**IBM Developer**
SKILLS NETWORK

# Winning Space Race with Data Science

Tanmayi Allu
1/20/2026

# 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.

# 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

- 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/AlluTanmayi-7/Applied-Data-Science-Capstone/blob/main/Mod1-jupyter-labs-spacex-data-collection-api.ipynb

**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:

```
In [9]:    static_json_url='https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DS0321EN-SkillsNetwork/datasets/API_
```

We should see that the request was successfull with the 200 status response code

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

```
In [11]:   response.status_code
```

```
Out[11]:   200
```

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

```
In [12]:   # 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

# Data Collection-Webscraping

- 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/AlluTanmayi-7/Applied-Data-Science-Capstone/blob/main/Mod1-jupyter-labs-webscraping.ipynb

## 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.

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

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

```
In [7]:   # Use BeautifulSoup() to create a BeautifulSoup object from a response text content
          soup = BeautifulSoup(response.text, 'html.parser')
```

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

```
In [8]:   # Use soup.title attribute
          print(soup.title)
```

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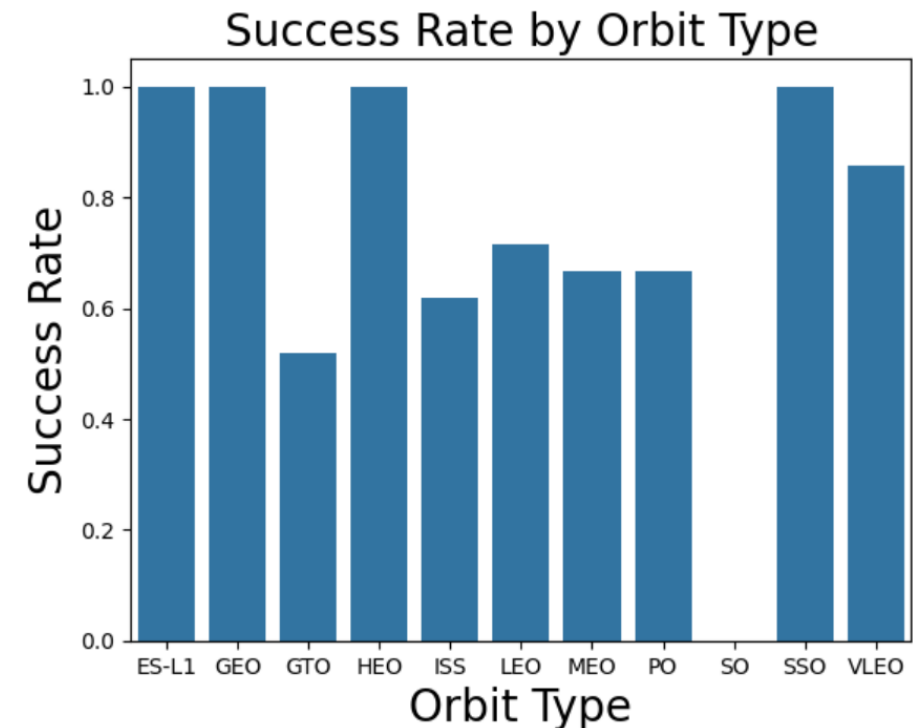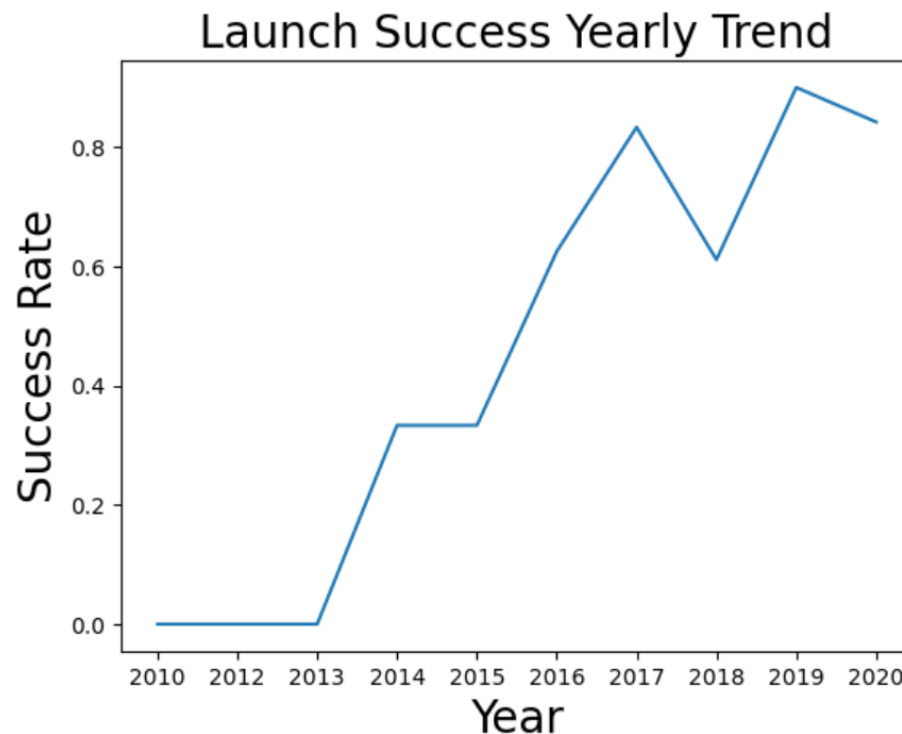
# Data Wrangling



- 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/AlluTanmayi-7/Applied-Data-Science-Capstone/blob/main/Mod1-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/AlluTanmayi-7/Applied-Data-Science-Capstone/blob/main/edadataviz.ipynb



Launch Success Yearly Trend



Success Rate by Orbit Type

# 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/AlluTanmayi-7/Applied-Data-Science-Capstone/blob/main/Mod2-jupyter-labs-eda-sql-coursera_sqllite.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/AlluTanmayi-7/Applied-Data-Science-Capstone/blob/main/Mod3-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/AlluTanmayi-7/Applied-Data-Science-Capstone/blob/main/Mod4-SpaceX_Machine%20Learning%20Prediction_Part_5.ipynb

# Results

- Exploratory data analysis results

- Interactive analytics demo in screenshots

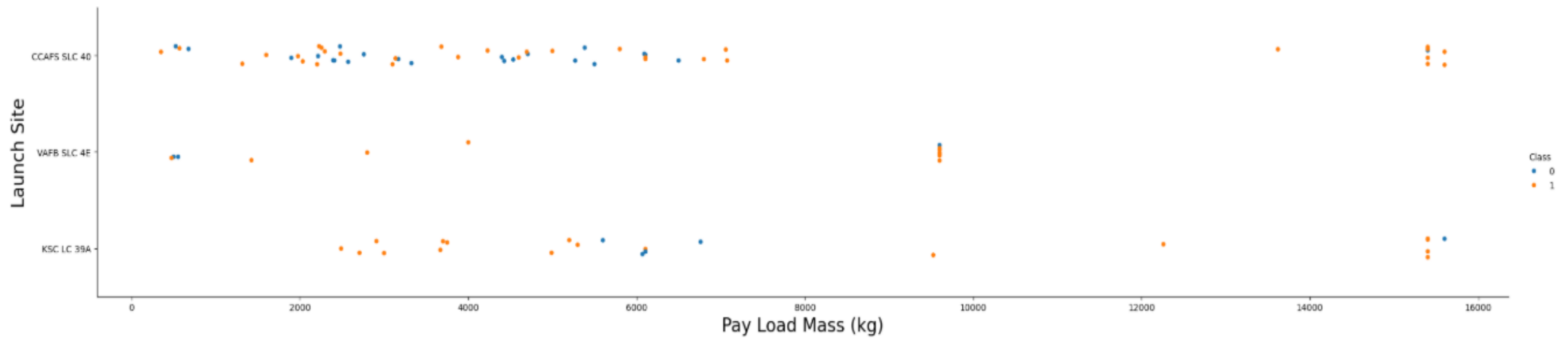- Predictive analysis results

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.
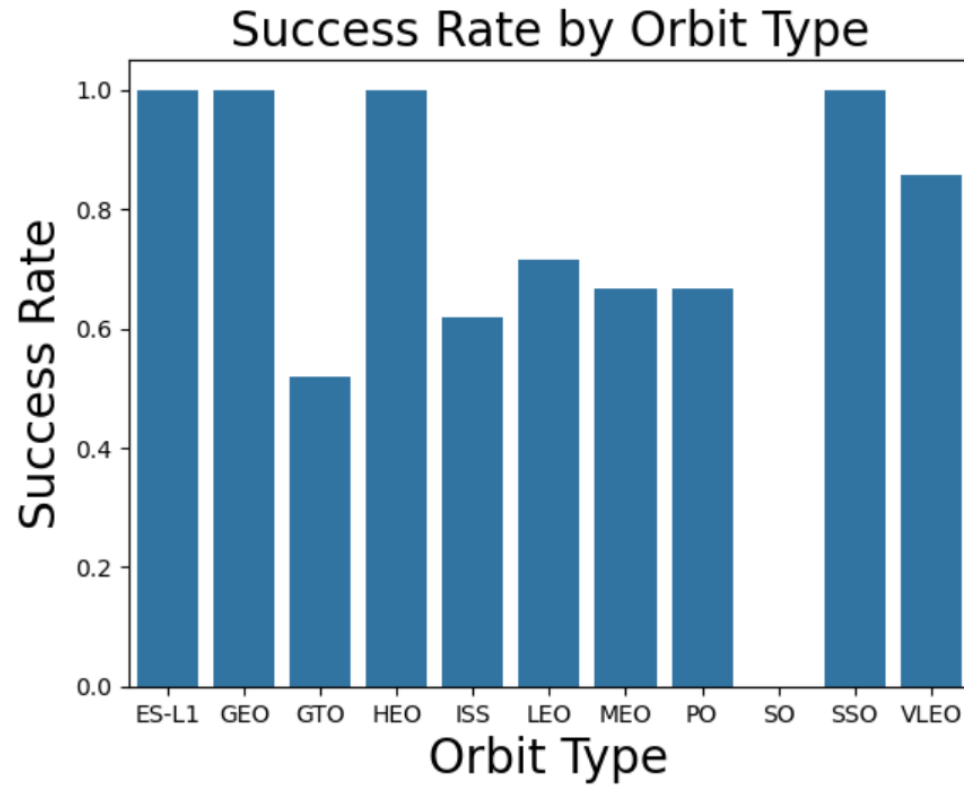
# Payload vs. Launch Site



The greater the payload was for launch site CCAFS SLC 40 the higher the succes rate for the rocket.
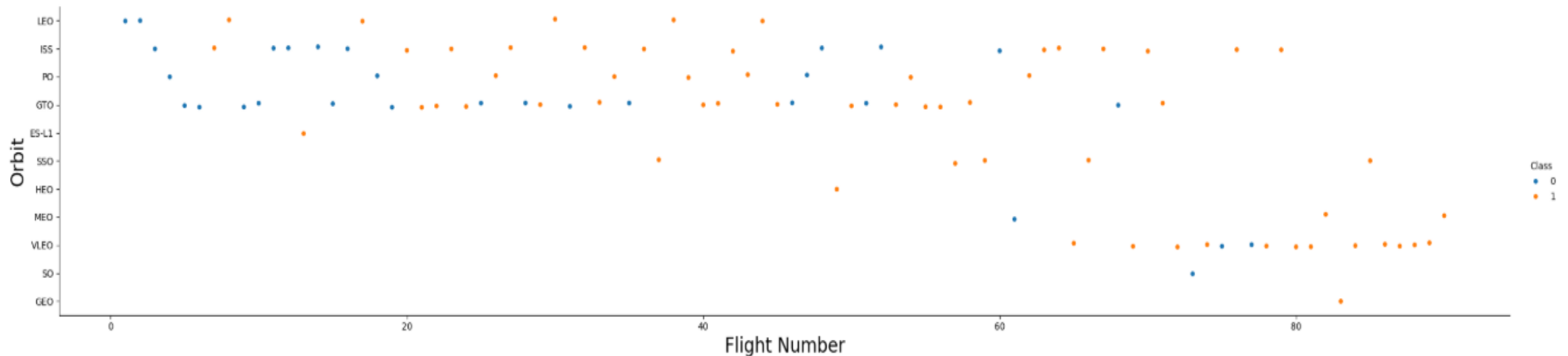
# Success Rate vs. Orbit Type



Success Rate by Orbit Type

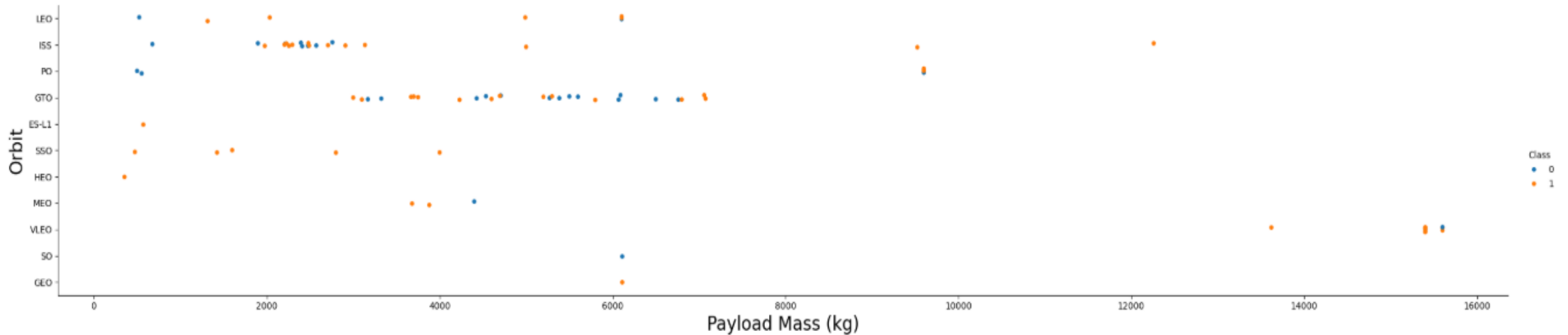- From the plot, we can see that ES-L1, GEO, HEO, SSO, VLEO had the most success rate.

# 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.
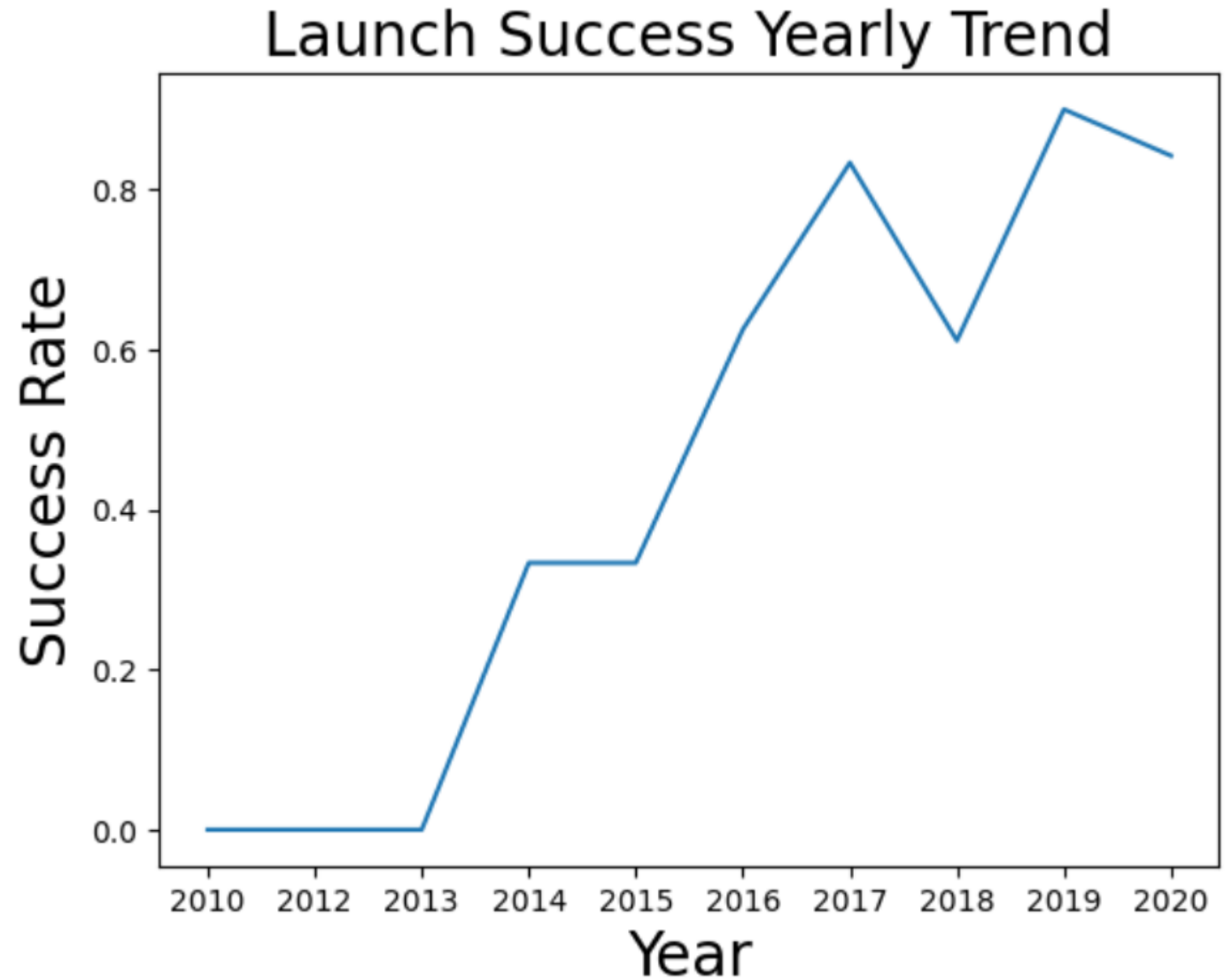
# Payload vs. Orbit Type

- 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.



Launch Success Yearly Trend

# First Successful Ground Landing Date

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

## Task 5

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

*Hint:Use min function*

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

* sqlite:///my_data1.db
Done.

**MIN("Date")**

2015-12-22

# Total Number of Successful and Failure Mission Outcomes

- This SQL query is designed to summarize the **Mission_Outcome** column from the SPACEXTABLE by counting how many times each specific outcome occurred.

## Task 7

List the total number of successful and failure mission outcomes

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

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

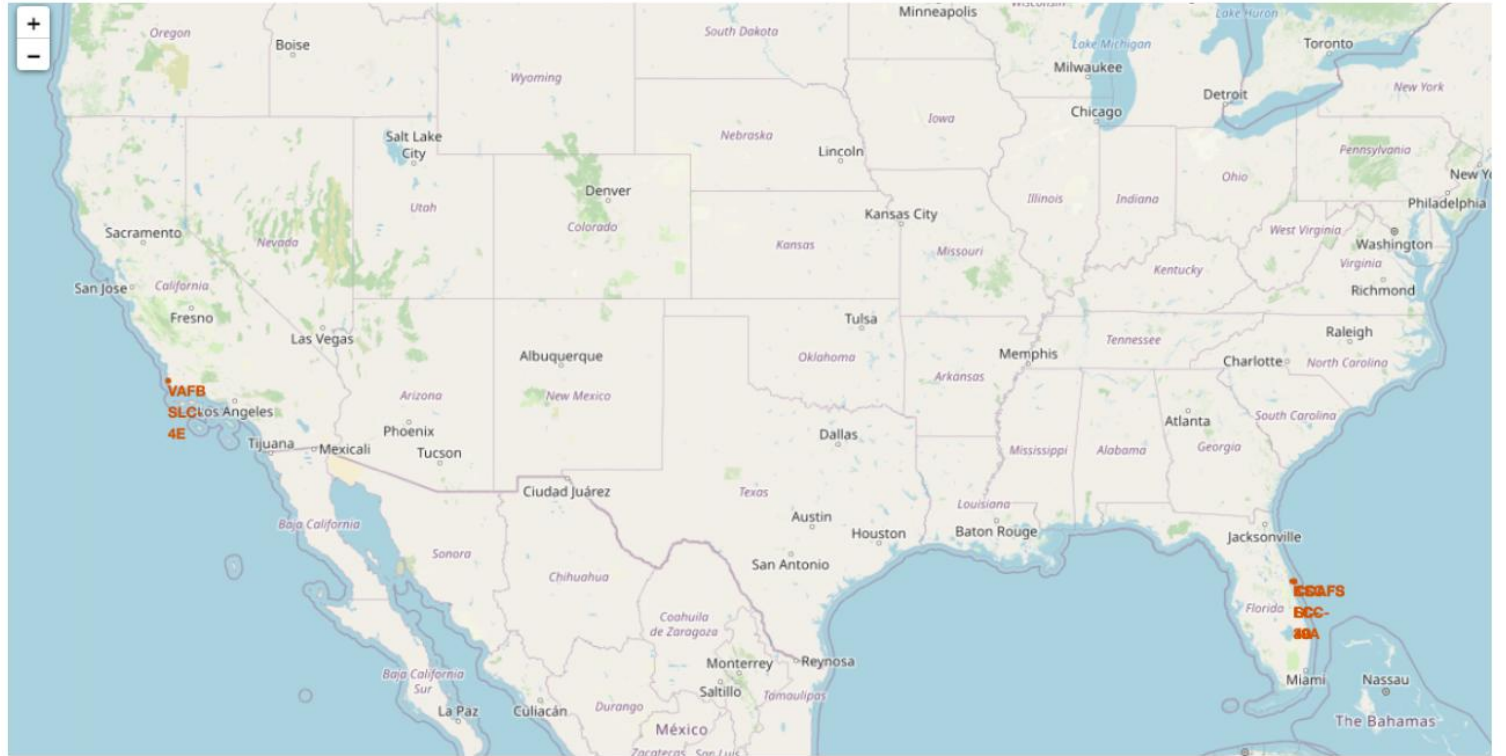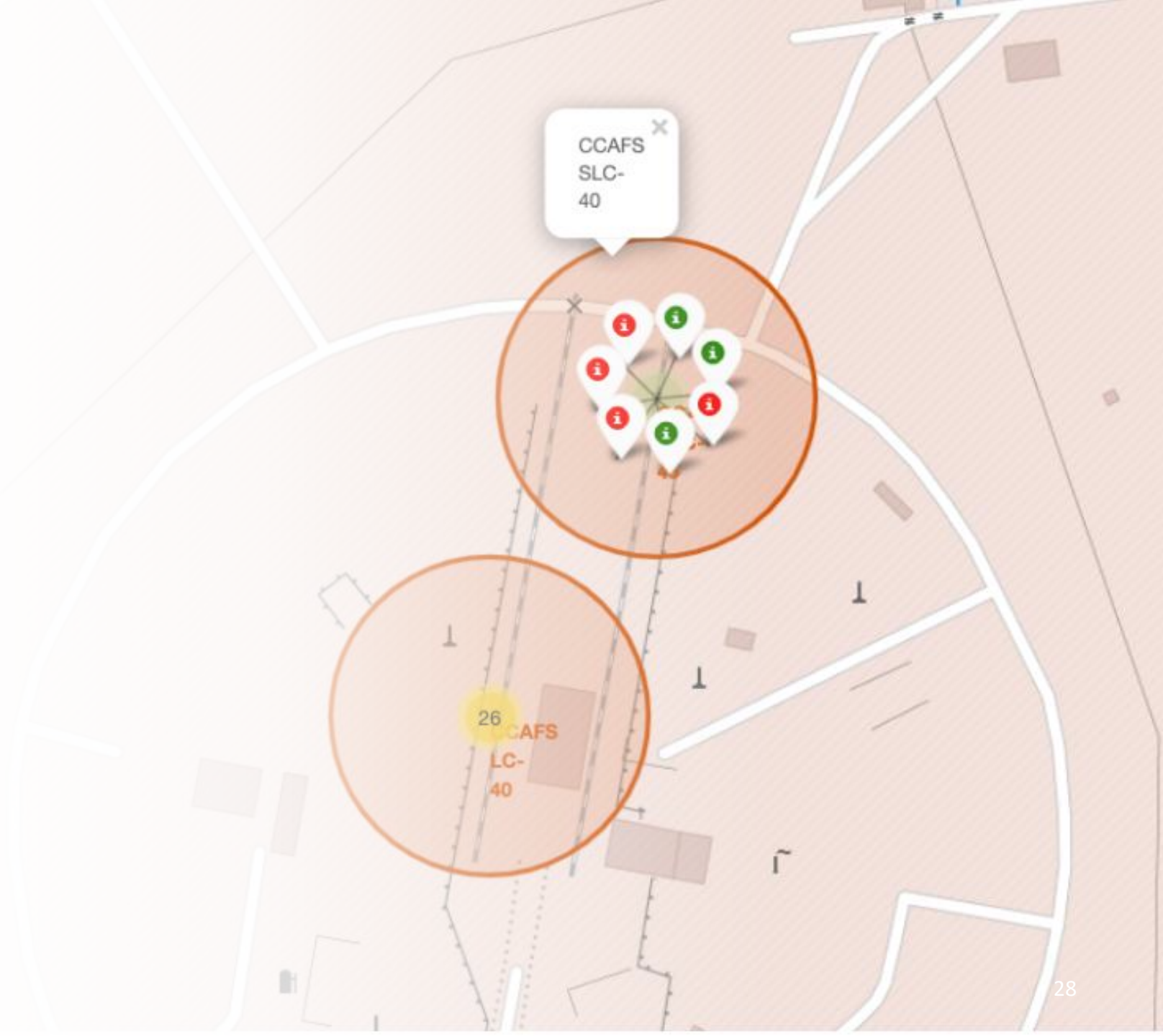| Mission_Outcome | TotalCount |
|---|---|
| Failure (in flight) | 1 |
| Success | 98 |
| Success | 1 |
| Success (payload status unclear) | 1 |

Section 4

# Launch Sites
# Proximities Analysis

# All launch sites global map markers

# Markers showing launch sites with color labels

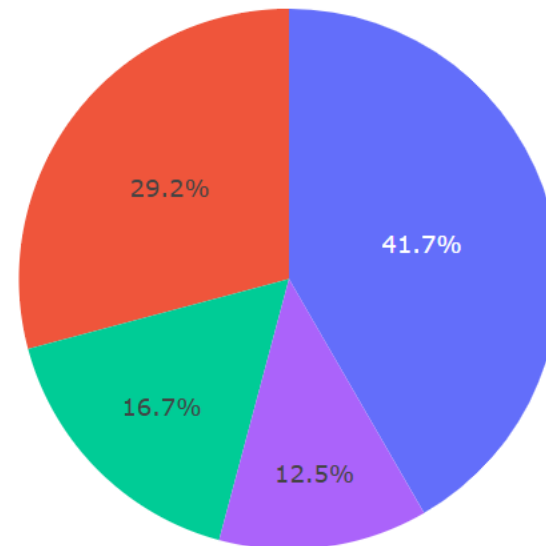# Build a Dashboard
# with Plotly Dash

# Pie chart showing the success percentage achieved by each launch site



## SpaceX Launch Records Dashboard

All Sites

### Total Success Launches By Site

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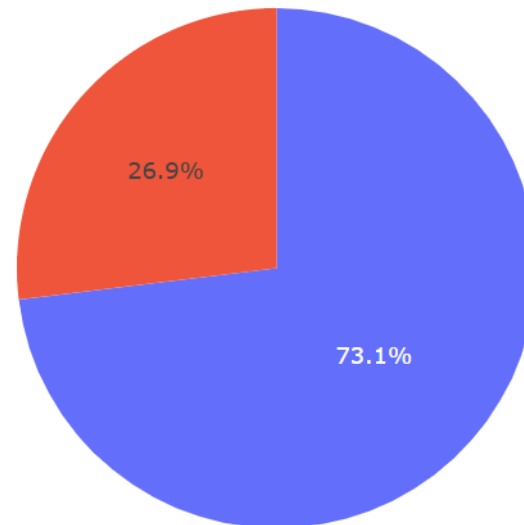- CCAFS LC-40
- VAFB SLC-4E
- CCAFS SLC-40

41.7%
29.2%
16.7%
12.5%

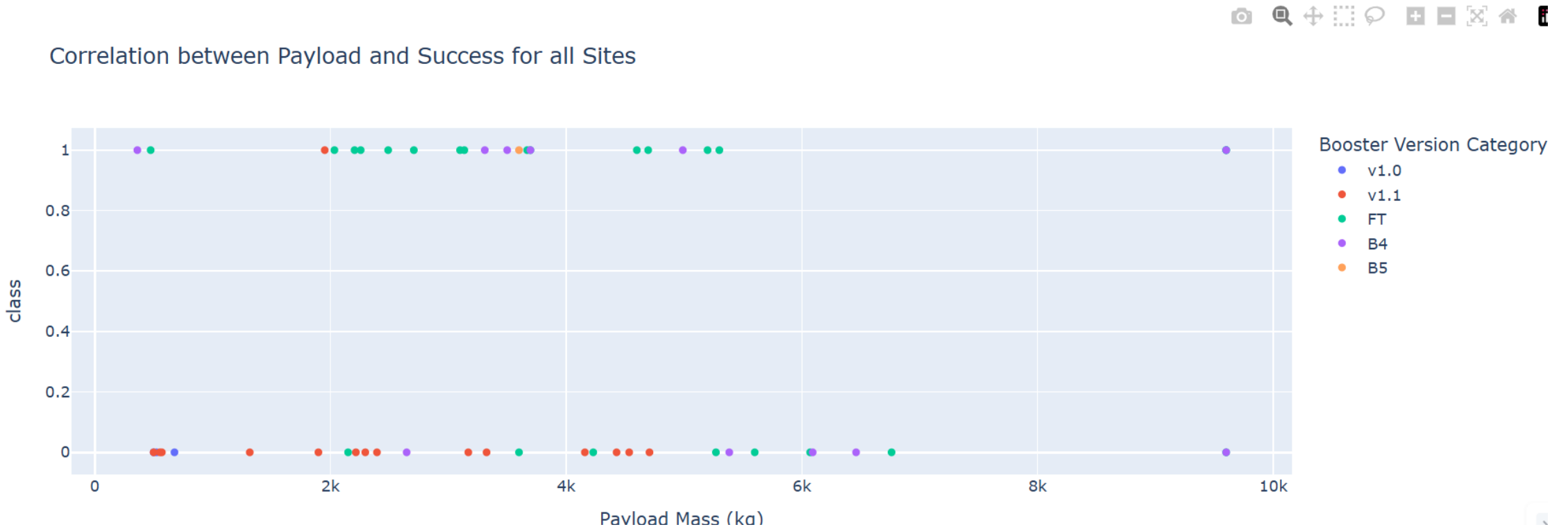# Pie chart showing the Launch site CCAFS LC-40

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



Correlation between Payload and Success for all Sites

Section 6

# Predictive Analysis (Classification)

# Classification Accuracy

- The Logistic Regression is the model with the highest classification accuracy

```python
accuracy_scores = {
    'Model': ['Logistic Regression', 'SVM', 'Decision Tree', 'KNN'],
    'Test Accuracy': [
        logreg_cv.score(X_test, Y_test),
        svm_cv.score(X_test, Y_test),
        tree_cv.score(X_test, Y_test),
        knn_cv.score(X_test, Y_test)
    ]
}
df_accuracy = pd.DataFrame(accuracy_scores)
print(df_accuracy)
```
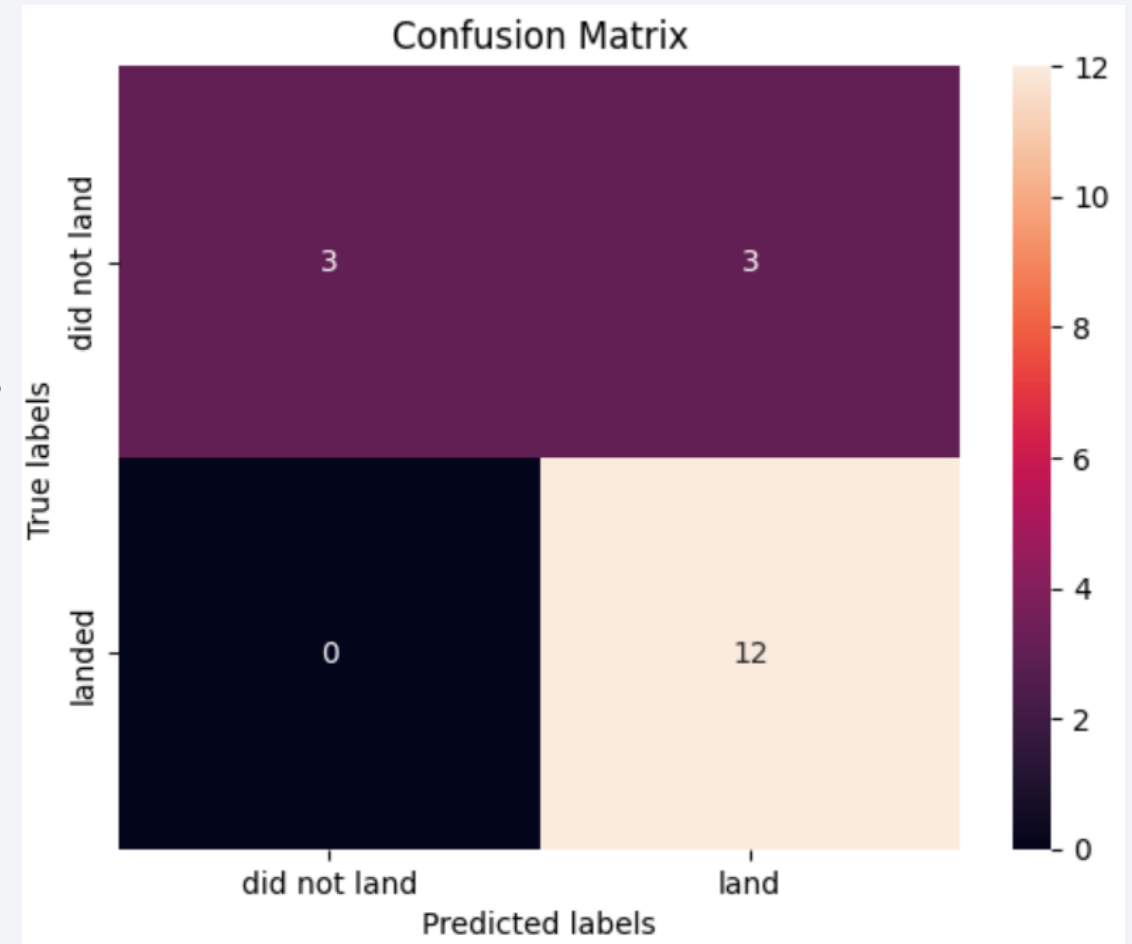
```
                 Model  Test Accuracy
0  Logistic Regression       0.833333
1                  SVM       0.833333
2        Decision Tree       0.777778
3                  KNN       0.833333
```

```python
best_model = df_accuracy.loc[df_accuracy['Test Accuracy'].idxmax()]
print(f"The best performing model is: {best_model['Model']} with an accuracy of {best_model['Test Accuracy']}")
```

```
The best performing model is: Logistic Regression with an accuracy of 0.8333333333333334
```

# 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 Logistic Regression is the best machine learning algorithm for this task.

Thank you!