



IBM Developer  
SKILLS NETWORK

# Winning Space Race with Data Science

Guru Raghavendra S  
02 May 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 from Machine Learning Lab

# Introduction

---

SpaceX has disrupted the space industry by offering cost-effective rocket launches, specifically through their Falcon 9 rocket which costs just 62 million dollars, as opposed to the higher price point of 165 million dollars offered by other providers. This is achievable due to SpaceX's innovative concept of reusing the initial launch stage by efficiently landing and renovating it for future missions, resulting in considerable savings. To compete with SpaceX in the bid for rocket launches, a startup intends to establish a machine learning process that can precisely forecast the first stage's landing outcome, which is essential in determining the appropriate bidding price.

The problems included:

- Identifying all factors that influence the landing outcome.
- The relationship between each variables and how it is affecting the outcome.
- The best condition needed to increase the probability of successful landing.

Section 1

# Methodology

# Methodology

---

## Executive Summary

- Data collection methodology:
  - Data was collected using SpaceX REST API and web scrapping from Wikipedia
  - Space X API (<https://api.spacexdata.com/v4/rockets/>)
- Perform data wrangling
- Perform exploratory data analysis (EDA) using visualization and SQL
- Perform interactive visual analytics using Folium and Plotly Dash
- Perform predictive analysis using classification models



# Data Collection

---

The process of data collection involves gathering and measuring specific information within a system to gain insights and evaluate outcomes. For this particular dataset, we utilized two methods - REST API and web scraping - to obtain the data from Wikipedia.

To obtain data through REST API, we initiated a GET request and then converted the response content into JSON format. We used the `json_normalize()` function to transform the data into a pandas dataframe. Afterward, we cleaned the data, identified missing values, and filled in the gaps as required.

For web scraping, we used BeautifulSoup to extract information from an HTML table of launch records, which we then parsed and converted into a pandas dataframe for further analysis.

# Data Collection - SpaceX API

---

```
spacex_url="https://api.spacexdata.com/v4/launches/past"
```

```
response = requests.get(spacex_url)
```

```
# Use json_normalize meethod to convert the json result into a dataframe  
data = pd.json_normalize(response.json())
```

```
# Lets take a subset of our dataframe keeping only the features we want a  
nd the flight number, and date_utc.  
data = data[['rocket', 'payloads', 'launchpad', 'cores', 'flight_number',  
'date_utc']]
```

```
# We will remove rows with multiple cores because those are falcon rocket  
s with 2 extra rocket boosters and rows that have multiple payloads in a  
single rocket.  
data = data[data['cores'].map(len)==1]  
data = data[data['payloads'].map(len)==1]
```

```
# Since payloads and cores are lists of size 1 we will also extract the s  
ingle value in the list and replace the feature.  
data['cores'] = data['cores'].map(lambda x : x[0])  
data['payloads'] = data['payloads'].map(lambda x : x[0])
```

```
# We also want to convert the date_utc to a datetime datatype and then ex  
tracting the date leaving the time  
data['date'] = pd.to_datetime(data['date_utc']).dt.date
```

```
# Using the date we will restrict the dates of the launches  
data = data[data['date'] <= datetime.date(2020, 11, 13)]
```



# Data Collection - Scraping

---

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

```
# Use BeautifulSoup() to create a BeautifulSoup object from a response text content
soup = BeautifulSoup(data, 'html.parser')
```

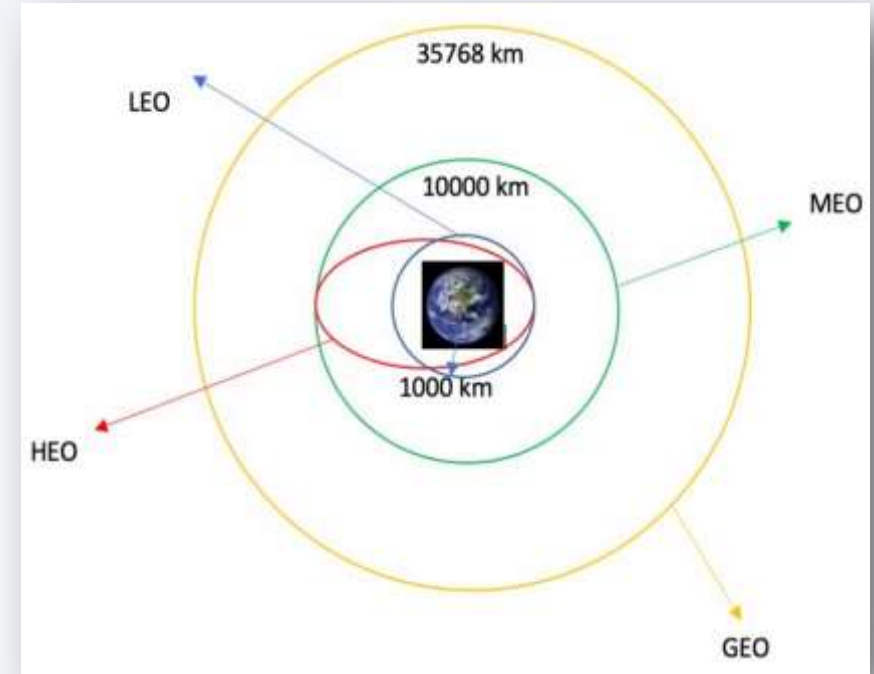
```
extracted_row = 0
#Extract each table
for table_number, table in enumerate(soup.find_all('table', "wikitable plainrowheaders collapsible")):
    # get table row
    for rows in table.find_all("tr"):
        #check to see if first table heading is as number corresponding to launch a number
        if rows.th:
            if rows.th.string:
                flight_number=rows.th.string.strip()
                flag=flight_number.isdigit()
            else:
                flag=False
```

# Data Wrangling

---

Data wrangling involves cleaning and organizing messy and complicated datasets to make them more convenient and suitable for exploratory data analysis.

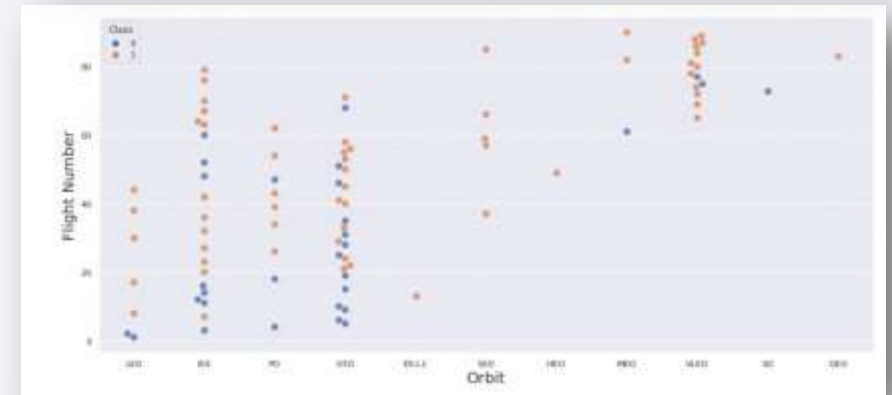
In order to accomplish this, we will begin by determining the quantity of launches from each launch site. After that, we will analyze the quantity and frequency of mission outcomes depending on orbit type. These actions will allow us to have a better grasp of the data and draw important conclusions from it.



# EDA with Data Visualization

We first started by using scatter graph to find the relationship between the attributes such as between:

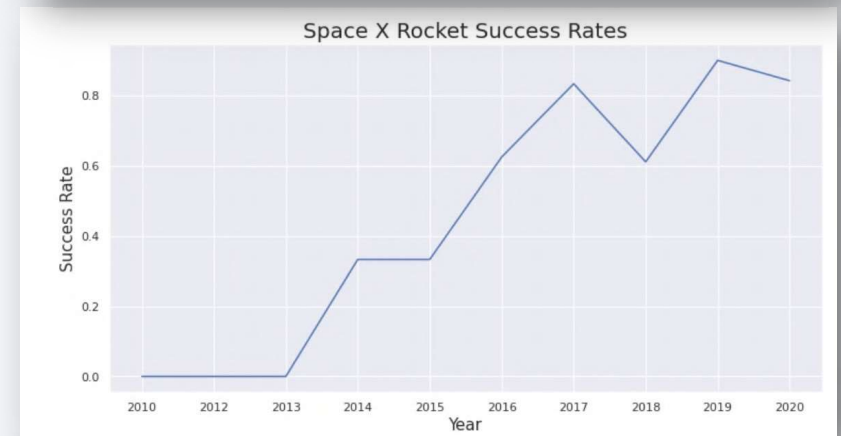
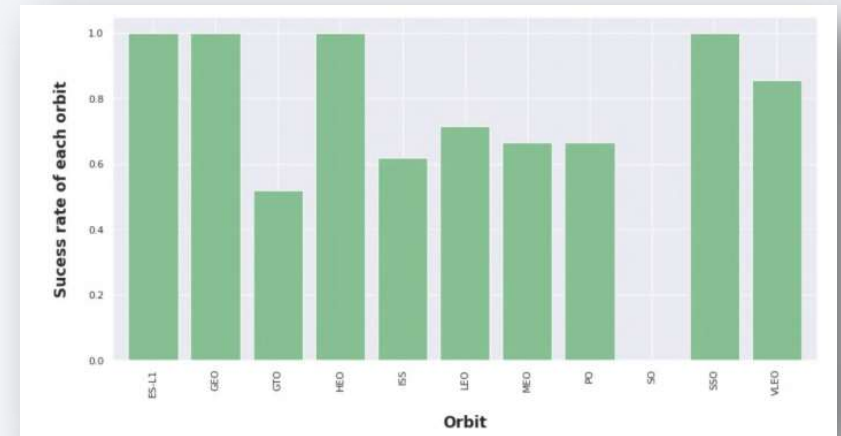
- Payload and Flight Number.
- Flight Number and Launch Site.
- Payload and Launch Site.
- Flight Number and Orbit Type.
- Payload and Orbit Type.



# EDA with Data Visualization

After exploring the relationships between different attributes using scatter plots, we will employ other types of visualizations, such as bar graphs and line plots, to perform additional analysis. Bar graphs provide a straightforward way to understand the connections between various attributes, and we will use them to determine which orbits have the highest probability of success.

Furthermore, we will employ line graphs to track trends or patterns over time, such as the yearly trend in launch success. To enable future success prediction, we will use feature engineering techniques, such as creating dummy variables for categorical columns. This allows us to extract useful information from the dataset and make more precise predictions about future outcomes.



# 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

---

To visualize the launch data into an interactive map. We took the latitude and longitude coordinates at each launch site and added a circle marker around each launch site with a label of the name of the launch site.

We then assigned the dataframe `launch_outcomes(failure,success)` to classes 0 and 1 with **Red** and **Green** markers on the map in `MarkerCluster()`.

We then used the Haversine's formula to calculate the distance of the launch sites to various landmarks to find answers to the questions of:

- How close the launch sites with railways, highways and coastlines?
- How close the launch sites with nearby cities?



# Build a Dashboard with Plotly Dash

---

- We built an interactive dashboard with Plotly dash which allowing the user to play around with the data as they need.
- We plotted pie charts showing the total launches by a certain sites.
- We then plotted scatter graph showing the relationship with Outcome and Payload Mass (Kg) for the different booster version.

# Results

---

The results will be categorized to 3 main results which is:

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

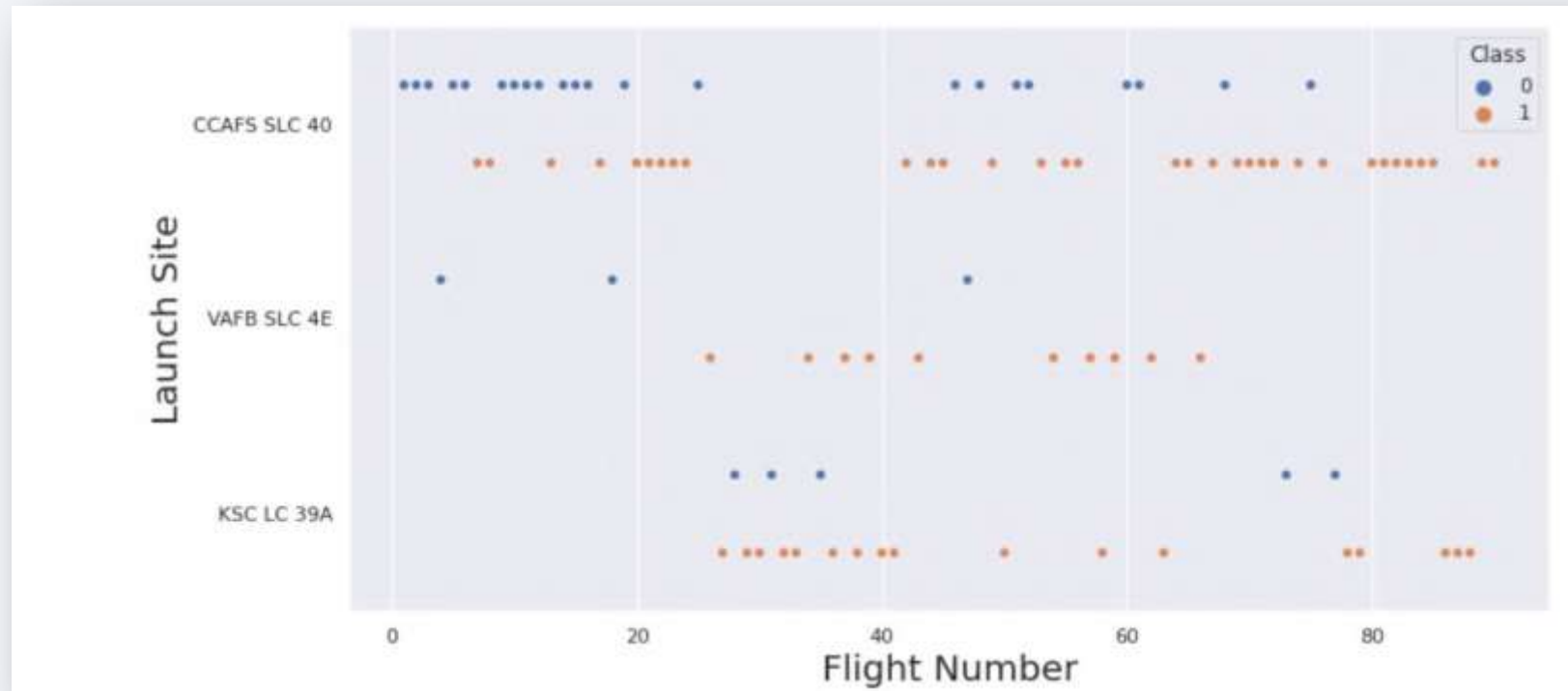
The background of the slide is an abstract composition. It features a solid blue area on the left side, which transitions into a dynamic pattern of diagonal streaks in shades of blue and red on the right. These streaks vary in thickness and intensity, creating a sense of motion and depth. A faint, light blue grid pattern is visible across the entire background, adding a technical or digital feel to the design.

Section 2

# Insights drawn from EDA

# Flight Number vs. Launch Site

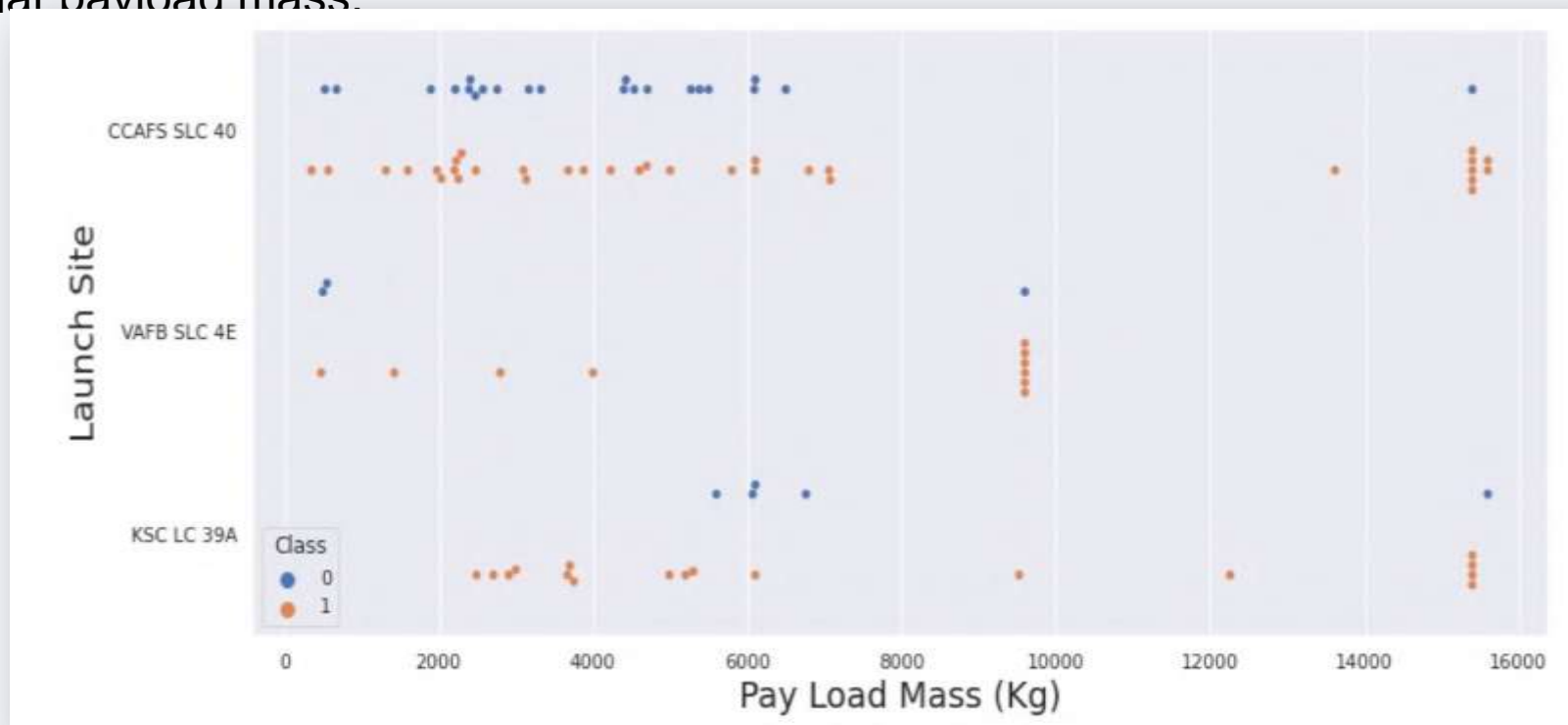
The scatter plot indicates that there is a positive correlation between the size of the launch site and the success rate of the flights. In other words, larger launch sites tend to have higher success rates. However, the launch site CCAFS SLC40 appears to deviate from this pattern and has a lower success rate despite its relatively large size.





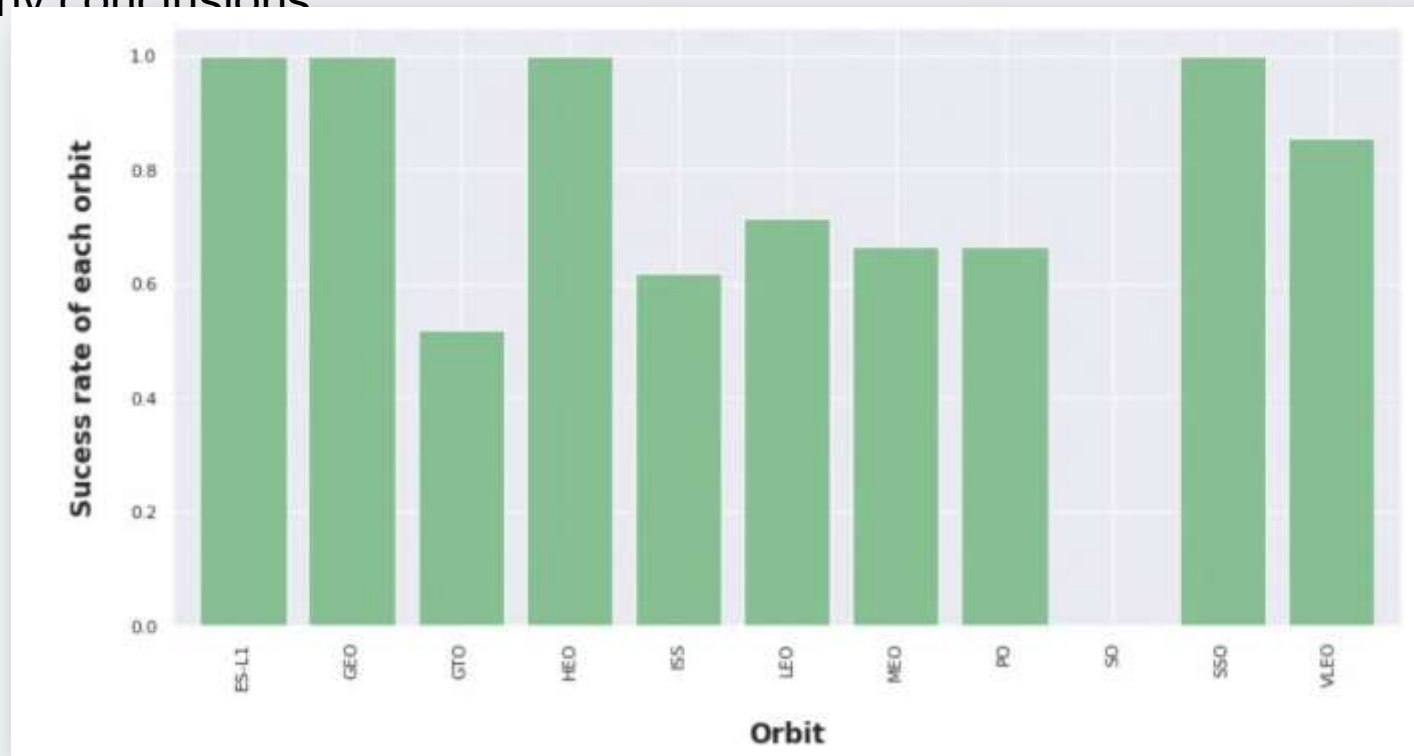
# Payload vs. Launch Site

The scatter plot displays that there is a strong correlation between payload mass and the success rate of a launch for payloads weighing over 7000kg. However, there is no discernible trend to suggest that the success rate is influenced by the launch site for a particular payload mass.



# Success Rate vs. Orbit Type

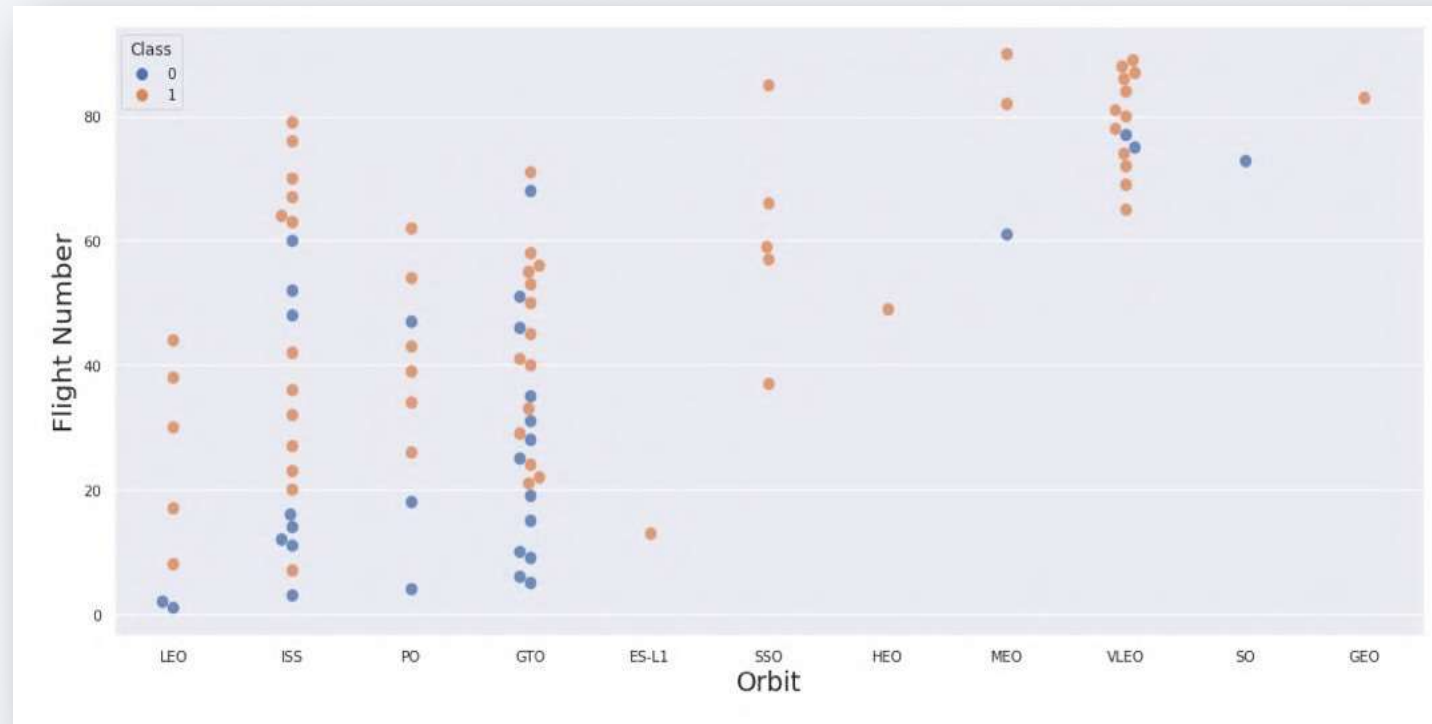
The figure illustrates that the type of orbit may have an impact on the landing outcomes, with certain orbits showing a 100% success rate, such as SSO, HEO, GEO, and ES-L1, while the SO orbit has a 0% success rate. However, further analysis indicates that some of these orbits only have one occurrence in the dataset, such as GEO, SO, HEO, and ES-L1. This means that more data is needed to identify any patterns or trends before drawing any conclusions.





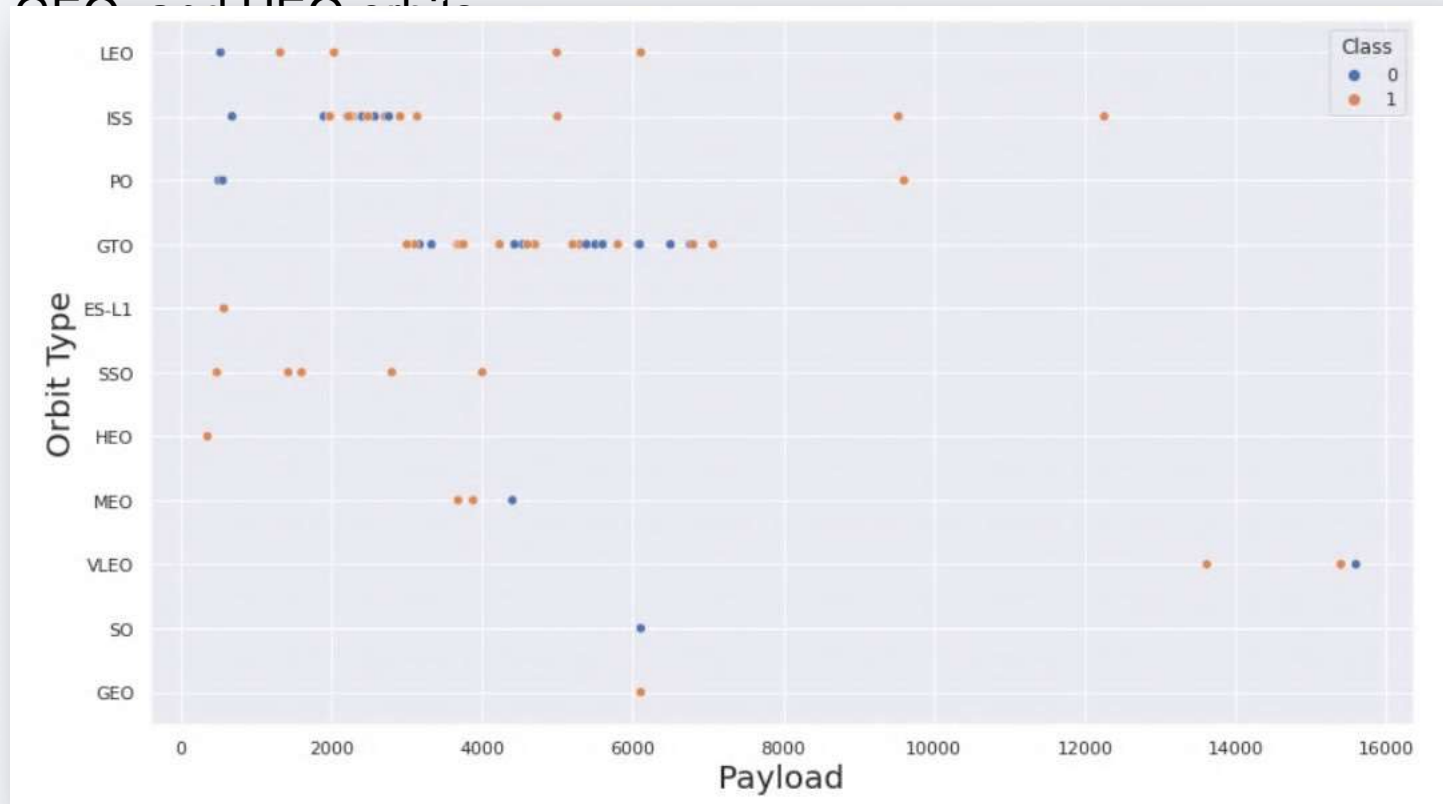
# Flight Number vs. Orbit Type

The scatter plot indicates that, in general, there is a positive correlation between the number of flights on each orbit and the success rate, particularly for the LEO orbit. However, for the GTO orbit, there is no clear relationship between these two attributes. It should be noted that any orbit with only one occurrence in the dataset should be excluded from the above statement since more data is needed to draw any conclusions.



# Payload vs. Orbit Type

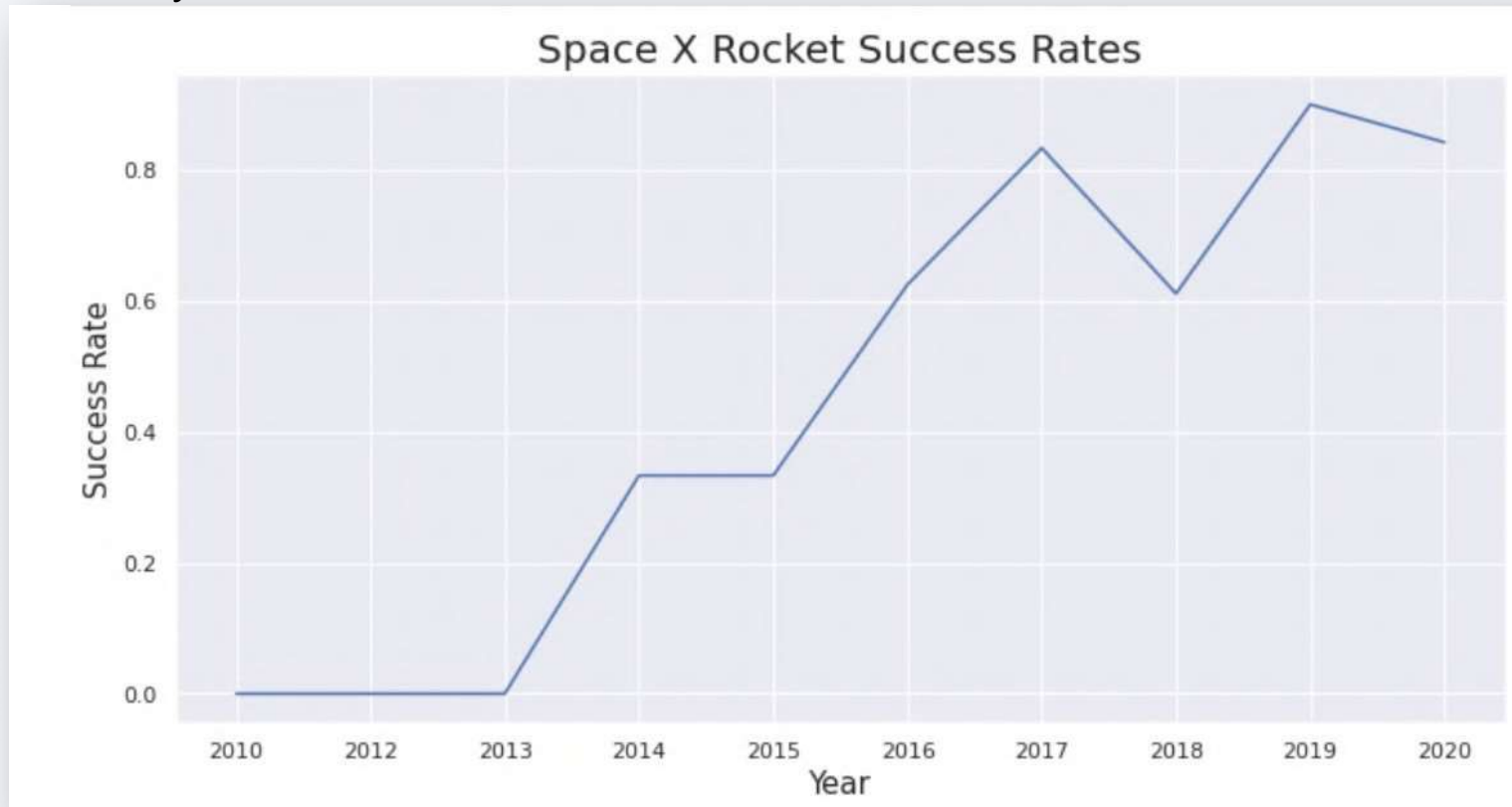
The weight of the payload has a positive impact on the success rate for the LEO, ISS, and PO orbits, while it has a negative impact on the MEO and VLEO orbits. However, there seems to be no clear relationship between the weight of the payload and the success rate for the GTO orbit. Additionally, further data is required to identify any patterns or trends for the SO, GEO, and HEO orbits.



# Launch Success Yearly Trend

---

These figures show a clear upward trend from the year 2013 to 2020, indicating an increase in the success rate over time. If this trend continues into the future, the success rate will steadily rise and eventually reach 100% success rate..



# All Launch Site Names

---

In [5]:

```
%sql SELECT DISTINCT LAUNCH_SITE as "Launch_Sites" FROM SPACEX;
```

```
* ibm_db_sa://zpw86771:***@fbd88901-ebdb-4a4f-a32e-9822b9fb237b.clogj3  
sd0tgtu0lqde00.databases.appdomain.cloud:32731/bludb  
Done.
```

Out[5]:

**Launch\_Sites**

CCAFS LC-40

CCAFS SLC-40

KSC LC-39A

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

# Total Payload Mass

---

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

```
%sql SELECT SUM(PAYLOAD_MASS__KG_) AS "Total Payload Mass by NASA (CRS)"
```

```
* ibm_db_sa://zpw86771:***@fbd88901-ebdb-4a4f-a32e-9822b9fb237b.clogj3  
sd0tgtu0lqde00.databases.appdomain.cloud:32731/bludb  
Done.
```

**Total Payload Mass by NASA (CRS)**

---

45596



# Average Payload Mass by F9 v1.1

---

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

```
%sql SELECT AVG(PAYLOAD_MASS__KG_) AS "Average Payload Mass by Booster  
WHERE BOOSTER_VERSION = 'F9 v1.1';
```

```
* ibm_db_sa://zpw86771:***@fbd88901-ebdb-4a4f-a32e-9822b9fb237b.clogj3  
sd0tgtu0lqde00.databases.appdomain.cloud:32731/bludb  
Done.
```

**Average Payload Mass by Booster Version F9 v1.1**

---

2928

# First Successful Ground Landing Date

---

We use the min() function to find the result

We observed that the dates of the first successful landing outcome on ground pad was 22<sup>nd</sup> December 2015

```
%sql SELECT MIN(DATE) AS "First Successful Landing Outcome in Ground Pad"  
WHERE LANDING__OUTCOME = 'Success (ground pad)';
```

```
* ibm_db_sa://zpw86771:***@fbd88901-ebdb-4a4f-a32e-9822b9fb237b.clogj3  
sd0tgtu0lqde00.databases.appdomain.cloud:32731/bludb  
Done.
```

**First Successful Landing Outcome in Ground Pad**

---

2015-12-22

## Successful Drone Ship Landing with Payload between 4000 and 6000

To narrow down our search, we utilized a WHERE clause to filter for boosters that had successfully landed on a drone ship. Additionally, we employed the AND condition to ensure that our results only included cases where the landing was successful and the payload mass was greater than 4000 but less than 6000.

```
%sql SELECT BOOSTER_VERSION FROM SPACEX WHERE LANDING__OUTCOME = 'Success (drone ship)' \
AND PAYLOAD_MASS__KG_ > 4000 AND PAYLOAD_MASS__KG_ < 6000;
```

```
* ibm_db_sa://zpw86771:***@fbd88901-ebdb-4a4f-a32e-9822b9fb237b.clogj3sd0tgtu0lqde00.databases.appdomain.cloud:32731/bludb
```

Done.

**booster\_version**

F9 FT B1022

F9 FT B1026

F9 FT B1021.2

F9 FT B1031.2

# Total Number of Successful and Failure Mission Outcomes

---

List the total number of successful and failure mission outcomes

```
%sql SELECT COUNT(MISSION_OUTCOME) AS "Successful Mission" FROM SPACEX WHERE MISSION_OUTCOME LIKE 'Success%';
```

```
* ibm_db_sa://zpw86771:***@fbd88901-ebdb-4a4f-a32e-9822b9fb237b.clogj3sd0tgtu0lqde00.databases.appdomain.cloud:32731/bludb  
Done.
```

Successful Mission
--------------------

100
-----

```
%sql SELECT COUNT(MISSION_OUTCOME) AS "Failure Mission" FROM SPACEX WHERE MISSION_OUTCOME LIKE 'Failure%';
```

```
* ibm_db_sa://zpw86771:***@fbd88901-ebdb-4a4f-a32e-9822b9fb237b.clogj3sd0tgtu0lqde00.databases.appdomain.cloud:32731/bludb  
Done.
```

Failure Mission
-----------------

1
---

# Boosters Carried Maximum Payload

```
mysql SELECT DISTINCT BOOSTER_VERSION AS "Booster Versions which carried the Maximum Payload Mass" FROM SPACEX  
WHERE PAYLOAD_MASS_KG_=(SELECT MAX(PAYLOAD_MASS_KG_) FROM SPACEX);
```

```
* ibm_db_sa://zpw86771:***@fbd88901-ebdb-4a4f-a32e-9822b9fb237b.clogj3sd0tgtu0lqde00.databases.appdomain.clou  
d:32731/bludb
```

Done.

## Booster Versions which carried the Maximum Payload Mass

F9 B5 B1048.4

F9 B5 B1048.5

F9 B5 B1049.4

F9 B5 B1049.5

F9 B5 B1049.7

F9 B5 B1051.3

F9 B5 B1051.4

F9 B5 B1051.6

F9 B5 B1056.4

F9 B5 B1058.3

F9 B5 B1060.2

F9 B5 B1060.3

# 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

```
%sql SELECT BOOSTER_VERSION, LAUNCH_SITE FROM SPACEX WHERE DATE LIKE '2015-%' AND \
LANDING__OUTCOME = 'Failure (drone ship)';
```

```
* ibm_db_sa://zpw86771:***@fbd88901-ebdb-4a4f-a32e-9822b9fb237b.clogj3sd0tgtu0lqde00.
databases.appdomain.cloud:32731/bludb
Done.
```

booster_version	launch_site
F9 v1.1 B1012	CCAFS LC-40
F9 v1.1 B1015	CCAFS LC-40



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

To extract specific information from the data, we selected the landing outcomes and the count of landing outcomes. Using the WHERE clause, we filtered the landing outcomes to include only those that occurred between June 4th, 2010 and March 20th, 2010. The GROUP BY clause was then used to group the landing outcomes, and the ORDER BY clause was applied to sort the grouped outcomes in descending order.

```
sql SELECT LANDING__OUTCOME as "Landing Outcome", COUNT(LANDING__OUTCOME) AS "Total Count" FROM SPACEX \
WHERE DATE BETWEEN '2010-06-04' AND '2017-03-20' \
GROUP BY LANDING__OUTCOME \
ORDER BY COUNT(LANDING__OUTCOME) DESC ;
```

```
* ibm_db_sa://zpw86771:***@fbd88901-ebdb-4a4f-a32e-9822b9fb237b.clogj3sd0tgtu0lqde00.databases.appdomain.c
loud:32731/bludb
Done.
```

Landing Outcome	Total Count
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

A satellite view of Earth from space, showing the curvature of the planet and city lights at night. The image is a composite of a solid blue background on the left and a satellite photograph of Earth on the right. The Earth's surface is dark, with numerous bright yellow and orange lights representing cities and urban areas. The horizon of the Earth is visible as a thin, curved line separating the dark surface from the deep blue of space.

Section 3

# Launch Sites Proximities Analysis

# Location of all the Launch Sites

---



# Markers showing launch sites with color labels





# Launch Sites Distance to Landmarks



- 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



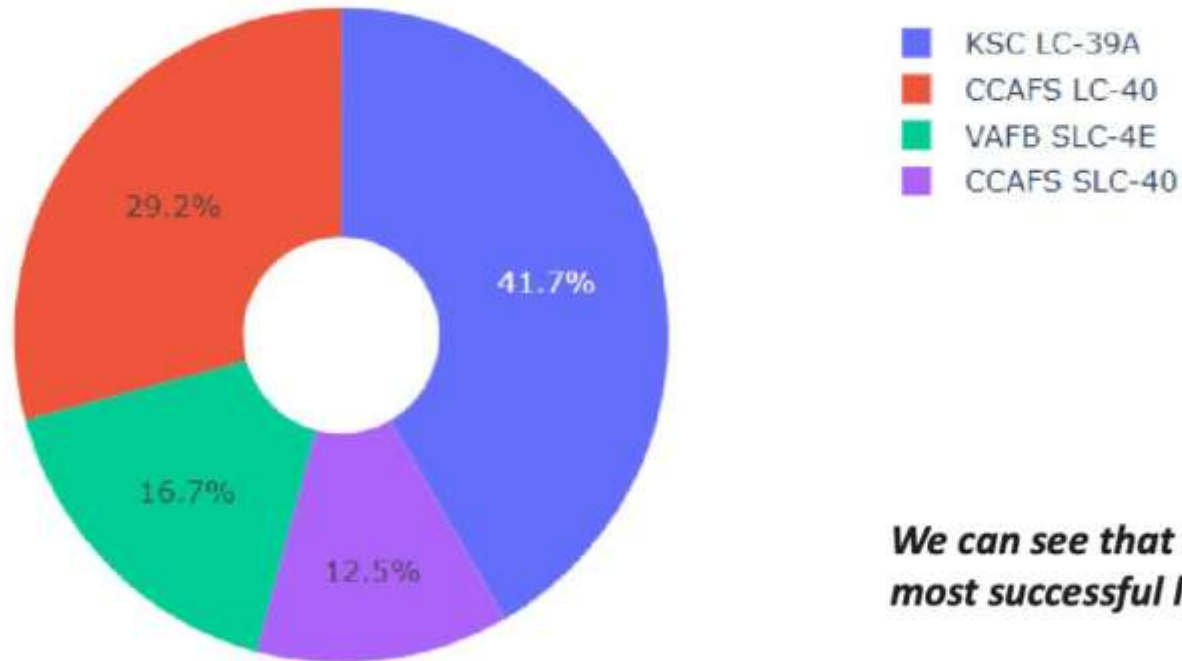
Section 4

# Build a Dashboard with Plotly Dash



# The success percentage by each sites.

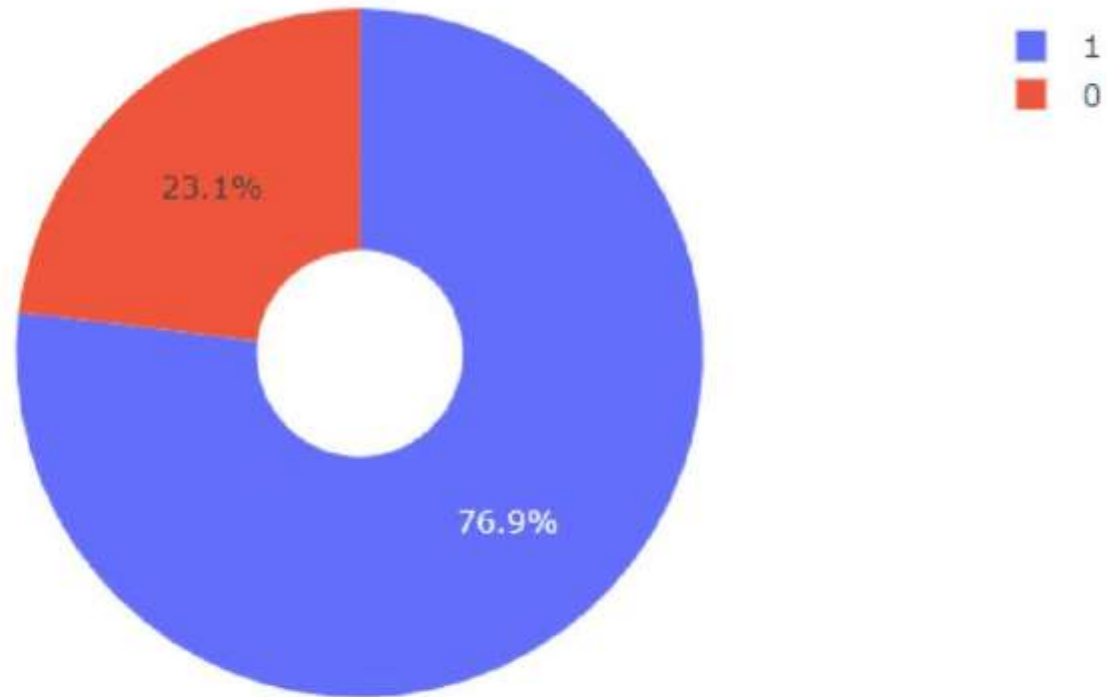
---



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

# The highest launch-success ratio: KSC LC-39A

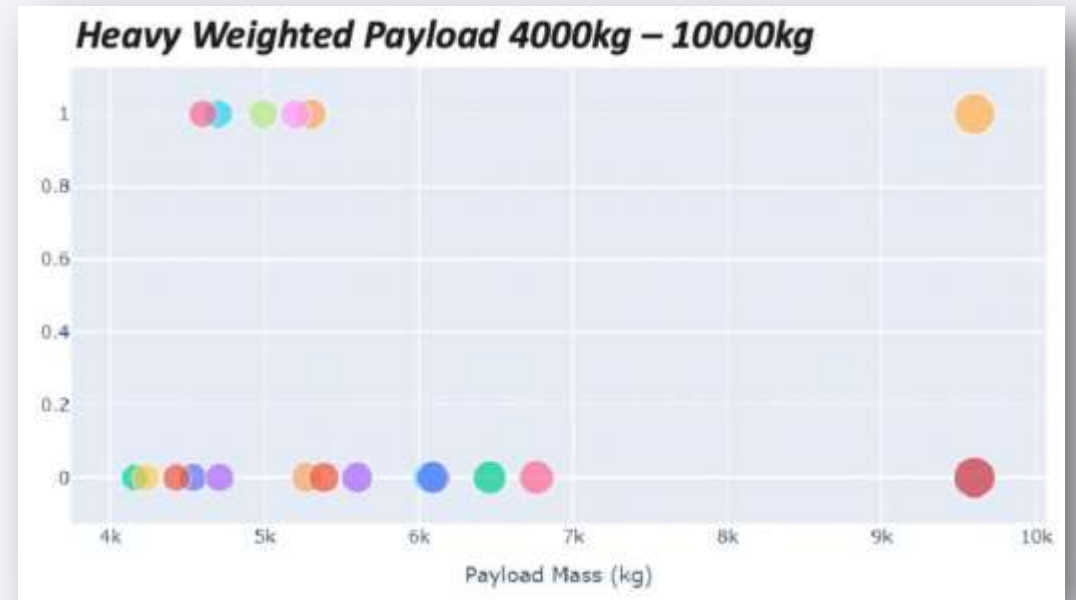
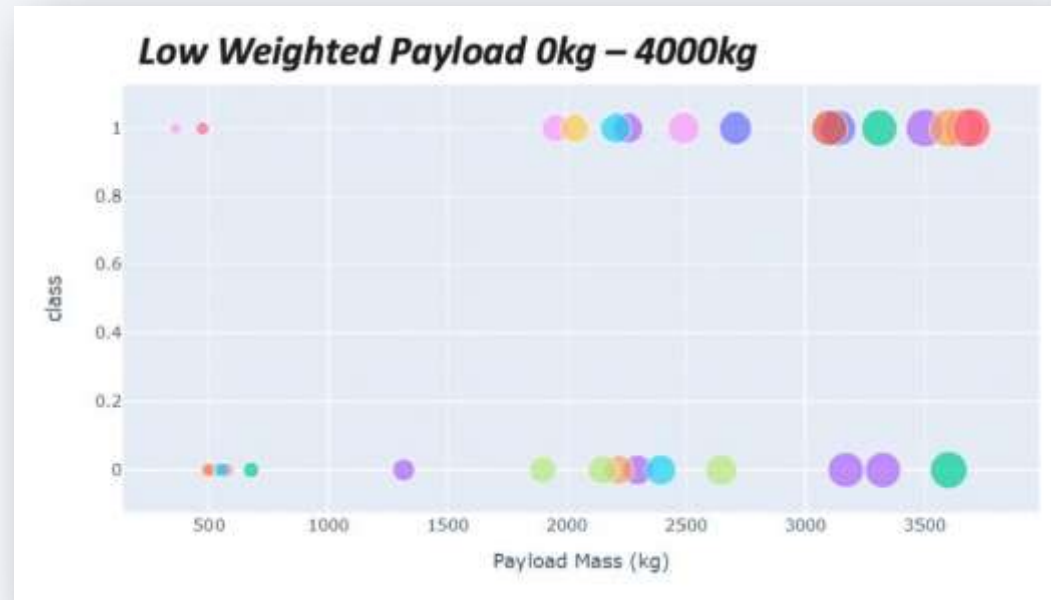
---



***KSC LC-39A achieved a 76.9% success rate while getting a 23.1% failure rate***

# Payload vs Launch Outcome Scatter Plot

---





Section 5

# Predictive Analysis (Classification)

# Classification Accuracy

---

Tree Algorithm with the highest classification accuracy.

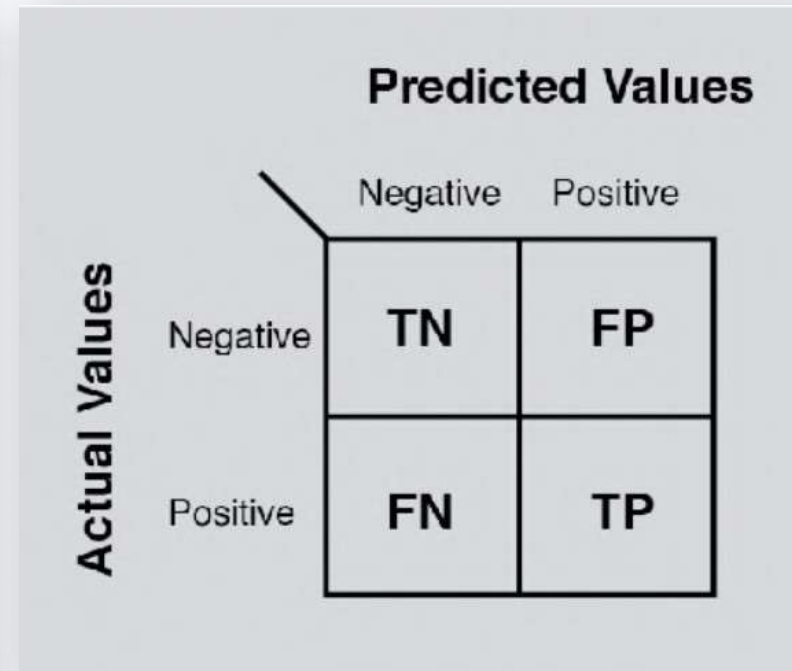
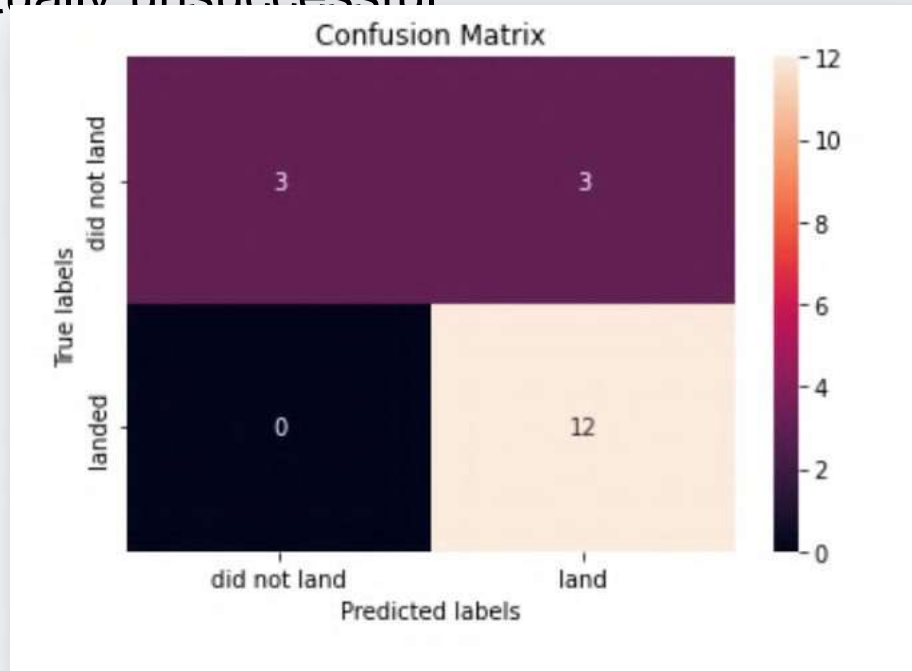
```
algorithms = {'KNN':knn_cv.best_score_, 'Tree':tree_cv.best_score_, 'LogisticRegression':logreg_cv.best_score_}
bestalgorithm = max(algorithms, key=algorithms.get)
print('Best Algorithm is',bestalgorithm,'with a score of',algorithms[bestalgorithm])
if bestalgorithm == 'Tree':
    print('Best Params is :',tree_cv.best_params_)
if bestalgorithm == 'KNN':
    print('Best Params is :',knn_cv.best_params_)
if bestalgorithm == 'LogisticRegression':
    print('Best Params is :',logreg_cv.best_params_)
```

Best Algorithm is Tree with a score of 0.9017857142857142

Best Params is : {'criterion': 'entropy', 'max\_depth': 10, 'max\_features': 'auto', 'min\_samples\_leaf': 2, 'min\_samples\_split': 10, 'splitter': 'random'}

# Confusion Matrix

The confusion matrix generated for the decision tree classifier indicates that the classifier is capable of differentiating between the various classes. However, the primary issue appears to be false positives, which are cases where the classifier incorrectly identifies a landing as successful when it was actually unsuccessful.





# Conclusions

---

Based on the analysis of the dataset, it was observed that the Tree Classifier Algorithm was the most effective Machine Learning approach. Interestingly, launches with a payload weight of 4000kg or less exhibited better performance compared to those with heavier payloads. The success rate of SpaceX launches has been steadily increasing since 2013, indicating the possibility of further improvement in the future. Among all launch sites, KSC LC-39A had the highest success rate of 76.9%. Similarly, the SSO orbit demonstrated the highest success rate, with all its launches being successful on multiple occasions.

Thank you!

