

Winning Space Race with Data Science

I Gusti Ngruah Sentana Putra 29th Maret 2024



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 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 you want to find answers

- 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.



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/chuksoo/IBM-Data-Science-Capstone-SpaceX/blob/main/Data%20Collect ion%20API.ipynb.

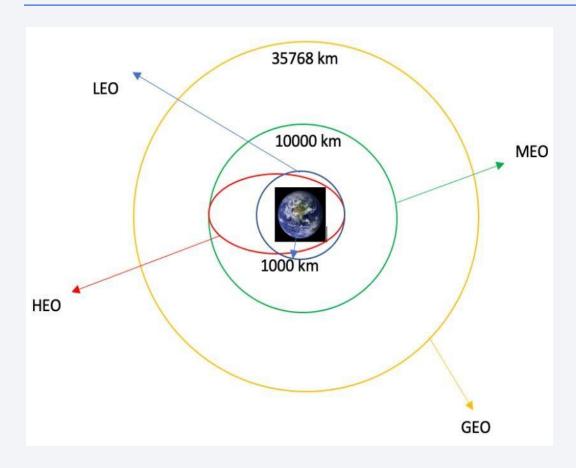
```
1. Get request for rocket launch data using API
          spacex url="https://api.spacexdata.com/v4/launches/past"
          response = requests.get(spacex url)
   2. Use json normalize method to convert json result to dataframe
In [12]:
           # Use json normalize method to convert the json result into a dataframe
           # decode response content as json
           static json df = res.json()
           # apply json normalize
           data = pd.json normalize(static json df)
   3. We then performed data cleaning and filling in the missing values
In [30]:
          rows = data falcon9['PayloadMass'].values.tolist()[0]
           df rows = pd.DataFrame(rows)
           df rows = df rows.replace(np.nan, PayloadMass)
          data falcon9['PayloadMass'][0] = df rows.values
           data falcon9
```

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/chuksoo/IBM-Data-Science-Capstone-SpaceX/blob/main/Data%20Collect ion%20with%20Web%20Scraping .ipynb.

```
1. Apply HTTP Get method to request the Falcon 9 rocket launch page
        static url = "https://en.wikipedia.org/w/index.php?title=List of Falcon 9 and Falcon Heavy launches&oldid=1927686922"
In [5]: # 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
Out[5]: 200
    2. Create a BeautifulSoup object from the HTML response
           # Use BeautifulSoup() to create a BeautifulSoup object from a response text content
           soup = BeautifulSoup(html_data.text, 'html.parser')
          Print the page title to verify if the BeautifulSoup object was created properly
           # Use soup.title attribute
           soup.title
          <title>List of Falcon 9 and Falcon Heavy launches - Wikipedia</title>
        Extract all column names from the HTML table header
          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 get a column name
          # Append the Non-empty column name ('if name is not None and Len(name) > \theta') into a list called column names
          element = soup.find all('th')
          for row in range(len(element)):
                 name = extract_column_from_header(element[row])
                 if (name is not None and len(name) > 0);
                     column names.append(name)
    4. Create a dataframe by parsing the launch HTML tables
    Export data to csv
```

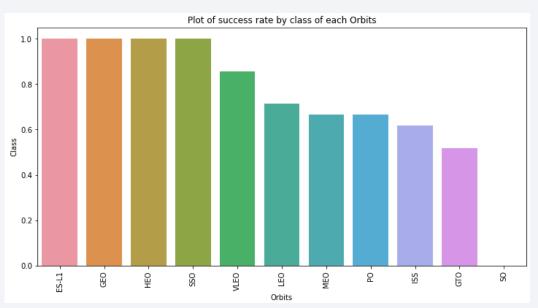
Data Wrangling

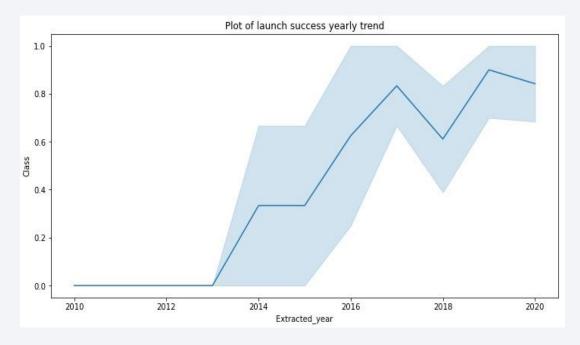


- We performed exploratory data analysis and determined the training labels.
- We calculated the number of launches at each site, and the number and occurrence of each orbits
- We created landing outcome label from outcome column and exported the results to csv.
- The link to the notebook is https://github.com/chuksoo/IBM-Data-Science-Capstone-SpaceX/blob/main/D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

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/chuksoo/IBM-Data-Science-Capstone-SpaceX/blob/main/EDA%20with%20D ata%20Visualization.ipynb

EDA with SQL

- We loaded the SpaceX dataset into a PostgreSQL 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/chuksoo/IBM-Data-Science-Capstone-SpaceX/blob/main/EDA%20with%20SQL.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.

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/chuksoo/IBM-Data-Science-Capstone-SpaceX/blob/main/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/chuksoo/IBM-Data-Science-Capstone-SpaceX/blob/main/Machine%20Learning%20Prediction.ipynb

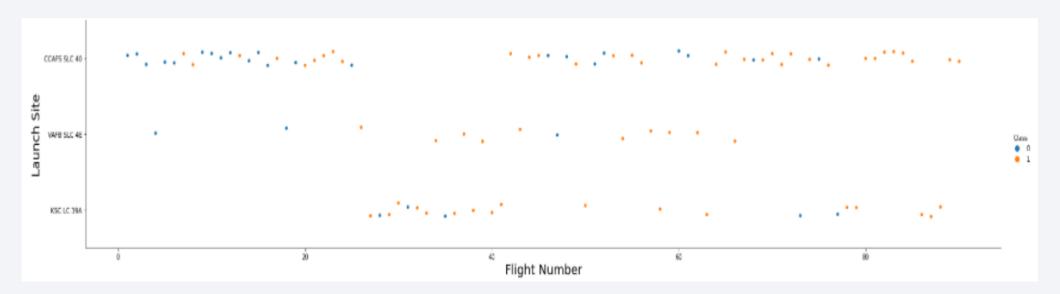
Results

- Exploratory data analysis results
- Interactive analytics demo in screenshots
- Predictive analysis results



Flight Number vs. Launch Site

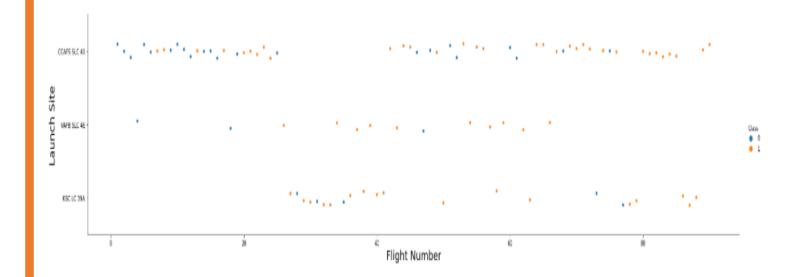
• From the plot, we found that the larger the flight amount at a launch site, the greater the success rate at a launch site.



Payload vs. Launch Site

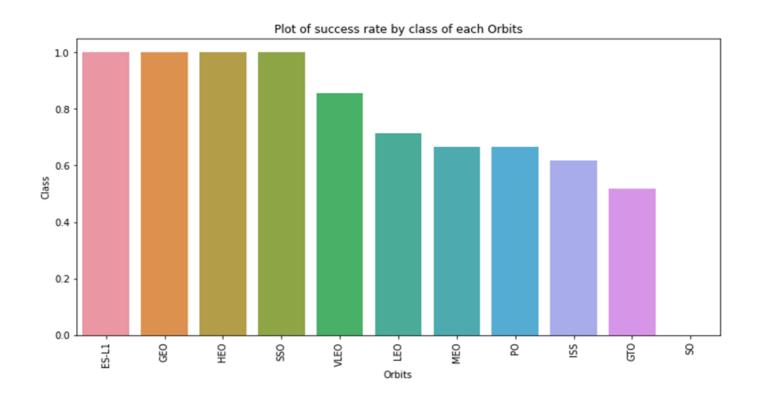


The greater the payload mass for launch site CCAFS SLC 40 the higher the success rate for the rocket.



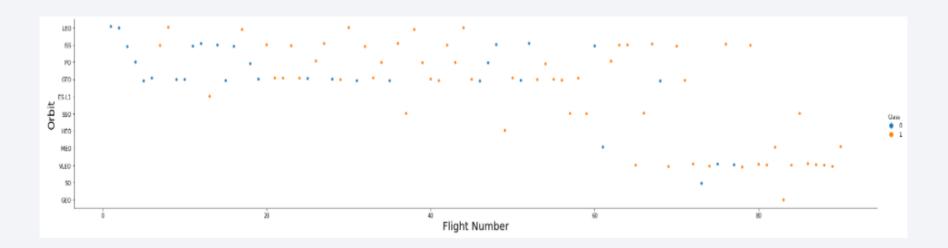
Success Rate vs. Orbit Type

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



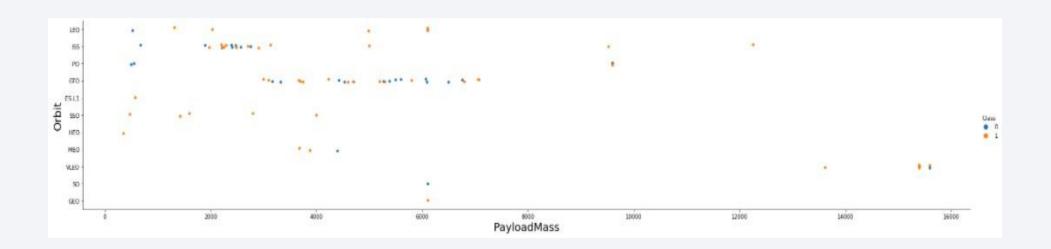
Flight Number vs. Orbit Type

 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.



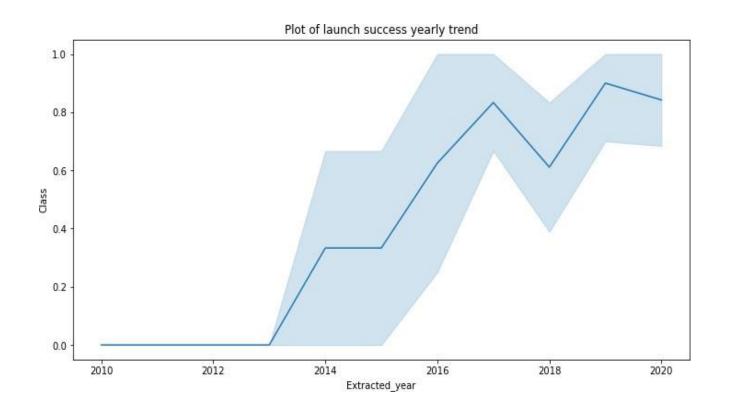
Payload vs. Orbit Type

 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.

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

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'

[11]:		FROM WHEN	ECT * 1 SpaceX RE Launc IT 5	hSite LIKE 'CC							
t[11]:		date	time	boosterversion	launchsite	payload	payloadmasskg	orbit	customer	missionoutcome	landingoutcom
	0	2010-04- 06	18:45:00	F9 v1.0 B0003	CCAFS LC- 40	Dragon Spacecraft Qualification Unit	0	LEO	SpaceX	Success	Failu (parachut
	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	Failu (parachut
	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 attem
	2	2012-05- 22 2012-08- 10	07:44:00	F9 v1.0 B0005		Dragon demo flight C2 SpaceX CRS-1	525 500	LEO (ISS) LEO (ISS)	NASA (COTS)	Success Success	No attem

• 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

 We calculated the average payload mass carried by booster version F9 v1.1 as 2928.4

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

```
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 22nd December 2015

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)
```

 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

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

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
                      100
         The total number of failed mission outcome is:
Out[16]:
            failureoutcome
```

 We used wildcard like '%' to filter for WHERE MissionOutcome was a success or a failure.

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

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

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



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))

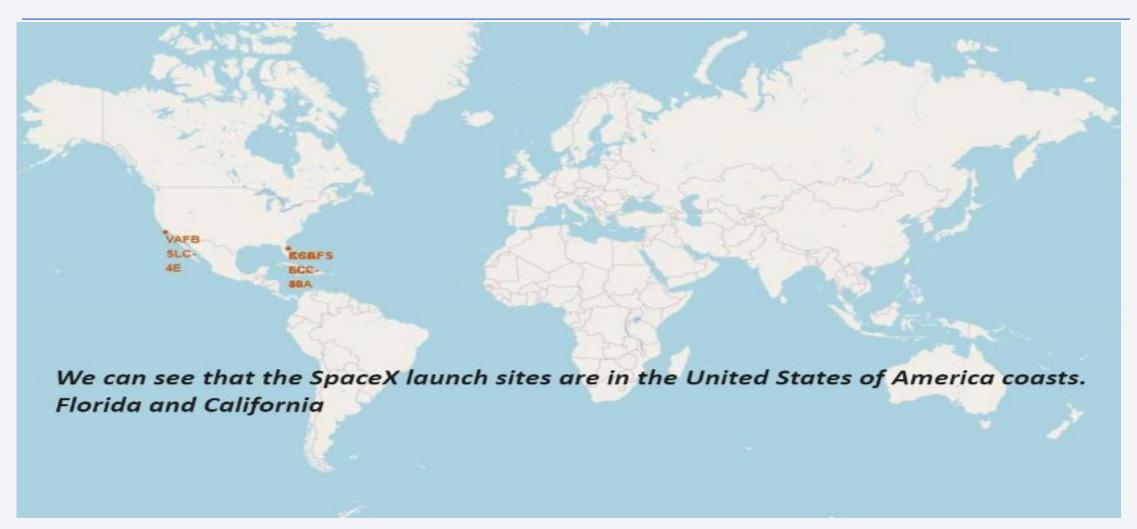
```
In [19]:
    task_10 = '''
        SELECT LandingOutcome, COUNT(LandingOutcome)
        FROM SpaceX
        WHERE DATE BETWEEN '2010-06-04' AND '2017-03-20'
        GROUP BY LandingOutcome
        ORDER BY COUNT(LandingOutcome) DESC
        '''
    create_pandas_df(task_10, database=conn)
```

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

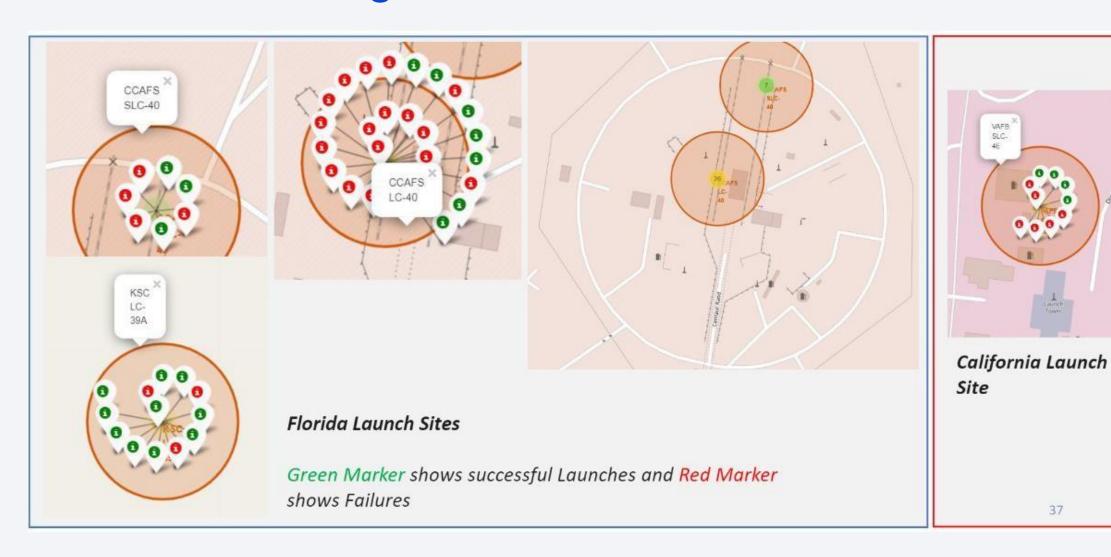
- We selected Landing outcomes and the COUNT of landing outcomes from the data and used the WHERE clause to filter for landing outcomes BETWEEN 2010-06-04 to 2010-03-20.
- We applied the GROUP BY clause to group the landing outcomes and the ORDER BY clause to order the grouped landing outcome in descending order.



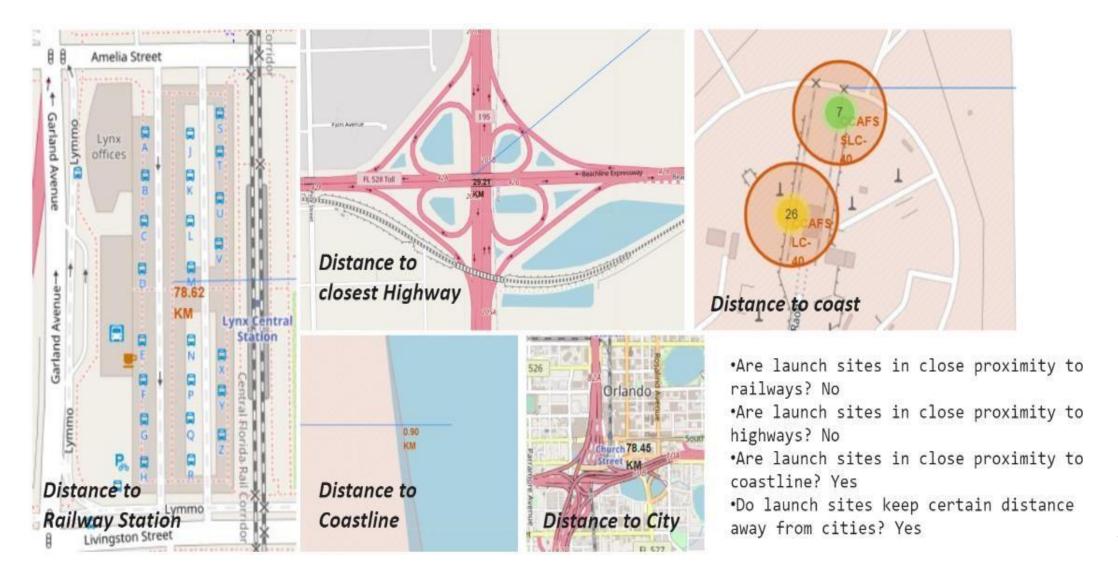
All launch sites global map markers



Markers showing launch sites with color labels



Launch Site distance to landmarks

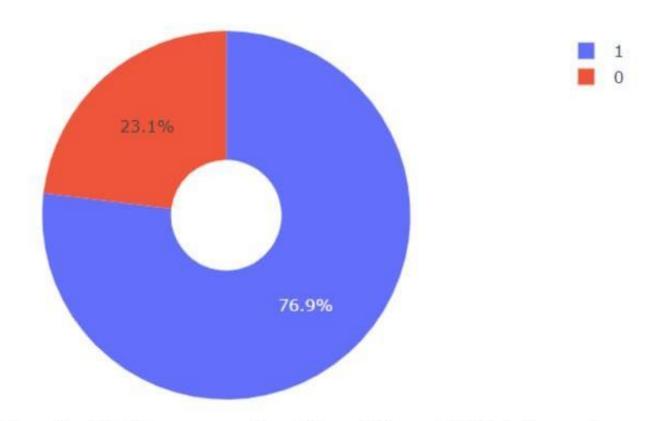




Pie chart showing the success percentage achieved by each launch site

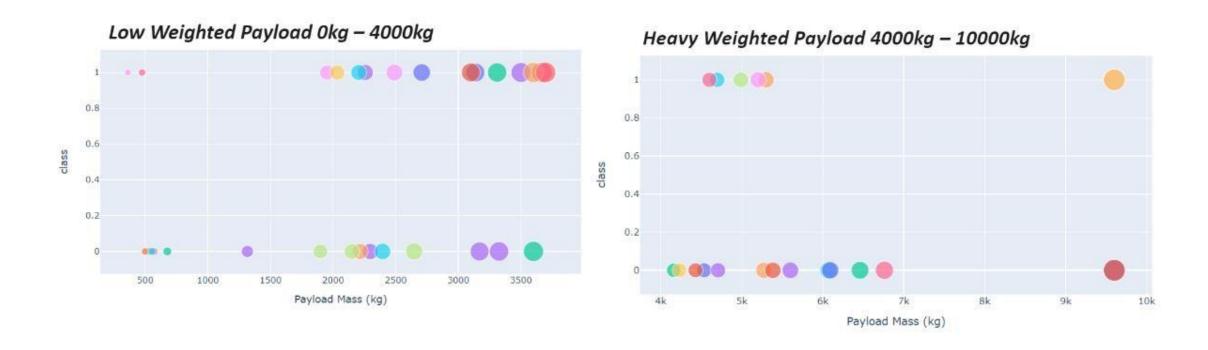


Pie chart showing the Launch site with the highest launch success ratio



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



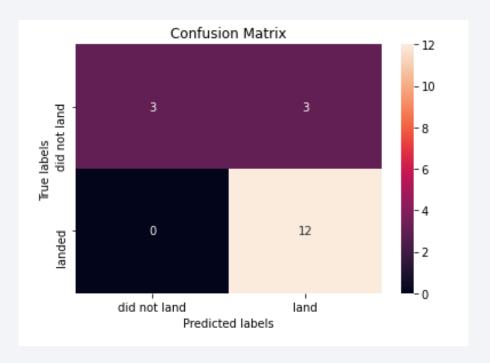
Classification Accuracy

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

```
models = { 'KNeighbors':knn cv.best score ,
               'DecisionTree': tree cv.best score ,
               'LogisticRegression':logreg cv.best score ,
               'SupportVector': svm cv.best score }
bestalgorithm = max(models, key=models.get)
print('Best model is', bestalgorithm,'with a score of', models[bestalgorithm])
if bestalgorithm == 'DecisionTree':
    print('Best params is :', tree cv.best params )
if bestalgorithm == 'KNeighbors':
    print('Best params is :', knn cv.best params )
if bestalgorithm == 'LogisticRegression':
    print('Best params is :', logreg cv.best params )
if bestalgorithm == 'SupportVector':
    print('Best params is :', svm cv.best params )
Best model is DecisionTree with a score of 0.8732142857142856
Best params is : {'criterion': 'gini', 'max_depth': 6, 'max_features': 'auto', 'min_samples_leaf': 2, 'min_samples_split': 5, 'splitter': 'random'}
```

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.
- KSCLC-39A had the most successful launches of any sites.
- The Decision tree classifier is the best machine learning algorithm for this task.

