

Winning Space Race with Data Science

Paulo Hiroaqui Ruiz Nakashima
November 29th, 2021



Outline

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

Executive Summary

- Summary of methodologies
 - Data collection
 - Data wrangling
 - EDA with data visualization
 - EDA with SQL
 - Building an interactive map with Folium
 - Building a Dashboard with Plotly Dash
 - Predictive analysis (Classification)
- Summary of all results
 - Exploratory data analysis results
 - Interactive analytics demo in screenshots
 - Predictive analysis results

Introduction

Project background and context:

We predicted if the Falcon 9 first stage will land successfully. SpaceX 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 SpaceX 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 SpaceX for a rocket launch.

Problems you want to find answers:

- What influences if the rocket will land successfully?
- The effect each relationship with certain rocket variables will impact in determining the success rate of a successful landing.
- What conditions does SpaceX have to achieve to get the best results and ensure the best rocket success landing rate.

Section 1

Methodology

Methodology

Data collection methodology:

- SpaceX Rest API
- Web Scrapping from Wikipedia

Perform data wrangling

- One Hot Encoding data fields for Machine Learning and dropping irrelevant columns

Perform exploratory data analysis (EDA) using visualization and SQL

- Plotting : Scatter Graphs, Bar Graphs to show relationships between variables to show patterns of data.

Perform interactive visual analytics using Folium and Plotly Dash

Perform predictive analysis using classification models

- How to build, tune, evaluate classification models

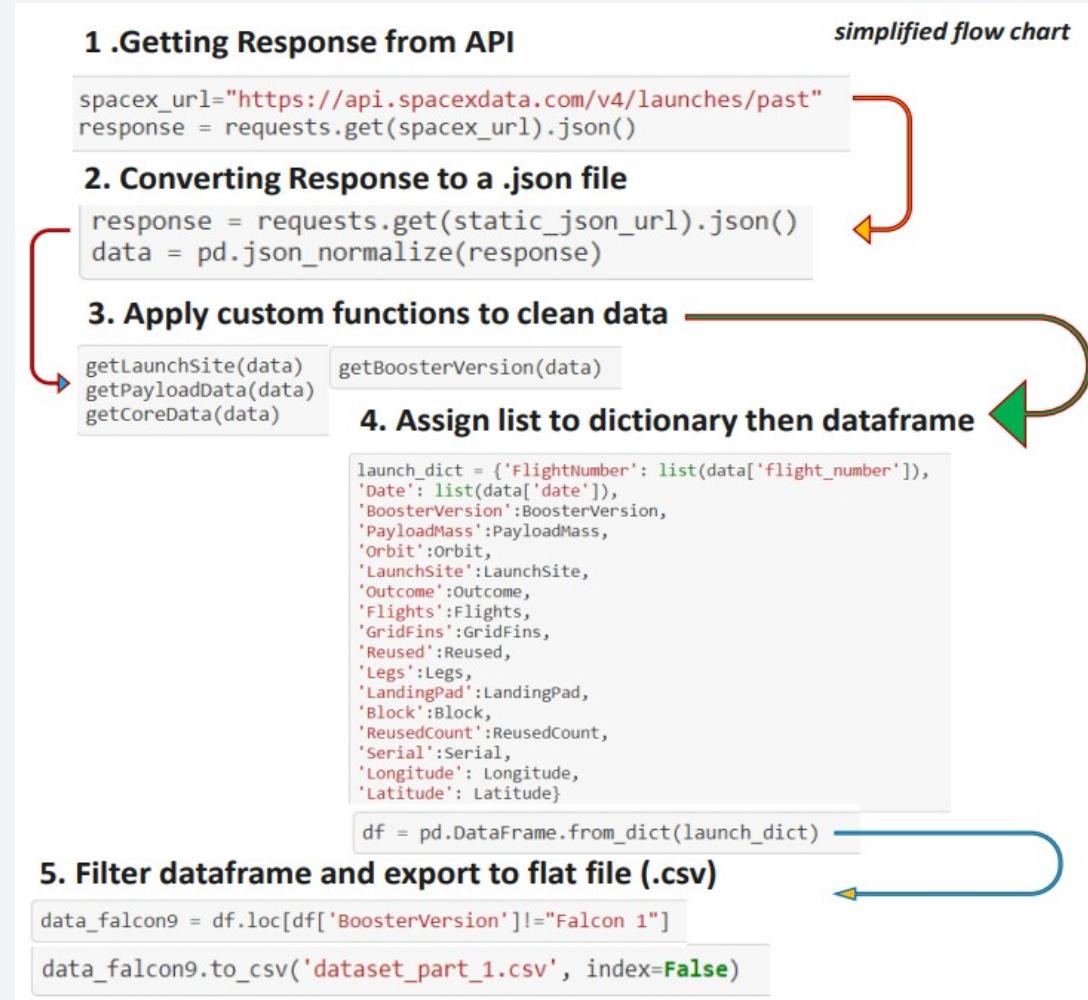
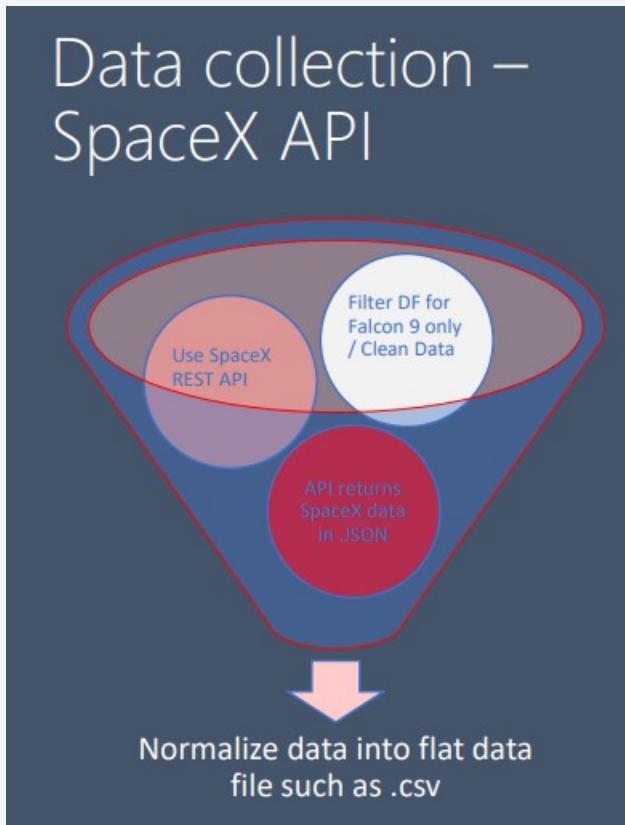
Data Collection

The following datasets was collected by

- We worked with SpaceX launch data that is gathered from the SpaceX REST API.
- This API will give us data about launches, including information about the rocket used, payload delivered, launch specifications, landing specifications, and landing outcome.
- Our goal is to use this data to predict whether SpaceX will attempt to land a rocket or not.
- The SpaceX REST API endpoints, or URL, starts with api.spacexdata.com/v4/.
- Another popular data source for obtaining Falcon 9 Launch data is web scraping Wikipedia using BeautifulSoup.

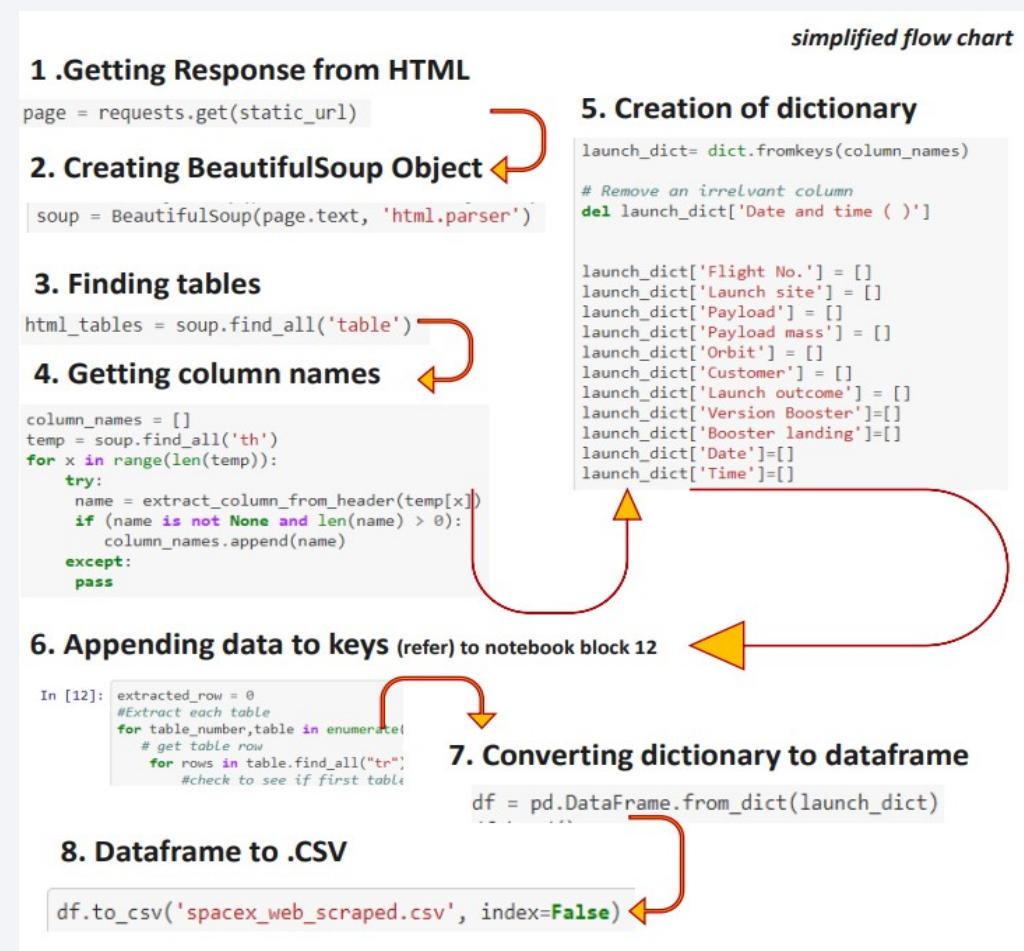
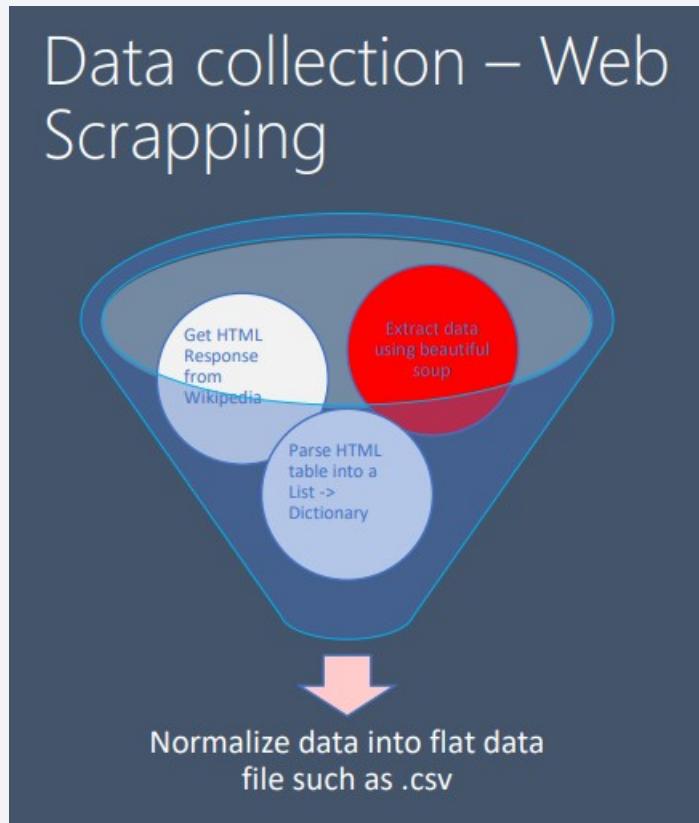
Data Collection – SpaceX API

<https://github.com/pnakashima/space-y/blob/master/Space%20Y.ipynb>



Data Collection - Scraping

<https://github.com/pnakashima/space-y/blob/master/Data%20Collection.ipynb>



Data Wrangling

In the data set, there are several different cases where the booster did not land successfully. Sometimes a landing was attempted but failed due to an accident; for example, True Ocean means the mission outcome was successfully landed to a specific region of the ocean while False Ocean means the mission outcome was unsuccessfully landed to a specific region of the ocean. True RTLS means the mission outcome was successfully landed to a ground pad False RTLS means the mission outcome was unsuccessfully landed to a ground pad. True ASDS means the mission outcome was successfully landed on a drone ship False ASDS means the mission outcome was unsuccessfully landed on a drone ship. We mainly convert those outcomes into Training Labels with 1 means the booster successfully landed 0 means it was unsuccessful.

Process

Perform Exploratory Data Analysis EDA on dataset

Calculate the number of launches at each site

Calculate the number and occurrence of each orbit

Calculate the number and occurrence of mission outcome per orbit type

Export dataset as .CSV

Create a landing outcome label from Outcome column

Work out success rate for every landing in dataset

Each launch aims to an dedicated orbit, and here are some common orbit types:

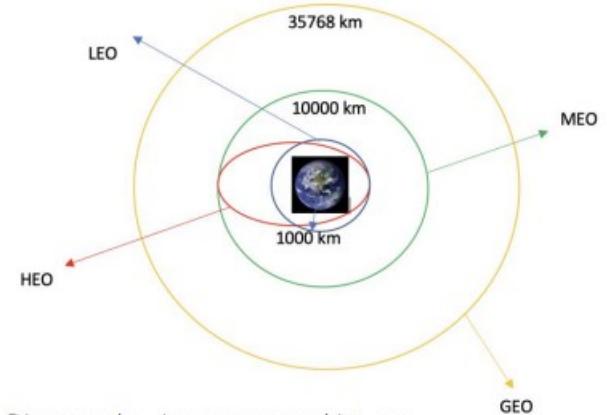


Diagram showing common orbit types
SpaceX uses

<https://github.com/pnakashima/space-y/blob/master/Data%20Wrangling.ipynb>

EDA with Data Visualization

Scatter Graphs being drawn:

- 1) Flight Number VS. Payload Mass
- 2) Flight Number VS. Launch Site
- 3) Payload VS. Launch Site
- 4) Orbit VS. Flight Number
- 5) Payload VS. Orbit Type
- 6) Orbit VS. Payload Mass

Scatter plots show how much one variable is affected by another. The relationship between two variables is called their correlation . Scatter plots usually consist of a large body of data.

<https://github.com/pnakashima/space-y/blob/master/Exploratory%20Analysis%20with%20Pandas%20and%20Matplotlib.ipynb>

Bar Graph being drawn:

- 1) Mean VS. Orbit

A bar diagram makes it easy to compare sets of data between different groups at a glance.

The graph represents categories on one axis and a discrete value in the other. The goal is to show the relationship between the two axes. Bar charts can also show big changes in data over time.

Line Graph being drawn:

- 2) Success Rate VS. Year

Line graphs are useful in that they show data variables and trends very clearly and can help to make predictions about the results of data not yet recorded

EDA with SQL

Performed SQL queries to gather information about the dataset:

- Displaying the names of the unique launch sites in the space mission
- Displaying 5 records where launch sites begin with the string 'KSC'
- Displaying the total payload mass carried by boosters launched by NASA (CRS)
- Displaying average payload mass carried by booster version F9 v1.1
- Listing the date where the successful landing outcome in drone ship was achieved.
- Listing the names of the boosters which have success in ground pad and have payload mass greater than 4000 but less than 6000
- Listing the total number of successful and failure mission outcomes
- Listing the names of the booster_versions which have carried the maximum payload mass.
- Listing the records which will display the month names, successful landing_outcomes in ground pad ,booster versions, launch_site for the months in year 2017
- Ranking the count of successful landing_outcomes between the date 2010-06-04 and 2017-03-20 in descending order.

<https://github.com/pnakashima/space-y/blob/master/Exploratory%20Analysis%20with%20SQL.ipynb>

Build an Interactive Map with Folium

To visualize the Launch Data into an interactive map, we took the Latitude and Longitude Coordinates at each launch site and added a Circle Marker around each launch site with a label of the name of the launch site.

We assigned the dataframe `launch_outcomes(failures, successes)` to classes 0 and 1 with Green and Red markers on the map in a `MarkerCluster()`

Using Haversine's formula we calculated the distance from the Launch Site to various landmarks to find various trends about what is around the Launch Site to measure patterns. Lines are drawn on the map to measure distance to landmarks.

Example of some trends in which the Launch Site is situated in:

- 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

<https://github.com/pnakashima/space-y/blob/master/Interactive%20Visual%20Analytics%20and%20Dashboard.ipynb>

Build a Dashboard with Plotly Dash

Graphs

- Pie Chart showing the total launches by a specific site/all sites
 - display relative proportions of multiple classes of data.
 - size of the circle can be made proportional to the total quantity it represents.
- Scatter Graph showing the relationship with Outcome and Payload Mass (Kg) for the different Booster Versions
 - It shows the relationship between two variables.
 - It is the best method to show you a non-linear pattern.
 - The range of data flow, i.e. maximum and minimum value, can be determined.
 - Observation and reading are straightforward.

https://github.com/pnakashima/space-y/blob/master/spacex_dash_app.py

Predictive Analysis (Classification)

BUILDING MODEL

- Load our dataset into NumPy and Pandas
- Transform Data
- Split our data into training and test data sets
- Check how many test samples we have
- Decide which type of machine learning algorithms we want to use
- Set our parameters and algorithms to GridSearchCV
- Fit our datasets into the GridSearchCV objects and train our dataset.

EVALUATING MODEL

- Check accuracy for each model
- Get tuned hyperparameters for each type of algorithms
- Plot Confusion Matrix

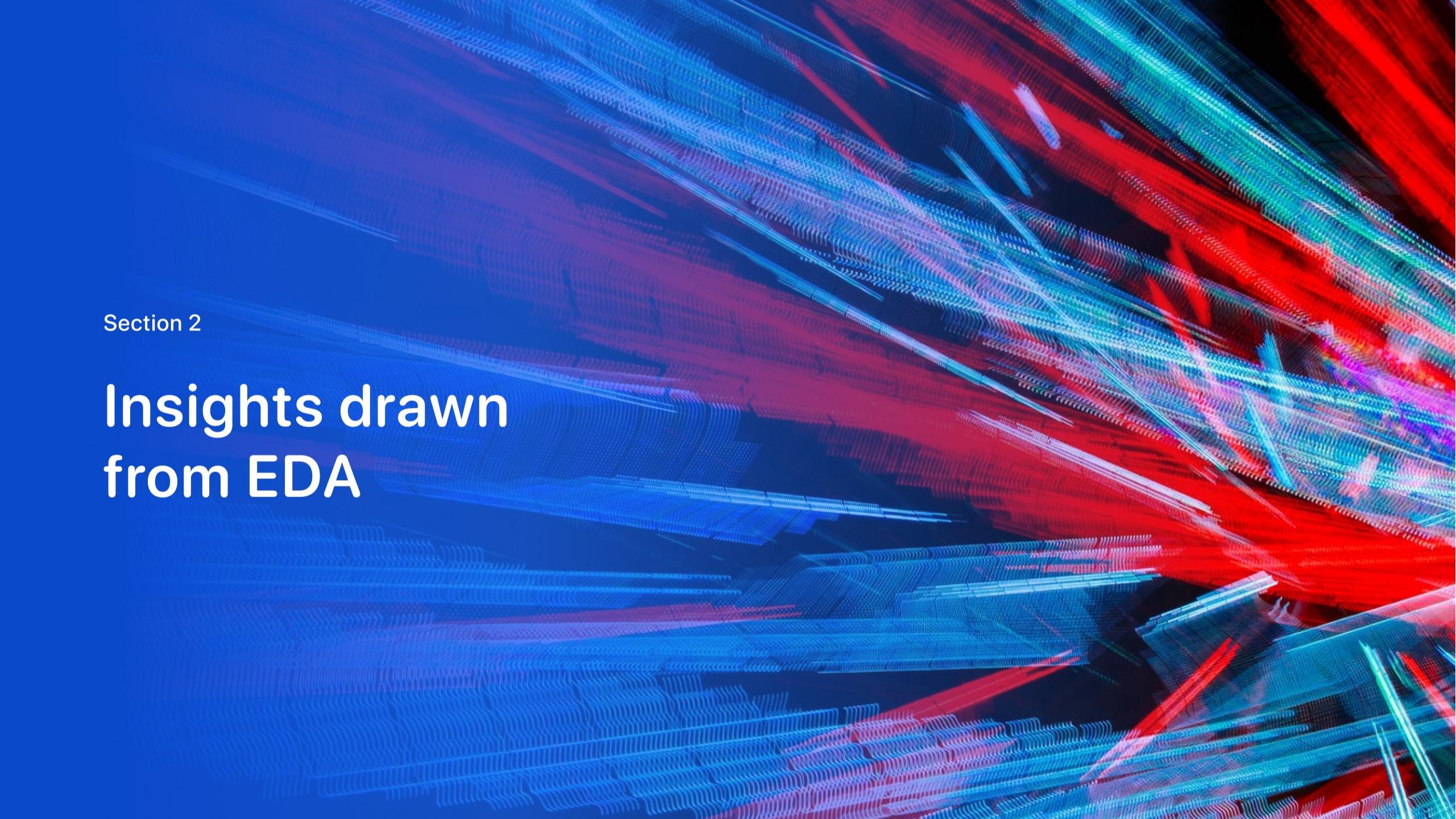
IMPROVING MODEL

- Feature Engineering
- Algorithm Tuning

FINDING THE BEST PERFORMING CLASSIFICATION MODEL

- The model with the best accuracy score wins the best performing model
- In the notebook there is a dictionary of algorithms with scores at the bottom of the notebook.

[https://github.com/pnakashima/space-y/blob/master/Predictive%20Analysis%20\(Classification\).ipynb](https://github.com/pnakashima/space-y/blob/master/Predictive%20Analysis%20(Classification).ipynb)

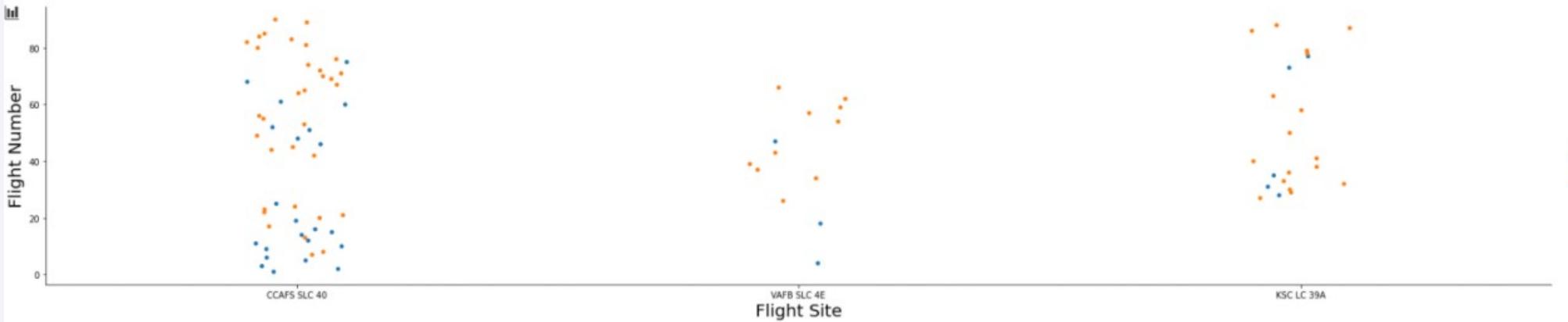
The background of the slide features a complex, abstract digital visualization. It consists of numerous thin, glowing lines that create a sense of depth and motion. The lines are primarily blue and red, with some green and purple highlights. They form a grid-like structure that curves and twists across the frame, resembling a three-dimensional space or a network of data points. The overall effect is futuristic and dynamic.

Section 2

Insights drawn from EDA

Results

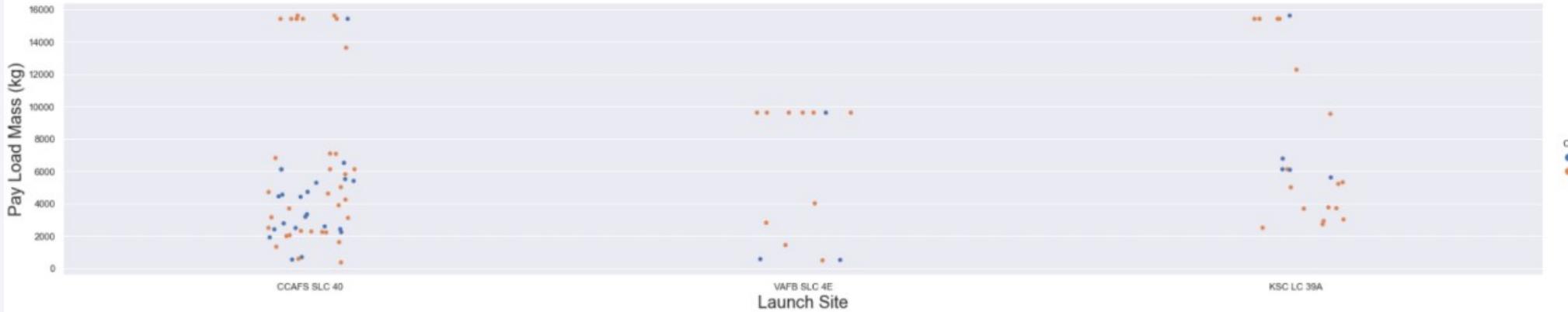
Flight Number vs. Flight Site



The more amount of flights at a launch site the greater the success rate at a launch site.

Results

Payload Mass vs. Launch Site

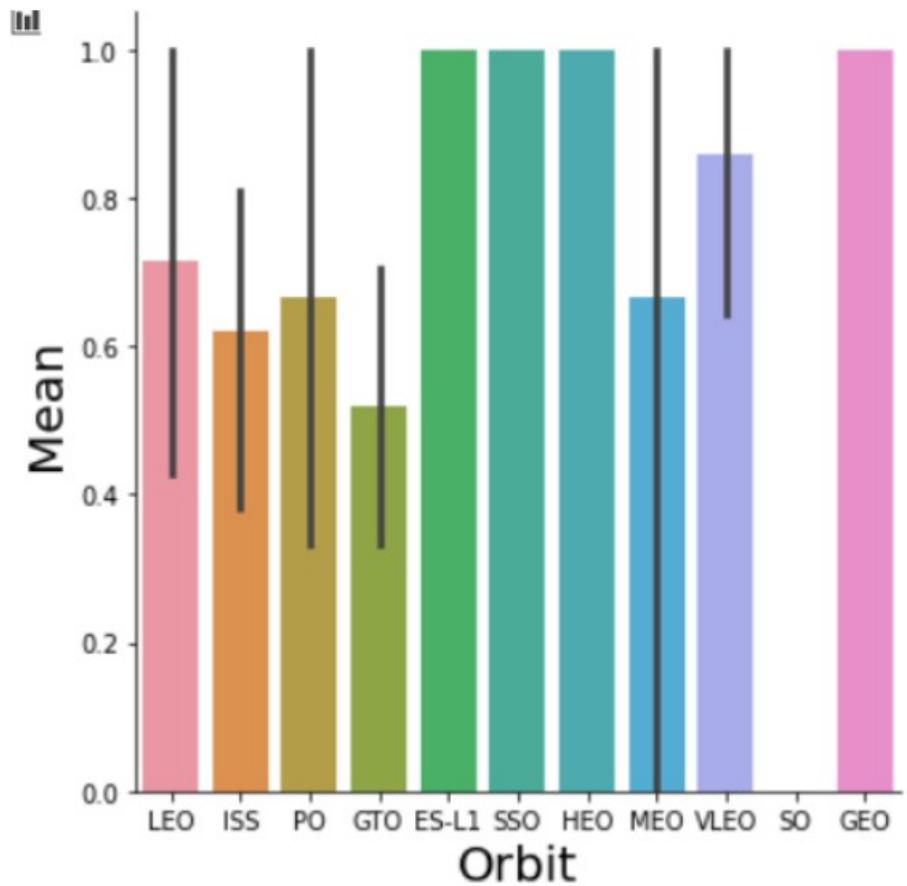


The greater the payload mass for Launch Site CCAFS SLC 40 the higher the success rate for the Rocket. There is not quite a clear pattern to be found using this visualization to make a decision if the Launch Site is dependant on Pay Load Mass for a success launch.

Results

Success rate vs. Orbit type

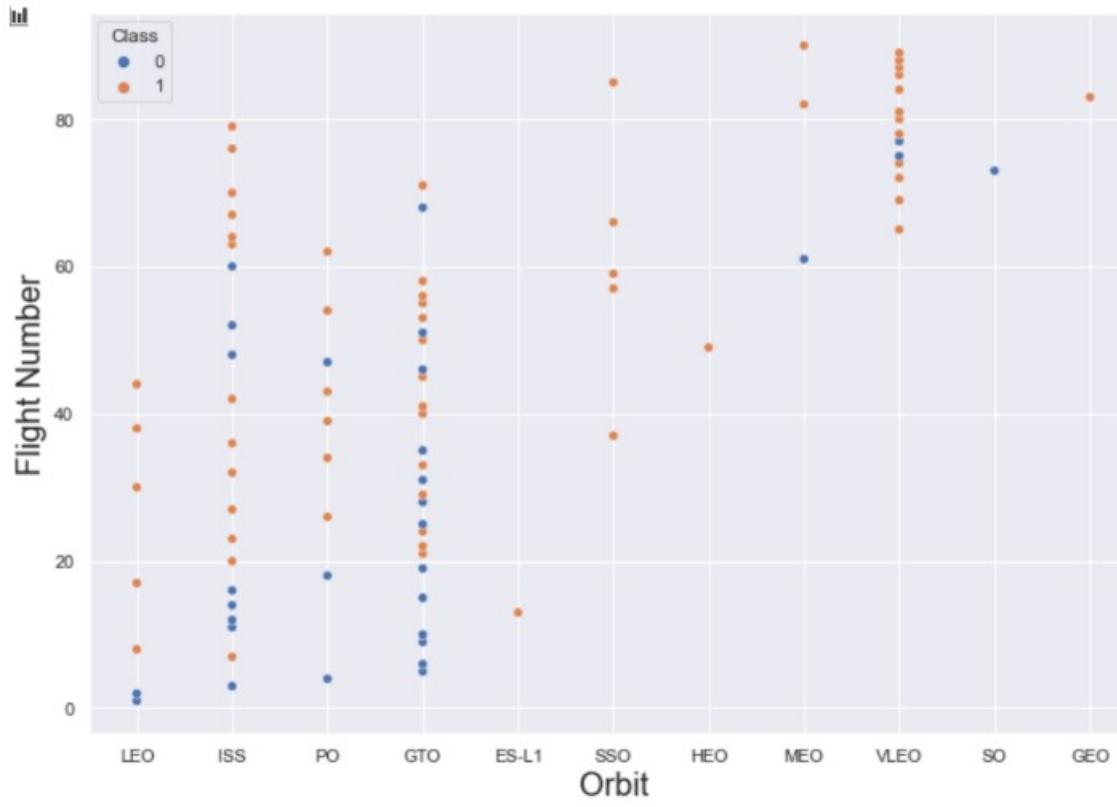
Orbit GEO, HEO, SSO, ES-L1 has the best Success Rate



Results

Flight Number vs. Orbit type

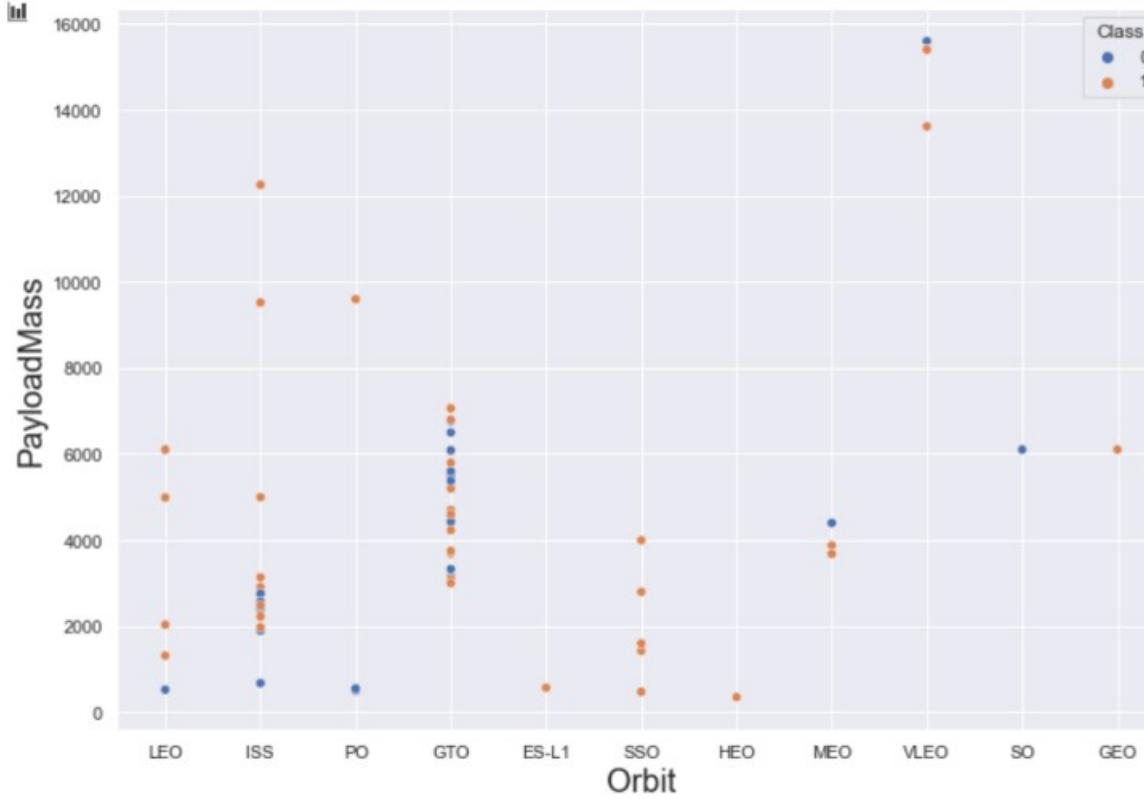
You should see that in the LEO orbit the Success appears related to the number of flights; on the other hand, there seems to be no relationship between flight number when in GTO orbit.



Results

Payload vs. Orbit type

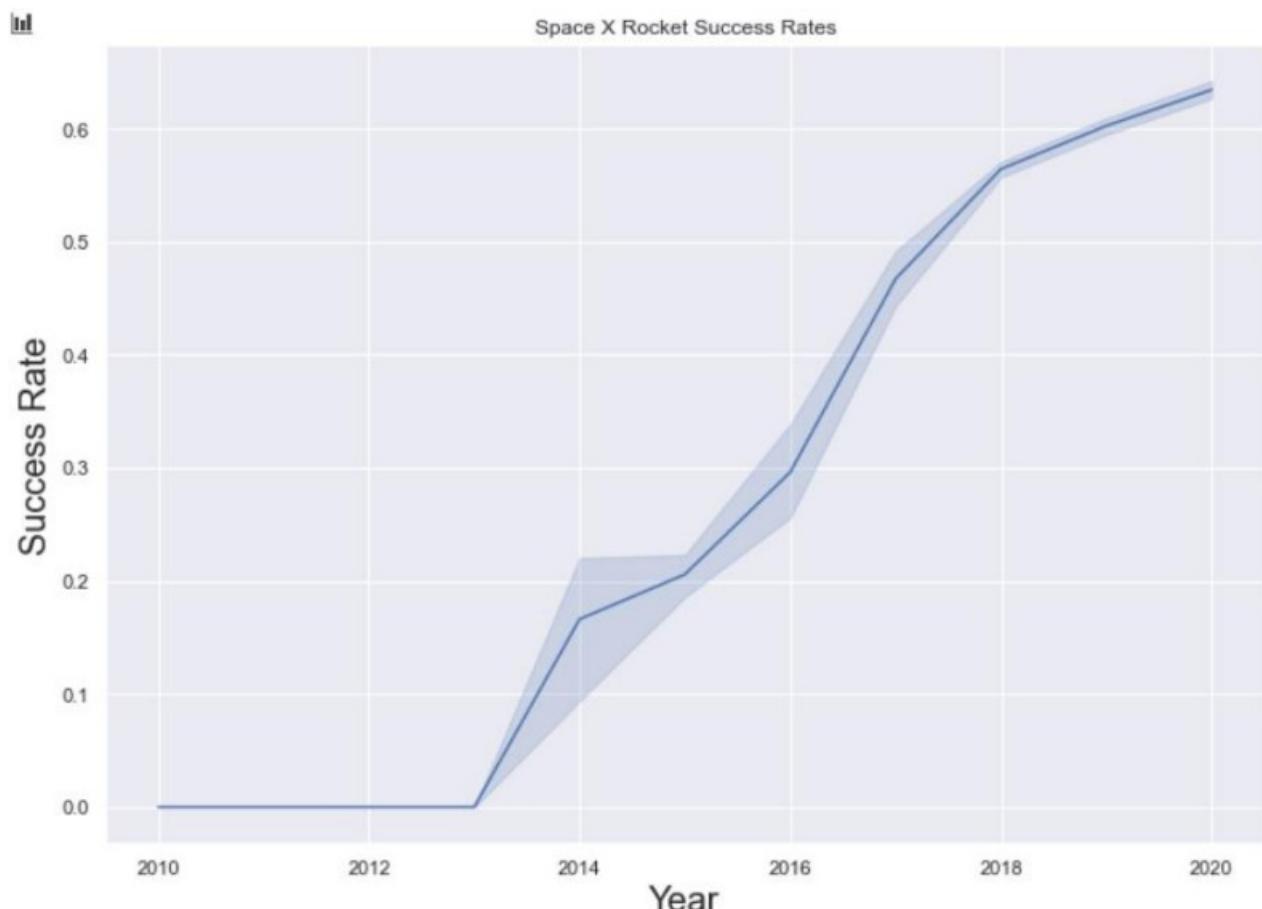
You should observe that Heavy payloads have a negative influence on GTO orbits and positive on GTO and Polar LEO (ISS) orbits.



Results

Launch success yearly trend

you can observe that
the success rate since
2013 kept increasing till
2020



All Launch Site Names

Task 1

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

```
In [6]: %sql select Unique(LAUNCH_SITE) from SPACEXTBL
* ibm_db_sa://rtz30224:***@b0aebb68-94fa-46ec-a1fc-1c999edb6187.c3n41cmd0nqnrk39u98g.databases.appdomain.cloud:31249/bludb
Done.
```

```
Out[6]: launch_site
CCAFS LC-40
CCAFS SLC-40
KSC LC-39A
VAFB SLC-4E
```

Launch Site Names Begin with 'CCA'

Task 2

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

In [7]:

```
%sql SELECT LAUNCH_SITE from SPACEXTBL where (LAUNCH_SITE) LIKE 'CCA%' LIMIT 5
```

```
* ibm_db_sa://rtz30224:***@b0aebb68-94fa-46ec-a1fc-1c999edb6187.c3n41cmd0nqnrk39u98g.databases.appdomain.cloud:31249/bludb
Done.
```

Out[7]: `launch_site`

CCAFS LC-40

CCAFS LC-40

CCAFS LC-40

CCAFS LC-40

CCAFS LC-40

Total Payload Mass

Task 3

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

In [8]:

```
%sql select sum(PAYLOAD_MASS__KG_) as payloadmass from SPACEXTBL  
* ibm_db_sa://rtz30224:***@b0aebb68-94fa-46ec-a1fc-1c999edb6187.c3n41cmd0nqnrk39u98g.databases.appdomain.cloud:31249/bludb  
Done.
```

Out[8]: payloadmass

```
619967
```

Average Payload Mass by F9 v1.1

Task 4

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

In [9]:

```
%sql select avg(PAYLOAD__MASS__KG_) as payloadmass from SPACEXTBL
```

```
* ibm_db_sa://rtz30224:***@b0aebb68-94fa-46ec-a1fc-1c999edb6187.c3n41cmd0nqnrk39u98g.databases.appdomain.cloud:31249/bludb
Done.
```

Out[9]: payloadmass

```
6138
```

First Successful Ground Landing Date

Task 5

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

Hint: Use min function

In [10]:

```
%sql select min(DATE) from SPACEXTBL
```

```
* ibm_db_sa://rtz30224:***@b0aebb68-94fa-46ec-a1fc-1c999edb6187.c3n41cmd0nqnrk39u98g.databases.appdomain.cloud:31249/bludb
Done.
```

Out[10]:

1

2010-06-04

Successful Drone Ship Landing with Payload between 4000 and 6000

Task 6

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

```
In [11]: %sql select BOOSTER_VERSION from SPACEXTBL where LANDING__OUTCOME='Success (drone ship)' and PAYLOAD_MASS__KG_ BETWEEN 4000 and 6000  
* ibm_db_sa://rtz30224:***@b0aebb68-94fa-46ec-a1fc-1c999edb6187.c3n41cmd0nqnrk39u98g.databases.appdomain.cloud:31249/bludb  
Done.  
Out[11]: booster_version  
F9 FT B1022  
F9 FT B1026  
F9 FT B1021.2  
F9 FT B1031.2
```

Total Number of Successful and Failure Mission Outcomes

Task 7

List the total number of successful and failure mission outcomes

```
In [12]: %sql select count(MISSION_OUTCOME) as missionoutcomes from SPACEXTBL GROUP BY MISSION_OUTCOME  
* ibm_db_sa://rtz30224:***@b0aebb68-94fa-46ec-a1fc-1c999edb6187.c3n41cmd0nqnrk39u98g.databases.appdomain.cloud:31249/bludb  
Done.  
Out[12]: missionoutcomes  
1  
99  
1
```

Boosters Carried Maximum Payload

Task 8

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

In [13]:

```
%sql select BOOSTER_VERSION as boosterversion from SPACEXTBL where PAYLOAD_MASS__KG_=(select max(PAYLOAD_MASS__KG_) from SPACEXTBL)  
  
* ibm_db_sa://rtz30224:***@b0aeabb68-94fa-46ec-a1fc-1c999edb6187.c3n41cmd0nqnrk39u98g.databases.appdomain.cloud:31249/bludb  
Done.
```

Out[13]: boosterversion

F9 B5 B1048.4
F9 B5 B1049.4
F9 B5 B1051.3
F9 B5 B1056.4
F9 B5 B1048.5
F9 B5 B1051.4
F9 B5 B1049.5
F9 B5 B1060.2
F9 B5 B1058.3
F9 B5 B1051.6
F9 B5 B1060.3
F9 B5 B1049.7

2015 Launch Records

Task 9

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

```
In [14]: %sql SELECT MONTH(DATE),MISSION_OUTCOME,BOOSTER_VERSION,LAUNCH_SITE FROM SPACEXTBL where EXTRACT(YEAR FROM DATE)='2015'  
* ibm_db_sa://rtz30224:***@b0aebb68-94fa-46ec-a1fc-1c999edb6187.c3n41cmd0nqnrk39u98g.databases.appdomain.cloud:31249/bludb  
Done.  
Out[14]: 1 mission_outcome booster_version launch_site  
1 Success F9 v1.1 B1012 CCAFS LC-40  
2 Success F9 v1.1 B1013 CCAFS LC-40  
3 Success F9 v1.1 B1014 CCAFS LC-40  
4 Success F9 v1.1 B1015 CCAFS LC-40  
4 Success F9 v1.1 B1016 CCAFS LC-40  
6 Failure (in flight) F9 v1.1 B1018 CCAFS LC-40  
12 Success F9 FT B1019 CCAFS LC-40
```

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

Task 10

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

```
In [15]: %sql SELECT LANDING__OUTCOME FROM SPACEXTBL WHERE DATE BETWEEN '2010-06-04' AND '2017-03-20' ORDER BY DATE DESC
* ibm_db_sa://rtz30224:***@b0aebb68-94fa-46ec-a1fc-1c999edb6187.c3n41cmd0nqnrik39u98g.databases.appdomain.cloud:31249/bludb
Done.
```

```
Out[15]: landing_outcome
No attempt
Success (ground pad)
Success (drone ship)
Success (drone ship)
Success (ground pad)
Failure (drone ship)
Success (drone ship)
Success (drone ship)
Success (drone ship)
Failure (drone ship)
Failure (drone ship)
Success (ground pad)
Precluded (drone ship)
No attempt
Failure (drone ship)
No attempt
Controlled (ocean)
Failure (drone ship)
Uncontrolled (ocean)
No attempt
No attempt
Controlled (ocean)
Controlled (ocean)
No attempt
No attempt
```

GitHub Link – EDA with SQL

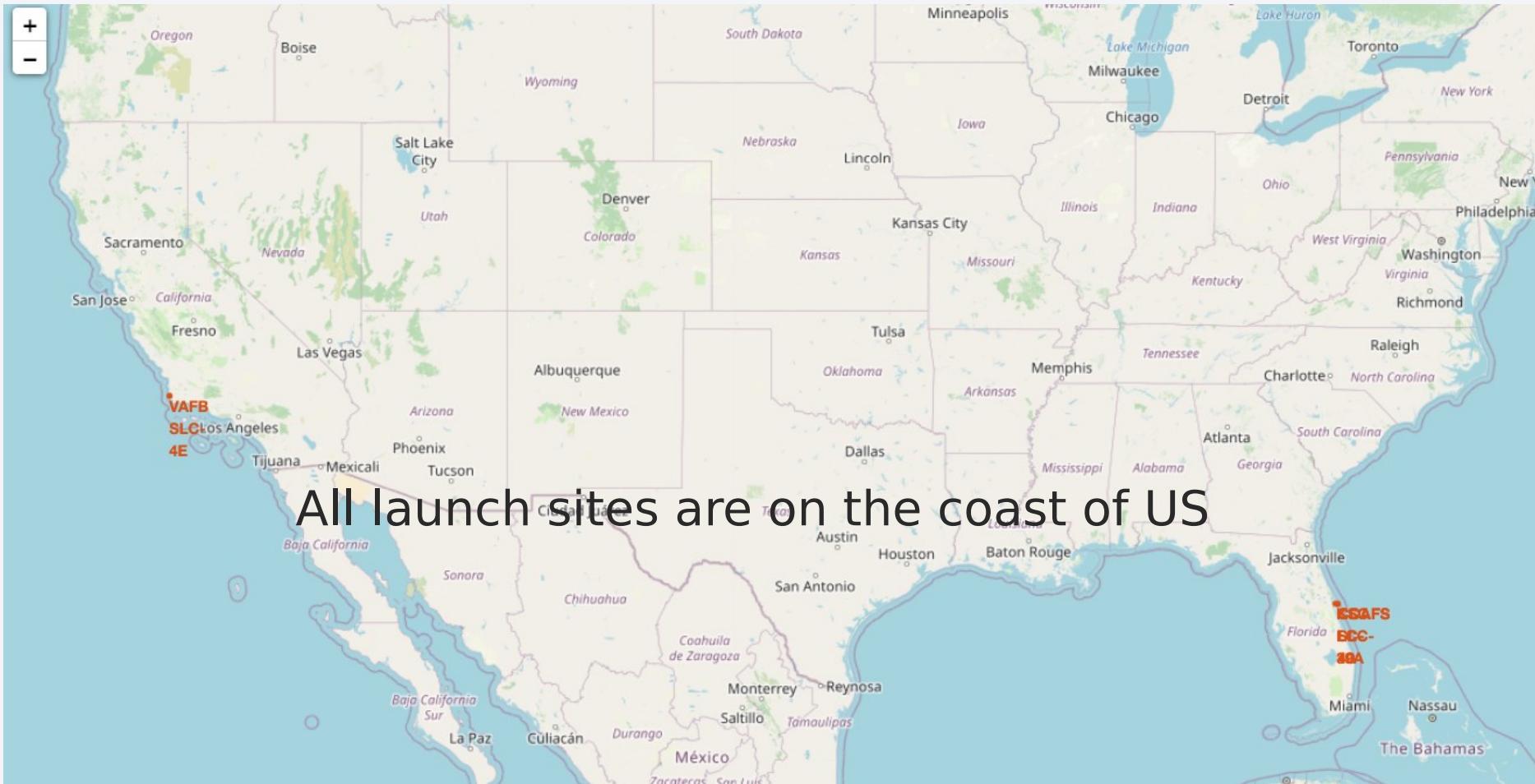
<https://github.com/pnakashima/space-y/blob/master/Exploratory%20Analysis%20with%20SQL.ipynb>

The background of the slide is a photograph taken from space at night. It shows the curvature of the Earth against the dark void of space. City lights are visible as numerous small white and yellow dots, primarily concentrated in the lower right quadrant where the United States appears. In the upper left quadrant, the green and blue glow of the aurora borealis (Northern Lights) is visible, appearing as horizontal bands of light.

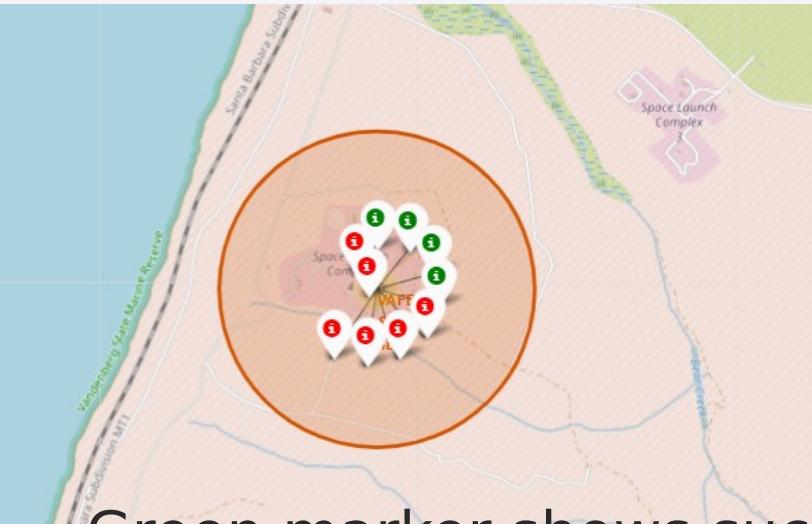
Section 4

Launch Sites Proximities Analysis

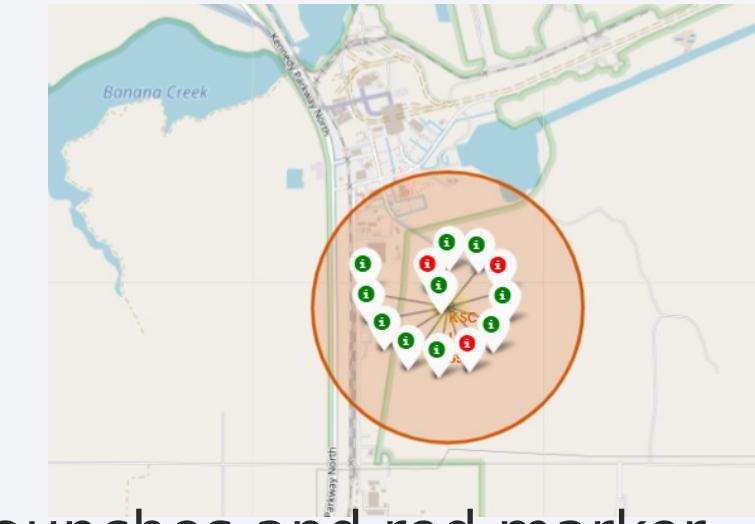
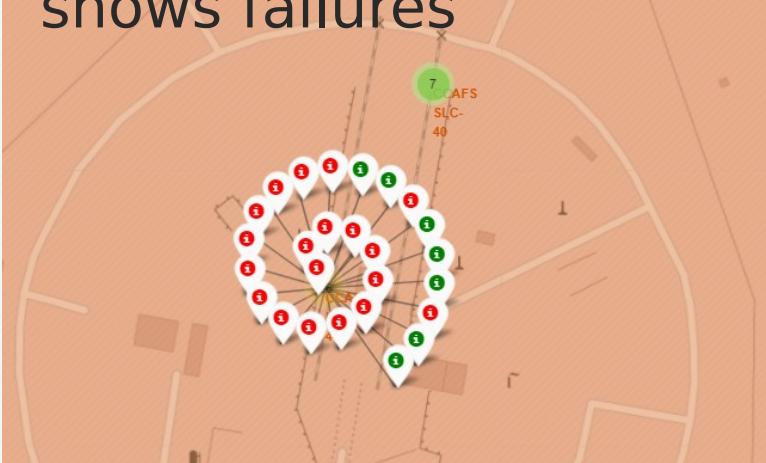
Launch Site Locations



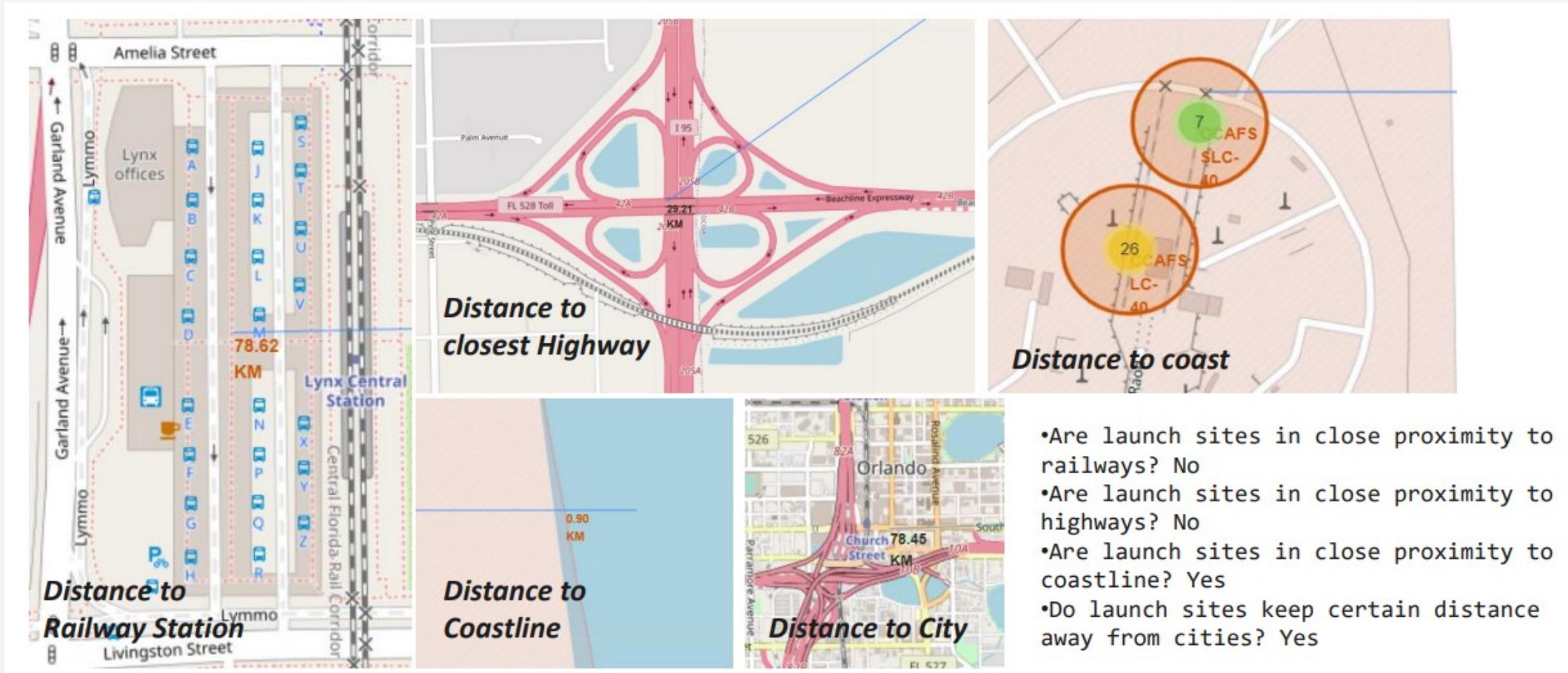
Color Labelled Markers



Green marker shows successful launches and red marker shows failures

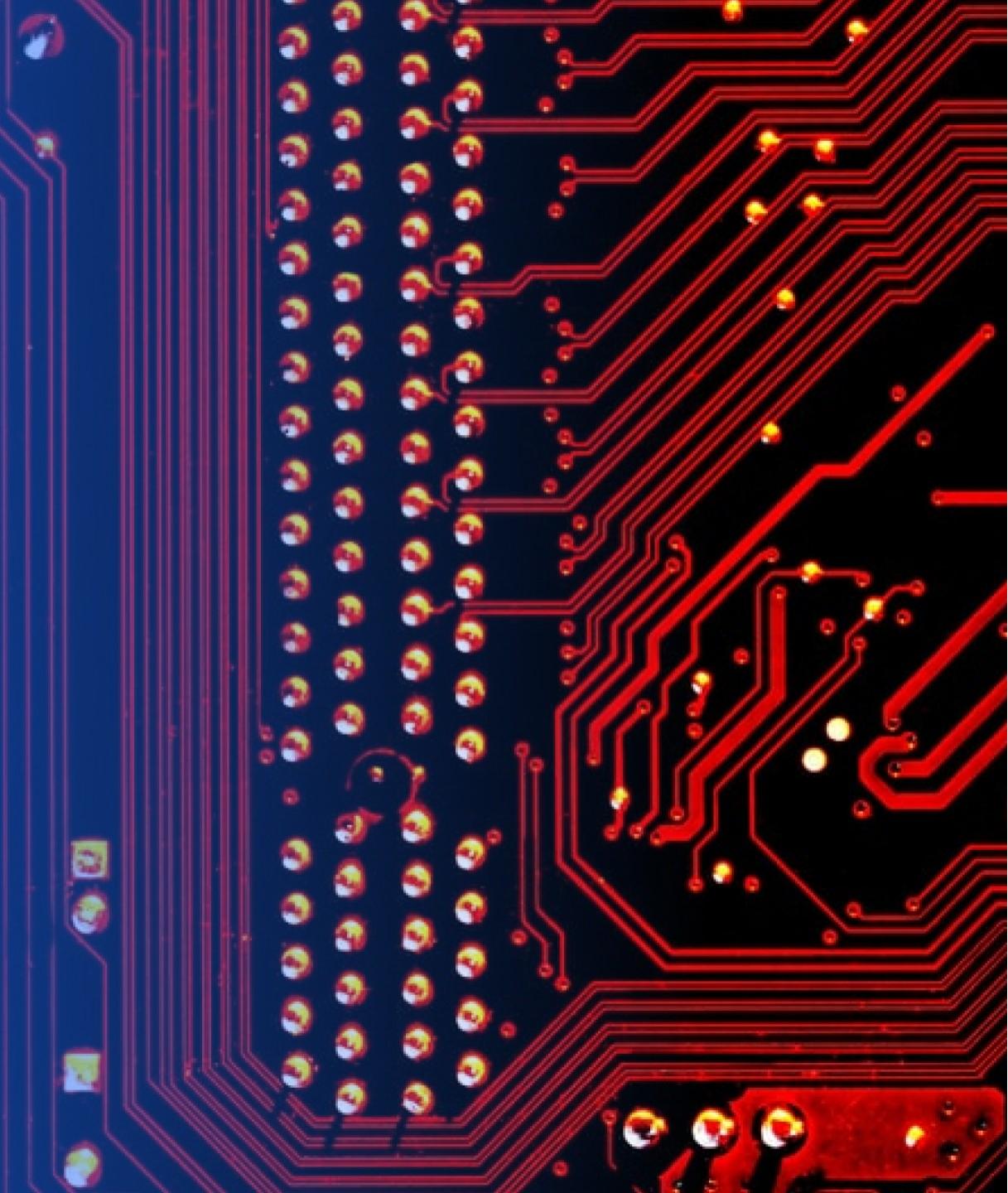


Launch Site Distances

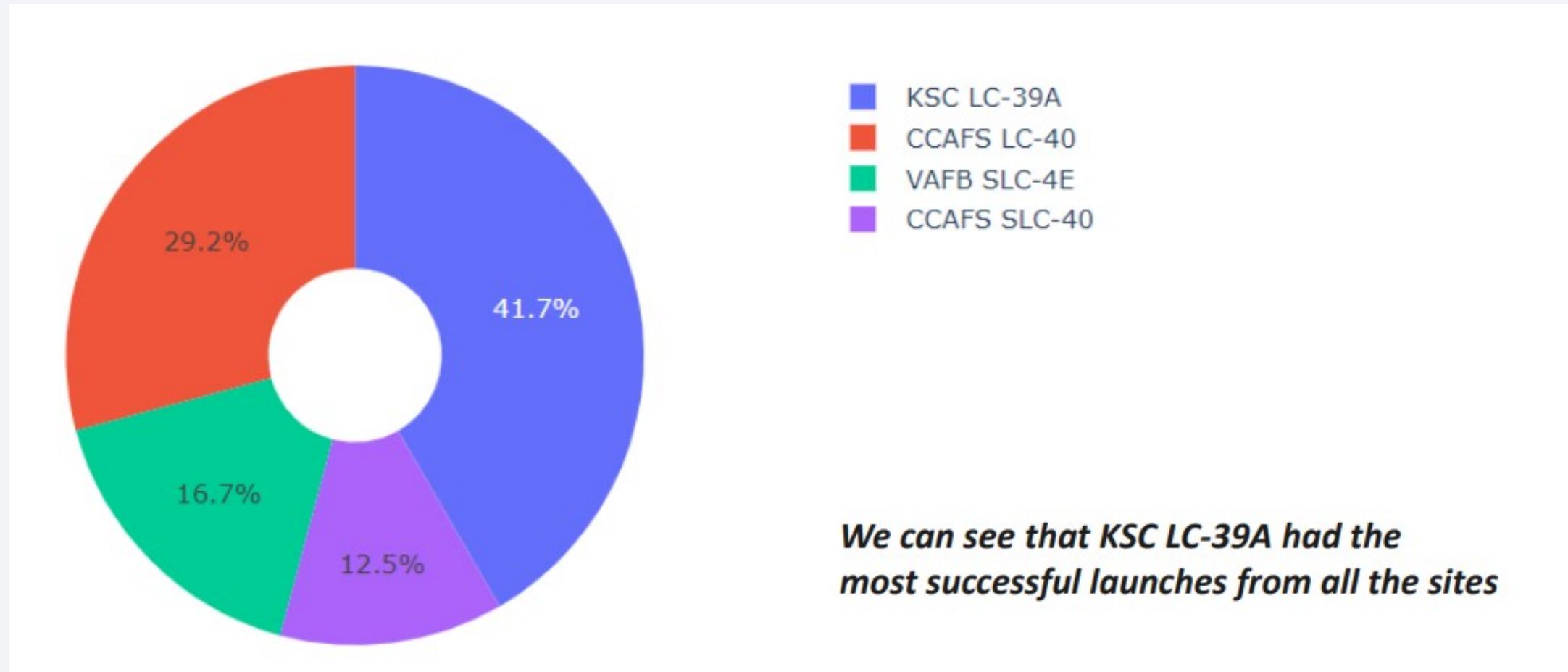


Section 5

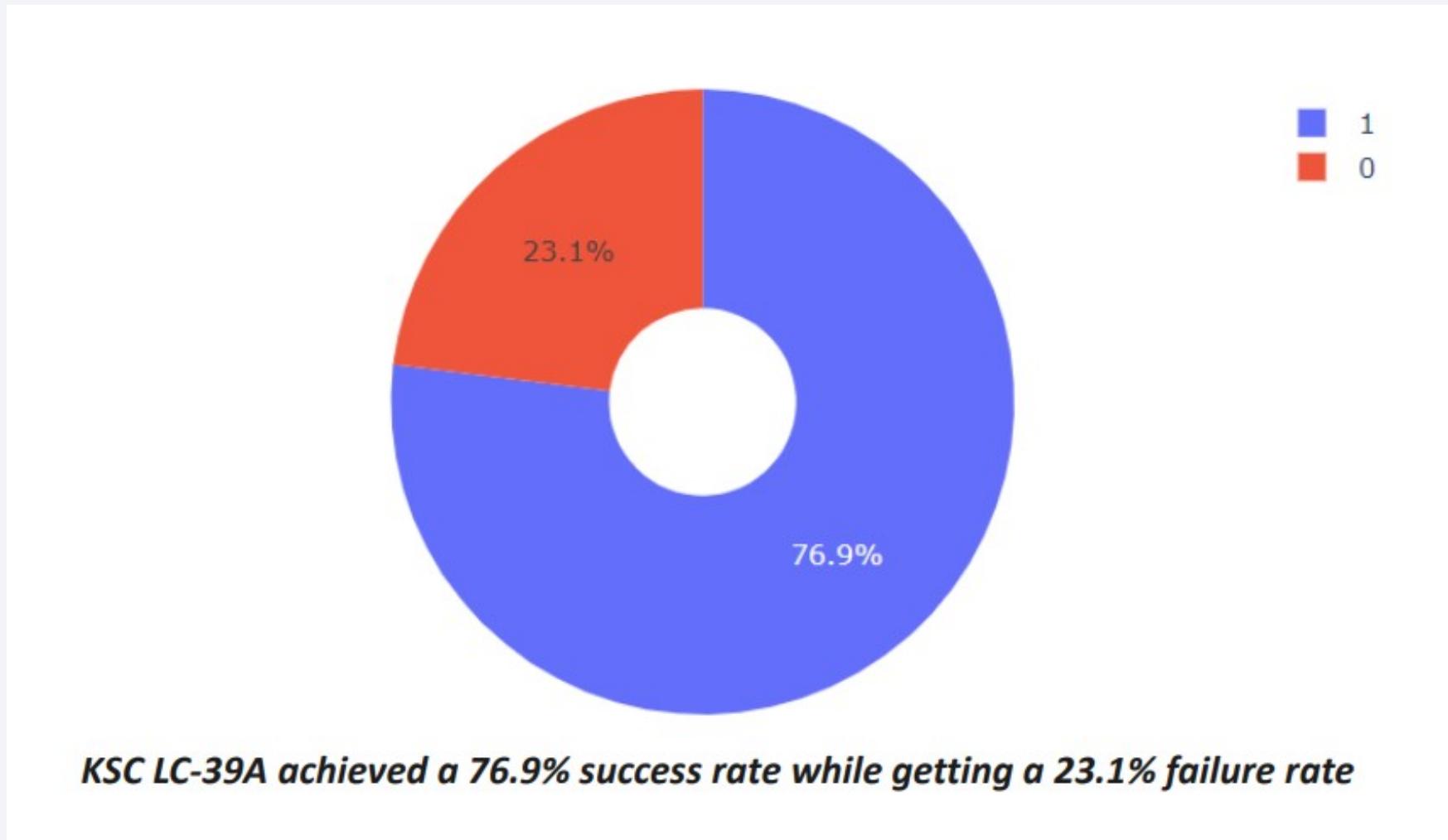
Build a Dashboard with Plotly Dash



Pie chart – Success Percentage by Launch Site

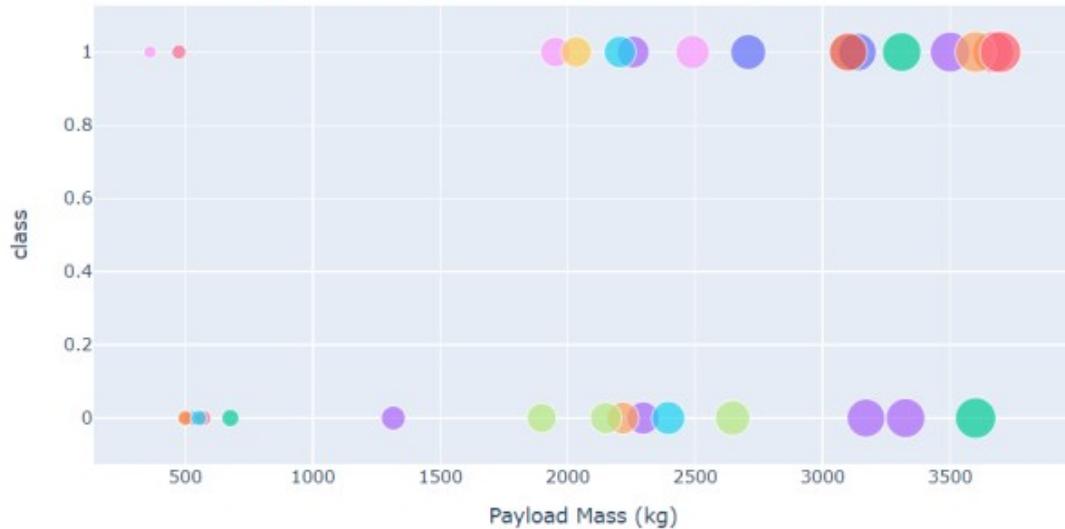


Pie Chart – Highest Success Ratio

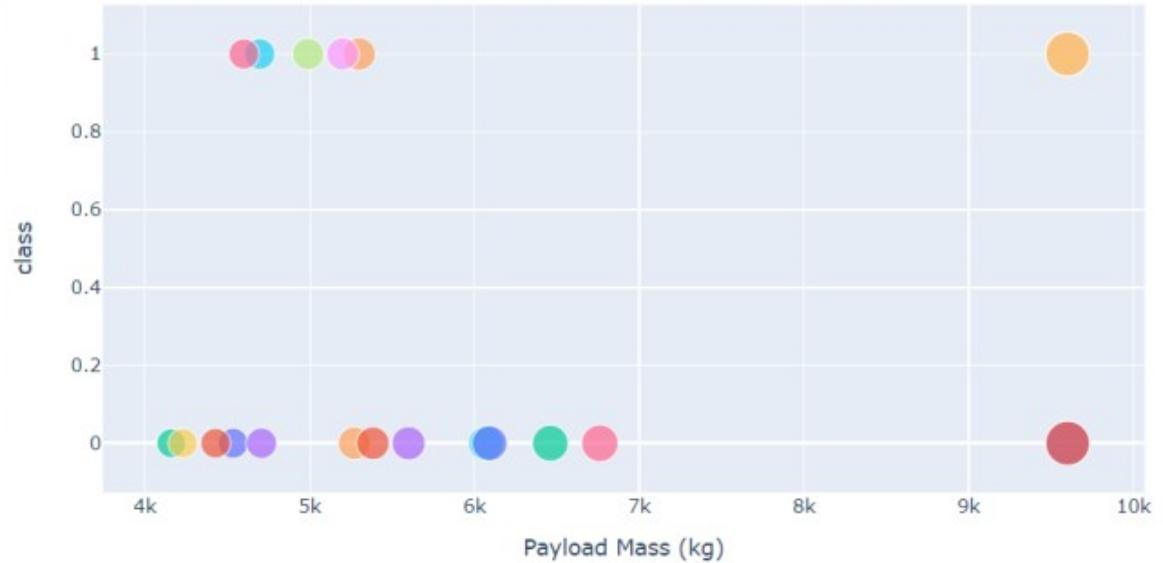


Payload vs Launch Outcome

Low Weighted Payload 0kg – 4000kg



Heavy Weighted Payload 4000kg – 10000kg



We can see the success rates for low weighted payloads is higher than the heavy weighted payloads

GitHub Link – Source Code

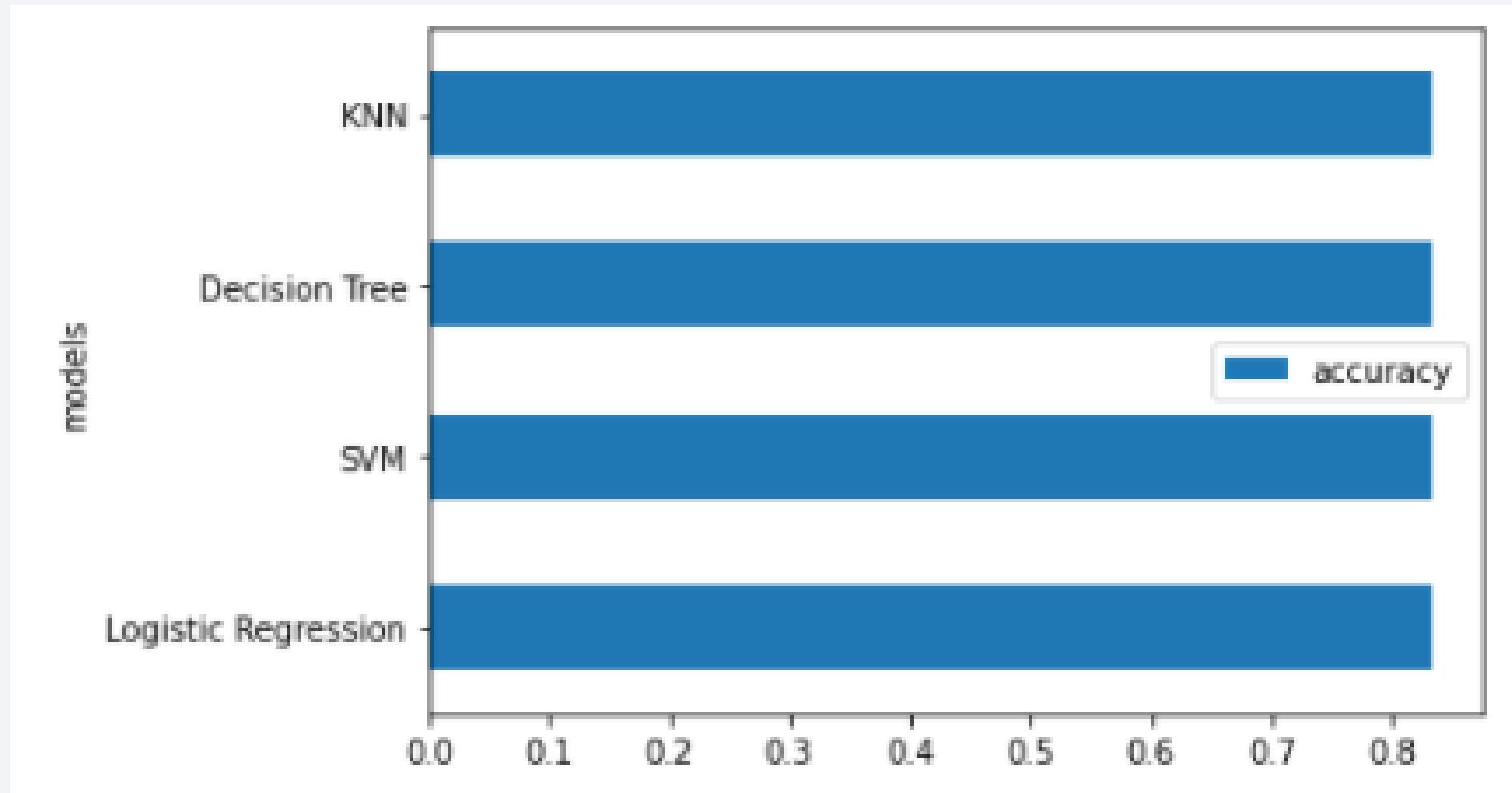
https://github.com/pnakashima/space-y/blob/master/spacex_dash_app.py

Section 6

Predictive Analysis (Classification)

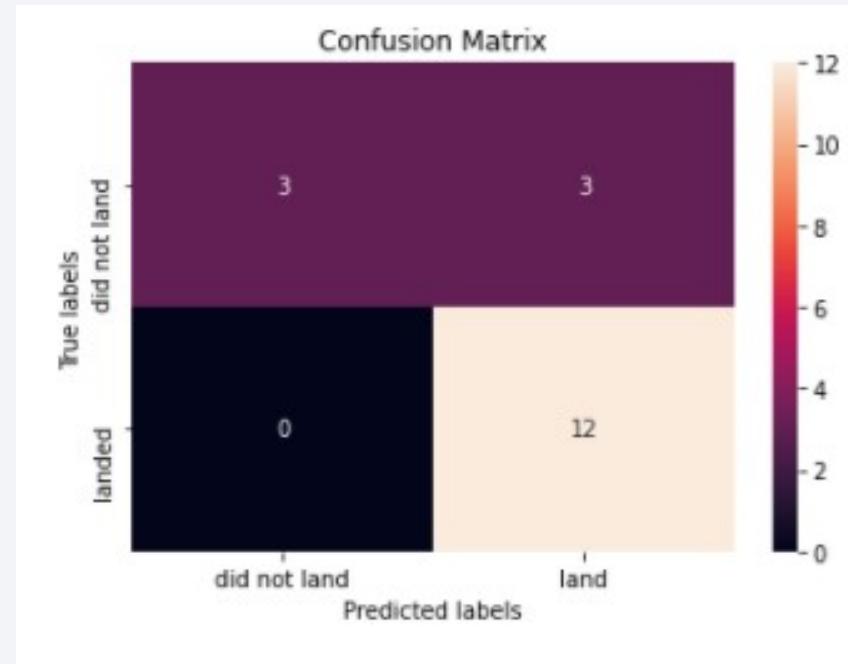
Classification Accuracy

All models have the same accuracy: 0.833333333333334



Confusion Matrix

- This is the SVM confusion matrix. The number of True Positives is 12, False Positives are 3, and True Negatives are also 3.



Conclusions

All Classifier Algorithms have the same performance for this dataset

Low weighted payloads perform better than the heavier payloads

The success rates for SpaceX launches is directly proportional time in years they will eventually perfect the launches

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

Orbit GEO,HEO,SSO,ES-L1 has the best Success Rate

Thank you!

