



IBM Developer
SKILLS NETWORK

Winning Space Race with Data Science

Prajwal Narayanaswamy
17th March 2023



Outline

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

Executive Summary

- Summary of methodologies
 - * Data Collection using API
 - * Data Collection with Web Scrapping
 - * Data Wrangling
 - * Exploratory data analysis using SQL and Data Visualization
 - * Interactive Visual Analytics using Folium
 - * Machine Learning Prediction
- Summary of all results
 - * Exploratory Data Analysis results
 - * Screenshots of Interactive Analysis
 - * Predictive Analysis Results

Introduction

- Project background and context

SpaceX advertises Falcon 9 rocket launch to cost 62 million dollars on its website; while other providers cost upward of 165 million dollars each, much of savings is because SpaceX can reuse the first stage. Therefore, if we determine if the first stage will land, we can determine the cost of the launch. This information can be sued if an alternate company wants to bud against SpaceX for a rocket launch. The 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

- Factors determining if the rocket will land successfully or not
- The interaction amongst various features that determine the success rate of a successful landing.
- Operating conditions required to ensure a successful landing program.

Section 1

Methodology

Methodology

Executive Summary

- Data collection methodology:
 - Data was collected using SpaceX API and web scrapping 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
 - The response was then decoded as JSON using `.json()` function call and convert it into pandas dataframe using `.json_normalize()`
 - The data was then cleaned, checked for missing values and fill in missing values wherever necessary.
 - In addition, we performed web scrapping 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 access the data from SpaceX API, clean it and did some basic wrangling and formatting.
- The link to the notebook is [https://github.com/Prajwal-Narayanaswamy/SpaceX---Project/tree/master/SpaceX Falcon 9 Lab1](https://github.com/Prajwal-Narayanaswamy/SpaceX---Project/tree/master/SpaceX_Falcon_9_Lab1).

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:

```
In [9]: static_json_url='https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DS0321EN-SkillsNetwork/datasets/API_call_spacex_api.json'
```

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

```
In [10]: response.status_code
```

```
Out[10]: 200
```

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

```
In [16]: # Use json_normalize method to convert the json result into a dataframe
jlist = requests.get(static_json_url).json()
df = pd.json_normalize(jlist)
```

Using the dataframe `data` print the first 5 rows

```
In [17]: # Get the head of the dataframe
df.head()
```

```
Out[17]:
```

	static_fire_date_utc	static_fire_date_unix	tbd	net	window	rocket	success	details	crew	ships	capsules	payloads
0	2006-03-17T00:00:00.000Z	1.142554e+09	False	False	0.0	5e9d0d95eda69955f709d1eb	False	Engine failure at 33 seconds and loss of vehicle				[5eb0e4b5b6c3bb0006eeb1e1] 5e9e4502f509f

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/Prajwal-Narayanaswamy/SpaceX---Project/tree/master/SpaceX_Falcon9_Lab1.

TASK 1: Request the Falcon9 Launch Wiki page from its URL

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

```
In [5]: # use requests.get() method with the provided static_url
# assign the response to a object
data = requests.get(static_url).text
```

Create a BeautifulSoup object from the HTML response

```
In [8]: # Use BeautifulSoup() to create a BeautifulSoup object from a response text content
soup = BeautifulSoup(data, 'lxml')
```

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

```
In [9]: # Use soup.title attribute
tag_title = soup.title
tag_string_tag_title = tag_title.string
tag_string_tag_title
```

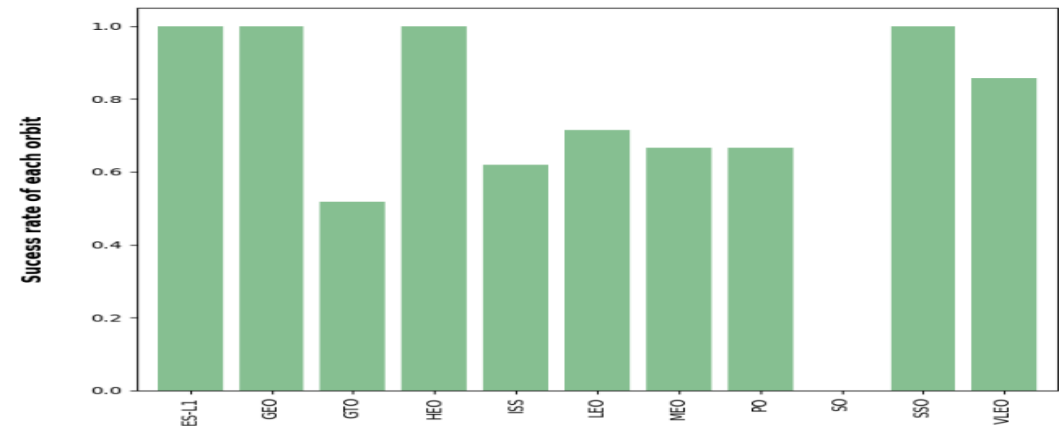
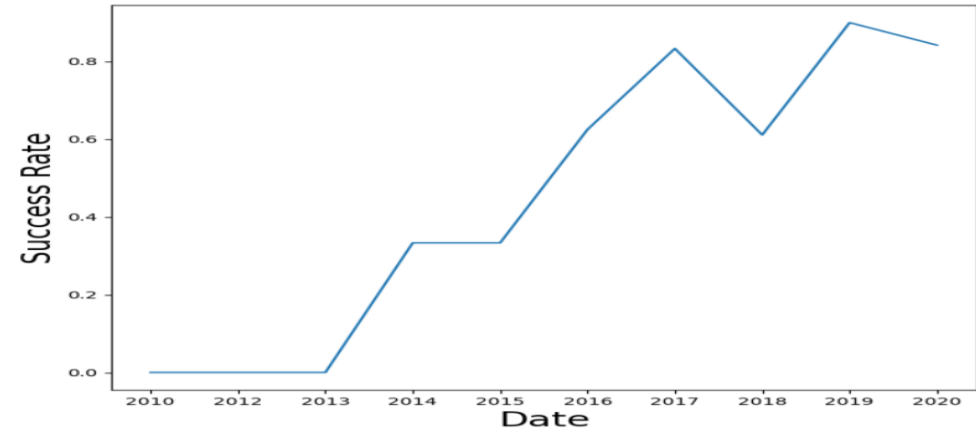
```
Out[9]: 'List of Falcon 9 and Falcon Heavy launches - Wikipedia'
```

Data Wrangling

- Exploratory data analysis was performed and training labels were determined.
- We calculated the number of launches at each site, and the number and occurrences of each orbits.
- We created landing outcome label from outcome column and exported the results to csv.
- The link to notebook is,
<https://github.com/Prajwal-Narayanaswamy/SpaceX---Project/blob/master/SpaceX%20-%20Notebook%20-%203.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/Prajwal-Narayanaswamy/SpaceX--Project/blob/master/SpaceX%20-%20Notebook%205.ipynb>



EDA with SQL

- We loaded the SpaceX dataset into PostgreSQL database without leaving the 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 F(v1.1
 - The total number of successful and failure mission outcomes.
 - The lading outcomes in drone ship, their booster version and launch site names.
- The link to notebook is <https://github.com/Prajwal-Narayanaswamy/SpaceX---Project/blob/master/SpaceX%20-%20Notebook%204.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 rates.
- We calculated the distances between a launch site to its proximities. We answered some questions 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

- An interactive dashboard was created using Plotly dash.
- We plotted pie charts showing total launches by 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/Prajwal-Narayanaswamy/SpaceX---Project/blob/master/SpaceX%20-%20Notebook%205.ipynb>

Predictive Analysis (Classification)

- Summarize how you built, evaluated, improved, and found the best performing classification model
- You need present your model development process using key phrases and flowchart
- Add the GitHub URL of your completed predictive analysis lab, as an external reference and peer-review purpose

Results

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

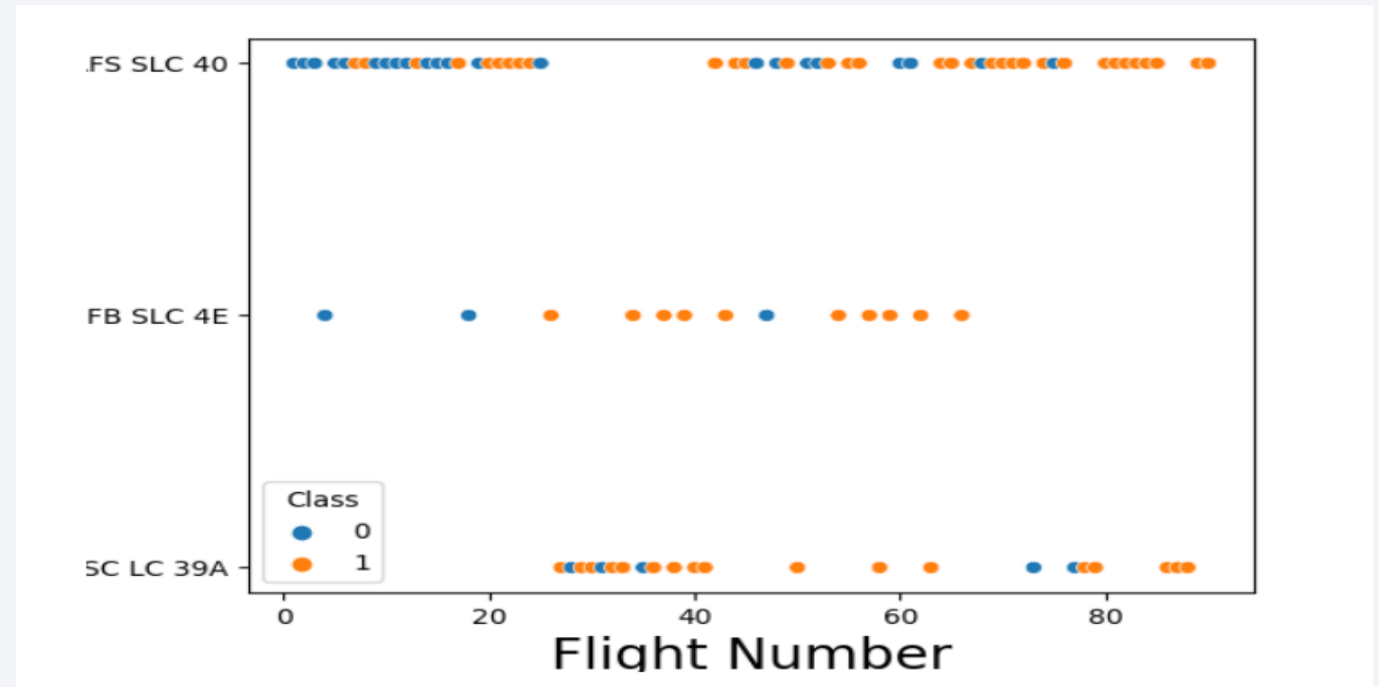
The background of the slide is an abstract composition. It features a dark blue base color. Overlaid on this are numerous diagonal streaks in shades of red and cyan. A faint, light blue grid pattern is also visible, particularly in the lower half of the image. The overall effect is dynamic and technological.

Section 2

Insights drawn from EDA

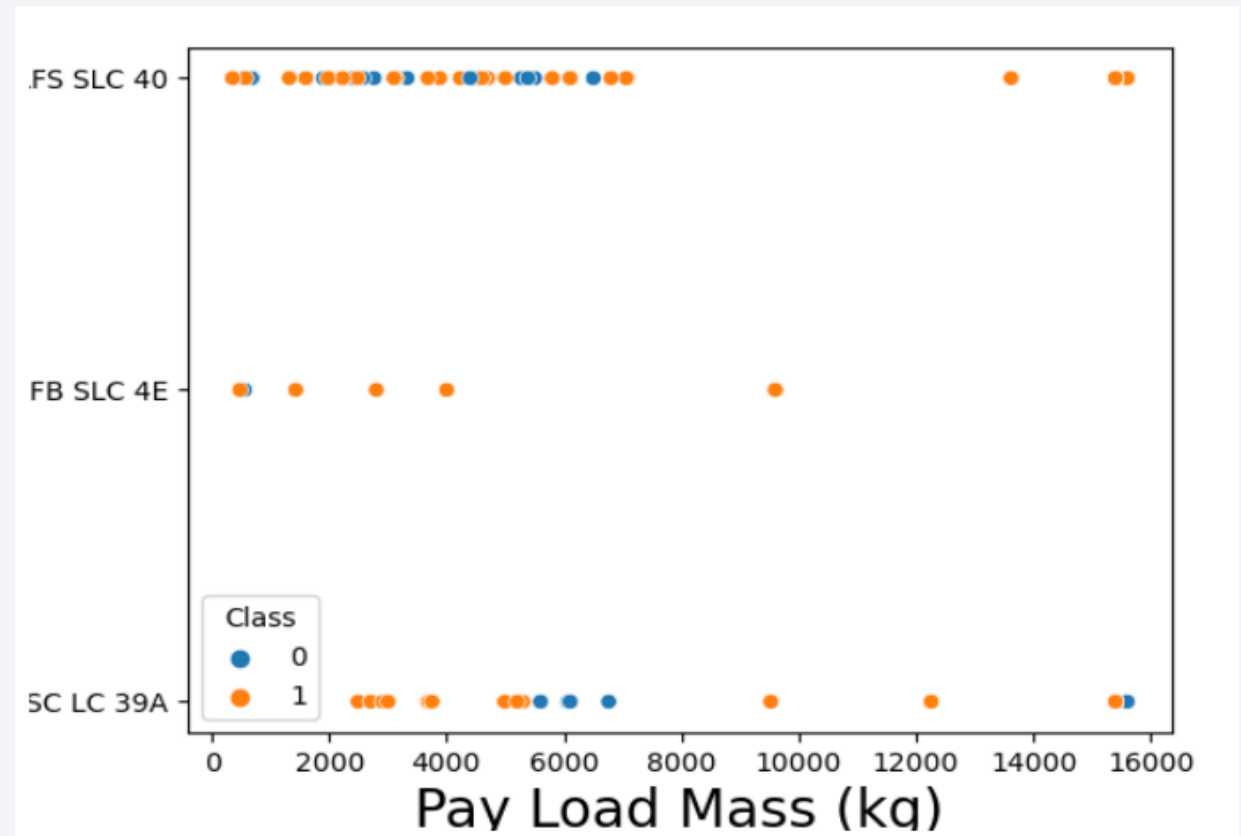
Flight Number vs. Launch Site

- From the plot, it was found that the larger the flight amount at a launch site, the greater the success rate at a launch site.



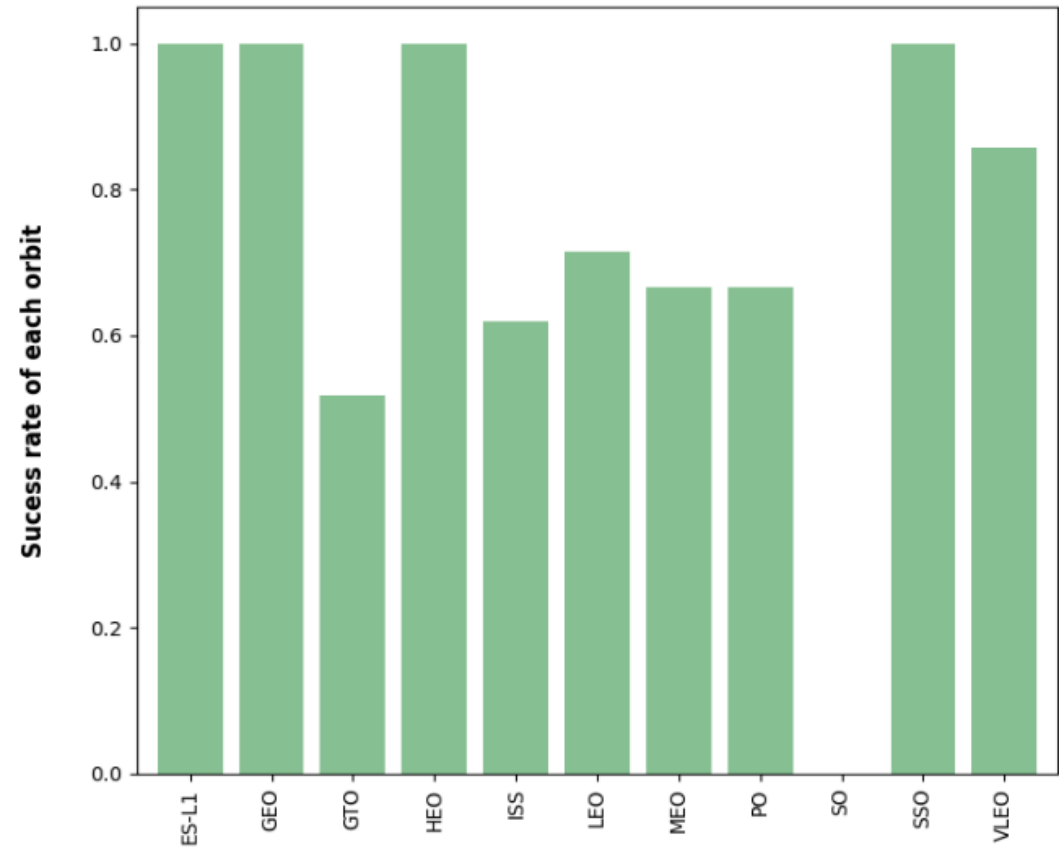
Payload vs. Launch Site

- The lower the payload mass for the launch site CCAFS SLC 40, the higher the success rate for the rocket.



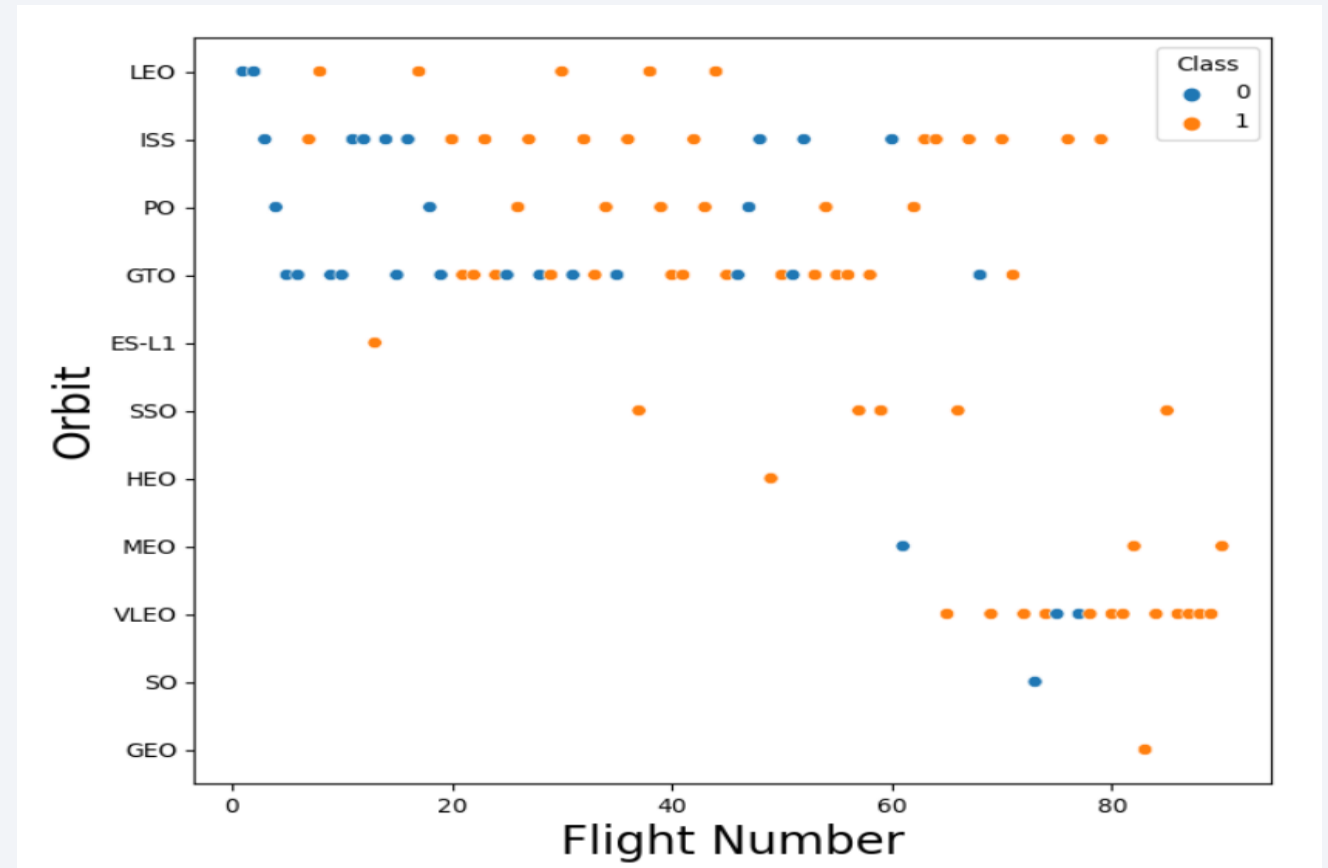
Success Rate vs. Orbit Type

- From the plot, we can see that ES-1, GEO, HEO, SSO, VLEO had most success rate.



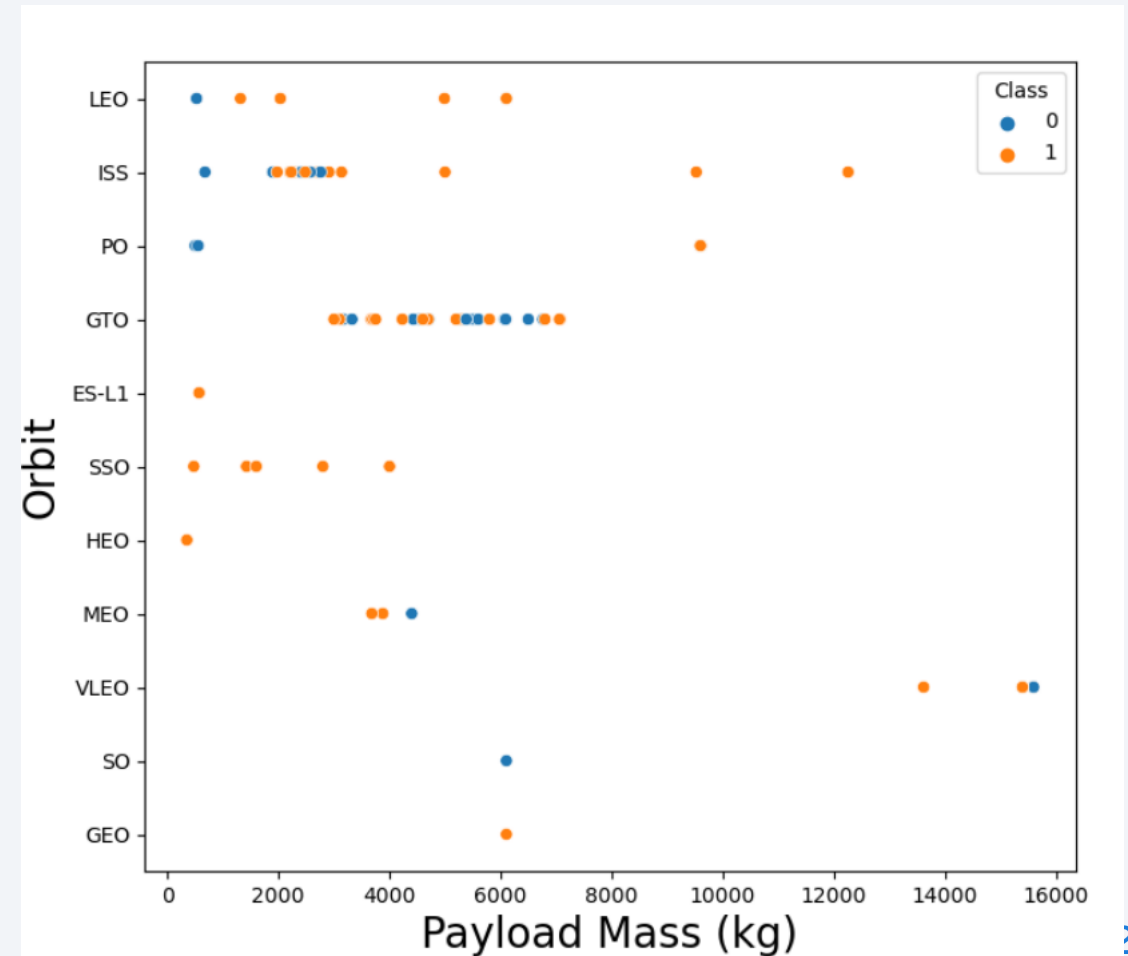
Flight Number vs. Orbit Type

- The plot indicates that 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.



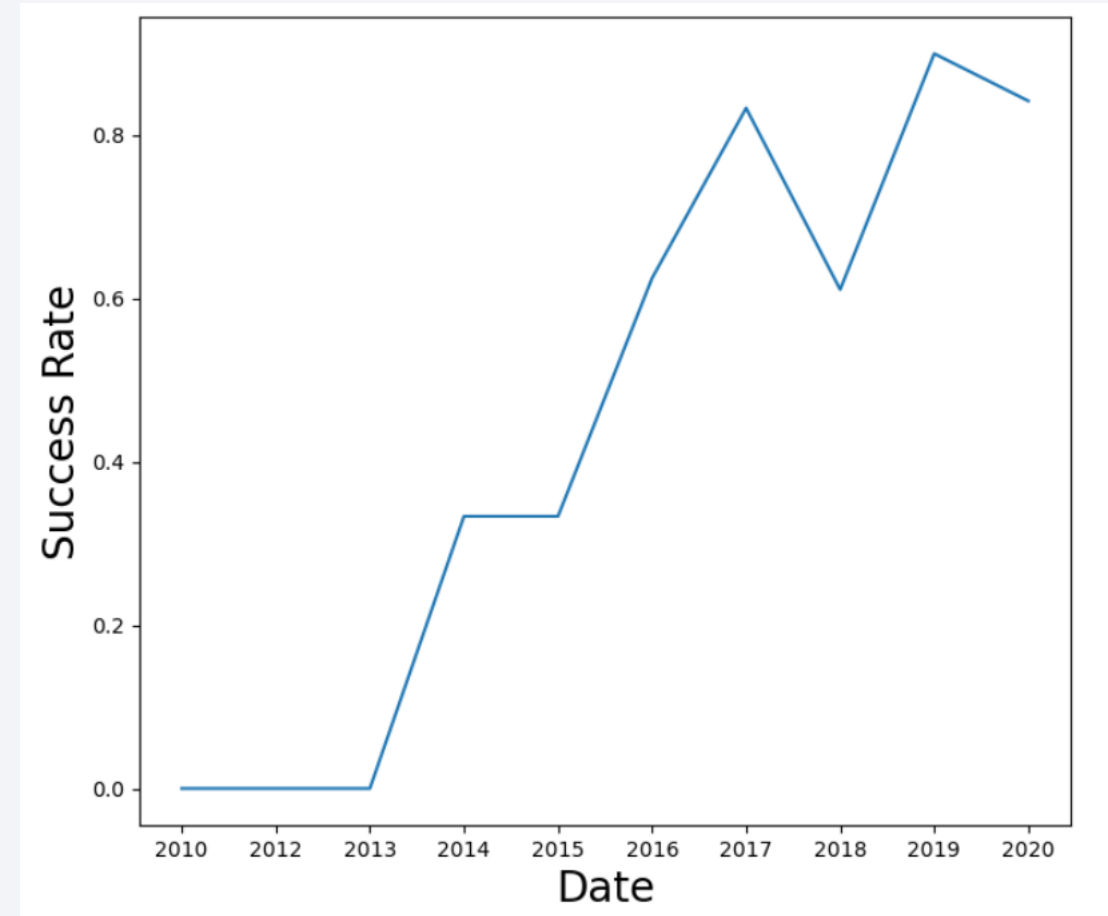
Payload vs. Orbit Type

- 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, it can be observed that success rate since 2013 kept on increasing till 2020.



All Launch Site Names

- We used the keyword 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'

- The below shown query was used to display 5 records where launch sites 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_pendas_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	12:43:00	F9 v1.0 B0004	CCAFS LC-40	Dragon demo flight C1, two CubeSats, barrel of...	0	LEO (ISS)	NASA (COTS) NRC	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

Total Payload Mass

- The total payload carried by boosters from NASA was 45596 and this was calculated 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

- The average payload mass carried by the Booster version F9 v1.1 was 2929.4

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 observe that the dates of the first successful landing outcome on the ground pad was 22nd 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

- WHERE clause was used 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 [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

Total Number of Successful and Failure Mission Outcomes

- We used wildcard like “%” to filter for WHERE mission outcome was a success or a failure.

```
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:
In [16]:  failureoutcome
0              1
```

Boosters Carried Maximum Payload

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

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

```
1x [17]: task_8 = '''
SELECT booster_version, PayloadMassKg
FROM SpaceX
WHERE PayloadMassKg = (
    SELECT MAX(PayloadMassKg)
    FROM SpaceX
)
ORDER BY Booster-Version
'''
create_pandas_df(task_8, database=conn)
```

Out[17]:

	booster_version	payload_mass_kg
0	F9 B5 B1040.4	15800
1	F9 B5 B1040.5	15800
2	F9 B5 B1049.1	15800
3	F9 B5 B1049.5	15800
4	F9 B5 B1049.7	15800
5	F9 B5 B1051.2	15800
6	F9 B5 B1051.4	15800
7	F9 B5 B1051.6	15800
8	F9 B5 B1056.4	15800
9	F9 B5 B1058.3	15800
10	F9 B5 B1060.2	15800
11	F9 B5 B1060.3	15800

2015 Launch Records

- We used a combinations of WHERE clause, LIKE, and BETWEEN conditions 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_3 = """
        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_3, database=conn)
```

```
Out[18]:
```

	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

- 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 2017-03-20
- We applied the GROUPBY clause to group the landing outcomes and the order by clause to order the grouped landing outcome in descending order.

Rank the count of landing outcomes (such as Failure (drone ship) or Success (ground pad))

```
In [18]: task_18 = """
SELECT LandingOutcomes, COUNT(LandingOutcomes)
FROM SpaceX
WHERE DATE BETWEEN '2010-06-04' AND '2017-03-20'
GROUP BY LandingOutcomes
ORDER BY COUNT(LandingOutcome) DESC
"""

create_pandas_df(task_18, database=conn)
```

Out[18]:

	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

A satellite view of Earth from space, showing the curvature of the planet and city lights at night. The background is a deep blue gradient.

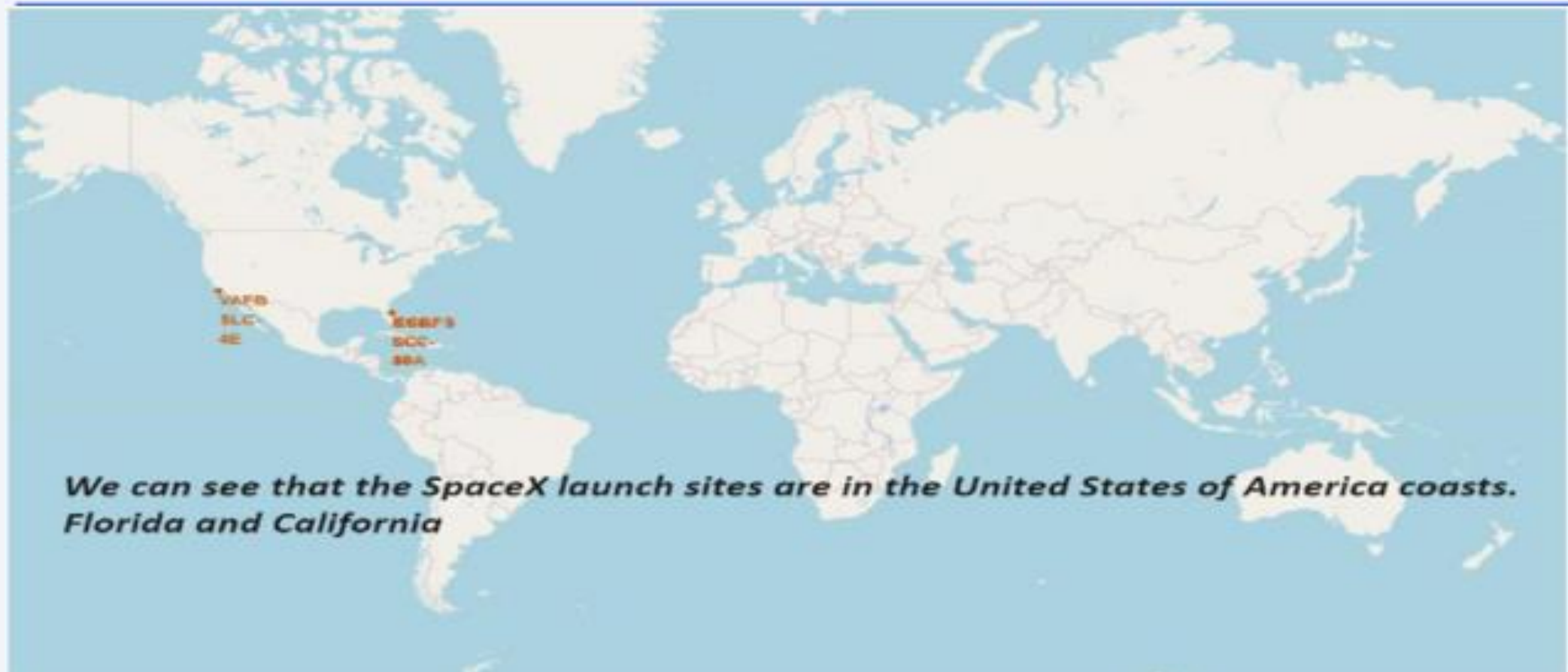
Section 3

Launch Sites Proximities Analysis

Markers showing launch sites with color labels



All launch sites global map markers



Launch Site distance to landmarks



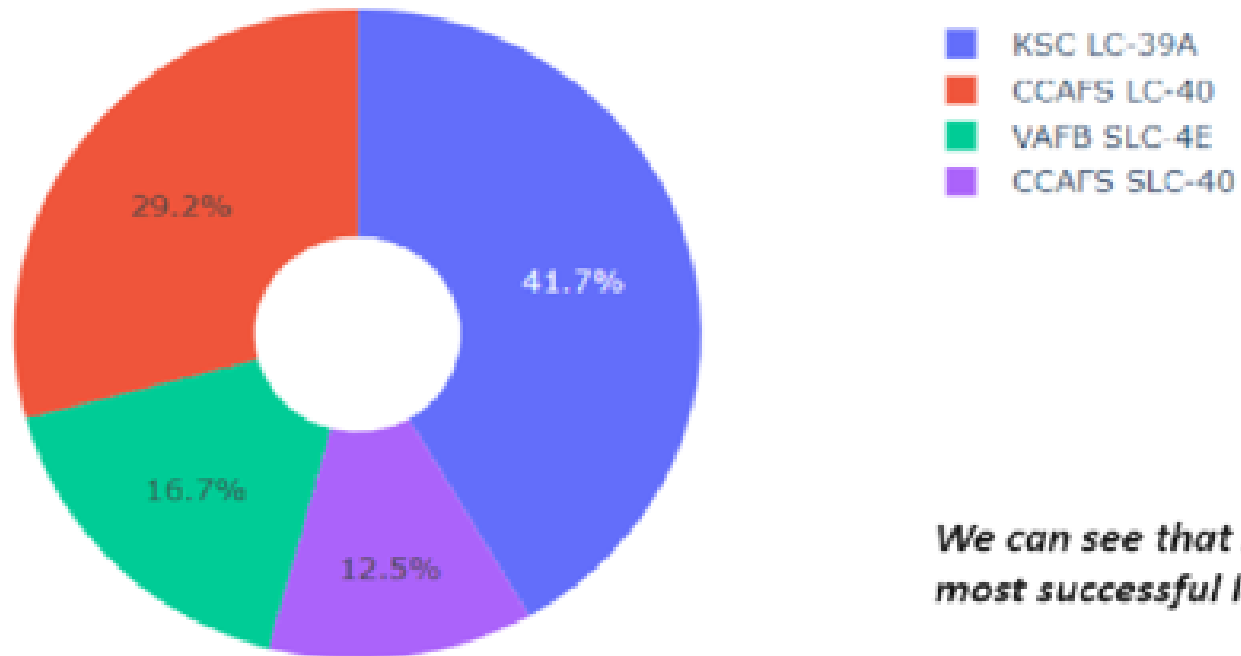


Section 4

Build a Dashboard with Plotly Dash

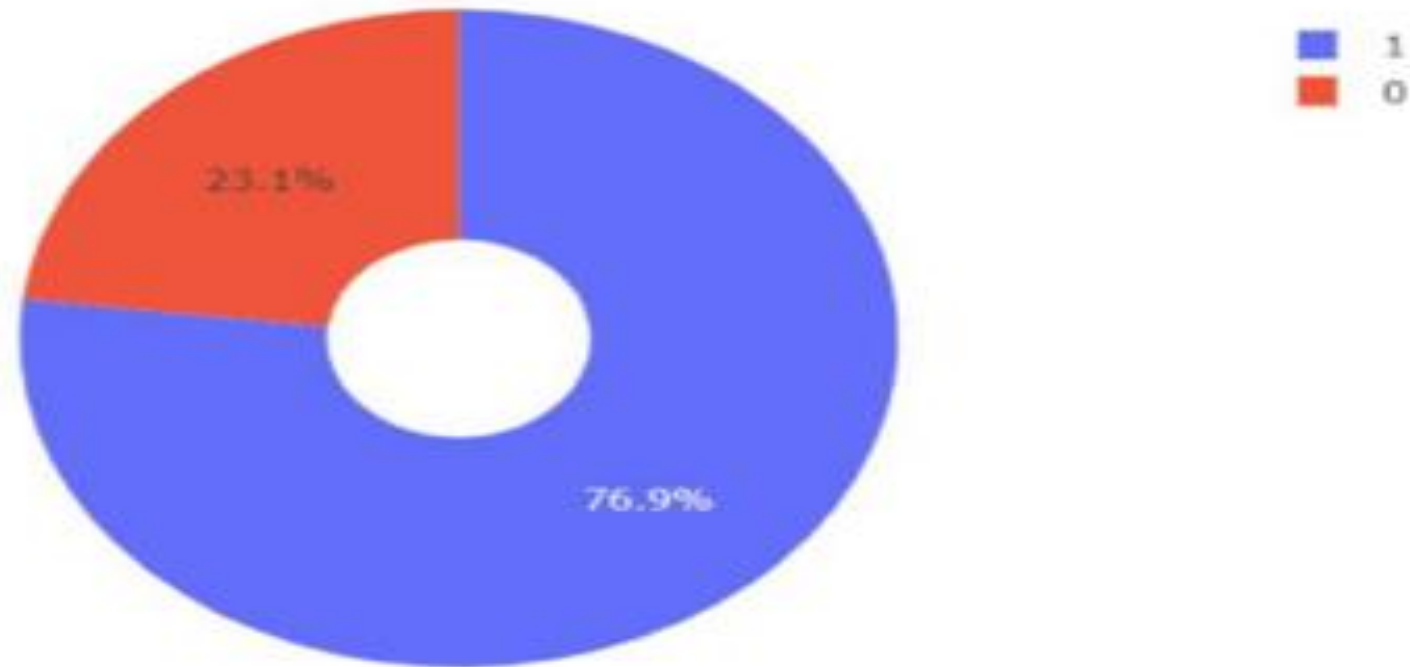
Pie Chart showing the success percentage achieved by each launch site

Total Success Launches By all sites



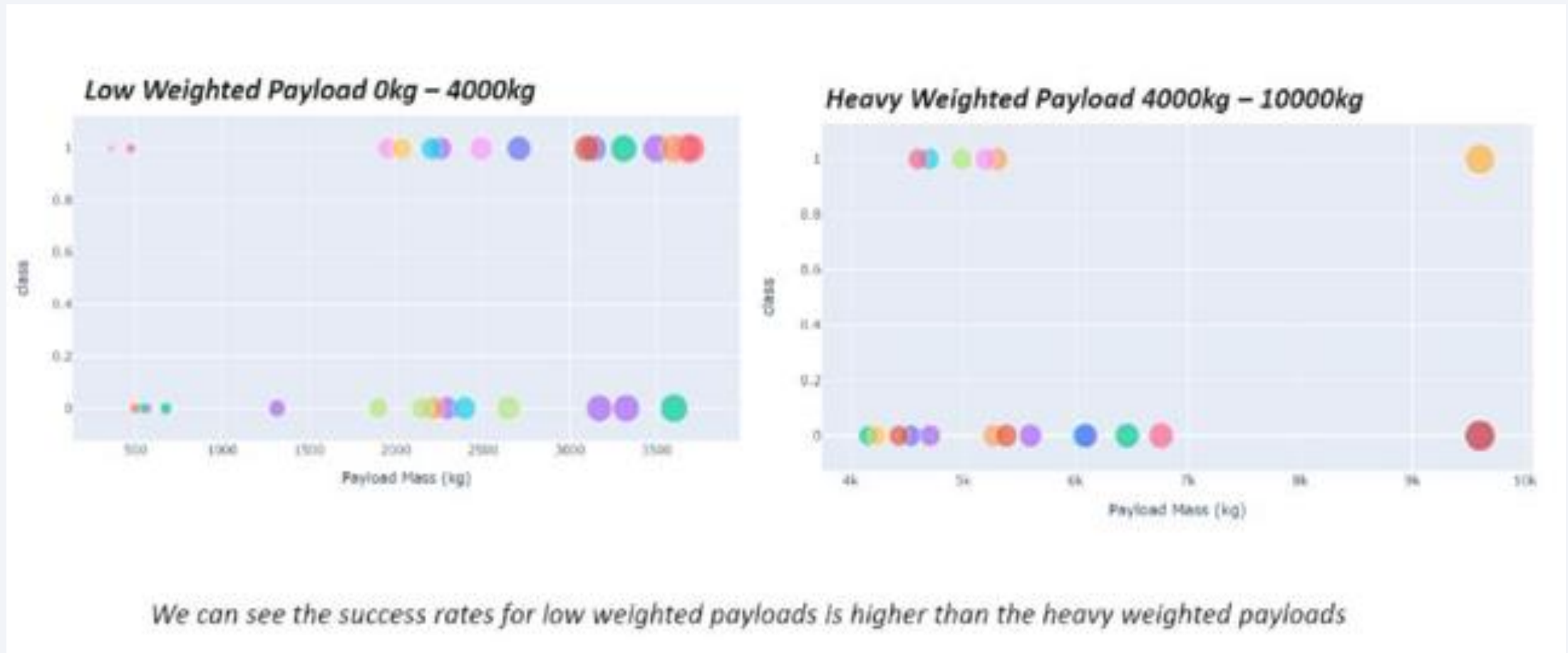
We can see that KSC LC-39A had the most successful launches from all the sites

Pie chart showing the launch sites with 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



Section 5

Predictive Analysis (Classification)

Classification Accuracy

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

```
models = {'KNeighbors': knn_cv.best_score_,
          'DecisionTree': tree_cv.best_score_,
          'LogisticRegression': logreg_cv.best_score_,
          'SupportVector': svm_cv.best_score_}

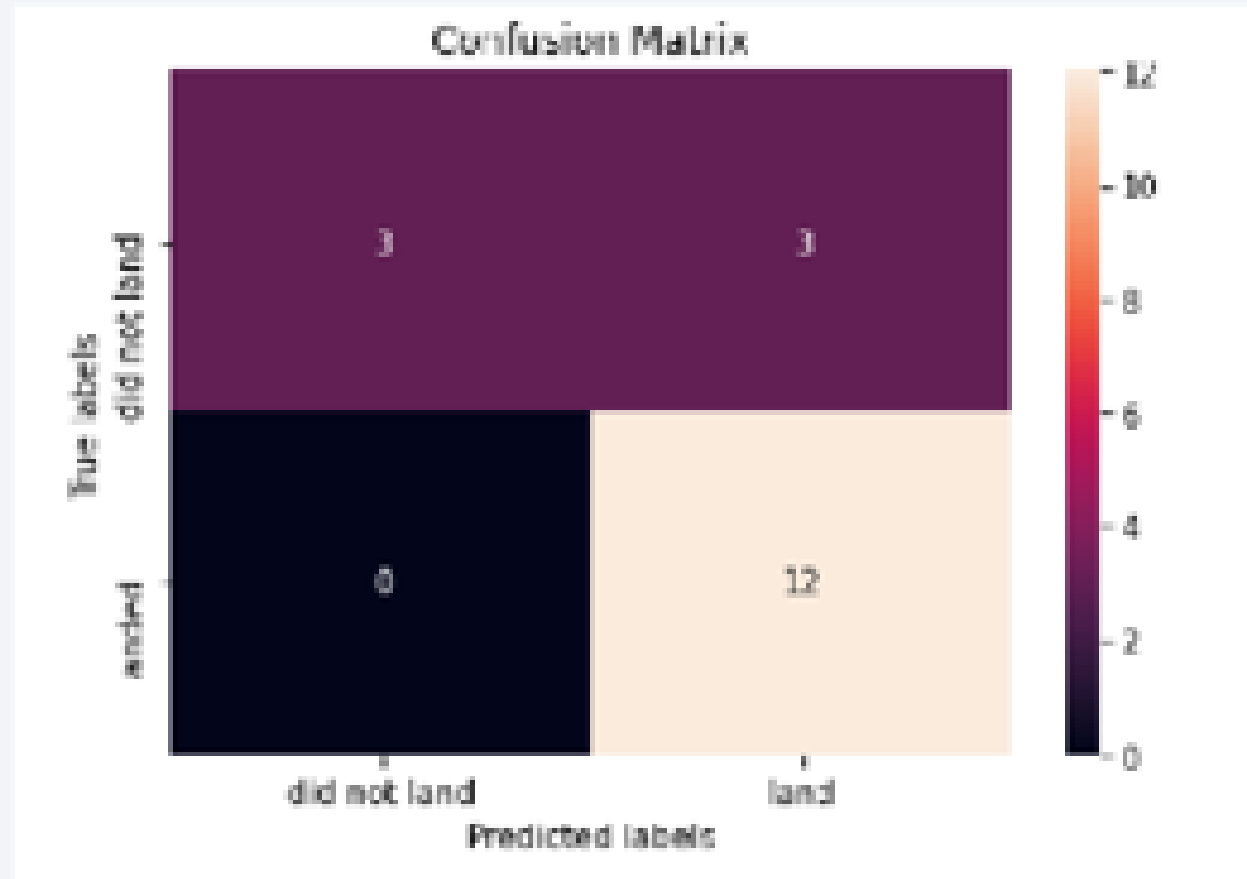
bestalgorithm = max(models, key=models.get)
print('Best model is', bestalgorithm, 'with a score of', models[bestalgorithm])
if bestalgorithm == 'DecisionTree':
    print('Best params is:', tree_cv.best_params_)
if bestalgorithm == 'KNeighbors':
    print('Best params is:', knn_cv.best_params_)
if bestalgorithm == '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': 5, '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

- 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 launch rate.
- The decision tree classifier is the best machine learning algorithm for this task.

Thank you!

