



IBM Developer  
SKILLS NETWORK

# Winning Space Race with Data Science

Robin Singh  
19/05/2025



# Outline

---

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

# Executive Summary

---

- Summary of methodologies
  - Data Collection through API
  - Data Collection with Web Scraping
  - Data Wrangling
  - Exploratory Data Analysis with SQL
  - Exploratory Data Analysis with Data Visualization
  - Interactive Visual Analytics with Folium
  - Machine Learning Prediction
- Summary of all results
  - Exploratory Data Analysis result
  - Interactive analytics in screenshots
  - Predictive Analytics result

# Introduction

---

- Project background and context

Space X advertises Falcon 9 rocket launches on its website with a cost of 62 million dollars; other providers cost upward of 165 million dollars each, much of the savings is because Space X can reuse the first stage. Therefore, if we can determine if the first stage will land, we can determine the cost of a launch. This goal of the project is to create a machine learning pipeline to predict if the first stage will land successfully.

- Questions to be answered

- How do variables such as payload mass, launch site, number of flights, and orbits affect the success of the first stage landing?
- Does the rate of successful landings increase over the years?
- What is the best algorithm that can be used for binary classification in this case ?



Section 1

# Methodology

# Methodology

---

## Executive Summary

- Data collection methodology:
  - Data was collected using SpaceX API and web scraping from Wikipedia.
- Perform data wrangling
  - One-hot encoding was applied to categorical features
- Perform exploratory data analysis (EDA) using visualization and SQL
- Perform interactive visual analytics using Folium and Plotly Dash
- Perform predictive analysis using classification models
  - How to build, tune, evaluate classification models

# Data Collection

---

- The data was collected using various methods
  - Data collection was done using get request to the SpaceX API.
  - Next, we decoded the response content as a Json using `.json()` function call and turn it into a pandas dataframe
  - We then cleaned the data, checked for missing values and fill in missing values where necessary.
  - In addition, we performed web scraping from Wikipedia for Falcon 9 launch records with BeautifulSoup.
  - The objective was to extract the launch records as HTML table, parse the table and convert it to a pandas dataframe for future analysis.

# Data Collection – SpaceX API

- We used the get request to the SpaceX API to collect data, clean the requested data and did some basic data wrangling and formatting.
- The link to the notebook is [https://github.com/RobinSingh410/IBM\\_Data\\_Science\\_Capstone.git](https://github.com/RobinSingh410/IBM_Data_Science_Capstone.git)

## 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:

```
[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

```
[10] response=requests.get(static_json_url)
```

```
[11] response.status_code
```

... 200

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

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

Using the dataframe `data` print the first 5 rows

```
[17] # Get the head of the dataframe
data.head()
```

	static_fire_date_utc	static_fire_date_unix	tbd	net	window	rocket	success	details	crew	ships	capsules	payloads	launchpad	auto_update
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]	5e9e4502f5090995de566f86		True



# 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/RobinSingh410/IBM\\_Data\\_Science\\_Capstone.git](https://github.com/RobinSingh410/IBM_Data_Science_Capstone.git)

To keep the lab tasks consistent, you will be asked to scrape the data from a snapshot of the [List of Falcon 9 and Falcon Heavy launches](#) Wikipage updated on [9th June 2021](#)

```
[4] static_url = "https://en.wikipedia.org/w/index.php?title=List_of_Falcon_9_and_Falcon_Heavy_launches&oldid=1027686922"
```

Next, request the HTML page from the above URL and get a `response` object

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.

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

Create a `BeautifulSoup` object from the HTML `response`

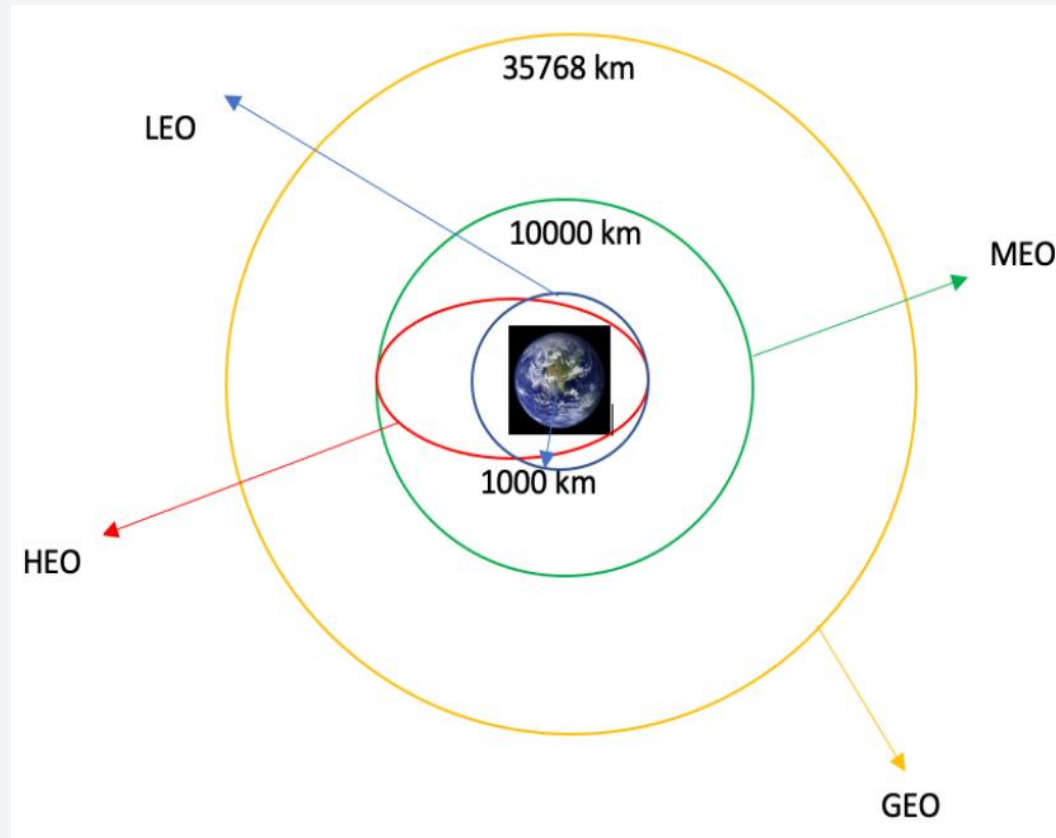
```
[6] # Use BeautifulSoup() to create a BeautifulSoup object from a response text content
soup = BeautifulSoup(response.content, 'html5lib')
```

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

```
[7] # Use soup.title attribut
soup.title

... <title>List of Falcon 9 and Falcon Heavy launches - Wikipedia</title>
```

# Data Wrangling



- We performed exploratory data analysis and determined the training labels.
- We calculated the number of launches at each site, and the number and occurrence of each orbits
- We created landing outcome label from outcome column and exported the results to csv.
- The link to the notebook is [https://github.com/RobinSingh410/IBM\\_Data\\_Science\\_Capstone.git](https://github.com/RobinSingh410/IBM_Data_Science_Capstone.git)

# 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, Launch success yearly trend
- For visualization we have used scatter plot , bar plot , line plot .
- Github link :  
[https://github.com/RobinSingh410/IBM\\_Data\\_Science\\_Capstone.git](https://github.com/RobinSingh410/IBM_Data_Science_Capstone.git)



# EDA with SQL

---

- We applied EDA with SQL to get insight from the data. We wrote queries to find out for instance:
  - The names of unique launch sites in the space mission.
  - The total payload mass carried by boosters launched by NASA (CRS)
  - The average payload mass carried by booster version F9 v1.1
  - The total number of successful and failure mission outcomes
  - The failed landing outcomes in drone ship, their booster version and launch site names.
- Github link :  
[https://github.com/RobinSingh410/IBM\\_Data\\_Science\\_Capstone.git](https://github.com/RobinSingh410/IBM_Data_Science_Capstone.git)

# Build an Interactive Map with Folium

---

- We plotted all launch sites on the Folium map and incorporated various map elements like markers, circles, and lines to visually represent the success or failure of launches at each location.
- The launch outcomes were categorized into two classes: 0 for failure and 1 for success.
- By utilizing color-coded marker clusters, we analyzed which launch sites demonstrated a higher success rate.
- We measured the distances between launch sites and their surrounding infrastructure. Additionally, we explored key questions such as:
  - Are launch sites located near railways, highways, or coastlines?
  - Do launch sites maintain a specific distance from cities?



# Build a Dashboard with Plotly Dash

---

- We developed an interactive dashboard using Plotly Dash.
- A pie chart was created to visualize the total number of launches at each site.
- A scatter plot was generated to examine the relationship between launch outcomes and payload mass (Kg) across different booster versions.
- Github link :  
[https://github.com/RobinSingh410/IBM\\_Data\\_Science\\_Capstone.git](https://github.com/RobinSingh410/IBM_Data_Science_Capstone.git)

# Predictive Analysis (Classification)

---

- We loaded and processed the data using NumPy and pandas, followed by data transformation and splitting into training and testing sets.
- Multiple machine learning models were built, and hyperparameters were optimized using GridSearchCV.
- Accuracy was used as the primary evaluation metric, and we enhanced the model through feature engineering and algorithm tuning.
- The best-performing classification model was identified based on its performance.
- Github link :  
[https://github.com/RobinSingh410/IBM\\_Data\\_Science\\_Capstone.git](https://github.com/RobinSingh410/IBM_Data_Science_Capstone.git)

# 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 solid blue area on the left side, which transitions into a dynamic pattern of diagonal streaks in shades of blue, red, and cyan on the right. These streaks are layered over a faint, grid-like pattern, creating a sense of depth and movement, reminiscent of a digital or data visualization theme.

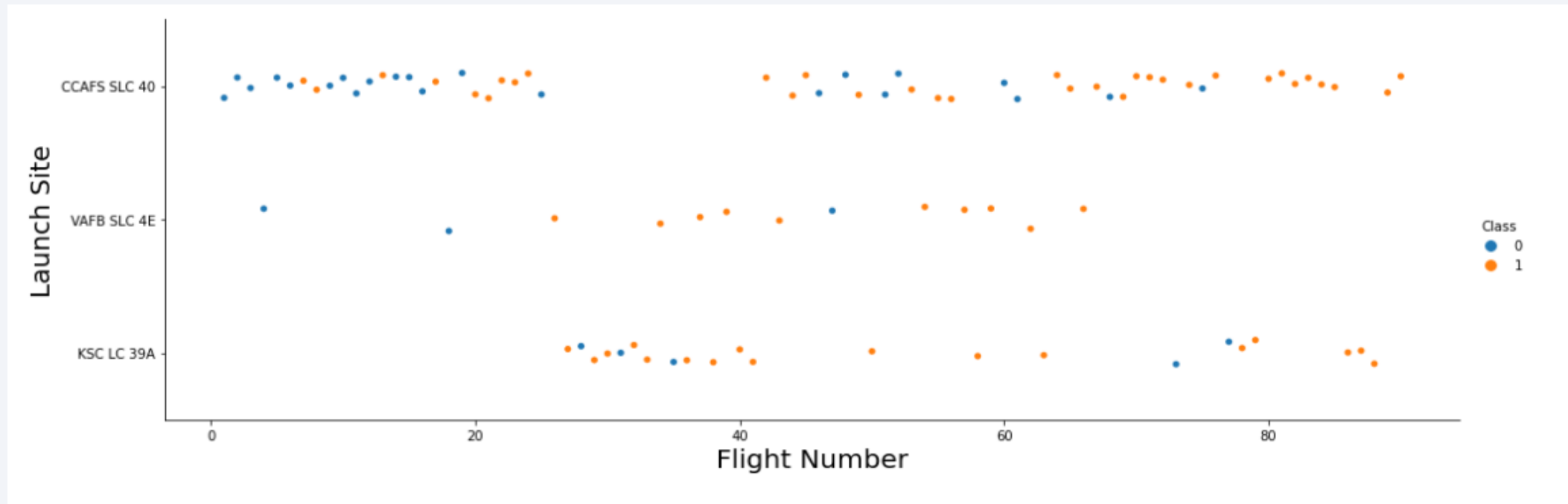
Section 2

# Insights drawn from EDA



# Flight Number vs. Launch Site

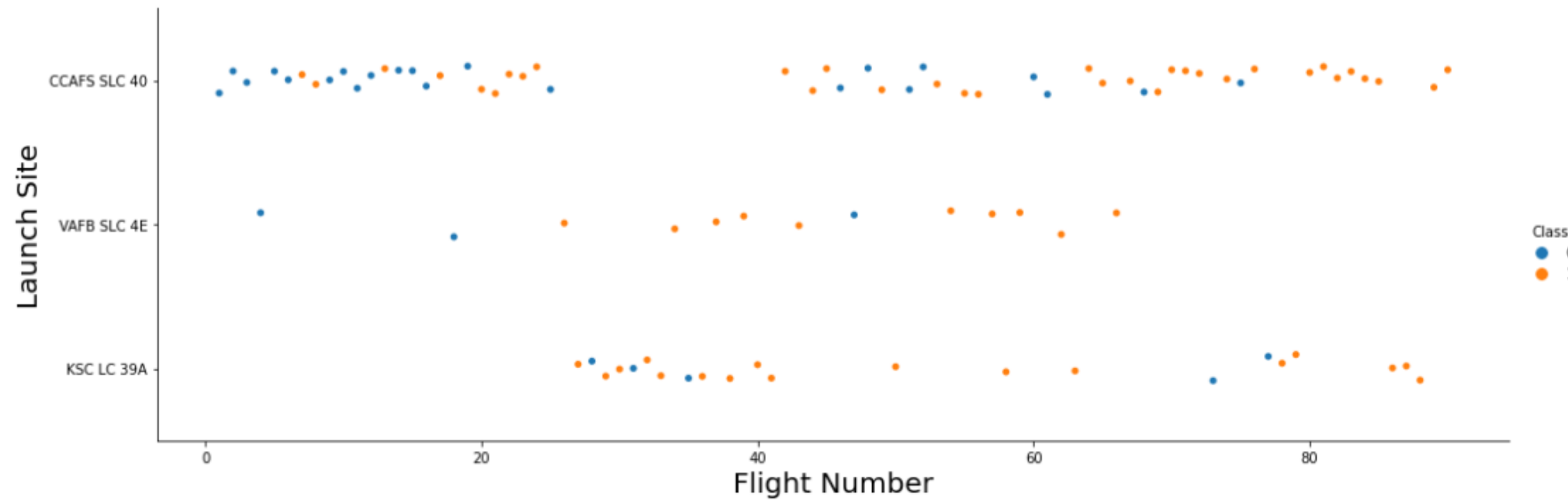
- The earliest flights all failed while the latest flights all succeeded.
- It can be assumed that each new launch has a higher rate of success.





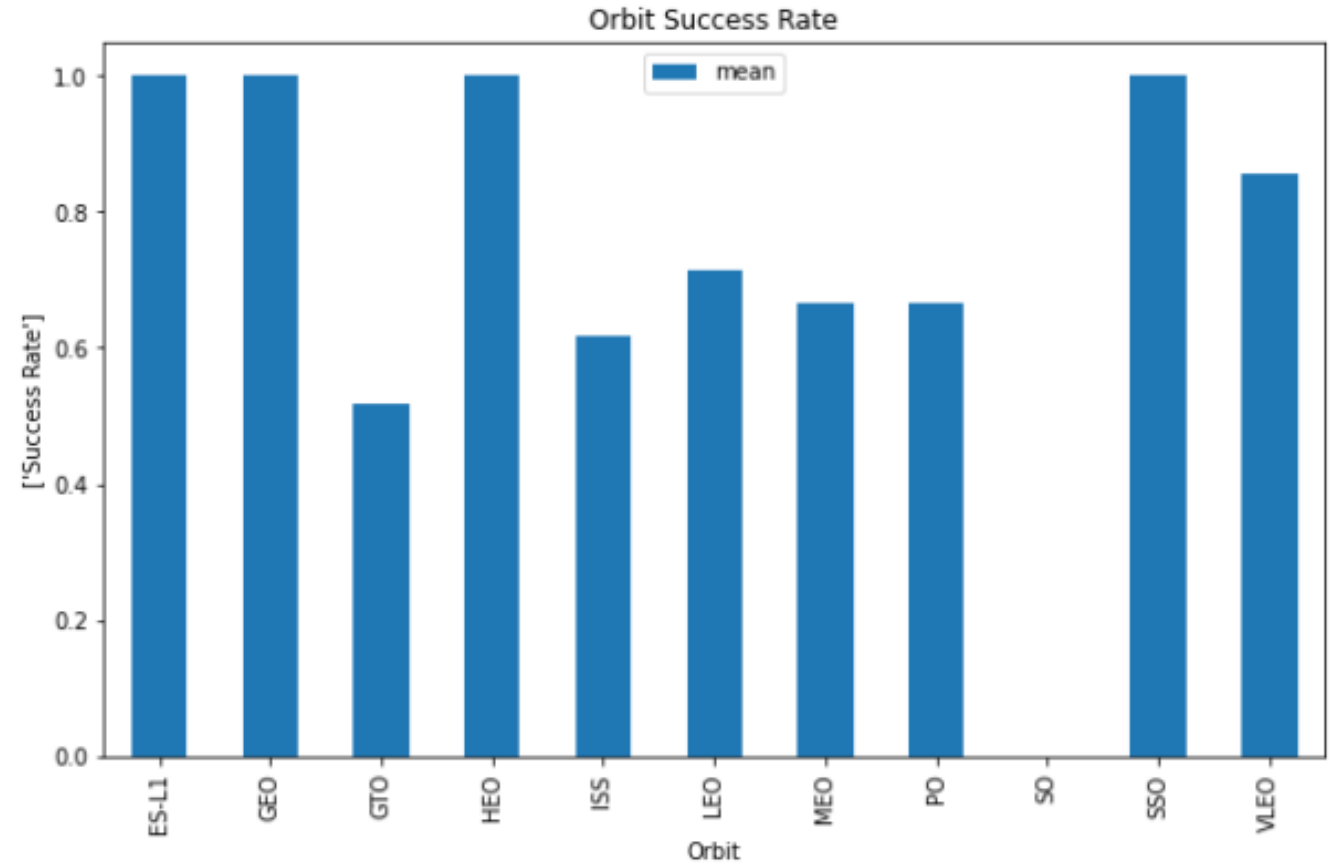
# Payload vs. Launch Site

- For every launch site the higher the payload mass, the higher the success rate



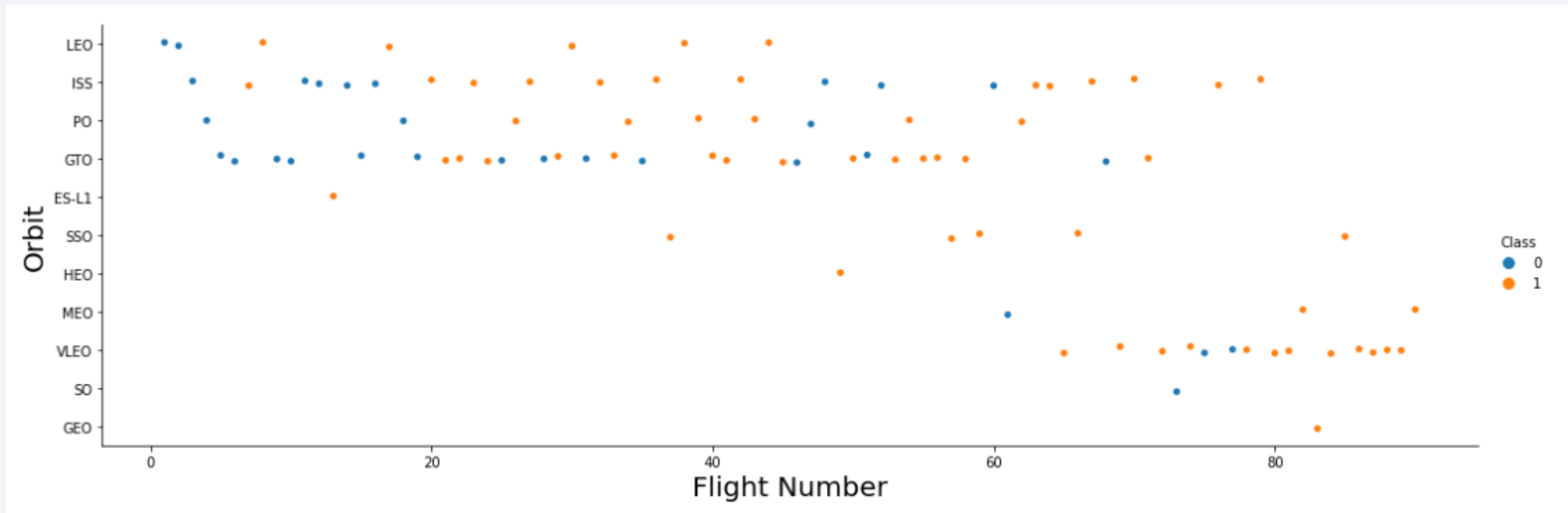
# Success Rate vs. Orbit Type

- Orbits with 100% success rate:
  - ES-L1, GEO, HEO, SSO
- Orbits with 0% success rate:
  - SO



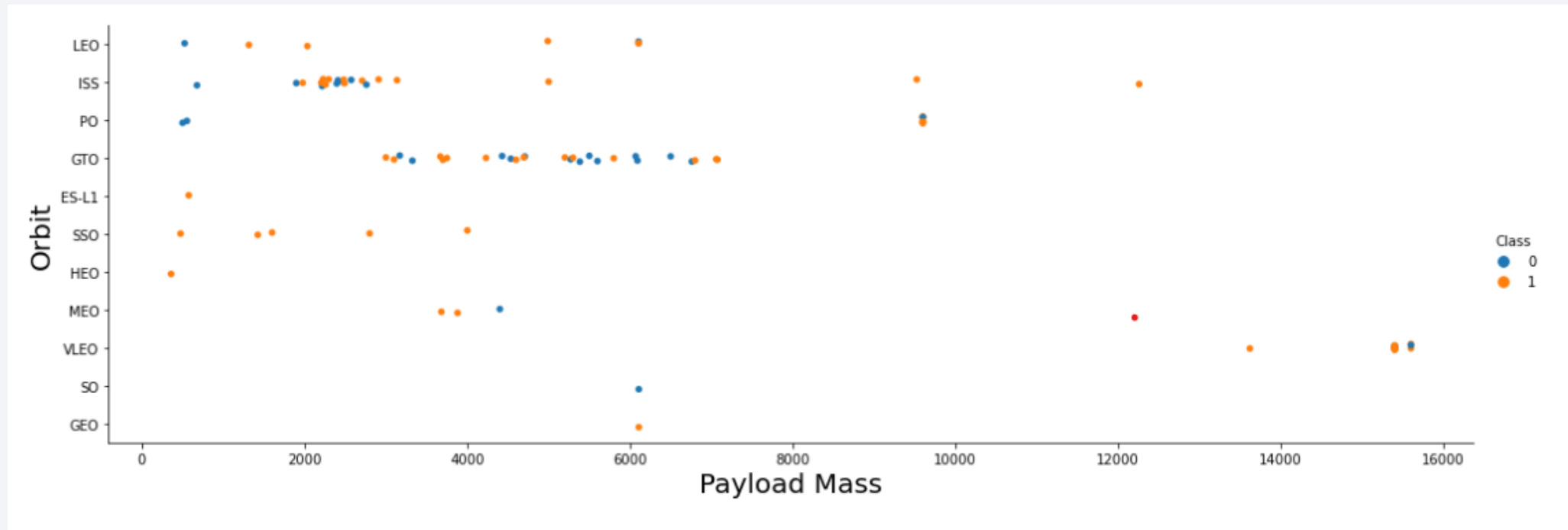
# Flight Number vs. Orbit Type

- 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.



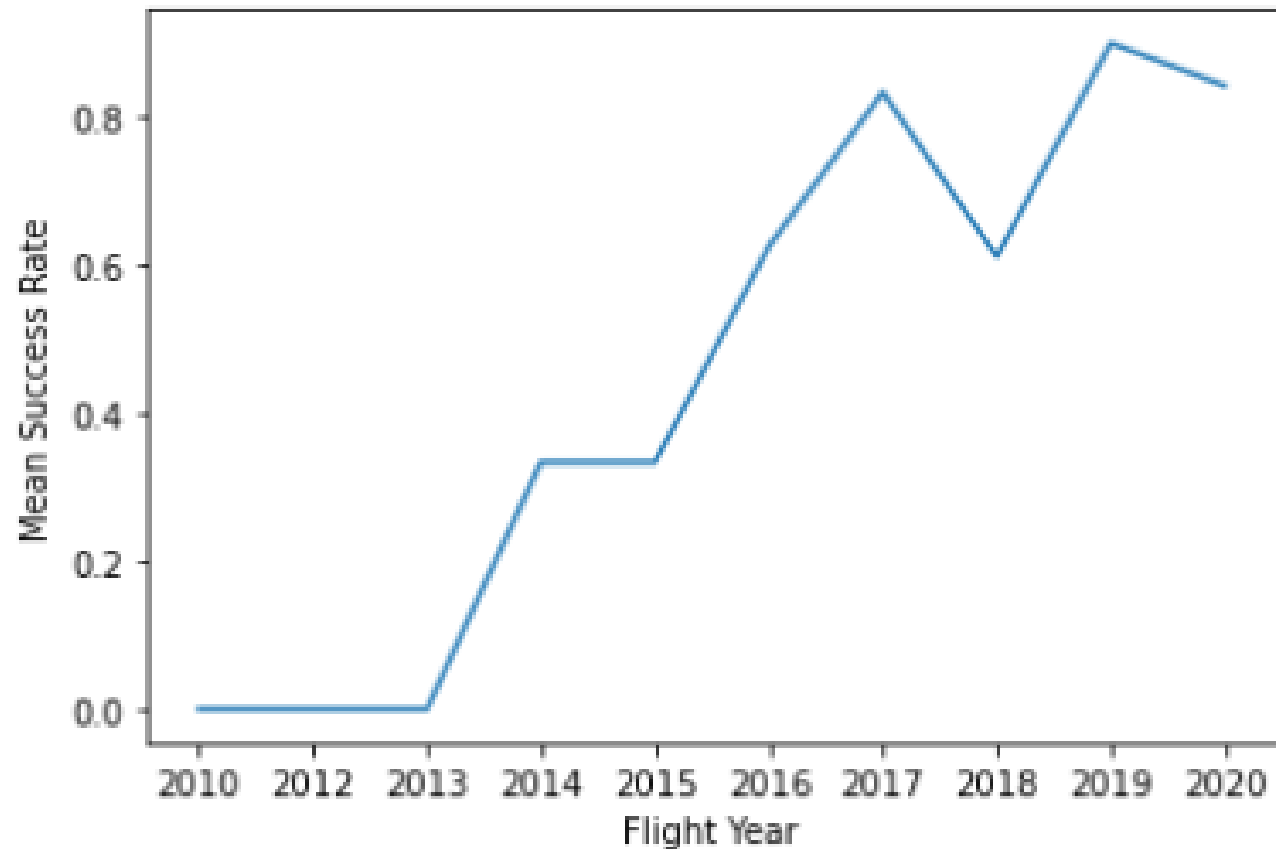
# Payload vs. Orbit Type

- We can observe that with heavy payloads, the successful landing are more for PO, LEO and ISS orbits.



# Launch Success Yearly Trend

- The success rate since 2013 kept increasing till 2020.





# All Launch Site Names

- Find the names of the unique launch sites

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

```
* ibm_db_sa://kcq64325:***@dashdb-txn-sbox-yp-dal09-04.  
Done.
```

launch\_site

CCAFS LC-40

CCAFS SLC-40

KSC LC-39A

VAFB SLC-4E

# Launch Site Names Begin with 'CCA'

---

```
%sql SELECT * FROM SPACEXTABLE WHERE "Launch_Site" LIKE 'CCA%' LIMIT 5;
```

\* [sqlite:///my\\_data1.db](#)  
Done.

Date	Time (UTC)	Booster_Version	Launch_Site	Payload	PAYLOAD_MASS_KG_	Orbit	Customer	Mission_Outcome	Landing_Outcome
2010-06-04	18:45:00	F9 v1.0 B0003	CCAFS LC-40	Dragon Spacecraft Qualification Unit	0	LEO	SpaceX	Success	Failure (parachute)
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 (parachute)
2012-05-22	7:44:00	F9 v1.0 B0005	CCAFS LC-40	Dragon demo flight C2	525	LEO (ISS)	NASA (COTS)	Success	No attempt
2012-10-08	0:35:00	F9 v1.0 B0006	CCAFS LC-40	SpaceX CRS-1	500	LEO (ISS)	NASA (CRS)	Success	No attempt
2013-03-01	15:10:00	F9 v1.0 B0007	CCAFS LC-40	SpaceX CRS-2	677	LEO (ISS)	NASA (CRS)	Success	No attempt

# Total Payload Mass

---

- Calculated the total payload carried by boosters from NASA

```
%sql SELECT SUM("PAYLOAD_MASS__KG_") AS Total_Payload_Mass FROM SPACEXTABLE WHERE "Customer" LIKE 'NASA (CRS)%';
```

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

```
Done.
```

Total_Payload_Mass
--------------------

48213
-------

# Average Payload Mass by F9 v1.1

- Calculated the average payload mass carried by booster version F9 v1.1

```
%sql SELECT AVG("PAYLOAD_MASS__KG_") AS Average_Payload_Mass FROM SPACE_TABLE WHERE "Booster_Version" = 'F9 v1.1';
```

\* [sqlite:///my\\_data1.db](#)

Done.

Average_Payload_Mass
----------------------

2928.4
--------

# First Successful Ground Landing Date

- The dates of the first successful landing outcome on ground pad

```
%sql SELECT MIN("Date") AS First_Successful_Landing FROM SPACEXTABLE WHERE "Landing_Outcome" = 'Success (ground pad)';
```

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

```
Done.
```

First_Successful_Landing
--------------------------

2015-12-22
------------



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

- Listing the names of boosters which have successfully landed on drone ship and had payload mass greater than 4000 but less than 6000

```
%sql SELECT DISTINCT "Booster_Version" FROM SPACEXTABLE WHERE "Landing_Outcome" = 'Success (drone ship)' AND "PAYLOAD_MASS_KG_" BETWEEN 4000 AND 6000;
```

\* [sqlite:///my\\_data1.db](#)

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

- Calculated the total number of successful and failure mission outcomes

```
%sql SELECT "Mission_Outcome", COUNT(*) AS Total_Count FROM SPACEXTABLE GROUP BY "Mission_Outcome";
```

\* [sqlite:///my\\_data1.db](#)

Done.

Mission_Outcome	Total_Count
Failure (in flight)	1
Success	98
Success	1
Success (payload status unclear)	1

# Boosters Carried Maximum Payload

- Listed the names of the booster which have carried the maximum payload mass

```
%sql SELECT "Booster_Version" FROM SPACEXTABLE WHERE "PAYLOAD_MASS__KG_" = (SELECT MAX("PAYLOAD_MASS__KG_") FROM SPACEXTABLE);
```

\* [sqlite:///my\\_data1.db](#)

Done.

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

# 2015 Launch Records

---

- Listed the failed landing outcomes in drone ship, their booster versions, and launch site names for in year 2015

```
%sql SELECT substr("Date", 6, 2) AS Month, "Landing_Outcome", "Booster_Version", "Launch_Site" FROM SPACEXTABLE WHERE "Landing_Outcome"  
LIKE '%Failure%' AND "Landing_Outcome" LIKE '%drone ship%' AND substr("Date", 0, 5) = '2015';|
```

\* [sqlite:///my\\_data1.db](#)

Done.

Month	Landing_Outcome	Booster_Version	Launch_Site
01	Failure (drone ship)	F9 v1.1 B1012	CCAFS LC-40
04	Failure (drone ship)	F9 v1.1 B1015	CCAFS LC-40

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

- Ranked 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

```
%sql SELECT "Landing_Outcome", COUNT("Landing_Outcome") AS Outcome_Count FROM SPACEXTABLE  
WHERE "Date" BETWEEN '2010-06-04' AND '2017-03-20' GROUP BY "Landing_Outcome" ORDER BY Outcome_Count DESC;
```

\* [sqlite:///my\\_data1.db](#)

Done.

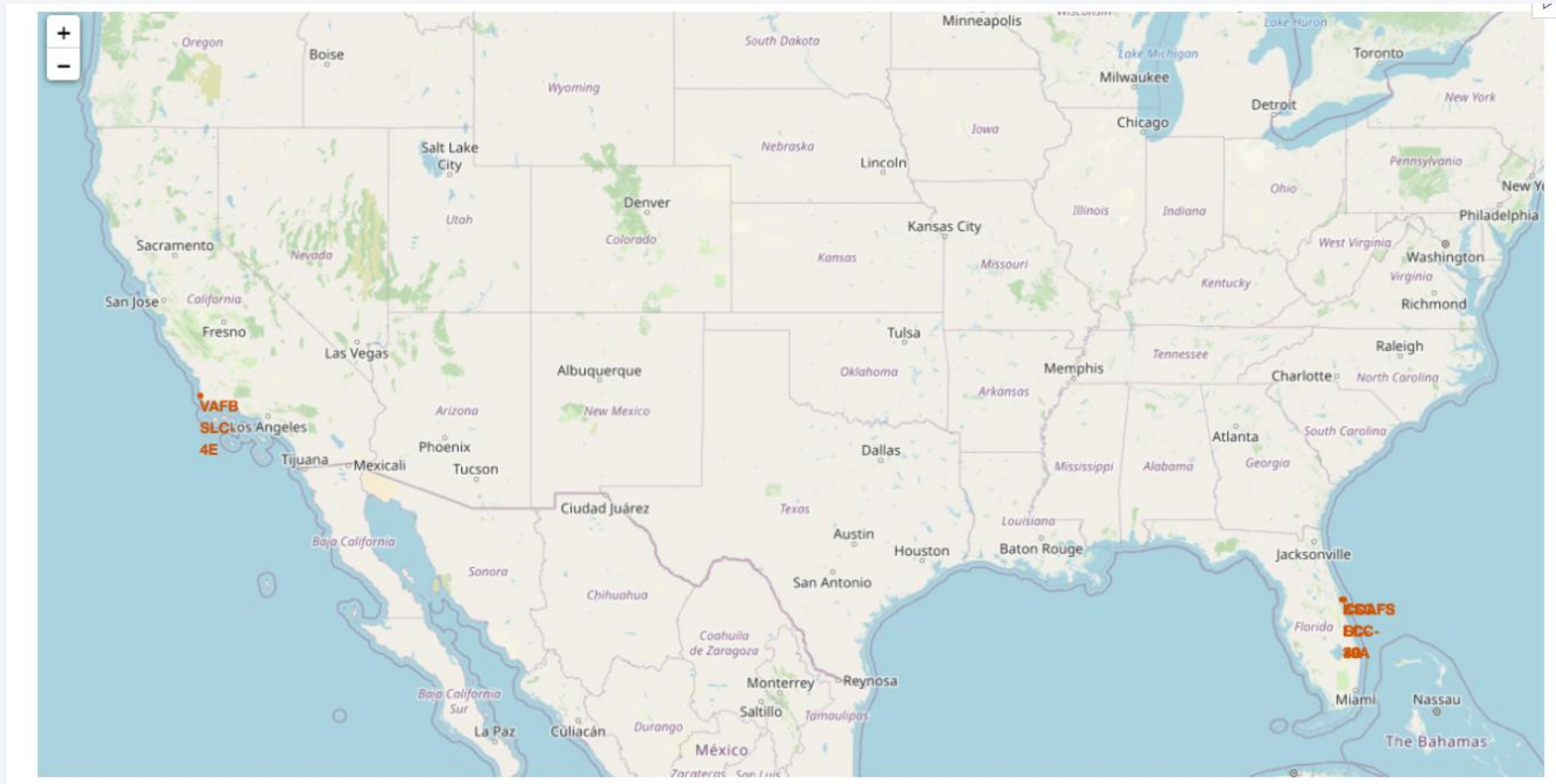
Landing_Outcome	Outcome_Count
No attempt	10
Success (drone ship)	5
Failure (drone ship)	5
Success (ground pad)	3
Controlled (ocean)	3
Uncontrolled (ocean)	2
Failure (parachute)	2
Precluded (drone ship)	1

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 upper half and a satellite photograph of the Earth's surface at night. The Earth's surface shows a dense network of city lights, primarily concentrated in the lower right quadrant, with a clear horizon line separating the dark blue sky from the illuminated land and sea.

Section 4

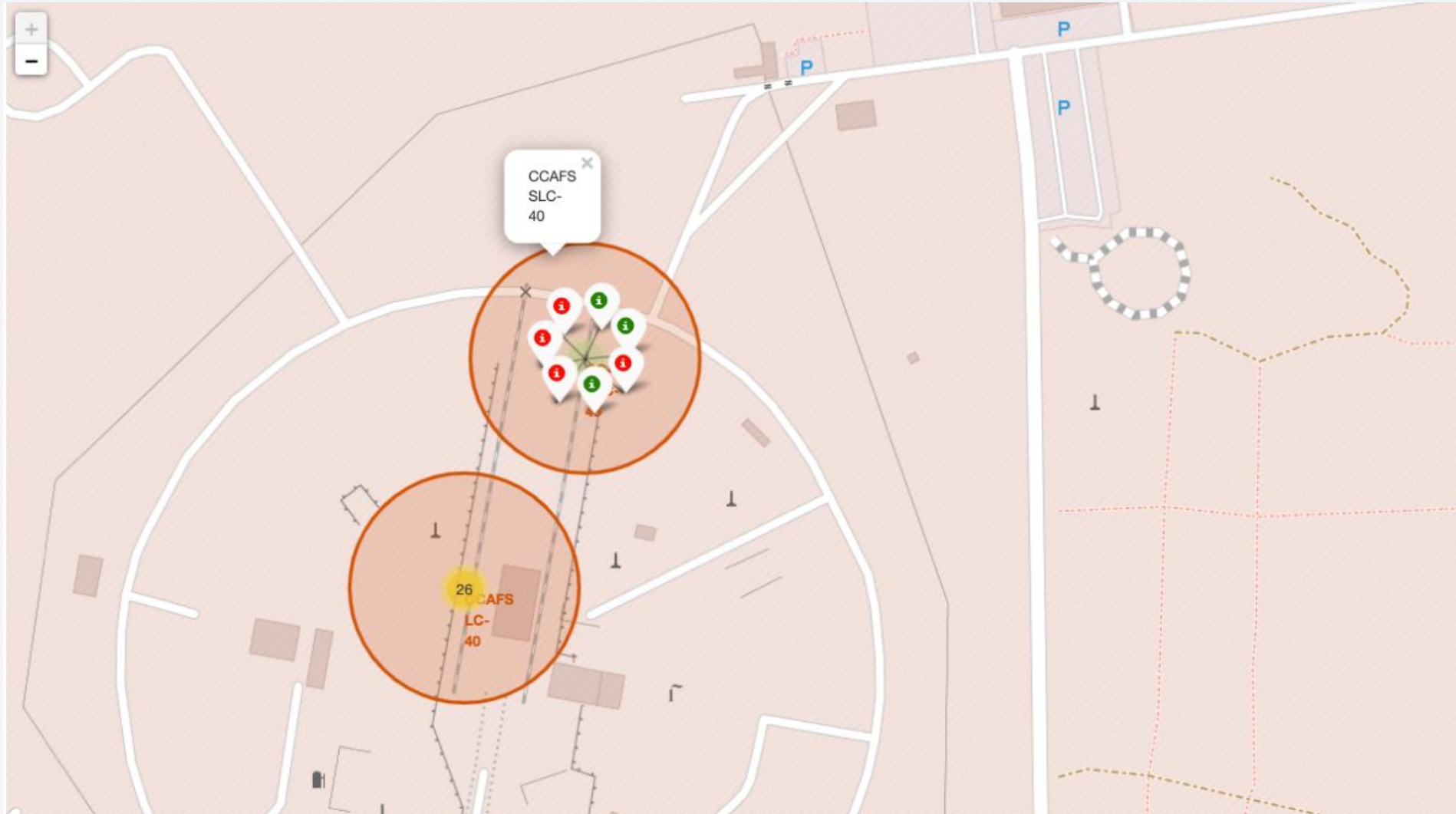
# Launch Sites Proximities Analysis

# All launch sites





# Markers showing launch sites with color labels





# Launch Site distance to landmarks





Section 6

# Predictive Analysis (Classification)

# Classification Accuracy

- Model 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_)
```

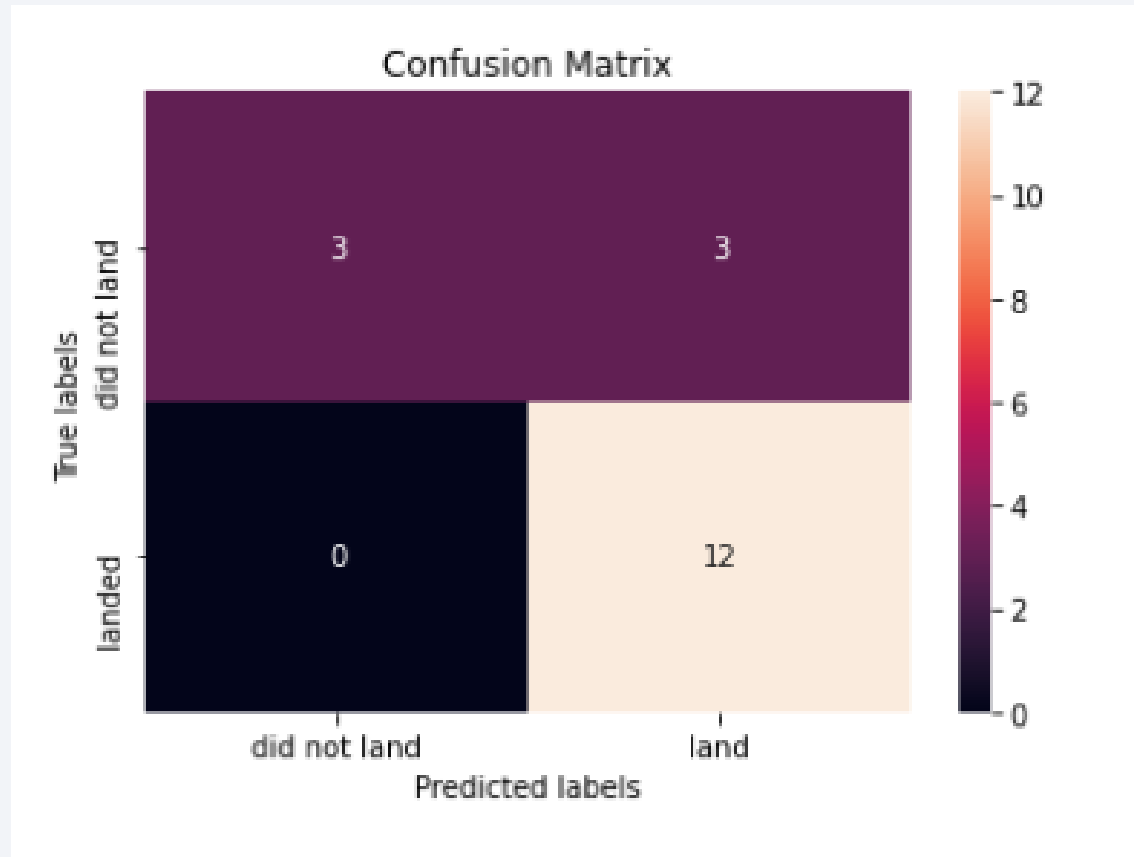
Best Algorithm is Tree with a score of 0.875

Best Params is : {'criterion': 'entropy', 'max\_depth': 2, 'max\_features': 'auto', 'min\_samples\_leaf': 1, 'min\_samples\_split': 10, 'splitter': 'best'}

# Confusion Matrix

---

- The confusion matrix for the decision tree



# Conclusions

---

We can conclude that:

- The Decision Tree model emerged as the best-performing algorithm for this dataset.
- Launches with lower payload mass tend to have higher success rates compared to those with heavier payloads.
- Most launch sites are strategically located near the Equator and in close proximity to the coast.
- The success rate of launches has shown a steady increase, particularly from 2013 to 2020.
- KSC LC-39A recorded the highest number of successful launches among all sites.
- Orbits such as ES-L1, GEO, HEO, and SSO demonstrated the highest success rates, with some achieving a 100% success rate.



Thank you!

