



IBM Developer
SKILLS NETWORK

Winning Space Race with Data Science

Subhash Somarouthu
15th April 2025



Index

- Executive Summary
- Introduction
- Methodology
- Results
- Conclusion
- Appendix

Executive Summary

This project investigates the success of SpaceX Falcon 9 first-stage landings using data-driven techniques.

Data was collected via the SpaceX API and web scraping from Wikipedia to build a comprehensive dataset.

After data wrangling and transformation, exploratory data analysis (EDA) was conducted using SQL, visualization tools, and interactive analytics.

An interactive map (Folium) and dashboard (Plotly Dash) were developed to visually explore launch sites, payloads, and outcomes.

Predictive classification models were built to determine whether a Falcon 9 first stage would land successfully.

The decision tree classifier achieved the highest accuracy, helping identify key factors that influence successful landings.

These insights can help competitors assess costs and strategies when bidding against SpaceX for future launches.

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

- 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
- [Falcon9 First Stage Land Prediction/SpaceX Machine Learning Prediction.ipynb at main · subhashsomarouthu/Falcon9 First Stage Land Prediction](#)

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
- [Falcon9 First Stage Land Prediction/webscraping.ipynb at main · subhashsomarouthu/Falcon9 First Stage Land Prediction](#)

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
- [Falcon9 First Stage Land Prediction/edadatavizualization.ipynb at main · subhashsomarouthu/Falcon9 First Stage Land Prediction](#)

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
- [Falcon9 First Stage Land Prediction/edadatavizualization.ipynb at main · subhashsomarouthu/Falcon9 First Stage Land Prediction](#)

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.

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
[Falcon9 First Stage Land Prediction/app.py at main · subhashsomarouthu/Falcon9 First Stage Land Prediction](#)

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 [Falcon9 First Stage Land Prediction/SpaceX Machine Learning Prediction.ipynb at main · subhashsomarouthu/Falcon9 First Stage Land Prediction](#)

Results

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

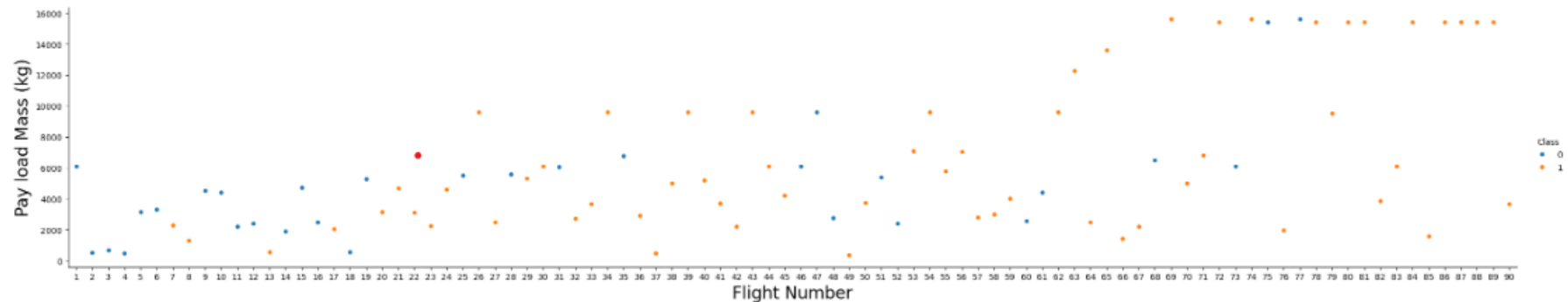
The background of the slide is a complex, abstract composition. It features a dark blue base color. Overlaid on this are numerous diagonal streaks in shades of red and cyan, creating a sense of motion and depth. A faint, white grid pattern is visible across the entire image, adding a technical or digital feel. The overall effect is modern and high-tech.

Section 2

Insights drawn from EDA

Flight Number vs. Launch Site

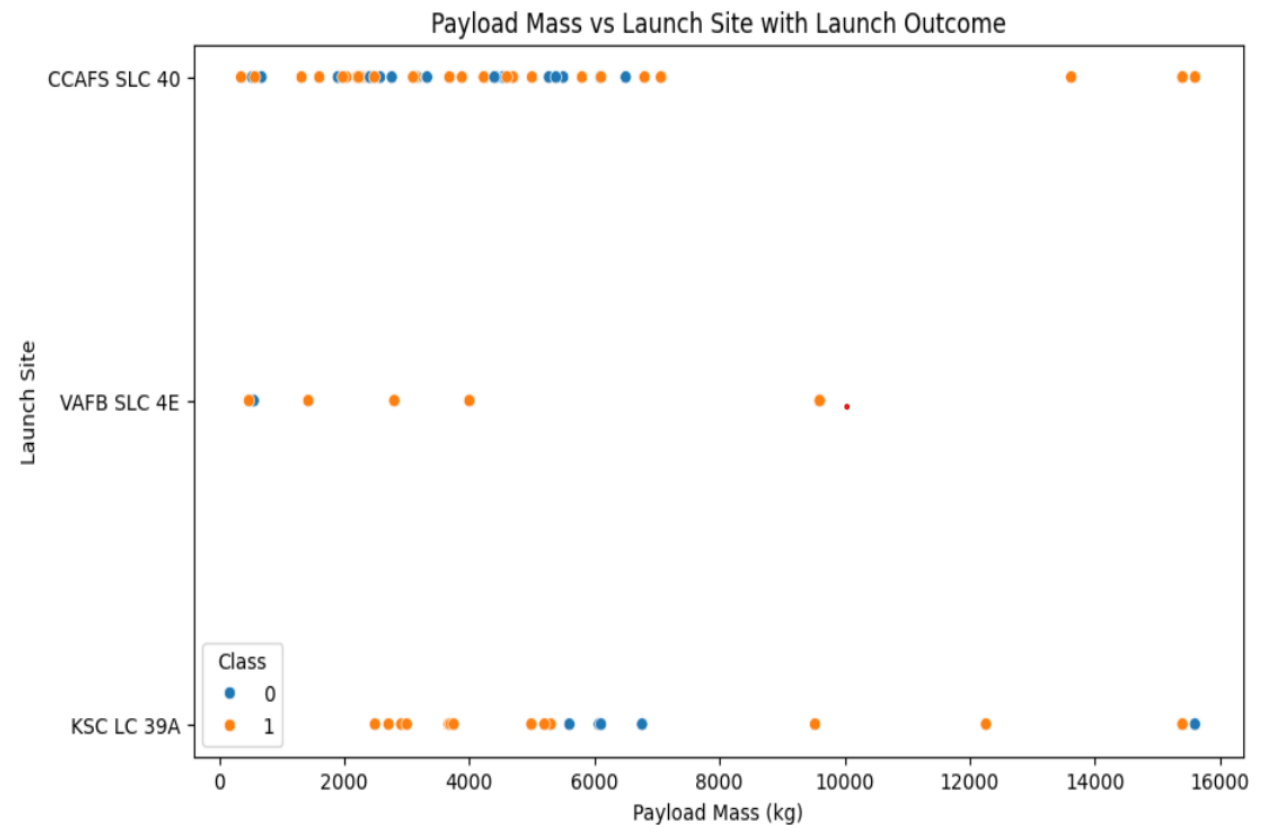
- From the plot, we found that the larger the flight amount at a launch site, the greater the success rate at a launch site.



Payload vs. Launch Site

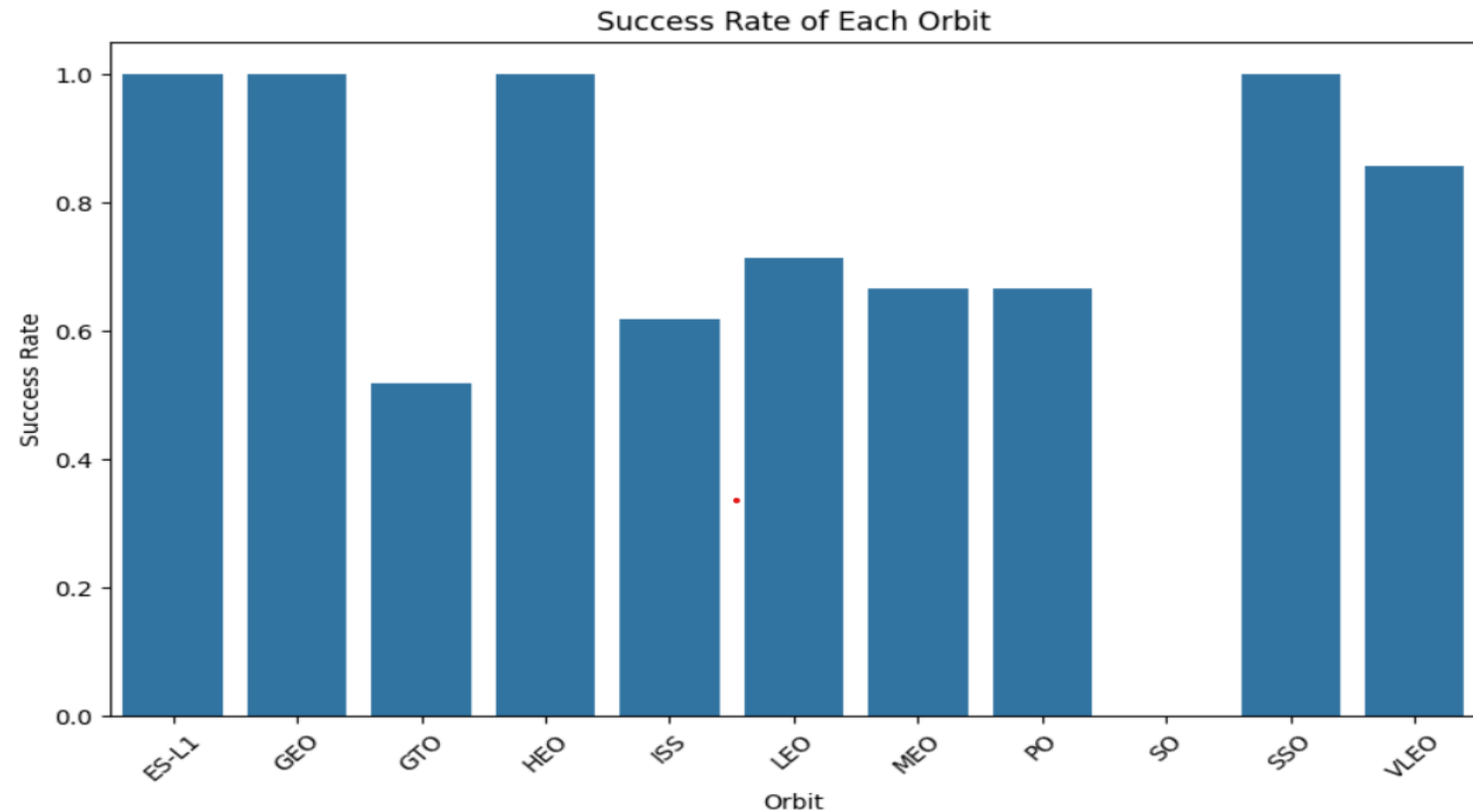


The greater the payload mass for launch site CCAFS SLC 40 the higher the success rate for the rocket.



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.



Total Payload Mass

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

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

```
In [12]: task_3 = '''
          SELECT SUM(PayloadMassKG) AS Total_PayloadMass
          FROM SpaceX
          WHERE Customer LIKE 'NASA (CRS)'
          '''
          create_pandas_df(task_3, database=conn)
```

Out[12]:

	total_payloadmass
0	45596

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 upper section and a photograph of the Earth's surface below. The Earth's surface shows a dark blue ocean with numerous bright yellow and orange lights from cities and towns, particularly concentrated along the coastlines and in the lower right quadrant. The horizon of the Earth is visible as a thin, curved line separating the dark surface from the blackness of space.

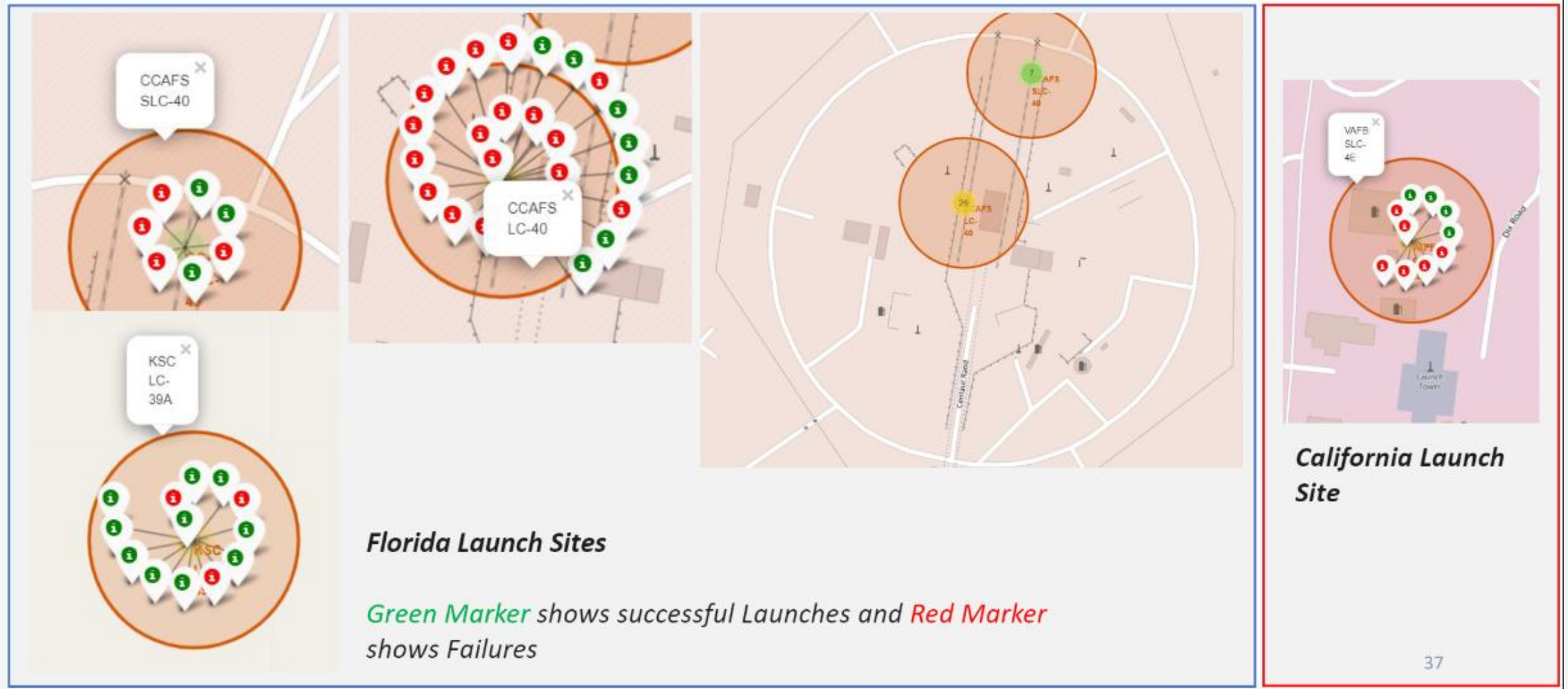
Section 4

Launch Sites Proximities Analysis

All launch sites global map markers



Markers showing launch sites with color labels





Section 5

Build a Dashboard with Plotly Dash

SpaceX Launch Records Dashboard

All Sites

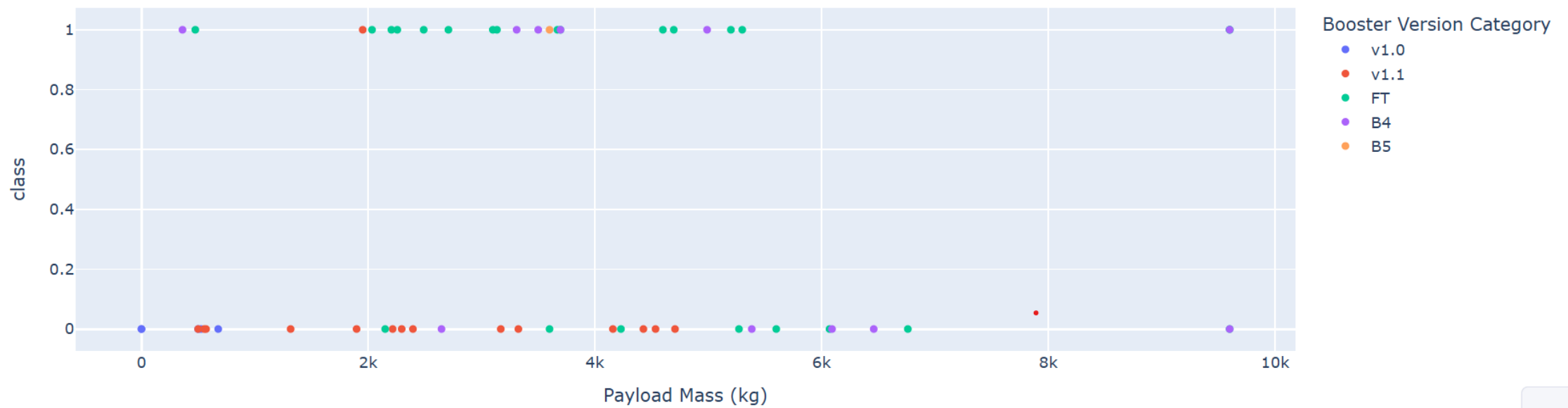
Total Success Launches by Site



Payload range (Kg):



Correlation between Payload and Launch Success



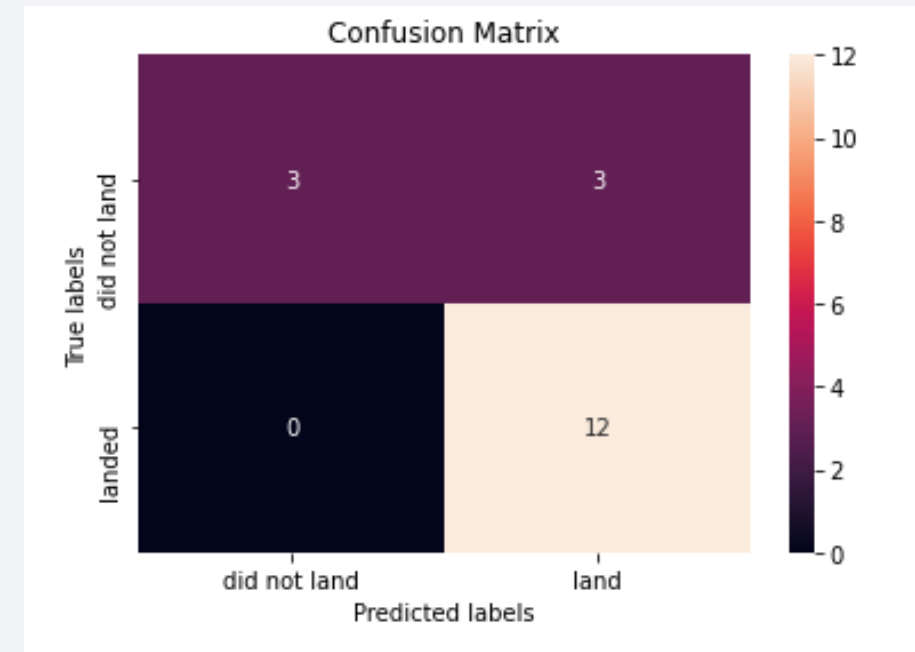


Section 6

Predictive Analysis (Classification)

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

- Thank you