**IBM Developer**
SKILLS NETWORK

# Winning Space Race with Data Science

Vilcacundo Alex
14th April 2023

# 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

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

  - 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

    - How to build, tune, evaluate classification models

# Data Collection

The data was collected using various methods

-Data collection was done using get request to the SpaceX API.

-Next, we decoded the response content as a Json using .json() function call and turn it into a pandas dataframe using .json_normalize().

-We then cleaned the data, checked for missing values and fill in missing values where necessary.

-In addition, we performed web scraping from Wikipedia for Falcon 9 launch records with BeautifulSoup.

-The objective was to extract the launch records as HTML table, parse the table and convert it to a pandas dataframe for future analysis.

# Data Collection – SpaceX API

- We used the get request to the SpaceX API to collect data, clean the requested data and did some basic data wrangling and formatting.

- The link to the notebook is https://github.com/alex-v505/IBM-DataScience-SpaceX-Capstone/blob/main/Spacex-data-collection-api.ipynb

# Data Collection - Scraping

- We applied web scrapping to webscrap Falcon 9 launch records with BeautifulSoup

- We parsed the table and converted it into a pandas dataframe.

- The link to the notebook is https://github.com/alex-v505/IBM-DataScience-SpaceX-Capstone/blob/main/Spacex_webscraping.ipynb

```
In [5]: static_url = "https://en.wikipedia.org/w/index.php?title=List_of_Falcon_9_and_Falco
```

```
In [13]: # use requests.get() method with the provided static_url
response = requests.get(static_url).text
# assign the response to a object
```

Create a BeautifulSoup object from the HTML response

```
In [16]: # Use BeautifulSoup() to create a BeautifulSoup object from a response text conten
soup = BeautifulSoup(response, 'html.parser')
```

Print the page title to verify if the BeautifulSoup object was created properly

```
In [17]: # Use soup.title attribute
soup.title
```
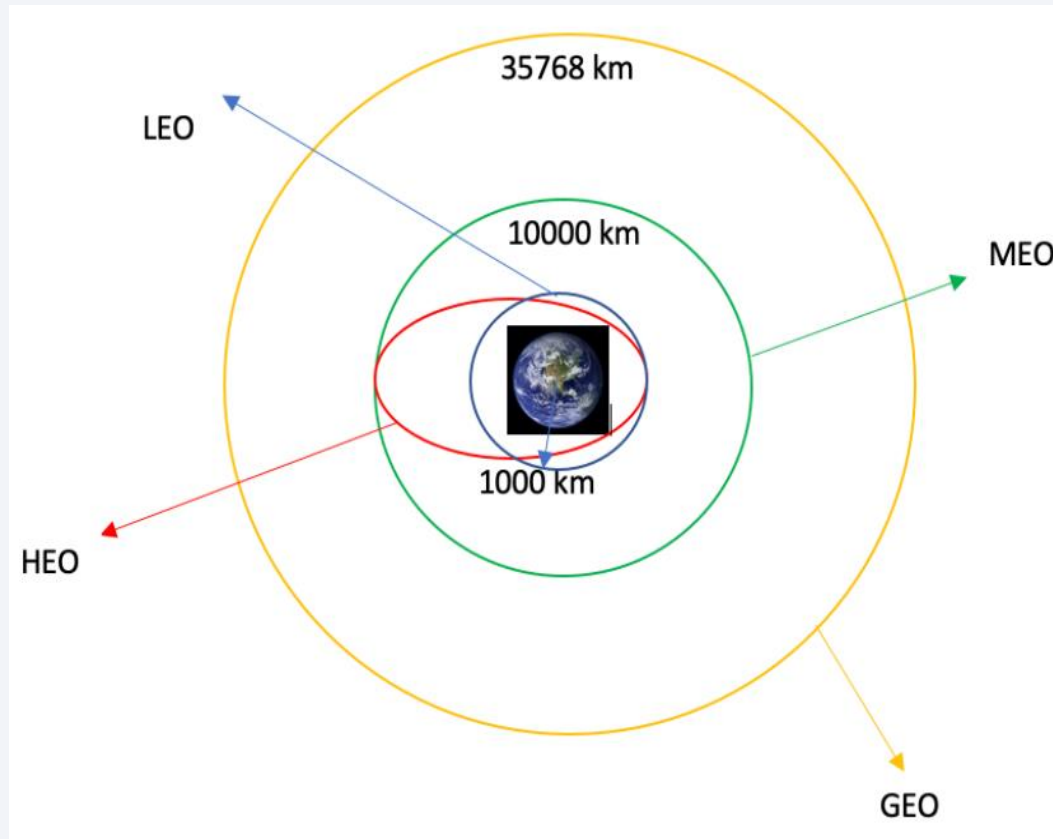
```
Out[17]: <title>List of Falcon 9 and Falcon Heavy launches - Wikipedia</title>
```

```
In [18]: # Use the find_all function in the BeautifulSoup object, with element type `table`
# Assign the result to a list called `html_tables`
html_tables = soup.find_all('table')
html_tables
```

```
In [20]: column_names = []

# Apply find_all() function with `th` element on first_launch_table
# Iterate each th element and apply the provided extract_column_from_header() to g
# Append the Non-empty column name (`if name is not None and len(name) > 0`) into
temp = soup.find_all('th')
for x in range(len(temp)):
    try:
        name = extract_column_from_header(temp[x])
        if(name is not None and len(name)>0):
            column_names.append(name)
    except:
        pass
```
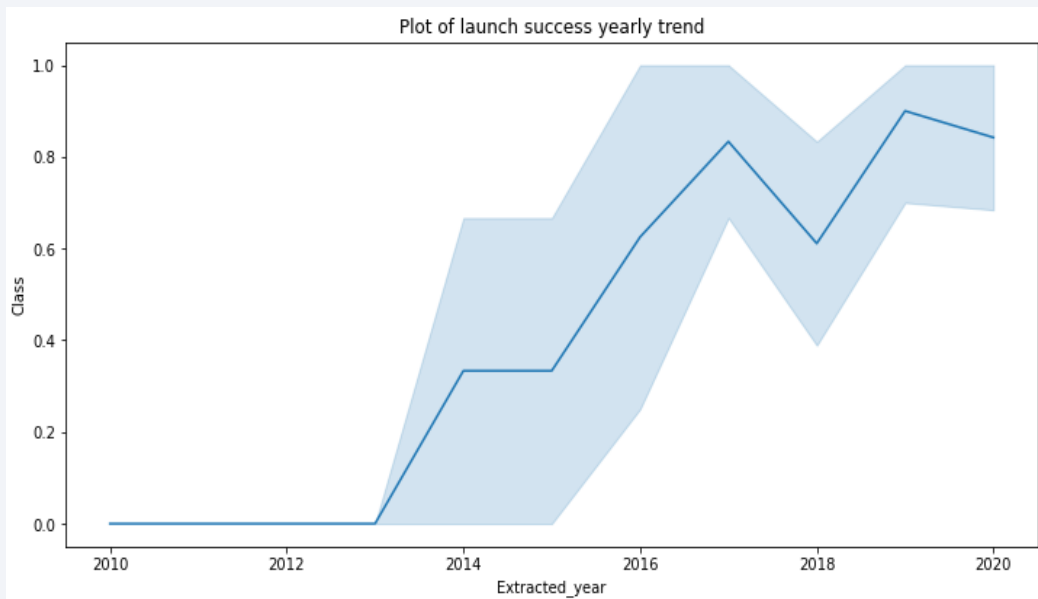
# Data Wrangling



- We performed exploratory data analysis and determined the training labels.
- We calculated the number of launches at each site, and the number and occurrence of each orbits
- We created landing outcome label from outcome column and exported the results to csv.
- The link to the notebook is https://github.com/alex-v505/IBM-DataScience-SpaceX-Capstone/blob/main/Spacex-data%20wrangling.ipynb

# 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.



The link to the notebook is
https://github.com/alex-v505/IBM-DataScience-SpaceX-Capstone/blob/main/Spacex-eda-dataviz.ipynb

# EDA with SQL

- We loaded the SpaceX dataset into a IBM DB2 database without leaving the jupyter notebook.

- We applied EDA with SQL to get insight from the data. We wrote queries to find out for instance:

  - 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.

- The link to the notebook is https://github.com/alex-v505/IBM-DataScience-SpaceX-Capstone/blob/main/Spacex-eda-sql.ipynb

# 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.

- The link to the notebook is https://github.com/alex-v505/IBM-DataScience-SpaceX-Capstone/blob/main/SpaceX_launch_site_location.ipynb

# 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.

- The link to the notebook is https://github.com/alex-v505/IBM-DataScience-SpaceX-Capstone/blob/main/spacex_dash_app.py

# 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.

- The link to the notebook is https://github.com/alex-v505/IBM-DataScience-SpaceX-Capstone/blob/main/SpaceX_Machine%20Learning%20Prediction_Part_5.ipynb

# Results

- Exploratory data analysis results

- Interactive analytics demo in screenshots
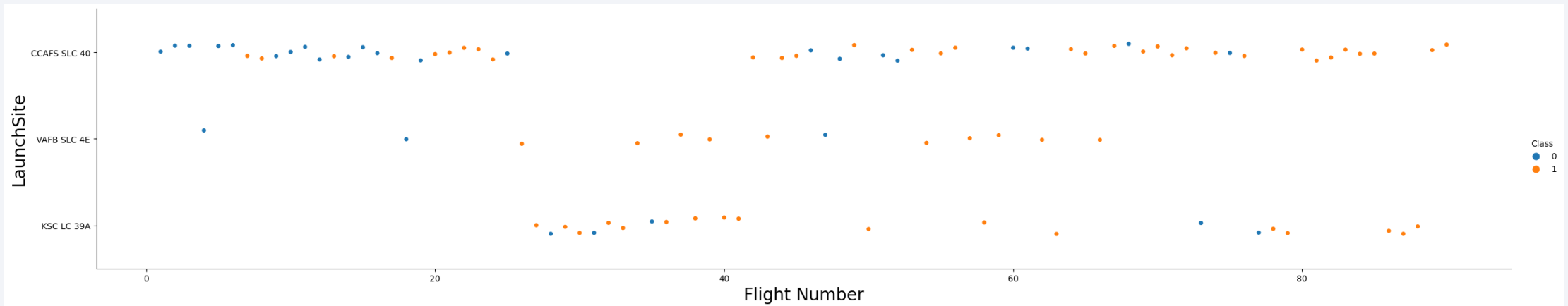
- Predictive analysis results
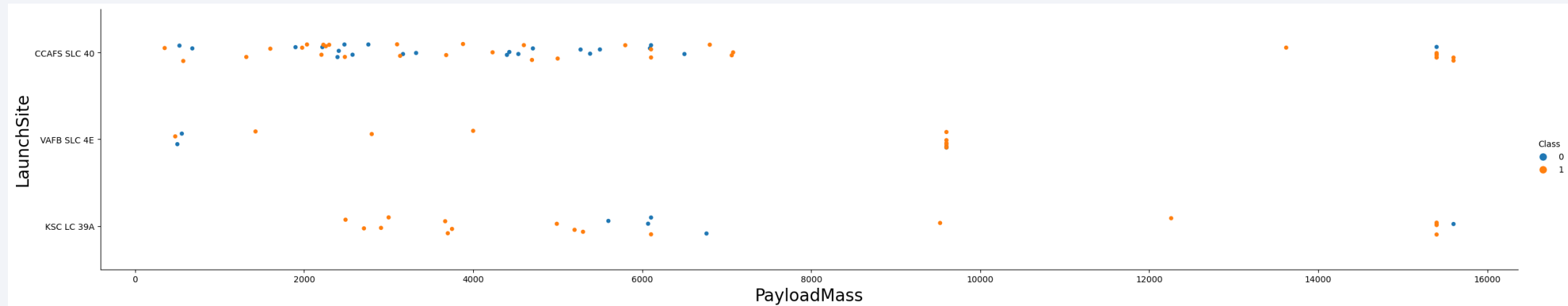
Section 2

# 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 a launch site.
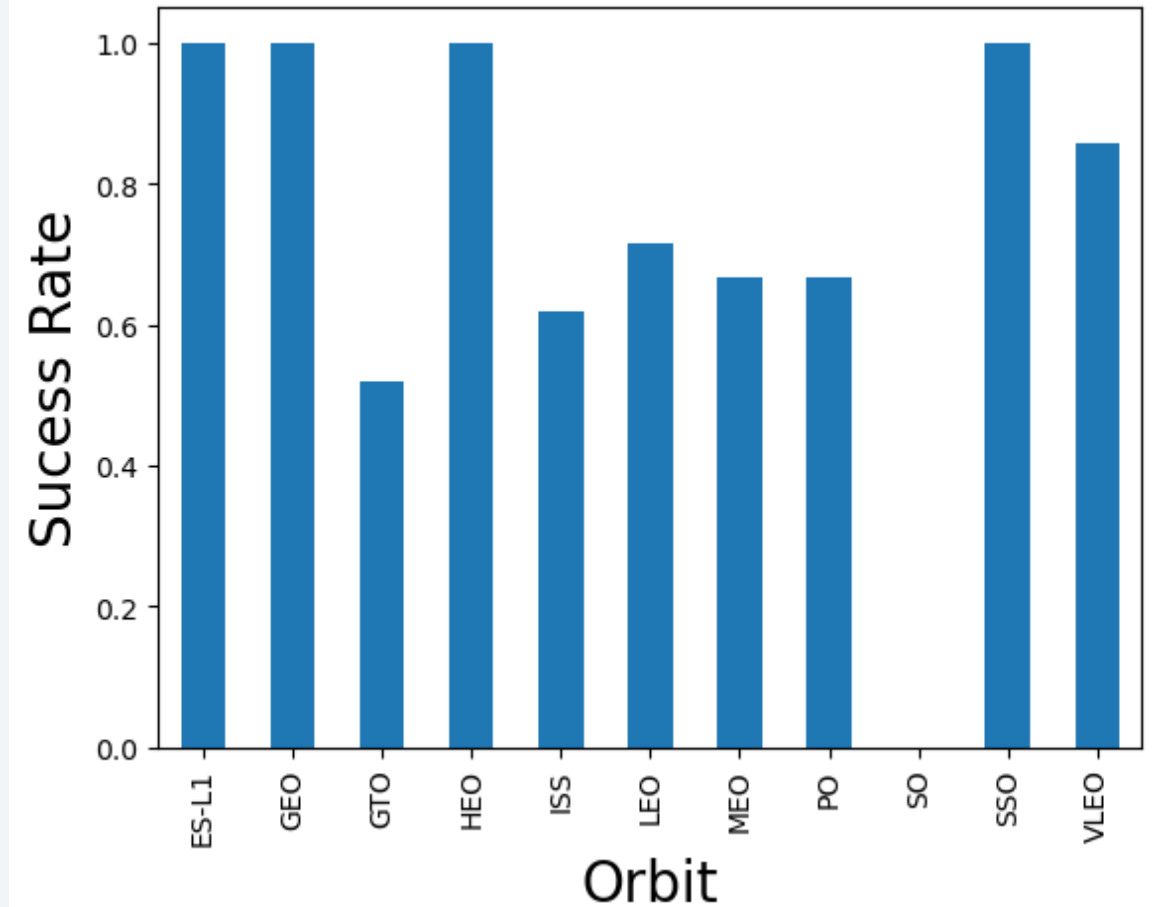
# Payload vs. Launch Site

- From the plot, we found that the VAFB-SLC launchsite there are no rockets launched for heavypayload mass(greater than 10000).
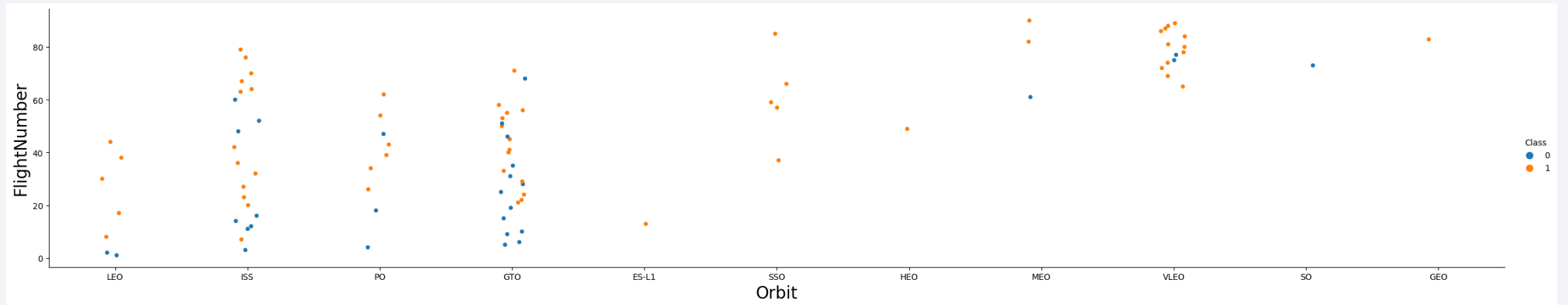
# Success Rate vs. Orbit Type

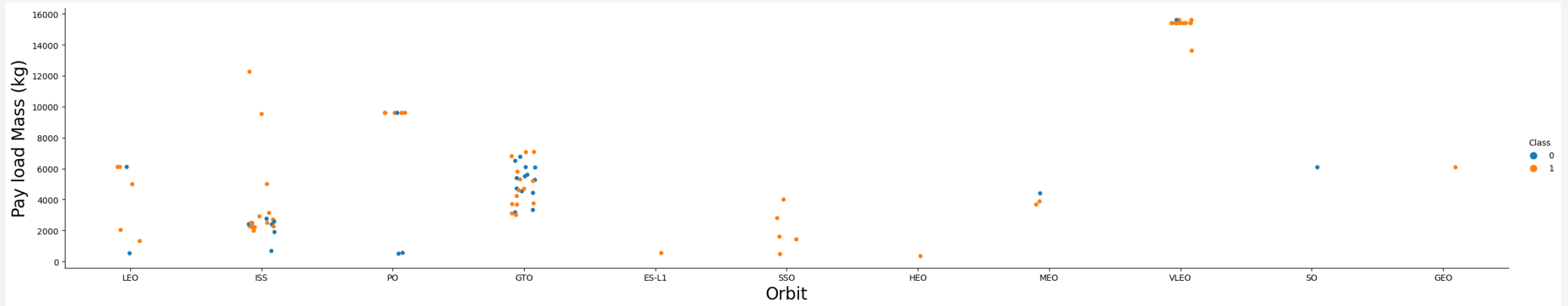- From the plot, we can see that ES-L1, GEO, HEO, SSO, VLEO had the most success rate.

# 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.
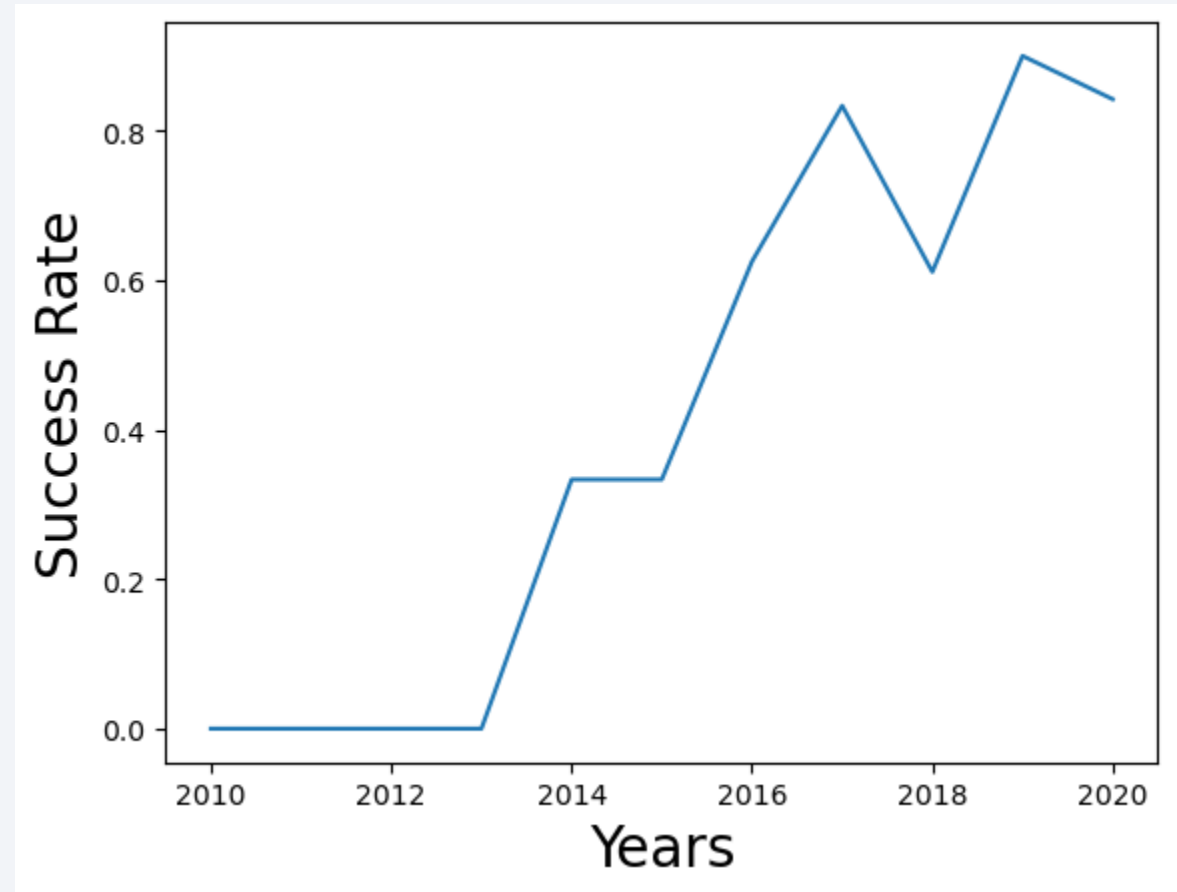
# Payload vs. Orbit Type

- 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

- We used the key word **DISTINCT** to show only unique launch sites from the SpaceX data.



```
In [7]:  %sql SELECT DISTINCT LAUNCH_SITE FROM SPACEX

          * ibm_db_sa://hsv46674:***@b0aebb68-94fa-46ec-a1fc-1c999edb6187.c3n41cmd0nqnrk39u
         98g.databases.appdomain.cloud:31249/BLUDB
         Done.

Out[7]:     launch_site

          CCAFS LC-40

          CCAFS SLC-40

          KSC LC-39A

          VAFB SLC-4E
```

# Launch Site Names Begin with 'CCA'

```
In [9]: %sql SELECT * FROM SPACEX WHERE LAUNCH_SITE LIKE 'CCA%' LIMIT 5
```

* ibm_db_sa://hsv46674:***@b0aebb68-94fa-46ec-a1fc-1c999edb6187.c3n41cmd0nqnrk39u98g.databases.appdomain.cloud:31249/BLUDB
Done.

Out[9]:

| DATE | time_utc | booster_version | launch_site | payload | payload_mass_kg_ | orbit | customer | mission_outcome | landing_outcome |
|---|---|---|---|---|---|---|---|---|---|
| 2010-06-04 | 18:45:00 | F9 v1.0 B0003 | CCAFS LC-40 | Dragon Spacecraft Qualification Unit | 0 | LEO | SpaceX | Success | Failure (parachute) |
| 2010-12-08 | 15:43:00 | F9 v1.0 B0004 | CCAFS LC-40 | Dragon demo flight C1, two CubeSats, barrel of Brouere cheese | 0 | LEO (ISS) | NASA (COTS) NRO | Success | Failure (parachute) |
| 2012-05-22 | 7:44:00 | F9 v1.0 B0005 | CCAFS LC-40 | Dragon demo flight C2 | 525 | LEO (ISS) | NASA (COTS) | Success | No attempt |
| 2012-10-08 | 0:35:00 | F9 v1.0 B0006 | CCAFS LC-40 | SpaceX CRS-1 | 500 | LEO (ISS) | NASA (CRS) | Success | No attempt |
| 2013-03-01 | 15:10:00 | F9 v1.0 B0007 | CCAFS LC-40 | SpaceX CRS-2 | 677 | LEO (ISS) | NASA (CRS) | Success | No attempt |

- We used the query above to display 5 records where launch sites begin with `CCA`

# Total Payload Mass

- We calculated the total payload carried by boosters from NASA as 45596 using the query below

```
In [16]: %sql SELECT sum(PAYLOAD_MASS_KG_) as "TOTAL_PAYLOAD_MASS_FROM_NASA_CRS" FROM SPACEX WHERE CUSTOMER = 'NASA (CRS)'
 * ibm_db_sa://hsv46674:***@b0aebb68-94fa-46ec-a1fc-1c999edb6187.c3n41cmd0nqnrk39u98g.databases.appdomain.cloud:31249/BLUDB
Done.
Out[16]: total_payload_mass_from_nasa_crs
                 45596
```

# Average Payload Mass by F9 v1.1

- We calculated the average payload mass carried by booster version F9 v1.1 as 2534

```
In [10]:  %sql SELECT AVG(PAYLOAD_MASS_KG_) as "AVERAGE_PAYLOAD_MASS_BOOSTER_F9 v1.1" FROM SPACEX WHERE BOOSTER_VERSION LIKE 'F9 v1.1%'
           * ibm_db_sa://hsv46674:***@b0aebb68-94fa-46ec-a1fc-1c999edb6187.c3n41cmd0nqnrk39u98g.databases.appdomain.cloud:31249/BLUDB
          Done.

Out[10]:  AVERAGE_PAYLOAD_MASS_BOOSTER_F9 v1.1
                          2534
```

# First Successful Ground Landing Date

- We observed that the dates of the first successful landing outcome on ground pad was 22[nd] December 2015

```
In [28]: %sql SELECT MIN(DATE) FROM SPACEX WHERE LANDING_OUTCOME = 'Success (ground pad)';
         * ibm_db_sa://hsv46674:***@b0aebb68-94fa-46ec-a1fc-1c999edb6187.c3n41cmd0nqnrk39u98g.databases.appdomain.cloud:31249/BLUDB
         Done.
Out[28]:         1
         2015-12-22
```

# Successful Drone Ship Landing with Payload between 4000 and 6000

- We used the **WHERE** clause to filter for boosters which have successfully landed on drone ship and applied the **AND** condition to determine successful landing with payload mass greater than 4000 but less than 6000

```
In [29]: %sql SELECT BOOSTER_VERSION FROM SPACEX WHERE LANDING_OUTCOME = 'Success (drone ship)' and PAYLOAD_MASS_KG_ between 4000 and 6000;

         * ibm_db_sa://hsv46674:***@b0aebb68-94fa-46ec-a1fc-1c999edb6187.c3n41cmd0nqnrk39u98g.databases.appdomain.cloud:31249/BLUDB
         Done.

Out[29]:   booster_version
              F9 FT B1022
              F9 FT B1026
              F9 FT B1021.2
              F9 FT B1031.2
```

# Total Number of Successful and Failure Mission Outcomes

- We used **WHERE** MissionOutcome was a success, failure or success (payload status unclear).

```
In [33]: %sql SELECT MISSION_OUTCOME, COUNT(MISSION_OUTCOME) as "Count" FROM SPACEX GROUP BY MISSION_OUTCOME
          * ibm_db_sa://hsv46674:***@b0aebb68-94fa-46ec-a1fc-1c999edb6187.c3n41cmd0nqnrk39u98g.databases.appdomain.cloud:31249/BLUDB
         Done.
Out[33]:
```

| mission_outcome | Count |
|---|---|
| Failure (in flight) | 1 |
| Success | 99 |
| Success (payload status unclear) | 1 |

# Boosters Carried Maximum Payload

- We determined the booster that have carried the maximum payload using a subquery in the **WHERE** clause and the **MAX()** function.

# 2015 Launch Records

- We used a combinations of the **WHERE** clause, **LIKE**, **AND**, and **BETWEEN** conditions to filter for failed landing outcomes in drone ship, their booster versions, and launch site names for year 2015

```
In [42]: %sql SELECT DATE,LANDING_OUTCOME, MISSION_OUTCOME, BOOSTER_VERSION, LAUNCH_SITE FROM SPACEX WHERE LANDING_OUTCOME = 'Failure (drone ship)' and YEAR(DATE) = '2015'
```
 * ibm_db_sa://hsv46674:***@b0aebb68-94fa-46ec-a1fc-1c999edb6187.c3n41cmd0nqnrk39u98g.databases.appdomain.cloud:31249/BLUDB
Done.

Out[42]:

| DATE | landing_outcome | mission_outcome | booster_version | launch_site |
|---|---|---|---|---|
| 2015-01-10 | Failure (drone ship) | Success | F9 v1.1 B1012 | CCAFS LC-40 |
| 2015-04-14 | Failure (drone ship) | Success | F9 v1.1 B1015 | CCAFS LC-40 |

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

- We selected Landing outcomes and the **COUNT** of landing outcomes from the data and used the **WHERE** clause to filter for landing outcomes **BETWEEN** 2010-06-04 to 2010-03-20.

- We applied the **GROUP BY** clause to group the landing outcomes and the **ORDER BY** clause to order the grouped landing outcome in descending order.

```
In [12]: %%sql
SELECT LANDING_OUTCOME, COUNT(*) as "TOTAL_LANDING_OUTCOMES" FROM SPACEX WHERE DATE BETWEEN '2010-06-04'
AND '2017-03-20' GROUP BY LANDING_OUTCOME ORDER BY TOTAL_LANDING_OUTCOMES DESC;
```

 * ibm_db_sa://hsv46674:***@b0aebb68-94fa-46ec-a1fc-1c999edb6187.c3n41cmd0nqnrk39u98g.databases.appdomain.cloud:31249/BLUDB
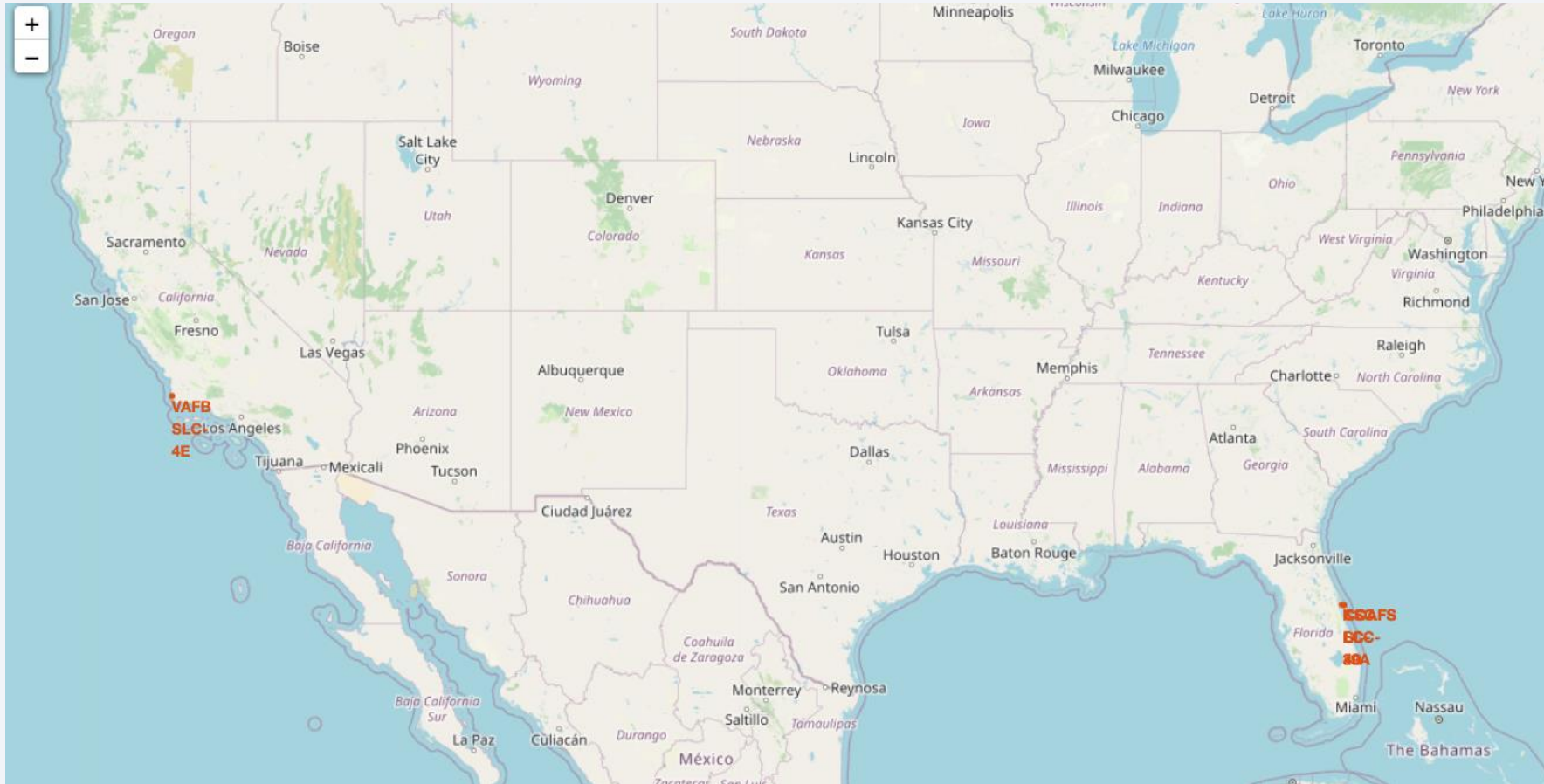Done.

Out[12]:

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

# Launch Sites Proximities Analysis

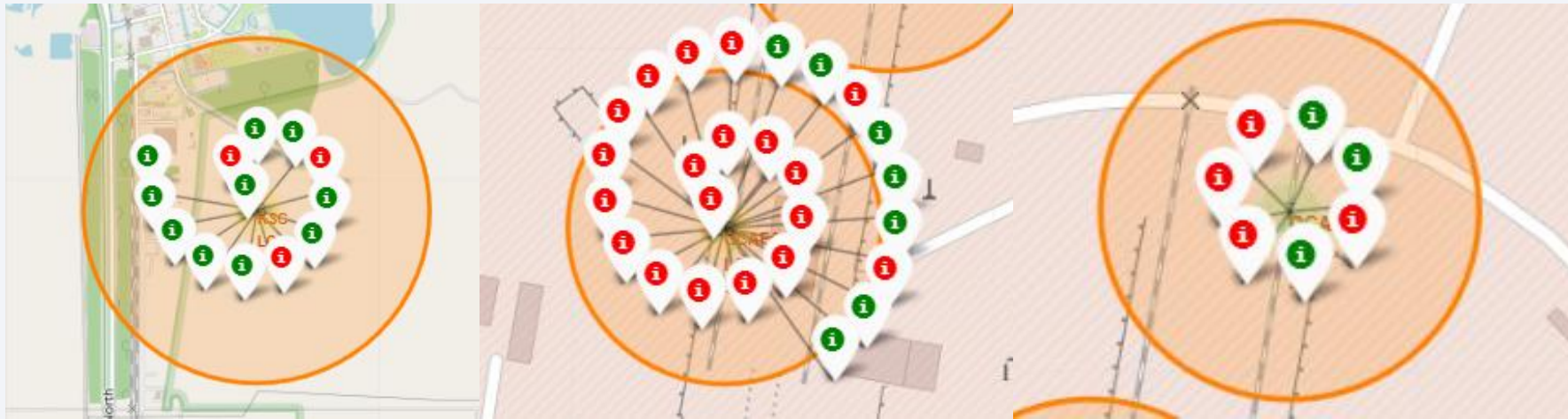# All launch sites global map markers



- We can see that the SpaceX launch sites are in the United States of America coasts. Florida and California

35

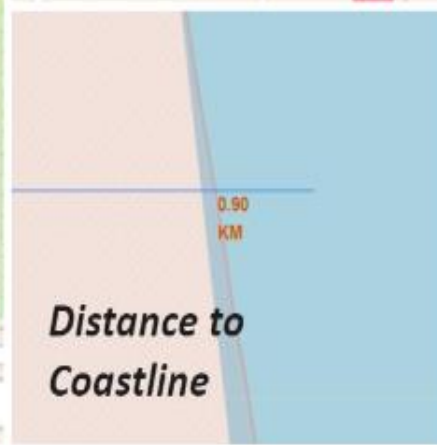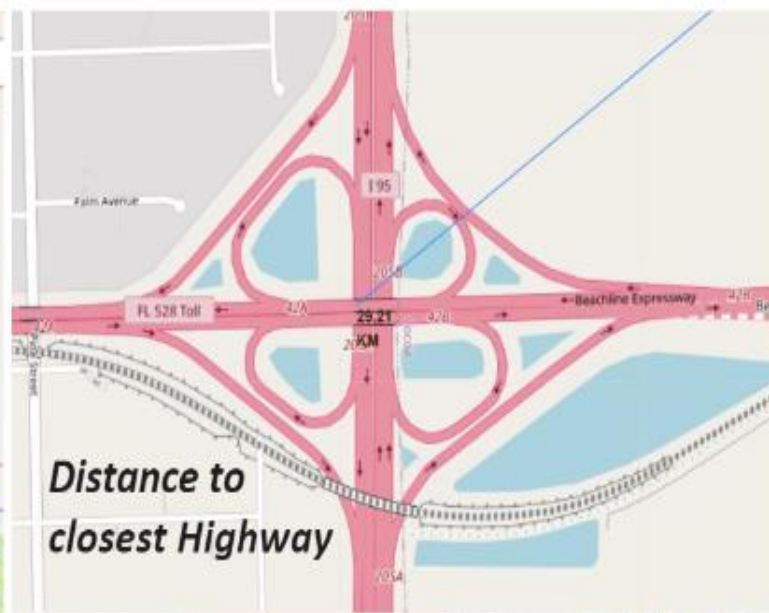# Markers showing launch sites with color labels



California Launch Site

Florida Launch Site

- Green Marker shows successful launches and Red Marker shows failures

36

# Launch Site distance to landmarks



Distance to Railway Station

Distance to closest Highway

Distance to coast

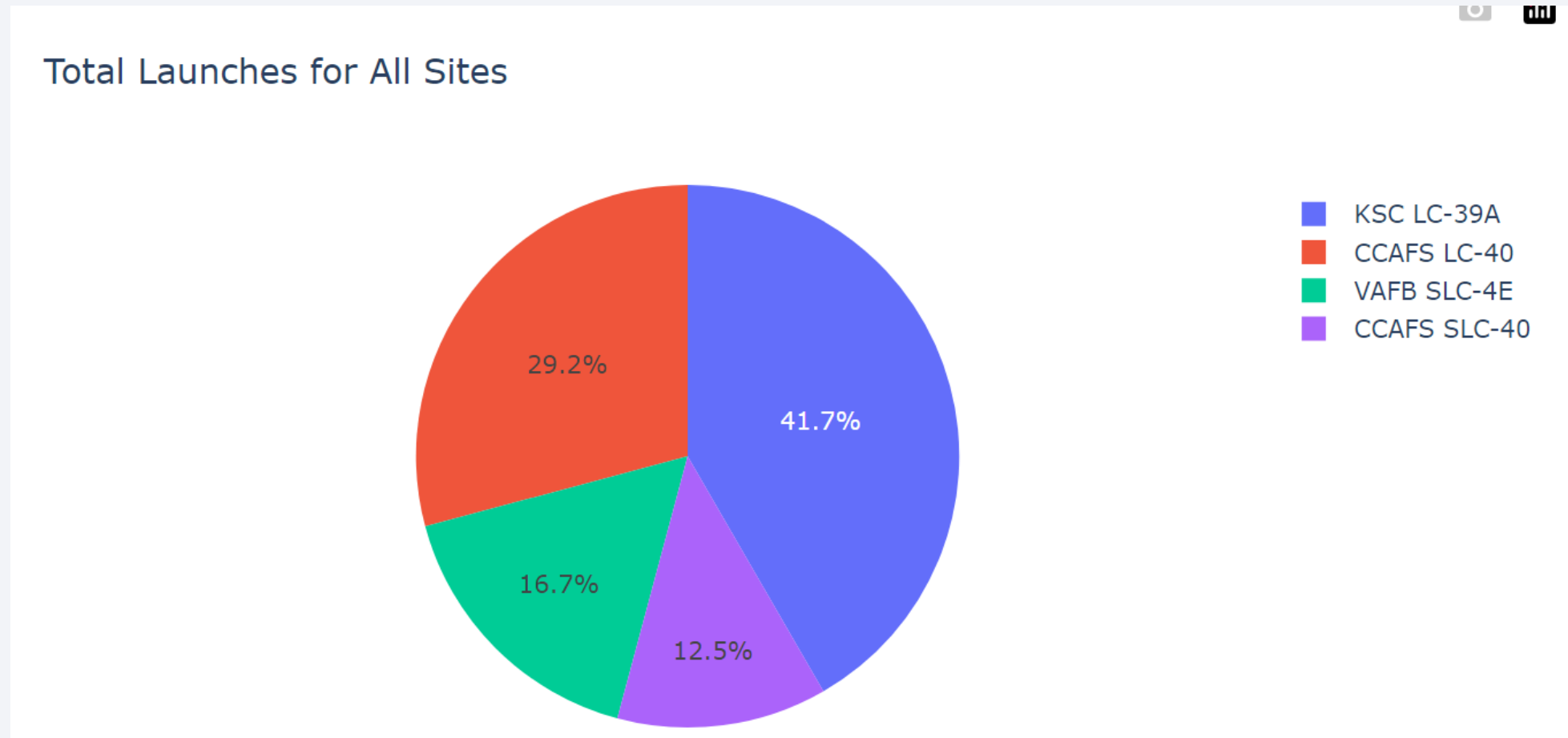Distance to Coastline

Distance to City

- 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
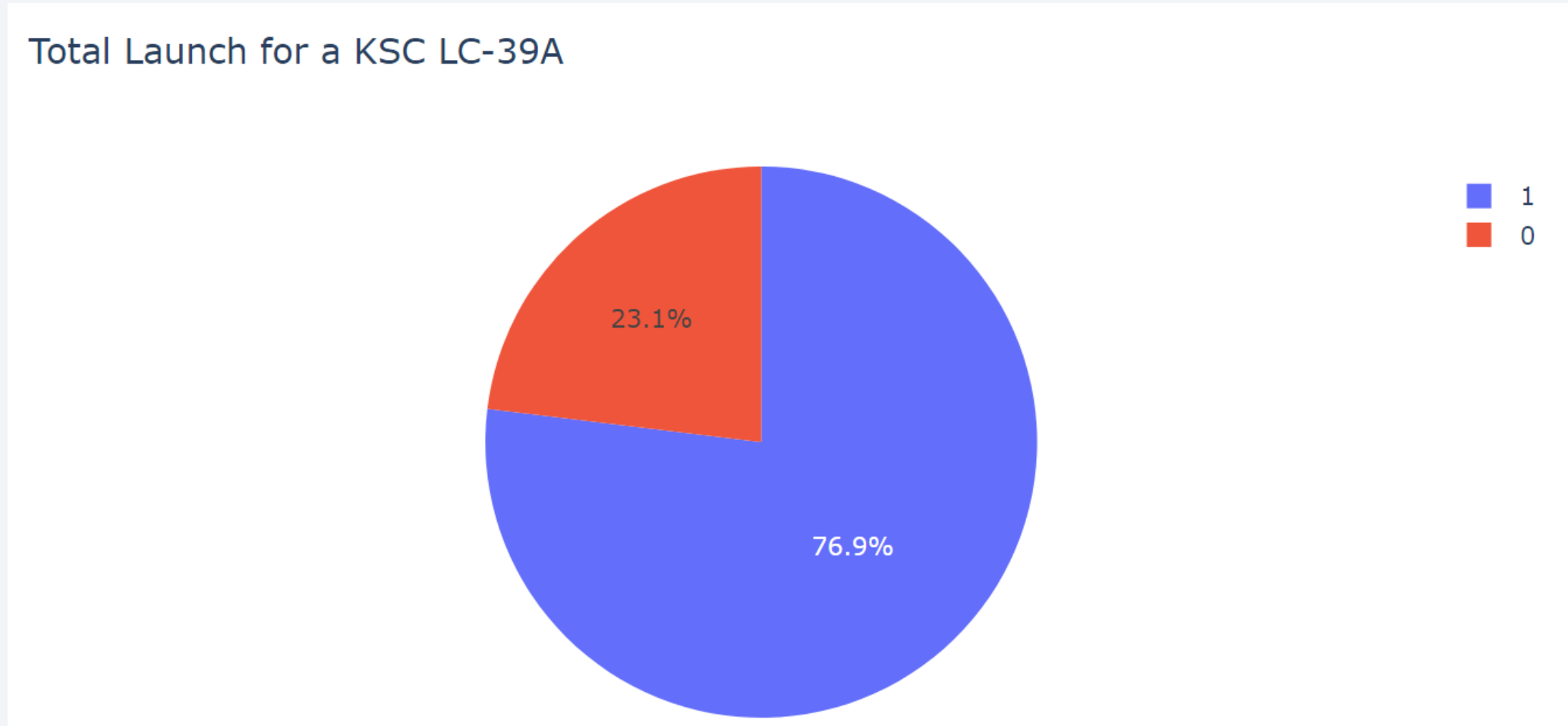- Do launch sites keep certain distance away from cities? Yes

# Build a Dashboard with Plotly Dash

# Pie chart showing the success percentage achieved by each launch site



**Total Launches for All Sites**

Legend:
- KSC LC-39A
- CCAFS LC-40
- VAFB SLC-4E
- CCAFS SLC-40
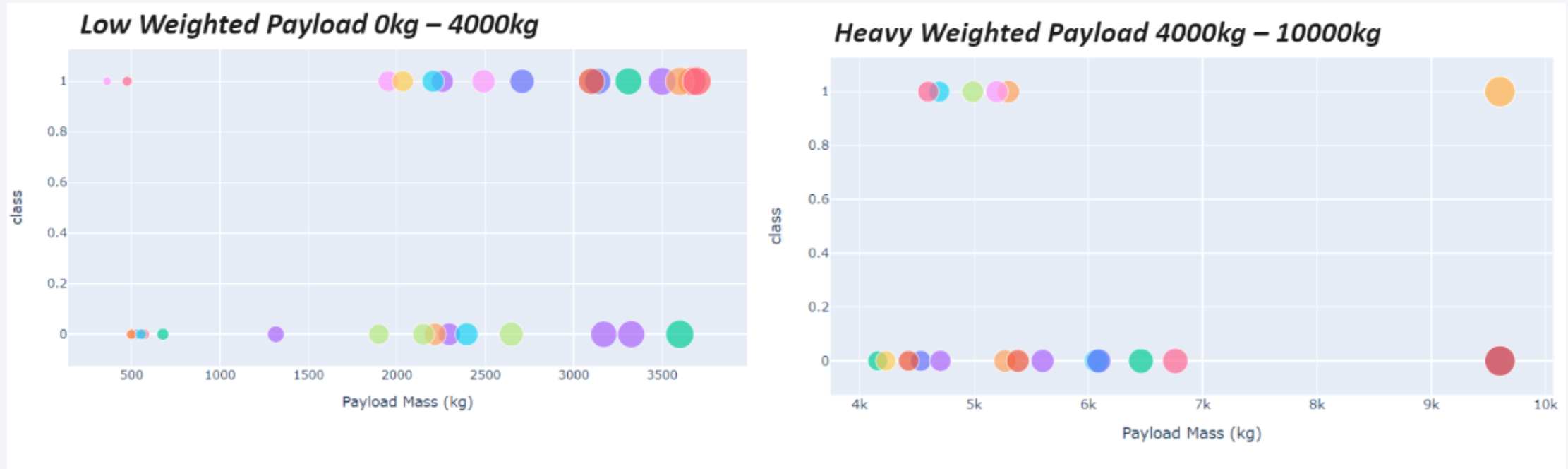
Pie chart values: 41.7%, 29.2%, 16.7%, 12.5%

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

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



Total Launch for a KSC LC-39A

23.1%

76.9%

1
0

KSC LC-39A archieved a 76,9% success rate while getting a 23,1% failure rate

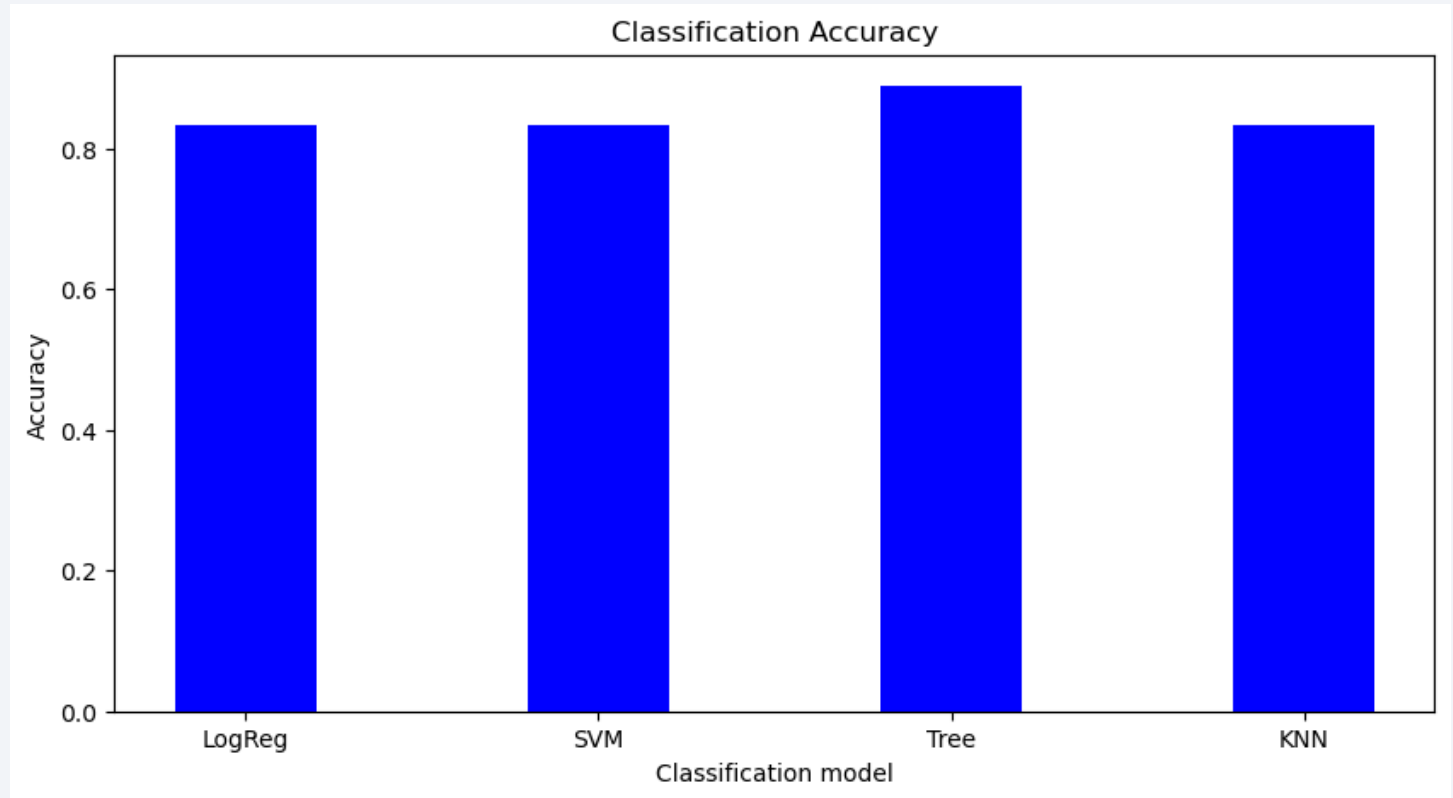# Scatter plot of Payload vs Launch Outcome for all sites, with different payload selected in the range slider



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

Section 5

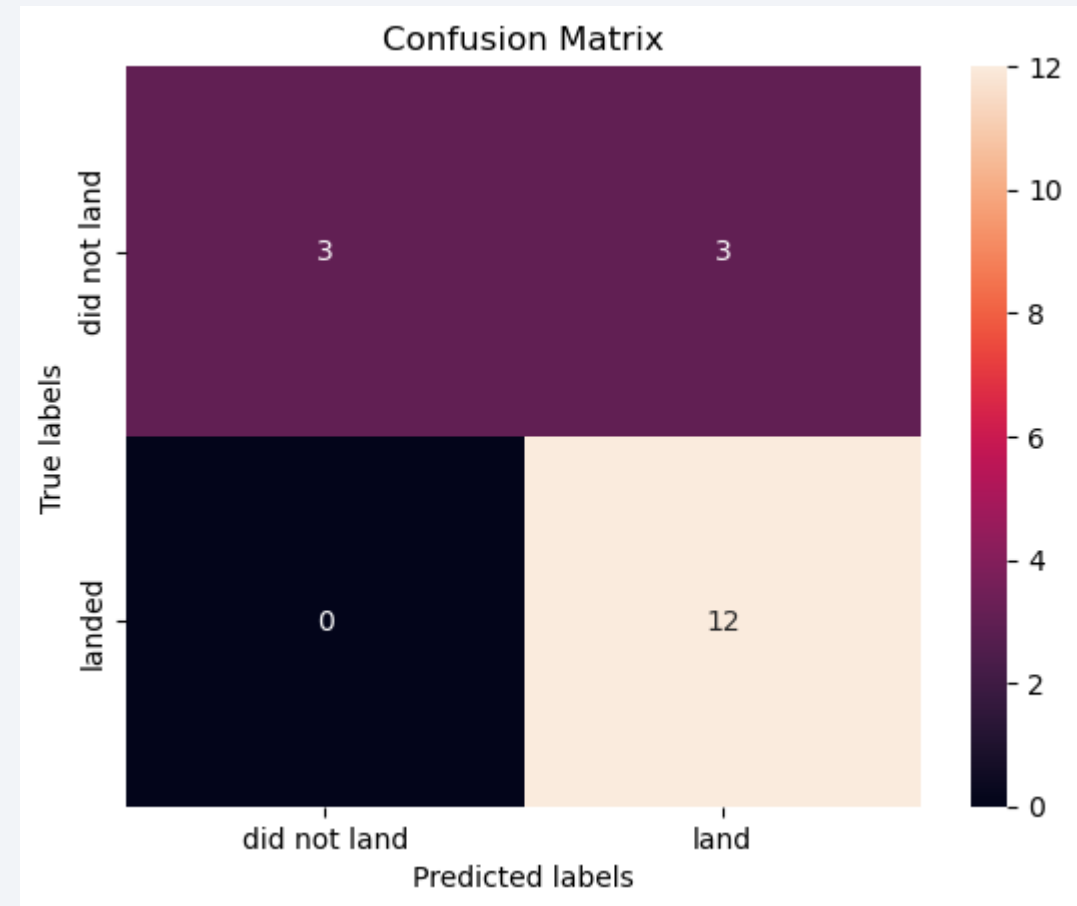# Predictive Analysis (Classification)

# Classification Accuracy

The decision tree classifier is the model with the highest classification accuracy with a 0.88

# Confusion Matrix

- The confusion matrix for the decision tree classifier shows that the classifier can distinguish between the different classes. The major problem is the false positives .i.e., unsuccessful landing marked as successful landing by the classifier.

# Conclusions

We can conclude that:

- The larger the flight amount at a launch site, the greater the success rate at a launch site.

- Launch success rate started to increase in 2013 till 2020.

- Orbits ES-L1, GEO, HEO, SSO, VLEO had the most success rate.

- KSC LC-39A had the most successful launches of any sites.

- The Decision tree classifier is the best machine learning algorithm for this task.

Thank you!