

Winning Space Race with Data Science

UFAQUE SHADAB
21 November 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 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, we decoded the response content as a Json using `.json()` function call and turn it into a pandas dataframe using `.json_normalize()`.
 - 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

- I 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/UFAQUE123/IB-M-Data-Science-Capstone-Project/blob/main/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)  
Check the content of the response  
[1]: # print(response.content)  
You should see the response contains massive information about SpaceX launches. Next, let's try to discover some more relevant information for this project.  
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  
[11]: 200  
Now we decode the response content as a json using .json() and turn it into a Pandas dataframe using .json_normalize()  
[12]: # Use json_normalize method to convert the json result into a dataframe  
data = pd.json_normalize(response.json())  
Using the dataframe data print the first 5 rows  
[13]: # Get the head of the dataframe  
data.head(2)  
[13]: static_fire_date_utc static_fire_date_unix tbd net window rocket success details crew ships capsules payloads launchpad auto_update failures flight_number name date_utc d  
  
0 2006-03-17T00:00:00Z 1142554e+09 False False 0.0 5e9d0d95eda6995f709d1eb False 33 Engine failure at 33 seconds and loss of vehicle None None 1 [Seb0e4b5b6c3bb0006eeb1e1] Se9e4502f5090995de566886 True 1 FalconSat 2006-03-24T22:30:00Z 114  
[[{"time": 33, "altitude": None, "reason": "merlin engine failure"}]]
```

Task 3: Dealing with Missing Values

Calculate below the mean for the `PayloadMass` using the `.mean()`. Then use the mean and the `.replace()` function to replace `np.nan` values in the data with the mean you calculated.

```
[28]: # Calculate the mean value of PayloadMass column
payloadmass_mean = data_falcon9['PayloadMass'].mean()
payloadmass_mean

# Replace the np.nan values with its mean value
data_falcon9['PayloadMass'] = data_falcon9['PayloadMass'].replace(np.nan, payloadmass_mean)
data_falcon9.isnull().sum()

/home/jupyterlab/conda/envs/python/lib/python3.7/site-packages/ipykernel_launcher.py:6: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy
```

```
[28]: FlightNumber      0
Date                  0
BoosterVersion       0
PayloadMass          0
Orbit                 0
LaunchSite            0
Outcome               0
Flights                0
GridFins              0
Reused                 0
Legs                   0
LandingPad           26
Block                  0
ReusedCount           0
Serial                  0
Longitude              0
Latitude                0
dtype: int64
```

You should see the number of missing values of the `PayloadMass` change to zero.

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/UFAQUE123/IBM-Data-Science-Capstone-Project/blob/main/jupyter-labs-webscraping.ipynb>

```
[4]: static_url = "https://en.wikipedia.org/w/index.php?title=List_of_Falcon_9_and_Falcon_Heavy_launches&oldid=1027686922"
headers = {
    "User-Agent": "Mozilla/5.0 (Windows NT 10.0; Win64; x64) "
                  "AppleWebKit/537.36 (KHTML, like Gecko) "
                  "Chrome/91.0.4472.124 Safari/537.36"
}
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 and headers
# assign the response to a object

# TASK 1: Request the Falcon9 Launch Wiki page
response = requests.get(static_url, headers=headers)

# Check status
response.status_code
[5]: 200
Create a BeautifulSoup object from the HTML response
```

```
[6]: # Use BeautifulSoup() to create a BeautifulSoup object from a response text content
from bs4 import BeautifulSoup

# Create a BeautifulSoup object from the response text
soup = BeautifulSoup(response.text, "html.parser")

# Optional: Check the type
type(soup)
[6]: bs4.BeautifulSoup
Print the page title to verify if the BeautifulSoup object was created properly
```

```
[7]: # Use soup.title attribute
# Print the page title to verify BeautifulSoup object
print(soup.title)
print(soup.title.string)
<title>List of Falcon 9 and Falcon Heavy launches - Wikipedia</title>
List of Falcon 9 and Falcon Heavy launches - Wikipedia
```

```
[10]: column_names = []

# Apply find_all() on <th> elements inside the table header
for th in first_launch_table.find_all("th"):
    name = extract_column_from_header(th)

    # Append only non-empty names
    if name is not None and len(name) > 0:
        column_names.append(name)

column_names

# Apply find_all() function with `th` element on first_launch_table
# Iterate each th element and apply the provided extract_column_from_header() to get a column name
# Append the Non-empty column name (`if name is not None and Len(name) > 0`) into a list called column_names
```

```
[10]: ['Flight No.',
       'Date and time ( )',
       'Launch site',
       'Payload',
       'Payload mass',
       'Orbit',
       'Customer',
       'Launch outcome']
```

Check the extracted column names

```
[11]: print(column_names)

['Flight No.', 'Date and time ( )', 'Launch site', 'Payload', 'Payload mass', 'Orbit', 'Customer', 'Launch outcome']
```

Data Wrangling

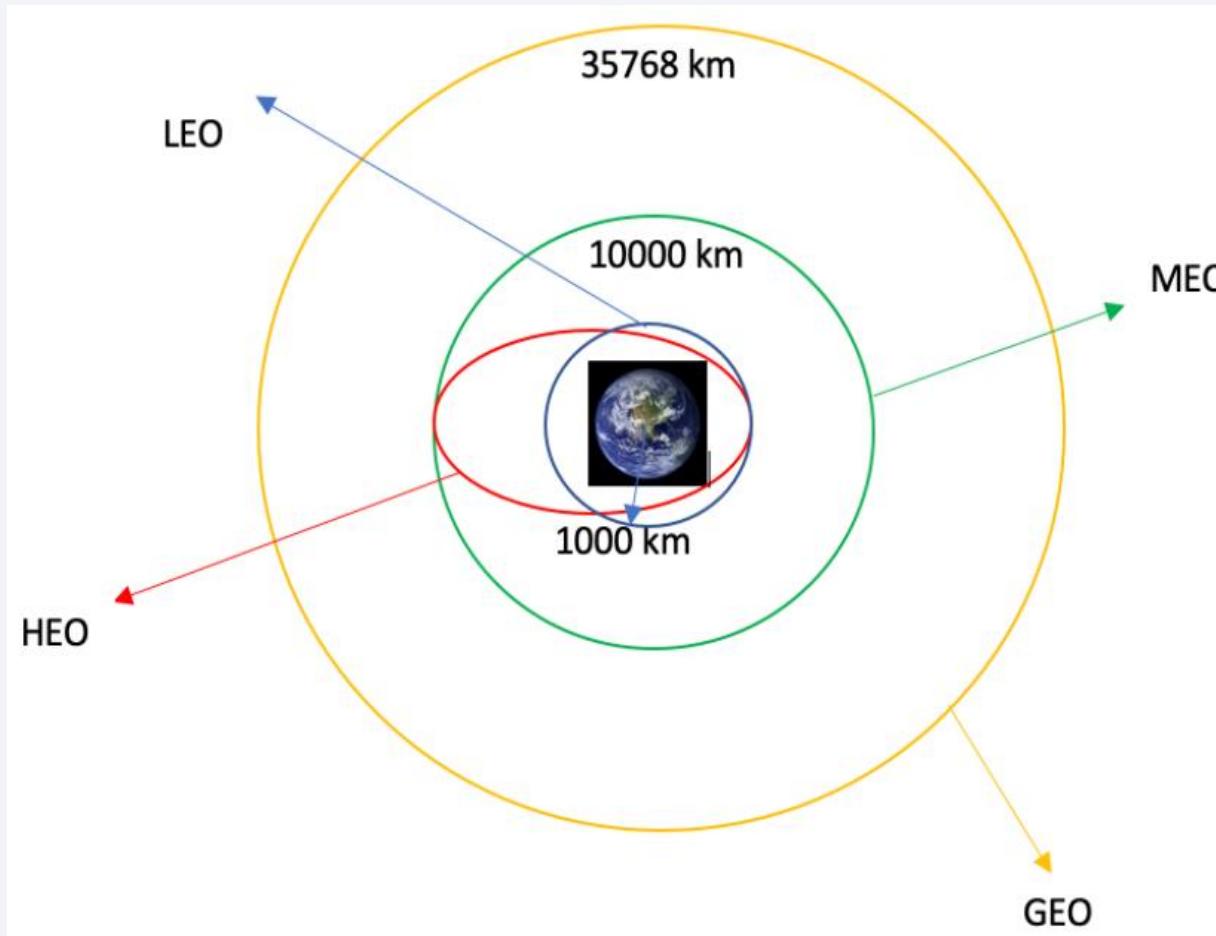
We performed exploratory data analysis and determined the training labels.

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.

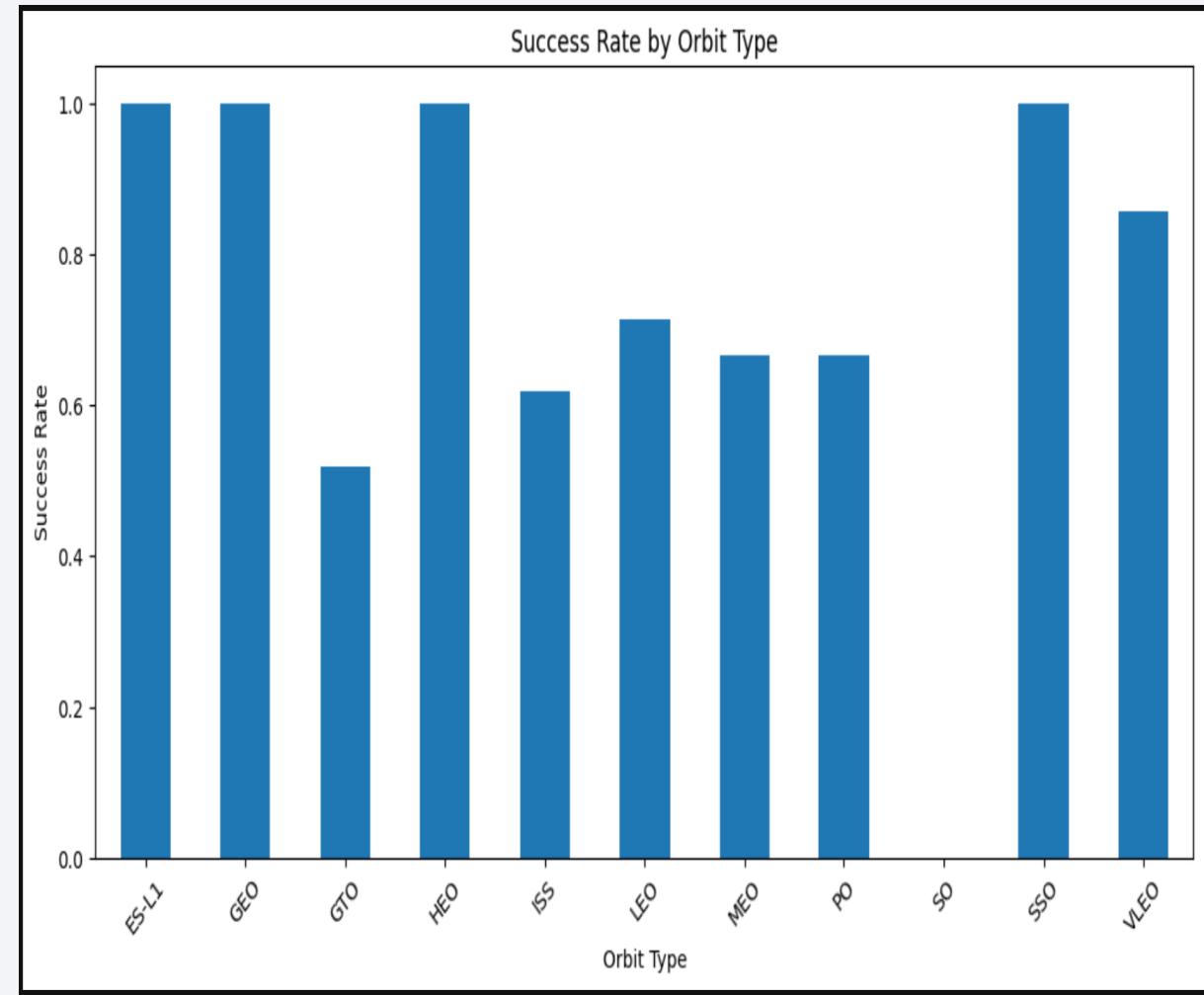
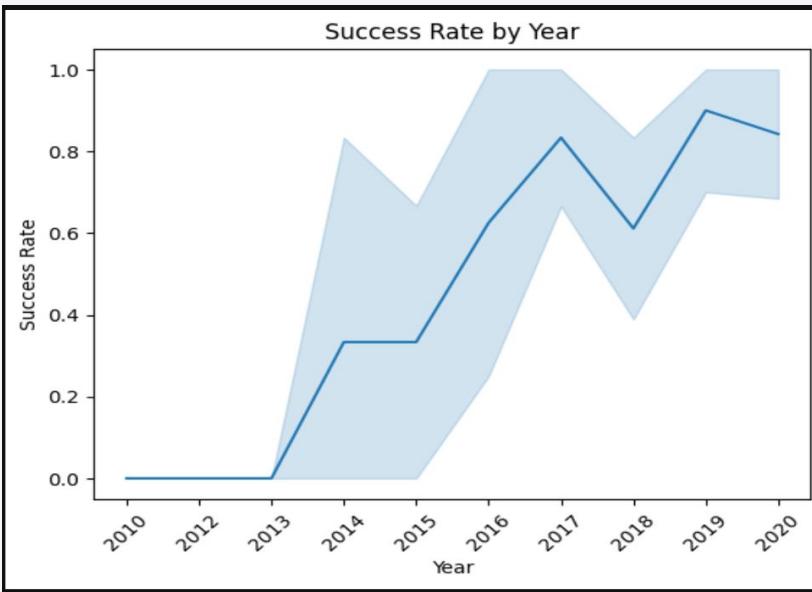
The link to the notebook is:

<https://github.com/UFAQUE123/IBM-Data-Science-Capstone-Project/blob/main/labs-jupyter-spacex-Data%20wrangling.ipynb>



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.
- The link to the notebook is:
<https://github.com/UFAQUE123/IBM-Data-Science-Capstone-Project/blob/main/edadataviz.ipynb>



EDA with SQL

- We loaded the SpaceX dataset into a PostgreSQL database without leaving the jupyter notebook.
- 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.
- The link to the notebook is: https://github.com/UFAQUE123/IBM-Data-Science-Capstone-Project/blob/main/jupyter-labs-eda-sql-coursera_sqlite.ipynb

Build an Interactive Map with Folium

- We 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.
- We 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, we identified which launch sites have relatively high success rate.
- We calculated the distances between a launch site to its proximities. We 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

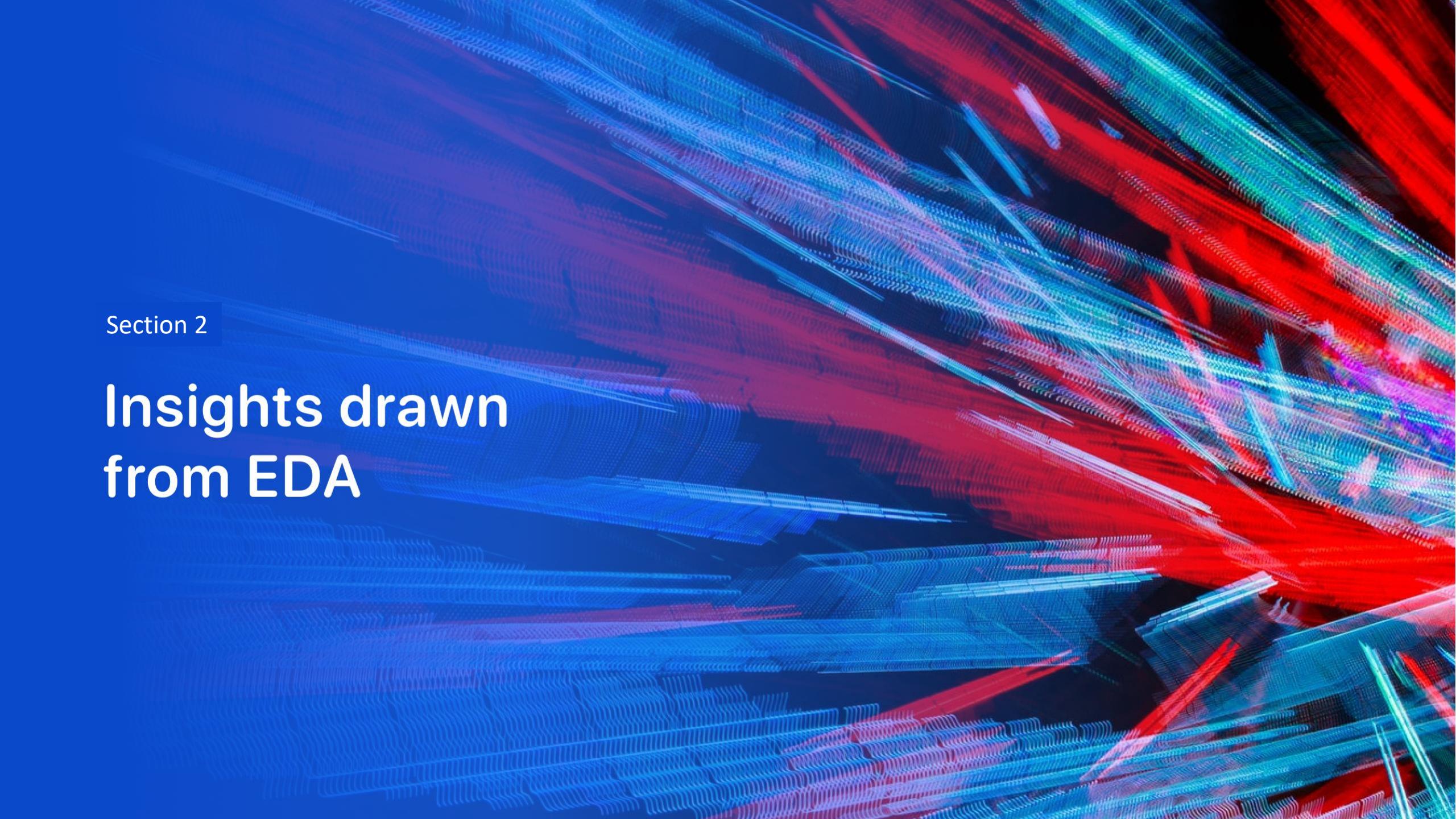
- We built an interactive dashboard with Plotly dash
- We plotted pie charts showing the total launches by a certain sites
- We plotted scatter graph showing the relationship with Outcome and Payload Mass (Kg) for the different booster version.
- The link to the notebook is <https://github.com/UFAQUE123/IBM-Data-Science-Capstone-Project/blob/main/spacex-dash-app.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

Results

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

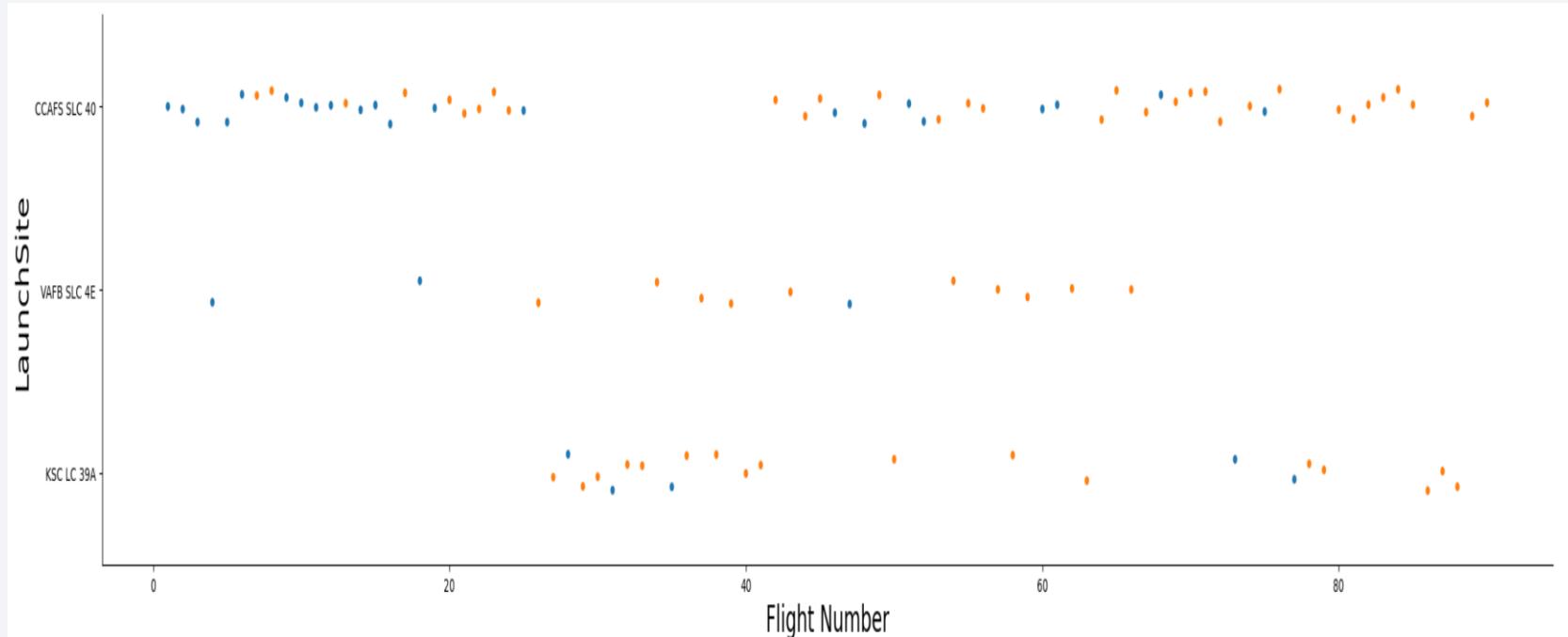
The background of the slide features a complex, abstract digital visualization. It consists of numerous thin, glowing lines that create a sense of depth and motion. The lines are primarily blue and red, with some green and purple highlights. They form a grid-like structure that curves and twists across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

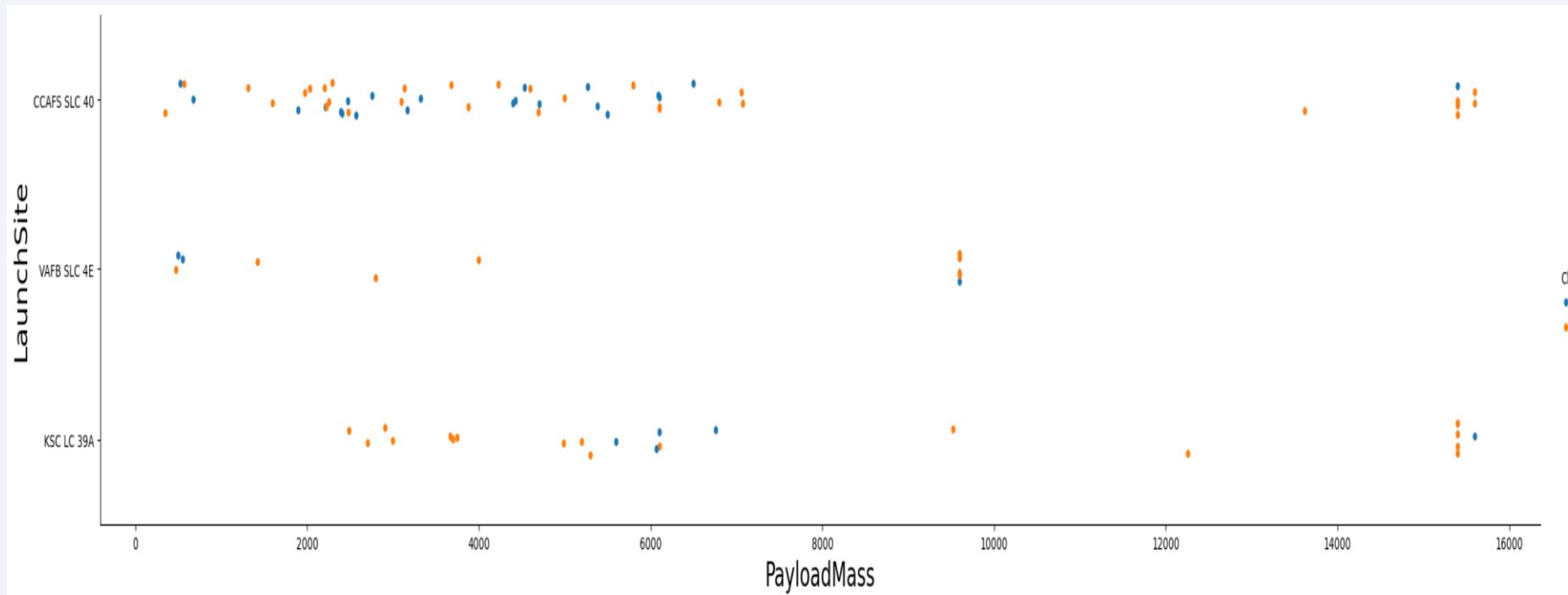
Flight Number vs. Launch Site

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



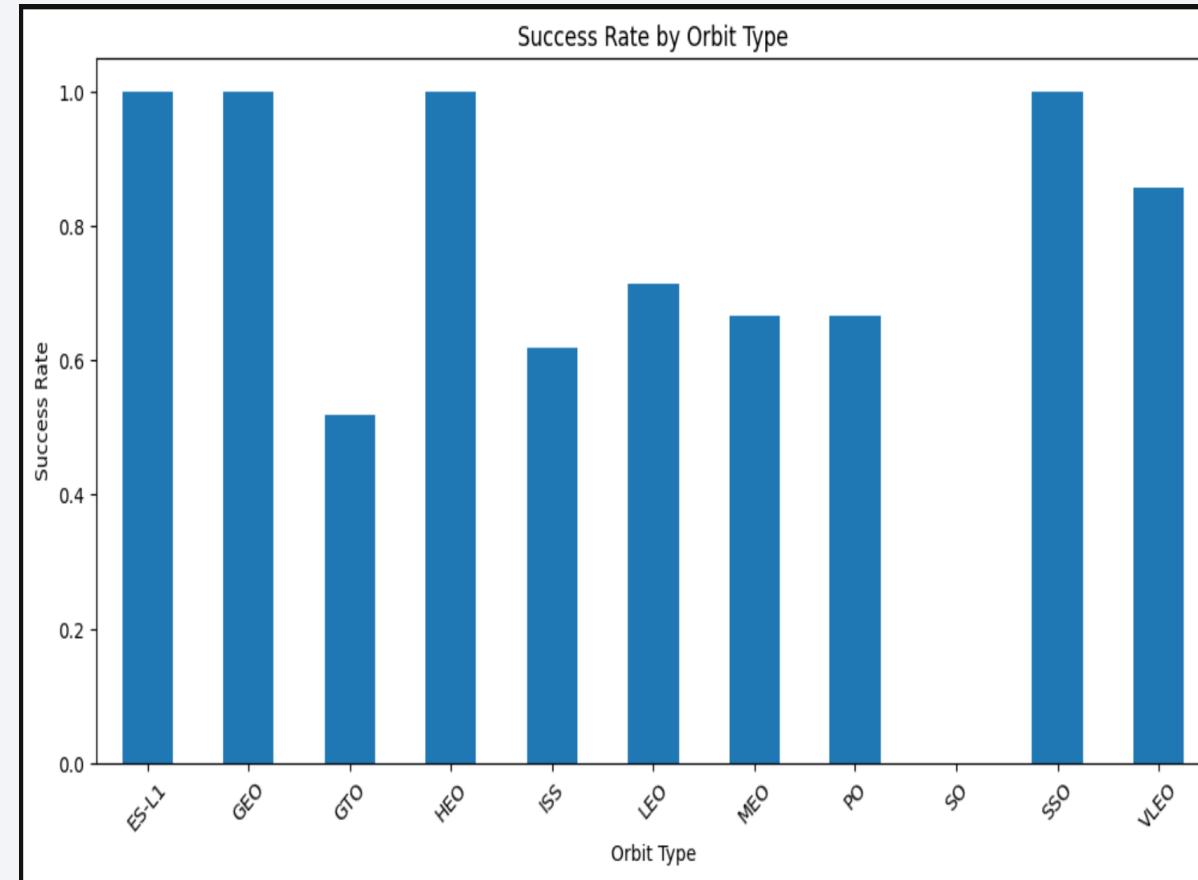
Payload vs. Launch Site

- The greater the Payload mass for Launch site 'CCAFS ALC 40' the higher the success rate for the rocket.



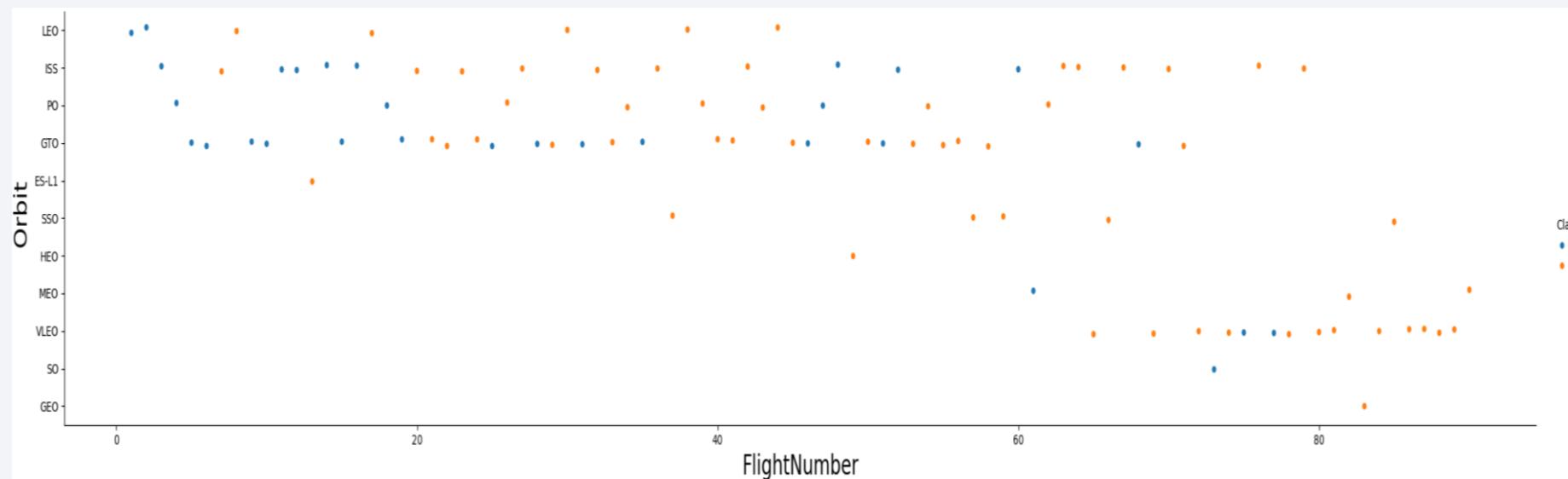
Success Rate vs. Orbit Type

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



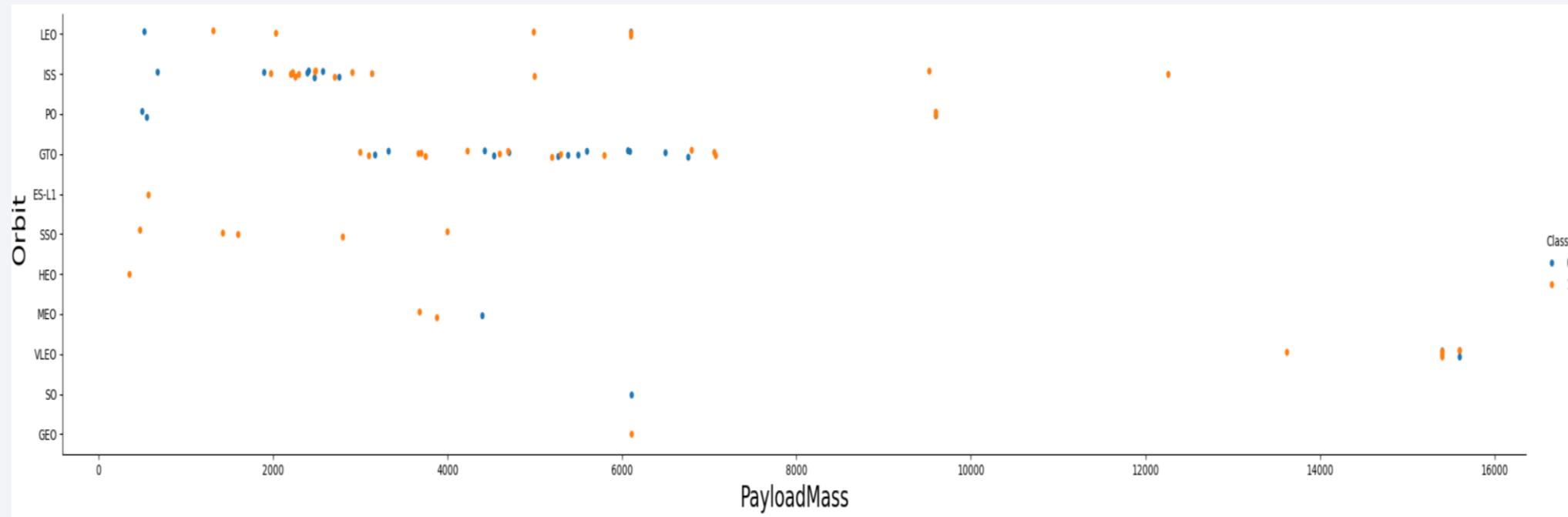
Flight Number vs. Orbit Type

- 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.



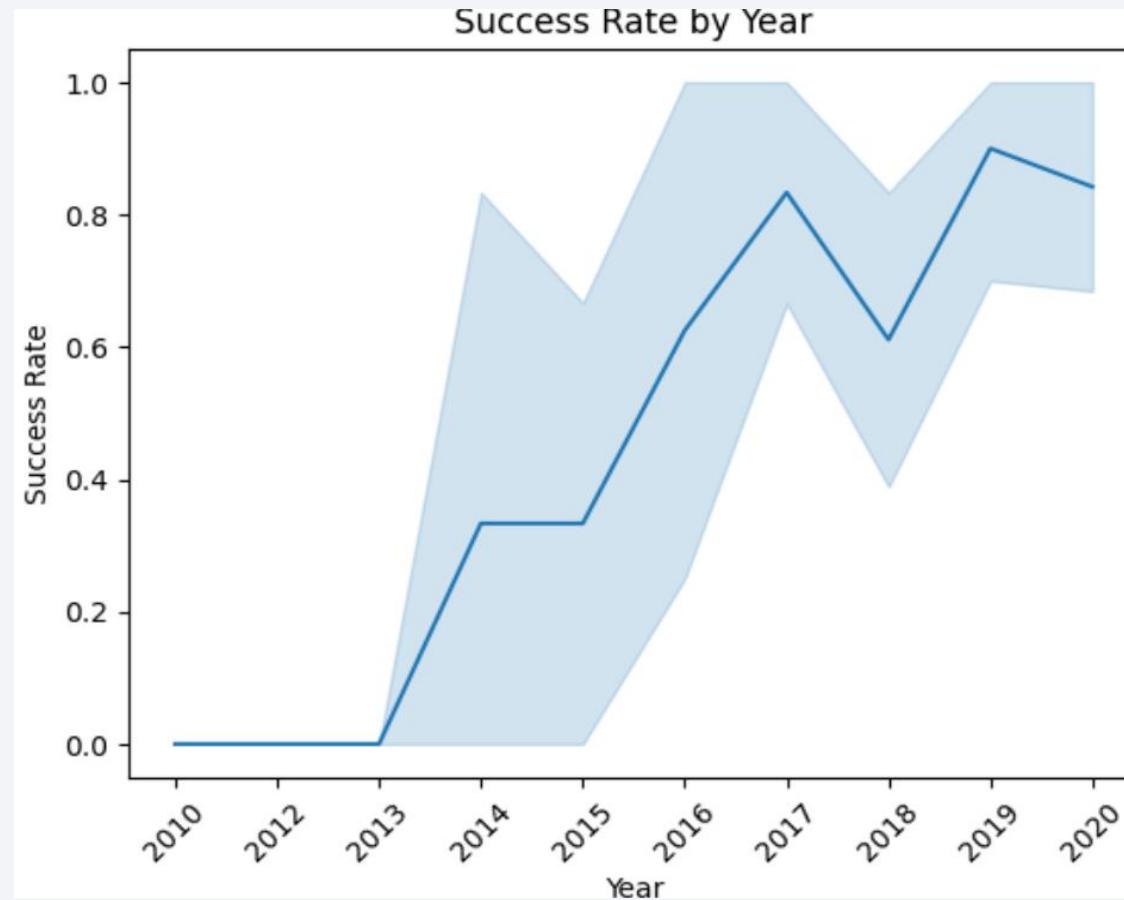
Payload vs. Orbit Type

- 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

```
[10]: %sql select Distinct(Launch_Site) from SPACEXTBL;
```

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

```
Done.
```

```
[10]: 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'

```
[11]: %sql select * from SPACEXTBL 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

- We calculated the total payload carried by boosters from NASA as 45596.0 using the query below.

Task 3

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

```
[12]: %sql select total(PAYLOAD_MASS__KG_) from SPACEXTBL where Customer='NASA (CRS)';
```

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

```
Done.
```

```
[12]: total(PAYLOAD_MASS_KG_)
```

```
45596.0
```

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

```
[13]: %sql select AVG(PAYLOAD_MASS__KG_) from SPACEXTBL where Booster_Version = 'F9 v1.1';
```

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

```
Done.
```

```
[13]: 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

```
[14]: %sql select min(Date) from SPACEXTBL where Landing_Outcome= 'Success (ground pad)';

* sqlite:///my_data1.db
Done.

[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.

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

```
[15]: %sql select Booster_Version from SPACEXTBL where Landing_Outcome = 'Success (drone ship)' AND PAYLOAD_MASS__KG_ > 4000 AND PAYLOAD_MASS__KG_ < 6000;
```

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

```
[15]: 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 ‘group by’ to filter whether MissionOutcome was a success or a failure.

Task 7

List the total number of successful and failure mission outcomes

```
[16]: %sql select Mission_Outcome, count(*) AS Total from SPACEXTBL GROUP BY Mission_Outcome;
```

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

Mission_Outcome	Total
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.

Task 8
List all the booster_versions that have carried the maximum payload mass, using a subquery with a suitable aggregate function.

```
[17]: %sql select Booster_Version from SPACEXTBL where PAYLOAD_MASS__KG_=(select max(PAYLOAD_MASS__KG_) from SPACEXTBL);  
  
* sqlite:///my_data1.db  
Done.  
[17]: 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

Task 9

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.

Note: SQLite does not support monthnames. So you need to use substr(Date, 6,2) as month to get the months and substr(Date,0,5)='2015' for year.

```
[18]: %sql select substr(Date, 6, 2) AS Month, Landing_Outcome,Booster_Version,Launch_Site from SPACEXTBL where substr(Date, 0, 5) = '2015' AND Landing_Outcome LIKE '%Failure%' AND Landing_Outcome LIKE '%drone ship%';
```

```
* 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

- 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. It can be seen in the code snippet attached in the next slide.

Task 10

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.

```
[19]: %sql select Landing_Outcome, count(*) as outcome_count from SPACEXTBL where Date between '2010-06-04' and '2017-03-20' group by Landing_Outcome order by outcome_count desc;
```

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

```
Done.
```

```
[19]:   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

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 a large, brightly lit urban area is visible. In the upper left quadrant, there are greenish-yellow bands of light, likely the Aurora Borealis or Australis, dancing across the atmosphere.

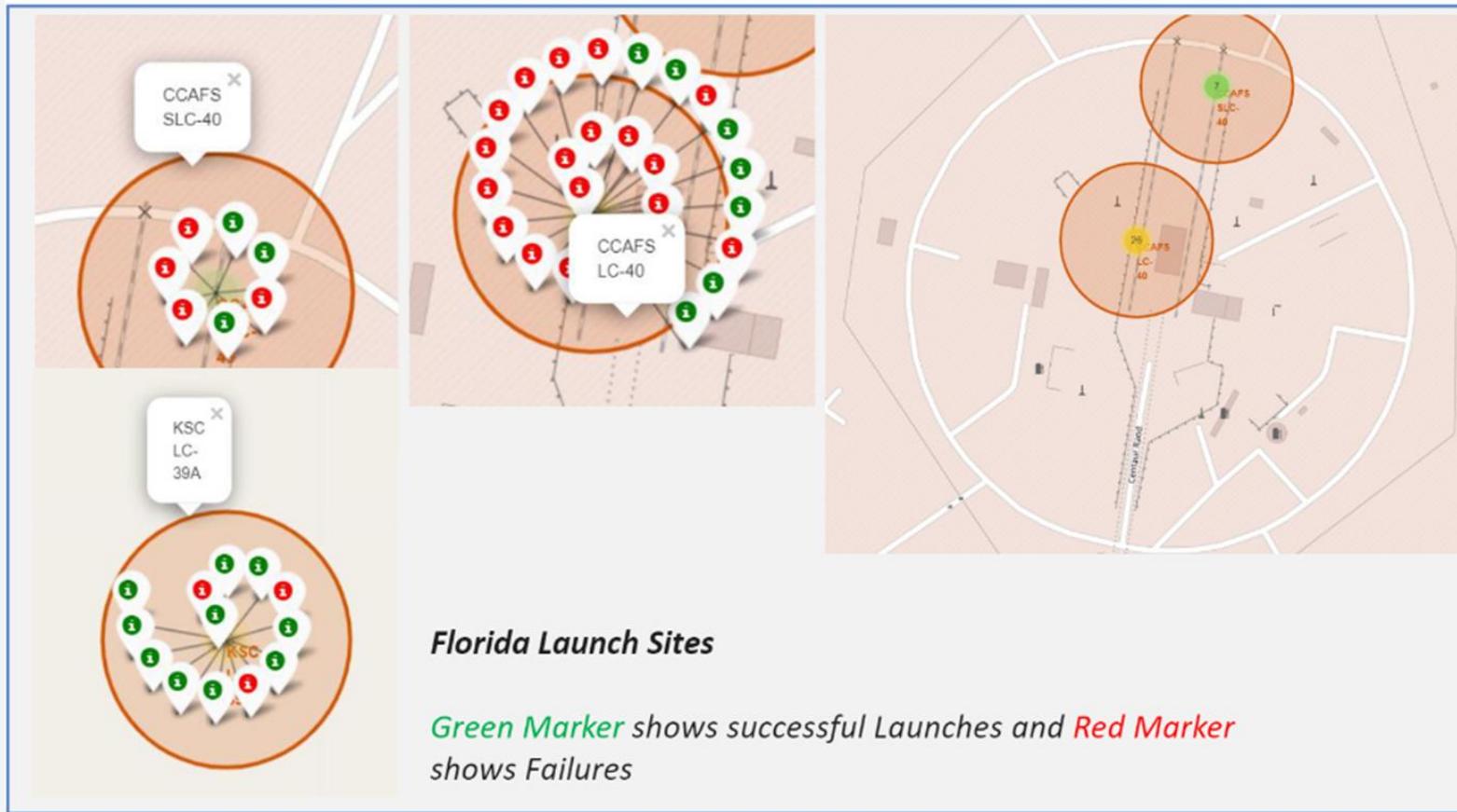
Section 3

Launch Sites Proximities Analysis

All Launch sites global map markers

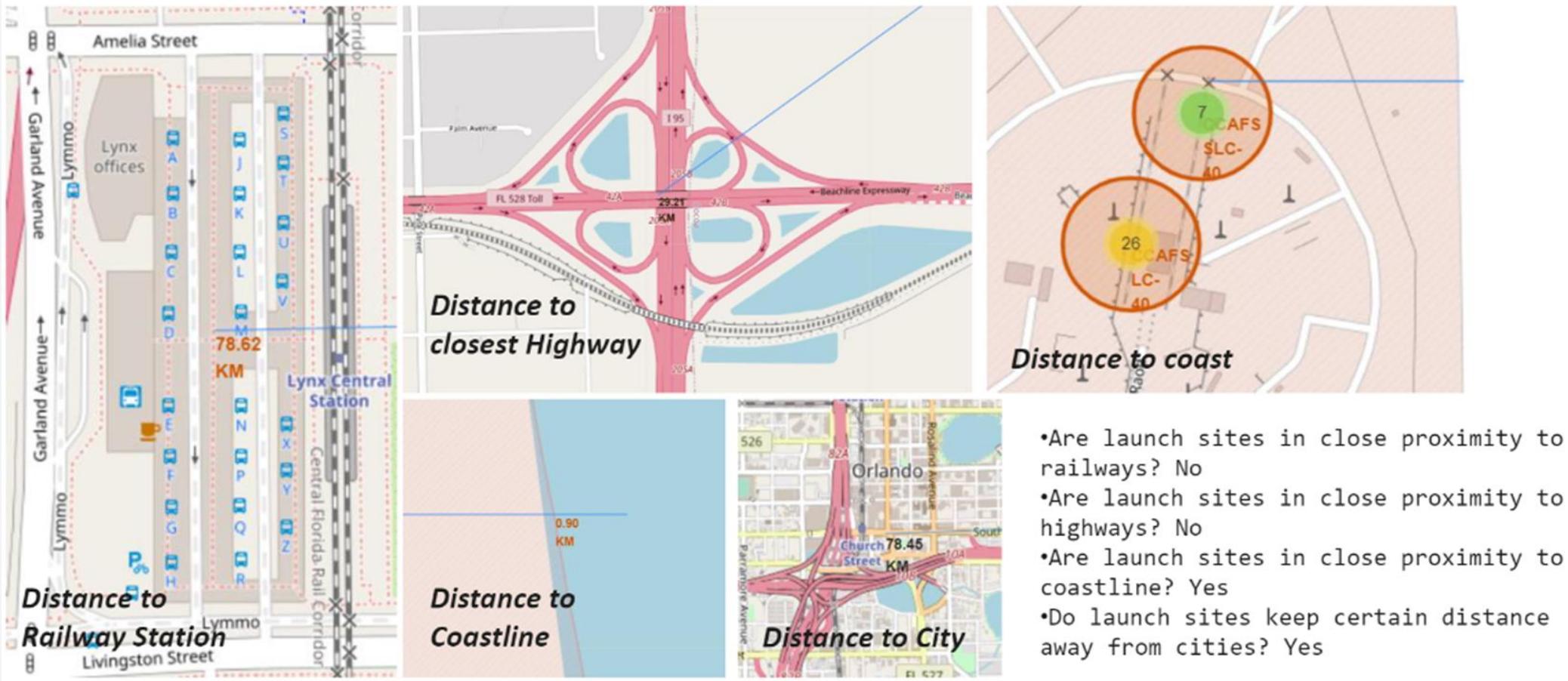


Markers showing Launch sites with color labels



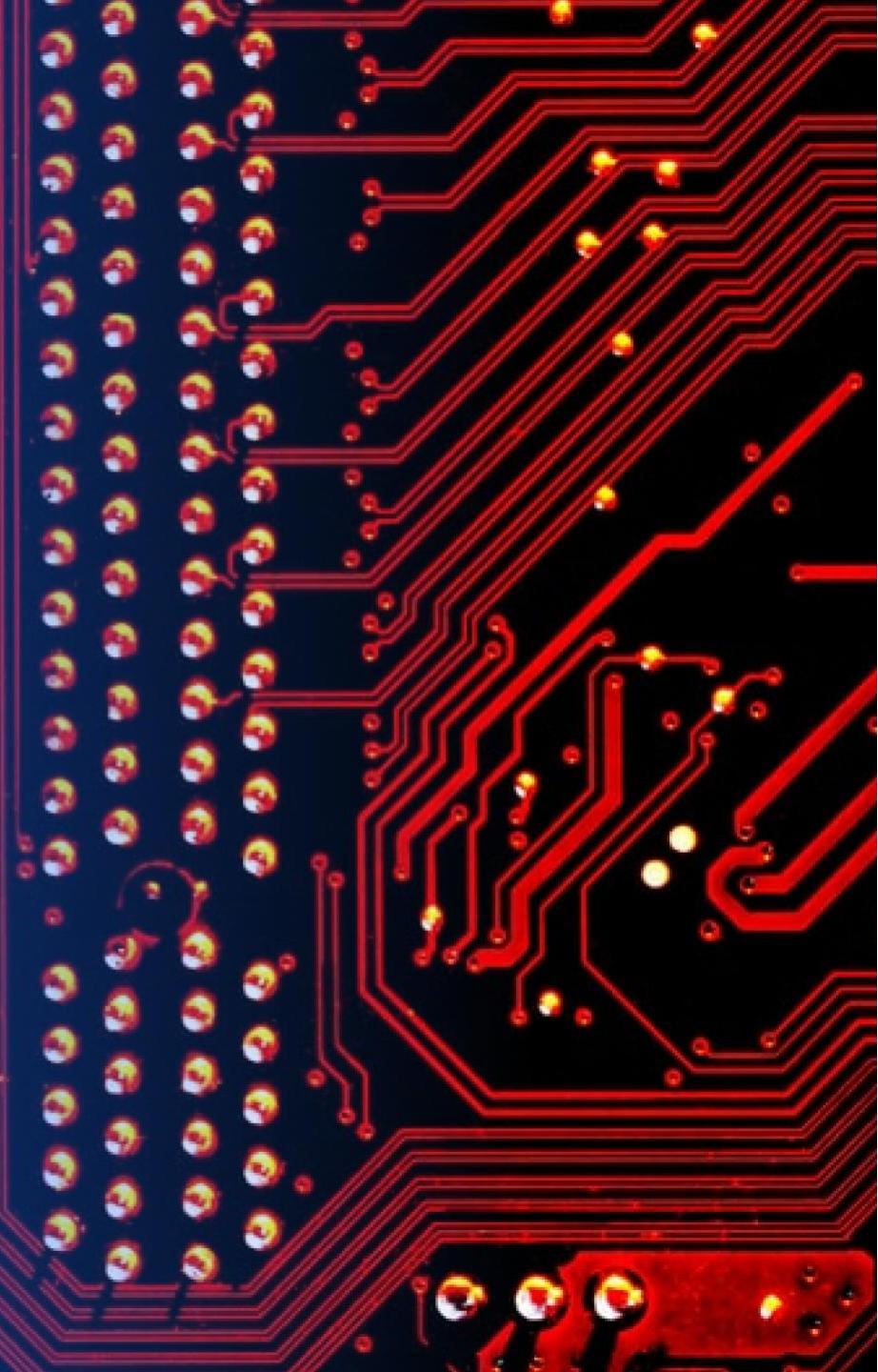
37

Launch Site Distance to Landmarks



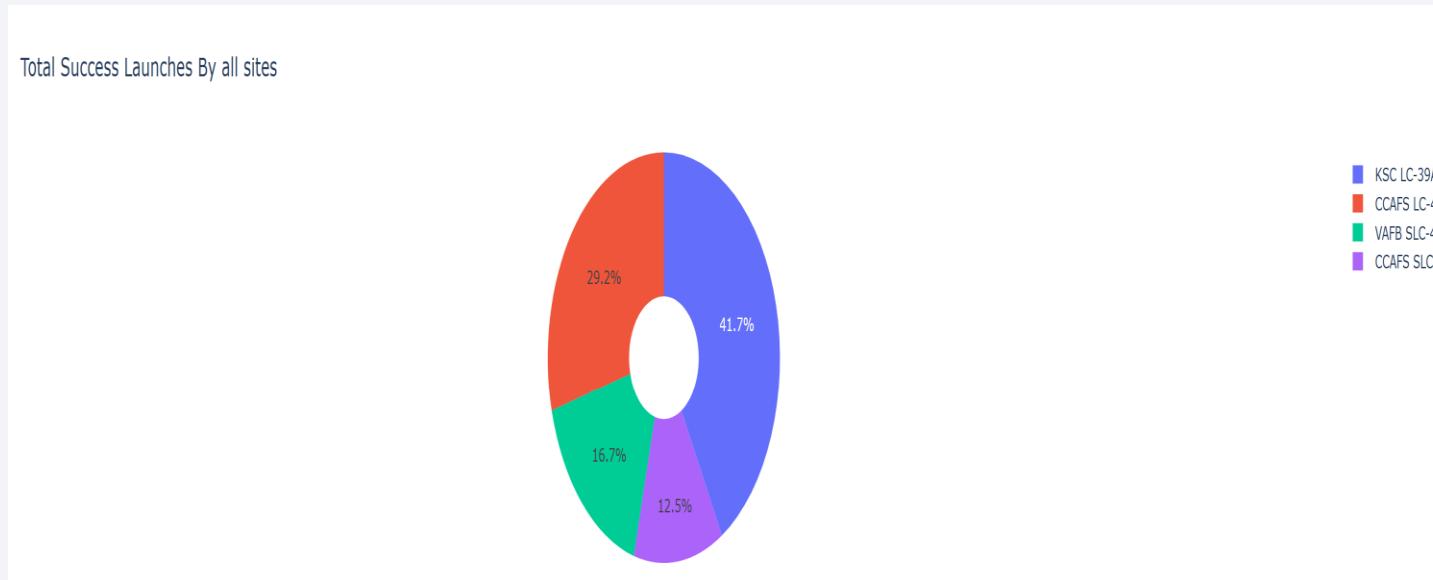
Section 4

Build a Dashboard with Plotly Dash



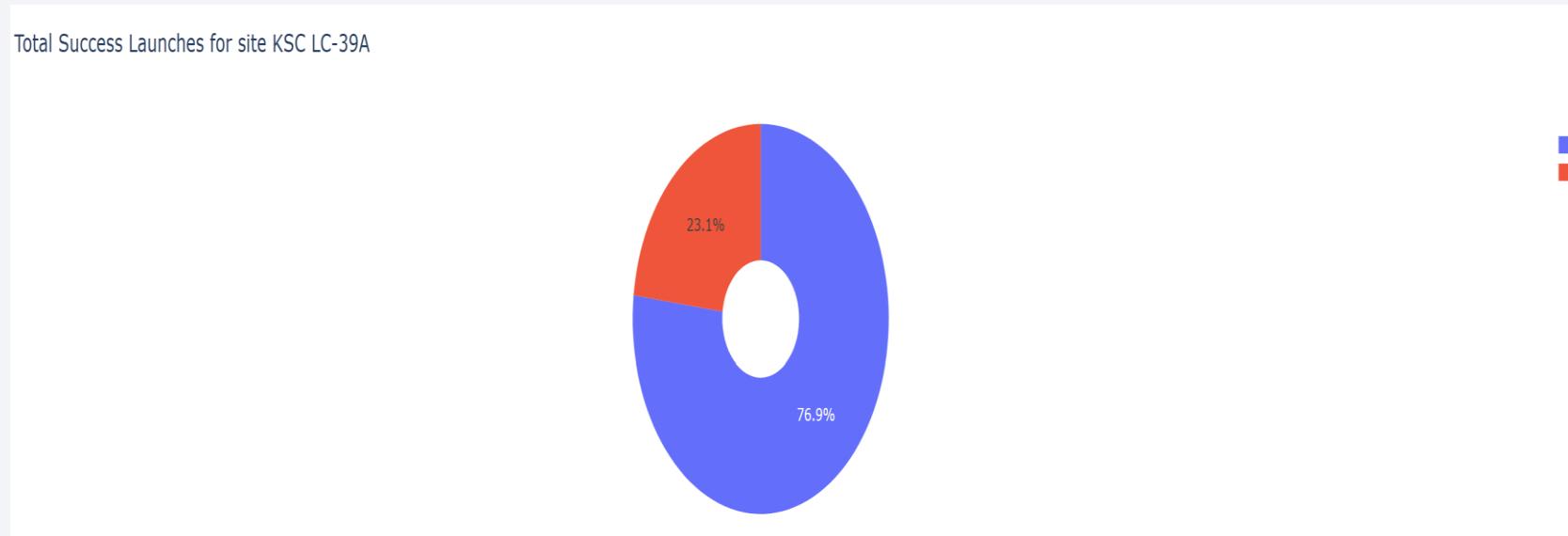
Pie chart showing the success percentage achieved by each launch site

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



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

- KSC LC-39A has achieved the success rate of 76.9%, while its failure rate is 23.1%, as shown in a pie chart below



Scatter plot of Payload vs Launch Outcome for all sites

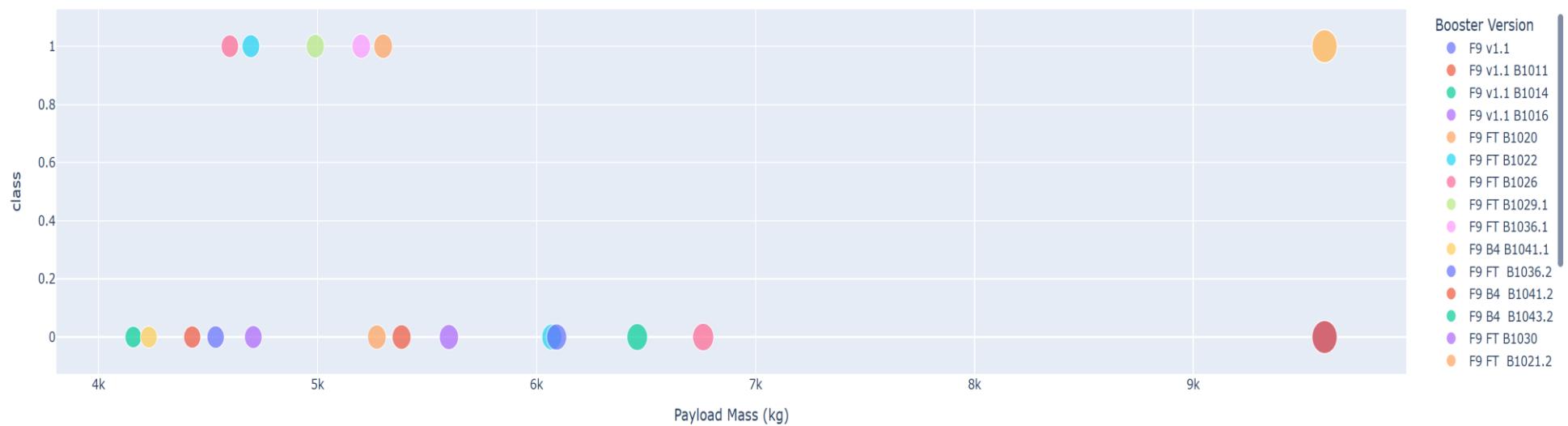
- It can be seen in the scatter plot below and the one in the next slide that:
“the success rate for low weighted payload is higher than the heavy weighted payloads”.

Low Weighted Payload Mass (0 kg – 4000 kg)



Heavy Weighted Payload Mass (4000 kg – 10,000 kg)

Payload range (Kg):

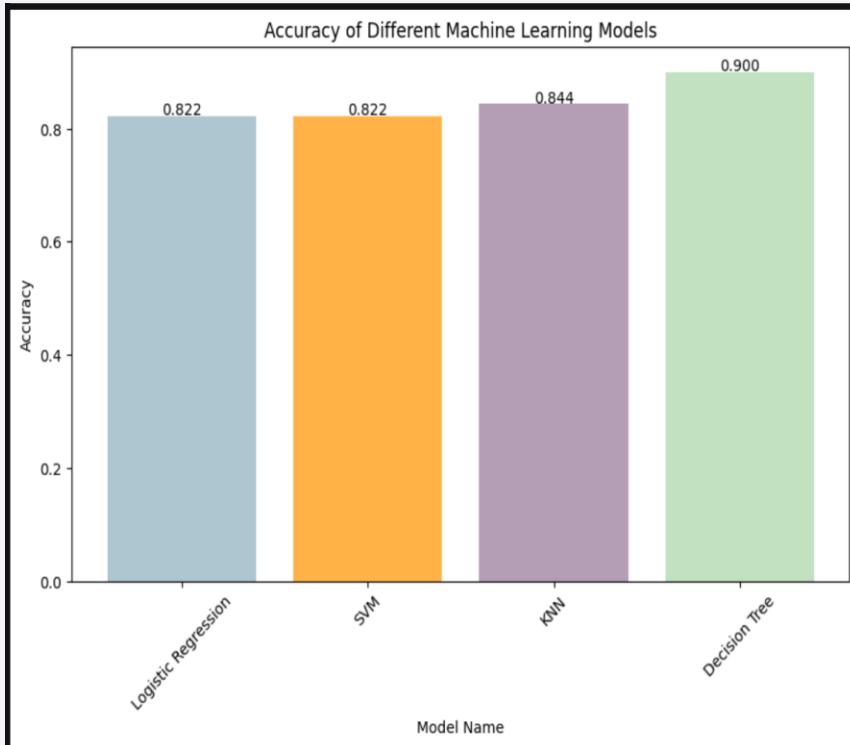


Section 5

Predictive Analysis (Classification)

Classification Accuracy

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



TASK 12

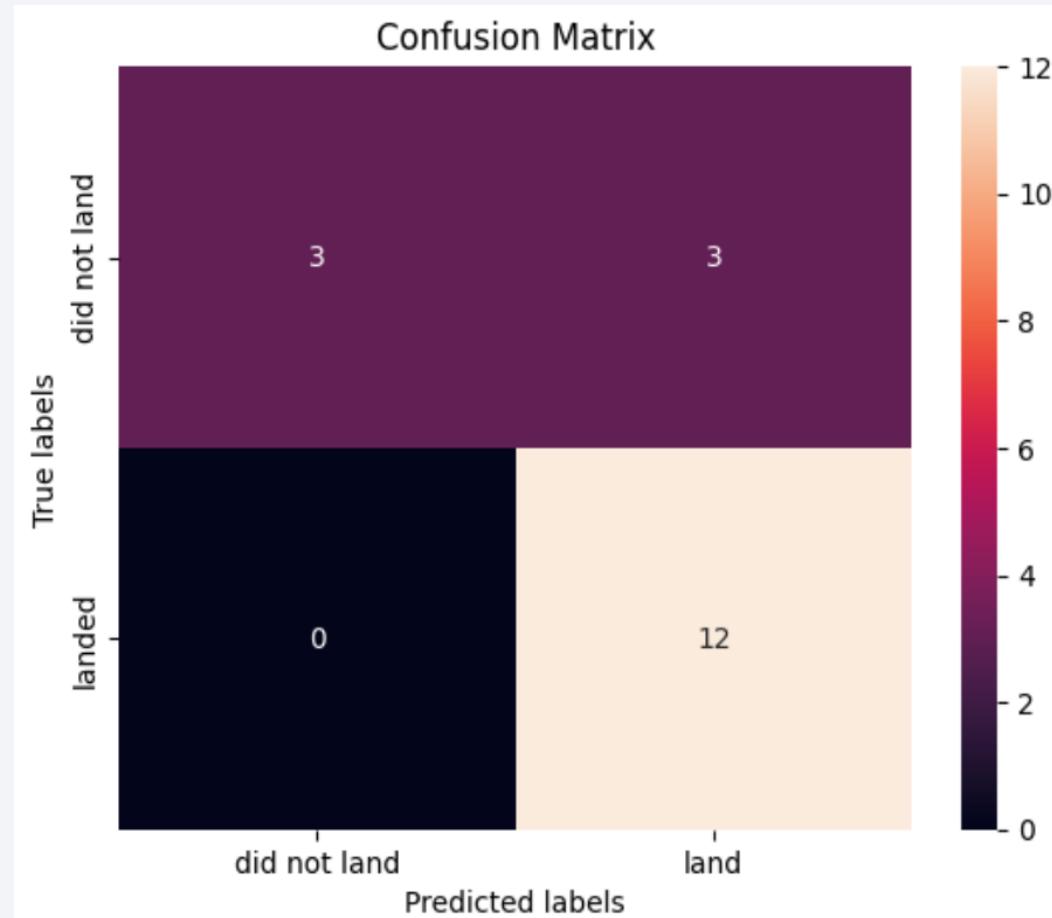
Find the method performs best:

```
[32]: models = {
    "Logistic Regression": logreg_cv.best_score_,
    "SVM": svm_cv.best_score_,
    "KNN": knn_cv.best_score_,
    "Decision Tree": tree_cv.best_score_
}
|
# Print best model
best_model = max(models, key=models.get)
print("Best performing model is:", best_model)
print("Accuracy:", models[best_model])
```

Best performing model is: Decision Tree
Accuracy: 0.9

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

It is 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.

Appendix

- Assets like Python code snippets, SQL queries, charts, Notebook outputs, or data sets that I have created during this project can be found in my github repository.
- The link of the repository is: <https://github.com/UFAQUE123/IBM-Data-Science-Capstone-Project>

Thank you!

