



IBM Developer
SKILLS NETWORK

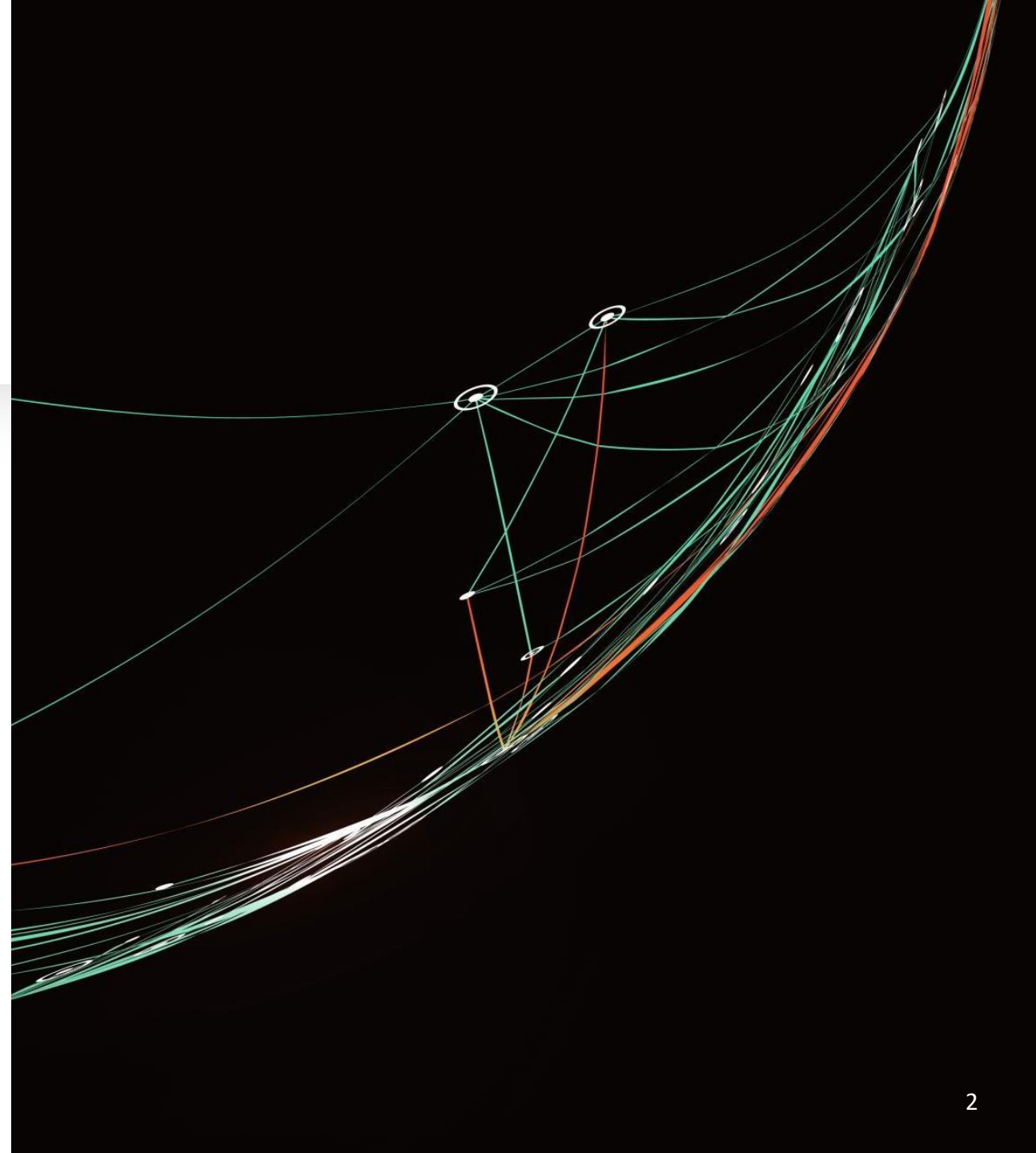
Winning Space Race with Data Science

Adithya Pagolu
23-06-2024



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

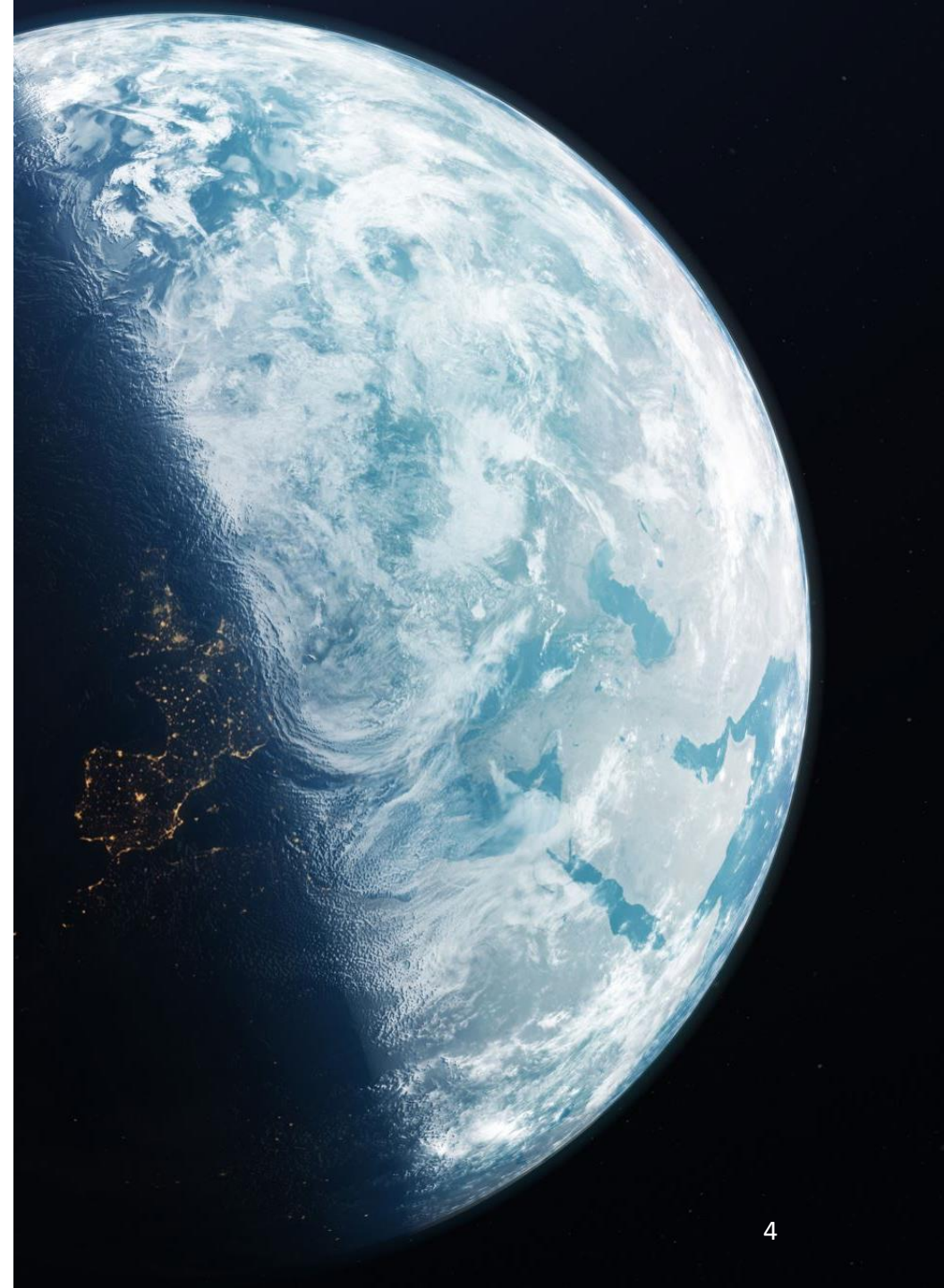
Introduction

- 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, 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

- To collect data for the project using SpaceX REST API calls, we begin with Authentication where the API key is used for secure access. Next, we proceed to Data Retrieval by sending GET requests to the SpaceX endpoints (e.g., /launchpads, /rockets) to fetch data on launches, rockets, payloads, etc. This data is then cleaned and used to extract relevant information and convert it into structured formats like JSON or CSV.

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

```
response = requests.get(spacex_url)
```

- GitHub URL of the completed SpaceX API calls notebook-[Adithya-Pagolu/Data-Science-and-Machine-Learning-Capstone-Project/Data Collection API.ipynb](#)

```
# Calculate the mean value of PayloadMass column  
payload_mass_mean = data_falcon9['PayloadMass'].mean()  
# Replace the np.nan values with its mean value  
data_falcon9['PayloadMass'].replace(np.nan, payload_mass_mean, inplace = True)
```


Data Collection - Scraping

- Here we applied web scrapping to webscrap Falcon 9 launch records with BeautifulSoup.

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

```
# Use soup.title attribute
soup.title
```

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

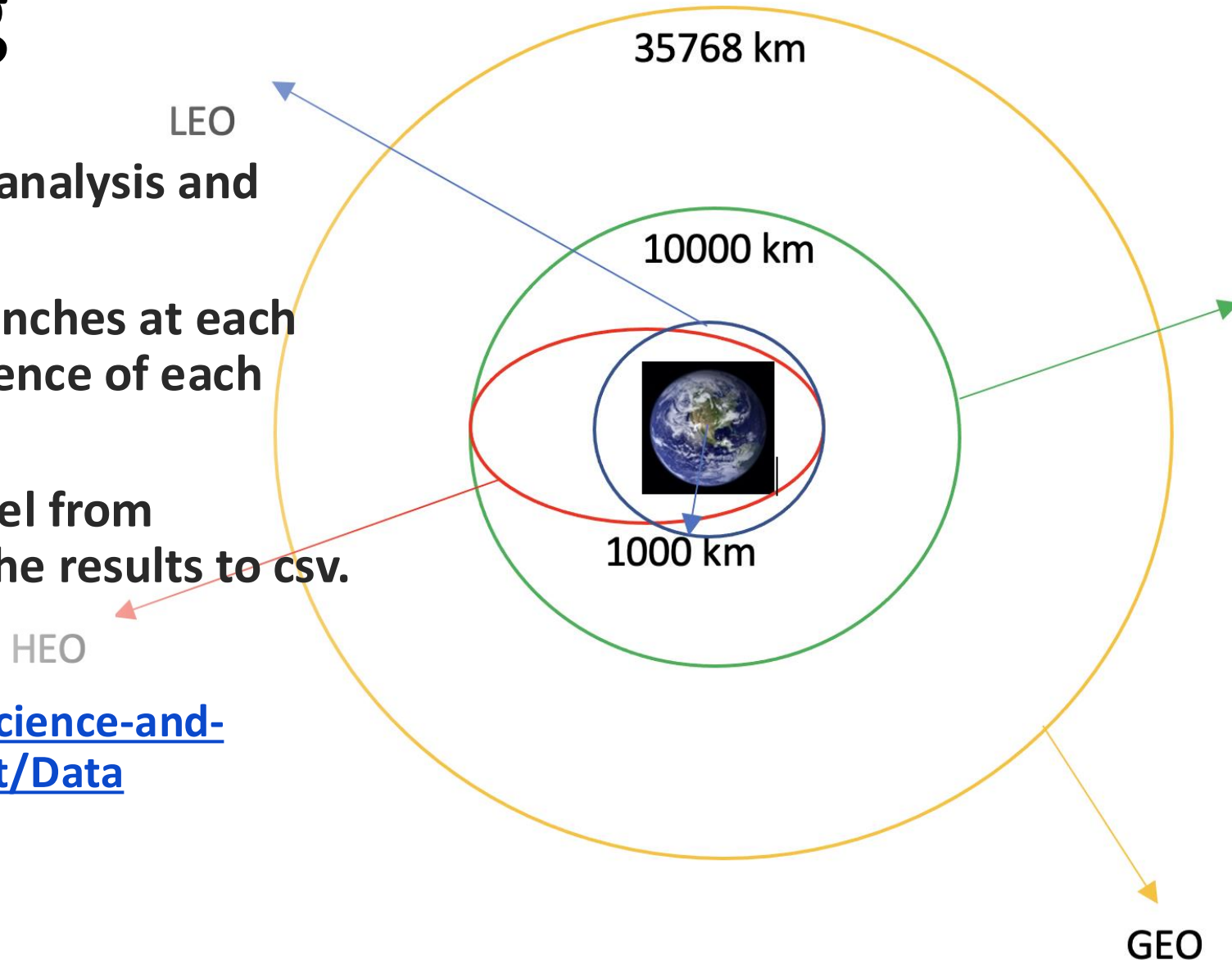
- We parsed the table and converted it into a pandas dataframe

```
df= pd.DataFrame({ key:pd.Series(value) for key, value in launch_dict.items() })
```

- GitHub URL-[Adithya-Pagolu/Data-Science-and-Machine-Learning-Capstone-Project/Data Collection with Web Scraping.ipynb](https://github.com/Adithya-Pagolu/Data-Science-and-Machine-Learning-Capstone-Project/Data-Collection-with-Web-Scraping.ipynb)

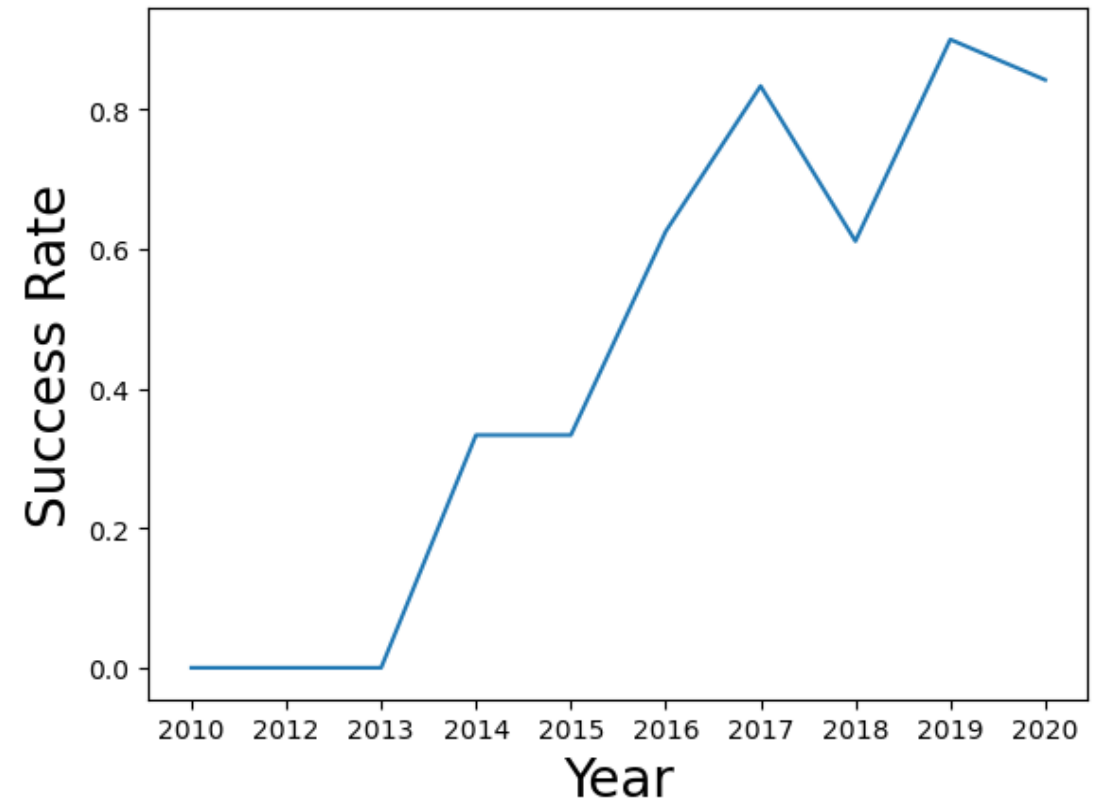
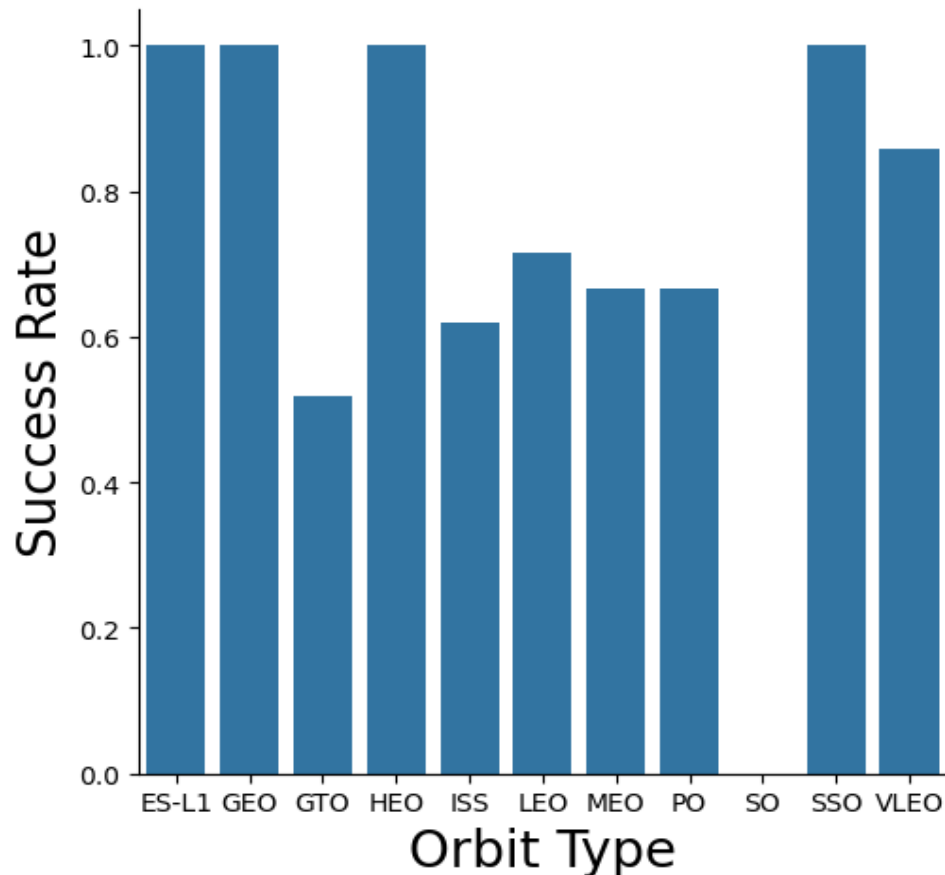
Data Wrangling

- 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.
- GitHub URL-[Adithya-Pagolu/Data-Science-and-Machine-Learning-Capstone-Project/DataWrangling.ipynb](https://github.com/Adithya-Pagolu/Data-Science-and-Machine-Learning-Capstone-Project/DataWrangling.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.
- GitHub URL-[Adithya-Pagolu/Data-Science-and-Machine-Learning-Capstone-Project/EDA with Visualization.ipynb](https://github.com/Adithya-Pagolu/Data-Science-and-Machine-Learning-Capstone-Project/EDA%20with%20Visualization.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.
- GitHub URL-[Adithya-Pagolu/Data-Science-and-Machine-Learning-Capstone-Project/EDA with SQL.ipynb](https://github.com/Adithya-Pagolu/Data-Science-and-Machine-Learning-Capstone-Project/EDA%20with%20SQL.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.
- GitHub URL-[Adithya-Pagolu/Data-Science-and-Machine-Learning-Capstone-Project/Interactive Visual Analytics with Folium.ipynb](https://github.com/Adithya-Pagolu/Data-Science-and-Machine-Learning-Capstone-Project/Interactive%20Visual%20Analytics%20with%20Folium.ipynb)

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.
- GitHub URL - [Adithya-Pagolu/Data-Science-and-Machine-Learning-Capstone-Project/spacex_dash_app.py](https://github.com/Adithya-Pagolu/Data-Science-and-Machine-Learning-Capstone-Project/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.
- GitHub URL-[Adithya-Pagolu/Data-Science-and-Machine-Learning-Capstone-Project/Machine Learning Prediction.ipynb](https://github.com/Adithya-Pagolu/Data-Science-and-Machine-Learning-Capstone-Project/blob/master/Machine%20Learning%20Prediction.ipynb)

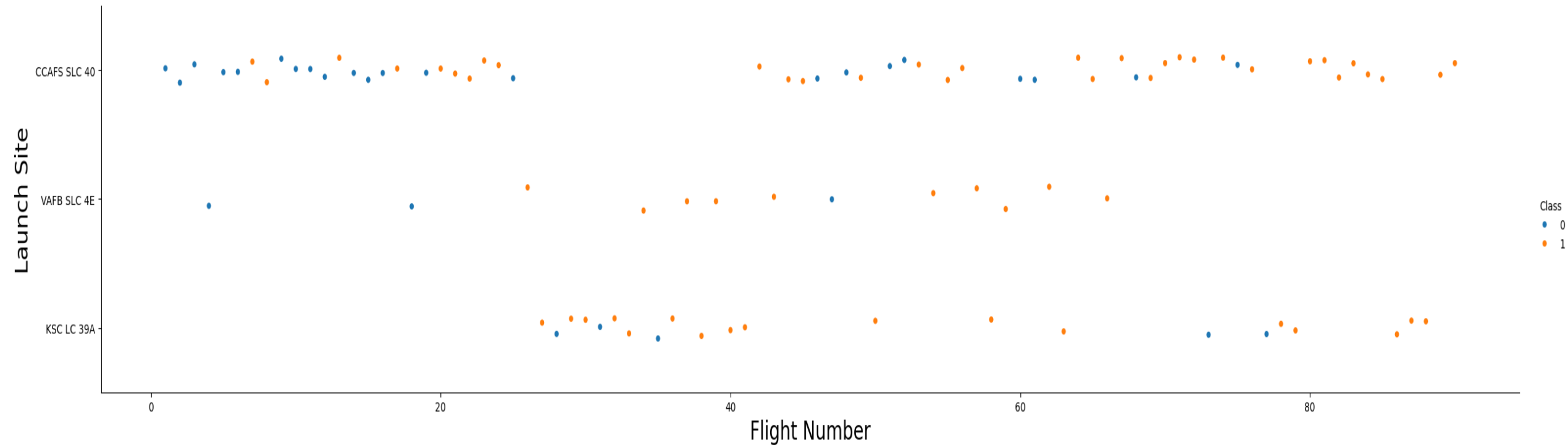
Results

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



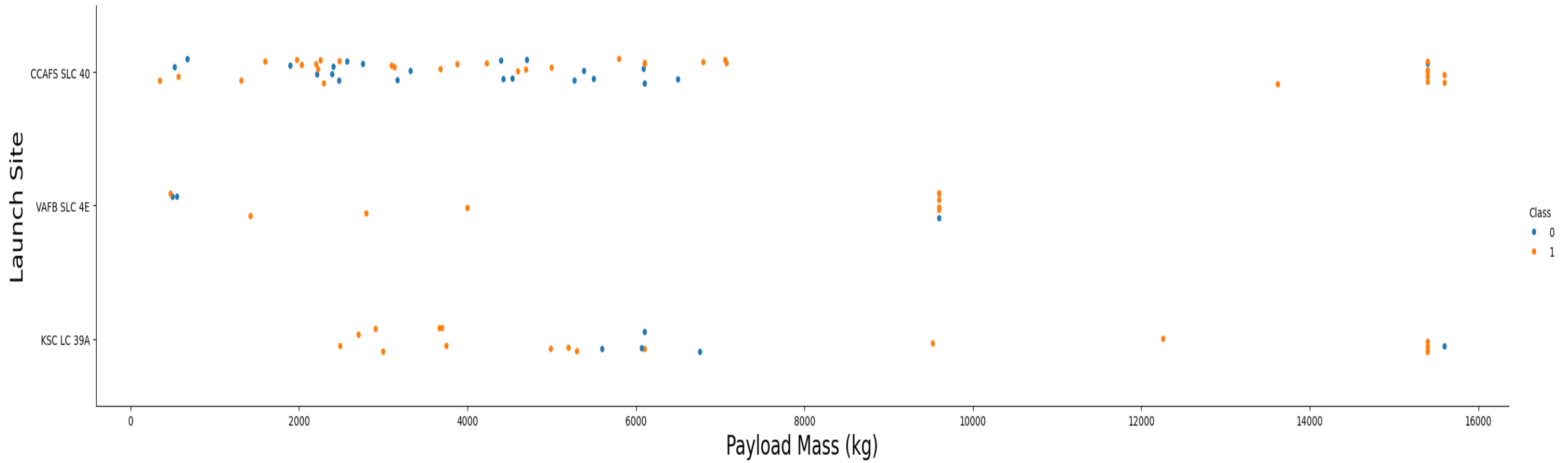
Section 2

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.

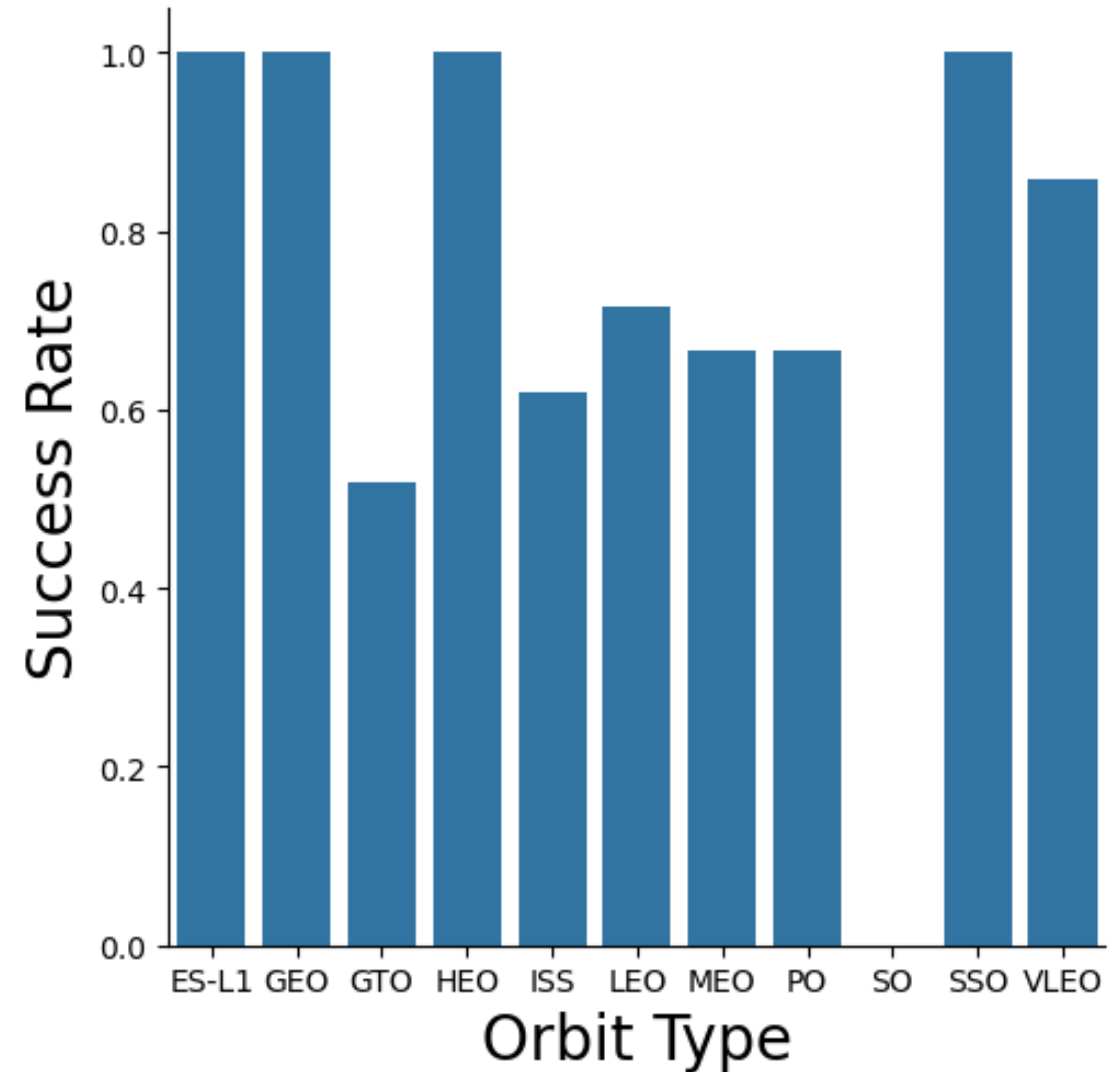


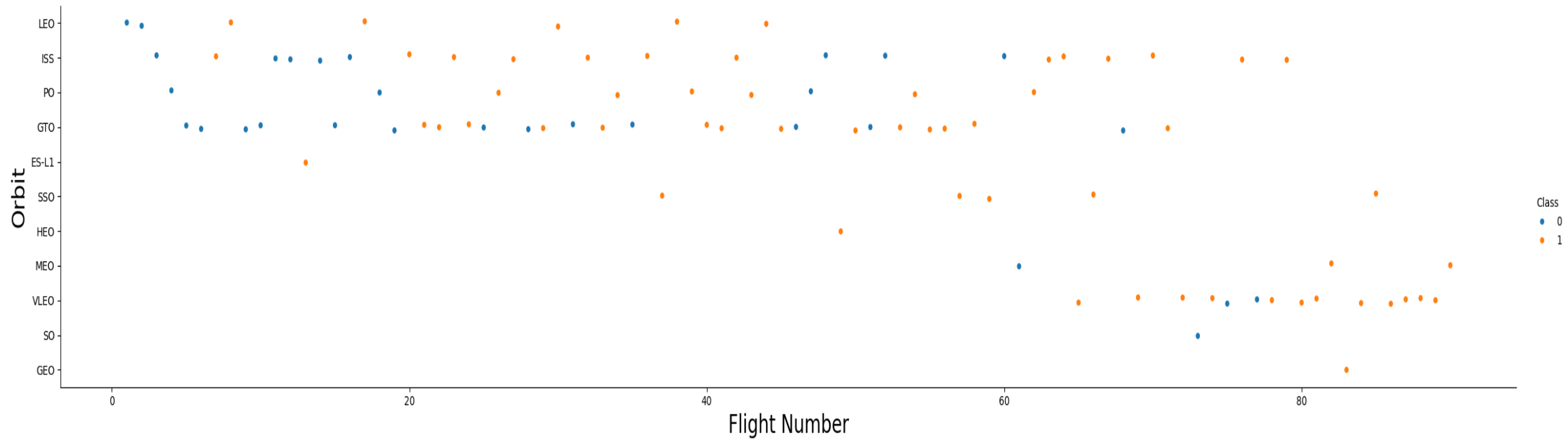
Payload vs. Launch Site

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

Success Rate vs. Orbit Type

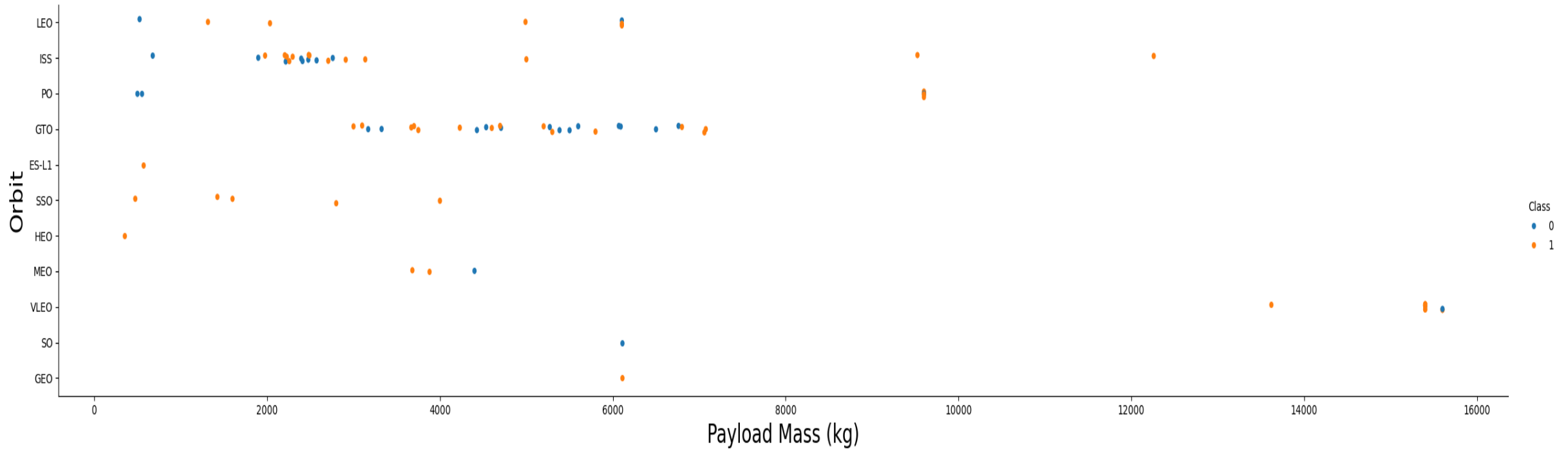
- From the plot, we can see that ES-L1, GEO, HEO, SSO, VLEO had the most success rate.





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.

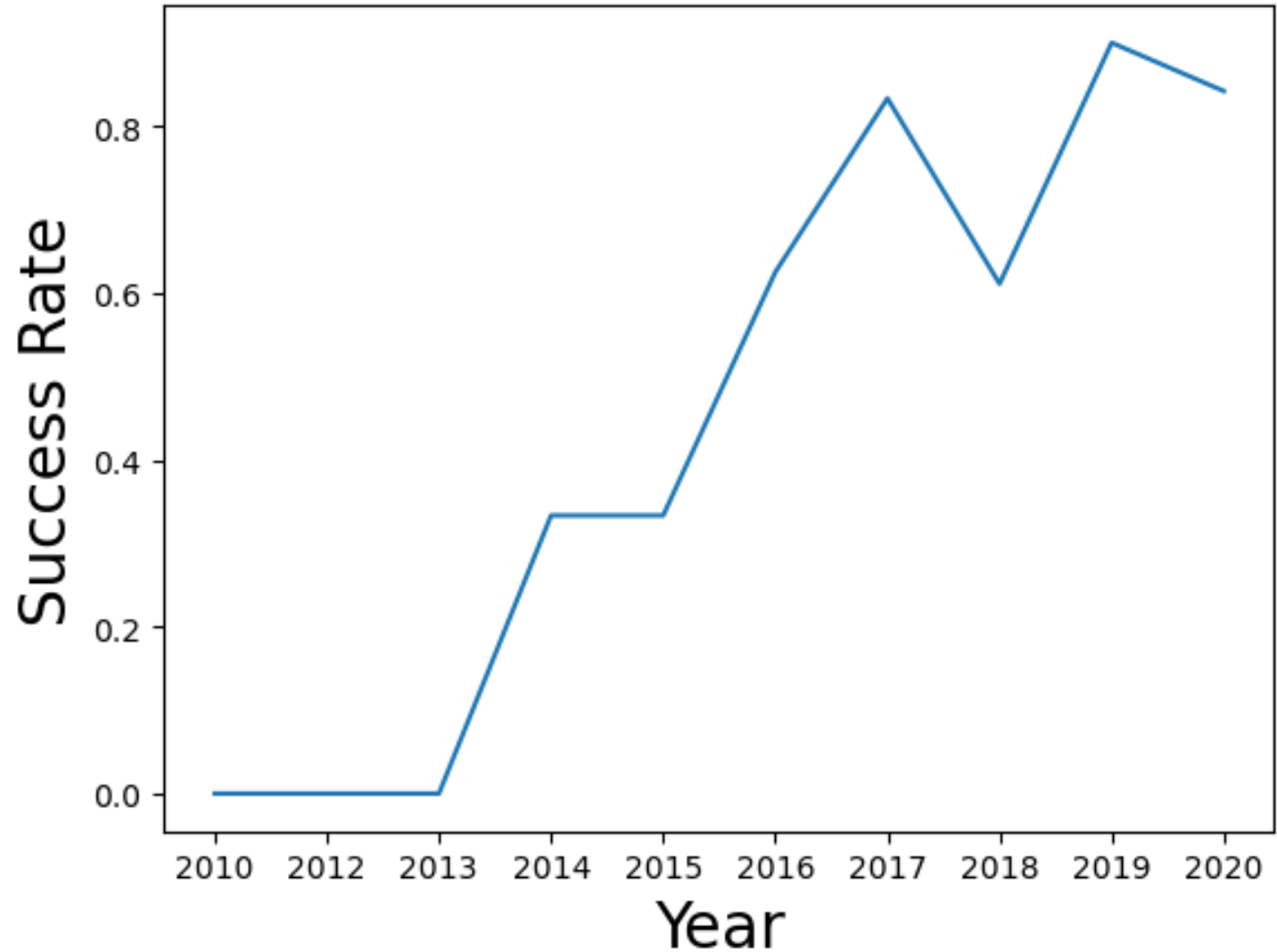


Payload vs. Orbit Type

- 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

- The SQL query `SELECT DISTINCT launch_site FROM SPACEXDATASET;` retrieves a unique list of all the different launch sites present in the SpaceX dataset. This query is useful for understanding the various locations from which SpaceX has conducted its launches.

```
%sql select distinct launch_site from SPACEXDATASET;
```

```
* ibm_db_sa://wzf08322:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:31198/bludb  
Done.
```

launch_site

CCAFS LC-40

CCAFS SLC-40

KSC LC-39A

VAFB SLC-4E

Launch Site Names Begin with 'KSC'

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

```
%sql select * from SPACEXDATASET where launch_site like 'CCA%' limit 5;
```

```
* ibm_db_sa://wzf08322:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:31198/bludb
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	07:44:00	F9 v1.0 B0005	CCAFS LC-40	Dragon demo flight C2	525	LEO (ISS)	NASA (COTS)	Success	No attempt
2012-10-08	00:35:00	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

- The SQL query `SELECT sum(payload_mass__kg_) AS total_payload_mass FROM SPACEXDATASET WHERE customer = 'NASA (CRS)';` is designed to calculate the total mass of payloads delivered by SpaceX for NASA's Commercial Resupply Services (CRS) missions.

```
%sql select sum(payload_mass__kg_) as total_payload_mass from SPACEXDATASET where customer = 'NASA (CRS)';
```

```
* ibm_db_sa://wzf08322:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:31198/bludb  
Done.
```

total_payload_mass

45596

Average Payload Mass by F9 v1.1

- The SQL query `SELECT AVG(payload_mass__kg_) AS average_payload_mass FROM SPACEXDATASET WHERE booster_version LIKE '%F9 v1.1%';` calculates the average payload mass for all records in the SPACEXDATASET table where the booster_version contains the substring 'F9 v1.1'.

```
%sql select avg(payload_mass__kg_) as average_payload_mass from SPACEXDATASET where booster_version like '%F9 v1.1%';
```

```
* ibm_db_sa://wzf08322:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90l08kqb1od8l1cg.databases.appdomain.cloud:31198/bludb
Done.
```

average_payload_mass
2534

First Successful Ground Landing Date

- The SQL query finds the earliest date when SpaceX achieved a successful landing on a ground pad.

```
%sql select min(date) as first_successful_landing from SPACEXDATASET where landing__outcome = 'Success (ground pad)';
```

```
* ibm_db_sa://wzf08322:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:31198/blddb  
Done.
```

first_successful_landing

2015-12-22

Successful Drone Ship Landing with Payload between 4000 and 6000

- The SQL query is designed to select the booster_version from a dataset named SPACEXDATASET where the landing outcome is 'Success (drone ship)' and the payload mass is between 4000 and 6000 kilograms.

**%sql select booster_version from SPACEXDATASET where
landing__outcome = 'Success (drone ship)' and payload_mass__kg_
between 4000 and 6000;**

booster_version

F9 FT B1022

F9 FT B1026

F9 FT B1021.2

F9 FT B1031.2

Total Number of Successful and Failure Mission Outcomes

- The SQL query is designed to aggregate the outcomes of SpaceX missions from the SPACEXDATASET table.

```
%sql select mission_outcome, count(*) as total_number from SPACEXDATASET group by mission_outcome;
```

```
* ibm_db_sa://wzf08322:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:31198/bludb  
Done.
```

mission_outcome	total_number
Failure (in flight)	1
Success	99
Success (payload status unclear)	1

Boosters Carried Maximum Payload

- The SQL query is designed to retrieve the booster_version from the SPACEXDATASET table where the payload_mass__kg_ is equal to the maximum payload_mass__kg_ in the entire dataset.

```
%sql select booster_version from SPACEXDATASET where payload_mass__kg_ = (select max(payload_mass__kg_) from SPACEXDATASET);
```

```
* ibm_db_sa://wzf08322:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:31198/bludb  
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

- SQL query is designed to retrieve data from the SPACEXDATASET table using SQL to analyze SpaceX missions in 2015 where the landing outcome was 'Failure (drone ship)'.

```
%%sql select monthname(date) as month, date, booster_version, launch_site, landing__outcome from SPACEXDATASET
       where landing__outcome = 'Failure (drone ship)' and year(date)=2015;
```

```
* ibm_db_sa://wzf08322:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90l08kqb1od8lcb.databases.appdomain.cloud:31198/blddb
Done.
```

MONTH	DATE	booster_version	launch_site	landing__outcome
January	2015-01-10	F9 v1.1 B1012	CCAFS LC-40	Failure (drone ship)
April	2015-04-14	F9 v1.1 B1015	CCAFS LC-40	Failure (drone ship)

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

- The SQL query is designed to count and display the occurrences of different landing outcomes from the SPACEXDATASET table within a specified date range.

```
%%sql select landing__outcome, count(*) as count_outcomes from SPACEXDATASET
      where date between '2010-06-04' and '2017-03-20'
      group by landing__outcome
      order by count_outcomes desc;
```

```
* ibm_db_sa://wzf08322:***@0c77d6f2-5da9-48a9-81f8-86b520b87518.bs2io90l08kqb1od8lcg.databases.appdomain.cloud:31198/bludb
Done.
```

landing__outcome	count_outcomes
No attempt	10
Failure (drone ship)	5
Success (drone ship)	5
Controlled (ocean)	3
Success (ground pad)	3
Failure (parachute)	2
Uncontrolled (ocean)	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 background on the left and a satellite photograph of Earth on the right. The Earth's surface is dark, with numerous bright yellow and orange lights representing cities and urban areas. The horizon of the Earth is visible as a curved line separating the dark surface from the deep blue of space.

Section 3

Launch Sites Proximities Analysis

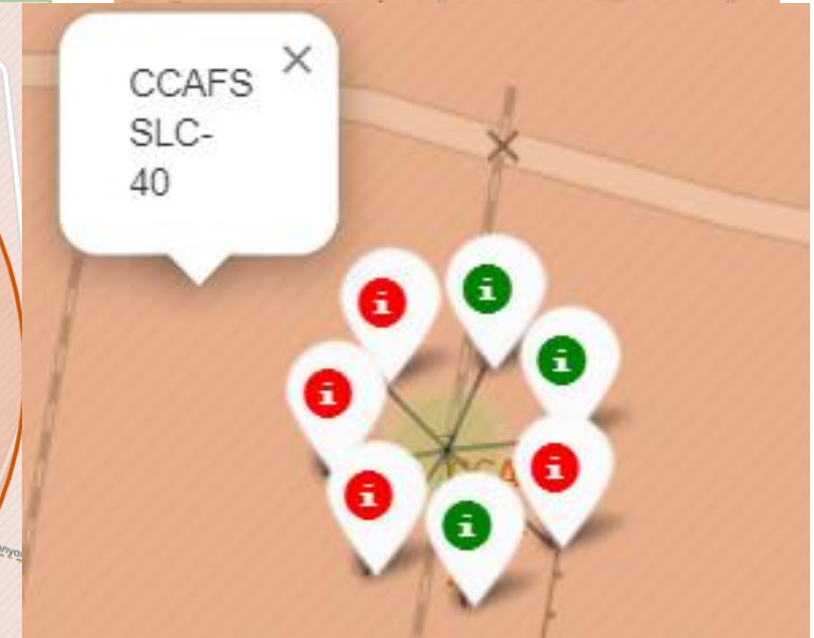
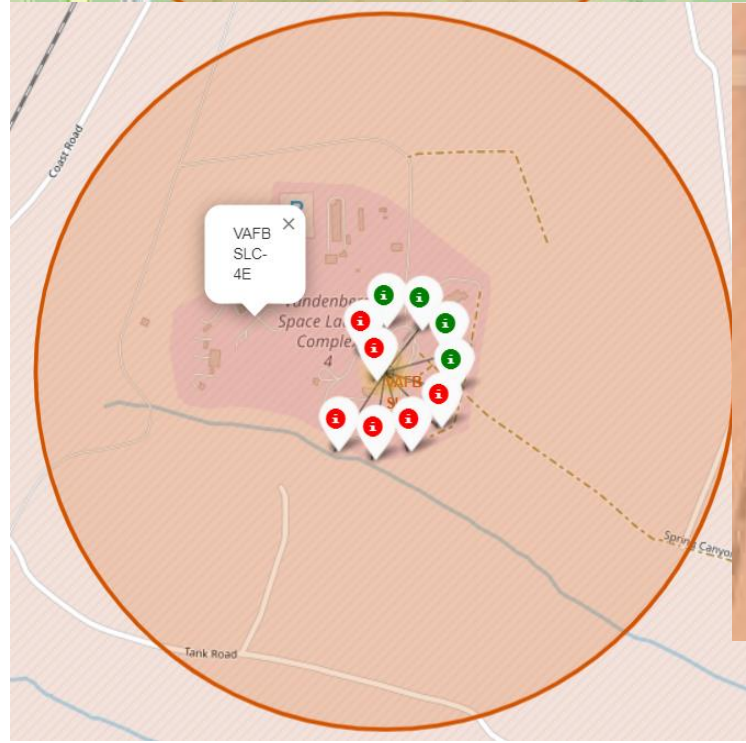
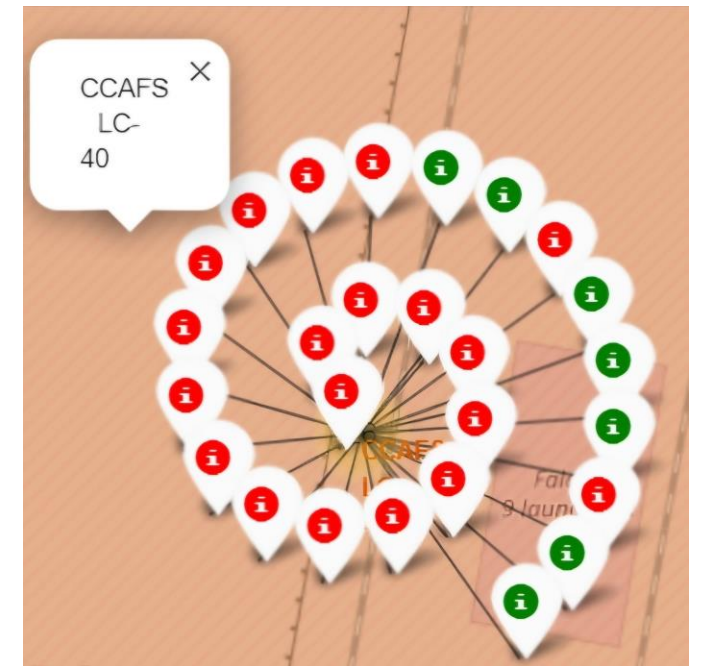
All launch sites global map markers



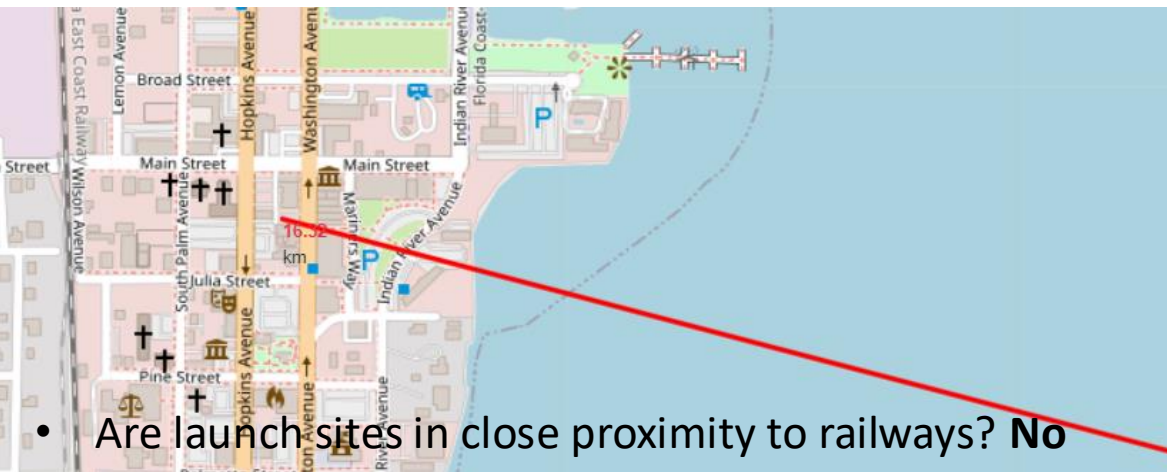
We can see that SpaceX launch sites are in the coasts of the United States of America.

Markers showing launch sites with color labels

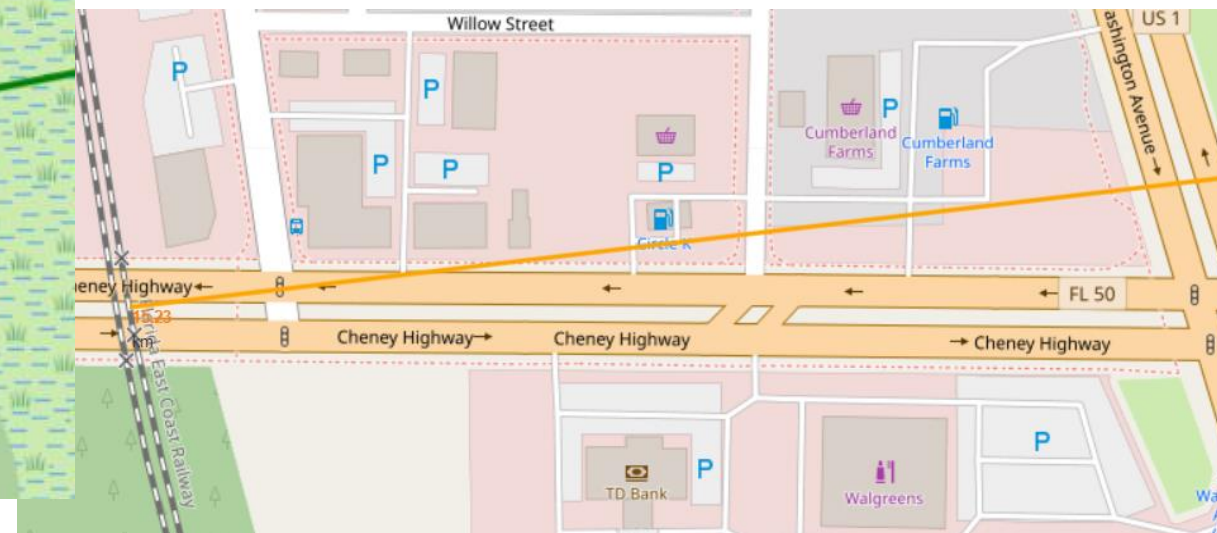
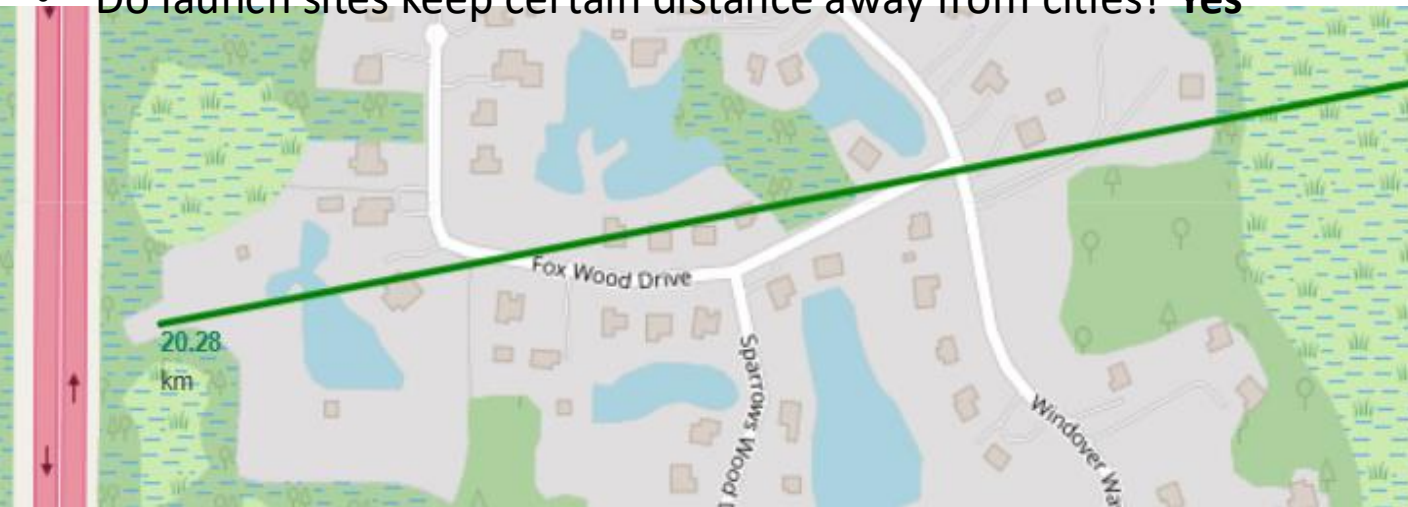
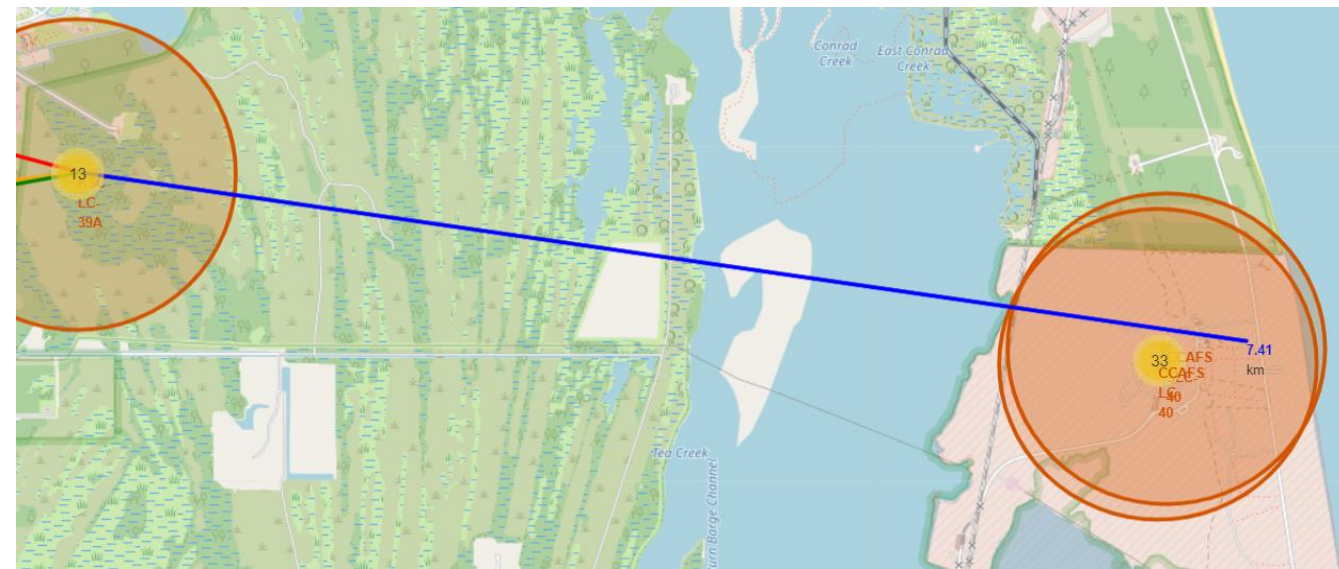
Green marker shows successful launches and Red marker shows failure launches



Launch Site distance to landmarks



- Are launch sites in close proximity to railways? **No**
- Are launch sites in close proximity to highways? **No**
- Are launch sites in close proximity to coastline? **Yes**
- Do launch sites keep certain distance away from cities? **Yes**





Section 4

Build a Dashboard with Plotly Dash

Pie chart showing the success percentage achieved by each launch site

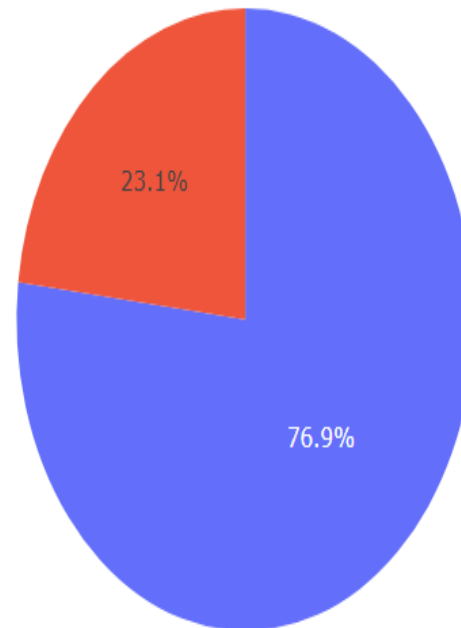


Pie chart showing the Launch site with the highest launch success ratio

KSC LC-39A

X ▼

Total Success Launches for Site KSC LC-39A



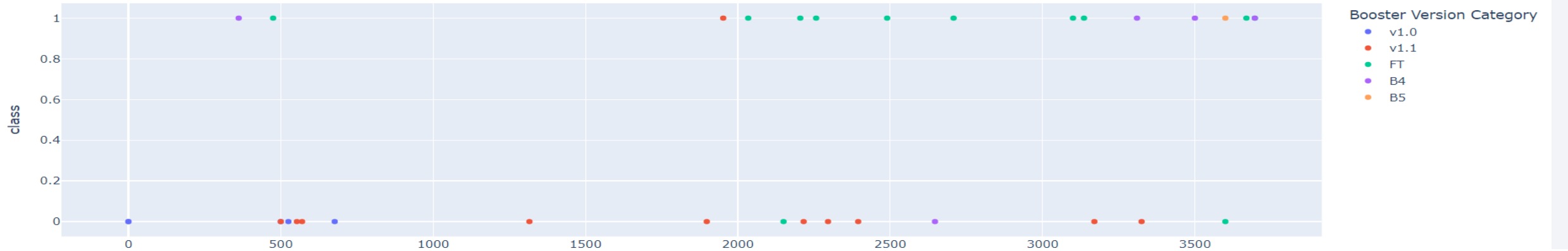
KSC LC-39A got 76.9% success rate while experiencing 23.1% failure rate.

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

Payload range (Kg):



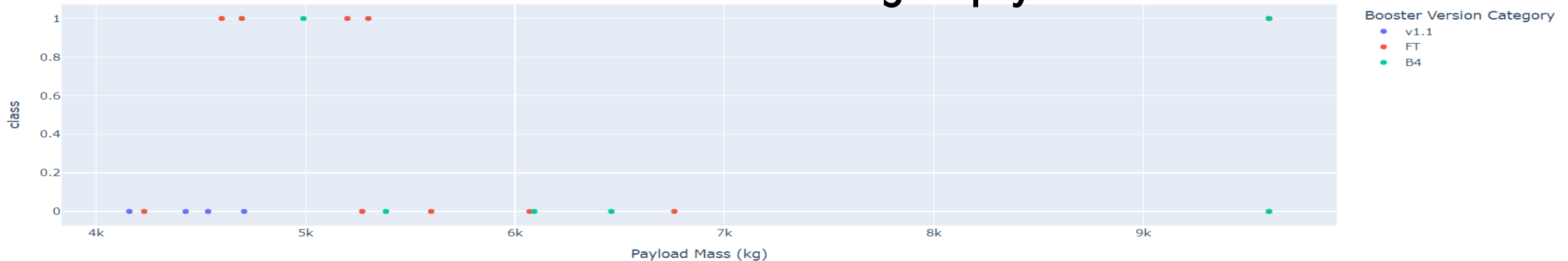
Correlation Between Payload and Success for All Sites



Payload range (Kg):



Correlation Between Payload and Success for All Sites



We can see that the success rate in lower payload is more than in higher payload.



Section 5

Predictive Analysis (Classification)

Classification Accuracy

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

```
models = {'KNeighbors': knn_cv.best_score_,
          'DecisionTree': tree_cv.best_score_,
          'LogisticRegression': logreg_cv.best_score_,
          'SupportVector': svm_cv.best_score_}

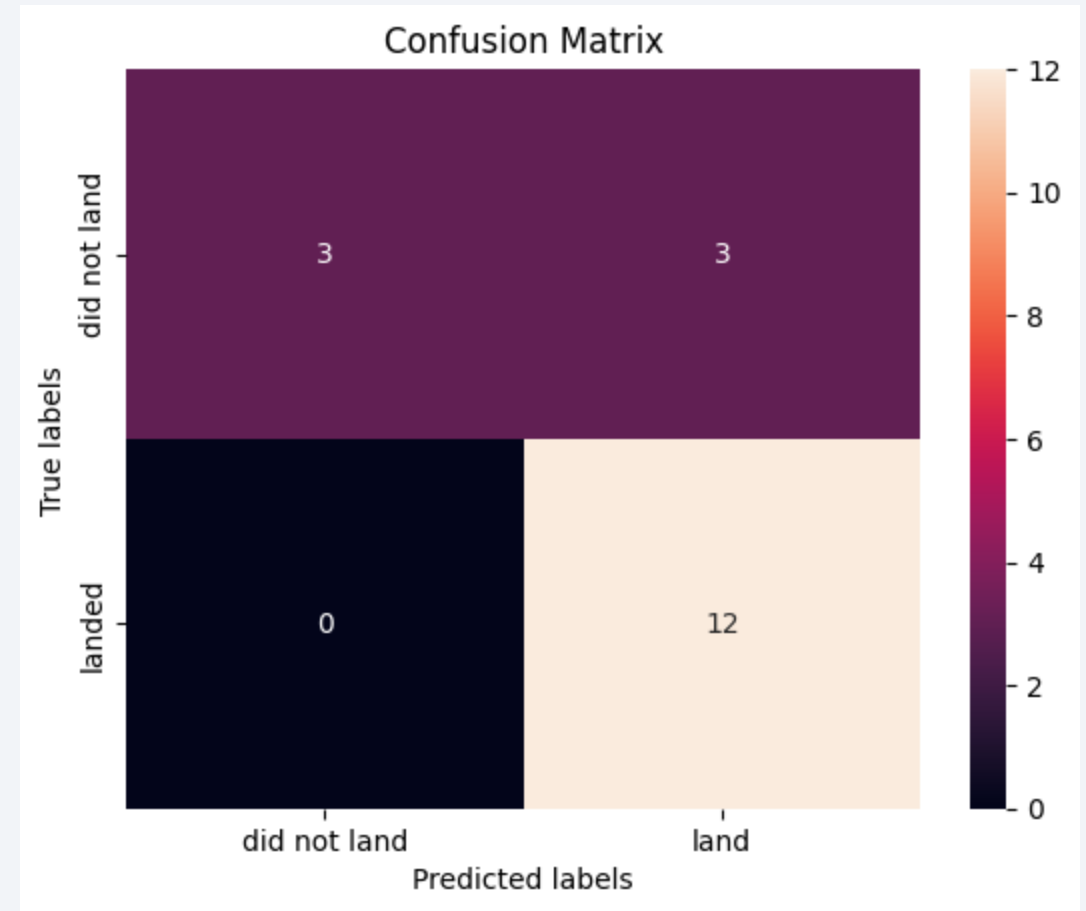
bestalgorithm = max(models, key=models.get)
print('Best model is', bestalgorithm, 'with a score of', models[bestalgorithm])
if bestalgorithm == 'DecisionTree':
    print('Best params is :', tree_cv.best_params_)
if bestalgorithm == 'KNeighbors':
    print('Best params is :', knn_cv.best_params_)
if bestalgorithm == 'LogisticRegression':
    print('Best params is :', logreg_cv.best_params_)
if bestalgorithm == 'SupportVector':
    print('Best params is :', svm_cv.best_params_)
```

Best model is DecisionTree with a score of 0.8732142857142856

Best params is : {'criterion': 'gini', 'max_depth': 6, 'max_features': 'sqrt', 'min_samples_leaf': 2, 'min_samples_split': 5, 'splitter': 'random'}

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.

Appendix

- SpaceX dataset- [Adithya-Pagolu/Data-Science-and-Machine-Learning-Capstone-Project/Spacex.csv](#)
- Dataset part-1- [Adithya-Pagolu/Data-Science-and-Machine-Learning-Capstone-Project/dataset part 1.csv](#)
- Dataset part-2- [Adithya-Pagolu/Data-Science-and-Machine-Learning-Capstone-Project/dataset part 2.csv](#)
- Dataset part-3- [Adithya-Pagolu/Data-Science-and-Machine-Learning-Capstone-Project/dataset part 3.csv](#)
- SpaceX launch dataset- [Adithya-Pagolu/Data-Science-and-Machine-Learning-Capstone-Project/spacex launch dash.csv](#)
- Project- [Adithya-Pagolu/Data-Science-and-Machine-Learning-Capstone-Project](#)

Thank you!

