# 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
  - SpaceX has transformed the space industry by dramatically reducing the cost of rocket launches. The Falcon 9, priced at $62 million, is significantly more affordable than competitors' rockets, which typically cost upwards of $165 million. This cost advantage is largely due to SpaceX's innovative approach of reusing the first stage of the rocket, which is successfully landed and repurposed for future missions. By reusing the first stage, SpaceX continues to reduce launch costs over time.
  - As a data scientist at a startup aiming to compete with SpaceX, the goal of this project is to develop a machine learning pipeline that predicts the success of the first stage landing. This prediction will help determine launch costs and enable competing companies to strategically bid against SpaceX.

- Problems you want to find answers
  - What factors influence the successful landing of the rocket's first stage?
  - How do variables such as payload mass, launch site, number of flights, and orbits impact the landing success?
  - What operational conditions maximize the likelihood of a successful landing?
  - How has the success rate of first-stage landings evolved over the years?
  - What is the most effective machine learning algorithm for predicting landing success in this binary classification task?

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 for this project was collected using two methods: REST API and web scraping. These approaches enabled the extraction of relevant data for analysis and decision-making.
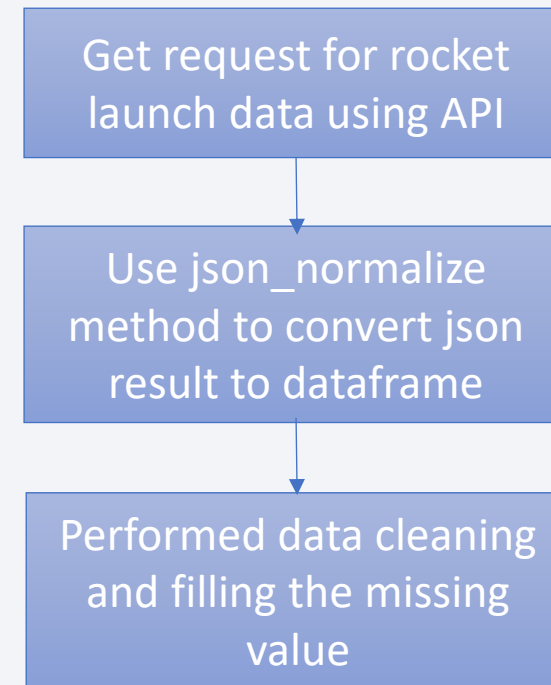
- REST API:

  Using the SpaceX API, data was gathered by sending GET requests. The response was decoded into JSON format using the .json() function and converted into a pandas dataframe with the json_normalize() method. The dataset was then cleaned, checked for missing values, and appropriately filled to ensure completeness.

- Web Scraping:

  Falcon 9 launch records were extracted from Wikipedia using BeautifulSoup. The data, retrieved as HTML tables, was parsed and transformed into pandas dataframes for further 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.

- GitHub URL: https://github.com/yash2412-d/IBM-Data-Science/blob/main/Applied%20Data%20Science%20Capstone/Week%201/jupyter-labs-spacex-data-collection-api.ipynb

Get request for rocket launch data using API

Use json_normalize method to convert json result to dataframe

Performed data cleaning and filling the missing value

# 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.

- GitHub URL: https://github.com/yash2412-d/IBM-Data-Science/blob/main/Applied%20Data%20Science%20Capstone/Week%201/jupyter-labs-webscraping.ipynb
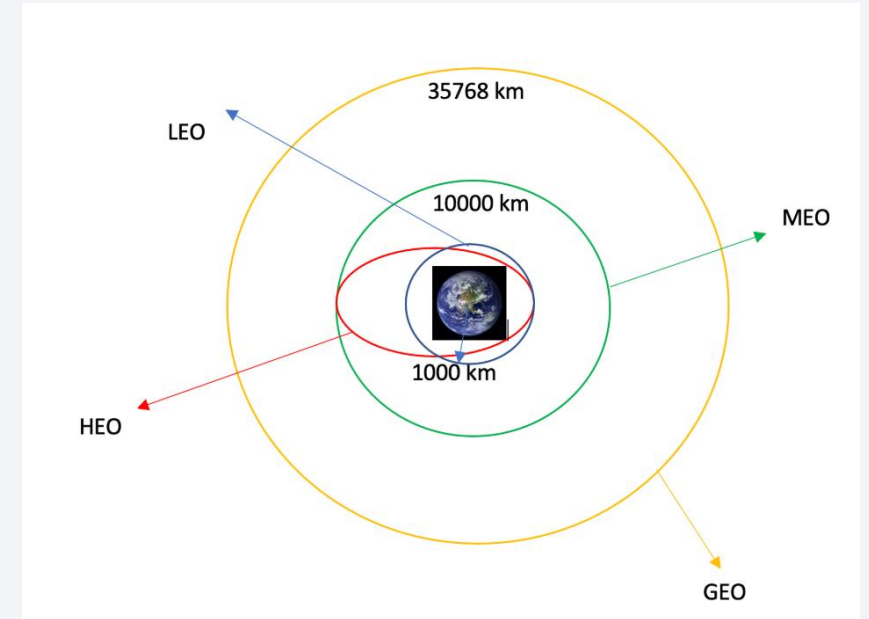
Request the Falcon9 Launch Wiki page from url

Create a BeautifulSoup from the HTML response

Extract all column/variable names from the HTML header

# Data Wrangling

Data Wrangling is the process of cleaning and unifying messy and complex data sets for easy access and Exploratory Data Analysis (EDA).

We will first calculate the number of launches on each site, then calculate the number and occurrence of mission outcome per orbit type.

We then create a landing outcome label from the outcome column. This will make it easier for further analysis, visualization, and ML. Lastly, we will export the result to a CSV.



GitHub URL: https://github.com/yash2412-d/IBM-Data-Science/blob/main/Applied%20Data%20Science%20Capstone/Week%201/labs-jupyter-spacex-Data%20wrangling.ipynb
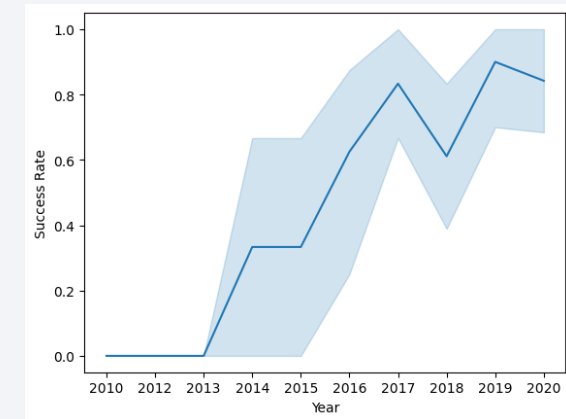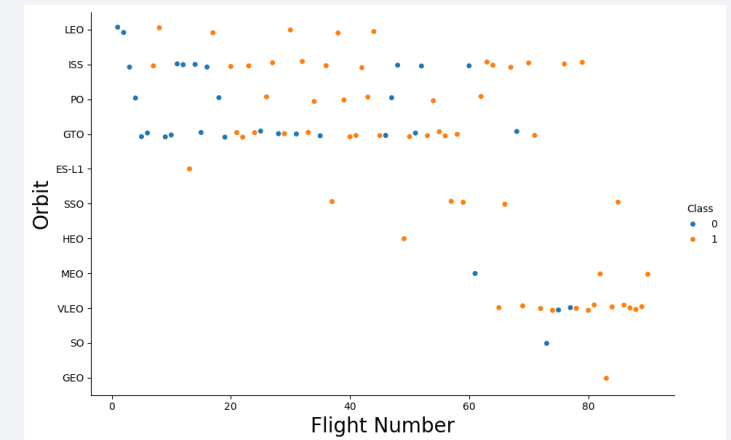
# EDA with Data Visualization

- 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: https://github.com/yash2412-d/IBM-Data-Science/blob/main/Applied%20Data%20Science%20Capstone/Week%202/edadataviz.ipynb

# EDA with SQL

Using SQL, we had performed many queries to get better understanding of the dataset, Ex:

- Displaying the names of the launch sites.

- Displaying 5 records where launch sites begin with the string 'CCA'.

- Displaying the total payload mass carried by booster launched by NASA (CRS).

- Displaying the average payload mass carried by booster version F9 v1.1.

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

- Listing the names of the boosters which have success in drone ship 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 failed landing_outcomes in drone ship, their booster versions, and launch sites names for in year 2015.

- Rank the count of landing outcomes or success between the date 2010-06-04 and 2017-03-20, in descending order.

# 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.

**Colored Markers of the launch outcomes for each Launch Site:**

- Added colored 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 colored Lines to show distances between the Launch Site KSC LC-39A (as an example) and its proximities like Railway, Highway, Coastline and Closest City.

GitHub URL: https://github.com/yash2412-d/IBM-Data-Science/blob/main/Applied%20Data%20Science%20Capstone/Week%203/lab_jupyter_launch_site_location.ipynb

# 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.
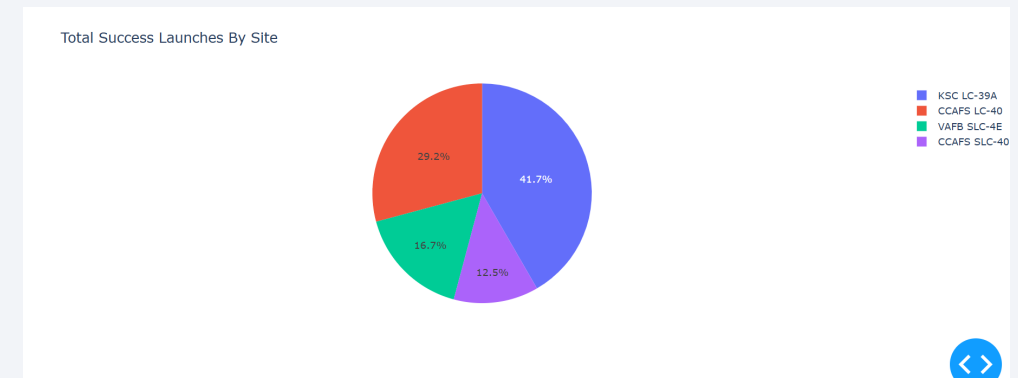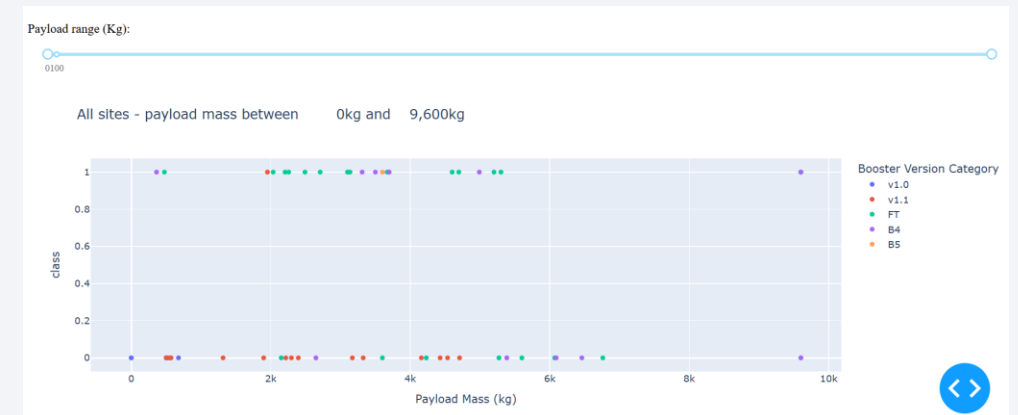
**Slider of Payload Mass Range:**

- 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: https://github.com/yash2412-d/IBM-Data-Science/blob/main/Applied%20Data%20Science%20Capstone/Week%203/spacex_dash_app.py

# Predictive Analysis (Classification)

**1. Building the Model**

- Load the dataset into
- NumPy and Pandas
- Transform the data and then split into training and test datasets
- Decide which type of ML to use
- Set the parameters and algorithms to GridSearchCV and fit it to dataset.

**2. Evaluating the Model**

- Check the accuracy for each model
- Get tuned hyperparameters for each type of algorithms
- Plot the confusion matrix.

**3. Improving the Model**

- Use Feature Engineering and Algorithm Tuning

**4. Find the Best Model**

- The model with the best accuracy score will be the best performing model.

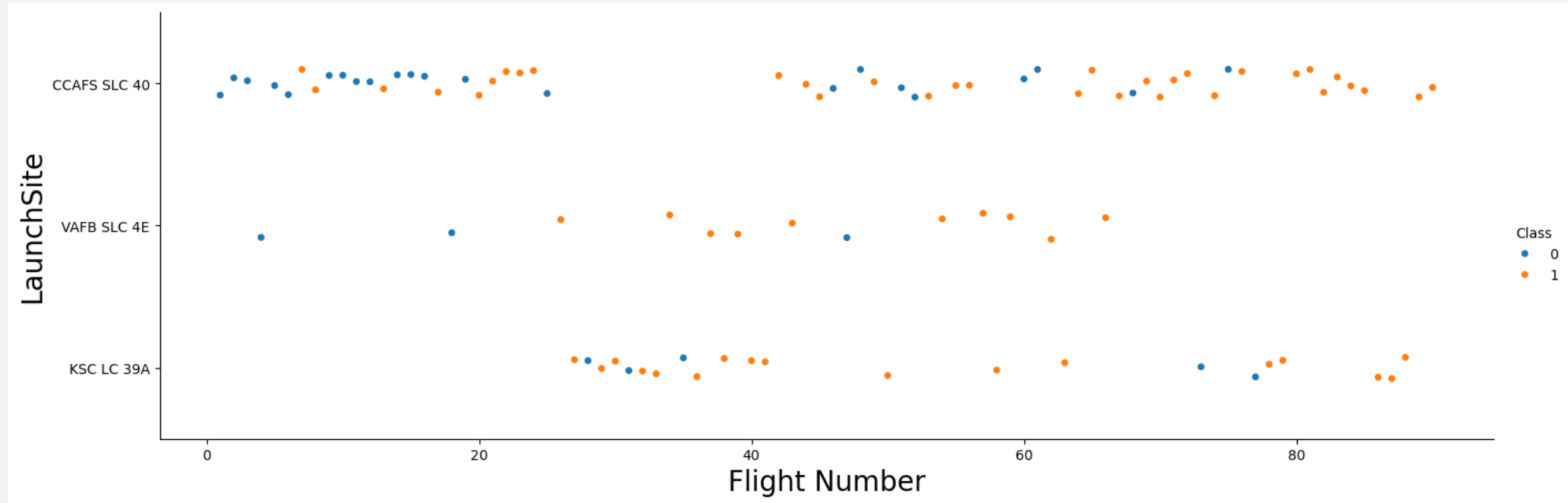# Results

- Exploratory data analysis results

- Interactive analytics demo in screenshots

- Predictive analysis results
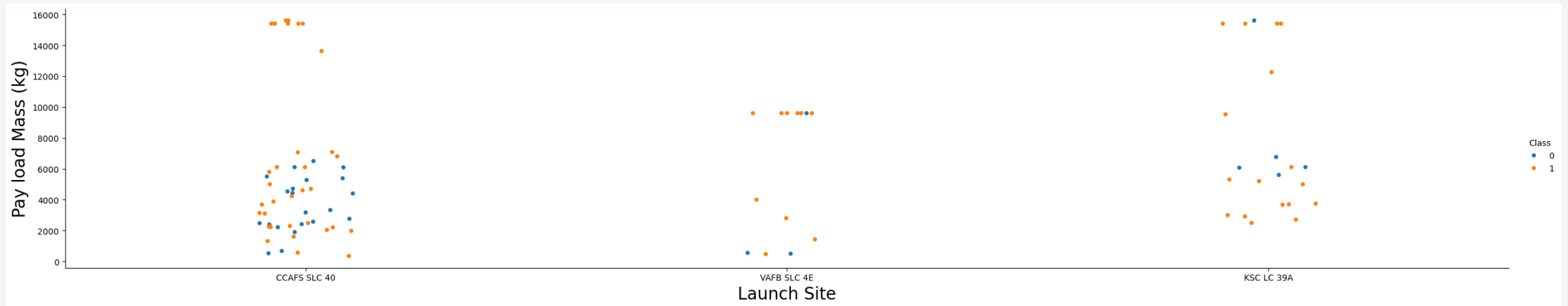
Section 2

# Insights drawn from EDA

# Flight Number vs. Launch Site



Explanation:

• The earliest flights all failed while the latest flights all succeeded.

• The CCAFS SLC 40 launch site has about a half of all launches.

• VAFB SLC 4E and KSC LC 39A have higher success rates.

• It can be assumed that each new launch has a higher rate of success.
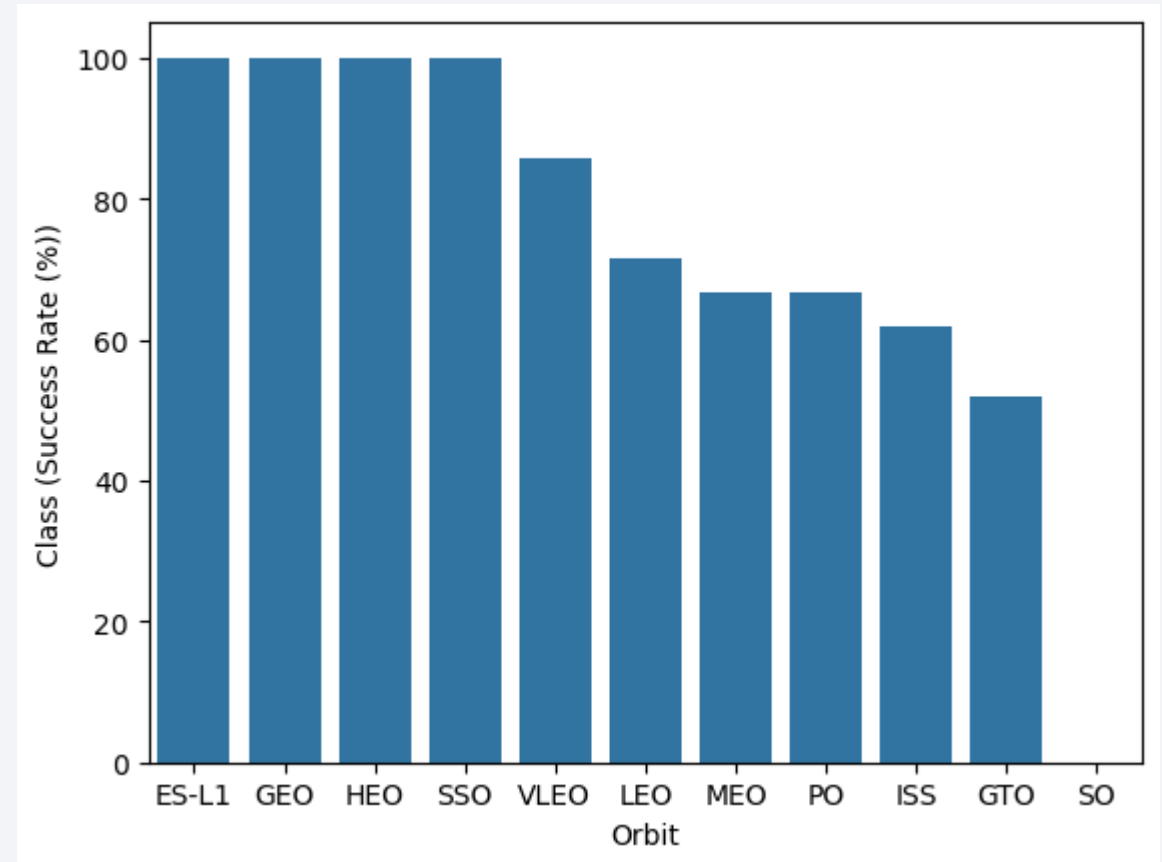
# Payload vs. Launch Site



Explanation:

• For every launch site the higher the payload mass, the higher the success rate.

• Most of the launches with payload mass over 7000 kg were successful.

• KSC LC 39A has a 100% success rate for payload mass under 5500 kg too.

# Success Rate vs. Orbit Type

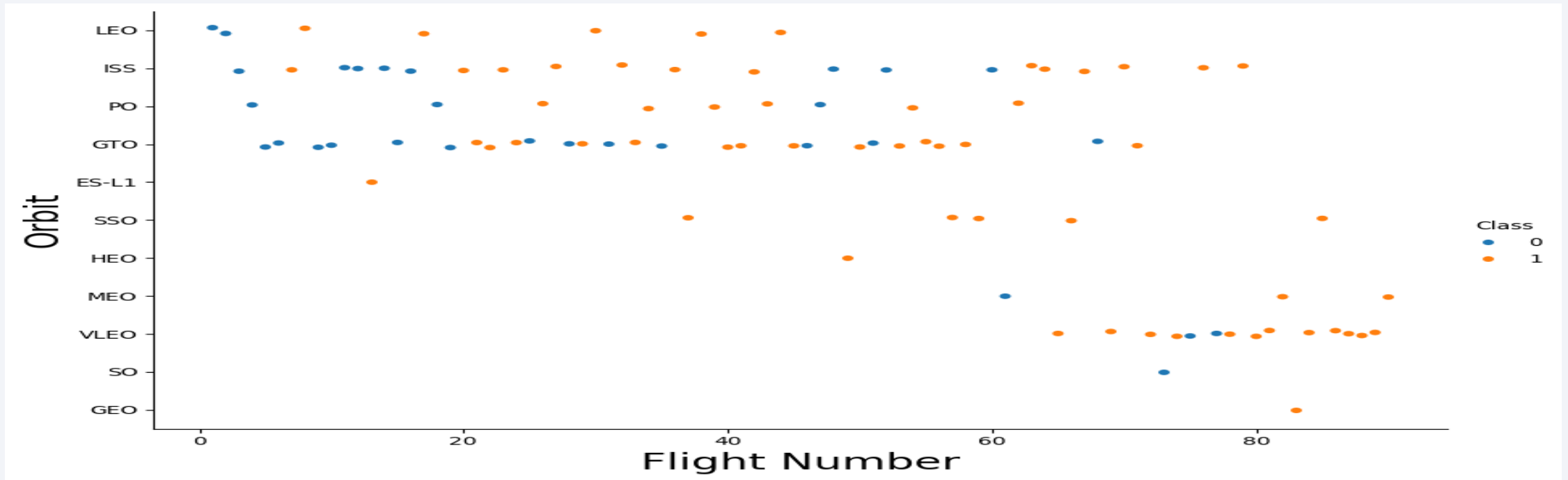Explanation:

- Orbits with 100% success rate:
  - ES-L1, GEO, HEO, SSO

- Orbits with 0% success rate:
  - SO

- Orbits with success rate between 50% and 85%:
  - GTO, ISS, LEO, MEO, PO

# Flight Number vs. Orbit Type



Explanation:

• 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.

# Payload vs. Orbit Type



Explanation:

• Heavy payloads have a negative influence on GTO orbits and positive on GTO and Polar LEO (ISS) orbits.

22

# Launch Success Yearly Trend

Explanation:

• The success rate since 2013 kept increasing till 2020.

# All Launch Site Names

Explanation:

• Displaying the names of the unique launch sites in the space mission.



```
In [11]:  %sql SELECT DISTINCT LAUNCH_SITE as "Launch_Sites" FROM SPACEXTBL;

 * sqlite:///my_data1.db
Done.
Out[11]:  Launch_Sites

          CCAFS LC-40

          VAFB SLC-4E

          KSC LC-39A

          CCAFS SLC-40
```

# Launch Site Names Begin with 'CCA'

Explanation:

• Displaying 5 records where launch sites begin with the string 'CCA'.

```
In [12]:  %sql SELECT * FROM 'SPACEXTBL' WHERE Launch_Site LIKE 'CCA%' LIMIT 5;
```

* sqlite:///my_data1.db
Done.

Out[12]:

| Date | Time (UTC) | Booster_Version | Launch_Site | Payload | PAYLOAD_MASS__KG_ | Orbit | Customer | Mission_Outcome |
|---|---|---|---|---|---|---|---|---|
| 2010-06-04 | 18:45:00 | F9 v1.0 B0003 | CCAFS LC-40 | Dragon Spacecraft Qualification Unit | 0 | LEO | SpaceX | Success |
| 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 |
| 2012-05-22 | 7:44:00 | F9 v1.0 B0005 | CCAFS LC-40 | Dragon demo flight C2 | 525 | LEO (ISS) | NASA (COTS) | Success |
| 2012-10-08 | 0:35:00 | F9 v1.0 B0006 | CCAFS LC-40 | SpaceX CRS-1 | 500 | LEO (ISS) | NASA (CRS) | Success |
| 2013-03-01 | 15:10:00 | F9 v1.0 B0007 | CCAFS LC-40 | SpaceX CRS-2 | 677 | LEO (ISS) | NASA (CRS) | Success |

# Total Payload Mass

```
In [13]:    %sql SELECT SUM(PAYLOAD_MASS__KG_) as "Total Payload Mass(Kgs)", Customer FROM 'SPACEXTBL' WHERE Customer = 'NASA (CRS)';

            * sqlite:///my_data1.db
            Done.
Out[13]:    Total Payload Mass(Kgs)    Customer

                              45596    NASA (CRS)
```

Explanation:

• Displaying the total payload mass carried by boosters launched by NASA (CRS).

# Average Payload Mass by F9 v1.1

Query:

%sql SELECT AVG(PAYLOAD_MASS__KG_) as "Payload Mass Kgs", Customer, Booster_Version FROM 'SPACEXTBL' WHERE Booster_Version LIKE 'F9 v1.1%';

```
In [14]:    %sql SELECT AVG(PAYLOAD_MASS__KG_) as "Payload Mass Kgs", Customer, Booster_Version FROM 'SPACEXTBL' WHERE Booster_Version L

         * sqlite:///my_data1.db
         Done.
Out[14]:    Payload Mass Kgs   Customer   Booster_Version

         2534.6666666666665       MDA      F9 v1.1 B1003
```

Explanation:

- Displaying average payload mass carried by booster version F9 v1.1.

# First Successful Ground Landing Date

```
In [24]:    %sql SELECT MIN(DATE) as Date FROM 'SPACEXTBL' WHERE "Landing_Outcome" = "Success (ground pad)";

            * sqlite:///my_data1.db
            Done.
Out[24]:        Date

            2015-12-22
```

Explanation:

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

# Successful Drone Ship Landing with Payload between 4000 and 6000

Query:

%sql SELECT DISTINCT Booster_Version, Payload FROM SPACEXTBL WHERE "Landing_Outcome" = "Success (drone ship)" AND PAYLOAD_MASS__KG_ > 4000 AND PAYLOAD_MASS__KG_ < 6000;

```
In [25]:   %sql SELECT DISTINCT Booster_Version, Payload FROM SPACEXTBL WHERE "Landing_Outcome" = "Success (drone ship)" AND PAYLOAD_MA

         * sqlite:///my_data1.db
         Done.
Out[25]:
```

| Booster_Version | Payload |
|---|---|
| F9 FT B1022 | JCSAT-14 |
| F9 FT B1026 | JCSAT-16 |
| F9 FT B1021.2 | SES-10 |
| F9 FT B1031.2 | SES-11 / EchoStar 105 |

Explanation:

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

# Total Number of Successful and Failure Mission Outcomes

**List the total number of successful and failure mission outcomes**

```
In [27]:  %sql SELECT "Mission_Outcome", COUNT(*) as Total FROM SPACEXTBL GROUP BY "Mission_Outcome" ORDER BY "Mission_Outcome";
```

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

Out[27]:

| Mission_Outcome | Total |
|---|---|
| Failure (in flight) | 1 |
| Success | 98 |
| Success | 1 |
| Success (payload status unclear) | 1 |

Explanation:

• Listing the total number of successful and failure mission outcomes.

# Boosters Carried Maximum Payload

Query:

%sql SELECT "Booster_Version",Payload, "PAYLOAD_MASS__KG_" FROM SPACEXTBL WHERE "PAYLOAD_MASS__KG_" = (SELECT MAX("PAYLOAD_MASS__KG_") FROM SPACEXTBL);

```
In [28]:  %sql SELECT "Booster_Version",Payload, "PAYLOAD_MASS__KG_" FROM SPACEXTBL WHERE "PAYLOAD_MASS__KG_" = (SELECT MAX("PAYLOAD_M
```

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

Out[28]:

| Booster_Version | Payload | PAYLOAD_MASS__KG_ |
|---|---|---|
| F9 B5 B1048.4 | Starlink 1 v1.0, SpaceX CRS-19 | 15600 |
| F9 B5 B1049.4 | Starlink 2 v1.0, Crew Dragon in-flight abort test | 15600 |
| F9 B5 B1051.3 | Starlink 3 v1.0, Starlink 4 v1.0 | 15600 |
| F9 B5 B1056.4 | Starlink 4 v1.0, SpaceX CRS-20 | 15600 |
| F9 B5 B1048.5 | Starlink 5 v1.0, Starlink 6 v1.0 | 15600 |
| F9 B5 B1051.4 | Starlink 6 v1.0, Crew Dragon Demo-2 | 15600 |
| F9 B5 B1049.5 | Starlink 7 v1.0, Starlink 8 v1.0 | 15600 |
| F9 B5 B1060.2 | Starlink 11 v1.0, Starlink 12 v1.0 | 15600 |
| F9 B5 B1058.3 | Starlink 12 v1.0, Starlink 13 v1.0 | 15600 |
| F9 B5 B1051.6 | Starlink 13 v1.0, Starlink 14 v1.0 | 15600 |
| F9 B5 B1060.3 | Starlink 14 v1.0, GPS III-04 | 15600 |
| F9 B5 B1049.7 | Starlink 15 v1.0, SpaceX CRS-21 | 15600 |

Explanation:

• Listing the names of the booster versions which have carried the maximum payload mass.

# 2015 Launch Records

Query:

%sql SELECT CASE WHEN substr(Date, 6, 2) = '01' THEN 'January' WHEN substr(Date, 6, 2) = '02' THEN 'February' WHEN substr(Date, 6, 2) = '03' THEN 'March' WHEN substr(Date, 6, 2) = '04' THEN 'April' WHEN substr(Date, 6, 2) = '05' THEN 'May' WHEN substr(Date, 6, 2) = '06' THEN 'June' WHEN substr(Date, 6, 2) = '07' THEN 'July' WHEN substr(Date, 6, 2) = '08' THEN 'August' WHEN substr(Date, 6, 2) = '09' THEN 'September' WHEN substr(Date, 6, 2) = '10' THEN 'October' WHEN substr(Date, 6, 2) = '11' THEN 'November' WHEN substr(Date, 6, 2) = '12' THEN 'December' END AS Month, substr(Date, 1, 4) AS Year, "Booster_Version", "Launch_Site", Payload, "PAYLOAD_MASS__KG_", "Mission_Outcome", "Landing_Outcome" FROM SPACEXTBL WHERE substr(Date, 1, 4) = '2015' AND "Landing_Outcome" = 'Failure (drone ship)';

```
In [37]:   %sql SELECT CASE WHEN substr(Date, 6, 2) = '01' THEN 'January' WHEN substr(Date, 6, 2) = '02' THEN 'February' WHEN substr(Da

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

Out[37]:

| Month | Year | Booster_Version | Launch_Site | Payload | PAYLOAD_MASS__KG_ | Mission_Outcome | Landing_Outcome |
|-------|------|-----------------|-------------|---------|-------------------|-----------------|-----------------|
| January | 2015 | F9 v1.1 B1012 | CCAFS LC-40 | SpaceX CRS-5 | 2395 | Success | Failure (drone ship) |
| April | 2015 | F9 v1.1 B1015 | CCAFS LC-40 | SpaceX CRS-6 | 1898 | Success | Failure (drone ship) |

Explanation:

• Listing the failed landing outcomes in drone ship, their booster versions and launch site names for the months in year 2015.

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

Query:
%sql SELECT * FROM SPACEXTBL WHERE "Landing_Outcome" LIKE 'Success%' AND (Date BETWEEN '2010-06-04' AND '2017-03-20') ORDER BY Date DESC;
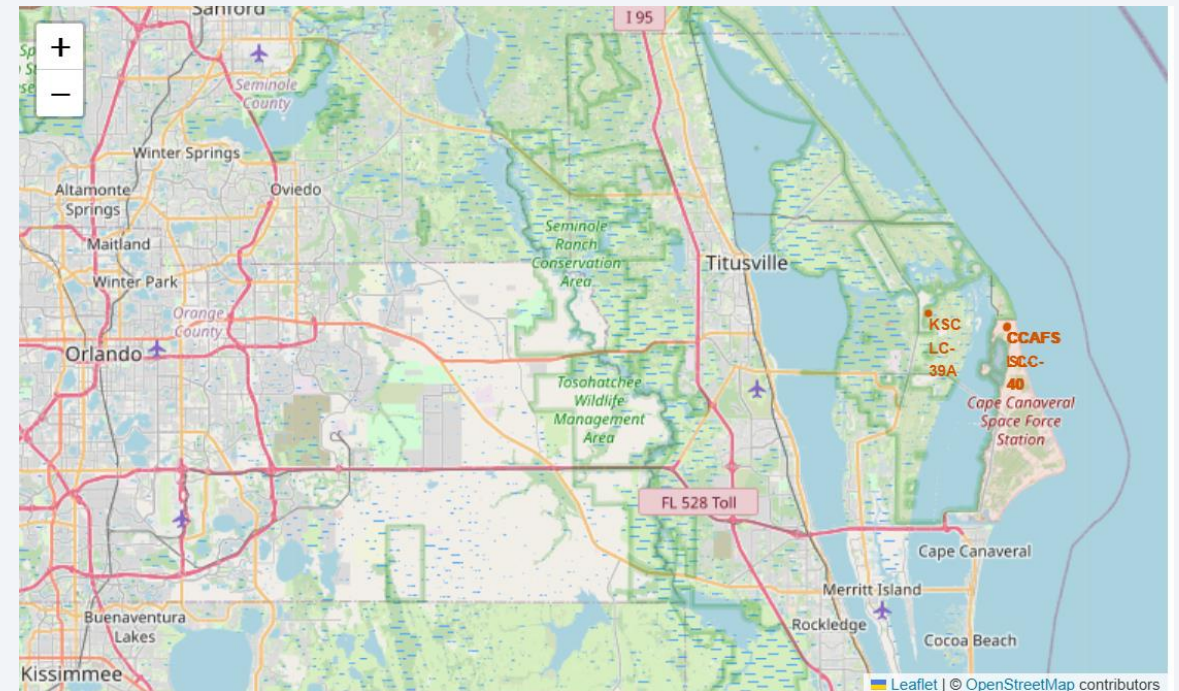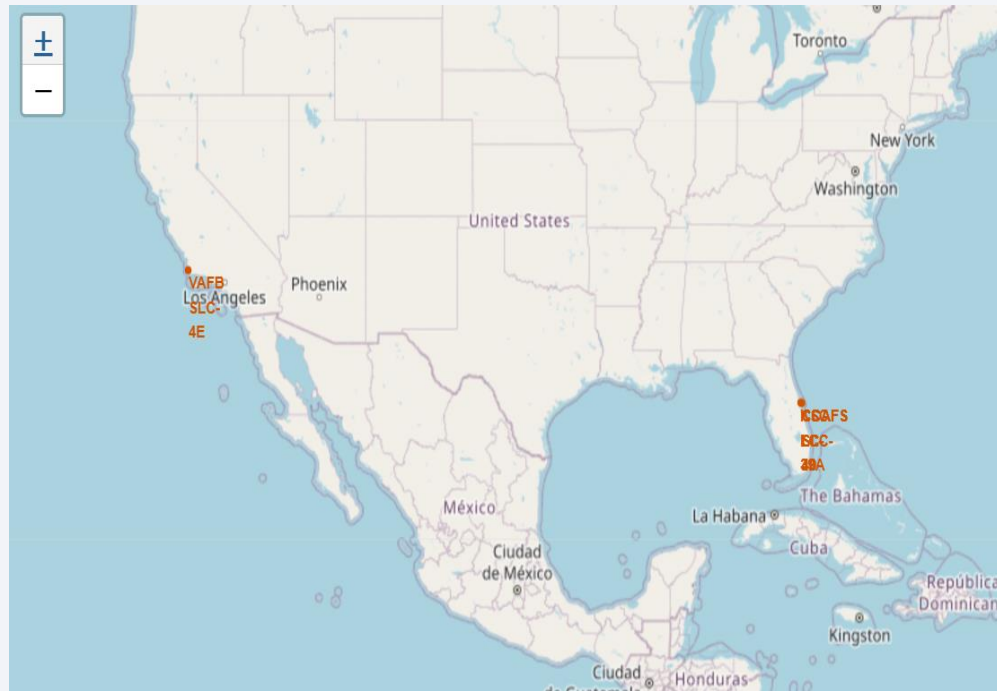
Explanation:

• Ranking 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 [40]: %sql SELECT * FROM SPACEXTBL WHERE "Landing_Outcome" LIKE 'Success%' AND (Date BETWEEN '2010-06-04' AND '2017-03-20') ORDER
```

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

Out[40]:

| Date | Time (UTC) | Booster_Version | Launch_Site | Payload | PAYLOAD_MASS_KG_ | Orbit | Customer | Mission_Outcome | Landing_Outco |
|---|---|---|---|---|---|---|---|---|---|
| 2017-02-19 | 14:39:00 | F9 FT B1031.1 | KSC LC-39A | SpaceX CRS-10 | 2490 | LEO (ISS) | NASA (CRS) | Success | Success (grou p |
| 2017-01-14 | 17:54:00 | F9 FT B1029.1 | VAFB SLC-4E | Iridium NEXT 1 | 9600 | Polar LEO | Iridium Communications | Success | Success (dro sh |
| 2016-08-14 | 5:26:00 | F9 FT B1026 | CCAFS LC-40 | JCSAT-16 | 4600 | GTO | SKY Perfect JSAT Group | Success | Success (dro sh |
| 2016-07-18 | 4:45:00 | F9 FT B1025.1 | CCAFS LC-40 | SpaceX CRS-9 | 2257 | LEO (ISS) | NASA (CRS) | Success | Success (grou p |
| 2016-05-27 | 21:39:00 | F9 FT B1023.1 | CCAFS LC-40 | Thaicom 8 | 3100 | GTO | Thaicom | Success | Success (dro sh |
| 2016-05-06 | 5:21:00 | F9 FT B1022 | CCAFS LC-40 | JCSAT-14 | 4696 | GTO | SKY Perfect JSAT Group | Success | Success (dro sh |
| 2016-04-08 | 20:43:00 | F9 FT B1021.1 | CCAFS LC-40 | SpaceX CRS-8 | 3136 | LEO (ISS) | NASA (CRS) | Success | Success (dro sh |
| 2015-12-22 | 1:29:00 | F9 FT B1019 | CCAFS LC-40 | OG2 Mission 2 11 Orbcomm-OG2 satellites | 2034 | LEO | Orbcomm | Success | Success (grou p |

Section 3

# Launch Sites Proximities Analysis
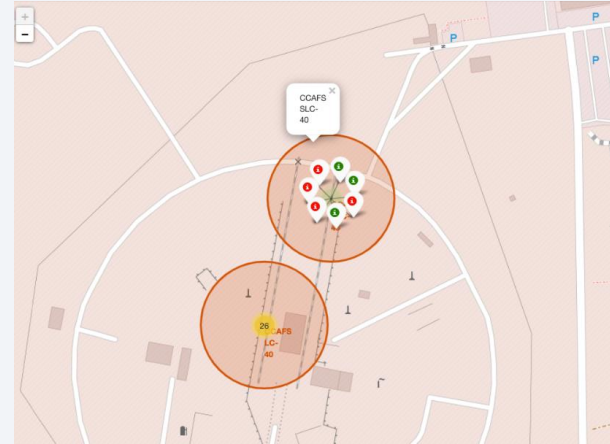
# Launch Site Locations



The left map shows all launch sites relative US map. The right map shows the two Florida launch sites since they are very close to each other. All launch sites are near the ocean.
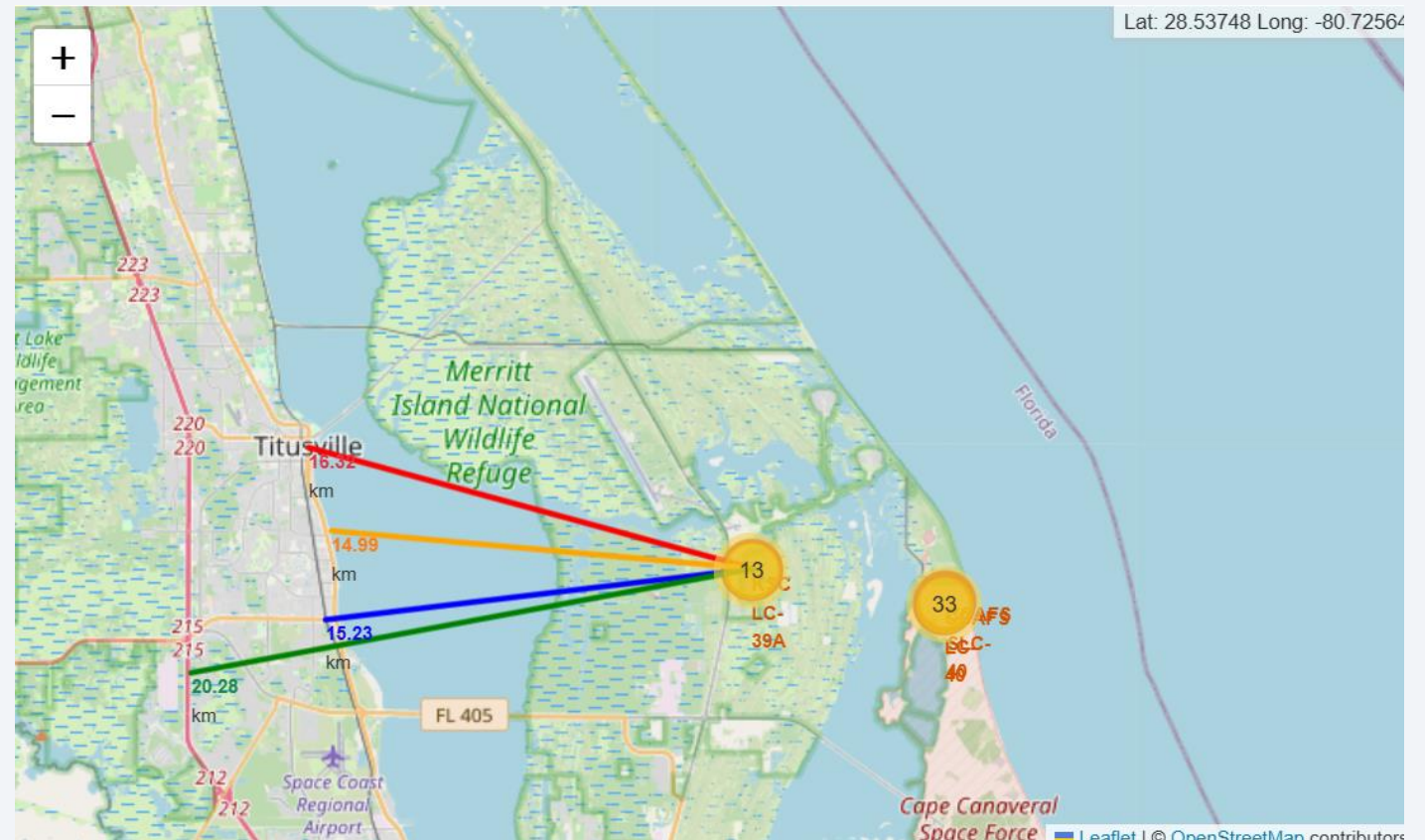
# Color-Coded Launch Markers

Explanation:

- From the colour-labeled markers we should be able to easily identify which launch sites have relatively high success rates.

    - Green Marker = Successful Launch

    - Red Marker = Failed Launch

- Launch Site KSC LC-39A has a very high Success Rate.

# Distance from the launch site KSC LC-39A to its proximities

Using KSC LC-39A as an example, launch sites are very close to railways for large part and supply  transportation. Launch sites are close to highways for human and supply transport. Launch sites  are also close to coasts and relatively far from cities so that launch failures can land in the sea to avoid rockets falling on densely populated areas.
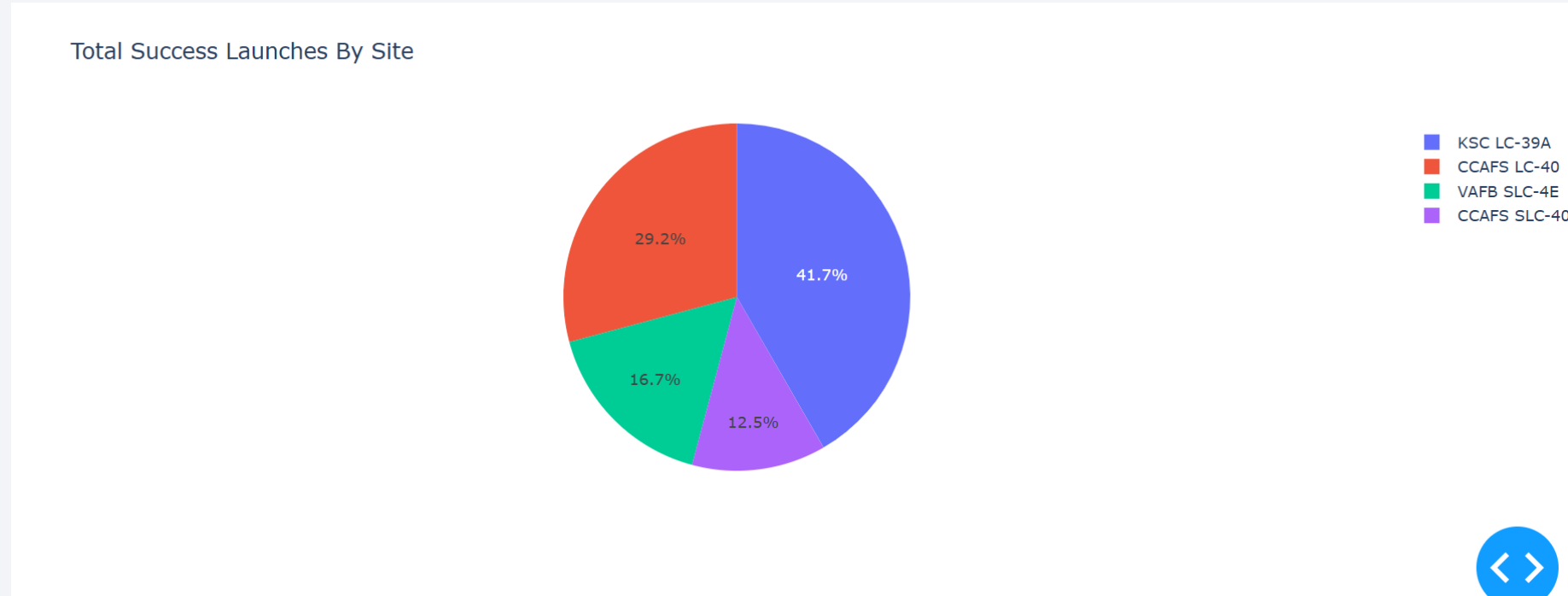
Section 4

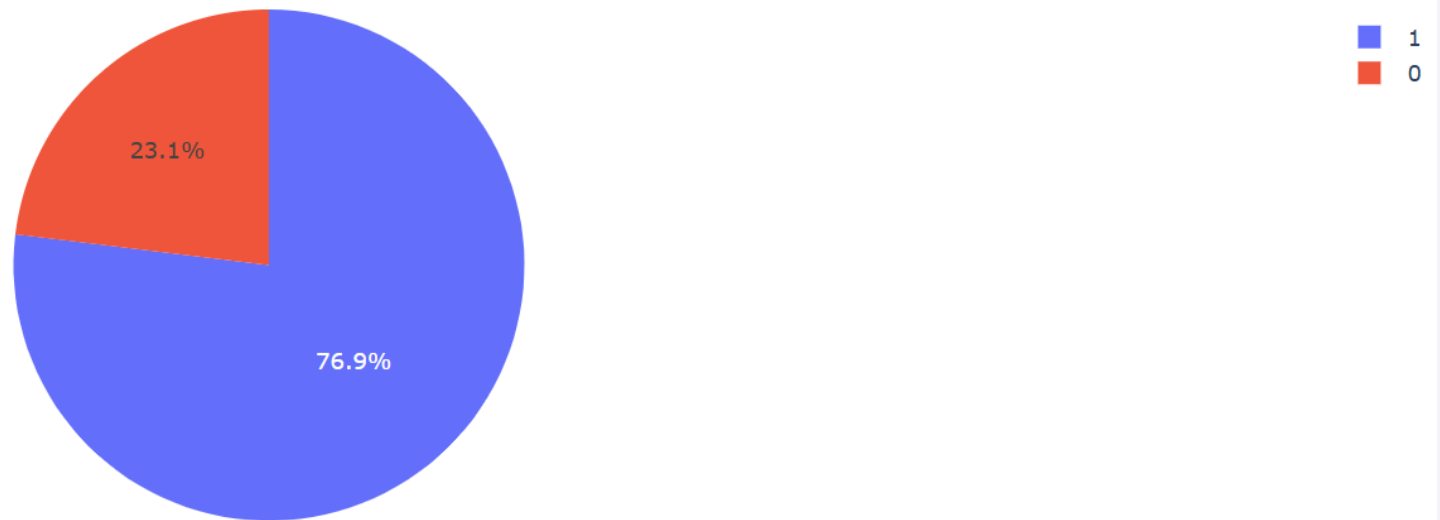# Build a Dashboard
# with Plotly Dash

# Launch success count for all sites



Total Success Launches By Site

KSC LC-39A 41.7%
CCAFS LC-40 29.2%
VAFB SLC-4E 16.7%
CCAFS SLC-40 12.5%

The chart clearly shows that from all the sites, KSC LC-39A has the most successful launches.

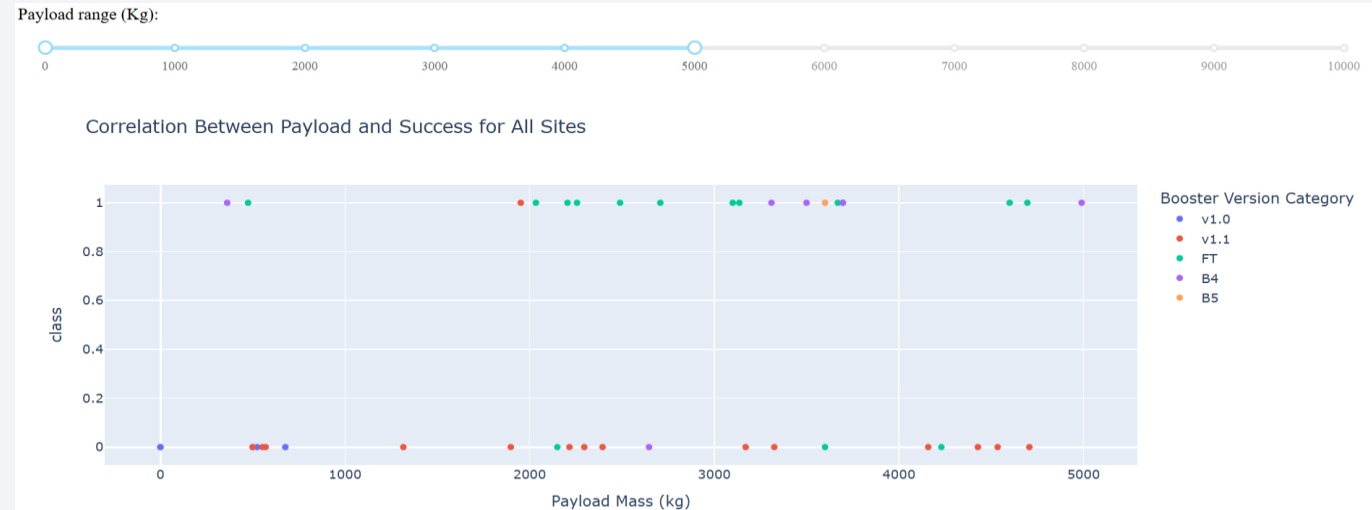# Launch site with highest launch success ratio

Total Launches for site KSC LC-39A



KSC LC-39A has the highest launch success rate (76.9%) with 10 successful and only 3 failed landings.

# Payload Mass vs. Launch Outcome for all sites

The charts show that payloads between 2000 and 5500 kg have the highest success rate.

Section 5

# Predictive Analysis (Classification)

# Classification Accuracy



```
Out[66]:            LogReg      SVM      Tree      KNN

Jaccard_Score     0.800000  0.800000  0.800000  0.800000

     F1_Score     0.888889  0.888889  0.888889  0.888889

     Accuracy     0.833333  0.833333  0.833333  0.833333
```

```
In [68]: algo_score = {'Logistic regresssion': [logreg_cv.best_score_], 'SVM': [svm_cv.best_score_], 'Decision tree': [tree_cv.best_s
         df = pd.DataFrame.from_dict(algo_score, orient='index', columns=['Best scores'])
         df
```

```
Out[68]:                          Best scores

        Logistic regresssion       0.846429

                        SVM        0.848214

               Decision tree       0.891071

                        KNN        0.848214
```
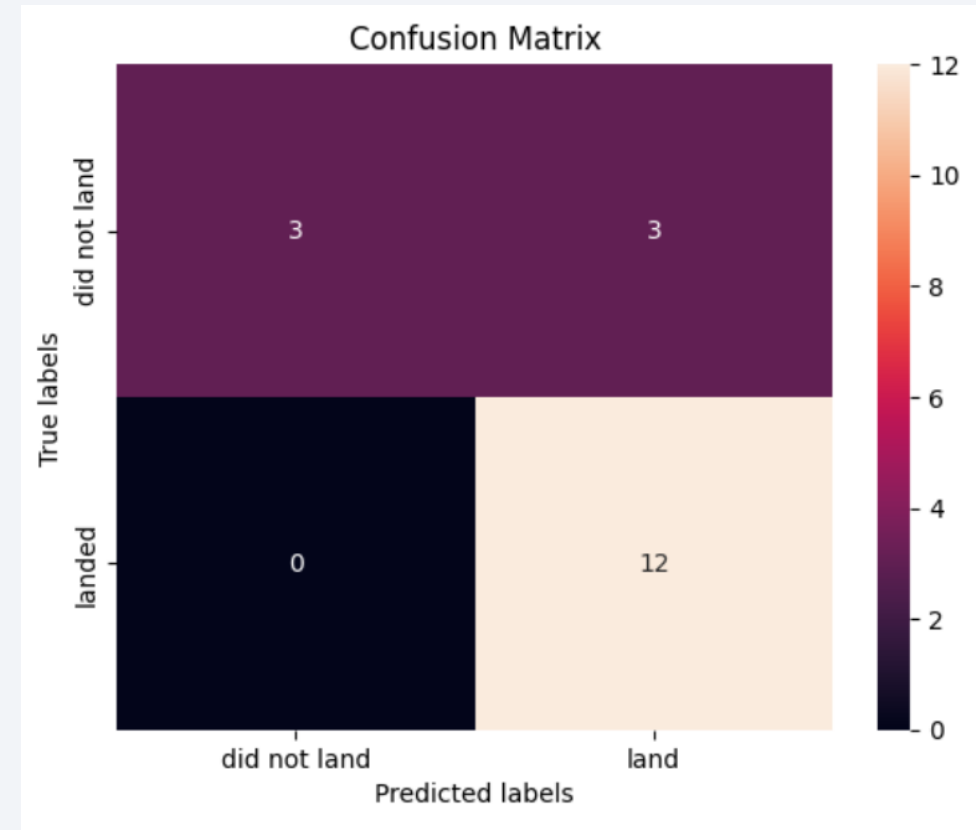
Accuracy of all the models is same. Since they are all same we picked best scores. So, based on their best scores Decision Tree is best model.

# Confusion Matrix

Since all models performed the same for the test set, the confusion matrix is the same across all models. The models predicted 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). Our models over predict successful landings.

# 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!