# Data-Driven Innovations in Commercial Spaceflight: A Python Project

Srikanth Potharla

# Executive Summary

This project embarked on a rigorous analysis of SpaceX's launch records to predict first-stage landing successes using machine learning models. Leveraging historical launch data, we applied logistic regression, support vector machines, and k-nearest neighbors' algorithms, achieving over 80% accuracy. Through hyperparameter tuning via GridSearchCV, model performance was significantly enhanced, underscoring the potency of AutoML tools in predictive analytics. The project also illuminated patterns within the data, such as no discernible correlation between payload mass and launch outcomes and revealed the robustness of SpaceX's operations across varying payload ranges. These findings can drive future improvements in space launch strategies.

Commercial space travel marks a new frontier, offering private-sector-led exploration beyond Earth, pioneered by innovative companies like SpaceX and Blue Origin

**Virgin Galactic:** Pioneering suborbital spaceflights for tourism.
**Rocket Lab:** Specializes in launching small satellites to orbit.
**Blue Origin:** Developing reusable rockets for suborbital and orbital flights.
**SpaceX:** Leading with multiple space innovations, including missions to the ISS, the Starlink project, and commercial manned missions

**Highlight on SpaceX:** Noteworthy for their cost-effective, reusable Falcon 9 rockets. A standard launch costs approximately $62 million, significantly lower than industry standards of up to $165 million.

# Introduction to Commercial Space Travel

# Role of a Data Scientist in Space Y

**Introduction to Space Y:** A burgeoning space company aiming to rival SpaceX.

**Role of Data Scientist:**

**Objective:** Analyze and optimize cost structures for each launch.

**Methods:** Collect and analyze data on SpaceX's operations; create dashboards to visualize cost metrics.

**Innovation:** Develop a machine learning model to predict the reusability of the Falcon 9's first stage based on publicly available data.

**Understanding Rocket Stages:** An analogy:

**First Stage:** Described as the "Power Booster," akin to the bottom, largest, and most powerful part of a pencil, responsible for the initial thrust.

**Second Stage:** Called the "Precision Guide," similar to the middle part of the pencil, ensuring payload delivery into precise orbits.

**Introduction to SpaceX API:**
    Accessing comprehensive data on SpaceX launches including rocket specifications and launch outcomes via the SpaceX REST API.

**API Details:**
    **Base URL: api.spacexdata.com/v4/**
    **Key Endpoints:**
        **/capsules** – Details about the capsules.
        **/cores** – Information on the rocket cores.
        **/launches/past** – Data on past launches, the primary focus for analysis.

**Data Retrieval Process:**
    Utilizing Python's **requests** library to send GET requests to **/launches/past** to fetch historical launch data.
    Handle JSON formatted API responses for processing.

**Additional Techniques:**
    Supplement API data with web scraping of Falcon 9 data from Wikipedia using Python's **BeautifulSoup** library.

# Data Collection Using SpaceX API  methodology

# Data Wrangling for SpaceX Launch Analysis  Methodology

**Processing JSON Data:**
　　Converting complex JSON data into a pandas DataFrame using **json_normalize** for easier analysis.
**Data Cleaning Steps:**
　　**Impute NULL Values:** Calculating and filling missing values in 'PayloadMass' with the mean.
　　**Remove Unused Columns:** Example, 'LandingPad' if not applicable.
　　**Data Filtering:** Exclude data related to non-Falcon 9 launches.
**Key Attributes for Analysis:**
　　Detailed list of attributes like Flight Number, Booster Version, Payload Mass, Outcome, etc., which are crucial for performance analysis.
**Outcome Classification:**
　　Converting landing outcomes into binary classification (0 for unsuccessful, 1 for successful) for simplified analysis.

**What is EDA?**
A critical, initial step in any data science workflow aimed at understanding the characteristics of data through statistical summaries and visual tools.
**Key Practices in EDA:**
Conducting initial data checks, examine basic statistics, and utilize visualizations to understand data structure and attributes.
Assessing data's predictive power for determining the success of Falcon 9 first stage landings using more detailed statistical tools.
**Key Attributes for Analysis:**
**Success Rate Over Time:** Tracking improvements in landing success rates, utilizing launch number as a key feature.
**Impact of Launch Sites on Success:** Analysis of how different launch sites influence launch outcomes with statistical evidence.
**Data Preparation for Machine Learning:**
Selecting features based on correlation with success, apply transformations like one-hot encoding for categorical data.

# Exploratory Data Analysis (EDA) Methodology

**Overview of Interactive Visual Analytics:**
Using tools that allow users to interact with data visualizations dynamically, enhancing the depth of analysis and understanding.
**Tools and Techniques:**
**Folium:** Used for creating interactive maps to visualize geographic data related to launch sites.
**Plotly Dash:** Framework for developing comprehensive web-based dashboards that incorporate interactive elements.
**Methodology:**
**Part 1:** Using Folium to map launch sites and analyze geographical advantages or patterns.
**Part 2:** Building an interactive dashboard using Plotly Dash, incorporating elements like dropdown lists and range sliders for dynamic data exploration.

**Interactive Visual Analytics and Dashboard Building- methodology**

**Predictive Analysis Overview:**
Applying machine learning to predict future events based on historical data.
**Data Preprocessing:**
**Purpose:** Prepares data for model training.
**Activities:**
Standardize numeric data to a uniform scale.
Address missing values.
Encode categorical variables for numerical analysis.
**Data Splitting with Train_test_split:**
**Training Data:** Used to train the machine learning models.
**Testing Data:** Used to test and validate model accuracy.

# Building the Machine Learning Pipeline- methodology

**Model Training Techniques:**
**Logistic Regression:** used for Predicting probabilities of binary outcomes.
**Support Vector Machines (SVM):** Classifies by finding the best boundary between outcomes.
**Decision Trees and K-Nearest Neighbors (KNN):** Decision-based and proximity-based approaches, respectively.
**Optimizing Model Performance:**
**Grid Search:** Automates the search for the best model parameters to improve accuracy.
**Model Evaluation:**
**Confusion Matrix:** Illustrates the accuracy of model predictions against actual outcomes, detailing true positives, false positives, true negatives, and false negatives.

**Model Training, Optimization, and Evaluation- methodology**

## Task 2

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

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

* sqlite:///my_data1.db
one.

| Date | Time (UTC) | Booster_Version | Launch_Site | Payload | PAYLOAD_MAS |
|---|---|---|---|---|---|
| 2010-06-04 | 18:45:00 | F9 v1.0 B0003 | CCAFS LC-40 | Dragon Spacecraft Qualification Unit | |
| 2010-12-08 | 15:43:00 | F9 v1.0 B0004 | CCAFS LC-40 | Dragon demo flight C1, two CubeSats, barrel of Brouere cheese | |
| 2012-05-22 | 7:44:00 | F9 v1.0 B0005 | CCAFS LC-40 | Dragon demo flight C2 | |
| 2012-10-08 | 0:35:00 | F9 v1.0 B0006 | CCAFS LC-40 | SpaceX CRS-1 | |
| 2013-03-01 | 15:10:00 | F9 v1.0 B0007 | CCAFS LC-40 | SpaceX CRS-2 | |

## Task 1

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

```sql
%sql SELECT DISTINCT "Launch_Site" FROM SPACEXTBL;
```

* sqlite:///my_data1.db
one.

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

## Task 4

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

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

* sqlite:///my_data1.db
one.

**Average_Payload_Mass**

2534.6666666666665

## Task 3

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

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

* sqlite:///my_data1.db
Done.

**Total_Payload_Mass**

48213

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
%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
ne.

| Booster_Version |
| --- |
| F9 FT B1022 |
| F9 FT B1026 |
| F9 FT B1021.2 |
| F9 FT B1031.2 |

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") AS First_Successful_Landing FROM SPACEXTBL WHERE "Landing_Outcome" = 'Success (ground pad)';
```

* sqlite:///my_data1.db
one.

| First_Successful_Landing |
| --- |
| 2015-12-22 |

## Task 7

List the total number of successful and failure mission outcomes

```
%sql SELECT CASE WHEN "Mission_Outcome" LIKE '%Success%' THEN 'Success' WHEN "Mission_Outcome" LIKE '%Failure%' THEN 'Failure' ELSE 'Other' END AS Outcome, COUNT(*) AS Total FRO
```

* sqlite:///my_data1.db
one.

| Outcome | Total |
|---------|-------|
| Failure | 1 |
| Success | 100 |

## Task 8

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

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

* sqlite:///my_data1.db
one.

| 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 |

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: SQLLite 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.**

`%sql SELECT substr(Date, 6, 2) AS Month, "Booster_Version", "Launch_Site", "Landing_Outcome" FROM SPACEXTBL WHERE substr(Date, 1, 4) = '2015' AND "Landing_Outcome" LIKE '%Failure (dr`

* sqlite:///my_data1.db
Done.

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

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.

`%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
one.

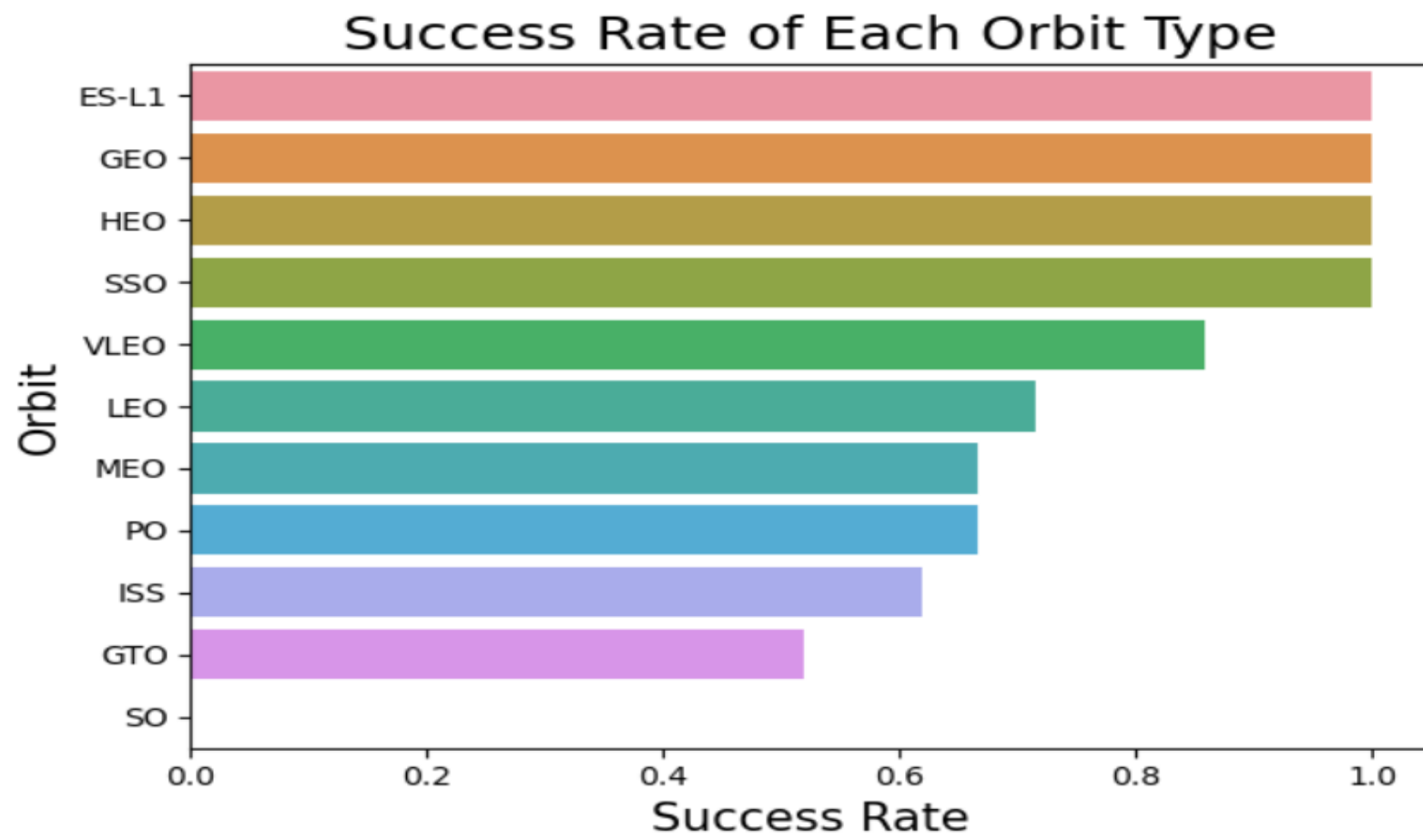| Landing_Outcome | Outcome_Count |
|-----------------|---------------|
| No attempt | 10 |
| Success (drone ship) | 5 |
| Failure (drone ship) | 5 |
| Success (ground pad) | 3 |
| Controlled (ocean) | 3 |
| Uncontrolled (ocean) | 2 |
| Failure (parachute) | 2 |
| Precluded (drone ship) | 1 |

EDA WITH VISUALIZATION RESULTS

Relationship between Flight Number and Launch Site by Class

Relationship between Payload Mass and Launch Site by Class

## Success Rate of Each Orbit Type
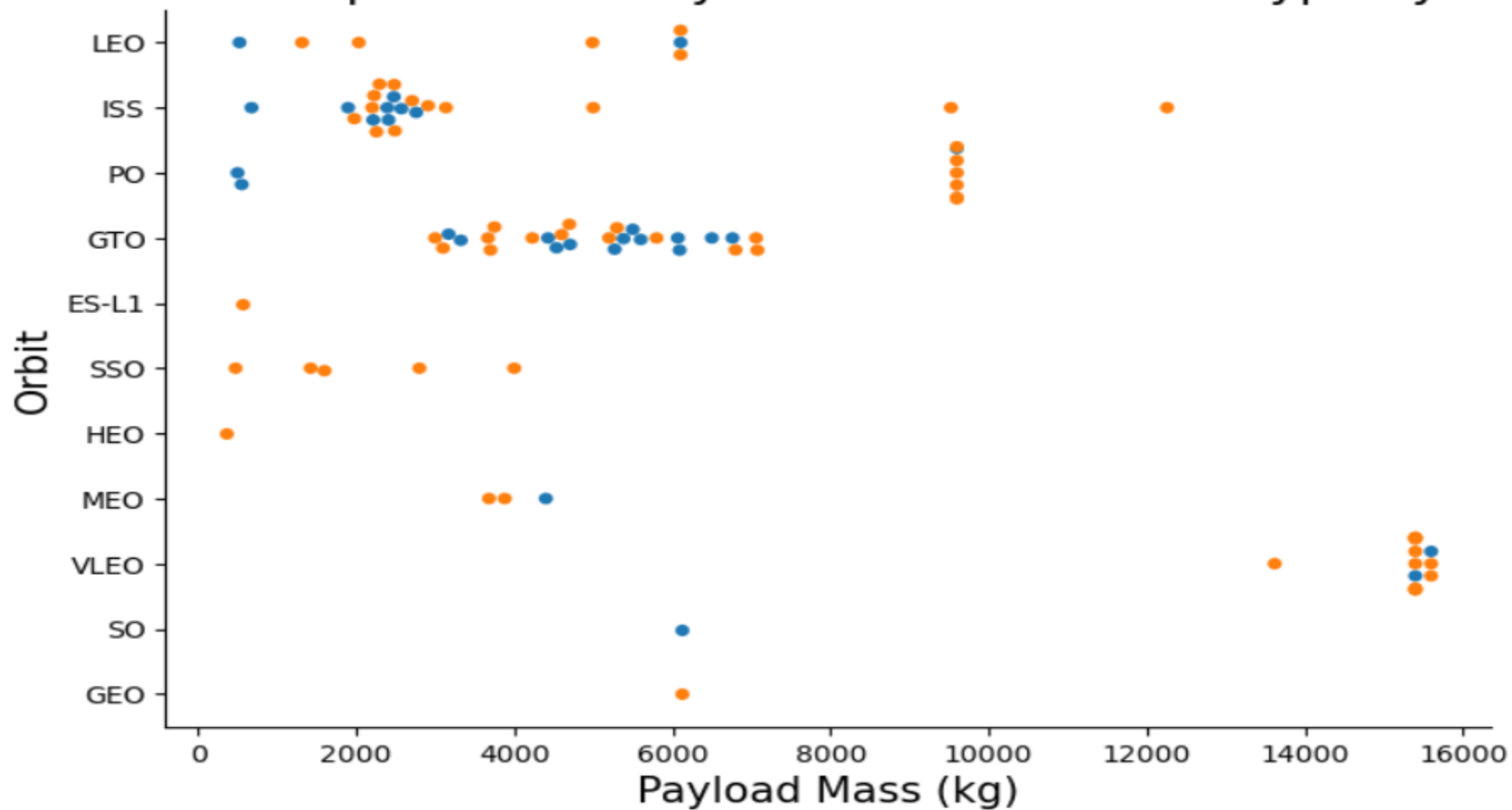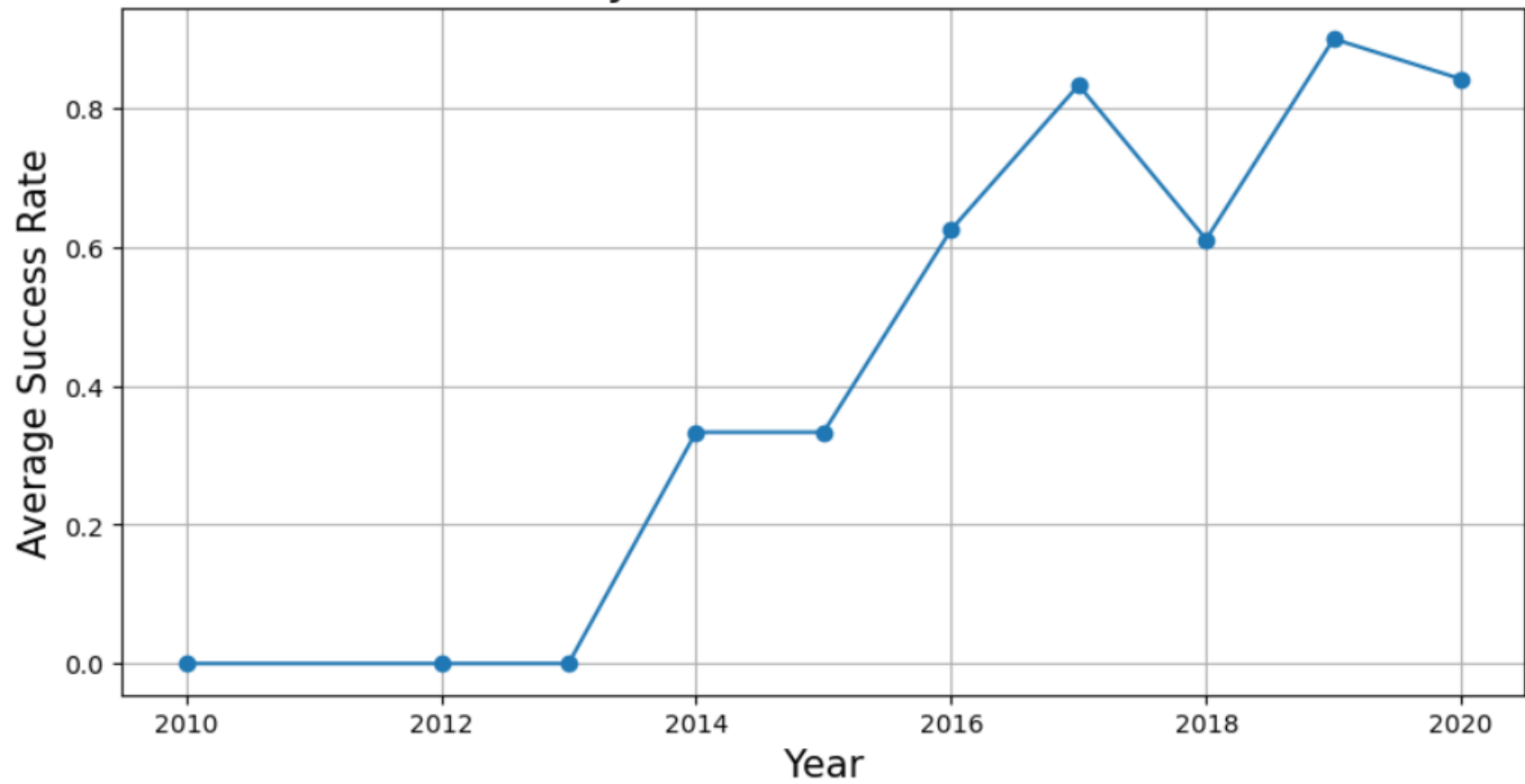
Relationship between Flight Number and Orbit Type by Class
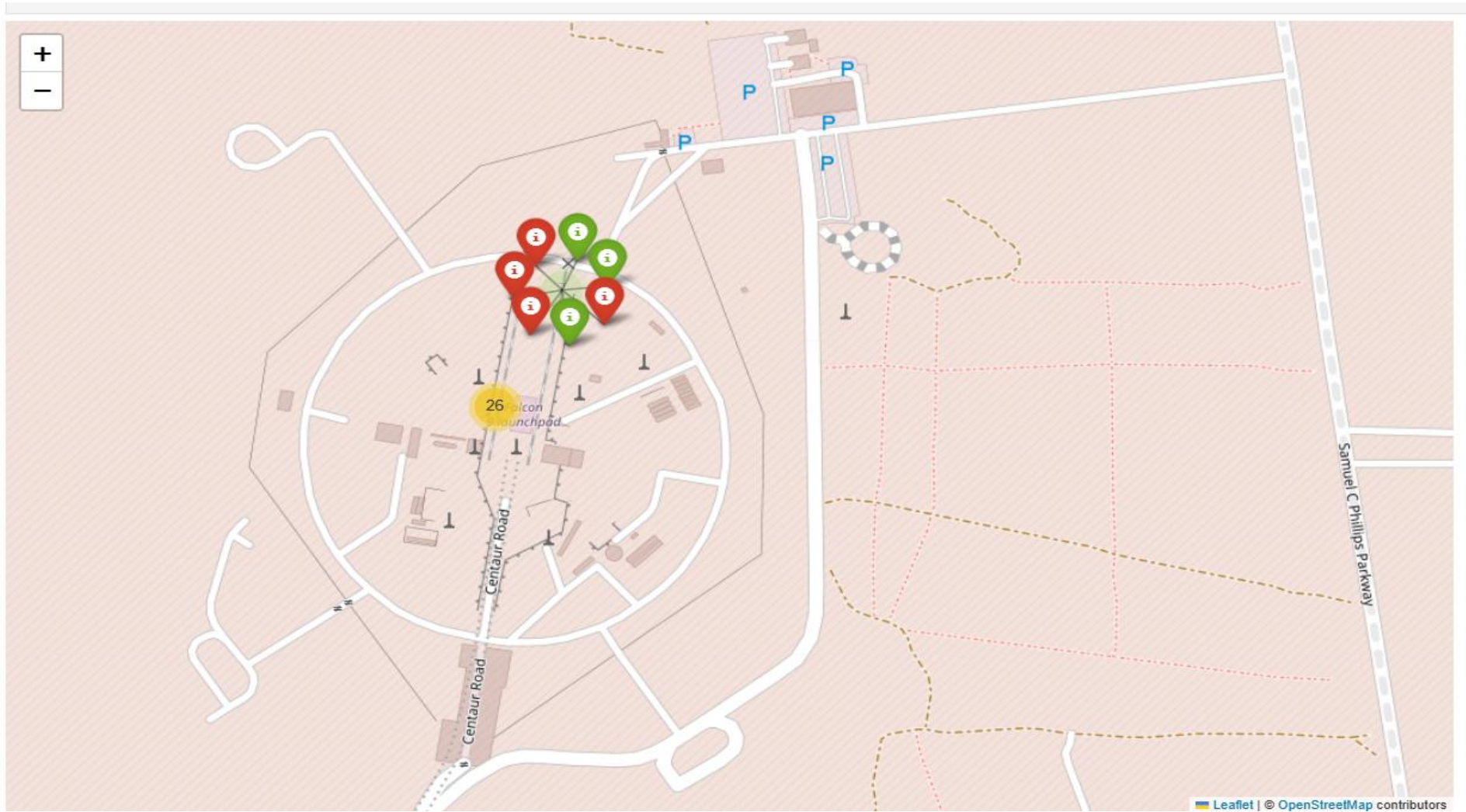
Relationship between Payload Mass and Orbit Type by Class

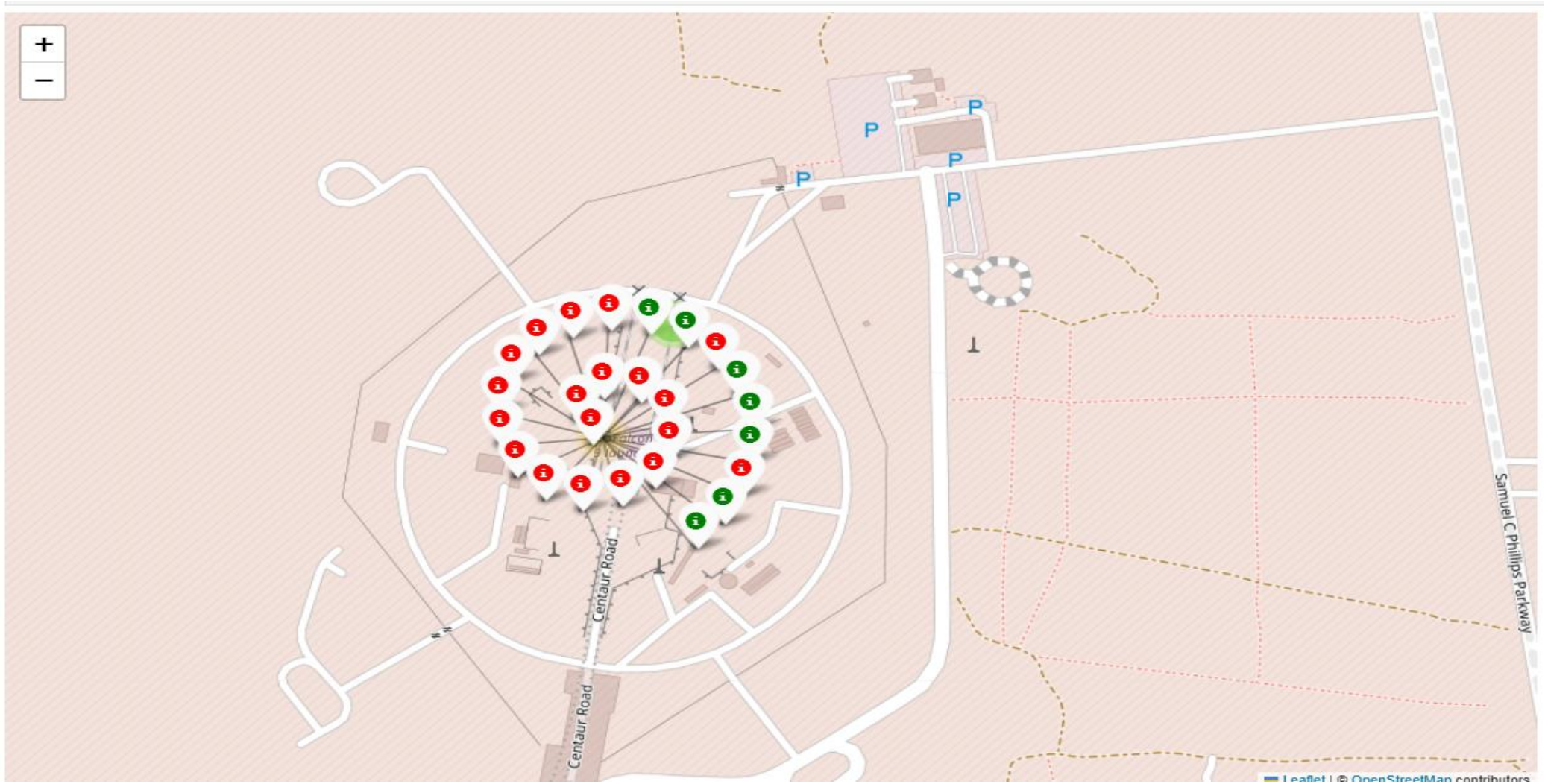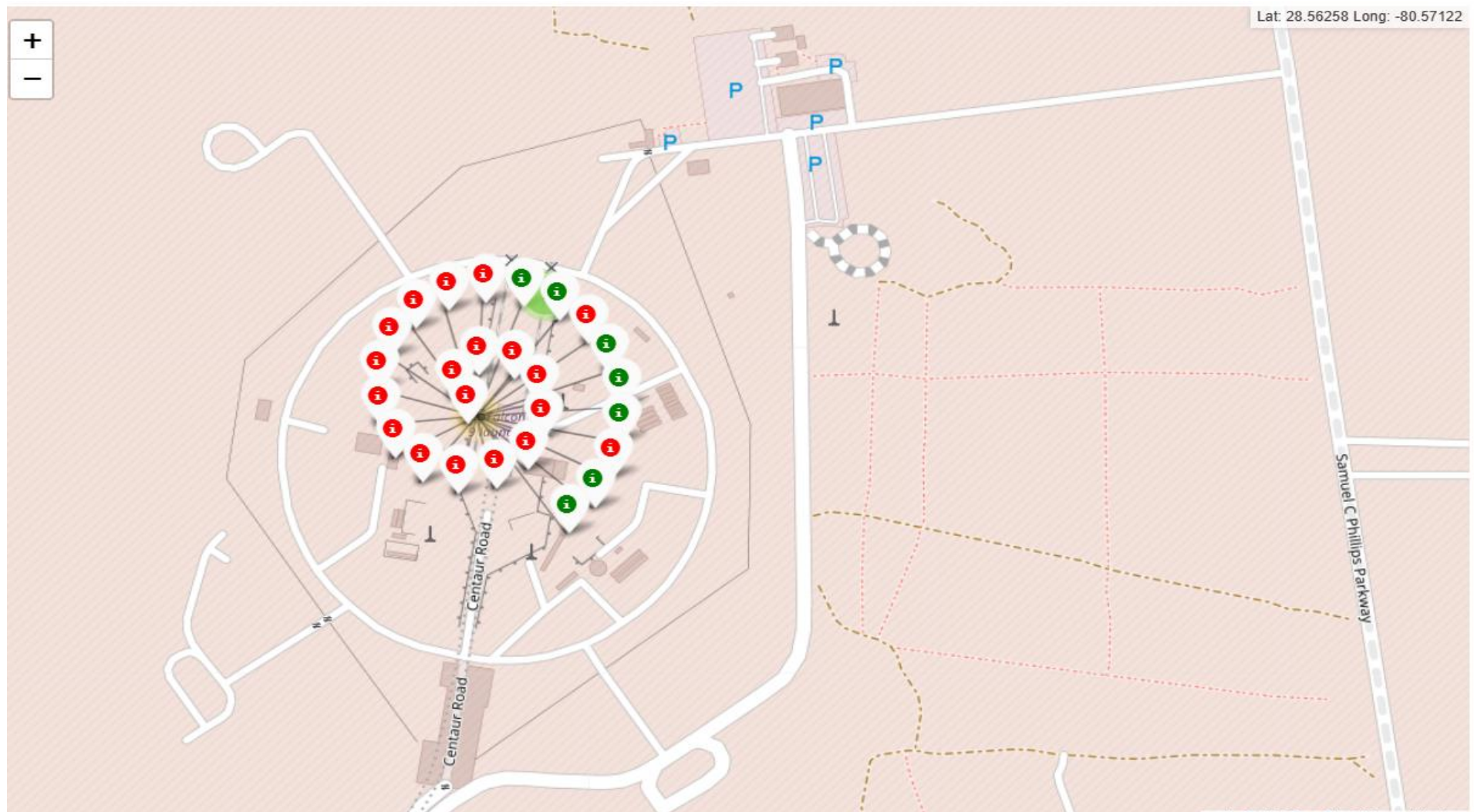Yearly Launch Success Trend

# INTERACTIVE MAPS WITH FOLIUM

# SpaceX Launch Records Dashboard

All Sites          × ▾

○————○————○————○————○————○————○————○————○————○————○
0 Kg    1000 Kg  2000 Kg  3000 Kg  4000 Kg  5000 Kg  6000 Kg  7000 Kg  8000 Kg  9000 Kg  10000 Kg

## Launch Successes and Failures for All Sites



■ 0
■ 1

42.9%

57.1%

## Correlation between Payload and Success for All Sites



**Booster Version**
- ● F9 v1.0  B0003
- ● F9 v1.0  B0004
- ● F9 v1.0  B0005
- ● F9 v1.0  B0006
- ● F9 v1.0  B0007
- ● F9 v1.1
- ● F9 v1.1 B1011
- ● F9 v1.1 B1010
- ● F9 v1.1 B1012
- ● F9 v1.1 B1013
- ● F9 v1.1 B1...
- ● F9 v1.1 B1...
- ● F9 v1.1 B10...

class

1
0.8
0.6
0.4
0.2
0

0        5k        10k

# SpaceX Launch Records Dashboard

CCAFS LC-40    ×  ▾

○————○————○————○————○————○————○————○————○————○————○
0 Kg    1000 Kg    2000 Kg    3000 Kg    4000 Kg    5000 Kg    6000 Kg    7000 Kg    8000 Kg    9000 Kg    10000 Kg

### Launch Successes and Failures for CCAFS LC-40



26.9%

73.1%

■ 0
■ 1

### Correlation between Payload and Success for CCAFS LC



**Booster Version**
- F9 v1.0  B0003
- F9 v1.0  B0004
- F9 v1.0  B0005
- F9 v1.0  B0006
- F9 v1.0  B0007
- F9 v1.1
- F9 v1.1 B1011
- F9 v1.1 B1010
- F9 v1.1 B1012
- F9 v1.1 B1013
- F9 v1.1 B1...
- F9 v1.1 B1...
- F9 v1.1 B10...

# SpaceX Launch Records Dashboard

VAFB SLC-4E ✕ ▼

○──────○──────○──────○──────○──────○──────○──────○──────○──────○
0 Kg    1000 Kg  2000 Kg  3000 Kg  4000 Kg  5000 Kg  6000 Kg  7000 Kg  8000 Kg  9000 Kg  10000
                                                                                        Kg

## Launch Successes and Failures for VAFB SLC-4E

■ 0
■ 1

40%

60%

## Correlation between Payload and Success for VAFB SLC·

**Booster Version**
- F9 v1.1  B1003
- F9 v1.1 B1017
- F9 FT B1029.1
- F9 FT B1036.1
- F9 FT B1038.1
- F9 B4 B1041.1
- F9 FT  B1036.2
- F9 FT  B1038.2
- F9 B4  B1041.2
- F9 B4  B1043.2

Payload Mass (kg)

# SpaceX Launch Records Dashboard

KSC LC-39A    ✕ ▾

0 Kg    1000 Kg    2000 Kg    3000 Kg    4000 Kg    5000 Kg    6000 Kg    7000 Kg    8000 Kg    9000 Kg    10000 Kg

## Launch Successes and Failures for KSC LC-39A



- 1
- 0

23.1%

76.9%

## Correlation between Payload and Success for KSC LC-39



**Booster Version**

- F9 FT B1031.1
- F9 FT B1030
- F9 FT  B1021.2
- F9 FT B1032.1
- F9 FT B1034
- F9 FT B1035.1
- F9 FT  B1029.2
- F9 FT B1037
- F9 B4 B1039.1
- F9 B4 B1040.1
- F9 FT  B10
- F9 B4 B10
- F9 B5  B104

# SpaceX Launch Records Dashboard

CCAFS SLC-40

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 Kg | 1000 Kg | 2000 Kg | 3000 Kg | 4000 Kg | 5000 Kg | 6000 Kg | 7000 Kg | 8000 Kg | 9000 Kg | 10000 Kg |

## Launch Successes and Failures for CCAFS SLC-40



- 0
- 1

42.9%

57.1%

## Correlation between Payload and Success for CCAFS SL



**Booster Version**
- F9 FT  B1035.2
- F9 B4 B1043.1
- F9 FT  B1032.2
- F9 B4 B1044
- F9 B4  B1039.2
- F9 B4 B1045.1
- F9 B4  B1040.2

# PREDICTIVE ANALYSIS RESULTS

# TASK 3

Use the function train_test_split to split the data X and Y into training and test data. Set the parameter test_size to 0.2 and random_state to 2. The training data and test data should be assigned to the following labels.

```
X_train, X_test, Y_train, Y_test
```

```python
from sklearn.model_selection import train_test_split

# Split the data into training and testing sets
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2, random_state=2)

# Optionally, display the shapes of the outputs to confirm the split
print("Training data shape (X_train):", X_train.shape)
print("Test data shape (X_test):", X_test.shape)
print("Training labels shape (Y_train):", Y_train.shape)
print("Test labels shape (Y_test):", Y_test.shape)
```

```
Training data shape (X_train): (72, 84)
Test data shape (X_test): (18, 84)
Training labels shape (Y_train): (72,)
Test labels shape (Y_test): (18,)
```

we can see we only have 18 test samples.

# TASK 4

Create a logistic regression object then create a GridSearchCV object `logreg_cv` with cv = 10. Fit the object to find the best parameters from the dictionary `parameters` .

We output the `GridSearchCV` object for logistic regression. We display the best parameters using the data attribute `best_params_` and the accuracy on the validation data using the data attribute `best_score_` .

```python
from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import LogisticRegression

# Define the logistic regression model
lr = LogisticRegression()

# Parameters dictionary for GridSearchCV
parameters = {'C': [0.01, 0.1, 1],
              'penalty': ['l2'],
              'solver': ['lbfgs']}

# Create a GridSearchCV object with the logistic regression model
logreg_cv = GridSearchCV(lr, parameters, cv=10)

# Fit GridSearchCV to the training data
logreg_cv.fit(X_train, Y_train)

# Display the best hyperparameters and the accuracy of the best model
print("Tuned hyperparameters (best parameters): ", logreg_cv.best_params_)
print("Accuracy on validation data: ", logreg_cv.best_score_)
```

```
Tuned hyperparameters (best parameters):  {'C': 0.01, 'penalty': 'l2', 'solver': 'lbfgs'}
Accuracy on validation data:  0.8464285714285713
```

# TASK 5

Calculate the accuracy on the test data using the method `score` :
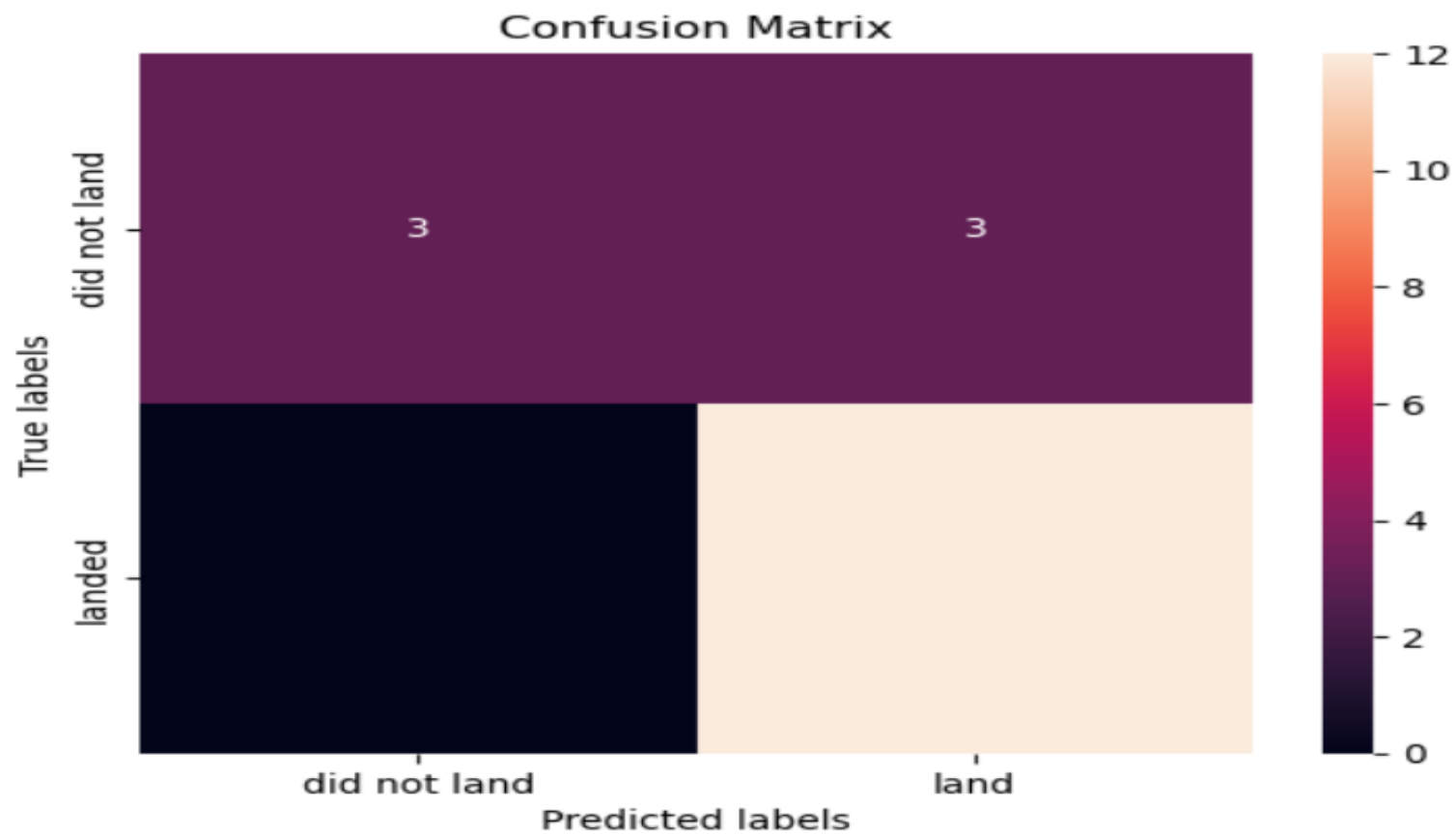
```python
# Calculate the accuracy on the test data
test_accuracy = logreg_cv.score(X_test, Y_test)


# Print the accuracy on the test data
print("Accuracy on test data: ", test_accuracy)
```

Accuracy on test data:  0.8333333333333334

```
[24]:  yhat=logreg_cv.predict(X_test)
        plot_confusion_matrix(Y_test,yhat)
```

## Confusion Matrix

# TASK 6

Create a support vector machine object then create a `GridSearchCV` object `svm_cv` with cv - 10. Fit the object to find the best parameters from the dictionary `parameters` .

```python
from sklearn.svm import SVC
from sklearn.model_selection import GridSearchCV
import numpy as np

# Define the support vector machine model
svm = SVC()

# Parameters dictionary for GridSearchCV
parameters = {
    'kernel': ('linear', 'rbf', 'poly', 'sigmoid'),
    'C': np.logspace(-3, 3, 5),
    'gamma': np.logspace(-3, 3, 5)
}

# Create a GridSearchCV object with the support vector machine model
svm_cv = GridSearchCV(svm, parameters, cv=10)

# Fit GridSearchCV to the training data
svm_cv.fit(X_train, Y_train)

# Display the best hyperparameters and the accuracy of the best model
print("Tuned hyperparameters (best parameters): ", svm_cv.best_params_)
print("Accuracy on validation data: ", svm_cv.best_score_)
```

```
Tuned hyperparameters (best parameters):  {'C': 0.03162277660168379, 'gamma': 0.001, 'kernel': 'linear'}
Accuracy on validation data:  0.8214285714285714
```

# TASK 7

Calculate the accuracy on the test data using the method `score`:

```python
# Calculate the accuracy on the test data
test_accuracy = svm_cv.score(X_test, Y_test)

# Print the accuracy on the test data
print("Accuracy on test data: ", test_accuracy)
```

Accuracy on test data:  0.8333333333333334

```python
# Expanding the parameter grid slightly more than the small test
expanded_parameters = {
    'criterion': ['gini', 'entropy'],
    'splitter': ['best', 'random'],
    'max_depth': [2, 4, 6],
    'max_features': ['auto', 'sqrt'],
    'min_samples_leaf': [1],
    'min_samples_split': [2, 5]
}


# Using more cross-validation folds than the smallest test, but fewer than the original large grid
expanded_tree_cv = GridSearchCV(tree, expanded_parameters, cv=5, verbose=3)  # Moderate verbosity
expanded_tree_cv.fit(X_train, Y_train)
print("Expanded Grid Search Complete.")
print("Best parameters: ", expanded_tree_cv.best_params_)
print("Accuracy on validation data: ", expanded_tree_cv.best_score_)
```
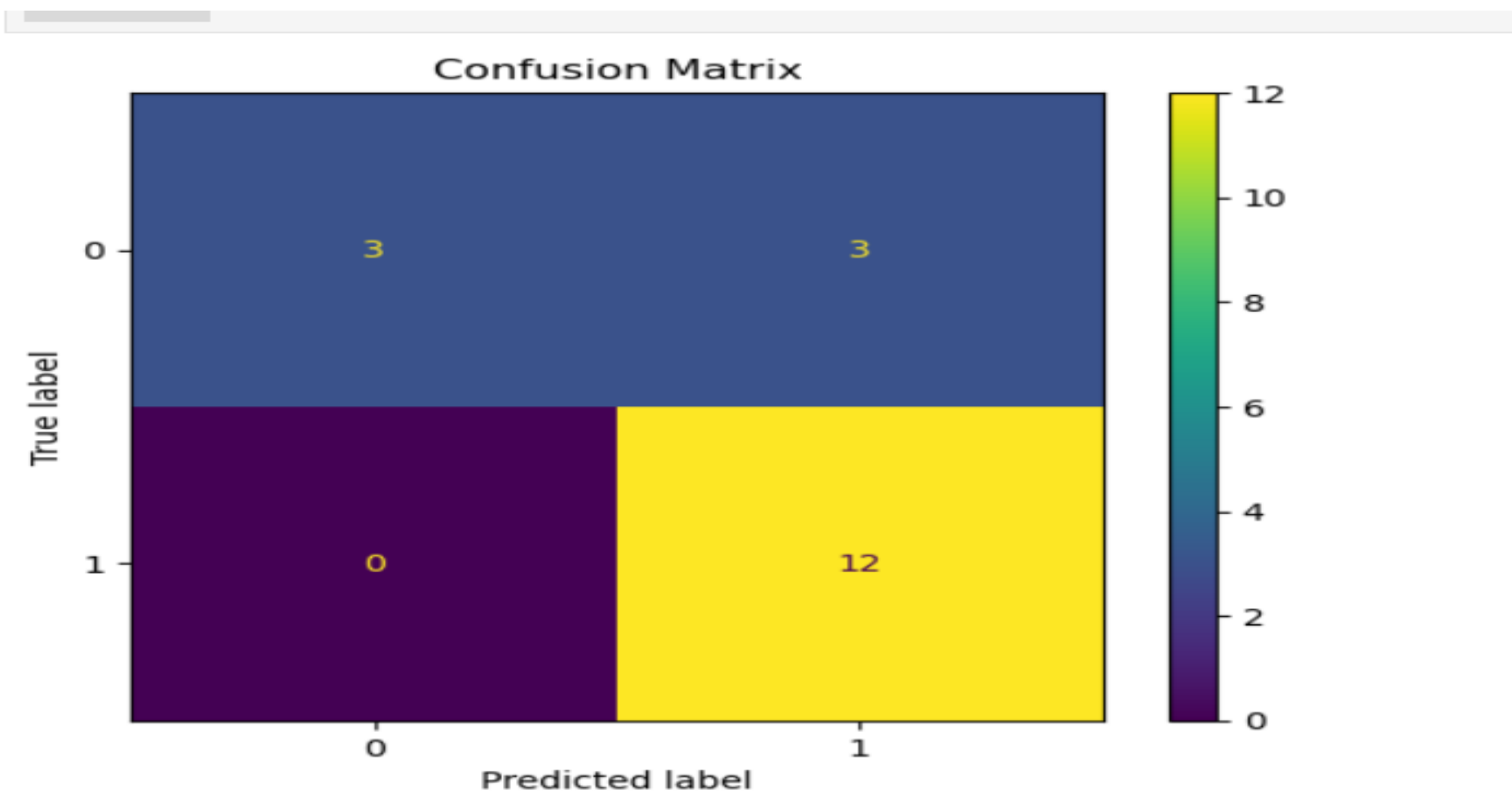
```
#Expanded Grid Search Complete.
#Best parameters:  {'criterion': 'gini', 'max_depth': 2, 'max_features': 'sqrt', 'min_samples_leaf': 1, 'min_samples_split': 5, 'splitter': 'best'}
Accuracy on validation data:  0.8885714285714286
```

Confusion Matrix

# TASK 10

Create a k nearest neighbors object then create a `GridSearchCV` object `knn_cv` with cv = 10. Fit the object to find the best parameters from the dictionary `parameters` .

```python
from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import GridSearchCV

# Define the K-Nearest Neighbors classifier
KNN = KNeighborsClassifier()

# Parameters dictionary for GridSearchCV
parameters = {
    'n_neighbors': list(range(1, 11)),  # Using a list of neighbors from 1 to 10
    'algorithm': ['auto', 'ball_tree', 'kd_tree', 'brute'],
    'p': [1, 2]  # 1 corresponds to Manhattan distance, 2 to Euclidean distance
}

# Create a GridSearchCV object with the KNN model
knn_cv = GridSearchCV(KNN, parameters, cv=10)

# Fit GridSearchCV to the training data
knn_cv.fit(X_train, Y_train)

# Display the best hyperparameters and the accuracy of the best model
print("Tuned hyperparameters (best parameters): ", knn_cv.best_params_)
print("Accuracy on validation data: ", knn_cv.best_score_)
```

```
Tuned hyperparameters (best parameters):  {'algorithm': 'auto', 'n_neighbors': 10, 'p': 1}
Accuracy on validation data:  0.8482142857142858
```

# TASK 11

Calculate the accuracy of knn_cv on the test data using the method `score` :
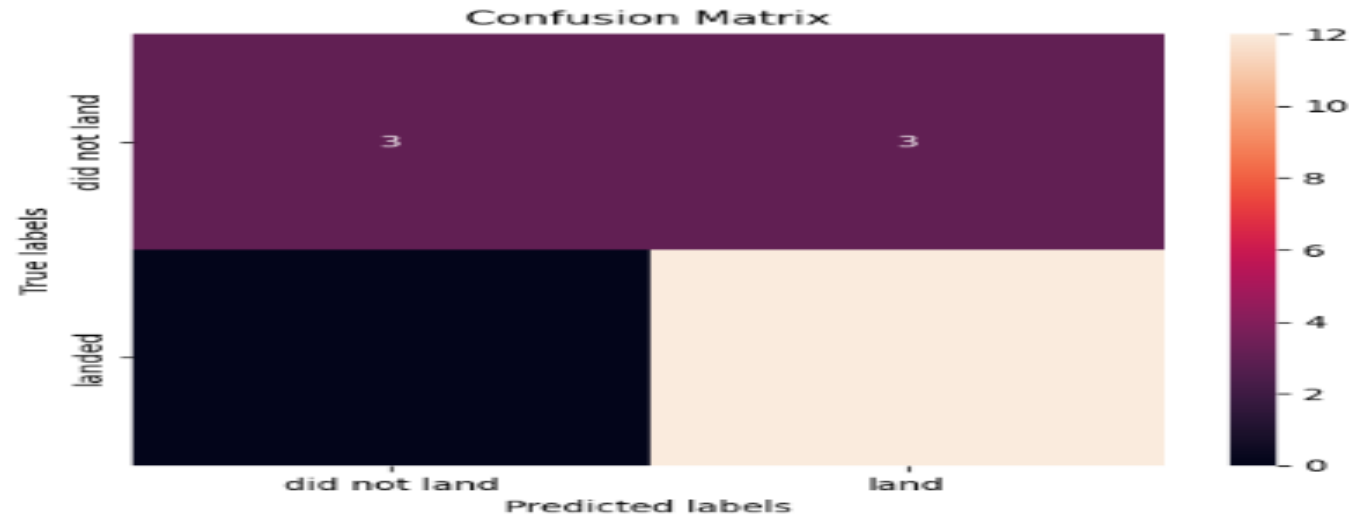
```
42]:   # Calculate the accuracy on the test data
       test_accuracy = knn_cv.score(X_test, Y_test)

       # Print the accuracy on the test data
       print("Accuracy on test data: ", test_accuracy)
```

```
Accuracy on test data:  0.8333333333333334
```

We can plot the confusion matrix

```
43]:   yhat = knn_cv.predict(X_test)
       plot_confusion_matrix(Y_test,yhat)
```

**Launch Site Variation**: Different launch sites have varying success rates; some sites show a notably higher likelihood of successful landings.

**Payload Influence**: Heavier payloads have been associated with successful outcomes, particularly at specific launch sites, suggesting a potential relationship between payload mass and successful landings.

**Technological Progress**: There is a clear trend showing an overall improvement in launch success rates over the years, indicating advancements in technology and operational experience.

**Orbit Type Success Rates**: Certain orbit types have a higher success rate; this may reflect the complexity of reaching different orbits, or it might indicate more focus and expertise on certain mission profiles.

**Flight Number Trends**: As the flight number increases, which correlates with more experience and technological refinements, there appears to be a higher success rate in landings, highlighting the benefits of iterative learning and process improvements.

# Conclusion – Exploratory Data Analysis(EDA)

**Strategic Launch Site Locations**: The launch sites are strategically placed near coastlines. This proximity to vast ocean areas is likely for safety, as it provides a clear flight path over unpopulated areas, minimizing risk in case of launch failures.

**Infrastructure Accessibility**: The maps show that launch sites are carefully situated with convenient access to critical infrastructure. Roads and railways are within close proximity, facilitating the transport of rocket components and supporting logistics efficiently, which is essential for launch preparations and operations.

**Conclusion – Interactive Folium Maps**

**Launch Outcomes by Site**: There is a variation in success rates across different launch sites, with some sites demonstrating a higher proportion of successful launches than others, which may suggest site-specific operational efficiencies or challenges.

**Payload and Success Correlation**: The scatter plots show no clear correlation between the payload mass and launch success, which suggests that SpaceX is capable of handling a wide range of payload masses with a consistent success rate, pointing towards the versatility and reliability of the launch vehicles.

**Conclusion – Plotly Dash-Boards**

**Model Accuracy**: The accuracy scores from different models are consistently above 80%, indicating a strong predictive capability for the test data. This suggests that the features selected for the model have a significant influence on the prediction of rocket landing successes or failures.

**Consistency Across Models**: Several models, including logistic regression, support vector machines, and k-nearest neighbors, yield similar accuracy scores, which suggests that the dataset is well-behaved, and the patterns are detectable by various algorithmic approaches.

**Confusion Matrix Interpretation**: The confusion matrices indicate a higher number of true positives and true negatives, which shows that the models are generally successful at predicting both successful and unsuccessful landings. However, the presence of false positives and false negatives indicates there is still room for model improvement.

**Hyperparameter Optimization**: The usage of GridSearchCV for hyperparameter tuning is effective in finding the optimal parameters for the models, which is evidenced by the improved accuracy scores upon tuning compared to default parameter settings. This demonstrates the importance of hyperparameter optimization in model performance.

# Conclusion – Predictive Analysis

Crafting a presentation with dynamic slide layouts and engaging backgrounds turns data into a visual narrative. Simple, clear bullet points and well-organized content ensure the information is digestible, making complex topics accessible.

This creative approach not only captivates the audience but also enhances comprehension and retention of the material.

**Creativity Applied in Presentation**

1. innovative and interesting insight from visualization of data as a part of EDA is the improvement in launch success rates over time, which suggests a learning curve effect at play. As the number of flights increases, the success rate also goes up, indicating that experience and technological advancements lead to better mission outcomes. This underscores the iterative nature of engineering in space travel where each launch provides valuable data to refine and enhance future missions.

2. An innovative and interesting insight from the predictive analysis is the effectiveness of the hyperparameter tuning process in improving model performance. Despite the different algorithms used, the GridSearchCV technique consistently found the best parameters that enhanced the accuracy of each model. This highlights the potential for automated machine learning (AutoML) tools to significantly streamline and optimize the model-building process, leading to robust predictive performance with less manual intervention.

## Innovative Insights from the Data Analysis

Thank you