



IBM Developer  
SKILLS NETWORK

# Winning Space Race with Data Science

Mitchell Brazell  
2023 January 18



# 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

SpaceX is the most successful private company in its industry. They advertise the 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 related to the fact that SpaceX 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 SpaceX 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.

Section 1

# Methodology

# Methodology

---

## Executive Summary

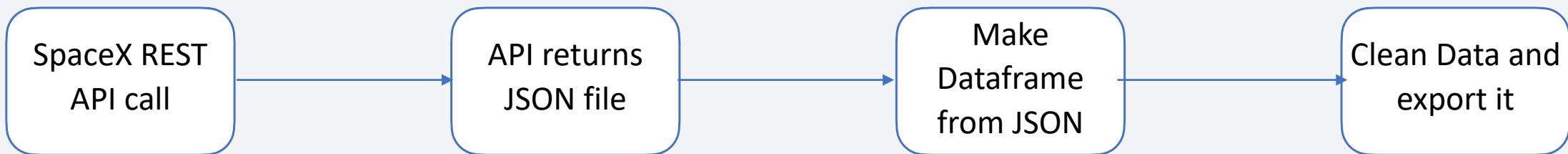
- Data collection methodology:
  - Data was collected using SpaceX API and web scraping from Wikipedia.
  - Data was collected using SpaceX API and web scraping from Wikipedia.
- Perform data wrangling
  - One-hot encoding was applied to categorical features
  - Dropping unnecessary columns
- Perform 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

---

- Datasets are collected using REST SpaceX API and web scraping Wikipedia

- The SpaceX REST API URL: [api.spacexdata.com/v4/](https://api.spacexdata.com/v4/)



- Wikipedia URL: [https://en.wikipedia.org/w/index.php?title=List\\_of\\_Falcon\\_9\\_and\\_Falcon\\_Heavy\\_launches&oldid=1027686922](https://en.wikipedia.org/w/index.php?title=List_of_Falcon_9_and_Falcon_Heavy_launches&oldid=1027686922)



# Data Collection – SpaceX API

## 1. Get Response from URL

```
spacex_url="https://api.spacexdata.com/v4/launches/past"  
  
response = requests.get(spacex_url)
```

## 2. Convert response to JSON

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

## 3. Transform Data

```
getBoosterVersion(data)  
getLaunchSite(data)  
getPayloadData(data)  
getCoreData(data)
```

## 4. Create Dictionary

```
launch_dict = {'FlightNumber': list(data['flight_number']),  
'Date': list(data['date']),  
'BoosterVersion':BoosterVersion,  
'PayloadMass':PayloadMass,  
'Orbit':Orbit,  
'LaunchSite':LaunchSite,  
'Outcome':Outcome,  
'Flights':Flights,  
'GridFins':GridFins,  
'Reused':Reused,  
'Legs':Legs,  
'LandingPad':LandingPad,  
'Block':Block,  
'ReusedCount':ReusedCount,  
'Serial':Serial,  
'Longitude': Longitude,  
'Latitude': Latitude}
```

## 5. Create Dataframe

```
# Create a data from launch_dict  
data_launch = pd.DataFrame.from_dict(launch_dict)
```

## 6. Filter Dataframe

```
# Hint data['BoosterVersion']!='Falcon 1'  
data_falcon9 = data_launch[data_launch['BoosterVersion'] != 'Falcon 1']
```

## 7. Export File

```
data_falcon9.to_csv('dataset_part_1.csv', index=False)
```

[Link to Notebook](#)

# Data Collection - Scraping

## 1. Getting Response from HTML

```
html_data = requests.get(static_url)
html_data.status_code
```

## 2. Create BeautifulSoup Object

```
soup = BeautifulSoup(html_data.text, 'html.parser')
```

## 3. Find all tables

```
html_tables=soup.find_all('table')
```

## 4. Get column names

```
for element in first_launch_table.find_all('th'):
    name = extract_column_from_header(element)
    if name is not None and len(name) > 0:
        column_names.append(name)
```

## 5. Create dictionary

```
launch_dict= dict.fromkeys(column_names)

# Remove an irrelevant column
del launch_dict['Date and time ( )']

# Let's initial the launch_dict with each
launch_dict['Flight No.']= []
launch_dict['Launch site']= []
launch_dict['Payload']= []
launch_dict['Payload mass']= []
launch_dict['Orbit']= []
launch_dict['Customer']= []
launch_dict['Launch outcome']= []
# Added some new columns
launch_dict['Version Booster']= []
launch_dict['Booster landing']= []
launch_dict['Date']= []
launch_dict['Time']= []
```

## 6. Add data to keys

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

## 7. Create dataframe

```
df=pd.DataFrame(launch_dict)
```

## 8. Export to file

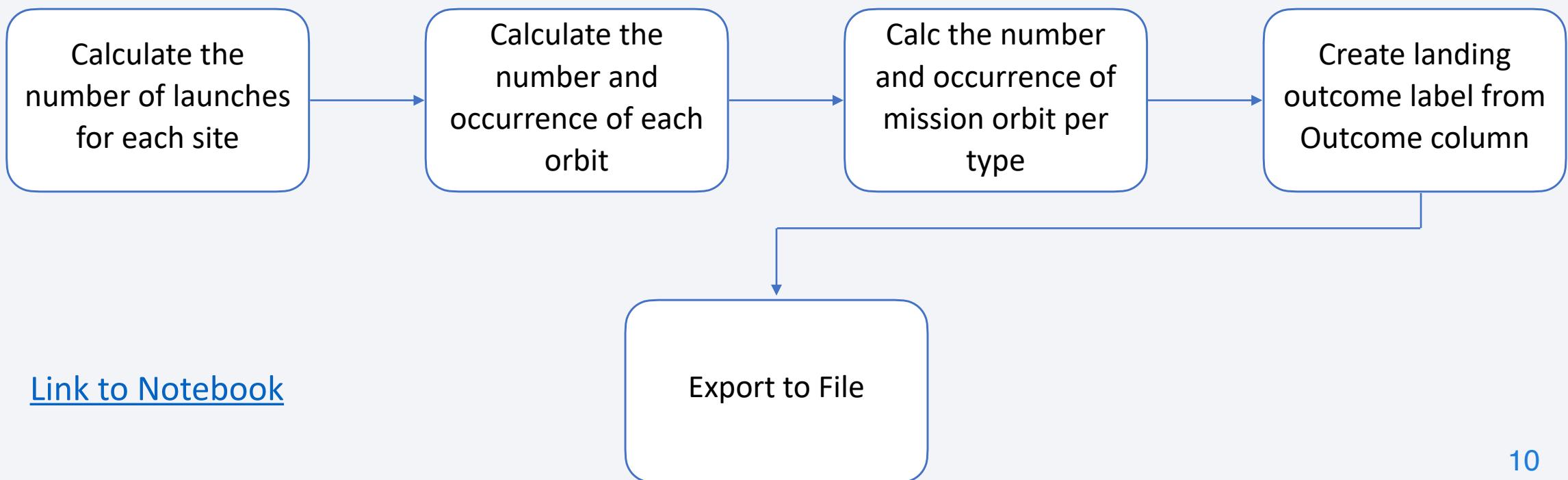
```
df.to_csv('spacex_web_scraped.csv', index=False)
```

[Link to Notebook](#)

# Data Wrangling

---

- The dataset included several cases where the booster did not land successfully
  - True Ocean, True RTLS, True ASDS means that the mission was successful
  - False Ocean, False RTLS, False ASDS means that the mission was a failure
- We needed to transform string variables into categorical variables where 1 means the mission has been successful and 0 means that the mission was a failure.



# EDA with Data Visualization

---

- Scatter Graphs
  - Flight Number vs. Payload Mass
  - Flight Number vs. Launch Site
  - Payload vs. Launch Site
  - Orbit vs. Flight Number
  - Payload vs. Orbit Type
  - Orbit vs Payload Mass
- Bar Graph
  - Success Rate vs. Orbit
- Line Graph
  - Success Rate vs. Year

[Link to Notebook](#)

# EDA with SQL

---

**We performed SQL queries to gather and understand data from dataset:**

- Displaying the names of the unique launch sites in the space mission.
- Display 5 records where launch sites begin with the string 'CCA'
- 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 when the first successful landing outcome in ground pad was achieved.
- List the names of the boosters which have success in drone ship 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.
- List the records which will display the month names, failure, landing\_outcomes in drone ship, booster versions, launch\_site for the months in year 2015.
- Rank the count of successful landing\_outcomes between the date 04-06-2010 and 20-03-2017 in descending order.

[Link to Notebook](#)

# Build an Interactive Map with Folium

---

Folium map object is a map centered on NASA Johnson Space Center at Houston, Texas

- Red circle at NASA Johnson Space Center's coordinate with label showing its name (folium.Circle, folium.map.Marker).
- Red circles at each launch site coordinates with label showing launch site name (folium.Circle, folium.map.Marker, folium.features.DivIcon).
- The grouping of points in a cluster to display multiple and different information for the same coordinates (folium.plugins.MarkerCluster).
- Markers to show successful and unsuccessful landings. **Green** for successful landing and **Red** for unsuccessful landing. (folium.map.Marker, folium.Icon).
- Markers to show distance between launch site to key locations (railway, highway, coastway, city) and plot a line between them . (folium.map.Marker, folium.PolyLine, folium.features.DivIcon)

These objects are created in order to understand better the problem and the data. We can show easily all launch sites, their surroundings and the number of successful and unsuccessful landings.

[Link to Notebook](#)

# Build a Dashboard with Plotly Dash

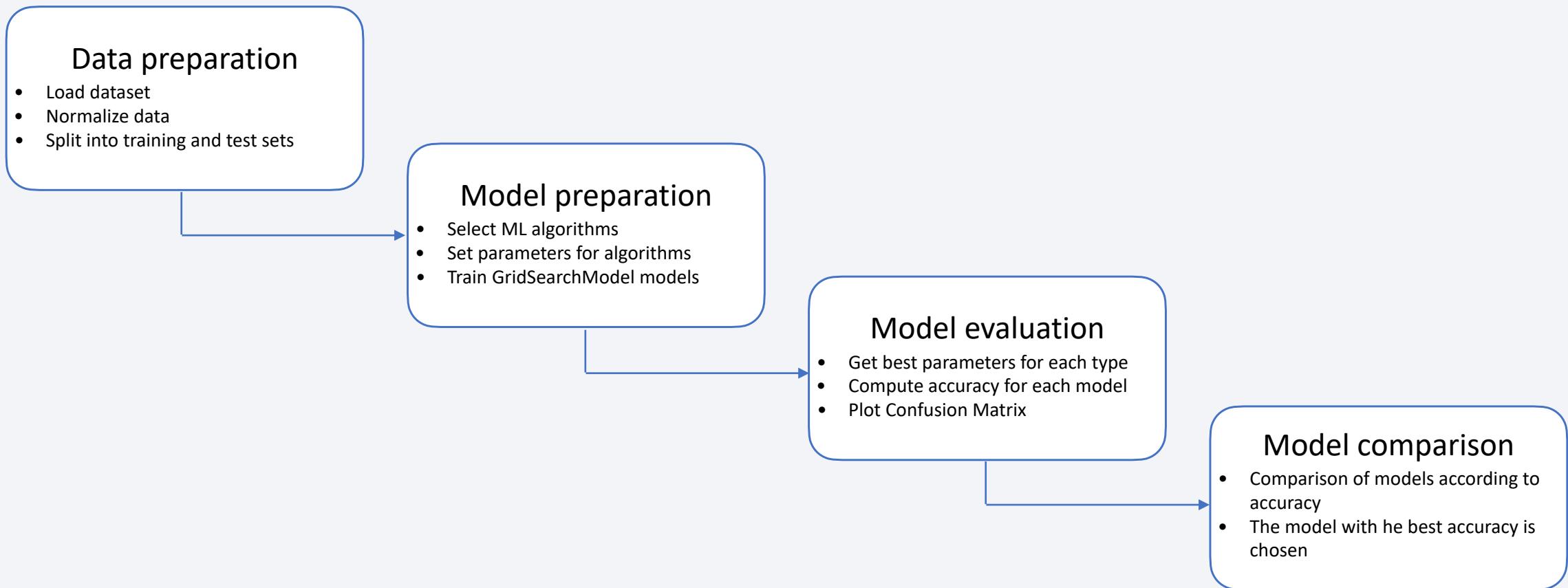
---

- The Dashboard has a dropdown, pie chart, range slider, and scatter plot components.
  - Dropdown allows the user to choose the launch site
  - Pie chart shows the total success and total failures for the launch site chosen above
  - Rangeslider allows a user to select a payload mass in a fixed range
  - Scatter chart shows the relationship between two variables

[Link to Code](#)

# Predictive Analysis (Classification)

---

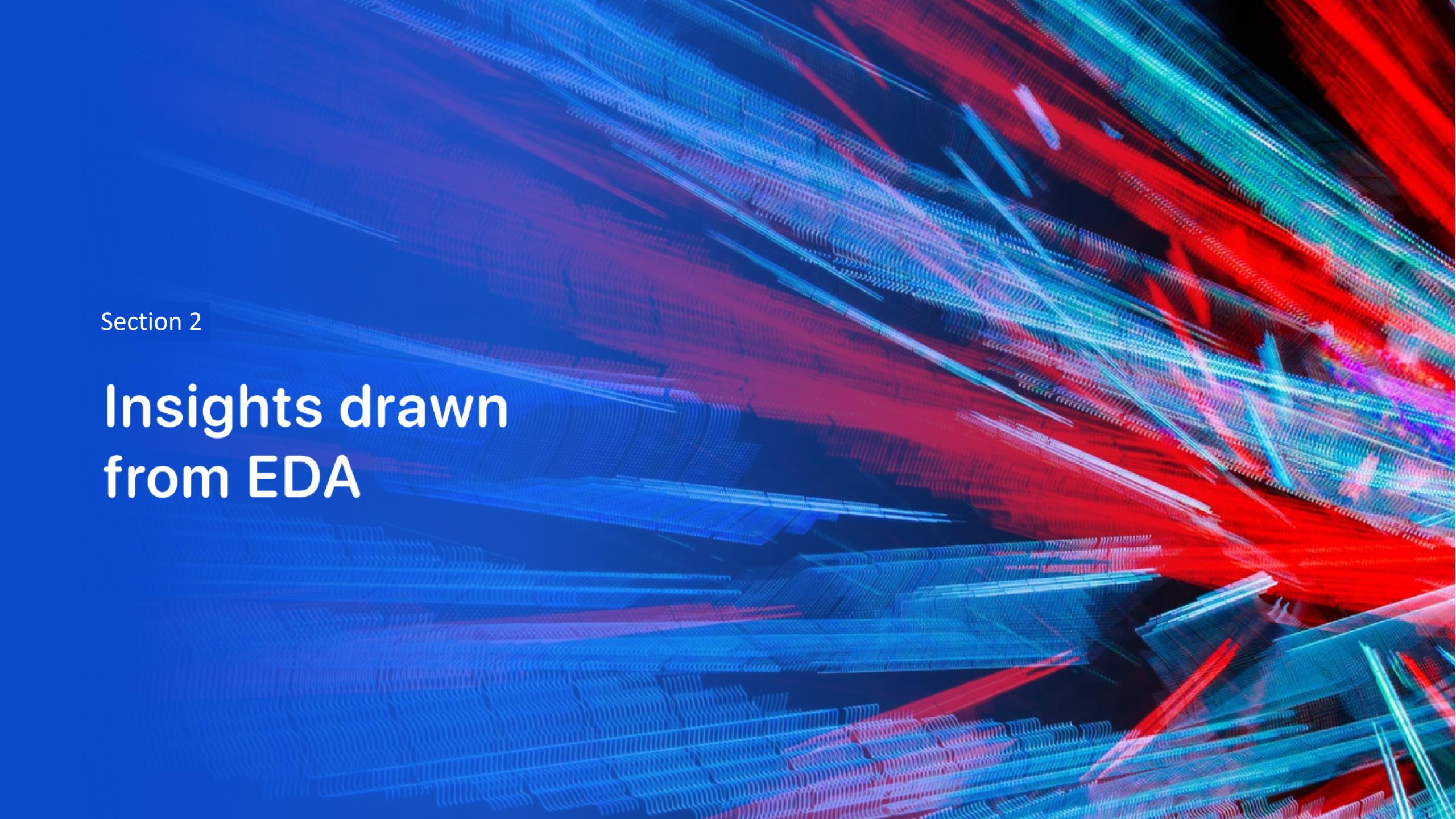


[Link to Notebook](#)

# Results

---

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

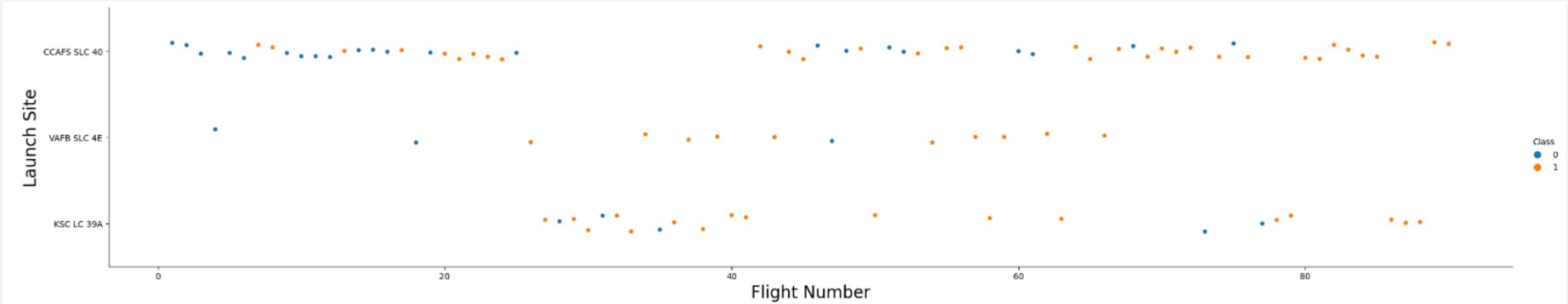
The background of the slide features a complex, abstract digital pattern. It consists of numerous thin, glowing lines that create a sense of depth and motion. The colors used are primarily shades of blue, red, and purple, which are bright against a dark, almost black, background. These lines form a grid-like structure that is more dense and vibrant towards the right side of the frame, while appearing more sparse and blurred towards the left.

Section 2

## Insights drawn from EDA

# Flight Number vs. Launch Site

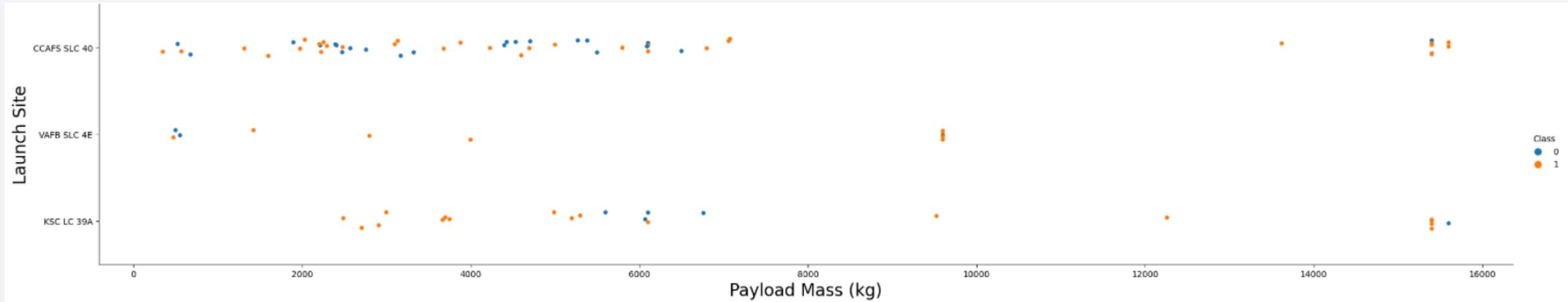
---



This graphic shows that for each site, the success rate is increasing.

# Payload vs. Launch Site

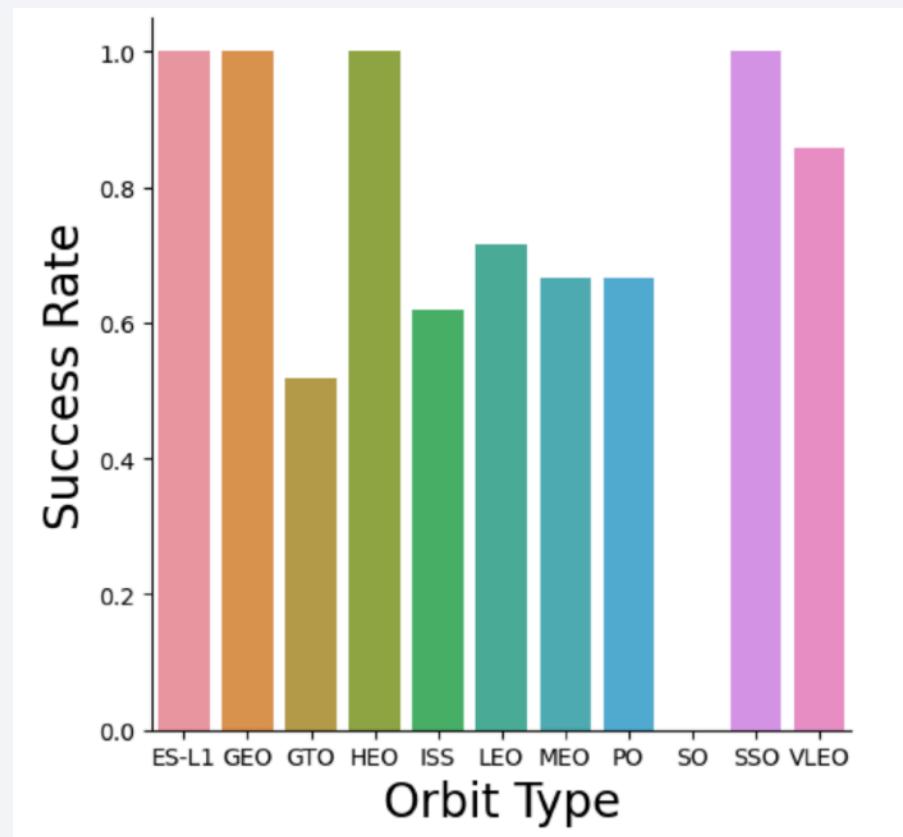
---



Depending on the launch site, a heavier payload may be a consideration for a successful landing, although a payload that is too heavy may result in failure.

# Success Rate vs. Orbit Type

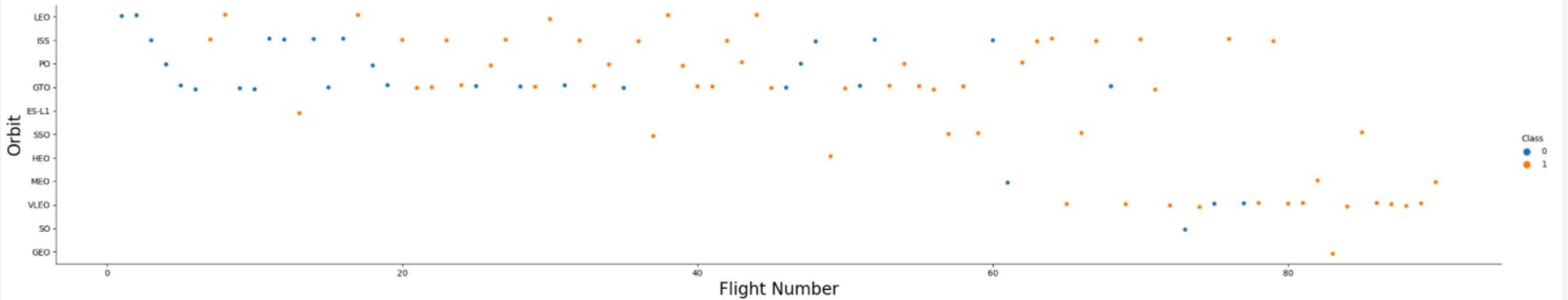
---



Here you can see the success rate for different orbit types. Note that ES-L1, GEO, HEP, SSO and the best success rates.

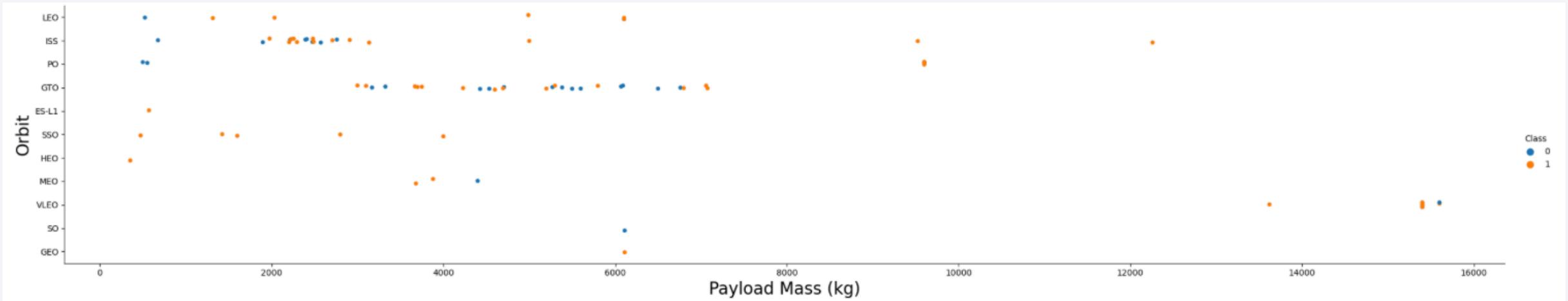
# Flight Number vs. Orbit Type

---



Notice that the success rate increases with the number of flights for the LEO orbit. For some orbits, like the GTO, there is no relation between the success rate and the number of flights.

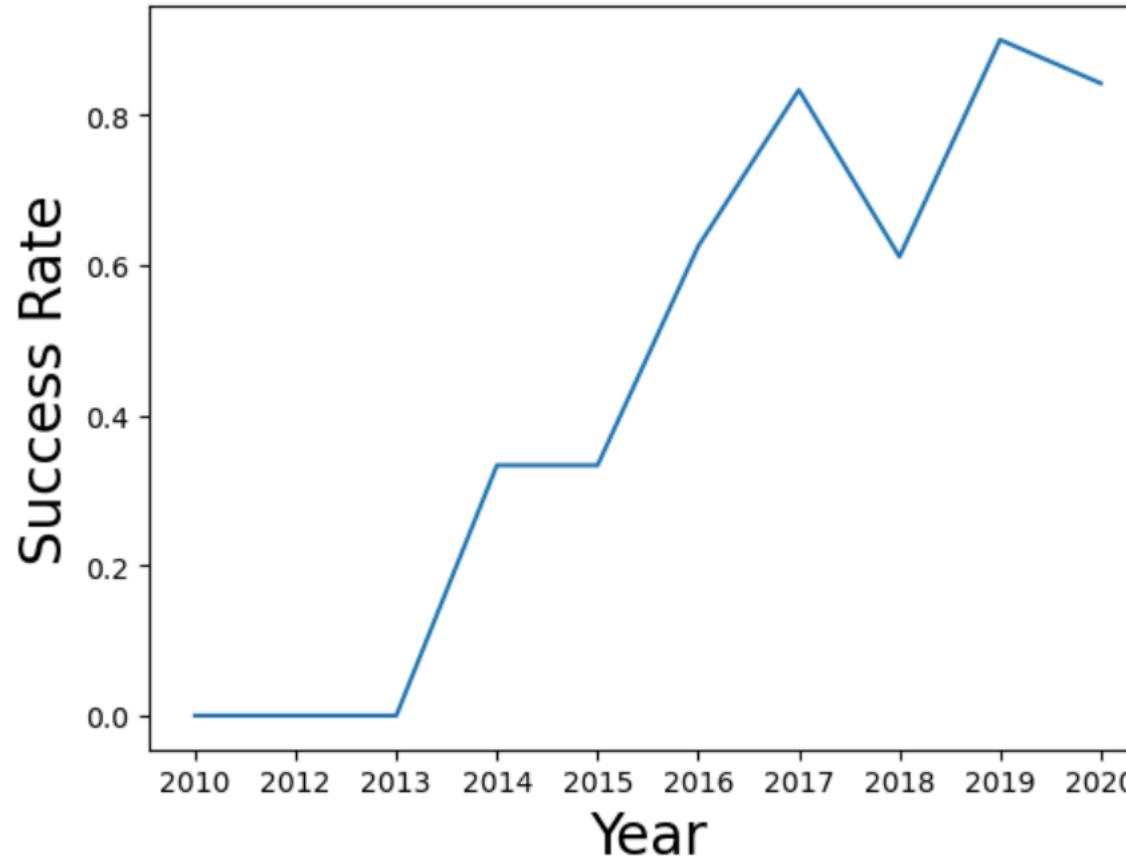
# Payload vs. Orbit Type



The weight of the payloads can have a great influence on the success rate of the launches in certain orbits. Heavier payloads improve the success rates of some orbits, like the LEO orbit. However, decreasing the payload of others, like the GTO orbit, improves the success of launch.

# Launch Success Yearly Trend

---



Beginning in 2013 there is a notable increase in the rocket success of SpaceX launches

# All Launch Site Names

---

SQL Query

```
SELECT DISTINCT "LAUNCH_SITE" FROM SPACEXTBL
```

Results

Launch_Site
CCAFS LC-40
VAFB SLC-4E
KSC LC-39A
CCAFS SLC-40

Explanation

The use of DISTINCT in the query removes duplicates in the column LAUNCH\_SITE.

# Launch Site Names Begin with 'CCA'

---

## SQL Query

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

## Results

Date	Time (UTC)	Booster_Version	Launch_Site	Payload	PAYLOAD_MASS_KG_	Orbit	Customer
04-06-2010	18:45:00	F9 v1.0 B0003	CCAFS LC-40	Dragon Spacecraft Qualification Unit		0 LEO	SpaceX
08-12-2010	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
22-05-2012	07:44:00	F9 v1.0 B0005	CCAFS LC-40	Dragon demo flight C2	525	LEO (ISS)	NASA (COTS)
08-10-2012	00:35:00	F9 v1.0 B0006	CCAFS LC-40	SpaceX CRS-1	500	LEO (ISS)	NASA (CRS)
01-03-2013	15:10:00	F9 v1.0 B0007	CCAFS LC-40	SpaceX CRS-2	677	LEO (ISS)	NASA (CRS)

## Explanation

The WHERE function followed by the LIKE function filters the launch sites column to only contain this that include CCA. LIMIT 5 shows 5 records from the table.

# Total Payload Mass

---

## SQL Query

```
SELECT SUM("PAYLOAD_MASS__KG_") FROM SPACEXTBL WHERE "CUSTOMER" = 'NASA (CRS)'
```

## Results

SUM("PAYLOAD_MASS__KG_")
45596

## Explanation

This query returns the sum of all payload masses where the customer is NASA (CRS).

# Average Payload Mass by F9 v1.1

---

## SQL Query

```
SELECT AVG("PAYLOAD_MASS_KG_") FROM SPACEXTBL WHERE "BOOSTER_VERSION" LIKE "%F9 v1.1%"
```

## Results

AVG("PAYLOAD_MASS_KG_")
2534.6666666666665

## Explanation

This query returns the average of all payload masses where the booster version contains the string F9 v1.1.

# First Successful Ground Landing Date

---

## SQL Query

```
SELECT MIN("DATE") FROM SPACEXTBL WHERE "Landing _Outcome" LIKE '%Success%'
```

## Results

MIN("DATE")
01-05-2017

## Explanation

This query selects the oldest successful landing. The WHERE function filters the dataset in order to keep only records where the landing was successful and the MIN function selects the record with the oldest date.

# Successful Drone Ship Landing with Payload between 4000 and 6000

---

## SQL Query

```
%sql SELECT "BOOSTER_VERSION" FROM SPACEXTBL WHERE "LANDING_OUTCOME" = 'Success (drone ship)' \
AND "PAYLOAD_MASS_KG_" > 4000 AND "PAYLOAD_MASS_KG_" < 6000;
```

## Results

Booster_Version
F9 FT B1022
F9 FT B1026
F9 FT B1021.2
F9 FT B1031.2

## Explanation

This query returns the booster version where the landing was successful and the payload was between 4000 and 6000 kg. The WHERE function and the AND function filter the dataset.

# Total Number of Successful and Failure Mission Outcomes

---

## SQL Query

```
%sql SELECT (SELECT COUNT("MISSION_OUTCOME") FROM SPACEXTBL WHERE "MISSION_OUTCOME" LIKE '%Success%') AS SUCCESS, \  
          (SELECT COUNT("MISSION_OUTCOME") FROM SPACEXTBL WHERE "MISSION_OUTCOME" LIKE '%Failure%') AS FAILURE
```

## Results

SUCCESS	FAILURE
100	1

## Explanation

The first subquery counts the successful missions. The second subquery counts the unsuccessful missions. The WHERE function and the LIKE function filters the mission outcome. And lastly, the COUNT function counts the filtered records.

# Boosters Carried Maximum Payload

---

## SQL Query

```
%sql SELECT DISTINCT "BOOSTER_VERSION" FROM SPACEXTBL \
WHERE "PAYLOAD_MASS_KG_" = (SELECT max("PAYLOAD_MASS_KG_") FROM SPACEXTBL)
```

## Results

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

## Explanation

The subquery filters the data by returning only the heaviest payloads using the MAX function. The main query uses these results and returns the unique booster version with the heaviest payloads.

# 2015 Launch Records

---

## SQL Query

```
%sql SELECT substr("DATE", 4, 2) AS MONTH, "BOOSTER_VERSION", "LAUNCH_SITE" FROM SPACEXTBL\\
WHERE "LANDING_OUTCOME" = 'Failure (drone ship)' and substr("DATE",7,4) = '2015'
```

## Results

MONTH	Booster_Version	Launch_Site
01	F9 v1.1 B1012	CCAFS LC-40
04	F9 v1.1 B1015	CCAFS LC-40

## Explanation

This query returns month, booster version, and launch site where the landing was unsuccessful and the landing date was before 2015.

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

---

## SQL Query

```
%sql SELECT "LANDING _OUTCOME", COUNT("LANDING _OUTCOME") FROM SPACEXTBL\
WHERE "DATE" >= '04-06-2010' and "DATE" <= '20-03-2017' and "LANDING _OUTCOME" LIKE '%Success%'\
GROUP BY "LANDING _OUTCOME" \
ORDER BY COUNT("LANDING _OUTCOME") DESC ;
```

## Results

Landing _Outcome	COUNT("LANDING _OUTCOME")
Success	20
Success (drone ship)	8
Success (ground pad)	6

## Explanation

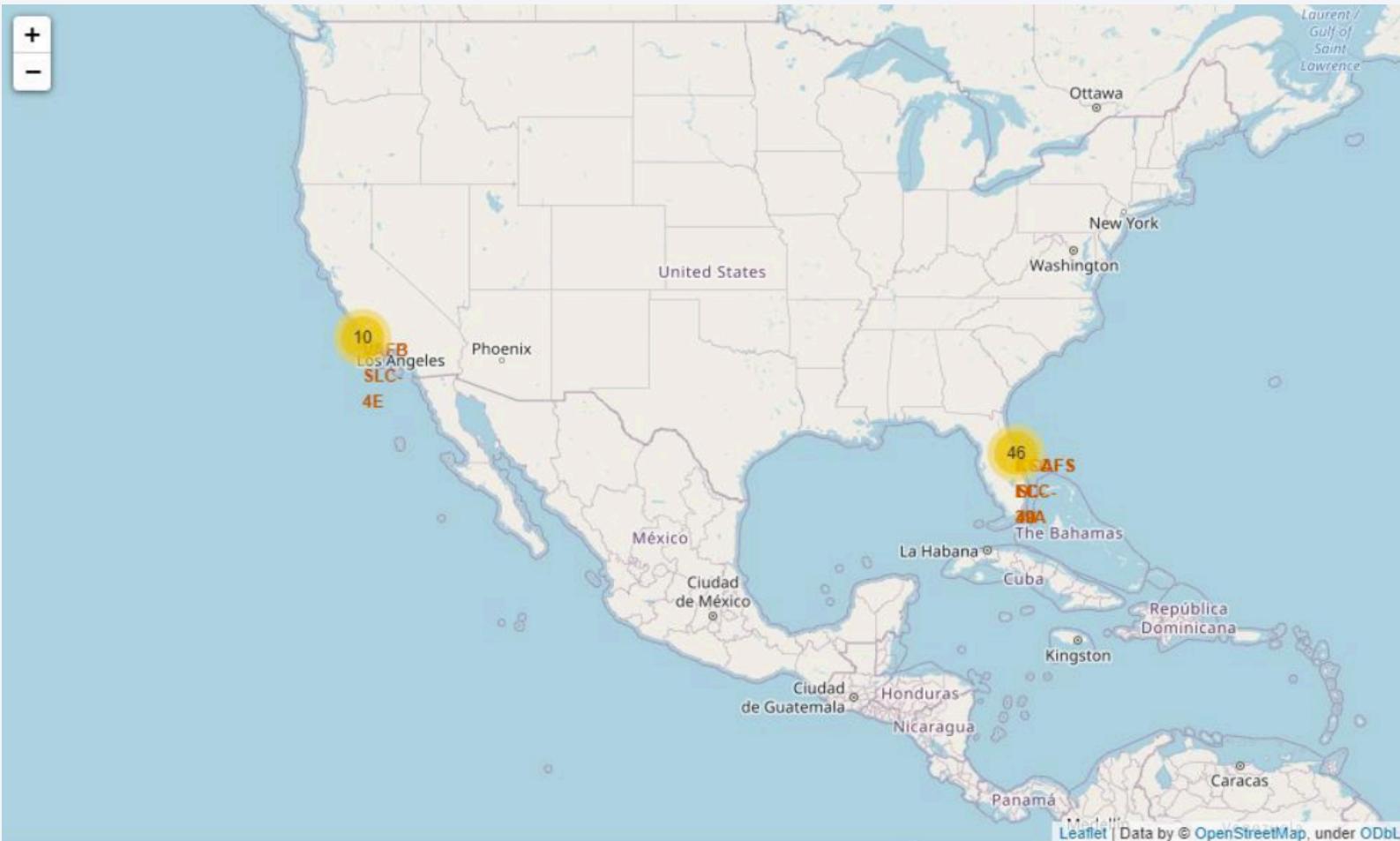
This query returns the landing outcomes and their count where the mission was successful and the date is between 06/04/10 and 03/20/17.

The background of the slide is a photograph taken from space at night. It shows the curvature of the Earth against a dark blue and black void of space. City lights are visible as numerous small white and yellow dots, primarily concentrated in the lower right quadrant where a large urban area is illuminated. In the upper right, there is a bright green and yellow glow, likely representing the Aurora Borealis or a similar atmospheric phenomenon.

Section 3

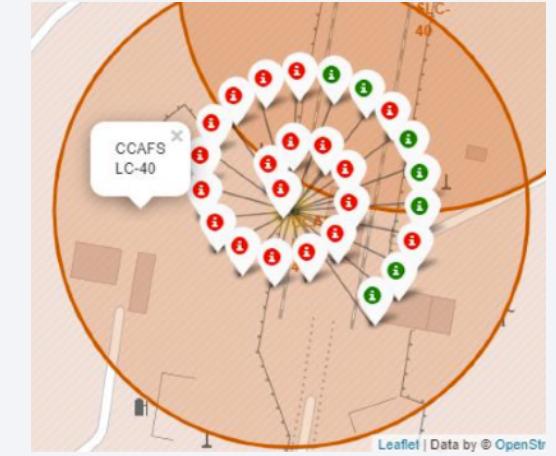
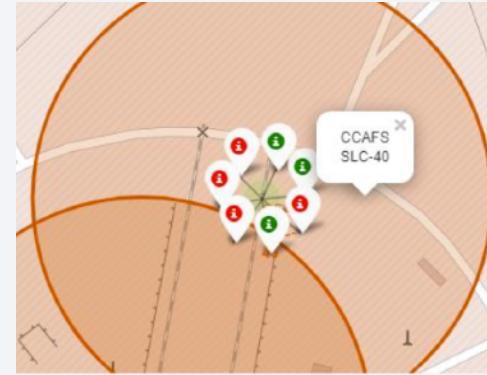
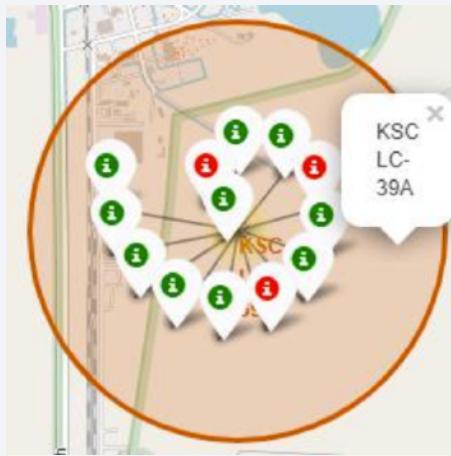
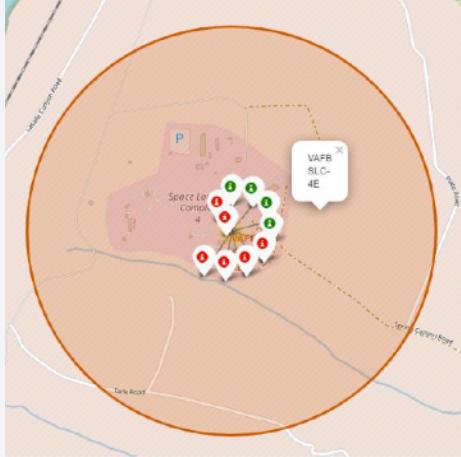
# Launch Sites Proximities Analysis

# Folium Map - Launch Sites



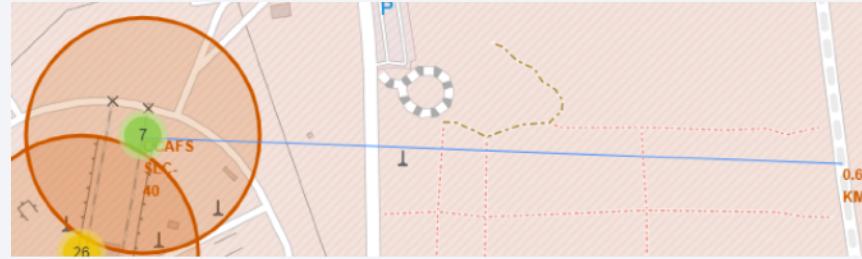
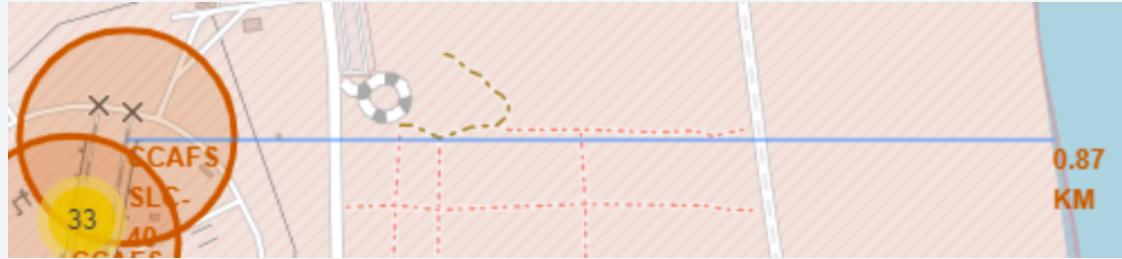
SpaceX launch sites are located on the coasts in the US.

# Folium Map - Labeled Markers

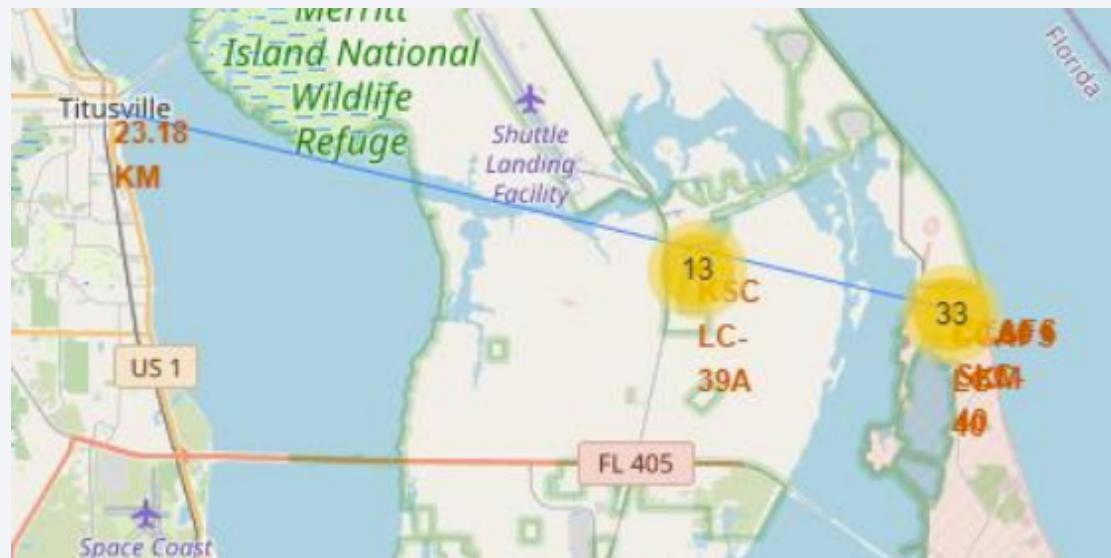


The **green** markers represent successful launches and the **red** markers represent unsuccessful launches.

# Folium Map - Distances between significant markers

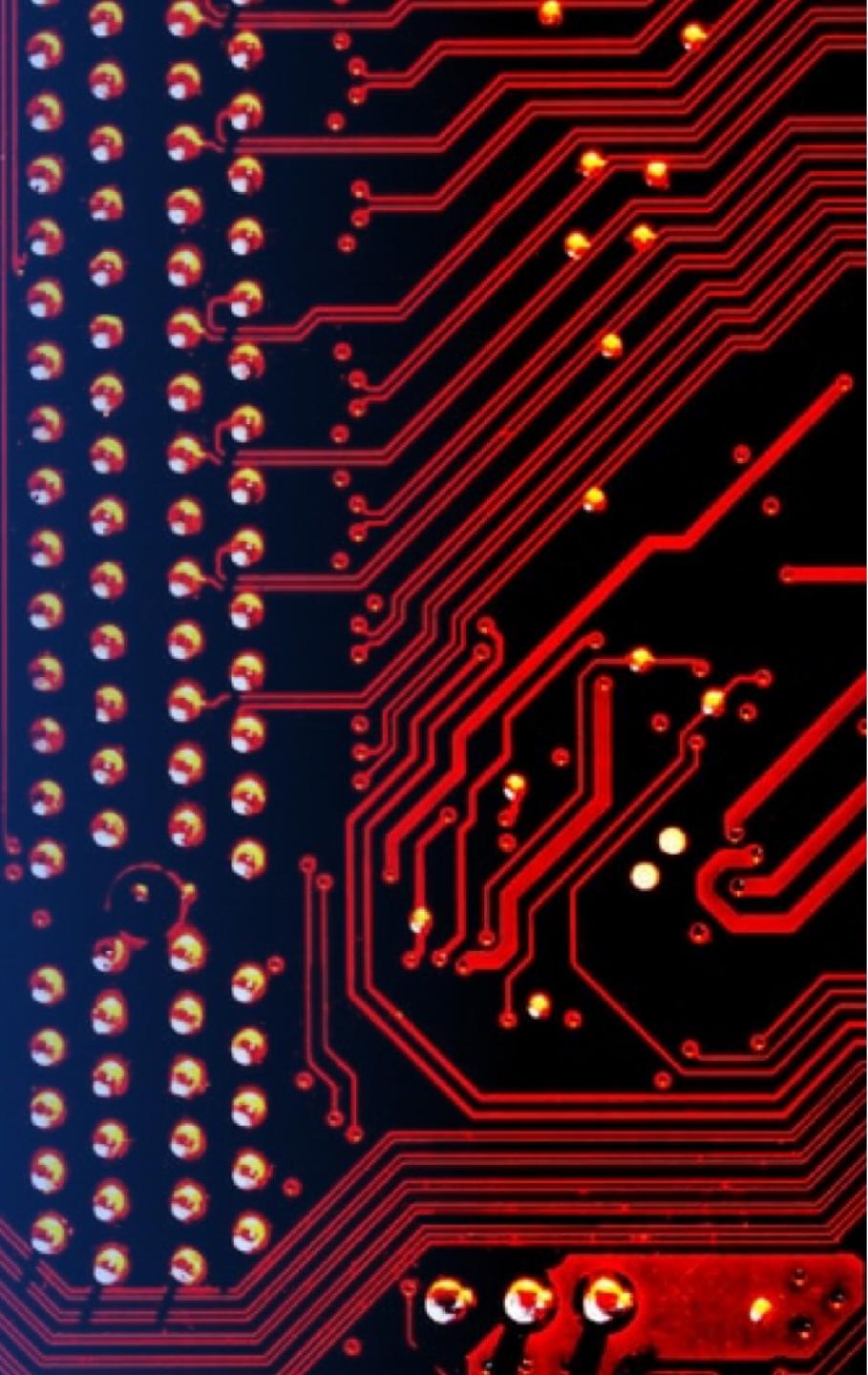


The distance markers allow SpaceX to measure the distance to major map elements and ensure the safety of residents in the area. Accurately measuring these distances ensures proper safety protocols can be followed.



Section 4

# Build a Dashboard with Plotly Dash



# Dashboard - Success by Site

---

Total Success Launches by Site

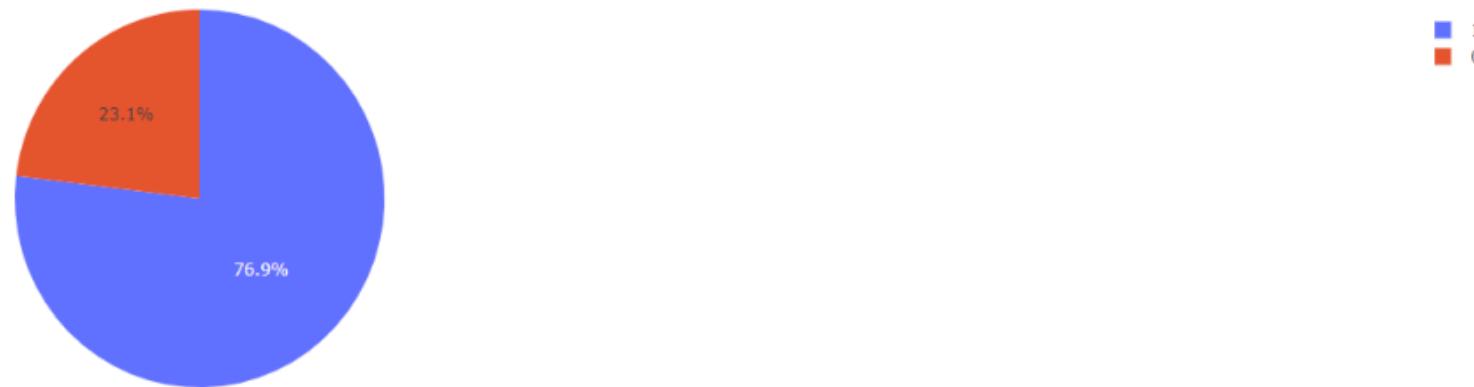


You can see that KSC-39A has the best success rate for launches

# Dashboard - Success Rate for launches at KSC LC-39A

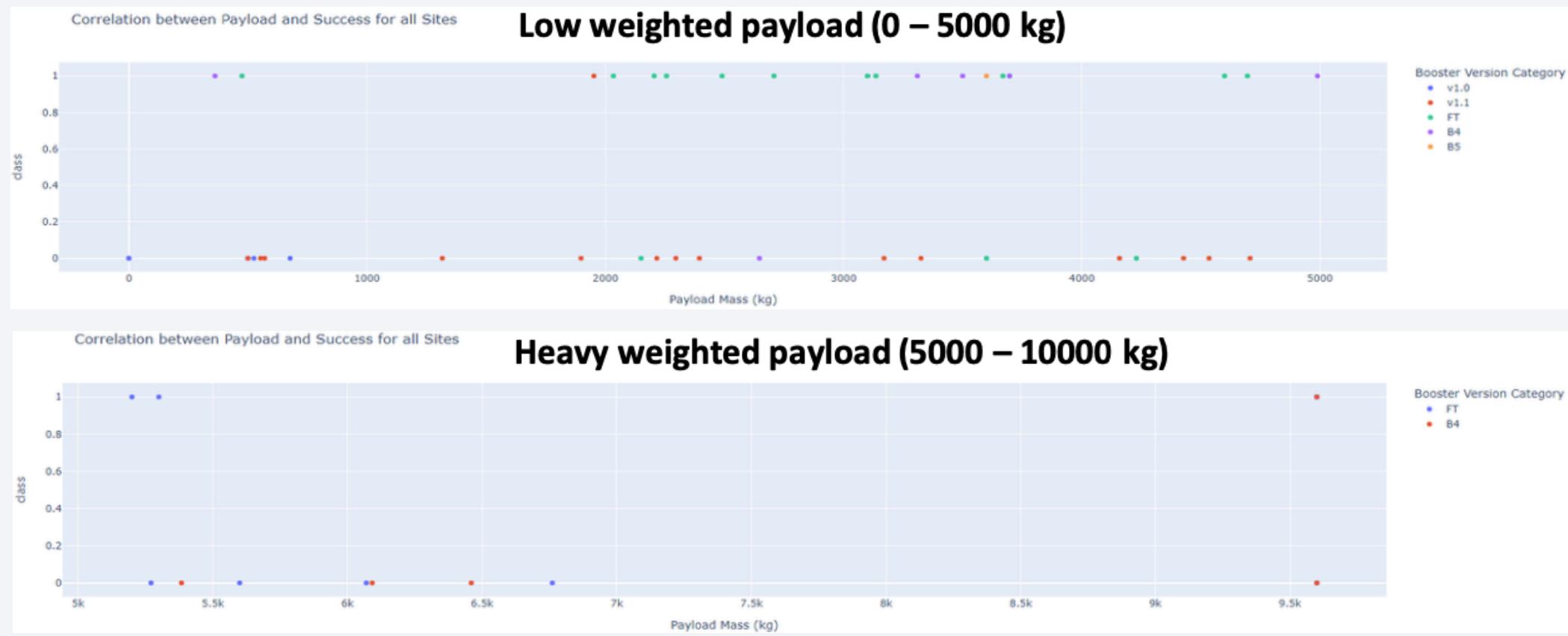
---

Total Success Launches for Site KSC LC-39A



The KSC LC-39A site has a success rate of 76.9% while having a failure rate of 23.1%.

# Dashboard - Payload Mass vs. Outcome for all sites



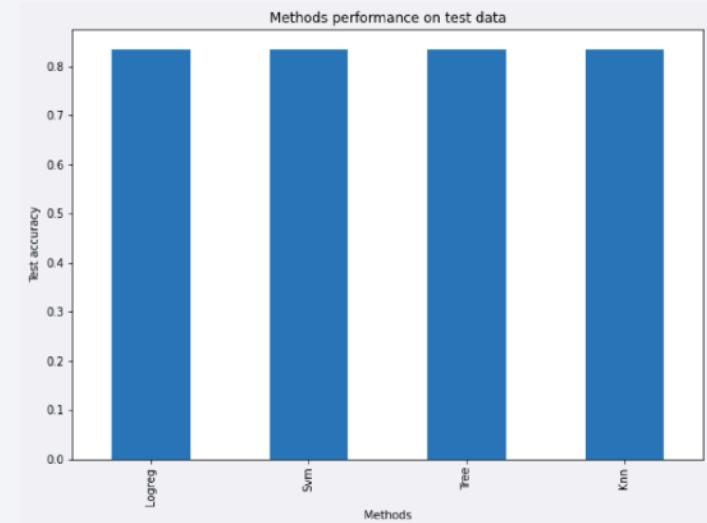
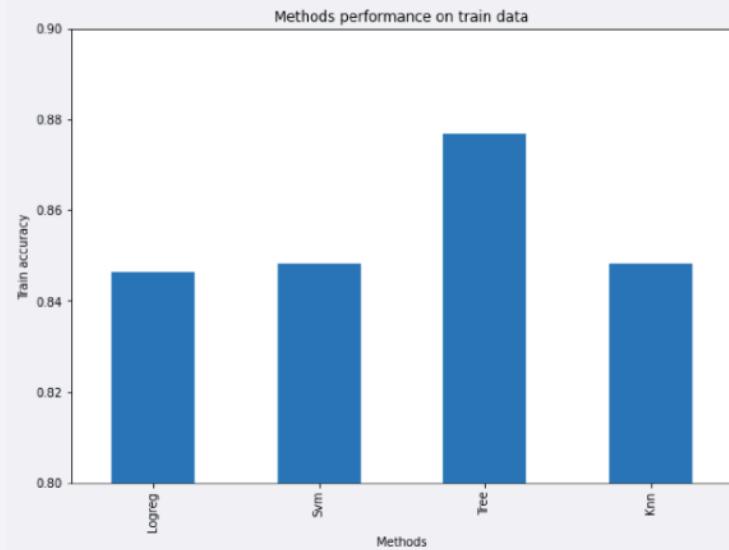
Low weighted payloads have better success rates than heave payloads.

Section 5

# Predictive Analysis (Classification)

# Classification Accuracy

	Accuracy Train	Accuracy Test
Tree	0.876786	0.833333
Knn	0.848214	0.833333
Svm	0.848214	0.833333
Logreg	0.846429	0.833333



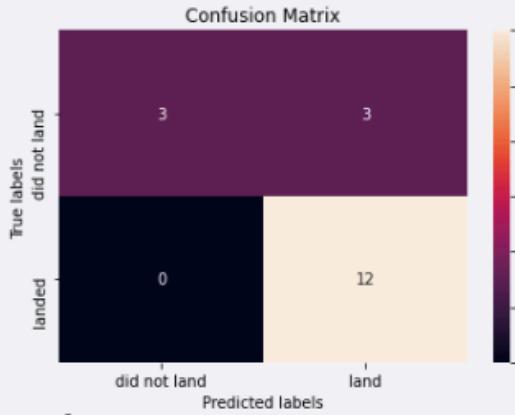
For accuracy test, all methods performed similarly.

## Decision tree parameters

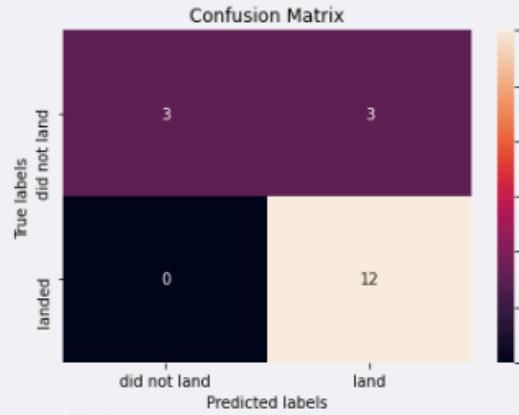
```
tuned hyperparameters :(best parameters)  {'criterion': 'entropy', 'max_depth': 12, 'max_features': 'sqrt', 'min_samples_leaf': 4, 'min_samples_split': 2, 'splitter': 'random'}
```

# Confusion Matrix

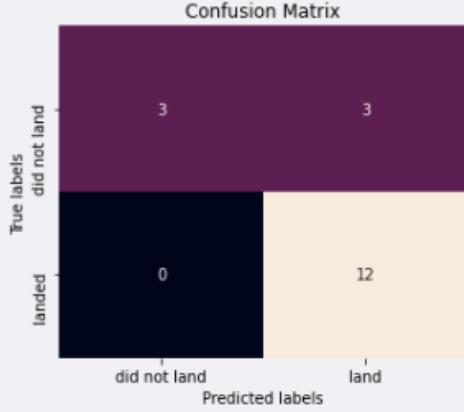
**Logistic regression**



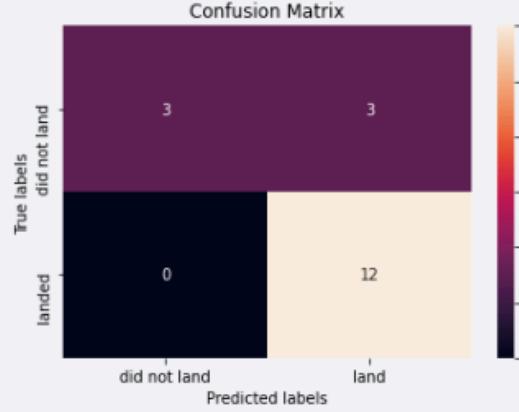
**Decision Tree**



**kNN**



**SVM**



The test accuracy was all equal, and so the confusion matrices are also identical. The primary challenge is false positives.

		Actual values	
		1	0
Predicted values	1	TP	FP
	0	FN	TN

# Conclusions

---

- The success of a mission can be explained by several factors such as the launch site, the orbit and especially the number of previous launches. Indeed, we can assume that there has been a gain in knowledge between launches that allowed to go from a launch failure to a success.
- The orbits with the best success rates are GEO, HEO, SSO, and ES-L1.
- Depending on the orbits, the payload mass can be a criterion to take into account for the success of a mission. Some orbits require a light or heavy payload mass. But generally low weighted payloads perform better than the heavy weighted payloads.
- With the current data, we cannot explain why some launch sites are better than others (KSCLC-39A is the best launch site). To get an answer to this problem, we could obtain atmospheric or other relevant data.
- For this dataset, we choose the Decision Tree Algorithm as the best model even if the test accuracy between all the models used is identical. We choose Decision Tree Algorithm because it has a better train accuracy.

Thank you!

