# Winning Space Race with Data Science

Jose Tavares
26 Jan 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

- Project background and context

  Space X's Falcon 9 rocket launches are more cost-effective than others due to the ability to reuse the first stage. By determining the likelihood of a successful first stage landing, the cost of a launch can be estimated. This information can be useful for other companies looking to compete with Space X for rocket launch contracts. The aim of this project is to develop a machine learning pipeline to predict the success of first stage landings.

- Problems you want to find answers

  - What factors determine if the rocket will land successfully?

  - What features that determine the success rate of a successful landing.

  - What conditions 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

- 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/zhead/AppliedDataScience-Capstone/blob/main/Data%20Collection%20API.ipynb



Task 1: Request and parse the SpaceX launch data using the GET request

To make the requested JSON results more consistent, we will use the following static response object for this project:

```python
static_json_url='https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DS0321EN-SkillsNetwo
```
[9]  ✓ 0.3s                                                                                    Python

We should see that the request was successfull with the 200 status response code

```python
response.status_code
```
[10]  ✓ 0.4s                                                                                   Python
···  200

Now we decode the response content as a Json using `.json()` and turn it into a Pandas dataframe using `.json_normalize()`

```python
# Use json_normalize meethod to convert the json result into a dataframe
data = pd.json_normalize(response.json())
```
[11]  ✓ 0.7s                                                                                   Python

# 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/zhead/Applied DataScience-Capstone/blob/mai n/Data%20Collection%20with%2 0Web%20Scraping.ipynb

First, let's perform an HTTP GET method to request the Falcon9 Launch HTML page, as an HTTP response.

```python
# use requests.get() method with the provided static_url
# assign the response to a object
html_data = requests.get(static_url)
html_data.status_code
```
[5] ✓ 8.3s                                              Python

... 200

Create a BeautifulSoup object from the HTML response

```python
# Use BeautifulSoup() to create a BeautifulSoup object from a response text content
soup = BeautifulSoup(html_data.text, 'html5lib')
```
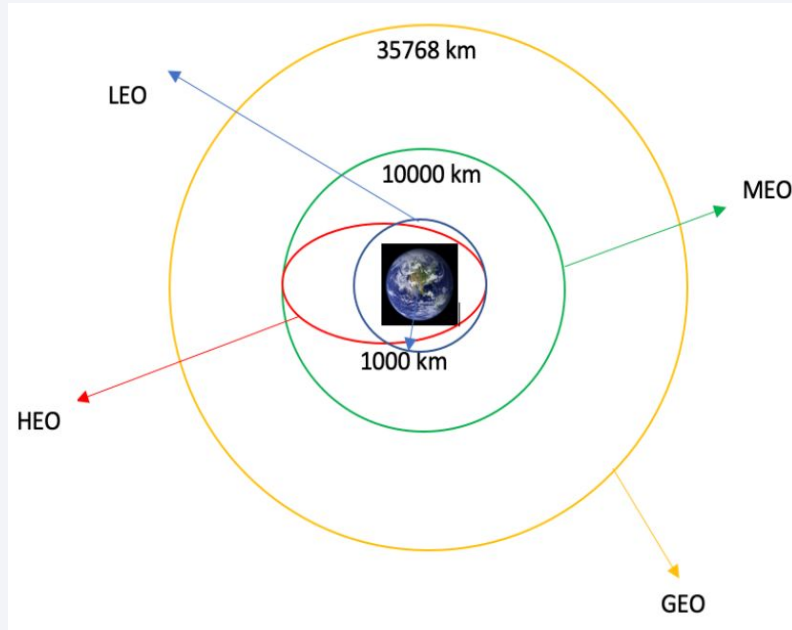[6] ✓ 1.4s                                              Python

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

```python
# Use soup.title attribute
soup.title
```
[7] ✓ 0.3s                                              Python

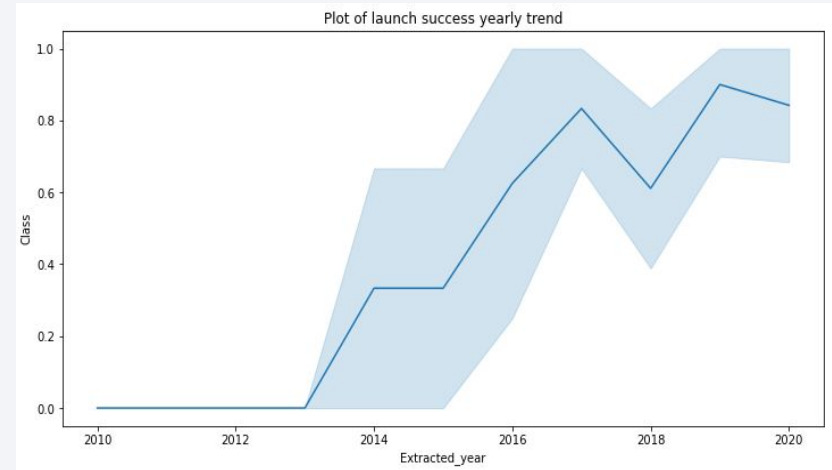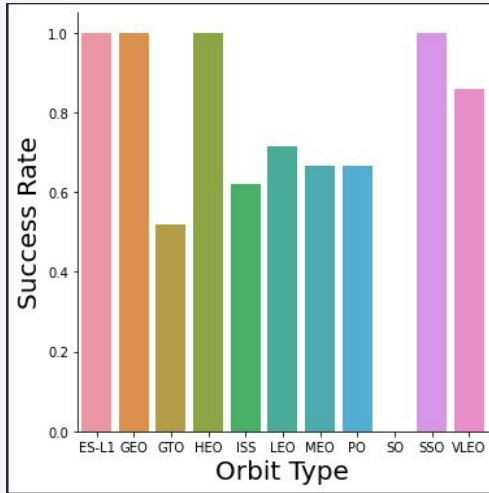... <title>List of Falcon 9 and Falcon Heavy launches - Wikipedia</title>

# Data Wrangling



- Performed exploratory data analysis and determined the training labels.

- Calculated the number of launches at each site, and the number and occurrence of each orbits

- Created landing outcome label from outcome column and exported the results to csv.

- The link to the notebook is https://github.com/zhead/AppliedDataScience-Capstone/blob/main/Data%20Wrangling.ipynb

# EDA with Data Visualization

- 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/zhead/AppliedDataScience-Capstone/blob/main/EDA%20with%20Data%20Visualization.ipynb

# EDA with SQL

- Loaded the SpaceX dataset into a PostgreSQL database without leaving the jupyter notebook.

- 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/zhead/AppliedDataScience-Capstone/blob/main/EDA%20with%20SQL.ipynb

# Build an Interactive Map with Folium

- 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.

- 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.

- 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.

# Build a Dashboard with Plotly Dash

- Built an interactive dashboard with Plotly dash

- Plotted pie charts showing the total launches by a certain sites

- 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/zhead/AppliedDataScience-Capstone/blob/main/space x_dash_app.py

# Predictive Analysis (Classification)

- Loaded the data using numpy and pandas, transformed the data, split our data into training and testing.

- Built different machine learning models and tune different hyperparameters using GridSearchCV.

- Used accuracy as the metric for our model, improved the model using feature engineering and algorithm tuning.

- Found the best performing classification model.

The link to the notebook is
https://github.com/zhead/AppliedDataScience-Capstone/blob/main/Machine%20Learning%20Prediction.ipynb

# Results

- Exploratory data analysis results

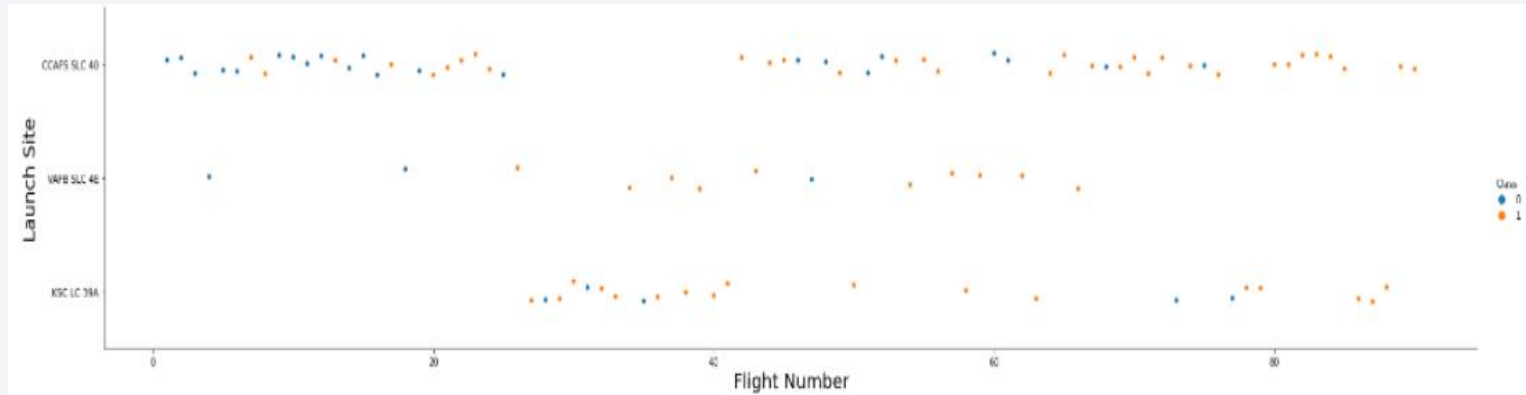- Interactive analytics demo in screenshots

- Predictive analysis results

Section 2
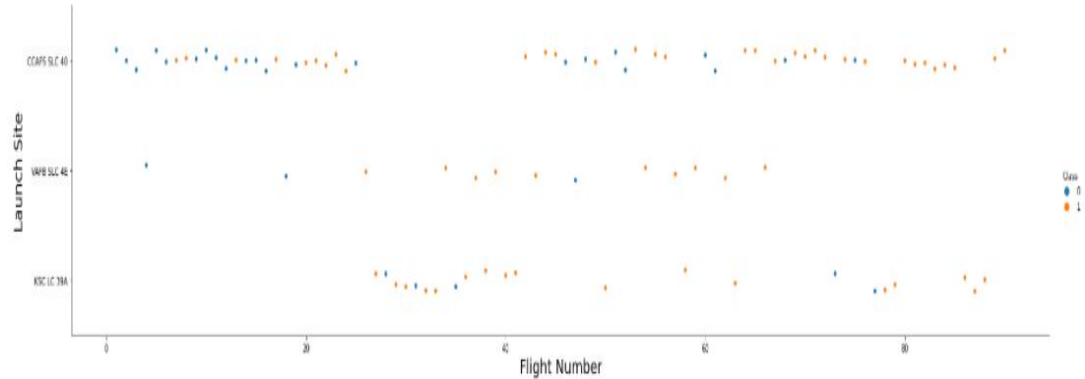
# Insights drawn
# from EDA

# Flight Number vs. Launch Site

- 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.
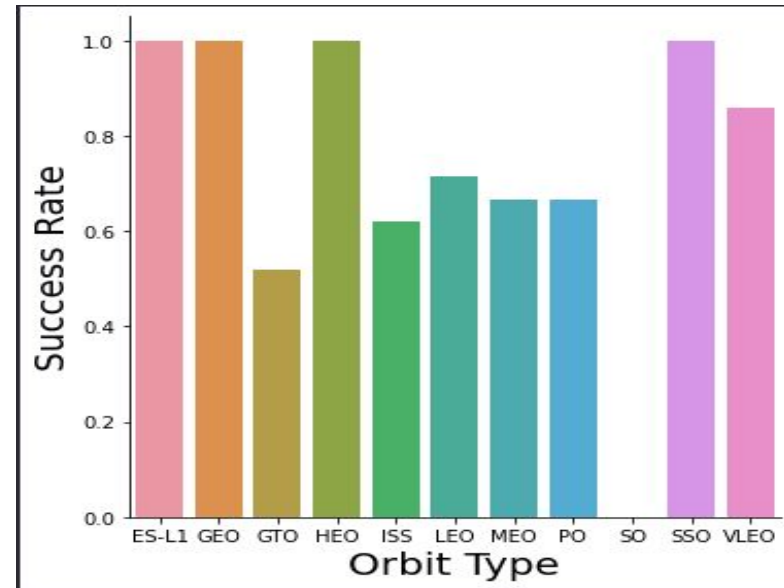- It can be assumed that each new launch has a higher rate of success.

# Payload vs. Launch Site

- 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.
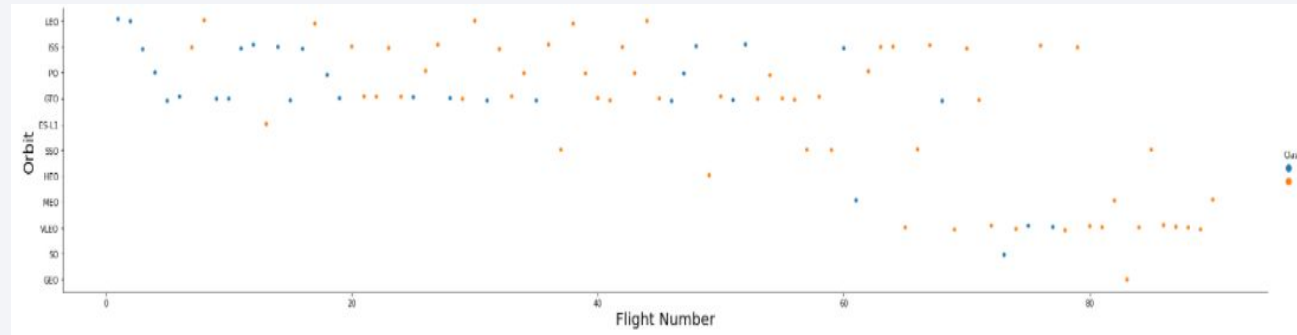
# Success Rate vs. Orbit Type

- ES-L1, GEO, HEO, SSO, VLEO had the most success rate.
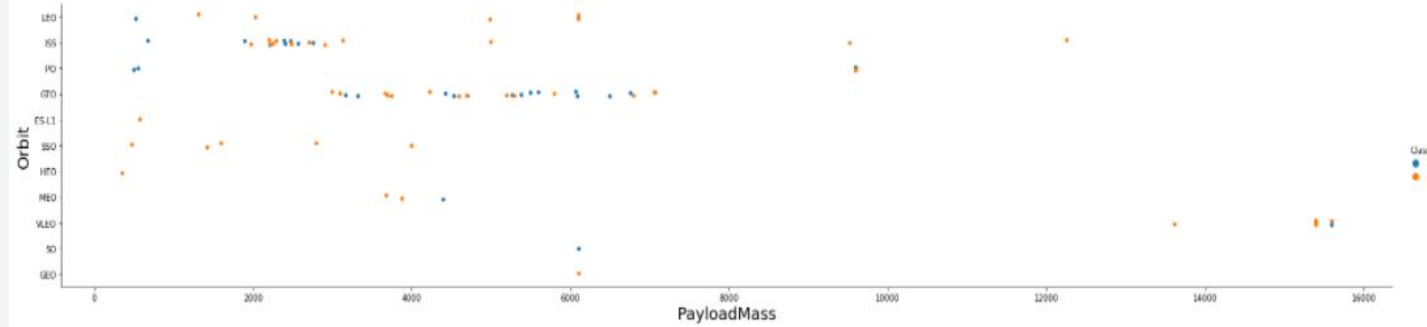
# Flight Number vs. Orbit Type

- There is no relationship between flight number and the orbit.

# Payload vs. Orbit Type

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

# Launch Success Yearly Trend

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



Plot of launch success yearly trend

# All Launch Site Names

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

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

```
In [10]:   task_1 = '''
                SELECT DISTINCT LaunchSite
                FROM SpaceX
           '''
           create_pandas_df(task_1, database=conn)
```

Out[10]:

| | launchsite |
|---|---|
| 0 | KSC LC-39A |
| 1 | CCAFS LC-40 |
| 2 | CCAFS SLC-40 |
| 3 | VAFB SLC-4E |

# Launch Site Names Begin with 'CCA'

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

```
In [11]:  task_2 = '''
              SELECT *
              FROM SpaceX
              WHERE LaunchSite LIKE 'CCA%'
              LIMIT 5
              '''
          create_pandas_df(task_2, database=conn)
```

Out[11]:

|   | date | time | boosterversion | launchsite | payload | payloadmasskg | orbit | customer | missionoutcome | landingoutcome |
|---|------|------|----------------|------------|---------|---------------|-------|----------|----------------|----------------|
| 0 | 2010-04-06 | 18:45:00 | F9 v1.0 B0003 | CCAFS LC-40 | Dragon Spacecraft Qualification Unit | 0 | LEO | SpaceX | Success | Failure (parachute) |
| 1 | 2010-08-12 | 15:43:00 | F9 v1.0 B0004 | CCAFS LC-40 | Dragon demo flight C1, two CubeSats, barrel of... | 0 | LEO (ISS) | NASA (COTS) NRO | Success | Failure (parachute) |
| 2 | 2012-05-22 | 07:44:00 | F9 v1.0 B0005 | CCAFS LC-40 | Dragon demo flight C2 | 525 | LEO (ISS) | NASA (COTS) | Success | No attempt |
| 3 | 2012-08-10 | 00:35:00 | F9 v1.0 B0006 | CCAFS LC-40 | SpaceX CRS-1 | 500 | LEO (ISS) | NASA (CRS) | Success | No attempt |
| 4 | 2013-01-03 | 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

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

```
In [12]:   task_3 = '''
               SELECT SUM(PayloadMassKG) AS Total_PayloadMass
               FROM SpaceX
               WHERE Customer LIKE 'NASA (CRS)'
               '''
           create_pandas_df(task_3, database=conn)
```

Out[12]:

| | total_payloadmass |
|---|---|
| 0 | 45596 |

# Average Payload Mass by F9 v1.1

- Average payload mass carried by booster version F9 v1.1 as 2928.4

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

```
In [13]:   task_4 = '''
               SELECT AVG(PayloadMassKG) AS Avg_PayloadMass
               FROM SpaceX
               WHERE BoosterVersion = 'F9 v1.1'
               '''
           create_pandas_df(task_4, database=conn)
```

Out[13]:

| | avg_payloadmass |
|---|---|
| 0 | 2928.4 |

# 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 [14]:  task_5 = '''
            SELECT MIN(Date) AS FirstSuccessfull_landing_date
            FROM SpaceX
            WHERE LandingOutcome LIKE 'Success (ground pad)'
            '''
          create_pandas_df(task_5, database=conn)
```

```
Out[14]:     firstsuccessfull_landing_date
          0                    2015-12-22
```

# Successful Drone Ship Landing with Payload between 4000 and 6000

```
In [15]:   task_6 = '''
               SELECT BoosterVersion
               FROM SpaceX
               WHERE LandingOutcome = 'Success (drone ship)'
                   AND PayloadMassKG > 4000
                   AND PayloadMassKG < 6000
               '''
           create_pandas_df(task_6, database=conn)
```

```
Out[15]:      boosterversion
           0     F9 FT B1022
           1     F9 FT B1026
           2    F9 FT B1021.2
           3    F9 FT B1031.2
```

- **WHERE** clause filters boosters which have successfully landed on drone ship. Applied the **AND** condition to determine successful landing with 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 [16]:   task_7a = '''
               SELECT COUNT(MissionOutcome) AS SuccessOutcome
               FROM SpaceX
               WHERE MissionOutcome LIKE 'Success%'
               '''

           task_7b = '''
               SELECT COUNT(MissionOutcome) AS FailureOutcome
               FROM SpaceX
               WHERE MissionOutcome LIKE 'Failure%'
               '''
           print('The total number of successful mission outcome is:')
           display(create_pandas_df(task_7a, database=conn))
           print()
           print('The total number of failed mission outcome is:')
           create_pandas_df(task_7b, database=conn)
```

The total number of successful mission outcome is:

|   | successoutcome |
|---|----------------|
| 0 | 100            |

The total number of failed mission outcome is:

Out[16]:

|   | failureoutcome |
|---|----------------|
| 0 | 1              |

- Wildcard like '%' to filter for **WHERE** MissionOutcome was a success or a failure.

30

# Boosters Carried Maximum Payload

Subquery in the **WHERE** clause and the **MAX()** function returns the booster that have carried the maximum payload

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

In [17]:
```
task_8 = '''
        SELECT BoosterVersion, PayloadMassKG
        FROM SpaceX
        WHERE PayloadMassKG = (
                                SELECT MAX(PayloadMassKG)
                                FROM SpaceX
                                )
        ORDER BY BoosterVersion
        '''
create_pandas_df(task_8, database=conn)
```

Out[17]:

| | boosterversion | payloadmasskg |
|---|---|---|
| 0 | F9 B5 B1048.4 | 15600 |
| 1 | F9 B5 B1048.5 | 15600 |
| 2 | F9 B5 B1049.4 | 15600 |
| 3 | F9 B5 B1049.5 | 15600 |
| 4 | F9 B5 B1049.7 | 15600 |
| 5 | F9 B5 B1051.3 | 15600 |
| 6 | F9 B5 B1051.4 | 15600 |
| 7 | F9 B5 B1051.6 | 15600 |
| 8 | F9 B5 B1056.4 | 15600 |
| 9 | F9 B5 B1058.3 | 15600 |
| 10 | F9 B5 B1060.2 | 15600 |
| 11 | F9 B5 B1060.3 | 15600 |

# 2015 Launch Records

- Filter for failed landing outcomes in drone ship, their booster versions, and launch site names for year 2015

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

```
In [18]:    task_9 = '''
                SELECT BoosterVersion, LaunchSite, LandingOutcome
                FROM SpaceX
                WHERE LandingOutcome LIKE 'Failure (drone ship)'
                    AND Date BETWEEN '2015-01-01' AND '2015-12-31'
                '''
            create_pandas_df(task_9, database=conn)
```

| | boosterversion | launchsite | landingoutcome |
|---|---|---|---|
| 0 | F9 v1.1 B1012 | CCAFS LC-40 | Failure (drone ship) |
| 1 | F9 v1.1 B1015 | CCAFS LC-40 | Failure (drone ship) |

# 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))

```python
In [19]:  task_10 = '''
              SELECT LandingOutcome, COUNT(LandingOutcome)
              FROM SpaceX
              WHERE DATE BETWEEN '2010-06-04' AND '2017-03-20'
              GROUP BY LandingOutcome
              ORDER BY COUNT(LandingOutcome) DESC
              '''
          create_pandas_df(task_10, database=conn)
```

Out[19]:

| | landingoutcome | count |
|---|---|---|
| 0 | No attempt | 10 |
| 1 | Success (drone ship) | 6 |
| 2 | Failure (drone ship) | 5 |
| 3 | Success (ground pad) | 5 |
| 4 | Controlled (ocean) | 3 |
| 5 | Uncontrolled (ocean) | 2 |
| 6 | Precluded (drone ship) | 1 |
| 7 | Failure (parachute) | 1 |

Landing outcomes and the **COUNT** of landing outcomes **BETWEEN** 2010-06-04 to 2010-03-20.

**GROUP BY** clause to group the landing outcomes and the **ORDER BY** clause to order the grouped landing outcome in descending order.
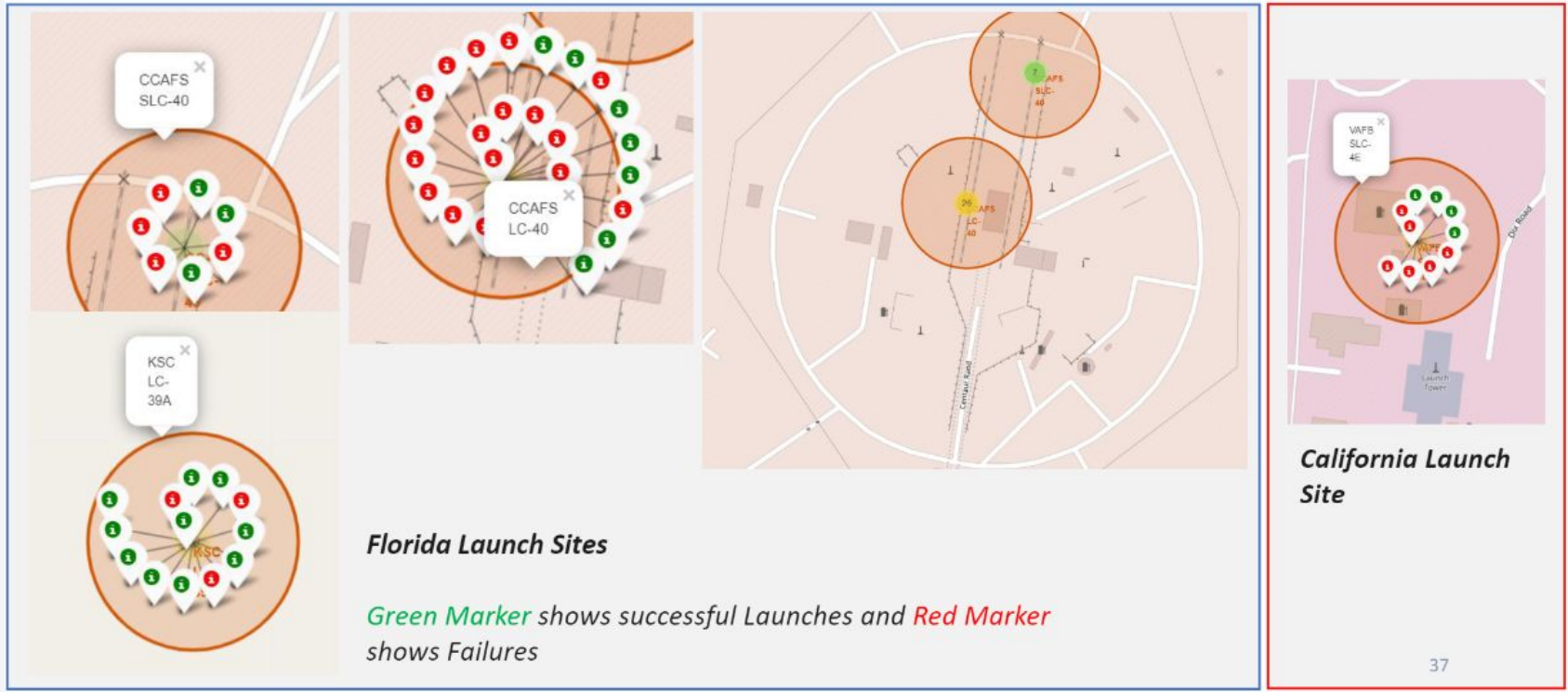
Section 4

# Launch Sites Proximities Analysis

# All launch sites global map markers



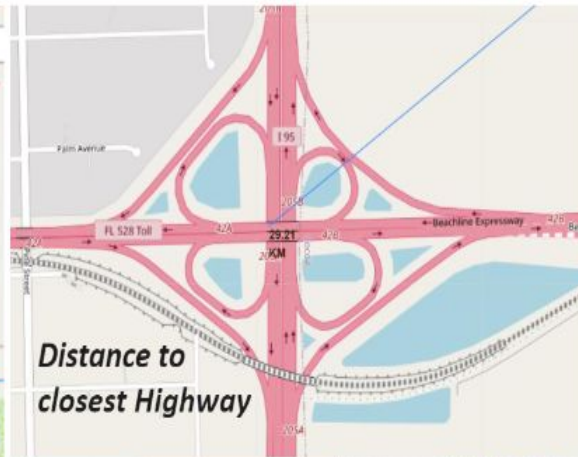We can see that the SpaceX launch sites are in the United States of America coasts. Florida and California

# Markers showing launch sites with color labels



**Florida Launch Sites**

Green Marker shows successful Launches and Red Marker shows Failures

**California Launch Site**

37

36

# Launch Site distance to landmarks



Distance to Railway Station

Distance to closest Highway

Distance to Coastline

Distance to City

Distance to coast

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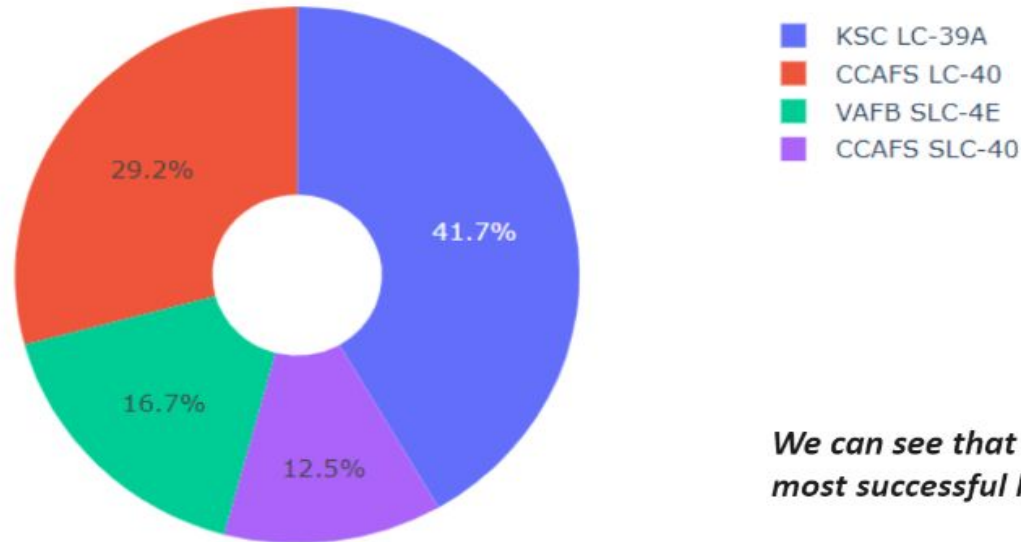•Do launch sites keep certain distance away from cities? Yes

Section 5

# Build a Dashboard
# with Plotly Dash

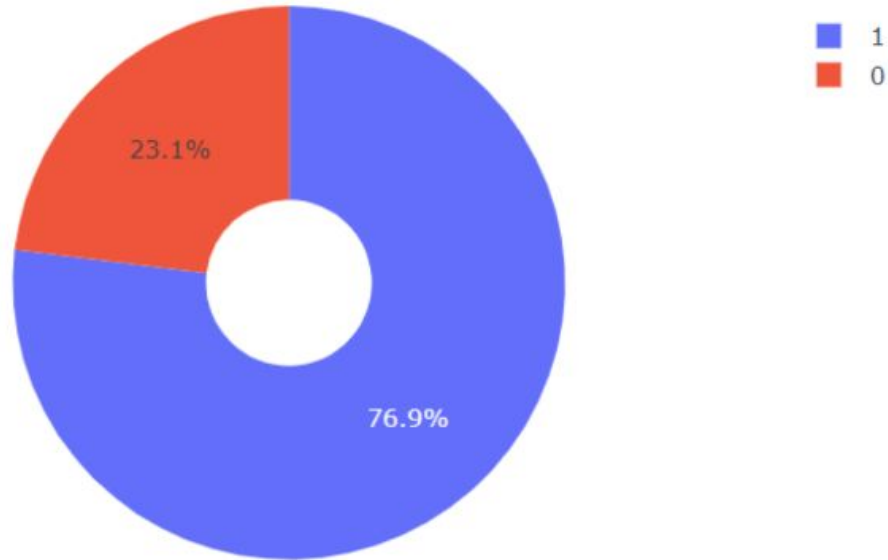# Pie chart showing the success percentage achieved by each launch site

**Total Success Launches By all sites**



Legend:
- KSC LC-39A
- CCAFS LC-40
- VAFB SLC-4E
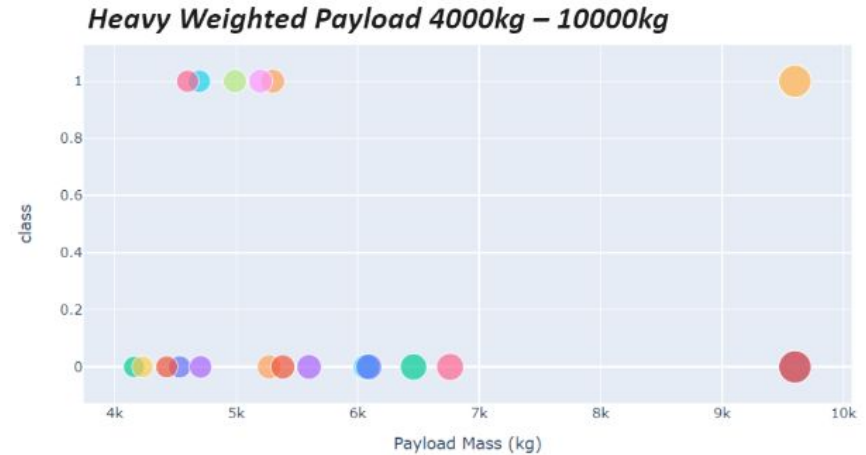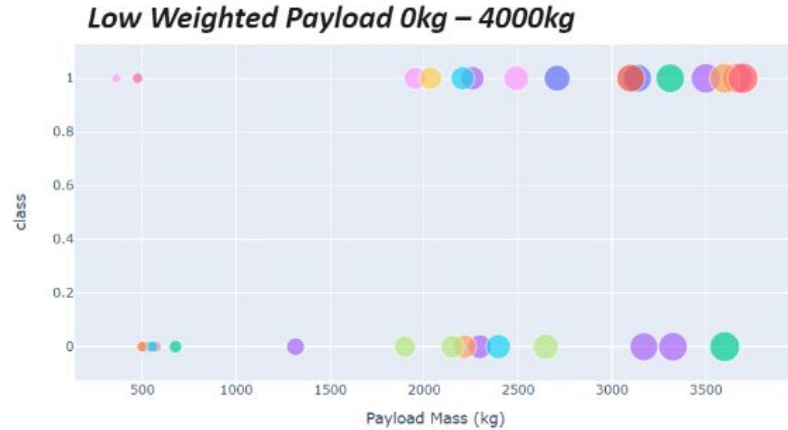- CCAFS SLC-40

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 the sites*

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



KSC LC-39A achieved 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

# Predictive Analysis (Classification)

# Classification Accuracy

- The decision tree classifier is the model with the highest classification accuracy

```python
# Examining the scores from the whole Dataset
jaccard_scores = [
                jaccard_score(Y, logreg_cv.predict(X), average='binary'),
                jaccard_score(Y, svm_cv.predict(X), average='binary'),
                jaccard_score(Y, tree_cv.predict(X), average='binary'),
                jaccard_score(Y, knn_cv.predict(X), average='binary'),
                ]

f1_scores = [
            f1_score(Y, logreg_cv.predict(X), average='binary'),
            f1_score(Y, svm_cv.predict(X), average='binary'),
            f1_score(Y, tree_cv.predict(X), average='binary'),
            f1_score(Y, knn_cv.predict(X), average='binary'),
            ]

accuracy = [logreg_cv.score(X, Y), svm_cv.score(X, Y), tree_cv.score(X, Y), knn_cv.score(X, Y)]

scores = pd.DataFrame(np.array([jaccard_scores, f1_scores, accuracy]),
                    index=['Jaccard_Score', 'F1_Score', 'Accuracy'],
                    columns=['LogReg', 'SVM', 'Tree', 'KNN'])
scores
```
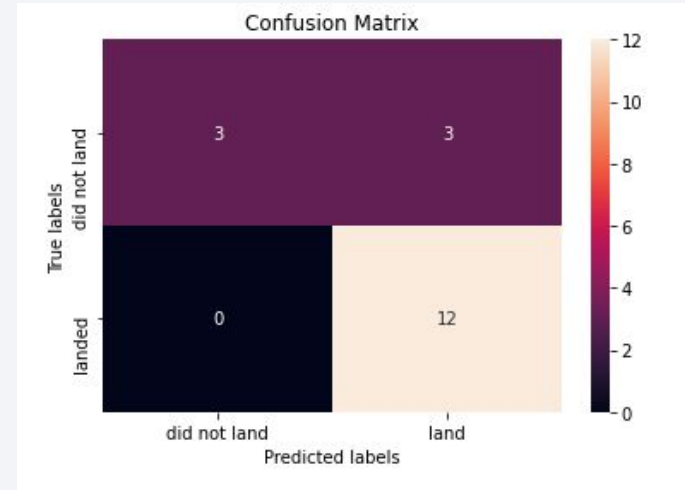
[31] ✓ 0.4s                                                                    Python

...

|               | LogReg   | SVM      | Tree     | KNN      |
|---------------|----------|----------|----------|----------|
| Jaccard_Score | 0.833333 | 0.845070 | 0.880597 | 0.819444 |
| F1_Score      | 0.909091 | 0.916031 | 0.936508 | 0.900763 |
| Accuracy      | 0.866667 | 0.877778 | 0.911111 | 0.855556 |

# 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

- 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!