

Winning Space Race with Data Science

Lucia Gonzalez
2024



Outline

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

Executive Summary

- **Summary of methodologies:**

- Data collection
- Data wrangling
- Exploratory Data Analysis with Data Visualization
- Exploratory Data Analysis 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:**

SpaceX is the most successful company of the commercial space age, making space travel affordable. The company 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. Based on public information and machine learning models, we are going to predict if SpaceX will reuse the first stage.

- **Questions to be answered:**

- How do variables such as payload mass, launch site, number of flights, and orbits affect the success of the first stage landing?
- Does the rate of successful landings increase over the years?
- What is the best algorithm that can be used for binary classification in this case?

Section 1

Methodology

Methodology

Executive Summary

- Data collection methodology:
 - Used SpaceX Rest API and web scrapping from wikipedia
- Perform data wrangling
 - Filtering the data, dealt with missing values and used One Hot Encoding to prepare the data for binary classification
- Perform exploratory data analysis (EDA) using visualization and SQL
- Perform interactive visual analytics using Folium and Plotly Dash
- Perform predictive analysis using classification models
 - Used machine learning classification models and evaluated them to ensure the best results.

Data Collection

Data collection process involved a combination of **API requests from SpaceX REST API** and **Web Scraping** data from a table in SpaceX's Wikipedia entry.

- Columns obtained using SpaceX REST API:

FlightNumber, Date, BoosterVersion, PayloadMass, Orbit, LaunchSite, Outcome, Flights, GridFins, Reused, Legs, LandingPad, Block, ReusedCount, Serial, Longitude, Latitude

- Columns obtained using Wikipedia Web Scraping:

Flight No., Launch site, Payload, PayloadMass, Orbit, Customer, Launch outcome, Version Booster, Booster landing, Date, Time

Data Collection – SpaceX API

- GitHub URL:
- [SpaceX Data Collection](#)



Data Collection - Scraping

- GitHub URL

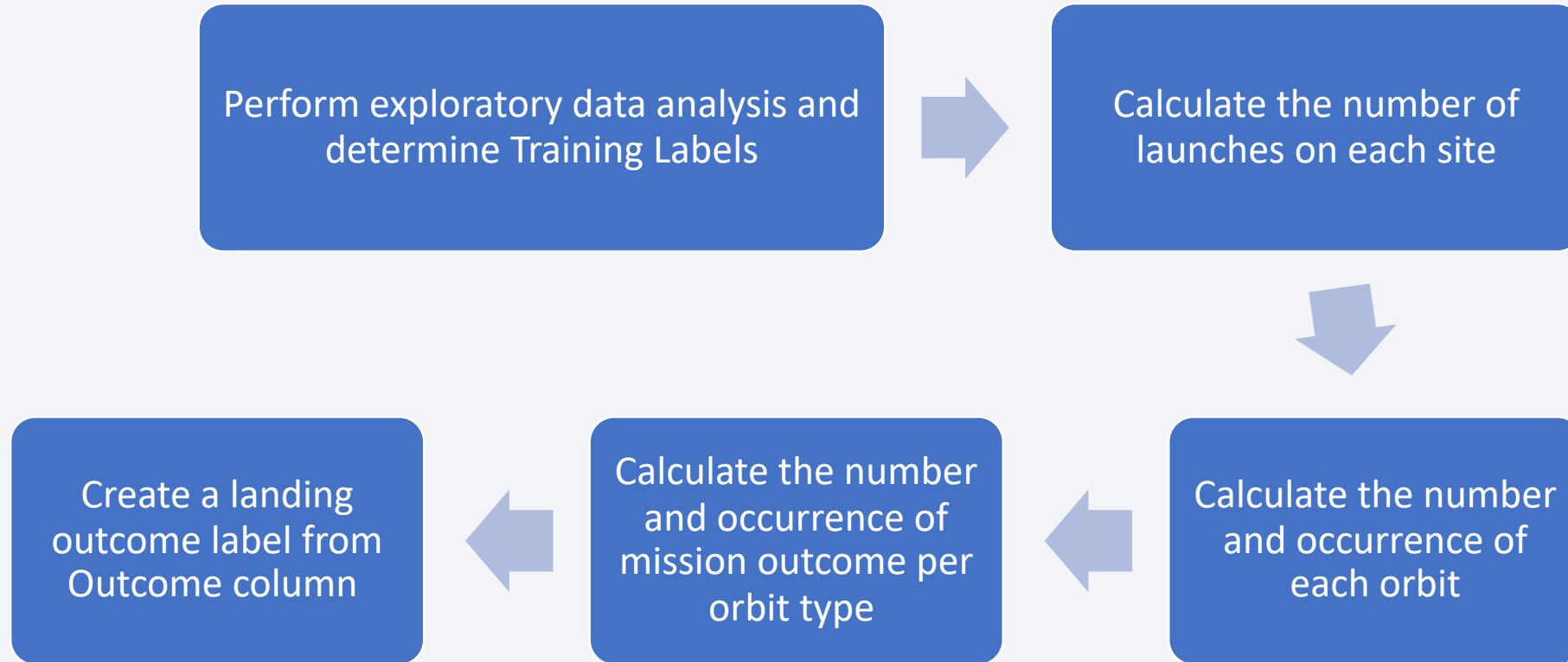
SpaceX Web Scrapping

Request the Falcon 9
launch Wikipedia page

Extract all column/variable
names from the HTML
table header

Create a data frame by
parsing the launch HTML
tables

Data Wrangling



[GitHub: SpaceX data wrangling](#)

EDA with Data Visualization

- Charts were plotted:
 - Flight Number vs. Payload Mass, Flight Number vs. Launch Site, Payload Mass vs. Launch Site, Orbit Type vs. Success Rate, Flight Number vs. Orbit Type, Payload Mass vs Orbit Type and Success Rate Yearly Trend.
- Scatter plots show the relationship between variables. If a relationship exists, they could be used in machine learning model.
- Bar charts show comparisons among discrete categories. The goal is to show the relationship between the specific categories being compared and a measured value.
- Line charts show trends in data over time (time series).

GitHub URL: [SpaceX EDA with data visualization](#)

EDA with SQL

The following SQL queries were performed:

- Names of the unique launch sites in the space mission;
- Top 5 launch sites whose name begin with the string 'CCA';
- Total payload mass carried by boosters launched by NASA (CRS);
- Average payload mass carried by booster version F9 v1.1;
- Date when the first successful landing outcome in ground pad was achieved;
- Names of the boosters which have success in drone ship and have payload mass between 4000 and 6000 kg;
- Total number of successful and failure mission outcomes;
- Names of the booster versions which have carried the maximum payload mass;
- Failed landing outcomes in drone ship, their booster versions, launch site names and month in the year 2015; and
- Rank of the count of landing outcomes (such as Failure (drone ship) or Success (ground pad)) between the date 2010-06-04 and 2017-03-20.

GitHub URL: [SpaceX EDA with SQL](#)

Build an Interactive Map with Folium

- **Markers of all Launch Sites:**
 - Added Marker with Circle, Popup Label and Text Label of NASA Johnson Space Center using its latitude and longitude coordinates as a start location.
 - Added Markers with Circle, Popup Label and Text Label of all Launch Sites using their latitude and longitude coordinates to show their geographical locations and proximity to Equator and coasts.
- **Coloured Markers of the launch outcomes for each Launch Site:**
 - Added coloured Markers of success (Green) and failed (Red) launches using Marker Cluster to identify which launch sites have relatively high success rates.
- **Distances between a Launch Site to its proximities:**
 - Added coloured Lines to show distances between the Launch Site CCAFS SLC - 40 (as an example) and its proximities like Railway, Highway, Coastline and Closest City.

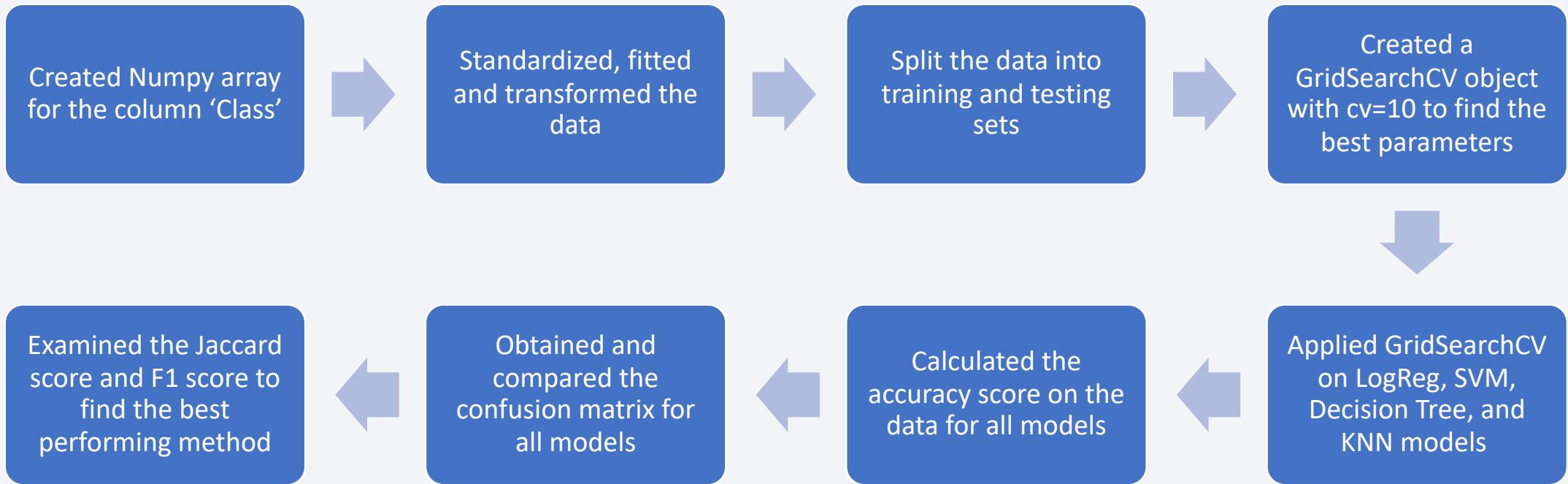
GitHub URL: [SpaceX Interactive Map](#)

Build a Dashboard with Plotly Dash

- Launch Sites Dropdown List:
 - Added a dropdown list to enable Launch Site selection.
- Pie Chart showing Success Launches (All Sites/Certain Site):
 - Added a pie chart to show the total successful launches count for all sites and the Success vs. Failed counts for the site, if a specific Launch Site was selected.
- Slider of Payload Mass Range:
 - Added a slider to select Payload range.
- Scatter Chart of Payload Mass vs. Success Rate for the different Booster Versions:
 - Added a scatter chart to show the correlation between Payload and Launch Success.

GitHub URL: [SpaceX Dashboard](#)

Predictive Analysis (Classification)



GitHub URL: [SpaceX predictive analysis](#)

Results

- Exploratory data analysis results:
 - Many Falcon 9 booster versions were successful at landing in drone ships having payload above the average;
 - Almost 100% of mission outcomes were successful;
 - Two booster versions failed at landing in drone ships in 2015: F9 v1.1 B1012 and F9 v1.1 B1015;
 - The number of landing outcomes became as better as years passed.
- Predictive analysis results:
 - All the methods' accuracy was almost the same, except for tree which fit train data slightly better but test data worse.
 - The launches occurred on the coast, mostly on the east one:

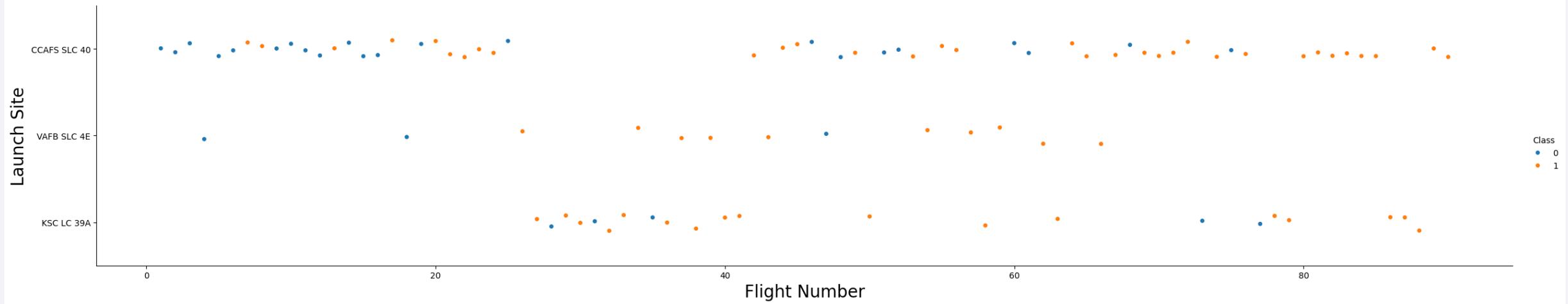


The background of the slide features a complex, abstract digital visualization. It consists of numerous thin, glowing lines that create a sense of depth and motion. The lines are primarily blue and red, with some green and purple highlights. They form a grid-like structure that curves and twists across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

Section 2

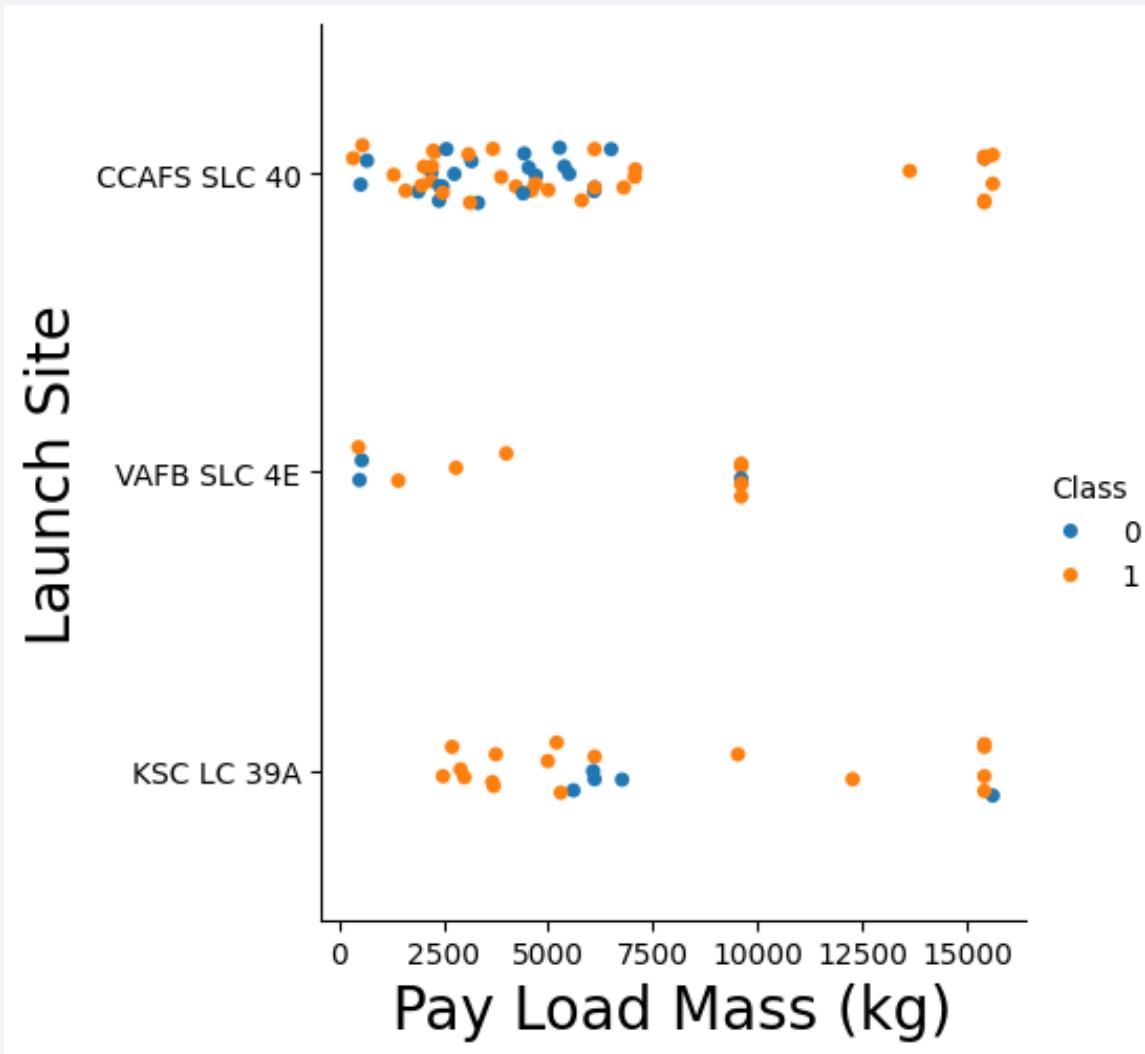
Insights drawn from EDA

Flight Number vs. Launch Site



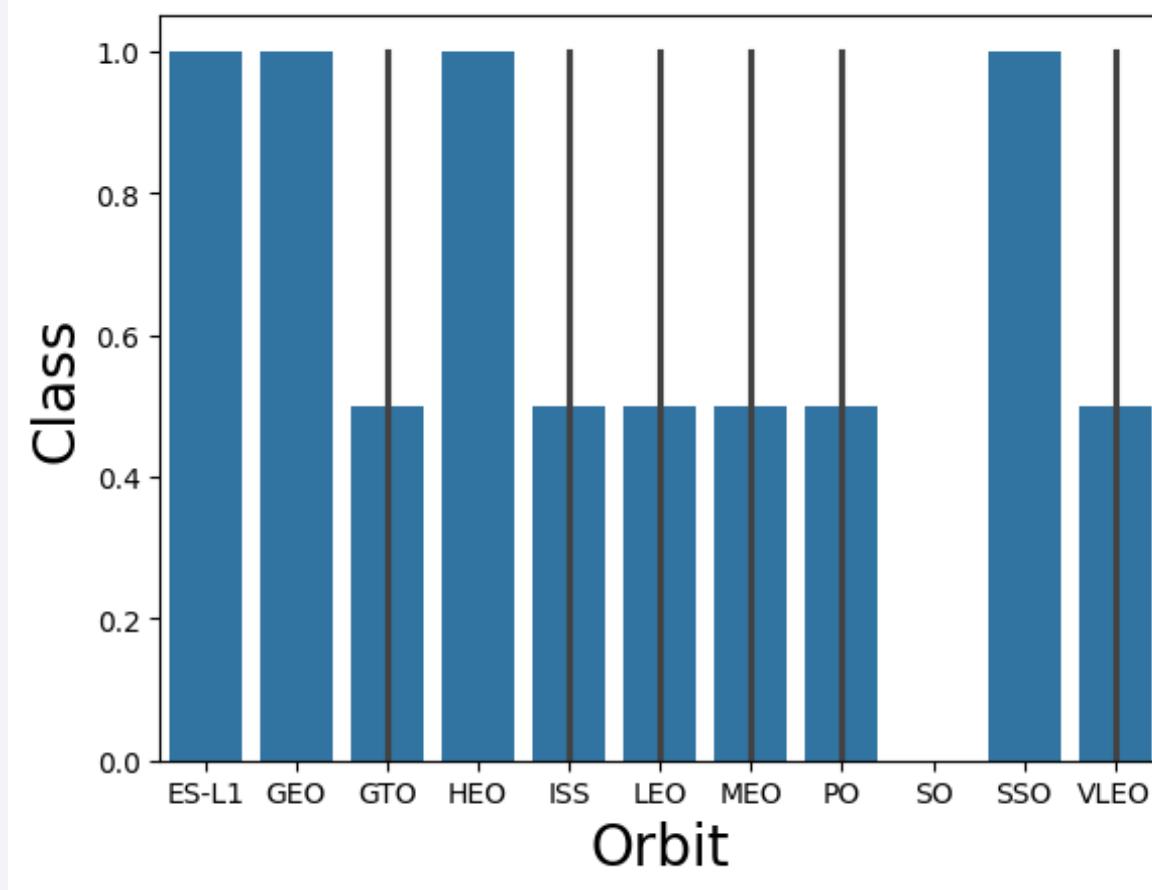
- According to the plot above, the best launch site is currently CCAF5 SLC 40, where most of recent launches were successful;
- In second place was VAFB SLC 4E and third place was KSC LC 39A;
- General success rate improved over time.

Payload vs. Launch Site



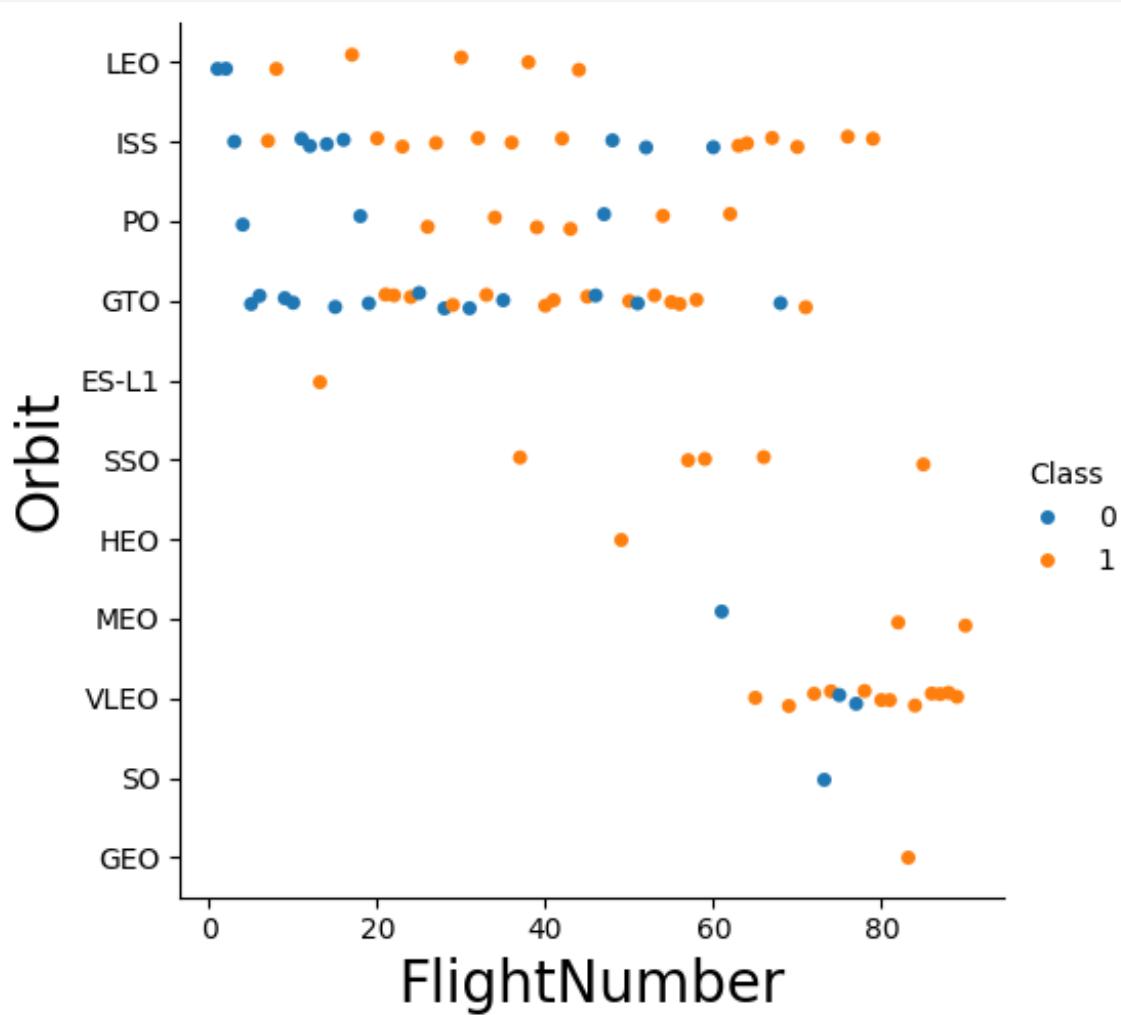
- Payloads over 9,000kg (about the weight of a school bus) have excellent success rate;
- Payloads over 12,000kg seems to be possible only on CCAFS SLC 40 and KSC LC 39A launch sites.

Success Rate vs. Orbit Type



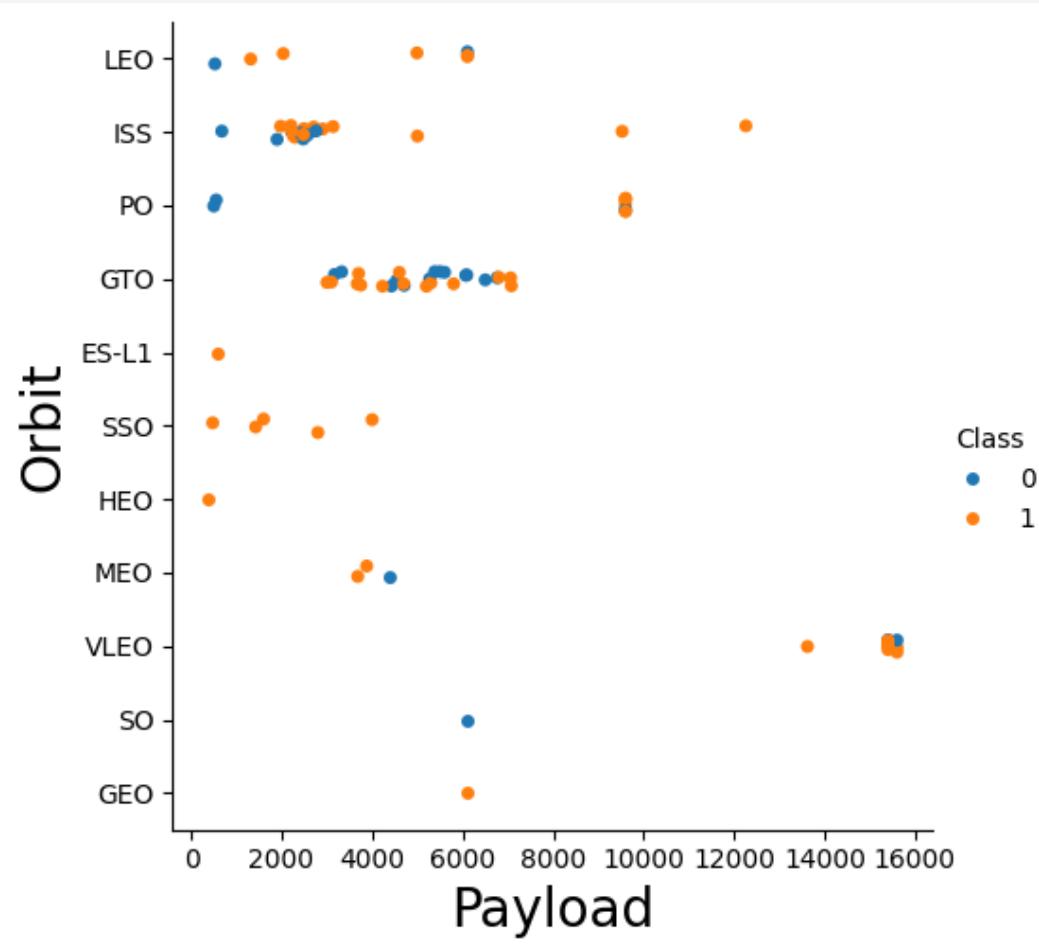
- The orbits with a success rate of 100% are: ES-L1, GEO, HEO and SSO
- The orbit SO has a 0% success rate

Flight Number vs. Orbit Type



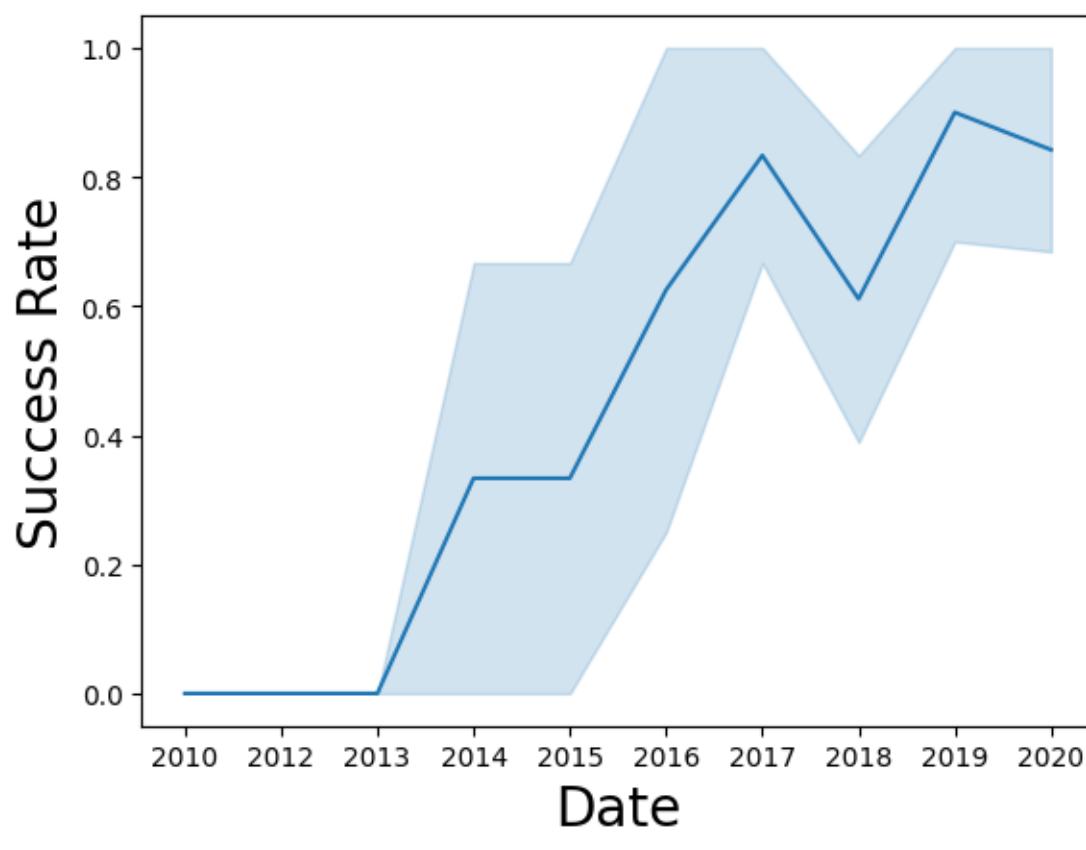
- In most orbits, the greater the number of flights the greater the chance for them to be successful
- GTO seems to be the exception and is both successful and unsuccessful across almost all number of flights

Payload vs. Orbit Type



- Apparently, there is no relation between payload and success rate to orbit GTO;
- ISS orbit has the widest range of payload and a good rate of success;
- There are few launches to the orbits ES-L1, HEO and GEO.

Launch Success Yearly Trend



- The success rate has been increasing since 2013 until 2020, with a minor decrease in 2018 but still following the increase trend

All Launch Site Names

```
%>sql
SELECT DISTINCT LAUNCH_SITE
FROM SPACEXTBL;

* sqlite:///my\_data1.db
Done.



| Launch_Site  |
|--------------|
| CCAFS LC-40  |
| VAFB SLC-4E  |
| KSC LC-39A   |
| CCAFS SLC-40 |


```

Launch Site Names Begin with 'CCA'

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

```
* sqlite:///my\_data1.db
Done.
```

Launch_Site
CCAFS LC-40

Total Payload Mass

```
%>sql
SELECT SUM(PAYLOAD_MASS__KG_)
FROM SPACEXTBL
WHERE Customer = 'NASA (CRS)';

* sqlite:///my\_data1.db
Done.

SUM(PAYLOAD_MASS__KG_)
45596
+ Co
```

Average Payload Mass by F9 v1.1

```
%%sql
SELECT AVG(PAYLOAD_MASS__KG_)
FROM SPACEXTBL
WHERE Booster_Version LIKE 'F9 v1.0%';

* sqlite:///my\_data1.db
Done.

AVG(PAYLOAD_MASS__KG_)
340.4
```

First Successful Ground Landing Date

```
%%sql
SELECT MIN(Date)
FROM SPACEXTBL
WHERE Landing_Outcome = 'Success (ground pad)';

* sqlite:///my\_data1.db
Done.

MIN(Date)
2015-12-22
```

Successful Drone Ship Landing with Payload between 4000 and 6000

```
%%sql
SELECT BOOSTER_VERSION
FROM SPACEXTBL
WHERE LANDING_OUTCOME = 'Success (drone ship)'
    AND 4000 < PAYLOAD_MASS__KG_ < 6000;

* sqlite:///my\_data1.db
Done.



| Booster_Version |
|-----------------|
| F9 FT B1021.1   |
| F9 FT B1022     |
| F9 FT B1023.1   |
| F9 FT B1026     |
| F9 FT B1029.1   |
| F9 FT B1021.2   |
| F9 FT B1029.2   |
| F9 FT B1036.1   |
| F9 FT B1038.1   |
| F9 B4 B1041.1   |
| F9 FT B1031.2   |
| F9 B4 B1042.1   |
| F9 B4 B1045.1   |
| F9 B5 B1046.1   |


```

Total Number of Successful and Failure Mission Outcomes

```
%%sql
SELECT MISSION_OUTCOME, COUNT(MISSION_OUTCOME) AS TOTAL_NUMBER
FROM SPACEXTBL
GROUP BY MISSION_OUTCOME;
```

```
* sqlite:///my_data1.db
Done.
```

Mission_Outcome	TOTAL_NUMBER
Failure (in flight)	1
Success	98
Success	1
Success (payload status unclear)	1

Boosters Carried Maximum Payload

```
%%sql
SELECT DISTINCT BOOSTER_VERSION
FROM SPACEXTBL
WHERE PAYLOAD_MASS__KG_ = (
    SELECT MAX(PAYLOAD_MASS__KG_)
    FROM SPACEXTBL);

* sqlite:///my_data1.db
Done.



| Booster_Version |
|-----------------|
| 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

- List the failed landing_outcomes in drone ship, their booster versions, and launch site names for in year 2015
- Present your query result with a short explanation here

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

```
%%sql
SELECT LANDING_OUTCOME, COUNT(LANDING_OUTCOME) AS TOTAL_NUMBER
FROM SPACEXTBL
WHERE DATE BETWEEN '2010-06-04' AND '2017-03-20'
GROUP BY LANDING_OUTCOME
ORDER BY TOTAL_NUMBER DESC
```

```
* sqlite:///my\_data1.db
Done.
```

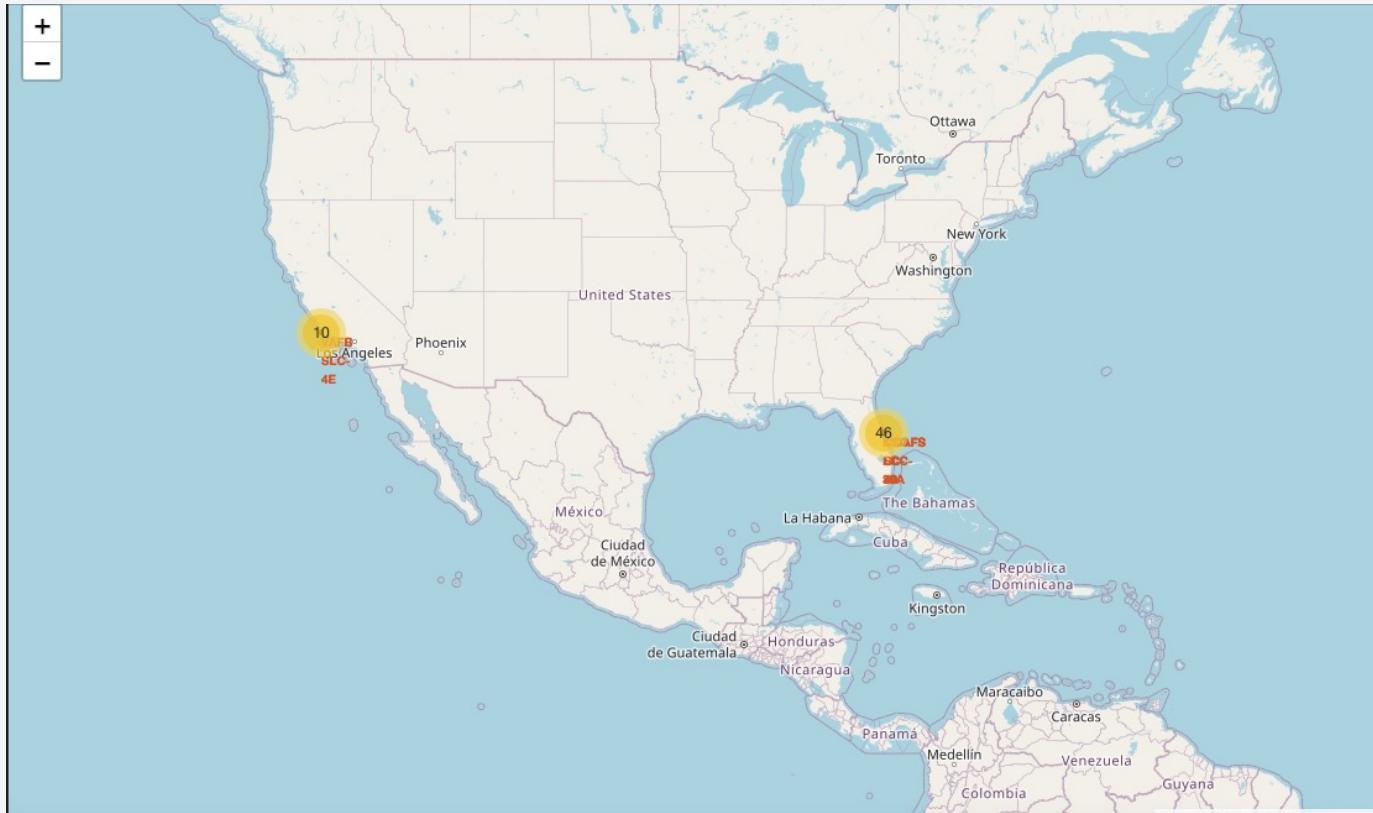
Landing_Outcome	TOTAL_NUMBER
No attempt	10
Success (drone ship)	5
Failure (drone ship)	5
Success (ground pad)	3
Controlled (ocean)	3
Uncontrolled (ocean)	2
Failure (parachute)	2
Precluded (drone ship)	1

The background of the slide is a photograph taken from space at night. It shows the curvature of the Earth against a dark blue-black 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 right, the green and yellow glow of the aurora borealis is visible. The atmosphere of the Earth is thin and hazy, appearing as a light blue band near the horizon.

Section 3

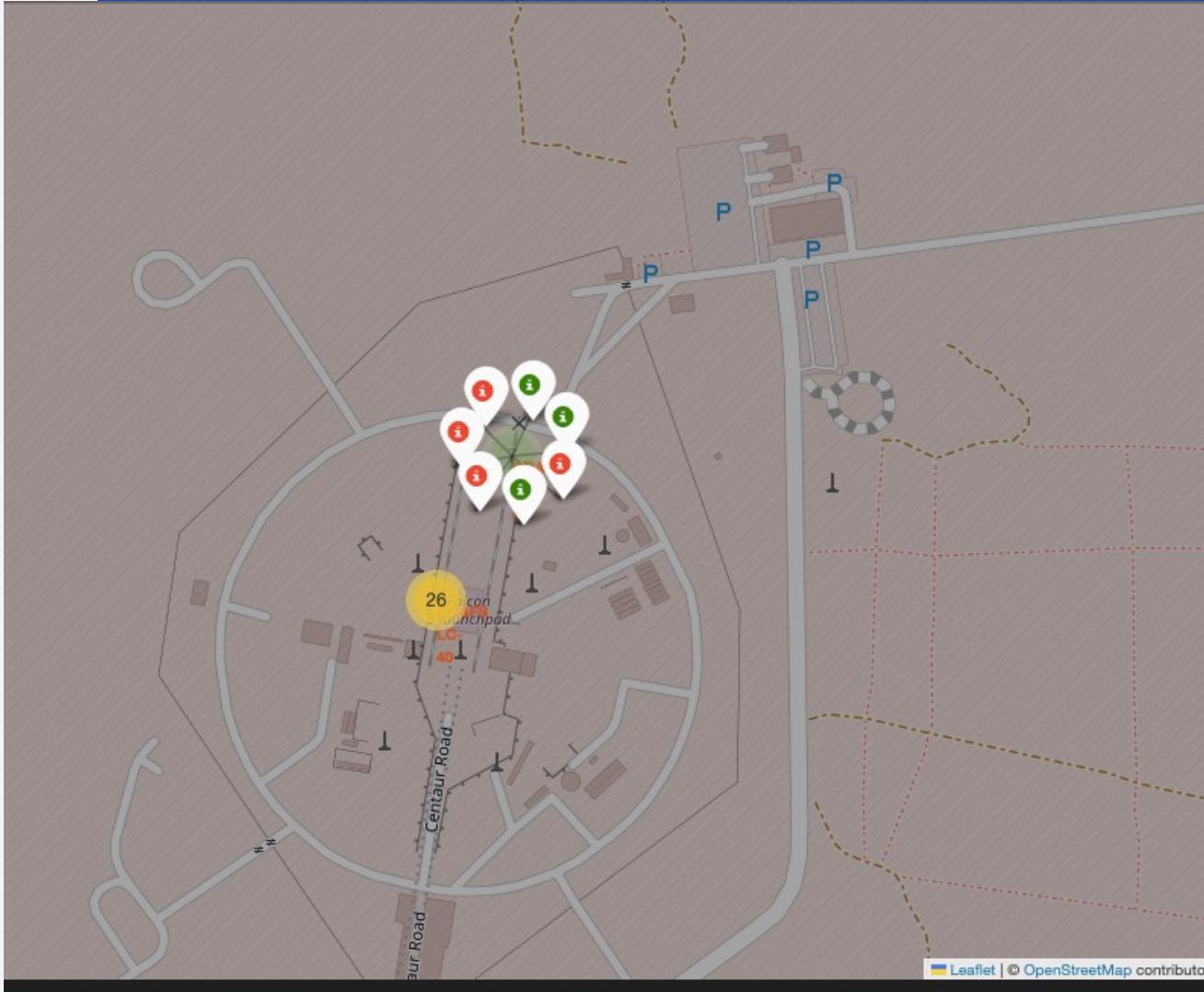
Launch Sites Proximities Analysis

Launch sites locations



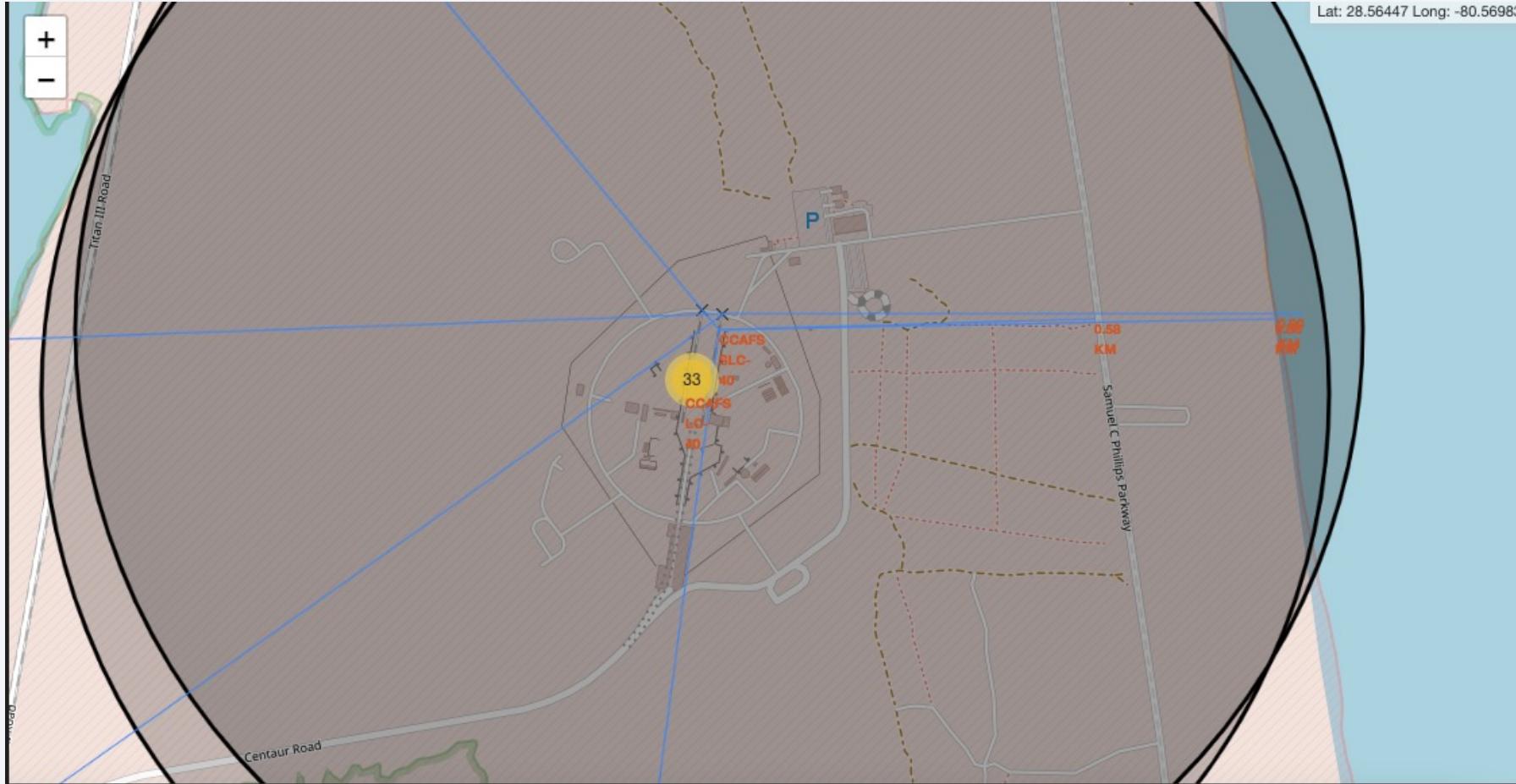
- All launch sites are in very proximity to the coast

Color-labeled markers



- A green marker indicates a successful launch, and a red marker indicates an unsuccessful one

Distance to proximity sites

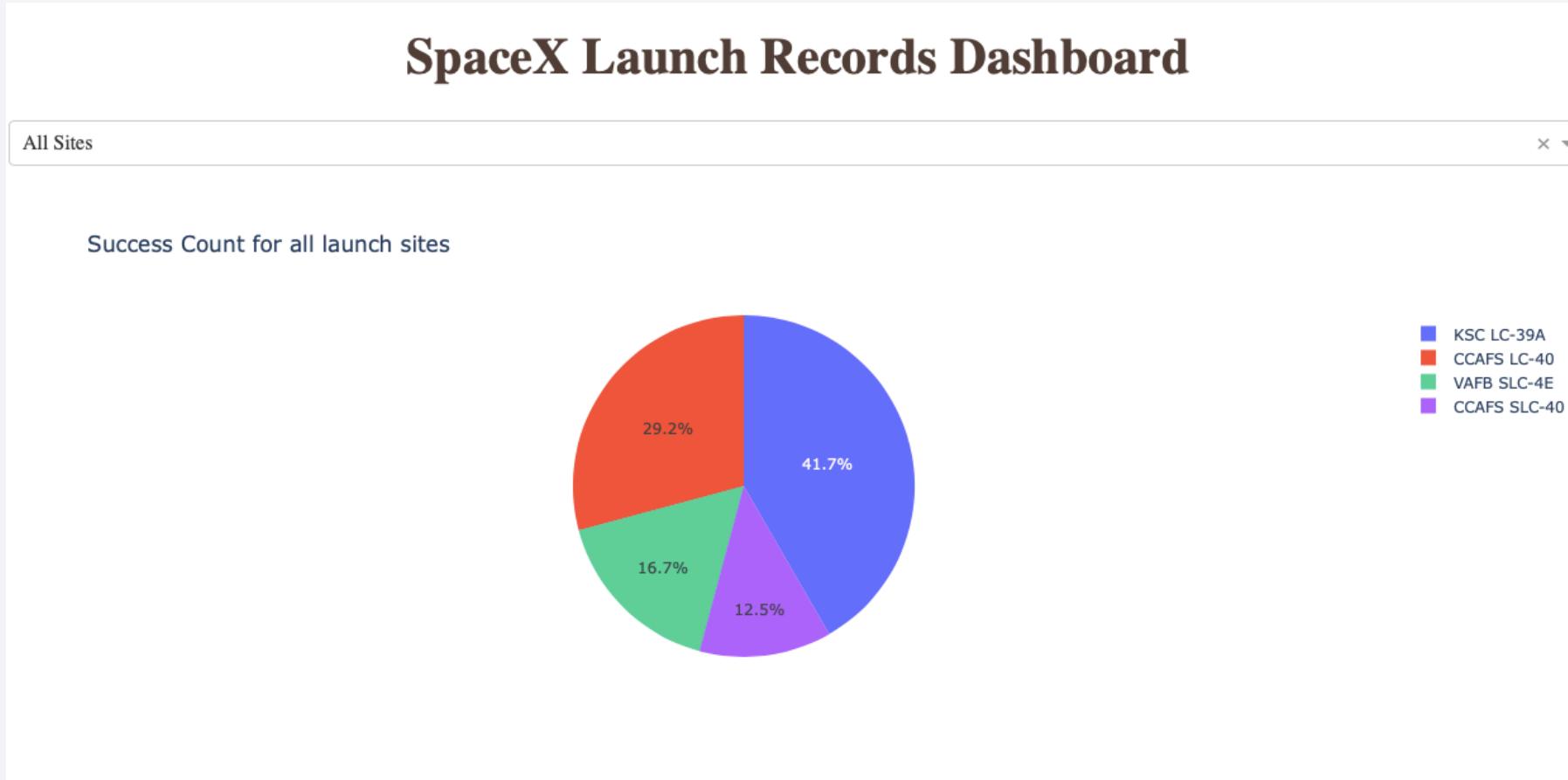


Section 4

Build a Dashboard with Plotly Dash

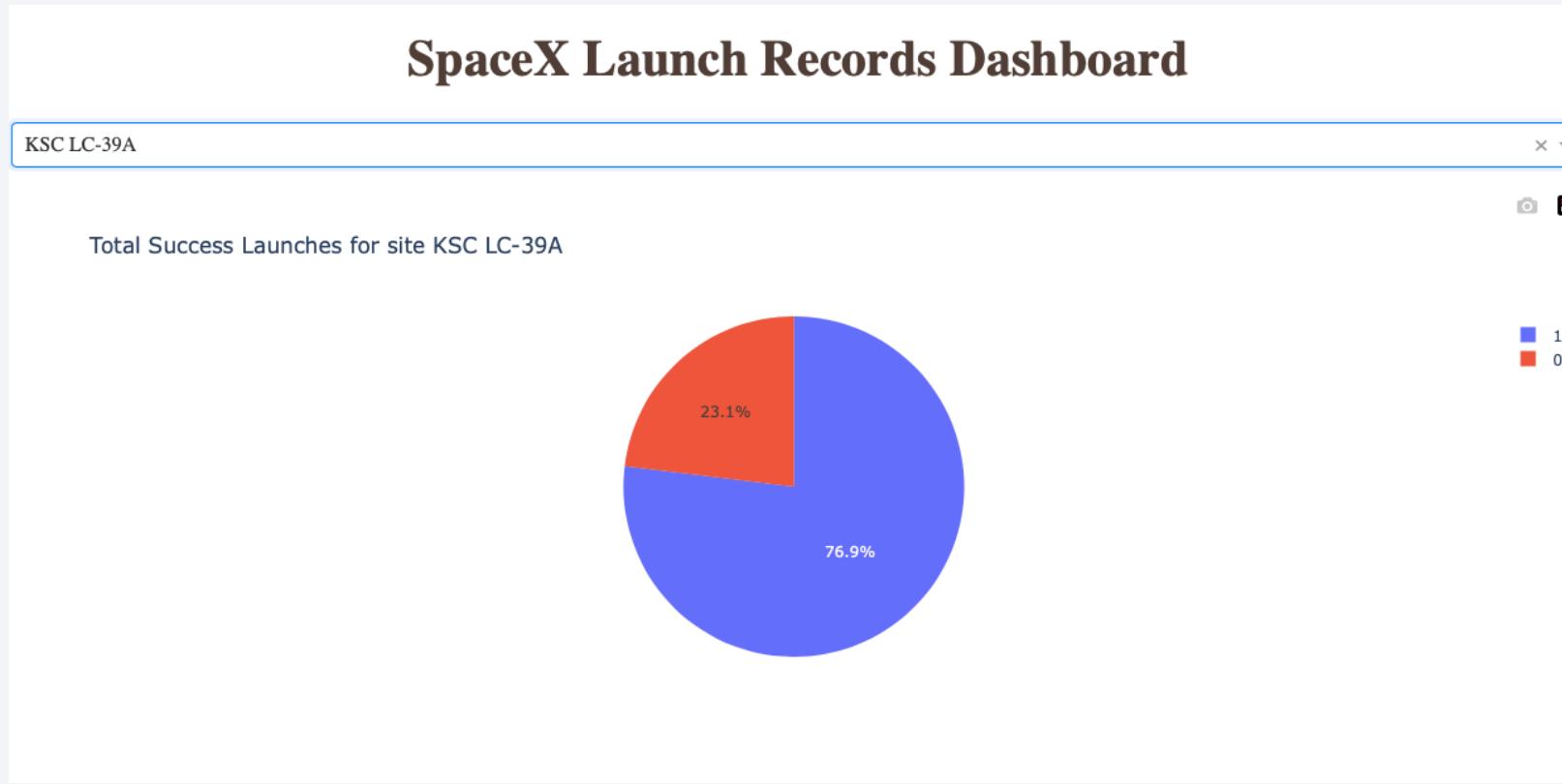


Launch success count for all sites



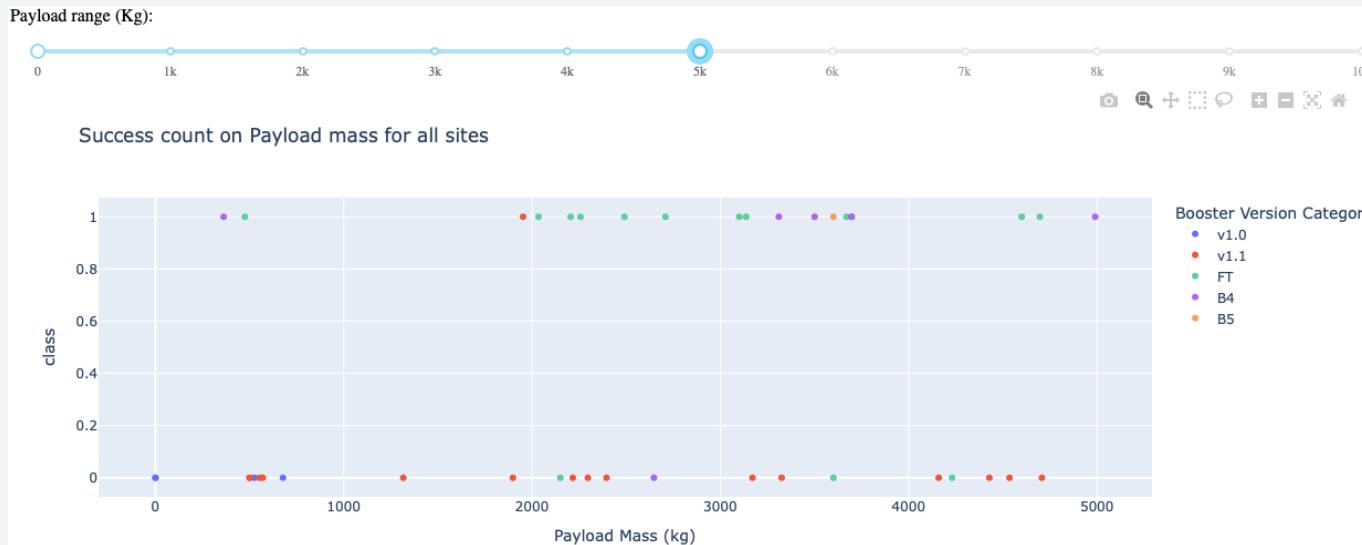
- The site KSC LC-39A is the most successful one

Launch success for the site KSC LC-39A

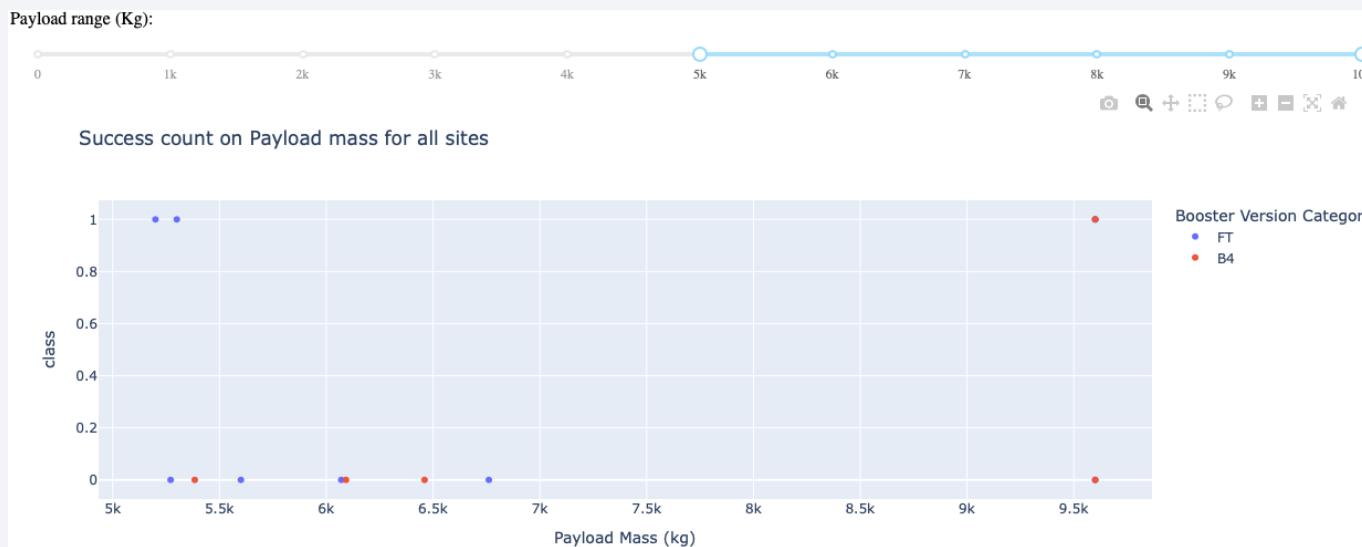


- It has a success rate of 76.9%

Payload mass vs. Launch outcome



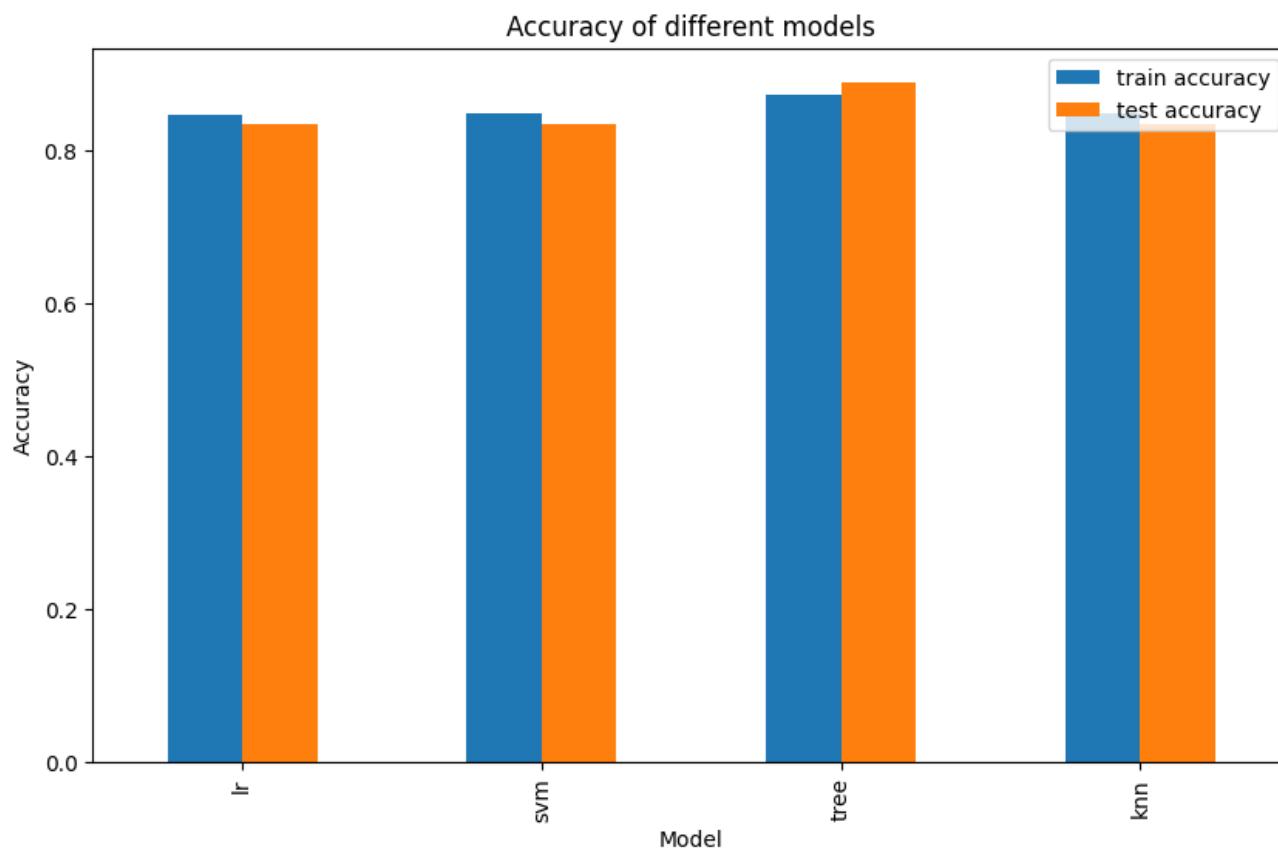
- Payloads between 2000kg and 5000kg show greater success than bigger payloads



Section 5

Predictive Analysis (Classification)

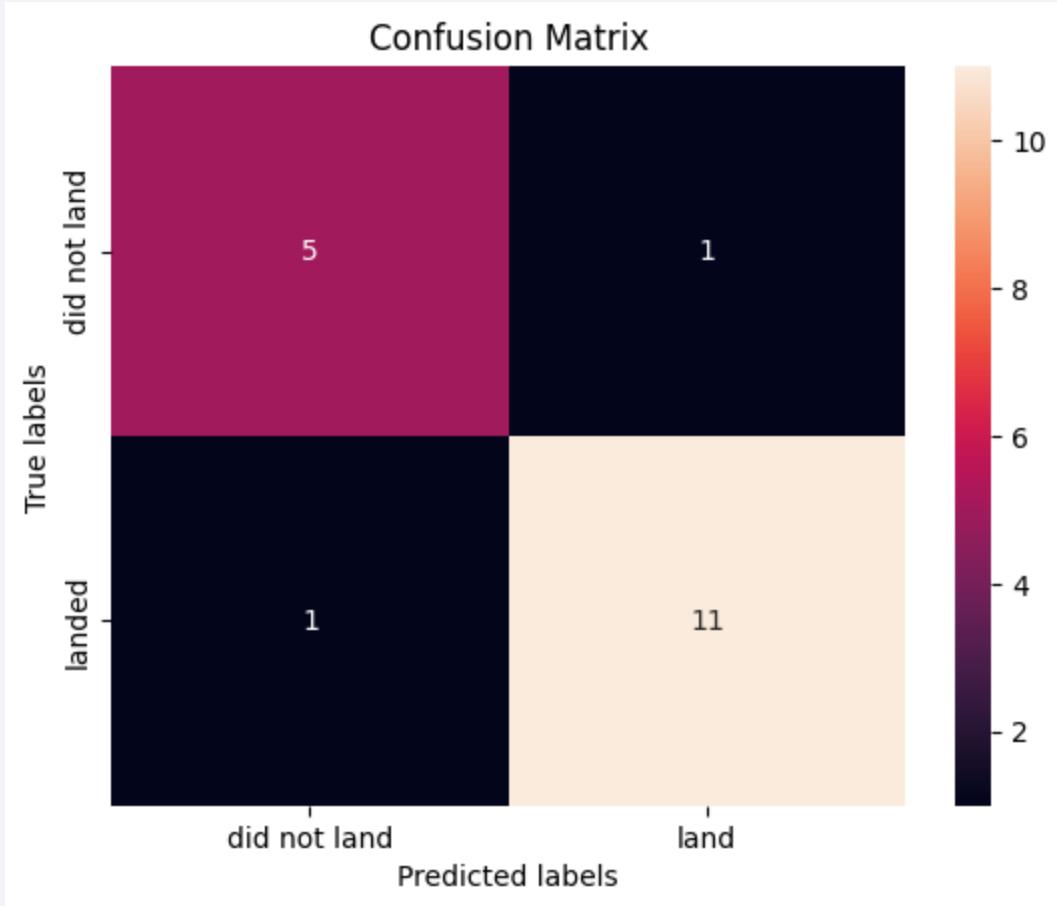
Classification Accuracy



- The model with the best accuracy is the decision tree classifier.

	train accuracy	test accuracy
lr	0.846429	0.833333
svm	0.848214	0.833333
tree	0.871429	0.888889
knn	0.848214	0.833333

Confusion Matrix



- Confusion matrix of Decision Tree Classifier proves its accuracy by showing the big numbers of true positive and true negative compared to the false ones.

Conclusions

- Different data sources were analysed, refining conclusions along the process;
- The best launch site is KSC LC-39A;
- Launches above 7,000kg are less risky;
- Although most of mission outcomes are successful, successful landing outcomes seem to improve over time, according the evolution of processes and rockets;
- Decision Tree Classifier can be used to predict successful landings and increase profits.

Appendix

- Include any relevant assets like Python code snippets, SQL queries, charts, Notebook outputs, or data sets that you may have created during this project

Thank you!

