

# Winning Space Race with Data Science

An Bui

11 May 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 from Machine Learning Lab

# Introduction

---

- In this capstone, we will predict if the Falcon 9 first stage will land successfully
- SpaceX 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 SpaceX can reuse the first stage.
- 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 SpaceX for a rocket launch.
- The problems
  - 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:
  - SpaceX launch data is gathered from SpaceX REST API
  - Perform data wrangling
  - Data was processed using one-hot encoding for 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

---

- SpaceX launch data is gathered from SpaceX REST API
  - API will give us data about launches, including information about the rocket used, payload delivered, launch specifications, landing specifications, and landing outcome.
- The dataset was collected by REST API and Web Scrapping from Wikipedia
  - For REST API, its started by using the get request.
  - We decoded the response content as Json and turn it into a pandas dataframe using `json_normalize()`.
  - We then cleaned the data, checked for missing values and fill with whatever needed.
- For web scrapping, we use the BeautifulSoup to extract the launch records as HTML table, parse the table and convert it to a pandas dataframe for further analysis

# Data Collection – SpaceX API

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

<https://github.com/abuialberta/spacex/blob/main/DataCollectionAPI.ipynb>

Now let's start requesting rocket launch data from SpaceX API with the following URL:

```
In [6]: spacex_url="https://api.spacexdata.com/v4/launches/past"
```

```
In [7]: response = requests.get(spacex_url)
```

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

```
In [15]: # Use json_normalize meethod to convert the json result into a dataframe
import pandas as pd
from pandas import json_normalize

data = pd.json_normalize(response.json())
```

We will now use the API again to get information about the launches using the IDs given for each launch. Specifically we will be using columns `rocket`, `payloads`, `launchpad`, and `cores`.

```
In [17]: # Lets take a subset of our dataframe keeping only the features we want and the flight number
data = data[['rocket', 'payloads', 'launchpad', 'cores', 'flight_number', 'date_utc']]

# We will remove rows with multiple cores because those are falcon rockets with 2 extra rocket cores
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 single value in the list
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 extracting the date and 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

Web scrap Falcon 9 launch records with *BeautifulSoup*



Extract a Falcon 9 launch records HTML table from Wikipedia



Parse the table and convert it into a Pandas data frame

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

In [9]:

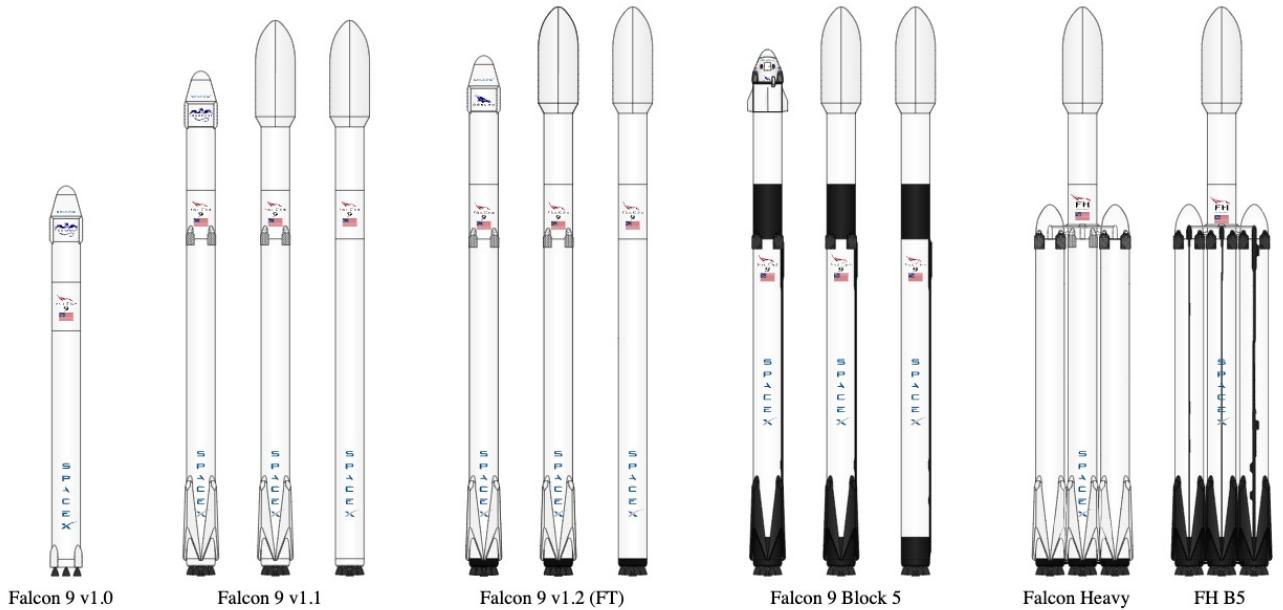
```
# Use BeautifulSoup() to create a BeautifulSoup object from a response text content  
soup = BeautifulSoup(data)
```

In [18]:

```
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  
        #get table element  
        row=rows.find_all('td')  
        #if it is number save cells in a dictionary  
        if flag:  
            extracted_row += 1  
            # Flight Number value  
            # TODO: Append the flight_number into launch_dict with key 'Flight No.'  
            #print(flight_number)  
            datatimelist=date_time(row[0])
```

# 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 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.

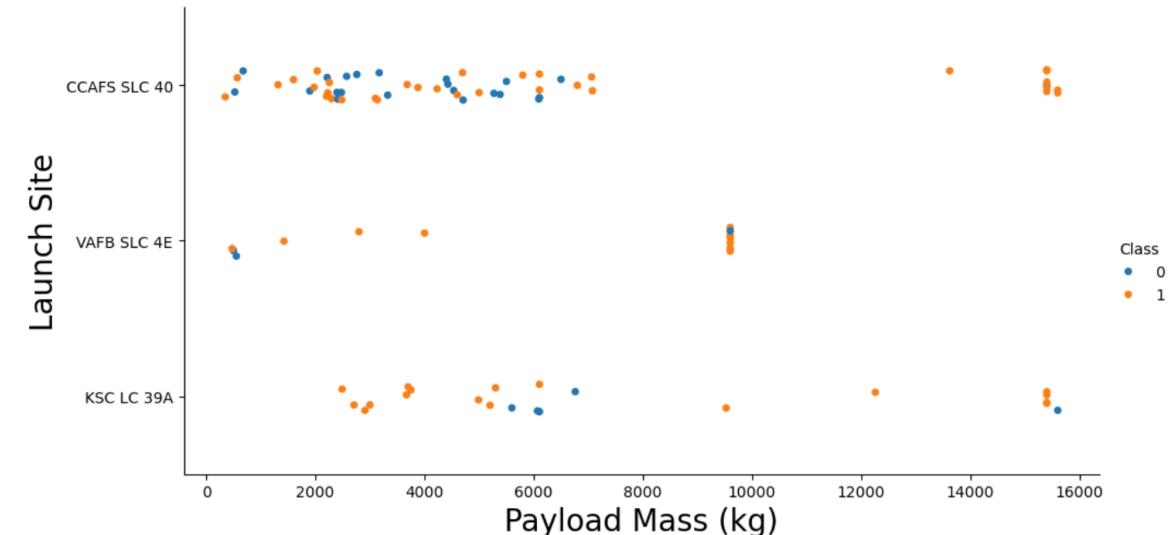
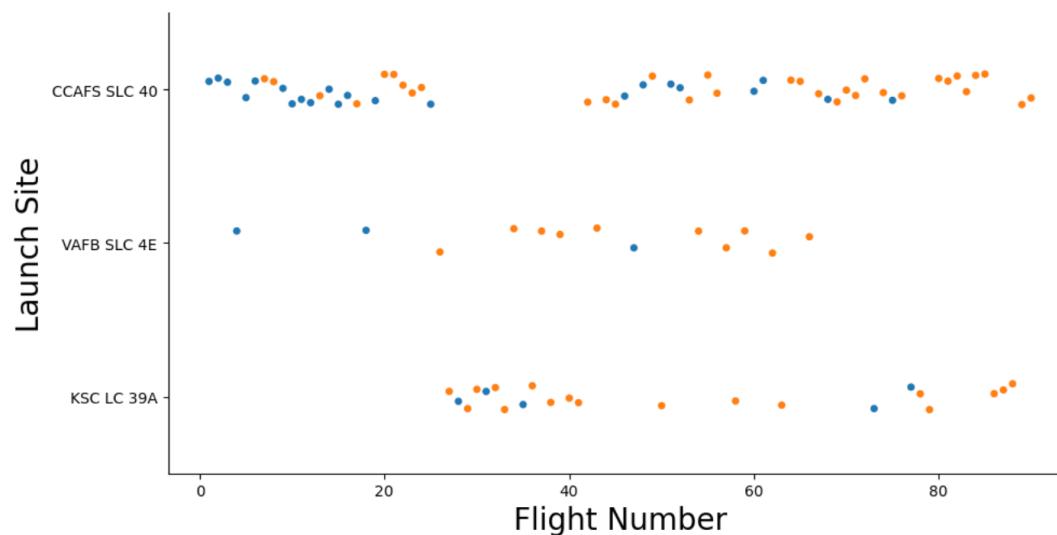


<https://github.com/abui-ualberta/spacex/blob/main/%20Data%20Wrangling.ipynb>

# EDA with Data Visualization

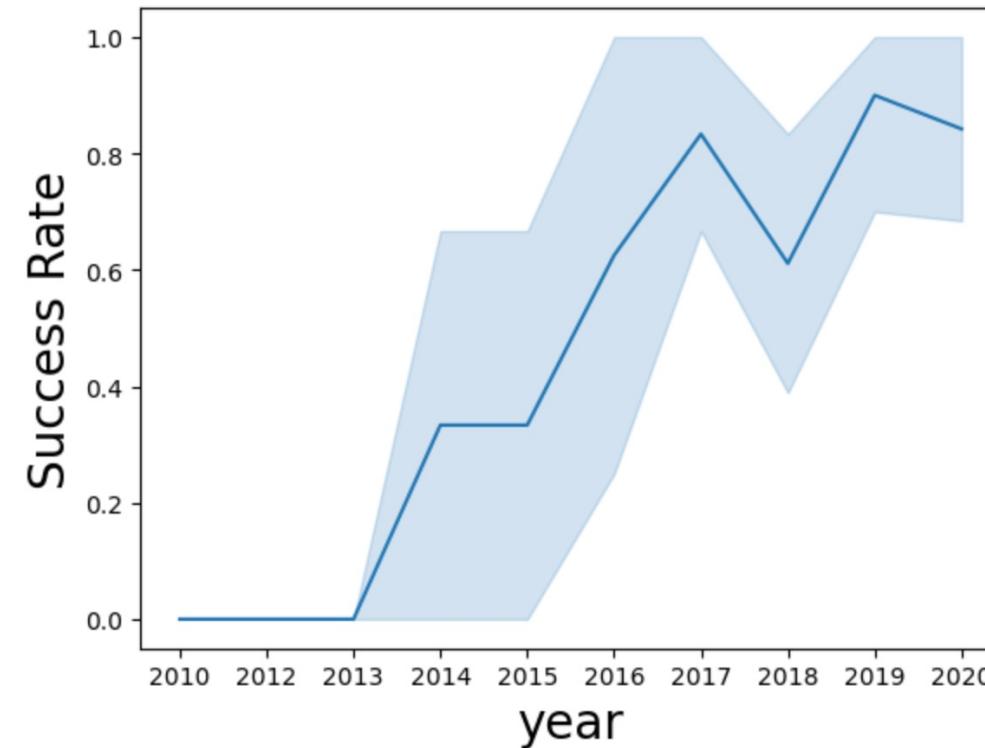
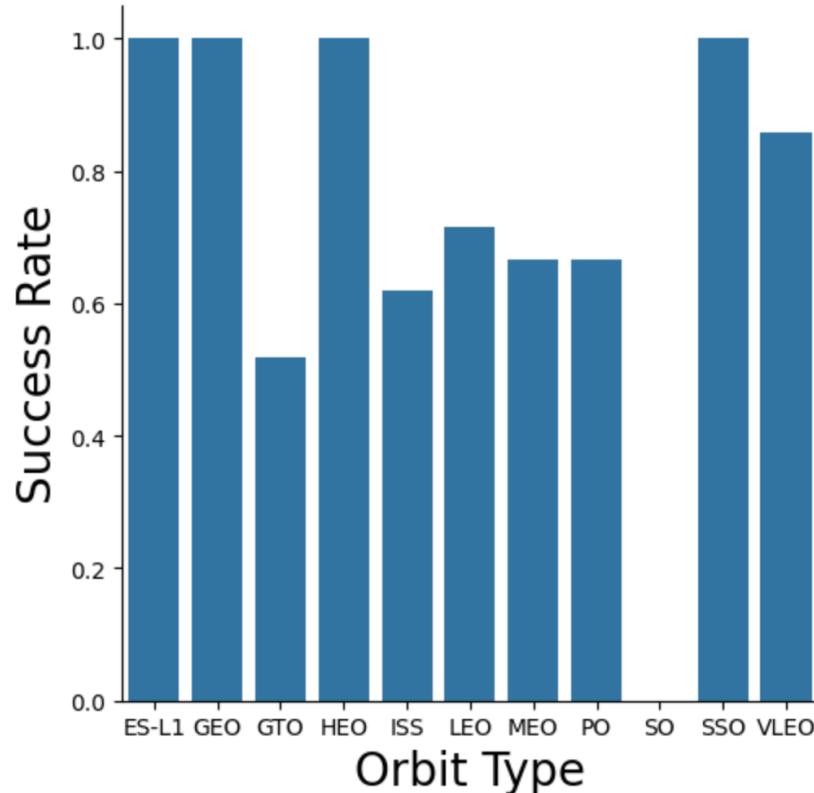
Scatter graph to find the relationship between the attributes

- Flight Number and Payload Mass
- Flight Number and Launch Site
- Launch Site and Payload Mass



# EDA with Data Visualization

- Bar graph to determine which orbits have the highest probability of success
- Line graph to show a trends of success rate over time



<https://github.com/abui-ulberta/spacex/blob/main/EDA%20Visualization.ipynb>

# EDA with SQL

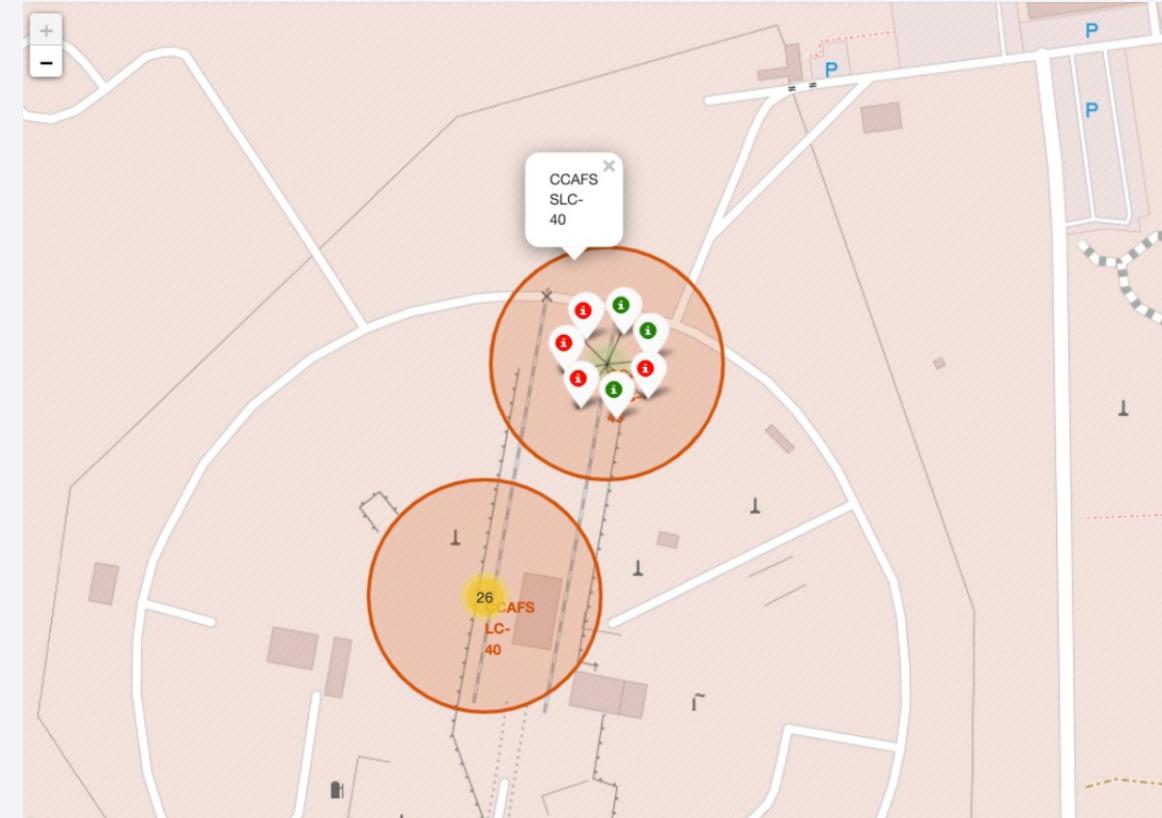
---

- Display names of the unique launch sites in the space mission
- Display 5 records where launch sites begin with the string 'KSC'
- Display the total payload mass carried by boosters launched by NASA (CRS)
- Display average payload mass carried by booster version F9 v1.1
- List the date where the successful landing outcome in drone ship was achieved.
- List the names of the boosters which have success in ground pad and have payload mass greater than 4000 but less than 6000
- List the total number of successful and failure mission outcomes
- List the names of the booster\_versions which have carried the maximum payload mass. Use a subquery
- List the records which will display the month names, successful landing\_outcomes in ground pad ,booster versions, launch\_site for the months in year 2017
- Rank 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

# Build an Interactive Map with Folium

---

- We mark latitude and longitude coordinates at each launch site with a label
- If a launch was successful (class=1), then we use a green marker and if a launch was failed, we use a red marker (class=0)
- We then used the Haversine's formula to calculated the distance of the launch sites to various landmark to answer
  - How close the launch sites with railways, highways and coastlines?
  - How close the launch sites with nearby cities?



# Predictive Analysis (Classification)

---

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.

Evaluating the model

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

Improving the model

- Use Feature Engineering and Algorithm Tuning

Find the best model

# Results

---

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

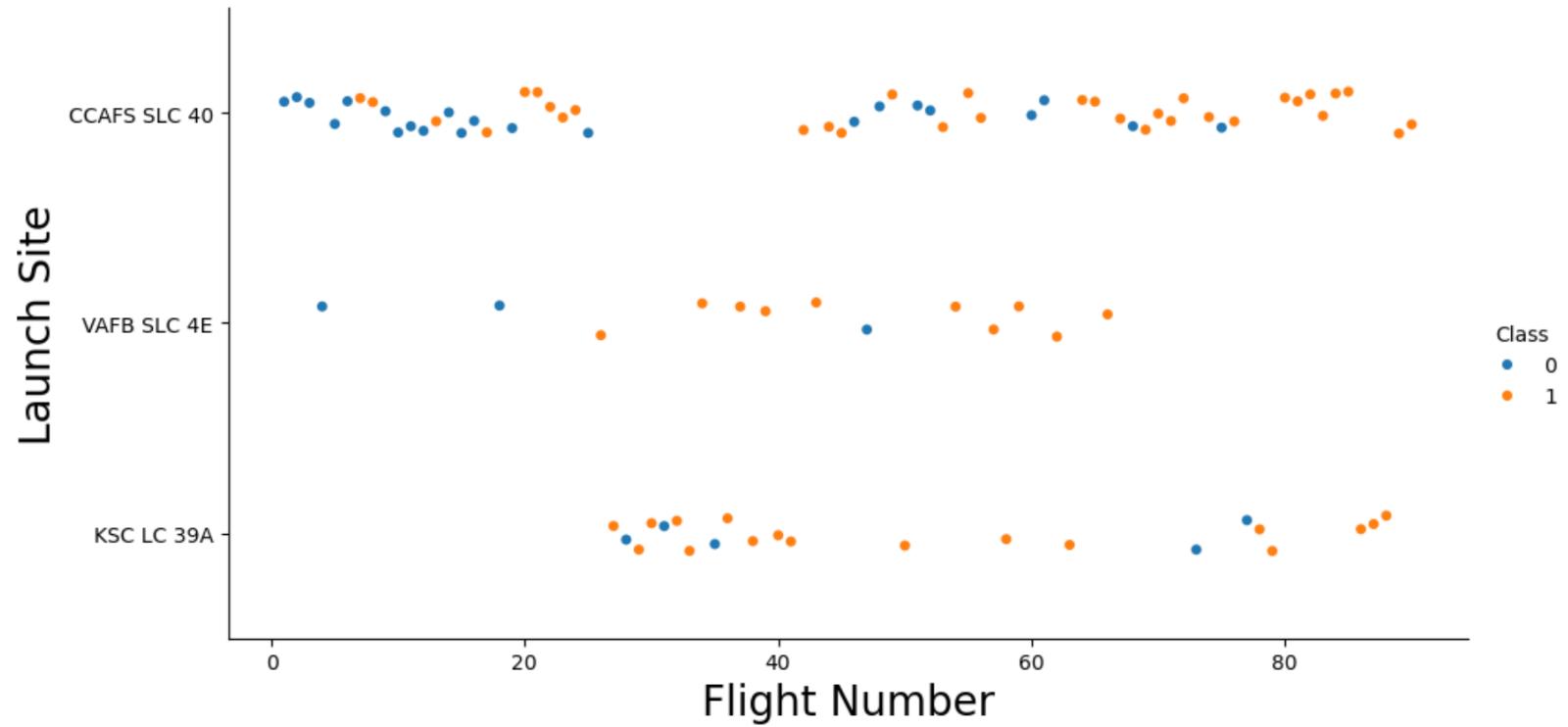
The background of the slide features a complex, abstract digital visualization. It consists of numerous thin, glowing lines that create a sense of depth and motion. The lines are primarily blue and red, with some green and purple highlights. They form a grid-like structure that curves and twists across the frame, resembling a three-dimensional space or a network of data points. The overall effect is futuristic and dynamic.

Section 2

## Insights drawn from EDA

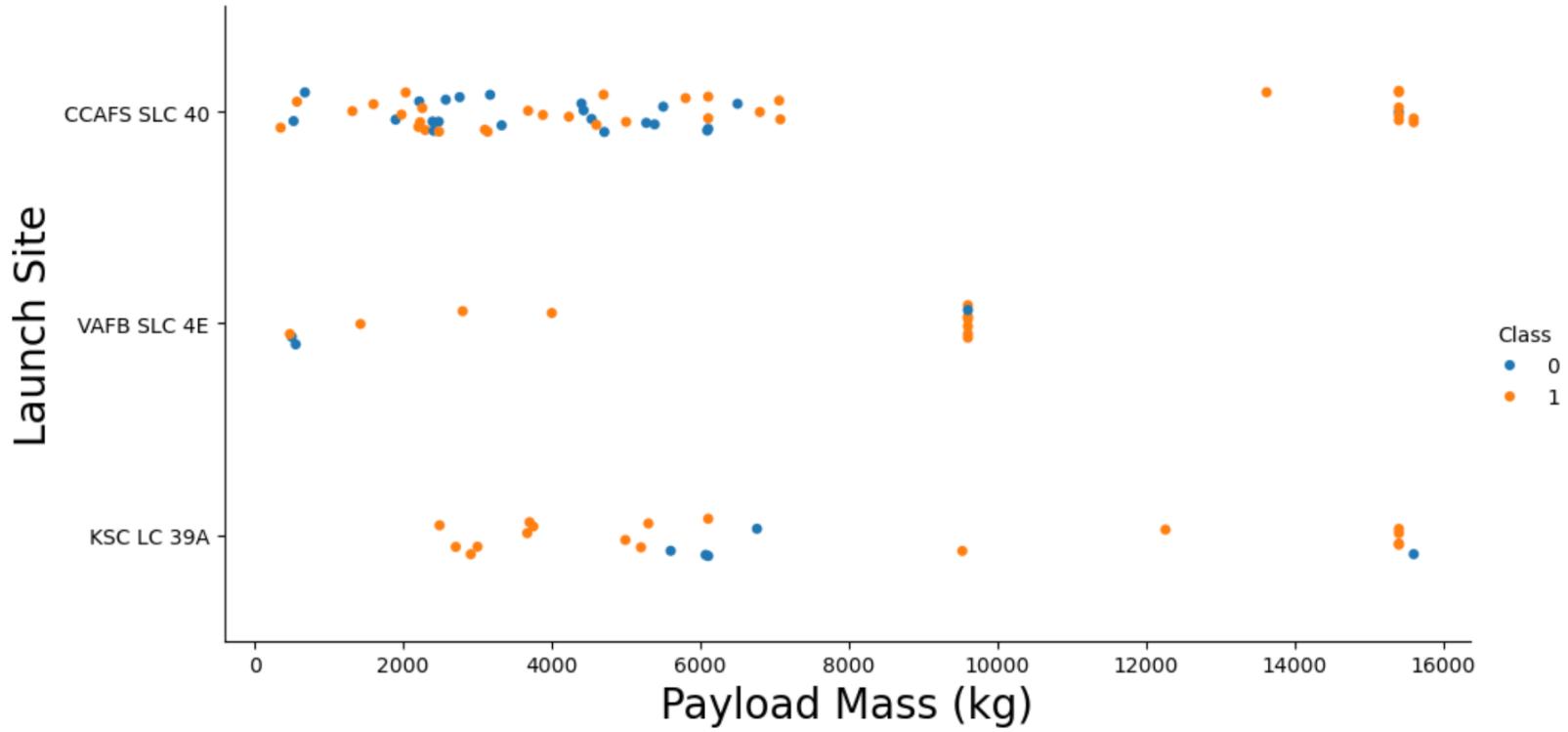
# Flight Number vs. Launch Site

- This scatter plot shows that the larger the flights amount of the launch site, the greater the success rate will be
- However, site CCAFS SLC40 shows the least pattern of this



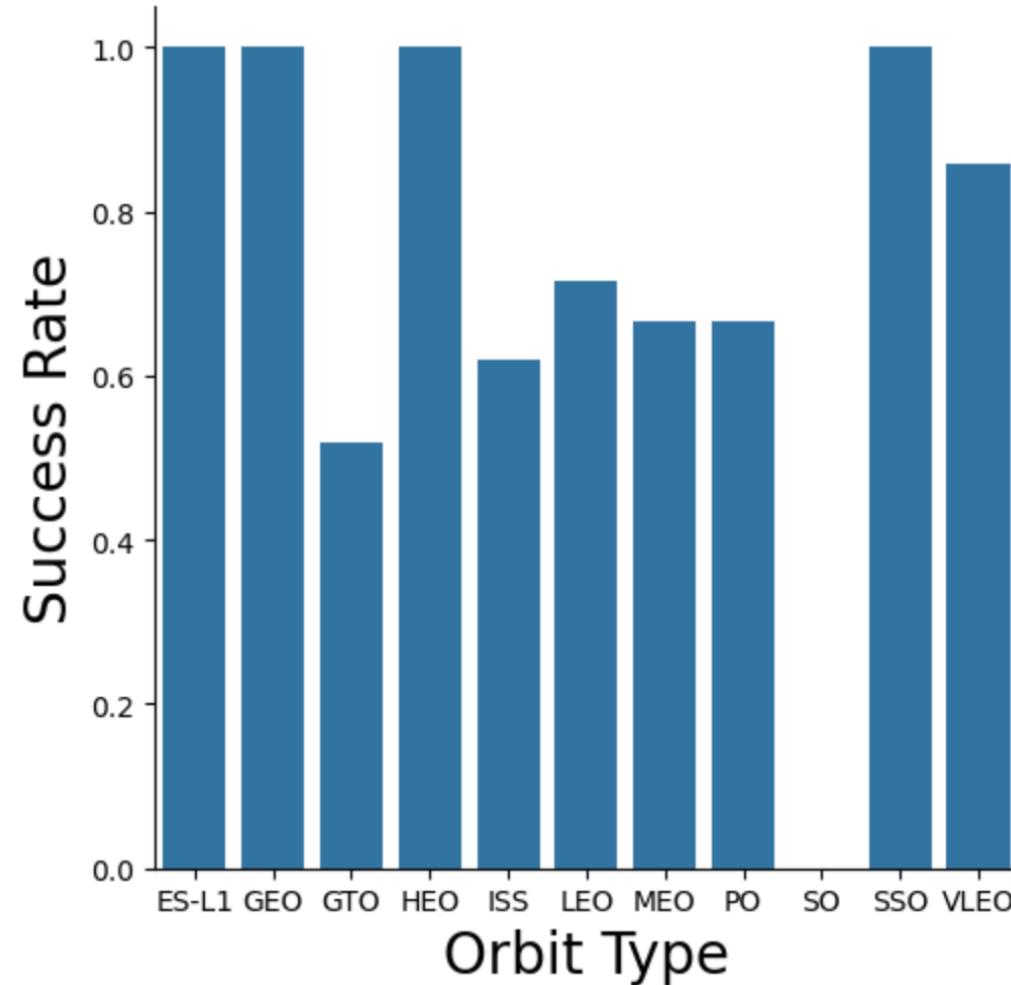
# Payload vs. Launch Site

- This scatter plot shows once the payload mass is greater than 7000kg, the probability of the success rate will be highly increased.
- However, there is no clear pattern to say the launch site is dependent to the payload mass for the success rate.



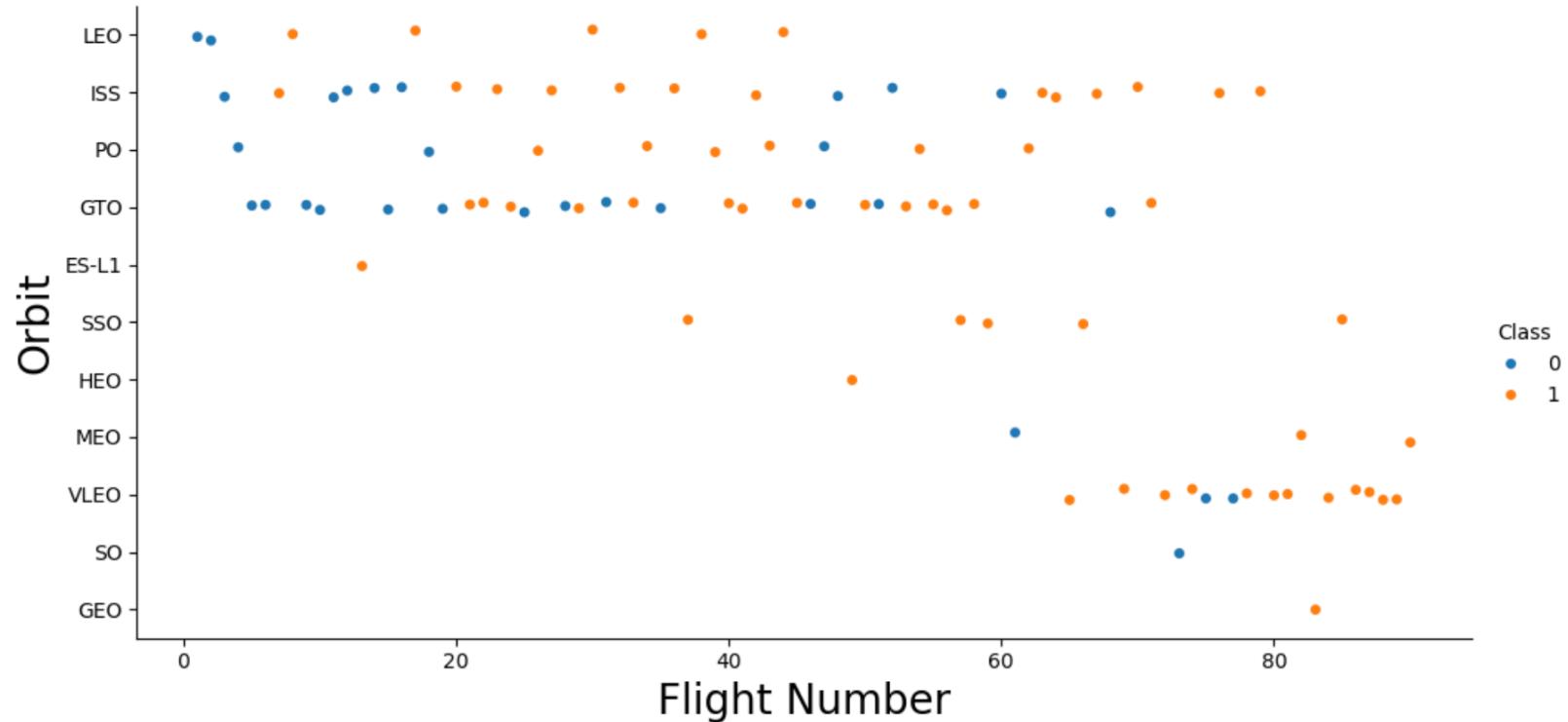
# Success Rate vs. Orbit Type

- This figure depicted the possibility of the orbits to influences the landing outcomes as some orbits has 100% success rate such as SSO, HEO, GEO AND ES-L1 while SO orbit produced 0% rate of success.
- However, deeper analysis show that some of this orbits has only 1 occurrence such as GEO, SO, HEO and ES-L1 which mean this data need more dataset to see pattern or trend before we draw any conclusion.



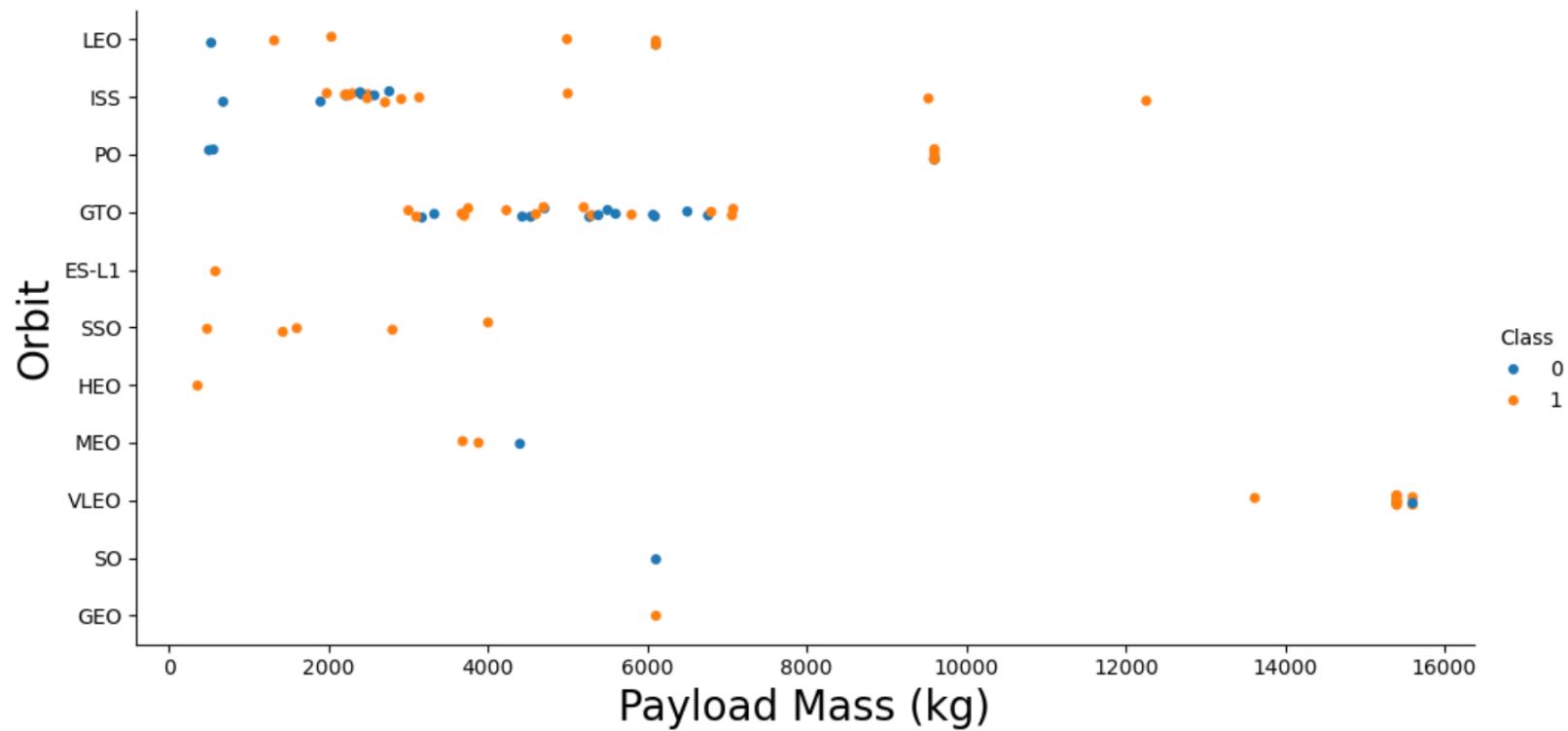
# Flight Number vs. Orbit Type

- This scatter plot shows that generally, the larger the flight number on each orbits, the greater the success rate (especially LEO orbit) except for GTO orbit which depicts no relationship between both attributes.
- Orbit that only has 1 occurrence should also be excluded from above statement as it's needed more dataset.



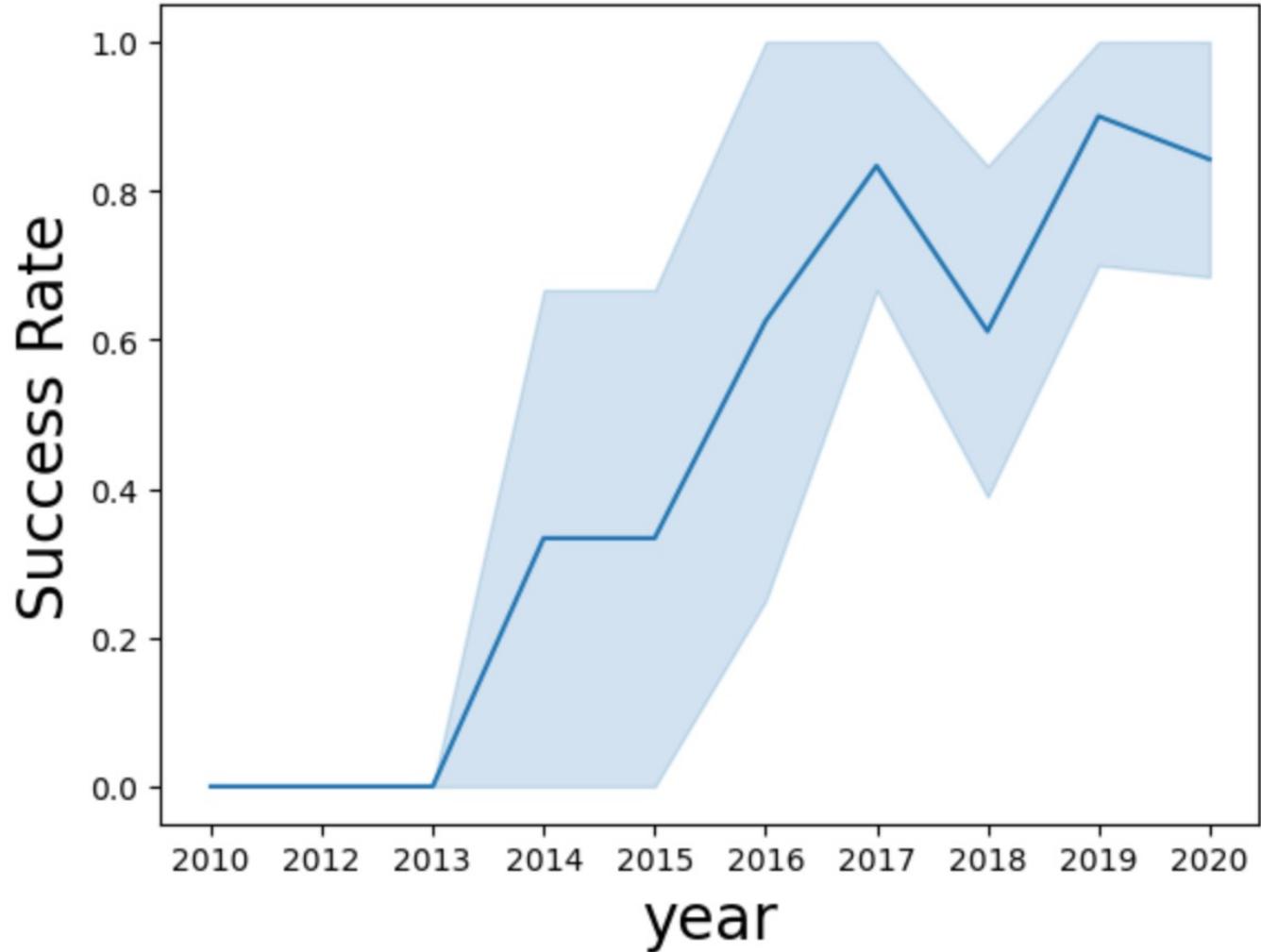
# Payload vs. Orbit Type

- Heavier payload has positive impact on LEO, ISS and PO orbit. However, it has negative impact on MEO and VLEO orbit.
- GTO orbit seem to depict no relation between the attributes.
- Meanwhile, again, SO, GEO and HEO orbit need more dataset to see any pattern or trend.



# Launch Success Yearly Trend

- This figure clearly depicted an increasing trend from the year 2013 until 2020.
- If this trend continues for the next year onward. The success rate will steadily increase until reaching 100% success rate.



# 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

In [9]:

```
%sql SELECT Distinct LAUNCH_SITE FROM SPACEXTBL
```

\* sqlite:///my\_data1.db

Done.

Out[9]:

Launch\_Site

CCAFS LC-40

VAFB SLC-4E

KSC LC-39A

CCAFS SLC-40

# Launch Site Names Begin with 'KSC'

Display 5 records where launch sites begin with the string 'KSC'

In [10]:

```
%sql SELECT * FROM SPACEXTBL WHERE LAUNCH_SITE LIKE 'CCA%' LIMIT 5
```

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

Out[10]:

| Date       | Time_(UTC) | Booster_Version | Launch_Site | Payload   | PAYLOAD_MASS__KG_ | Orbit     | Customer        | Mission_Outcome | Landing_Site |
|------------|------------|-----------------|-------------|---|-------------------|-----------|-----------------|-----------------|--------------|
| 2010-06-04 | 18:45:00   | F9 v1.0 B0003   | CCAFS LC-40 | Dragon Spacecraft Qualification Unit                          | 0                 | LEO       | SpaceX          | Success         | Failure      |
| 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         | Failure      |
| 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

- 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 [11]: %sql SELECT SUM(PAYLOAD_MASS__KG_) FROM SPACEXTBL WHERE CUSTOMER='NASA (CRS)'
```

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

```
Out[11]: SUM(PAYLOAD_MASS__KG_)
```

---

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

In [12]:

```
%sql SELECT AVG(PAYLOAD_MASS__KG_) FROM SPACEXTBL WHERE BOOSTER_VERSION='F9 v1.1'
```

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

Out[12]: AVG(PAYLOAD\_MASS\_\_KG\_)

2928.4

# 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 22nd December 2015

List the date where the succesful landing outcome in drone ship was acheived.

*Hint:Use min function*

In [14]:

```
%sql SELECT MIN(DATE) FROM SPACEXTBL WHERE LANDING_OUTCOME = 'Success (ground pad)'
```

\* sqlite:///my\_data1.db

Done.

Out[14]: **MIN(DATE)**

---

2015-12-22

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

- 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

List the names of the boosters which have success in ground pad and have payload mass greater than 4000 but less than 6000

```
In [16]: %sql SELECT BOOSTER_VERSION FROM SPACEXTBL WHERE PAYLOAD_MASS__KG_ BETWEEN 4000 AND 6000 AND LANDING_OUTCOME = 'S  
* sqlite:///my_data1.db  
Done.
```

```
Out[16]: Booster_Version  
F9 FT B1022  
F9 FT B1026  
F9 FT B1021.2  
F9 FT B1031.2
```

# Total Number of Successful and Failure Mission Outcomes

- We used wildcard like ‘%’ to filter for WHERE MissionOutcome

List the total number of successful and failure mission outcomes

```
In [17]: %sql SELECT COUNT(*) FROM SPACEXTBL WHERE MISSION_OUTCOME LIKE '%Success%' OR MISSION_OUTCOME LIKE '%Failure%'  
* sqlite:///my_data1.db  
Done.
```

```
Out[17]: COUNT(*)
```

101

# Boosters Carried Maximum Payload

List the names of the booster\_versions which have carried the maximum payload mass. Use a subquery

In [18]:

```
%sql SELECT BOOSTER_VERSION FROM SPACEXTBL WHERE PAYLOAD_MASS__KG_ = (SELECT MAX(PAYLOAD_MASS__KG_) FROM SPACEXTB
```

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

```
Done.
```

Out[18]: **Booster\_Version**

F9 B5 B1048.4

F9 B5 B1049.4

F9 B5 B1051.3

F9 B5 B1056.4

F9 B5 B1048.5

F9 B5 B1051.4

F9 B5 B1049.5

F9 B5 B1060.2

F9 B5 B1058.3

F9 B5 B1051.6

F9 B5 B1060.3

F9 B5 B1049.7

# 2017 Launch Records

In [28]:

```
%sql SELECT substr(Date,6,2) as month, DATE, BOOSTER_VERSION, LAUNCH_SITE, Landing_Outcome \
FROM SPACEXTBL \
WHERE Landing_Outcome = 'Success (ground pad)' \
AND substr(Date,1,4) = '2017';
```

\* sqlite:///my\_data1.db

Done.

Out[28]:

| month | Date       | Booster_Version | Launch_Site  | Landing_Outcome      |
|-------|------------|-----------------|--------------|----------------------|
| 02    | 2017-02-19 | F9 FT B1031.1   | KSC LC-39A   | Success (ground pad) |
| 05    | 2017-05-01 | F9 FT B1032.1   | KSC LC-39A   | Success (ground pad) |
| 06    | 2017-06-03 | F9 FT B1035.1   | KSC LC-39A   | Success (ground pad) |
| 08    | 2017-08-14 | F9 B4 B1039.1   | KSC LC-39A   | Success (ground pad) |
| 09    | 2017-09-07 | F9 B4 B1040.1   | KSC LC-39A   | Success (ground pad) |
| 12    | 2017-12-15 | F9 FT B1035.2   | CCAFS SLC-40 | Success (ground pad) |

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 2017

# 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)) between the date 2010-06-04 and 2017-03-20, in descending order

```
[31]: %sql SELECT Landing_Outcome, count(*) as count_outcomes \
FROM SPACEXTBL \
WHERE DATE between '2010-06-04' and '2017-03-20' \
GROUP BY Landing_Outcome \
ORDER BY count_outcomes DESC;
```

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

```
[31]: 

| Landing_Outcome        | count_outcomes |
|------------------------|----------------|
| No attempt             | 10             |
| Success (drone ship)   | 5              |
| Failure (drone ship)   | 5              |
| Success (ground pad)   | 3              |
| Controlled (ocean)     | 3              |
| Uncontrolled (ocean)   | 2              |
| Failure (parachute)    | 2              |
| Precluded (drone ship) | 1              |


```

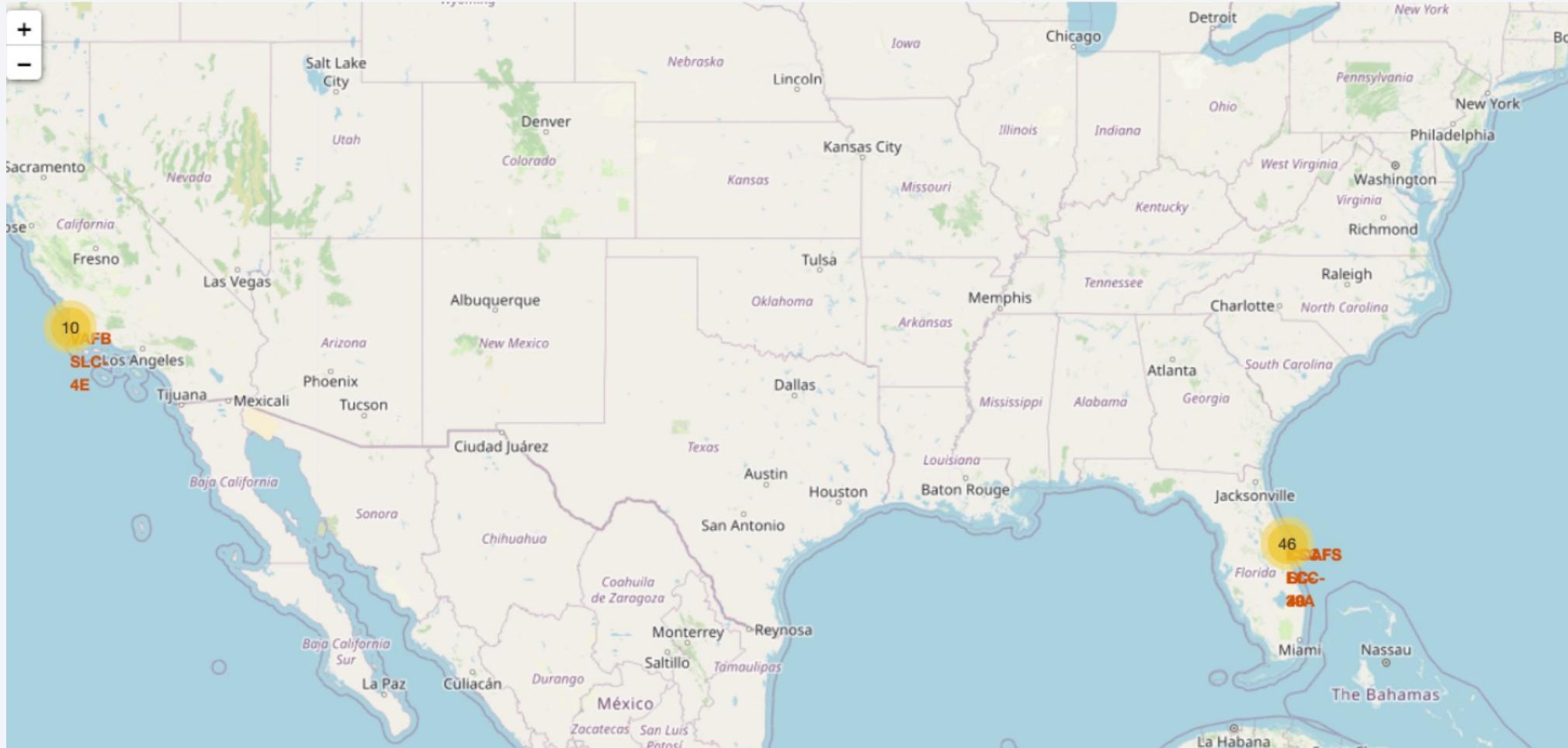
| Landing_Outcome        | count_outcomes |
|------------------------|----------------|
| No attempt             | 10             |
| Success (drone ship)   | 5              |
| Failure (drone ship)   | 5              |
| Success (ground pad)   | 3              |
| Controlled (ocean)     | 3              |
| Uncontrolled (ocean)   | 2              |
| Failure (parachute)    | 2              |
| Precluded (drone ship) | 1              |

The background of the slide is a photograph taken from space at night. It shows the curvature of the Earth against a dark blue-black void of space. City lights are visible as numerous small white and yellow dots, primarily concentrated in the lower right quadrant where the United States appears. In the upper right, the green and yellow glow of the aurora borealis is visible. The atmosphere of the Earth is thin and hazy, appearing as a light blue band near the horizon.

Section 3

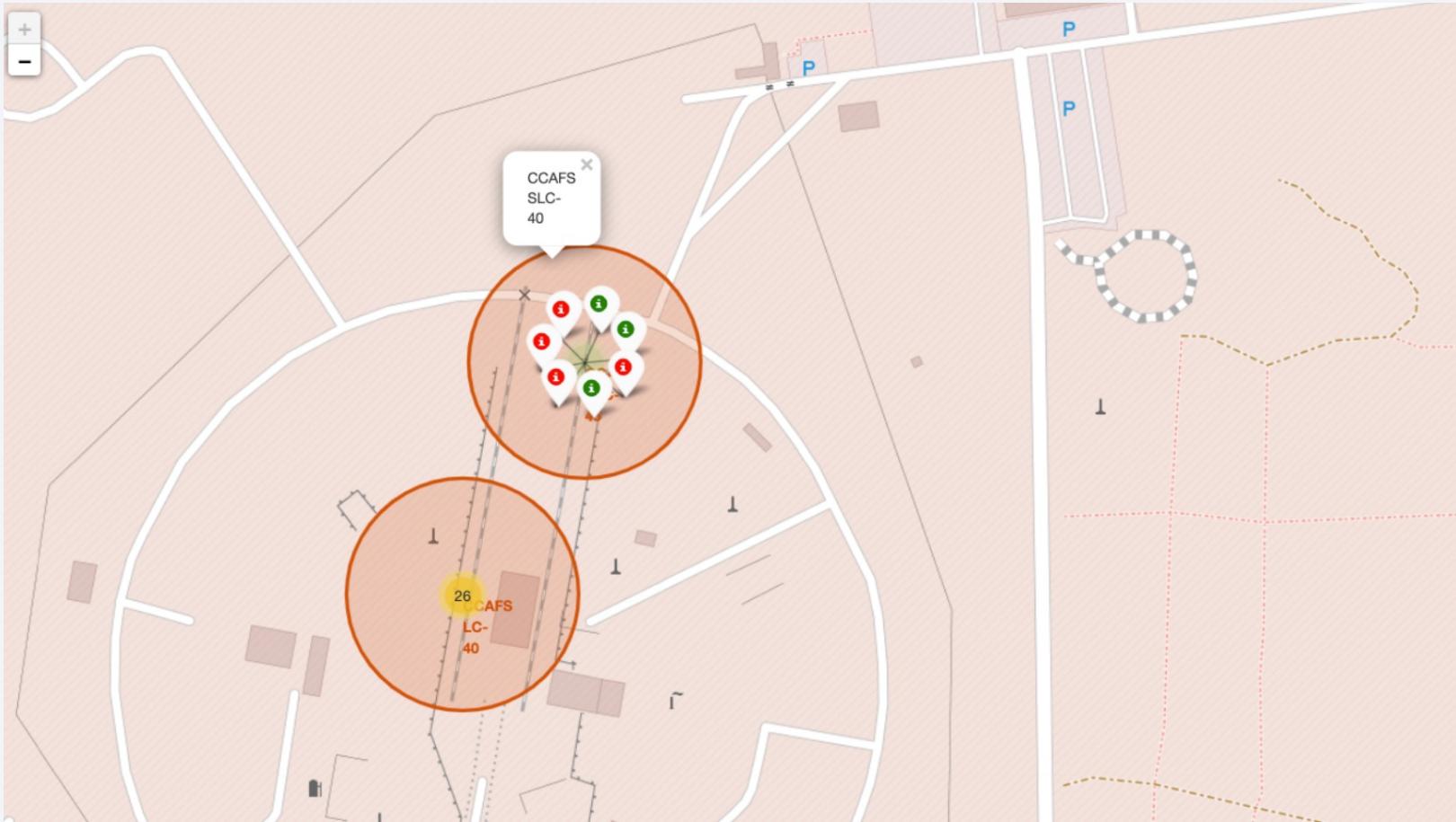
# Launch Sites Proximities Analysis

# Location of all the Launch Sites



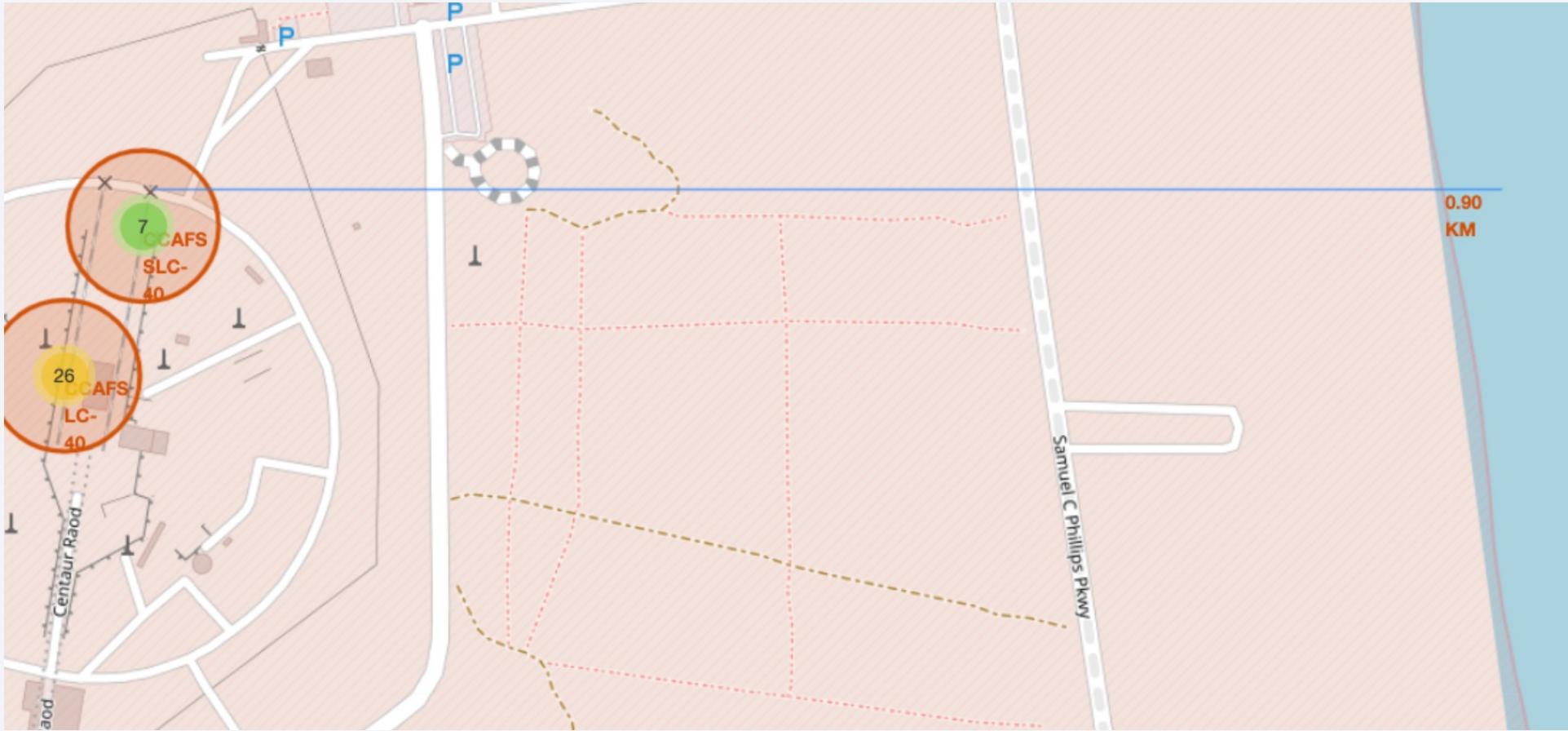
We can see that all the SpaceX launch sites are located inside the United States

# Detailed launch sites



# <Folium Map Screenshot 3>

---

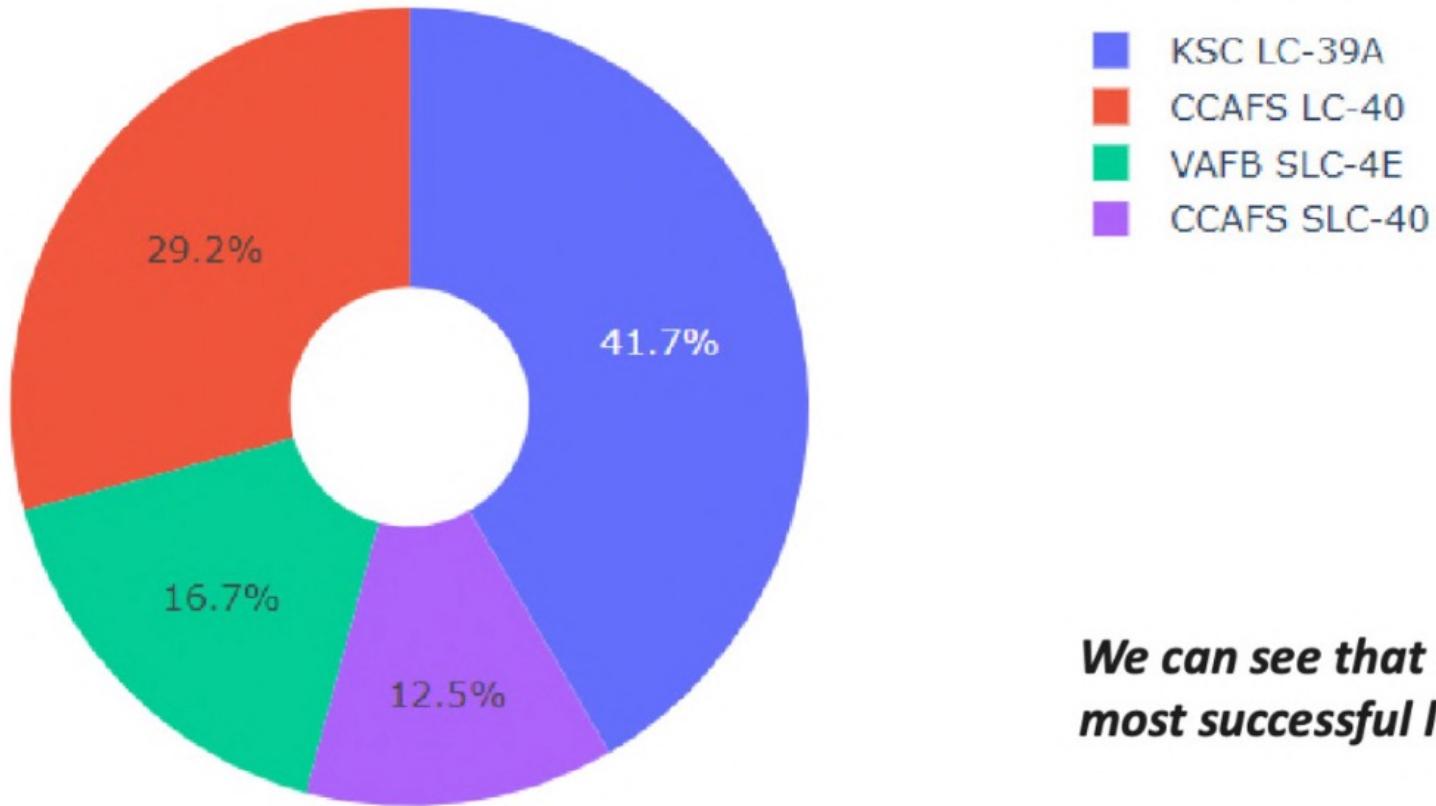


Section 4

# Build a Dashboard with Plotly Dash

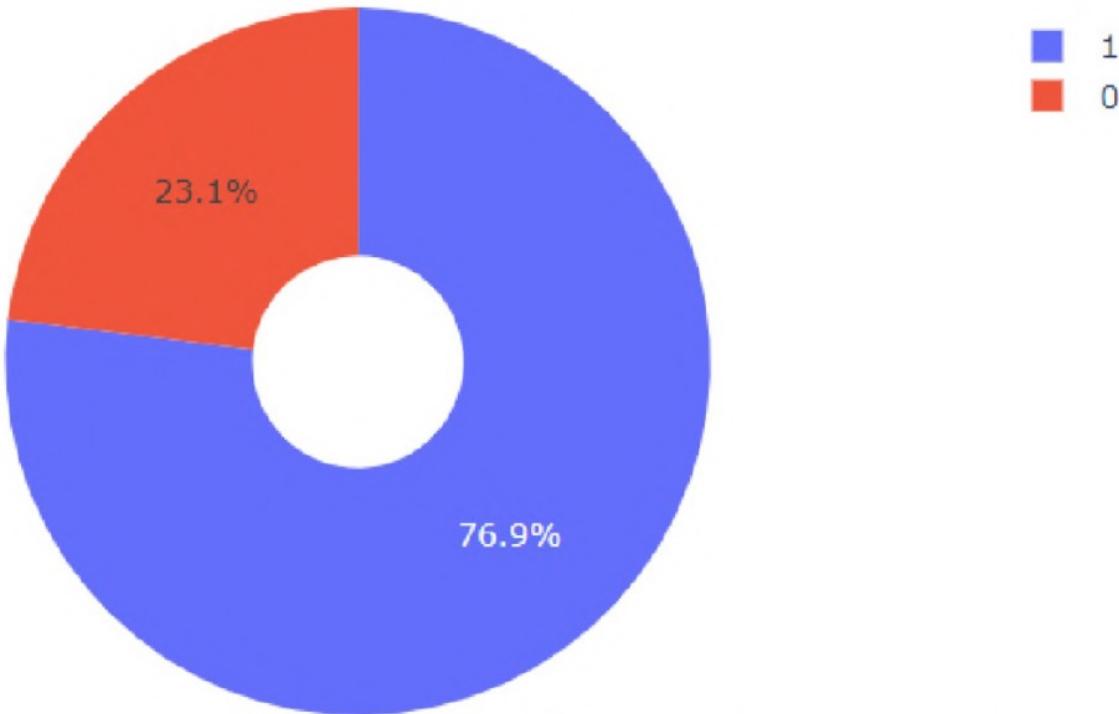


# The success percentage by each sites.



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

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

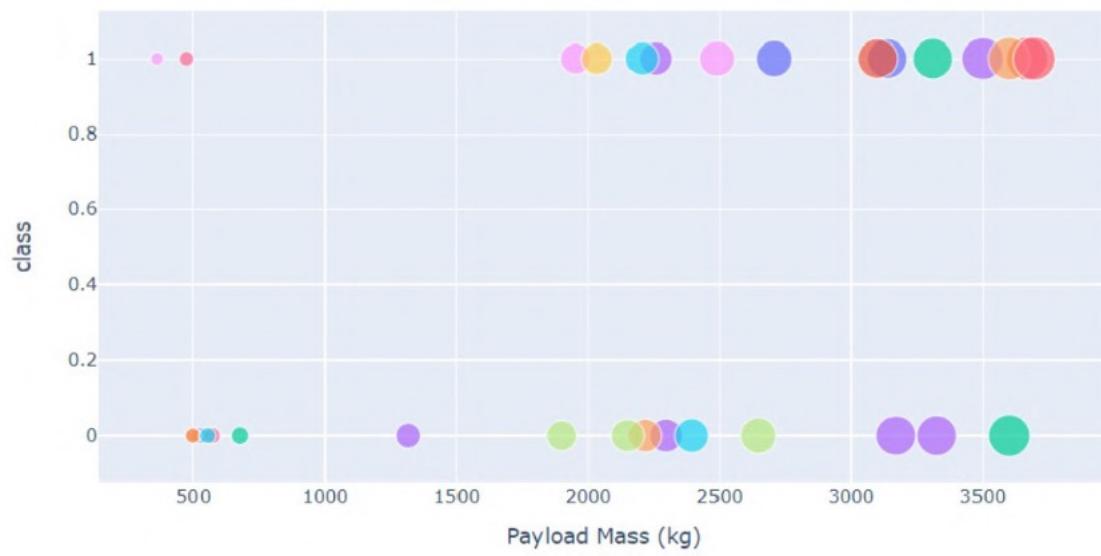


*KSC LC-39A achieved a 76.9% success rate while getting a 23.1% failure rate*

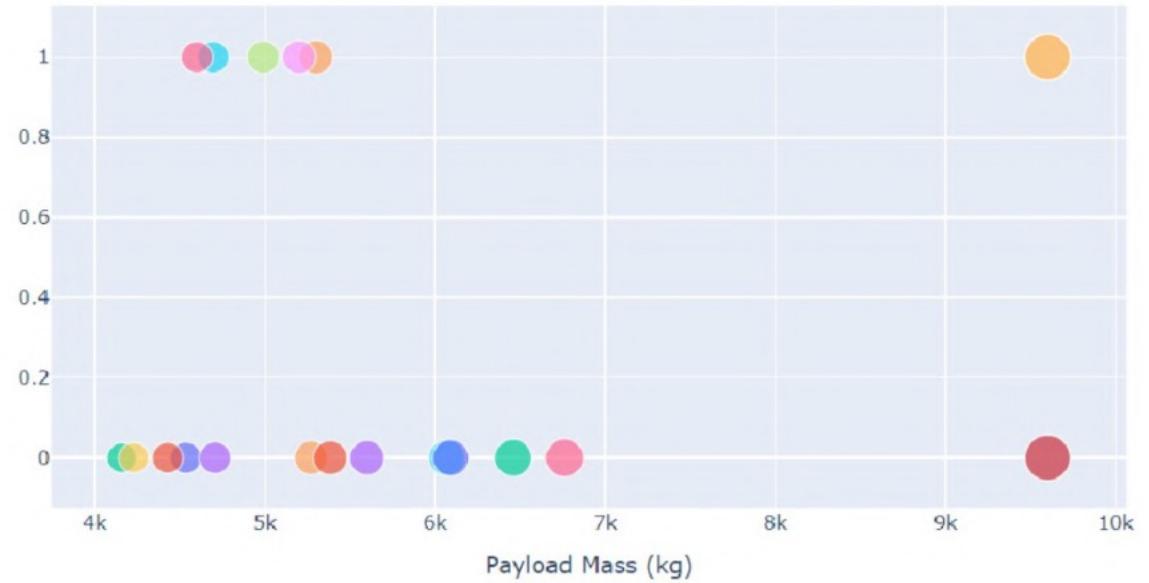
# Payload vs Launch Outcome Scatter Plot

- We can see that all the success rate for low weighted payload is higher than heavy weighted

***Low Weighted Payload 0kg – 4000kg***



***Heavy Weighted Payload 4000kg – 10000kg***



The background of the slide features a dynamic, abstract design. It consists of several thick, curved lines that transition from a bright yellow at the top right to a deep blue at the bottom left. These lines create a sense of motion and depth, resembling a tunnel or a stylized road. The overall effect is modern and professional.

Section 5

# Predictive Analysis (Classification)

# Classification Accuracy

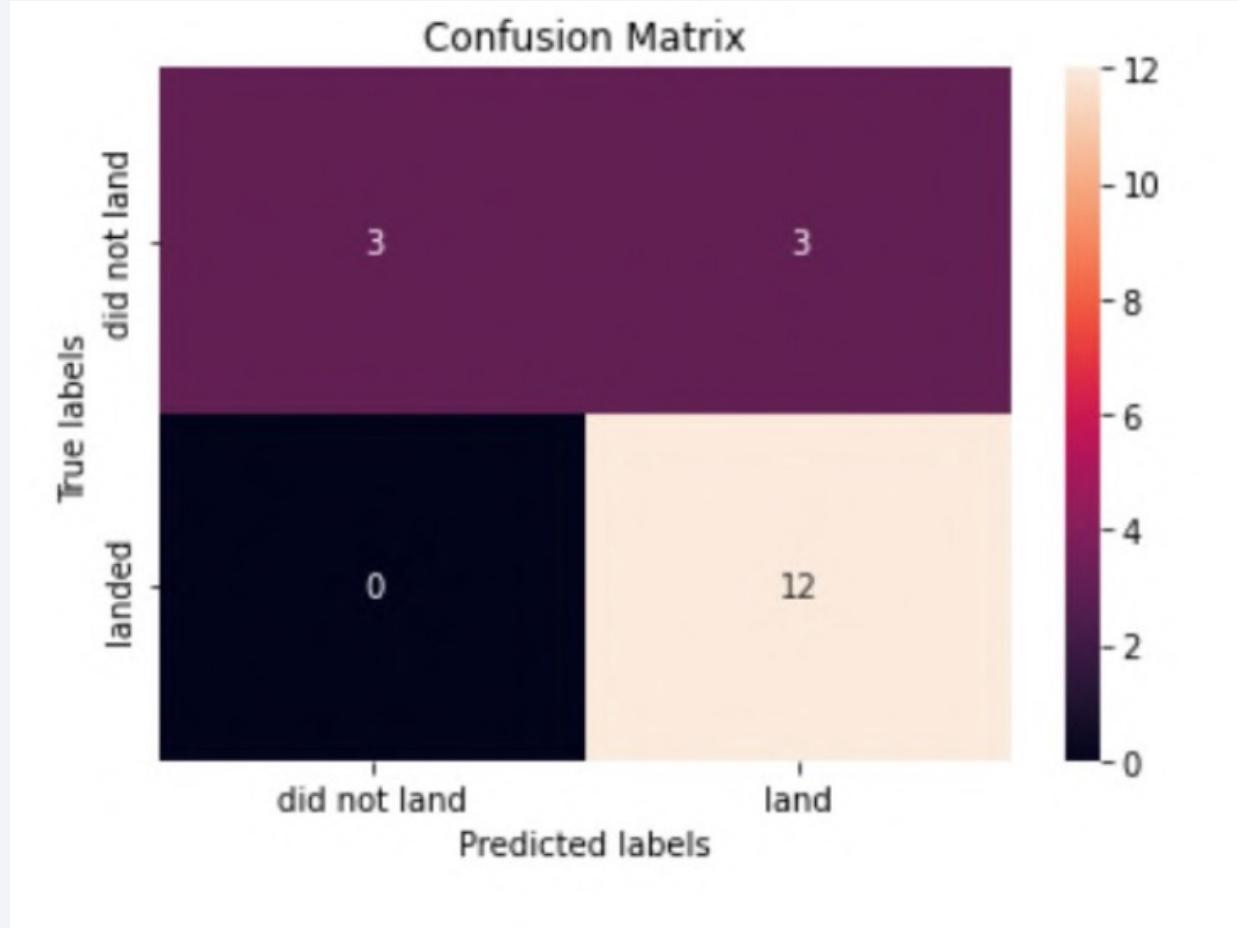
---

```
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

---



# Conclusions

---

- The Tree Classifier Algorithm is the best Machine Learning approach for this dataset.
- The low weighted payloads (which define as 4000kg and below) performed better than the heavy weighted payloads.
- Starting from the year 2013, the success rate for SpaceX launches is increased, directly proportional time in years to 2020, which it will eventually perfect the launches in the future.
- KSC LC-39A have the most successful launches of any sites; 76.9%
- SSO orbit have the most success rate; 100% and more than 1 occurrence.

# Appendix

---

- <https://github.com/abui-ualberta/spacex/>

Thank you!

