

Winning Space Race with Data Science

Vagisha Ojha
30th 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

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 information can be used if an alternate company wants to bid against space X 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.
- 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, decoded the response content as a Json using `.json()` function call and turn it into a pandas dataframe using `.json_normalize()`.
 - Cleaned the data, checked for missing values and fill in missing values where necessary.
 - In addition, 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

- Utilized 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/vagishao/IBM-Data-Science-Professional-Certificate/blob/main/Applied%20Data%20Science%20Capstone/module1/jupyter-labs-spacex-data-collection-api.ipynb>

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

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

```
[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()`

```
[33] import json  
# Use json_normalize method to convert the json result into a dataframe  
data = requests.get(static_json_url).json()  
df= pd.json_normalize(data)
```

Finally we will remove the Falcon 1 launches keeping only the Falcon 9 launches. Filter the data dataframe using the `BoosterVersion` column to only keep the Falcon 9 launches. Save the filtered data to a new dataframe called `data_falcon9`.

```
[54] # Hint data['BoosterVersion']!='Falcon 1'  
data_falcon9= df_launch_dict[df_launch_dict['BoosterVersion']!='Falcon 1']  
data_falcon9.head()
```

Python

...	FlightNumber	Date	BoosterVersion	PayloadMass	Orbit	LaunchSite	Outcome	Flights	GridFins	Reused	Legs	Landi
4	6	2010-06-04	Falcon 9	NaN	LEO	CCSFS SLC 40	None None	1	False	False	False	
5	8	2012-05-22	Falcon 9	525.0	LEO	CCSFS SLC 40	None None	1	False	False	False	

Data Collection - Scraping

- Applied web scrapping to webscrap Falcon 9 launch records with BeautifulSoup
- Parsed the table and converted it into a pandas dataframe.
- The link to the notebook is <https://github.com/vagishao/IBM-Data-Science-Professional-Certificate/blob/main/Applied%20Data%20Science%20Capstone/module1/jupyter-labs-webscraping.ipynb>

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 Wikipedia updated on 9th June 2021

```
[5] static_url = "https://en.wikipedia.org/w/index.php?title=List_of_Falcon_9_and_Falcon_Heavy_launches&oldid=1027686" Python
```

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.

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

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

```
[7] # Use BeautifulSoup() to create a BeautifulSoup object from a response text content  
soup = BeautifulSoup(data, "html.parser") Python
```

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

```
[8] # Use soup.title attribute  
soup.title Python
```

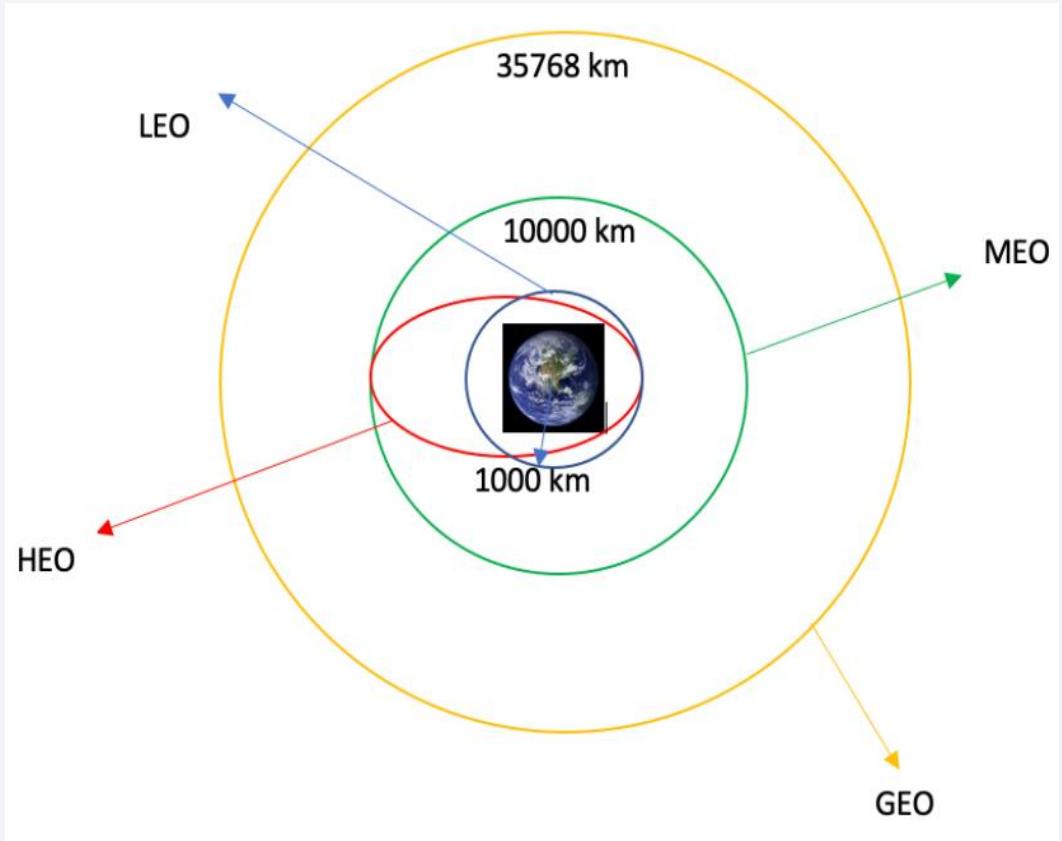
```
[9] ... <title>List of Falcon 9 and Falcon Heavy launches - Wikipedia</title>
```

After you have fill in the parsed launch record values into `launch_dict`, you can create a dataframe from it.

```
[10] df= pd.DataFrame({ key:pd.Series(value) for key, value in launch_dict.items() })  
df.head() Python
```

Flight No.	Launch site	Payload	Payload mass	Orbit	Customer	Launch outcome	Version Booster	Booster landing	Date	Time
0	1	CCAFS	Dragon Spacecraft Qualification Unit	0	LEO	SpaceX	Success\n	F9 v1.1	Failure	4 June 2010 18:45

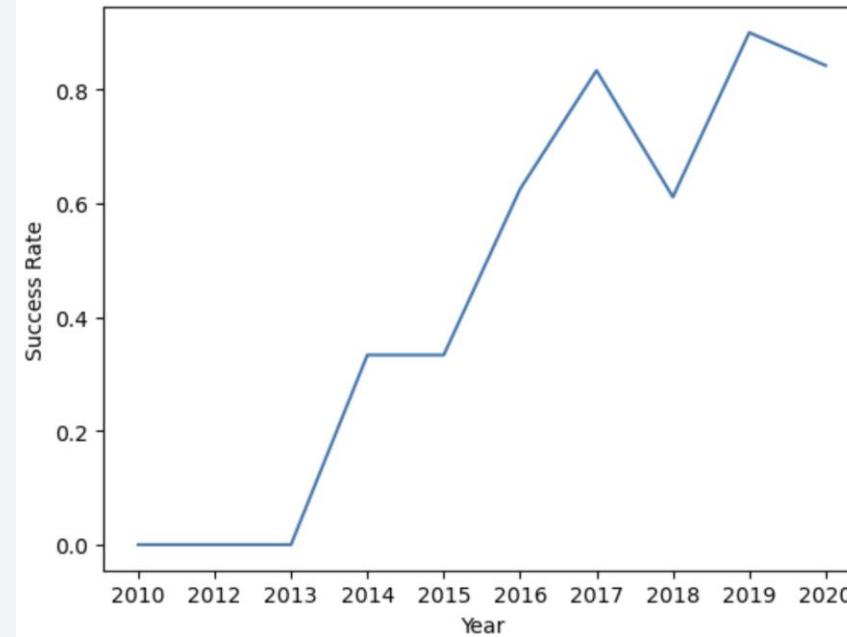
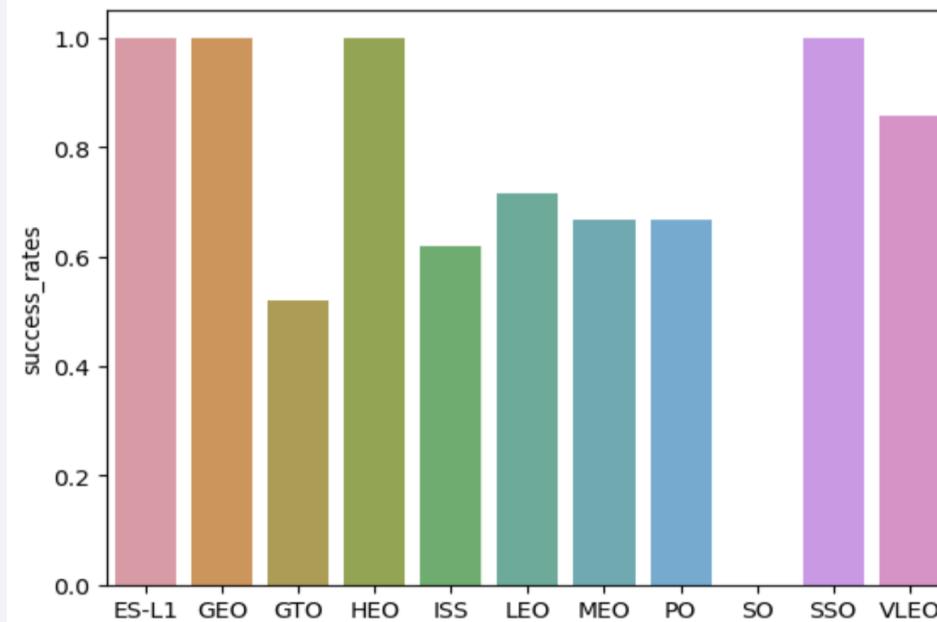
Data Wrangling



- Performed exploratory data analysis and determined the training labels.
- Calculated number of launches at each site, and the number and occurrence of each orbits
- Created landing outcome label from outcome column and exported the results to csv.
- The link to the notebook is
<https://github.com/vagishao/IBM-Data-Science-Professional-Certificate/blob/main/Applied%20Data%20Science%20Capstone/module1/labs-jupyter-spacex-Data%20wrangling.ipynb>

EDA with Data Visualization

- 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/vagishao/IBM-Data-Science-Professional-Certificate/blob/main/Applied%20Data%20Science%20Capstone/module2/jupyter-labs-eda-dataviz.ipynb>

EDA with SQL

- Loaded SpaceX dataset into a database.
- Applied EDA with SQL to get insight from the data. Wrote queries to determine:
 - 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.
- The link to the notebook is <https://github.com/vagishao/IBM-Data-Science-Professional-Certificate/blob/main/Applied%20Data%20Science%20Capstone/module2/jupyter-labs-eda-dataviz.ipynb>

Build an Interactive Map with Folium

- 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.
- 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, identified which launch sites have relatively high success rate.
- Calculated the distances between a launch site to its proximities. Answered some question 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

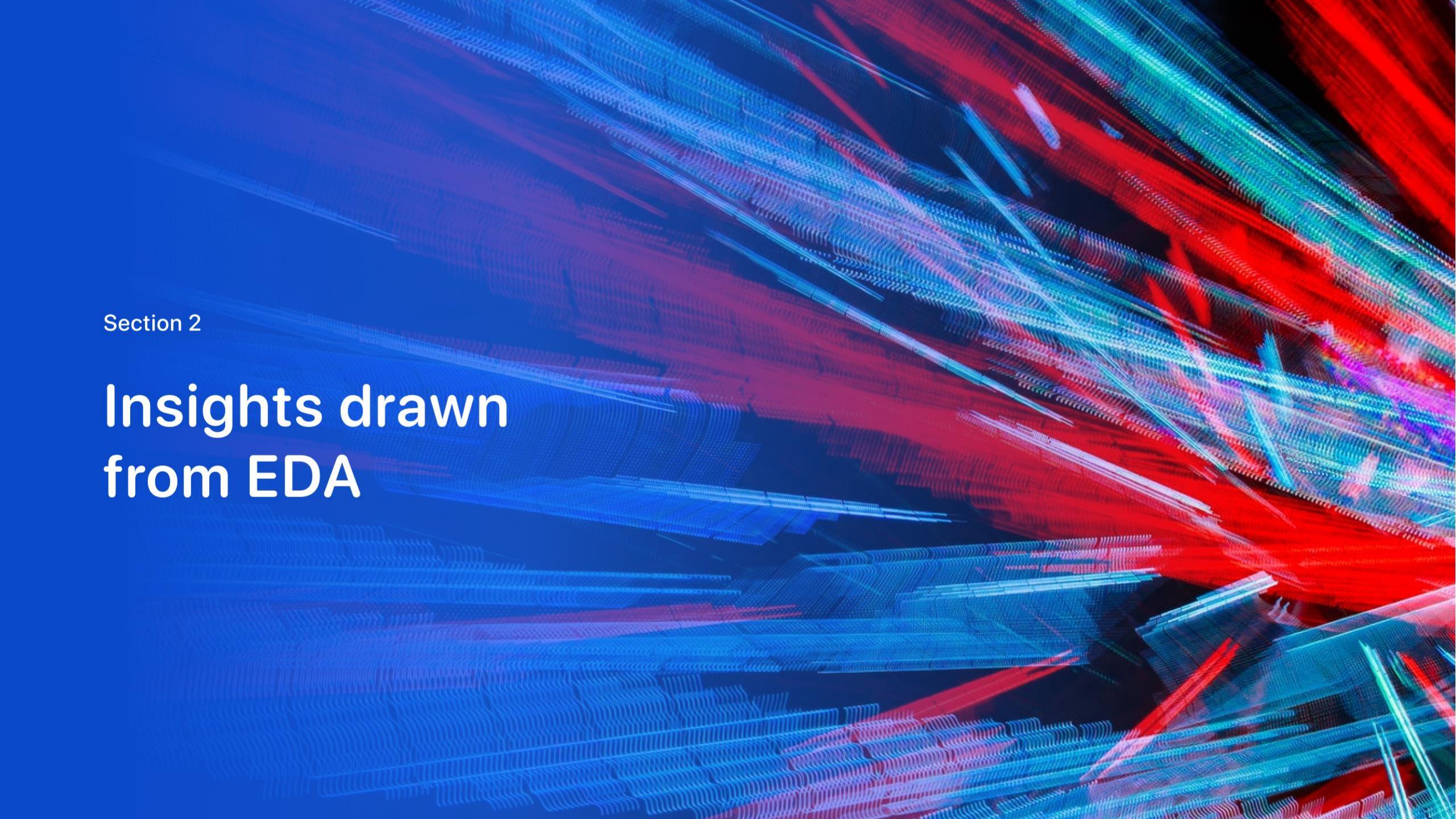
- Built an interactive dashboard with Plotly dash
- Plotted pie charts showing the total launches by a certain sites
- Plotted scatter graph showing the relationship with Outcome and Payload Mass (Kg) for the different booster version.
- The link to the notebook is https://github.com/vagishao/IBM-Data-Science-Professional-Certificate/blob/main/Applied%20Data%20Science%20Capstone/module3/spacex_dash_app.py

Predictive Analysis (Classification)

- We loaded the data using numpy and pandas, transformed the data, split our data into training and testing.
- We built different machine learning models and tune different hyperparameters using GridSearchCV.
- We used accuracy as the metric for our model, improved the model using feature engineering and algorithm tuning.
- We found the best performing classification model.
- The link to the notebook is https://github.com/vagishao/IBM-Data-Science-Professional-Certificate/blob/main/Applied%20Data%20Science%20Capstone/module4/SpaceX_Machine_Learning_Prediction_Part_5.jupyterlite.ipynb

Results

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

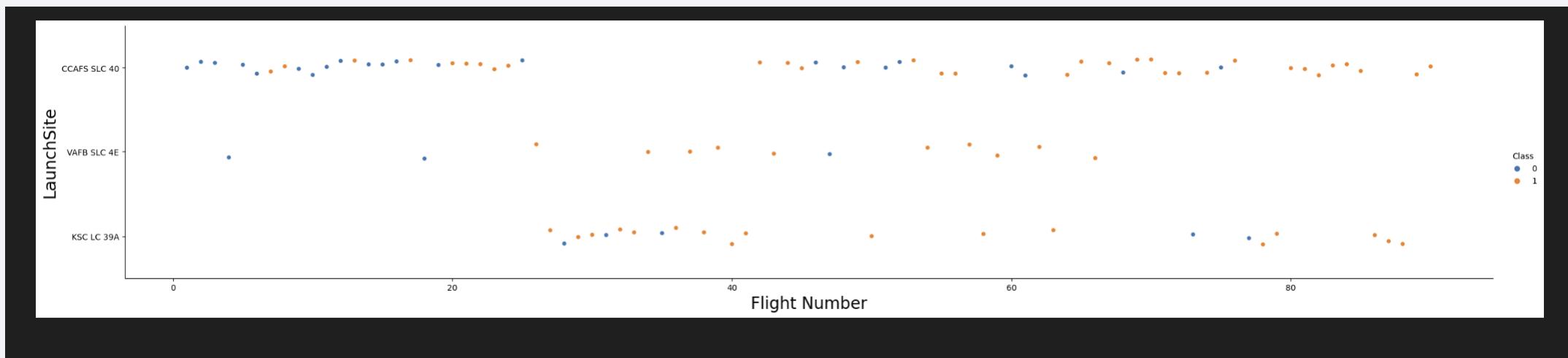
The background of the slide features a complex, abstract digital visualization. It consists of numerous thin, glowing lines that create a sense of depth and motion. The lines are primarily blue and red, with some green and purple highlights. They form a grid-like structure that curves and twists across 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

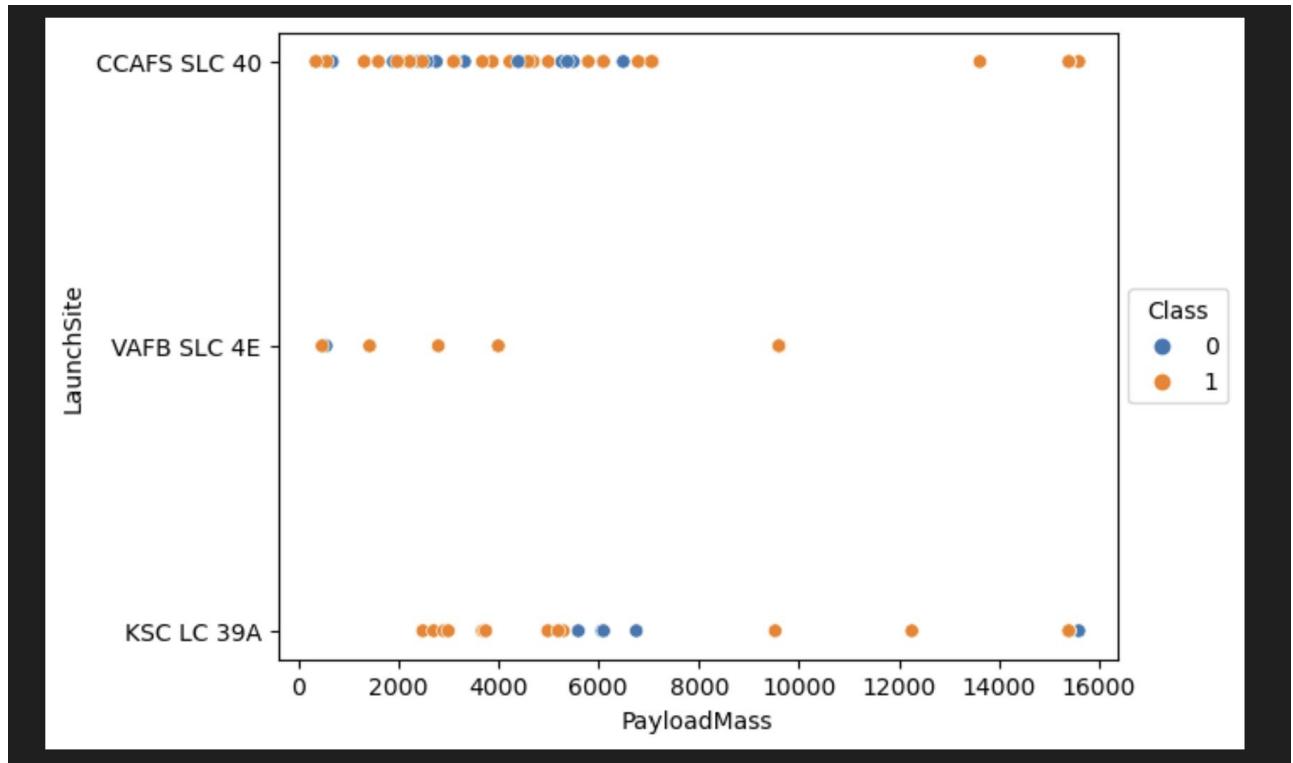
- From the plot, we found that the larger the flight amount at a launch site, the greater the success rate at a launch site.



Payload vs. Launch Site

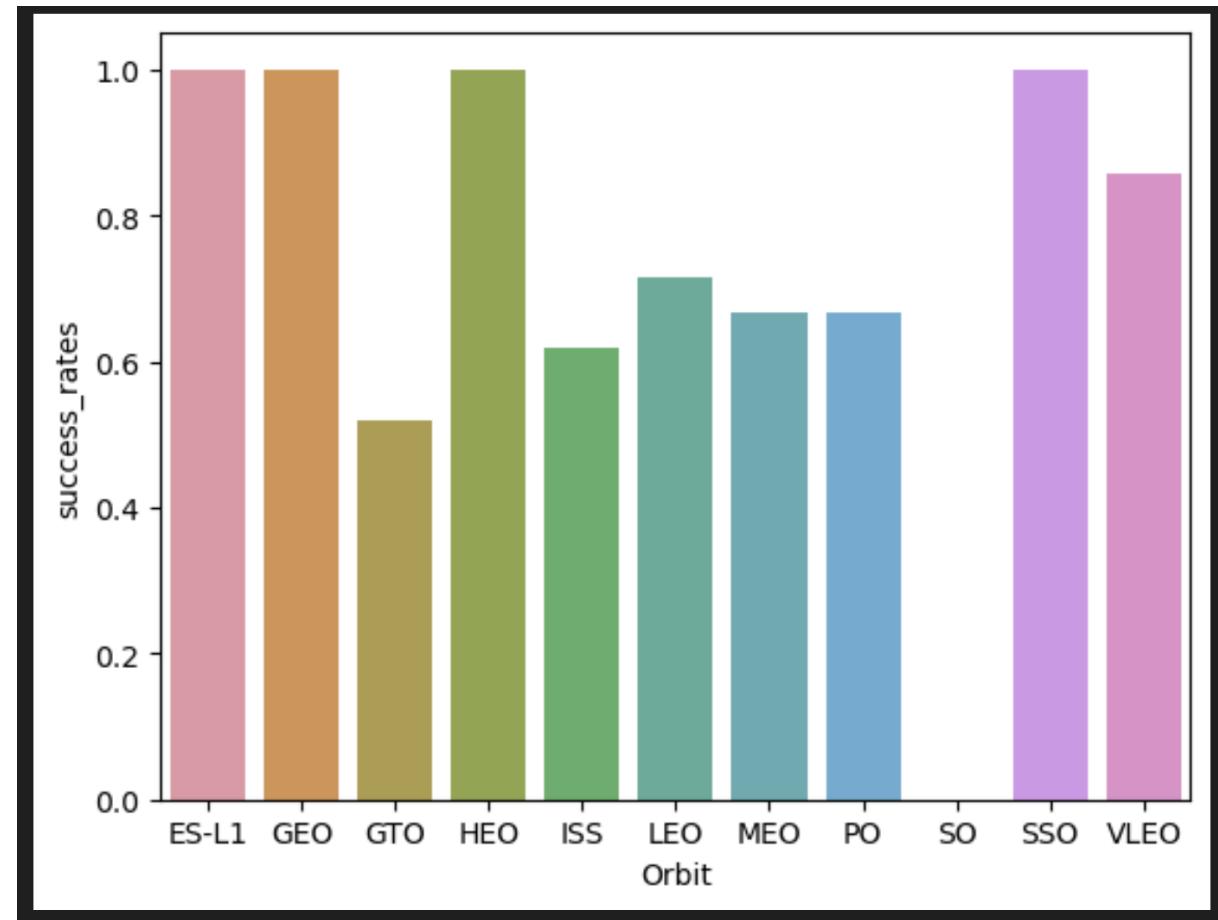


The greater the payload mass for launch site CCAFS SLC 40 the higher the success rate for the rocket.



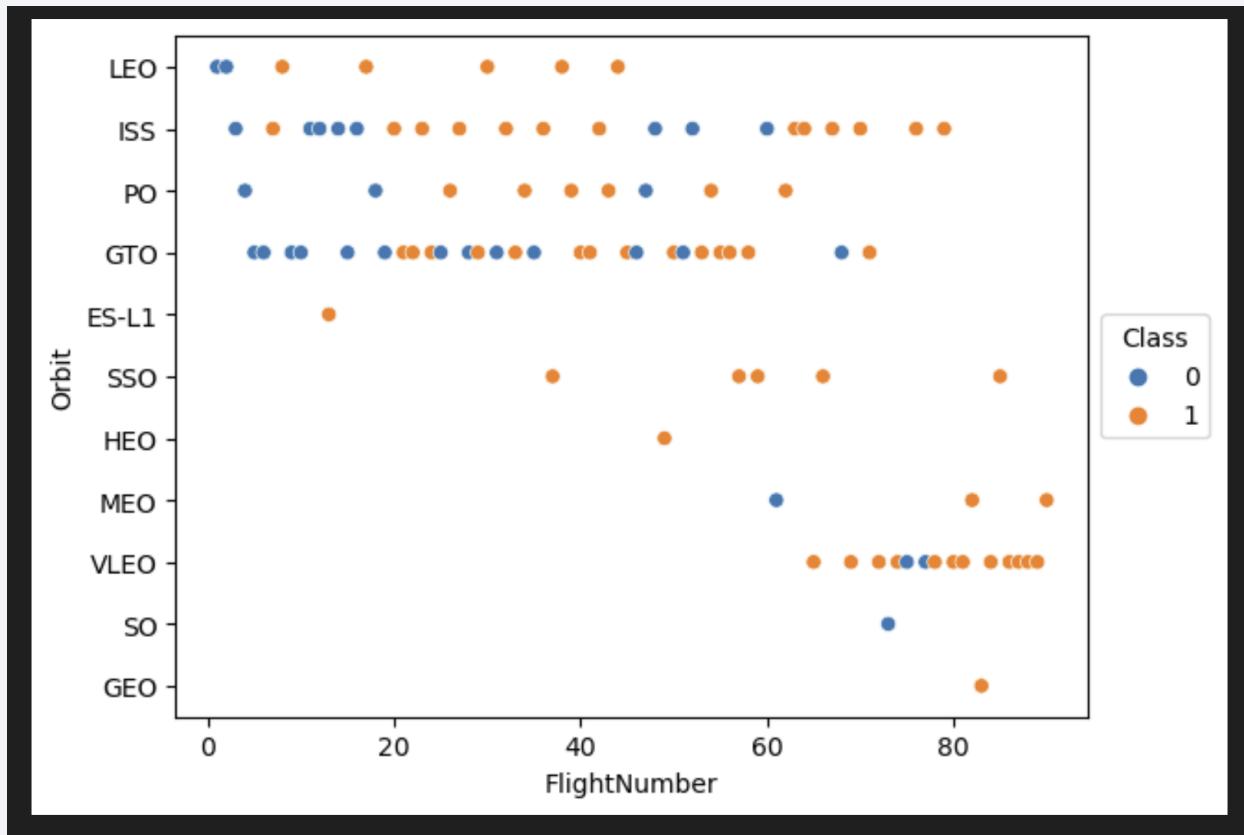
Success Rate vs. Orbit Type

- From the plot, we can see that ES-L1, GEO, HEO, SSO, VLEO had the most success rate.



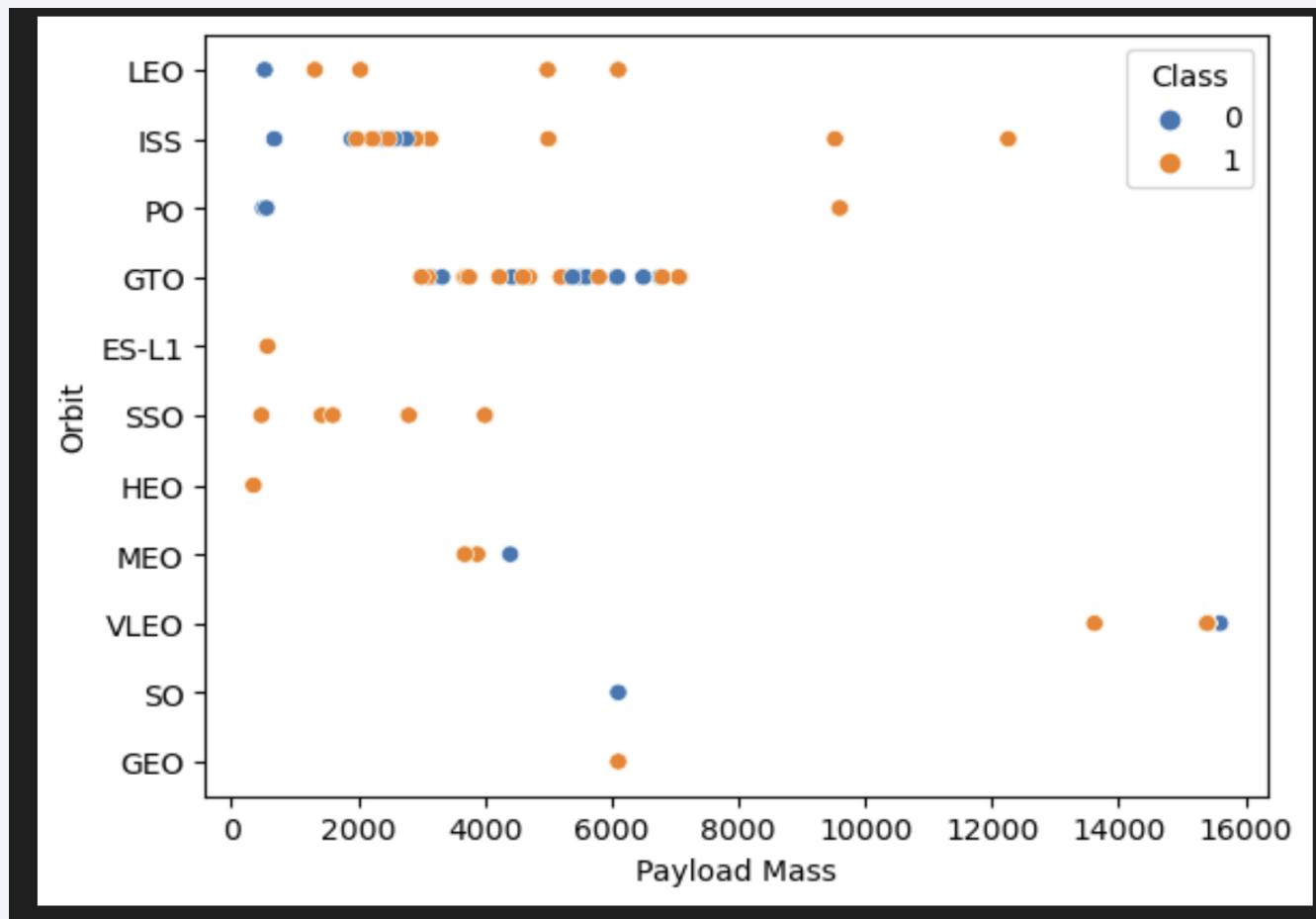
Flight Number vs. Orbit Type

- The plot below shows 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.



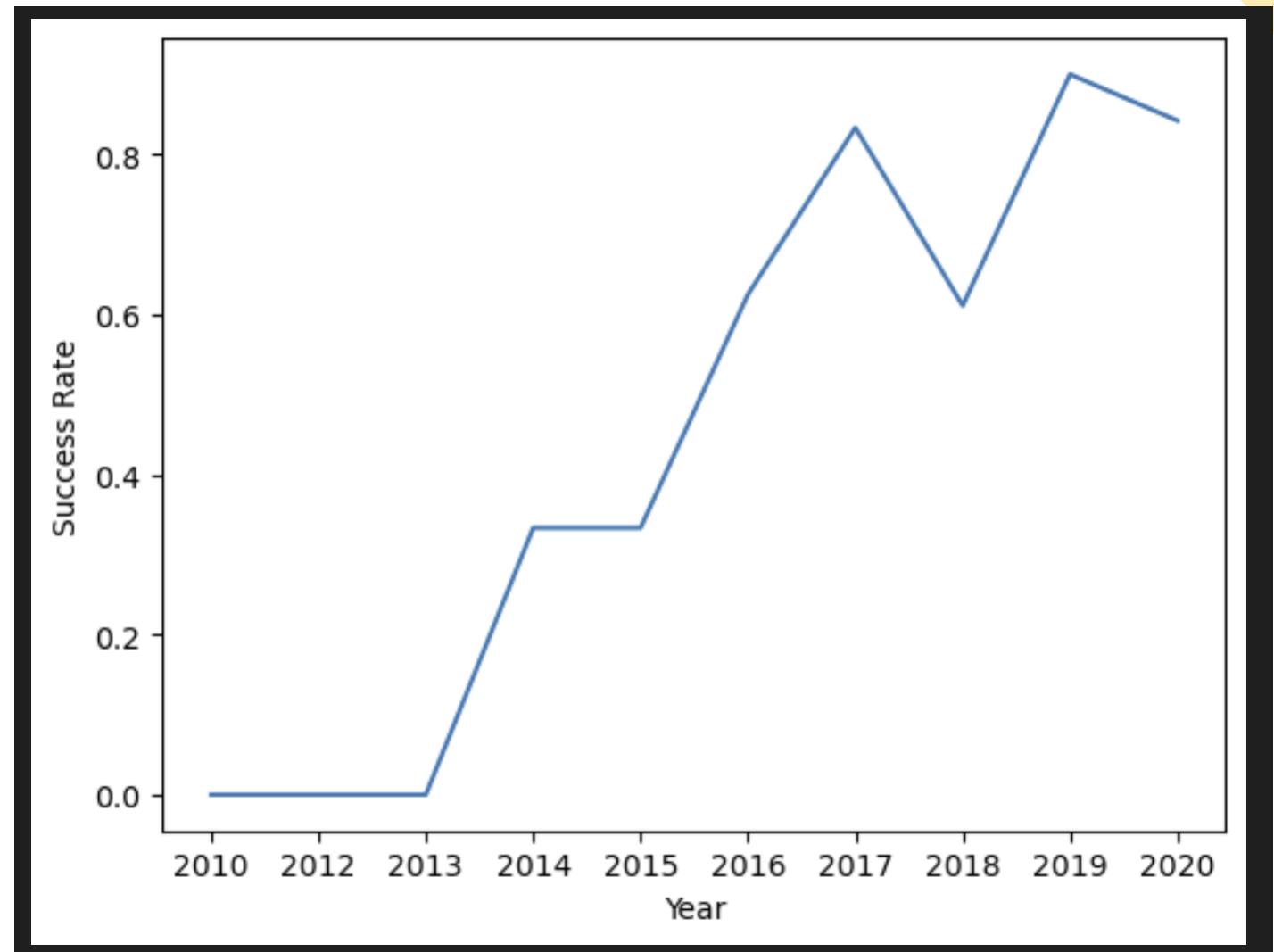
Payload vs. Orbit Type

- 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, we can observe that success rate since 2013 kept on increasing till 2020.



All Launch Site Names

- We used the key word **DISTINCT** to show only unique launch sites from the SpaceX data.

Task 1

Display the names of the unique launch sites in the space mission

```
%%sql  
select distinct "Launch_Site" from SPACEXTABLE
```

[10]

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

...

Launch_Site

CCAFS LC-40

VAFB SLC-4E

KSC LC-39A

CCAFS SLC-40

Launch Site Names Begin with 'CCA'

- We used the query above to display 5 records where launch sites begin with 'CCA'

▼ Task 2

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

```
%%sql
select * from SPACEXTABLE where "Launch_Site" like 'CCA%' LIMIT 5
[11] Python
...
* sqlite:///my_data1.db
Done.
```

	Date	Time (UTC)	Booster_Version	Launch_Site	Payload	PAYLOAD_MASS_KG_	Orbit	Customer	Mission_Outcome	Landin
...	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	0	LEO (ISS)	NASA (COTS) NRO	Success	Failure

Total Payload Mass

- We calculated the total payload carried by boosters from NASA as 45596 using the query below

Task 3

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

```
%%sql  
select sum(PAYLOAD_MASS__KG_)from SPACEXTABLE where Customer ='NASA (CRS)'
```

2]

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

Done.

```
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

Task 4

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

```
%%sql
select avg(PAYLOAD_MASS_KG_)from SPACEXTABLE where "Booster_Version" ='F9 v1.1'
13]
.. * sqlite:///my\_data1.db
Done.

.. avg(PAYLOAD_MASS_KG_)
2928.4
```

First Successful Ground Landing Date

- We observed that the dates of the first successful landing outcome on ground pad was 22nd December 2015

Task 5

List the date when the first succesful landing outcome in ground pad was acheived.

Hint:Use min function

```
%%sql  
select Date from SPACEXTABLE where "Landing_Outcome"='Success (ground pad)' order by Date asc limit 1
```

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

Done.

Date

2015-12-22

Successful Drone Ship Landing with Payload between 4000 and 6000

Task 6

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

```
%%sql
select "Booster_Version" from SPACEXTABLE where "Landing_Outcome"='Success (drone ship)' and PAYLOAD_MASS_KG_ BE
15] * sqlite:///my\_data1.db
Done.
```

Booster_Version

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

- 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

Total Number of Successful and Failure Mission Outcomes

- We used wildcard like '%' to filter for WHERE MissionOutcome was a success or a failure.

Task 7

List the total number of successful and failure mission outcomes

```
%%sql
select "Mission_Outcome" , count("Mission_Outcome") from SPACEXTABLE group by "Mission_Outcome"
16]
.. * sqlite:///my\_data1.db
Done.
```

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

Boosters Carried Maximum Payload

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

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

```
[17] %%sql  
select DISTINCT "Booster_Version" from SPACEXTABLE where PAYLOAD_MASS__KG_ = (select MAX(PAYLOAD_MASS__KG_) from S  
Python  
... * 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

- We used 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

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.

```
%%sql
select Landing_Outcome, count("Landing_Outcome") from SPACEXTABLE
where Landing_Outcome in ('Failure (drone ship)', 'Success (ground pad)')
and Date between '2010-06-04' and '2017-03-20'
group by Landing_Outcome
order by Date desc
```

[52]

Python

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

```
...   Landing_Outcome  count("Landing_Outcome")
...   Success (ground pad)      3
...   Failure (drone ship)      5
```

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.

```
%%sql
select Landing_Outcome, count("Landing_Outcome") from SPACEXTABLE
where Landing_Outcome in ('Failure (drone ship)', 'Success (ground pad)')
and Date between '2010-06-04' and '2017-03-20'
group by Landing_Outcome
order by Date desc
```

52] * [sqlite:///my_data1.db](#)
Done.

Landing_Outcome	count("Landing_Outcome")
Success (ground pad)	3
Failure (drone ship)	5

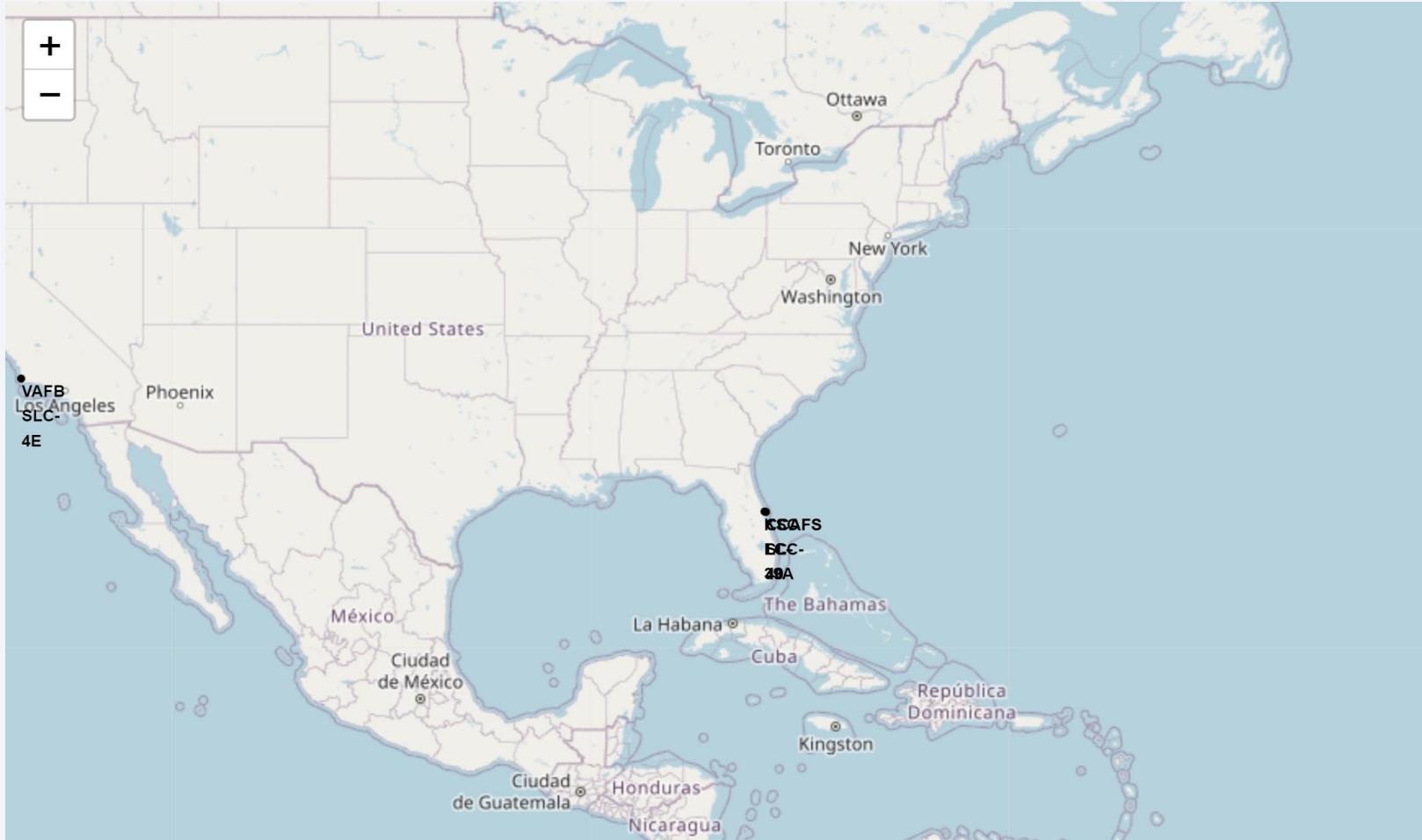
- We selected Landing outcomes and the **COUNT** of landing outcomes from the data and used the **WHERE** clause to filter for landing outcomes **BETWEEN** 2010-06-04 to 2010-03-20.
- We applied the **GROUP BY** clause to group the landing outcomes and the **ORDER BY** clause to order the grouped landing outcome in descending order.

The background of the slide is a photograph taken from space at night. It shows the curvature of the Earth's horizon against a dark blue sky. 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 corner, there is a faint, greenish glow of the aurora borealis or a similar atmospheric phenomenon.

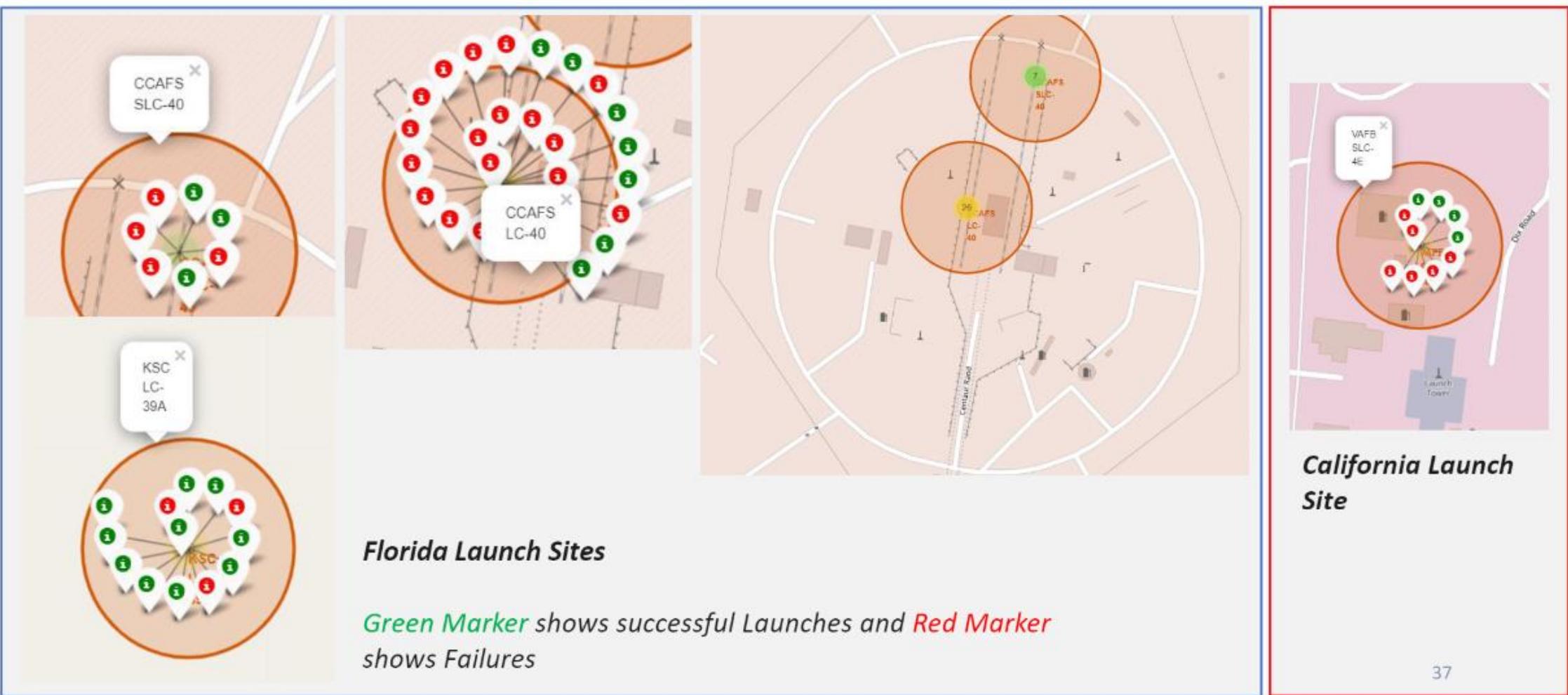
Section 4

Launch Sites Proximities Analysis

All launch sites global map markers



Markers showing launch sites with color labels

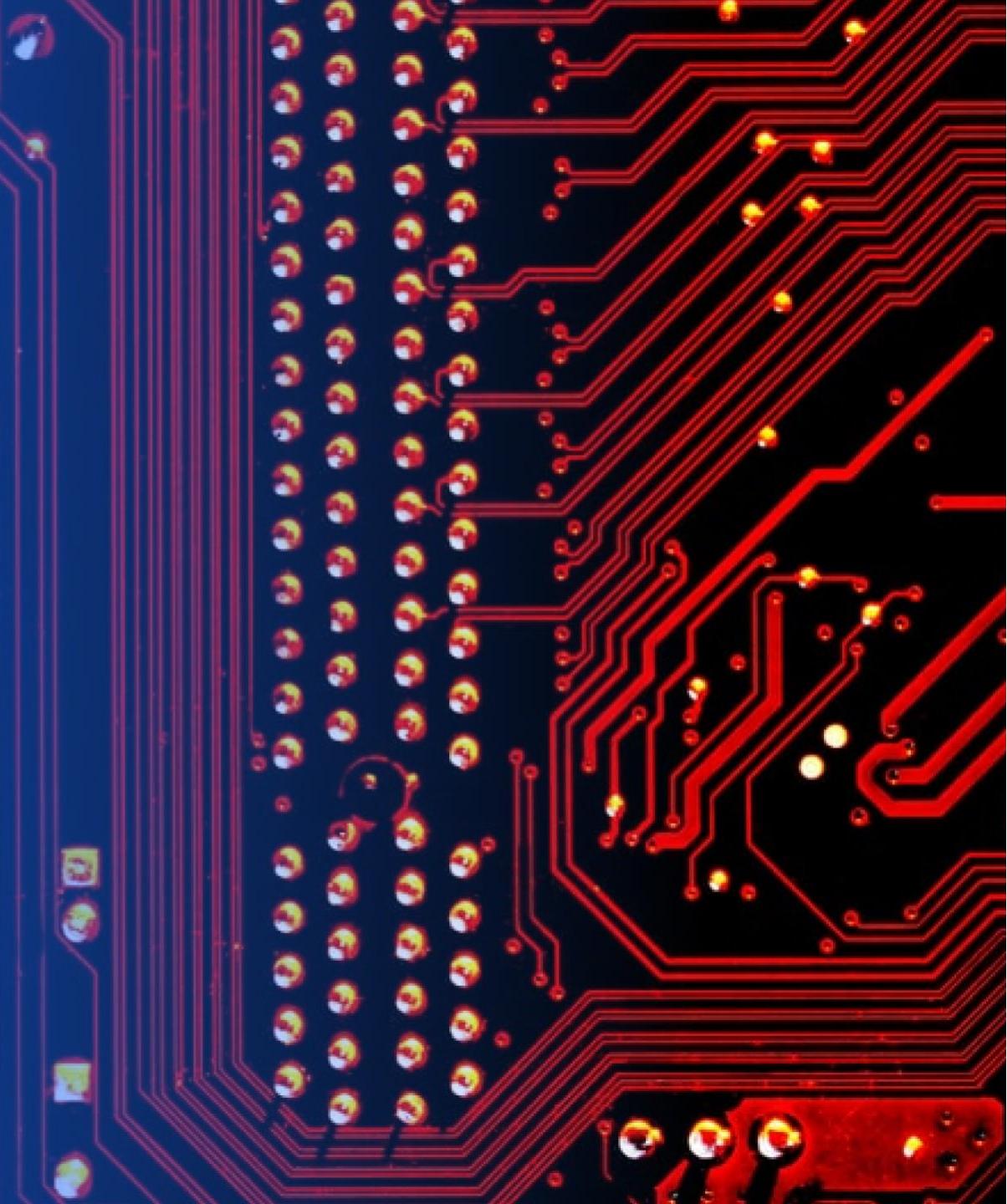


37

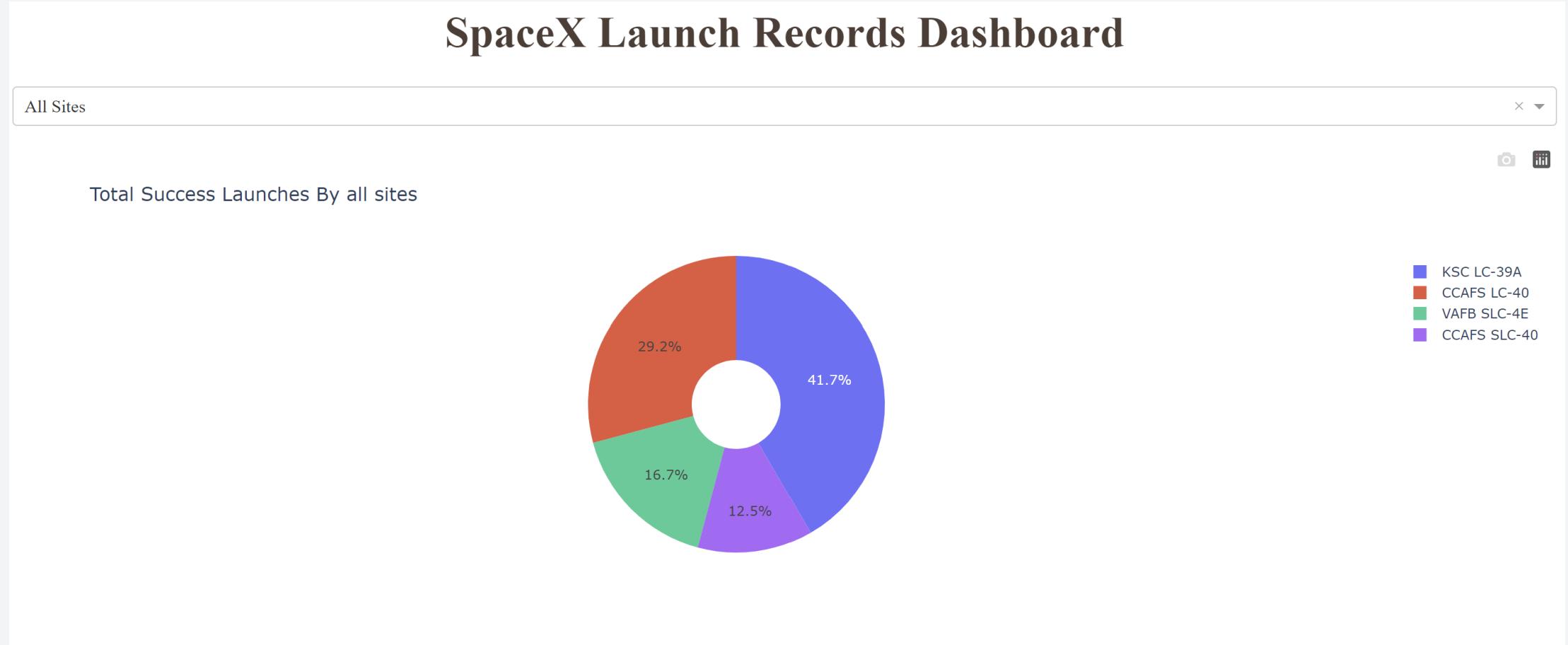
36

Section 5

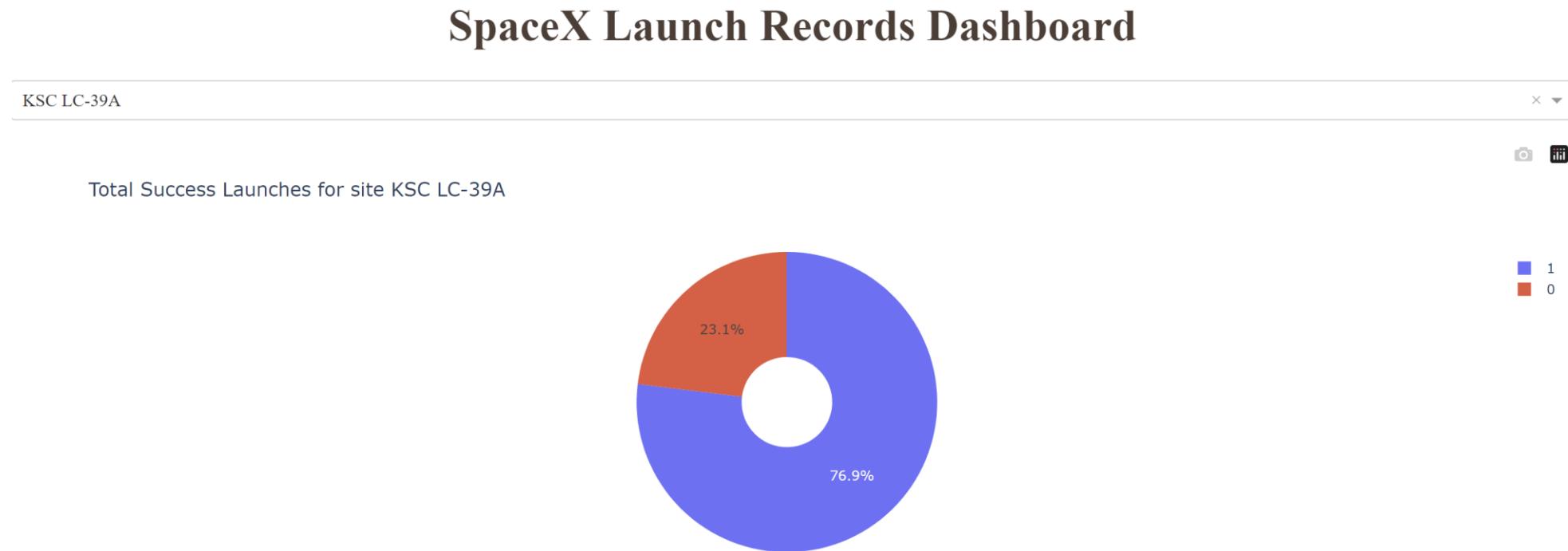
Build a Dashboard with Plotly Dash



Pie chart showing the success percentage achieved by each launch site

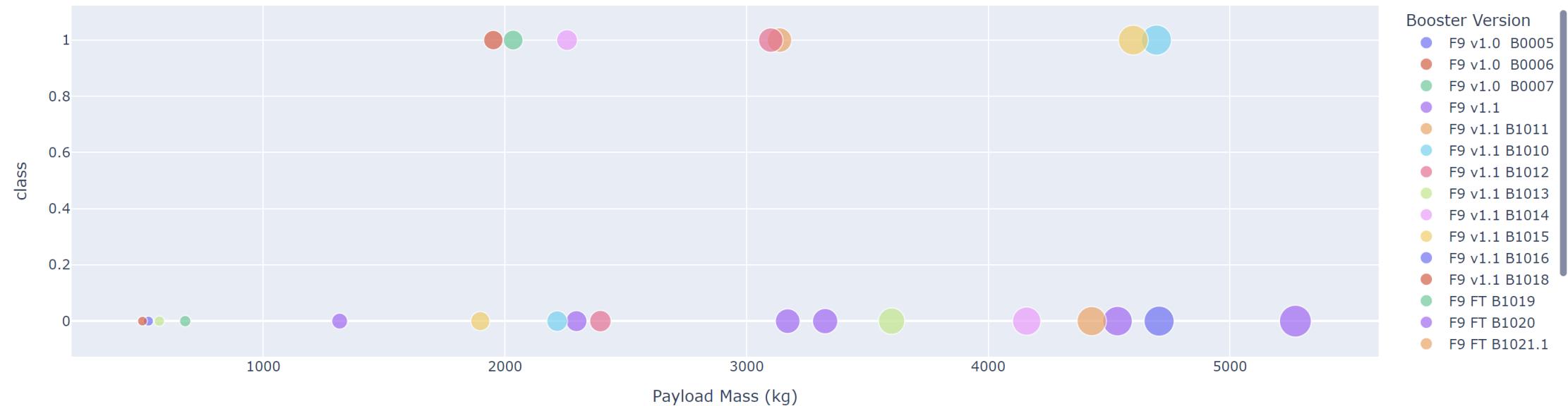


Pie chart showing the Launch site with the highest launch success ratio



Scatter plot of Payload vs Launch Outcome for all sites, with different payload selected in the range slider

Payload range (Kg):



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 landscape. The overall effect is modern and professional.

Section 6

Predictive Analysis (Classification)

Classification Accuracy

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

Find the method performs best:

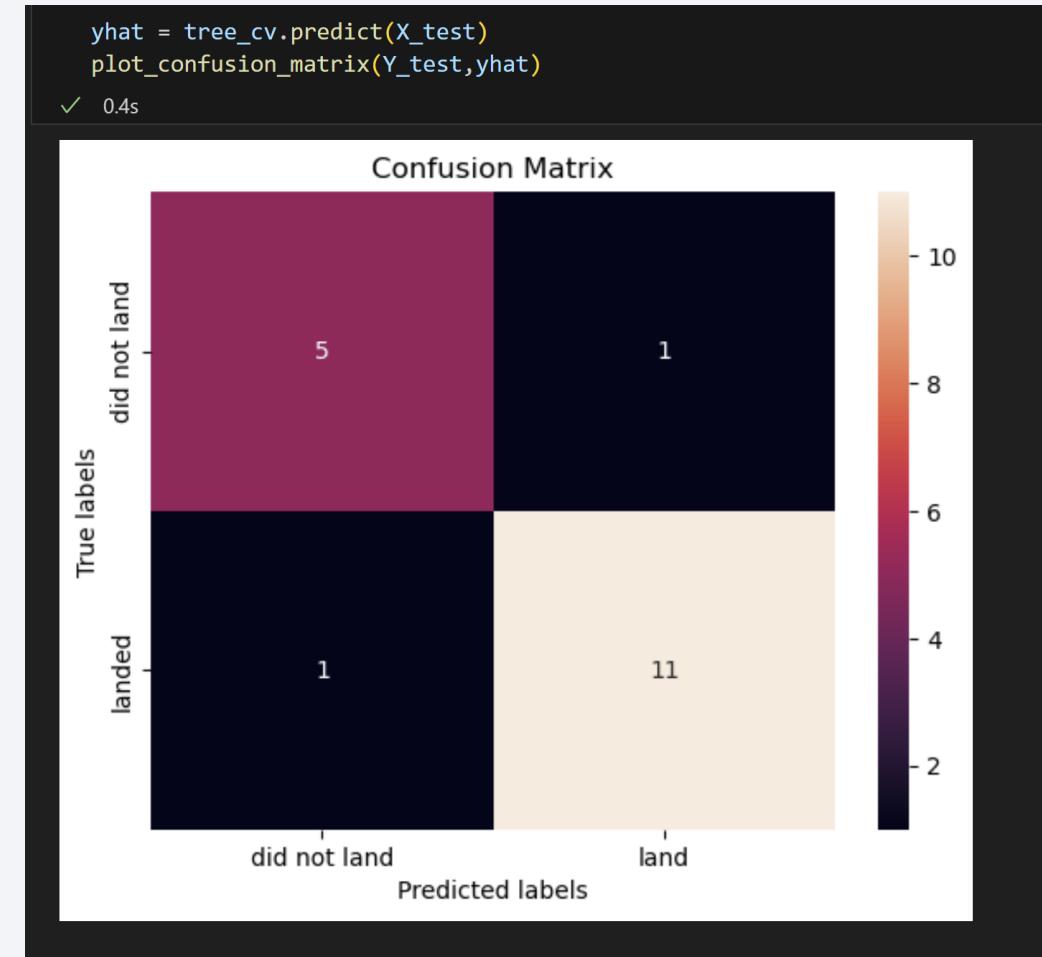
```
Report = pd.DataFrame({'Method' : ['Test Data Accuracy']})  
  
knn_accuracy=knn_score  
Decision_tree_accuracy=tree_score  
SVM_accuracy=svm_score  
Logistic_Regression=lrscore  
  
Report['Logistic_Reg'] = [Logistic_Regression]  
Report['SVM'] = [SVM_accuracy]  
Report['Decision Tree'] = [Decision_tree_accuracy]  
Report['KNN'] = [knn_accuracy]  
  
Report.transpose()
```

✓ 0.0s

Method	Test Data Accuracy
Logistic_Reg	0.833333
SVM	0.833333
Decision Tree	0.888889
KNN	0.833333

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

We can conclude that:

- 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 launches of any sites.
- The Decision tree classifier is the best machine learning algorithm for this task.

Thank you!

