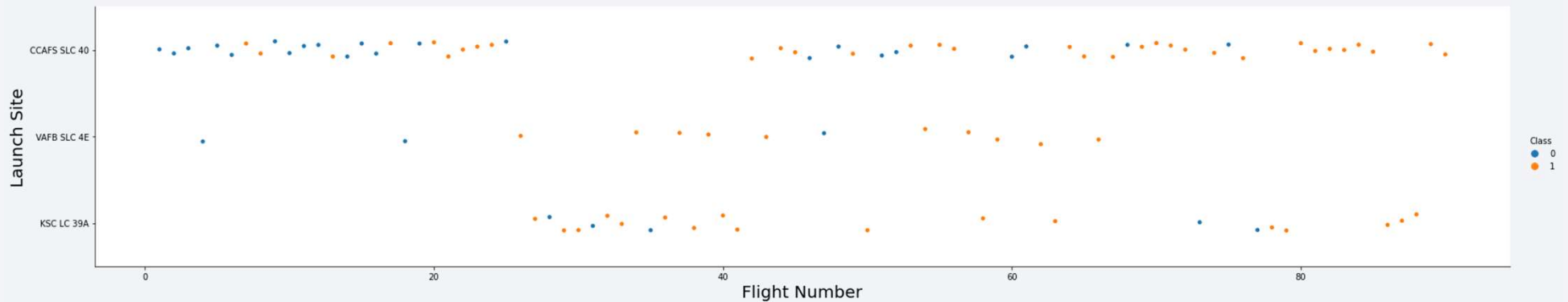
- Exploratory data analysis results
- Interactive analytics demo in screenshots
- Predictive analysis results



Section 2

Insights drawn from EDA

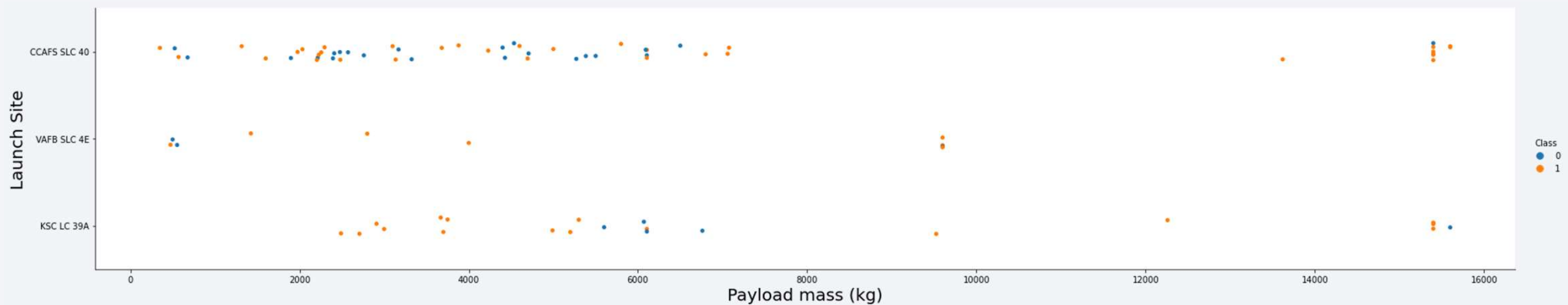
Flight Number vs. Launch Site



- From the plot, we found that the larger the flight amount at a launch site, the greater the success rate at the launch site.

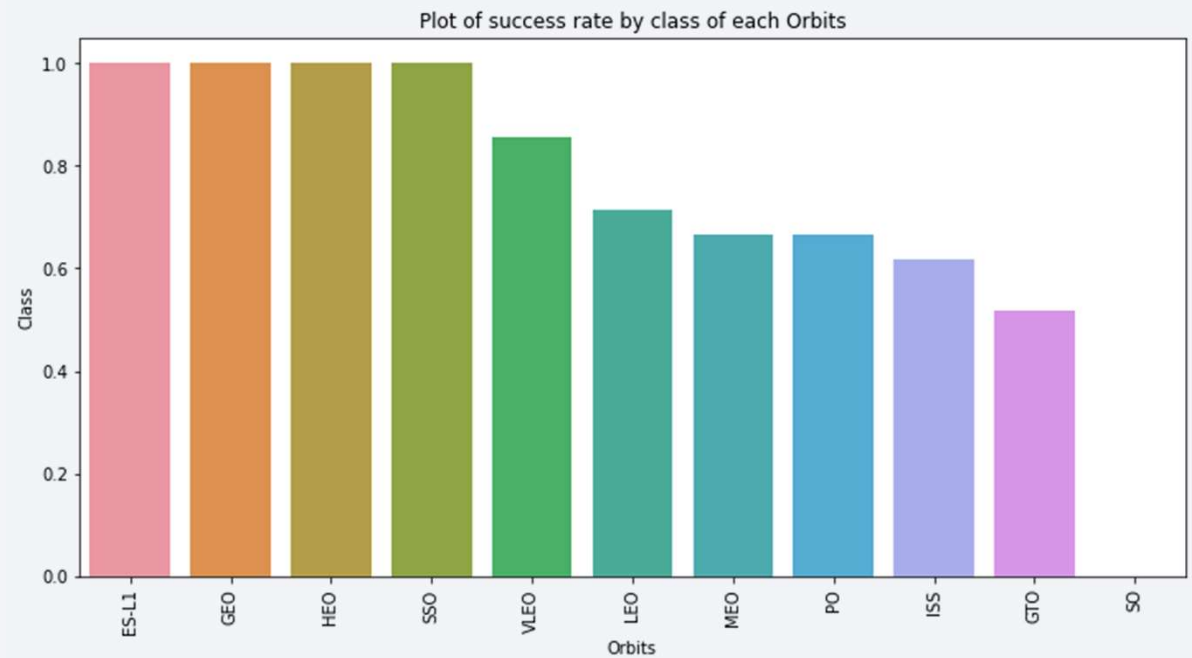
Payload vs. Launch Site

- From the plot, we found that the greater the payload mass is at a launch site, the higher the success rate at the launch site.

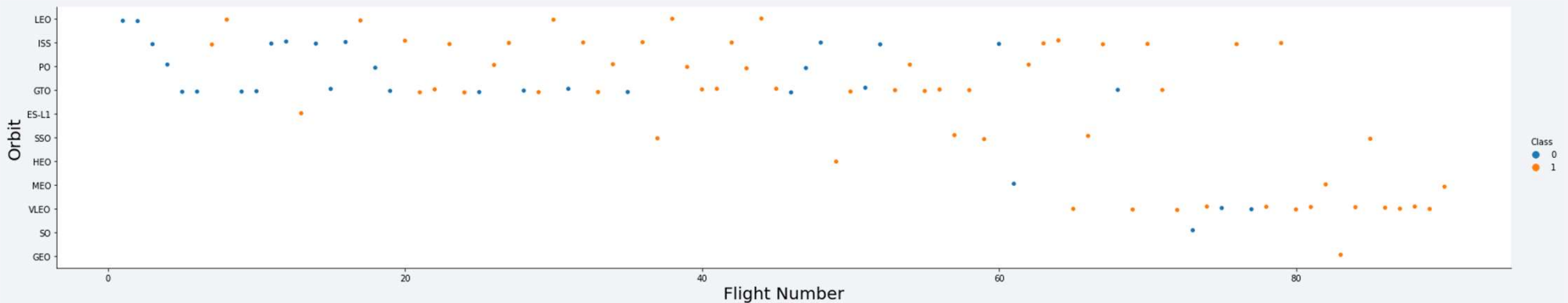


Success Rate vs. Orbit Type

- From the plot, we can see that ES-L1, GEO, HEO, SSO, VLEO had the most success rate of 100%.



Flight Number vs. Orbit Type



- The plot below shows the Flight Number vs. Orbit type. We observe that in the LEO orbit, success is related to the number of flights whereas in the GTO orbit, there is no relationship between flight number and the orbit.

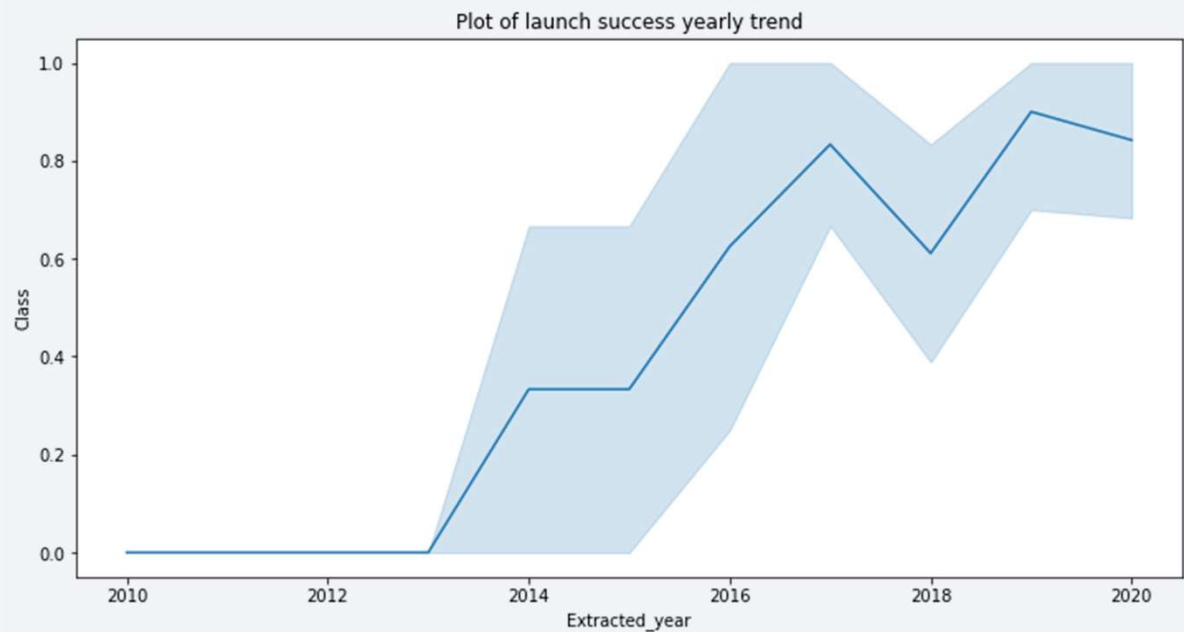
Payload vs. Orbit Type



- We can observe that with heavy payloads, the successful landing are more for PO, LEO and ISS orbits.

Launch Success Yearly Trend

- From the plot, we can observe that success rate since 2013 kept on increasing till 2020.



All Launch Site Names

Display the names of the unique launch sites in the space mission

```
%sql SELECT DISTINCT launch_site FROM SPACEXTBL;
```

```
* ibm_db_sa://xvs97022:***@ea286ace-86c7-4d5b-8580-3fbfa40
Done.
```

launch_site

CCAFS LC-40

CCAFS SLC-40

KSC LC-39A

VAFB SLC-4E

Launch Site Names Begin with 'CCA'

Display 5 records where launch sites begin with the string 'CCA'

```
%sql SELECT launch_site FROM SPACEXTBL WHERE launch_site LIKE 'CCA%' LIMIT 5;
```

```
* ibm_db_sa://xvs97022:***@ea286ace-86c7-4d5b-8580-3fbfa46b1c66.bs2io90l08kqb1c  
Done.
```

launch_site

CCAFS LC-40

CCAFS LC-40

CCAFS LC-40

CCAFS LC-40

CCAFS LC-40

Total Payload Mass

Display the total payload mass carried by boosters launched by NASA (CRS)

```
%sql SELECT SUM(payload_mass__kg_) FROM SPACEXTBL WHERE customer LIKE 'NASA (CRS)';
```

```
* ibm_db_sa://xvs97022:***@ea286ace-86c7-4d5b-8580-3fbfa46b1c66.bs2io90l08kqb1od8lc  
Done.
```

```
1
```

```
45596
```

Average Payload Mass by F9 v1.1

Display average payload mass carried by booster version F9 v1.1

```
%sql SELECT AVG(payload_mass__kg_) FROM SPACEXTBL WHERE booster_version LIKE '%F9 v1.1%';
```

```
* ibm_db_sa://xvs97022:***@ea286ace-86c7-4d5b-8580-3fbfa46b1c66.bs2io90l08kqb1od8lcg.databa  
Done.
```

```
1
```

```
2534
```

First Successful Ground Landing Date

List the date when the first successful landing outcome in ground pad was achieved.

Hint: Use min function

```
%sql SELECT MIN(DATE) FROM SPACEXTBL WHERE landing__outcome LIKE 'Success%';
```

```
* ibm_db_sa://xvs97022:***@ea286ace-86c7-4d5b-8580-3fbfa46b1c66.bs2io90l08kqb1  
Done.
```

1

2015-12-22

Successful Drone Ship Landing with Payload between 4000 and 6000

List the names of the boosters which have success in drone ship and have payload mass greater than 4000 but less than 6000

```
%sql SELECT booster_version FROM SPACEXTBL WHERE mission_outcome='Success' AND landing__outcome LIKE '%drone ship%' AND payload_mass__kg_ BETWEEN 4000 .
```

```
* ibm_db_sa://xvs97022:***@ea286ace-86c7-4d5b-8580-3fbfa46b1c66.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:31505/bludb
Done.
```

booster_version

F9 FT B1020

F9 FT B1022

F9 FT B1026

F9 FT B1021.2

F9 FT B1031.2

Total Number of Successful and Failure Mission Outcomes

List the total number of successful and failure mission outcomes

```
%sql SELECT mission_outcome, COUNT(*) FROM SPACEXTBL GROUP BY mission_outcome;
```

```
* ibm_db_sa://xvs97022:***@ea286ace-86c7-4d5b-8580-3fbfa46b1c66.bs2io90l08kqb1od
Done.
```

| mission_outcome | 2 |
|-----------------|---|
|-----------------|---|

| | |
|---------------------|---|
| Failure (in flight) | 1 |
|---------------------|---|

| | |
|---------|----|
| Success | 99 |
|---------|----|

| | |
|----------------------------------|---|
| Success (payload status unclear) | 1 |
|----------------------------------|---|

Boosters Carried Maximum Payload

List the names of the booster_versions which have carried the maximum payload mass. Use a subquery

```
%sql SELECT booster_version, payload_mass__kg_ FROM SPACEXTBL WHERE payload_mass__kg_=(SELECT MAX(payload_mass__kg_) FROM SPACEXTBL);
```

```
* ibm_db_sa://xvs97022:***@ea286ace-86c7-4d5b-8580-3fbfa46b1c66.bs2io90108kqb1od8lcg.databases.appdomain.cloud:31505/bludb  
Done.
```

| booster_version | payload_mass__kg_ |
|-----------------|-------------------|
| F9 B5 B1048.4 | 15600 |
| F9 B5 B1049.4 | 15600 |
| F9 B5 B1051.3 | 15600 |
| F9 B5 B1056.4 | 15600 |
| F9 B5 B1048.5 | 15600 |
| F9 B5 B1051.4 | 15600 |
| F9 B5 B1049.5 | 15600 |
| F9 B5 B1060.2 | 15600 |
| F9 B5 B1058.3 | 15600 |
| F9 B5 B1051.6 | 15600 |
| F9 B5 B1060.3 | 15600 |
| F9 B5 B1049.7 | 15600 |

2015 Launch Records

List the failed landing_outcomes in drone ship, their booster versions, and launch site names for in year 2015

```
%sql SELECT landing__outcome, booster_version, launch_site FROM SPACEXTBL WHERE YEAR(DATE)=2015 AND landing__outcome LIKE '%Fail%drone ship%';
```

```
* ibm_db_sa://xvs97022:***@ea286ace-86c7-4d5b-8580-3fbfa46b1c66.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:31505/bludb
Done.
```

| landing__outcome | booster_version | launch_site |
|----------------------|-----------------|-------------|
| Failure (drone ship) | F9 v1.1 B1012 | CCAFS LC-40 |
| Failure (drone ship) | F9 v1.1 B1015 | CCAFS LC-40 |

Rank Landing Outcomes Between 2010-06-04 and 2017-03-20

Rank the count of landing outcomes (such as Failure (drone ship) or Success (ground pad)) between the date 2010-06-04 and 2017-03-20, in descending order

```
%sql SELECT landing__outcome, COUNT(*) FROM SPACEXTBL WHERE (DATE BETWEEN '2010-06-04' AND '2017-03-20') GROUP BY landing__outcome ORDER BY 1 DESC;
```

```
* ibm_db_sa://xvs97022:***@ea286ace-86c7-4d5b-8580-3fbfa46b1c66.bs2io90l08kqb1od8l1cg.databases.appdomain.cloud:31505/bludb
Done.
```

| landing__outcome | 2 |
|------------------------|----|
| Uncontrolled (ocean) | 2 |
| Success (ground pad) | 3 |
| Success (drone ship) | 5 |
| Precluded (drone ship) | 1 |
| No attempt | 10 |
| Failure (parachute) | 2 |
| Failure (drone ship) | 5 |
| Controlled (ocean) | 3 |

A satellite view of Earth from space, showing the curvature of the planet and city lights at night. The image is a composite of a solid blue gradient on the left and a satellite photograph of Earth on the right. The Earth's surface is dark blue, with numerous bright yellow and orange lights representing city lights at night. The horizon line of the Earth is visible, separating the dark blue of the planet from the blackness of space.

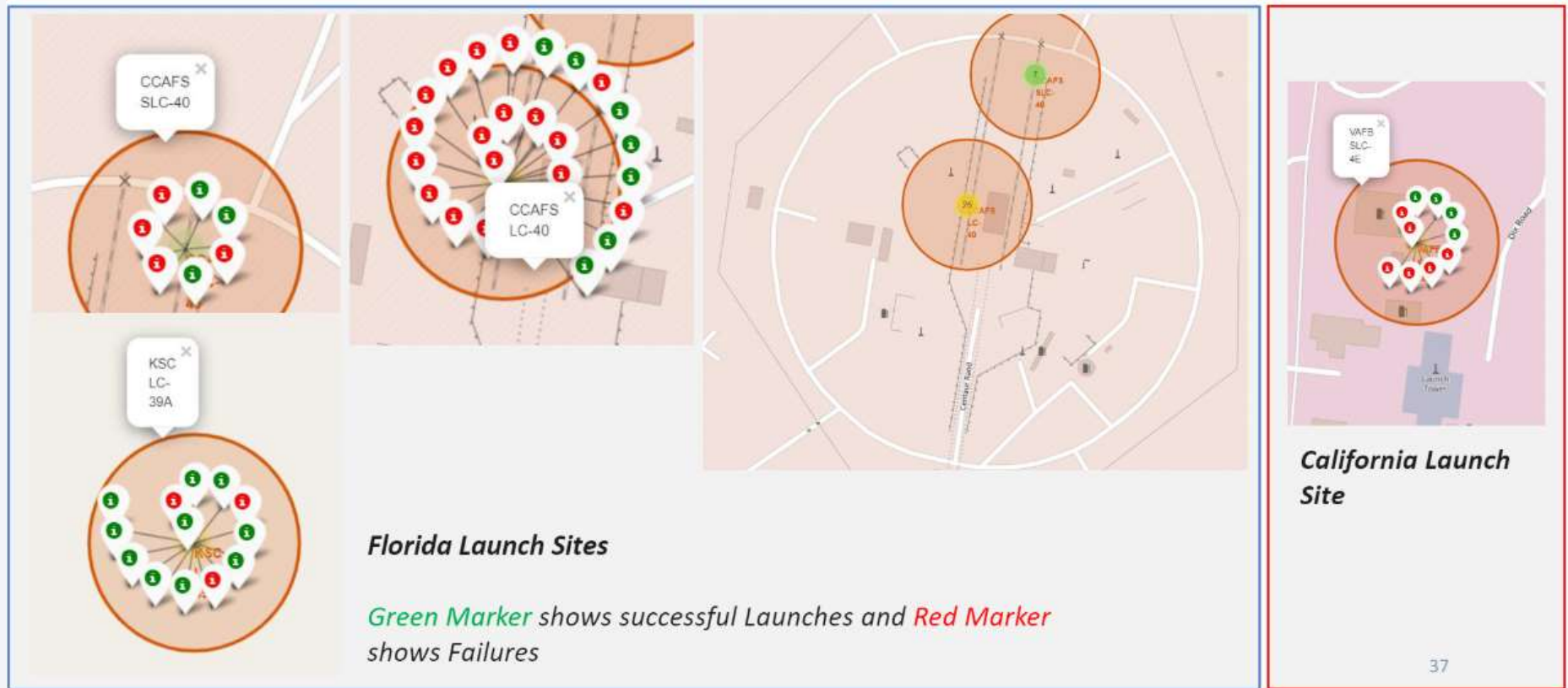
Section 3

Launch Sites Proximities Analysis

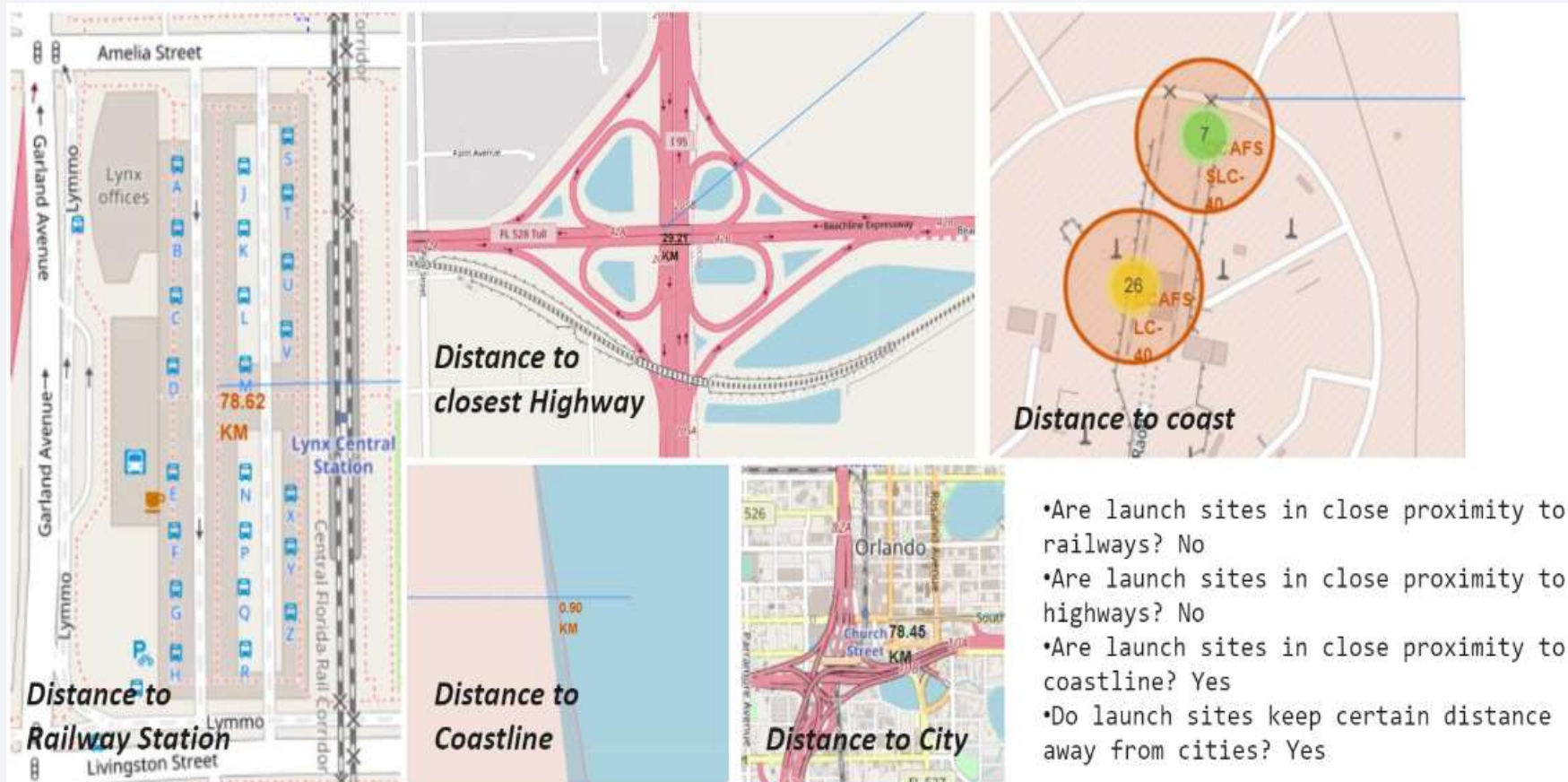
All SpaceX launch sites



Markers showing launch sites with color labels



Launch Site distance to landmarks



- Are launch sites in close proximity to railways? No
- Are launch sites in close proximity to highways? No
- Are launch sites in close proximity to coastline? Yes
- Do launch sites keep certain distance away from cities? Yes

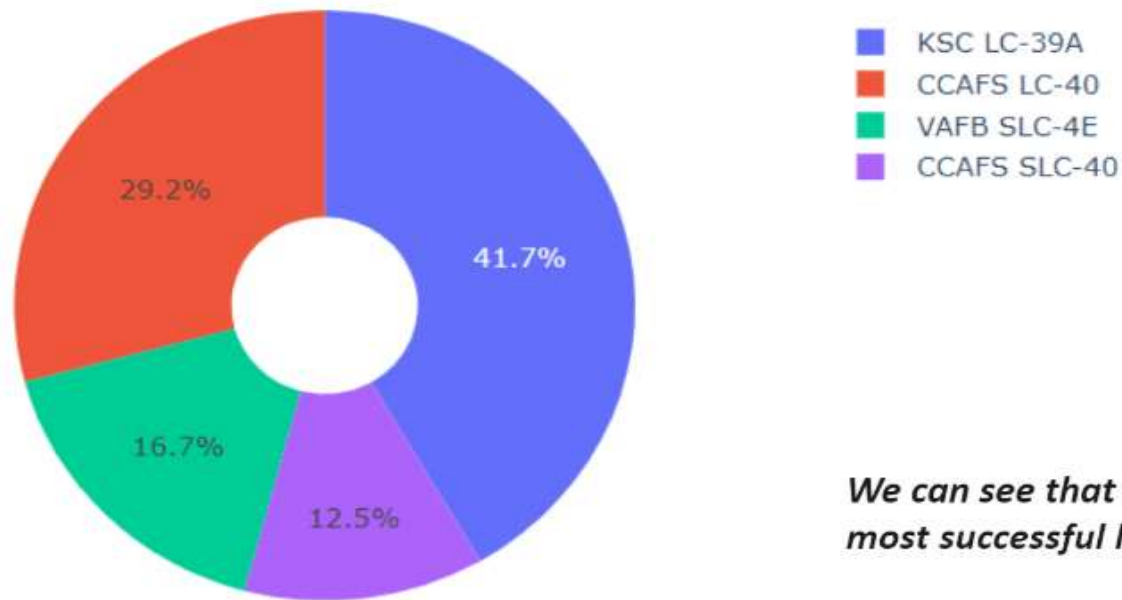


Section 4

Build a Dashboard with Plotly Dash

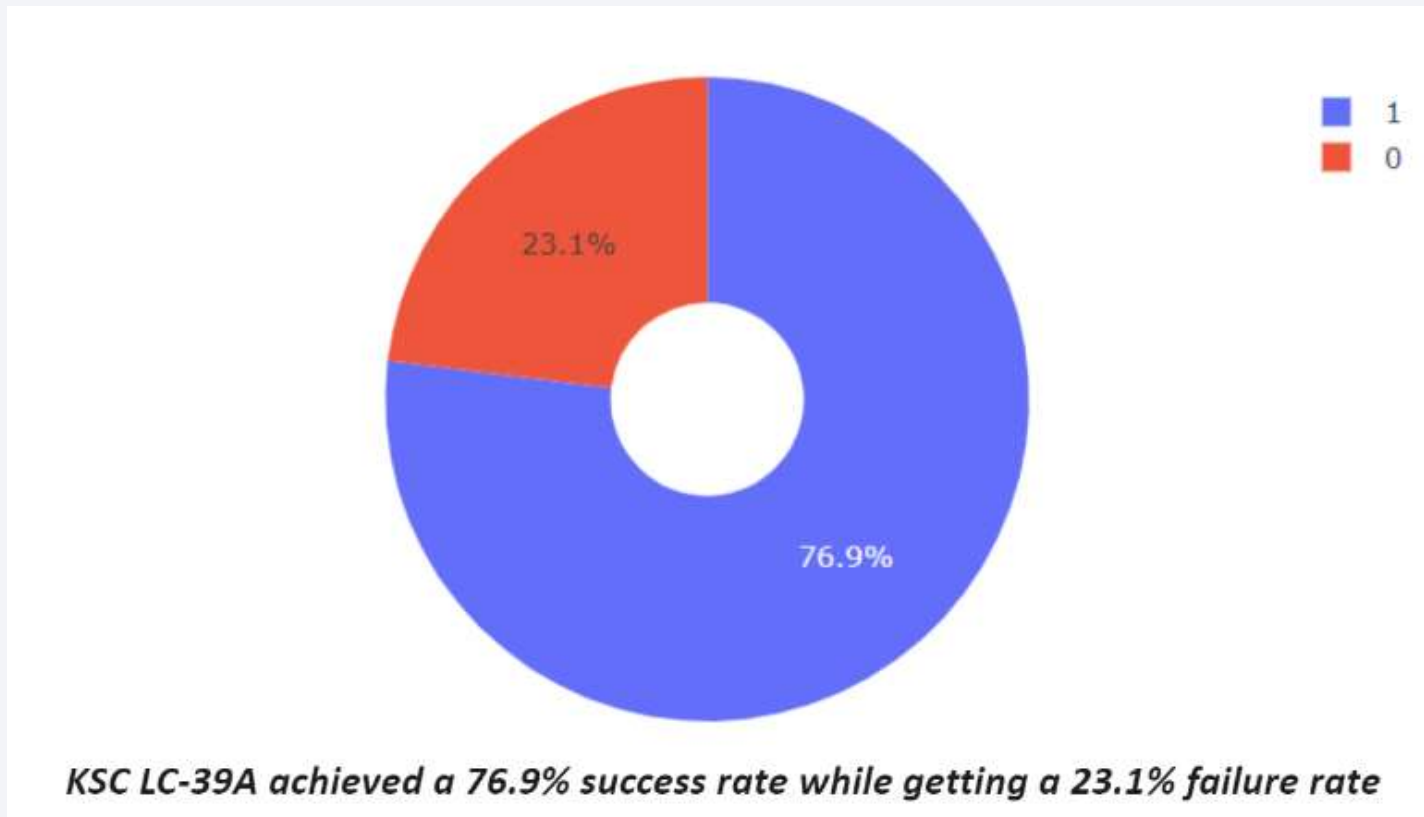
Pie chart showing the success percentage achieved by each launch site

Total Success Launches By all sites

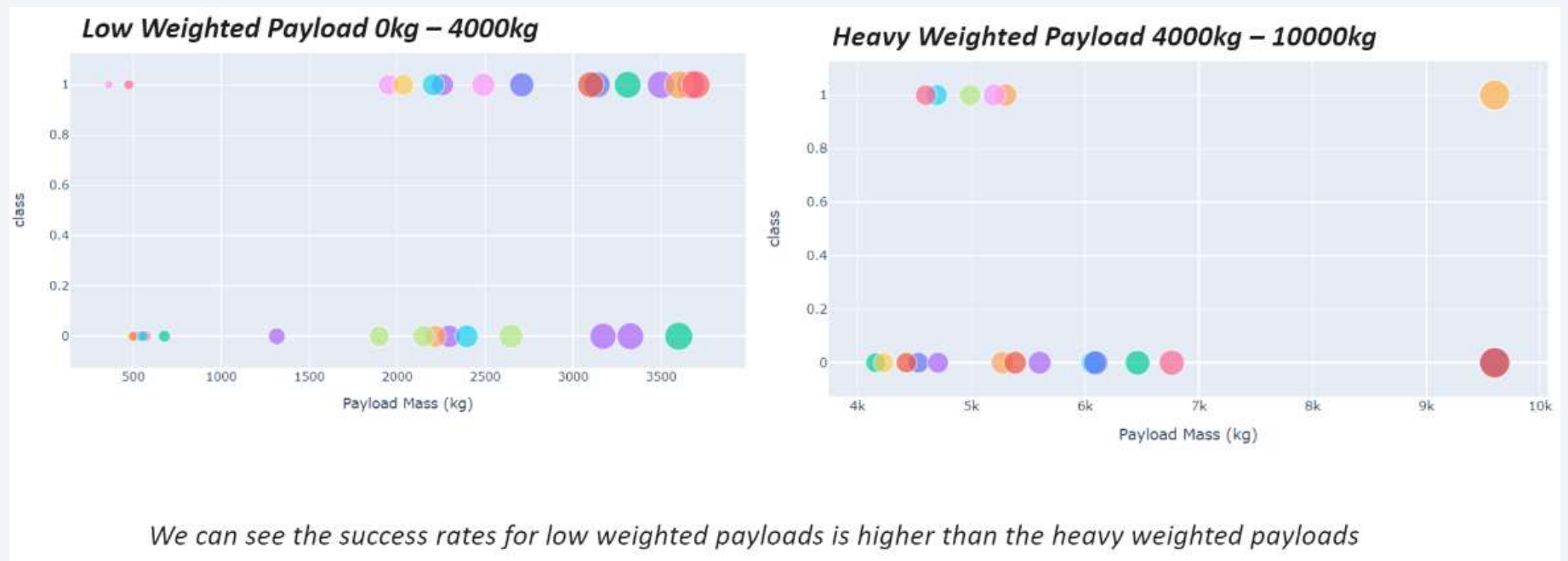


We can see that KSC LC-39A had the most successful launches from all the sites

Pie chart showing the launch site with the highest launch success ratio



Scatter plot of Payload vs Launch Outcome for all sites, with different payload selected in the range slider





Section 5

Predictive Analysis (Classification)

Classification Accuracy

```
models = {'KNeighbors': knn_cv.best_score_,
          'DecisionTree': tree_cv.best_score_,
          'LogisticRegression': logreg_cv.best_score_,
          'SupportVector': svm_cv.best_score_}

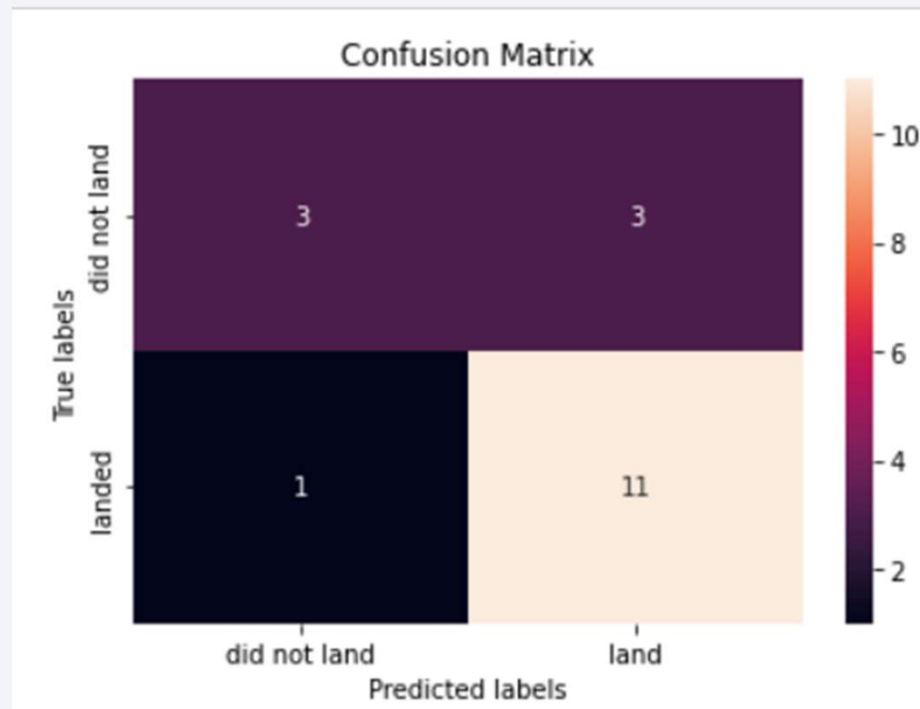
bestalgorithm = max(models, key=models.get)
print('Best model is', bestalgorithm, 'with a score of', models[bestalgorithm])
if bestalgorithm == 'DecisionTree':
    print('Best params is :', tree_cv.best_params_)
if bestalgorithm == 'KNeighbors':
    print('Best params is :', knn_cv.best_params_)
if bestalgorithm == 'LogisticRegression':
    print('Best params is :', logreg_cv.best_params_)
if bestalgorithm == 'SupportVector':
    print('Best params is :', svm_cv.best_params_)
```

Best model is DecisionTree with a score of 0.8892857142857145

Best params is : {'criterion': 'gini', 'max_depth': 4, 'max_features': 'auto', 'min_samples_leaf': 1, 'min_samples_split': 2, 'splitter': 'random'}

- Decision tree is the best classifier with parameters of {'criterion': 'gini', 'max_depth': 4, 'max_features': 'auto', 'min_samples_leaf': 1, 'min_samples_split': 2, 'splitter': 'random'}
- The accuracy score is 88.93%

Confusion Matrix



- The confusion matrix for the decision tree classifier shows that the classifier can distinguish between the different classes. However, the false positives is still quite high .i.e., failed landing marked as successful landing by the classifier.

Conclusions

We can conclude that:

- Higher amount of flight at a launch site is positively correlated with the success rate at the launch site.
- Launch success rate started to increase in 2013 till 2018.
- Orbits ES-L1, GEO, HEO, SSO, VLEO had the most success rate at 100%.
- KSC LC-39A had the most successful launches of any sites.
- The Decision tree classifier is the best machine learning algorithm for this project.

Suggestions

My suggestions for SpaceY:

1. Repeatedly use KSC LC-39A launch site to launch our spacecraft
2. Launches rocket with high amount of payload into LEO orbit would improve success rate
3. Landing on ground pad is very likely to success

Thank you!

