



# DATA SCIENCE FOR SPACE RACE

*Submitted by:*  
*Hadiya Ashraf*

# CONTENTS

- Executive Summary
- Introduction
- Methodology
- Results
- Conclusions

# 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 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.

## RESEARCH QUESTIONS

- 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

## *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 And Scrapping

- We used the get request to the SpaceX API to collect data, clean the requested data and do some basic data wrangling and formatting.

- The link to the notebook is

[https://github.com/Hadiya29/Data\\_Science\\_Capstone\\_Project/blob/main/jupyter-labs-spacex-data-collection-api.ipynb](https://github.com/Hadiya29/Data_Science_Capstone_Project/blob/main/jupyter-labs-spacex-data-collection-api.ipynb)

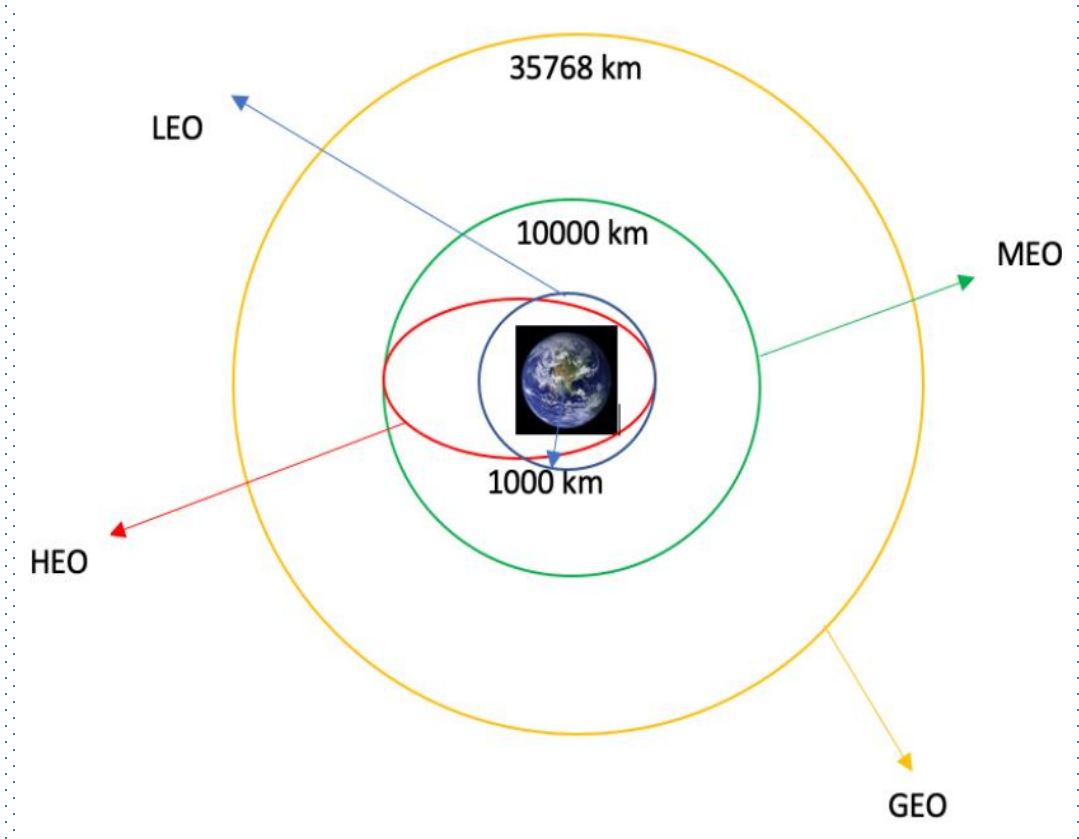
- We applied web scrapping to web scrap Falcon 9 launch records with BeautifulSoup
- We parsed the table and converted it into a pandas data frame.

- The link to the notebook is

[https://github.com/Hadiya29/Data\\_Science\\_Capstone\\_Project/blob/main/jupyter-labs-webscraping.ipynb](https://github.com/Hadiya29/Data_Science_Capstone_Project/blob/main/jupyter-labs-webscraping.ipynb)

# Data Wrangling

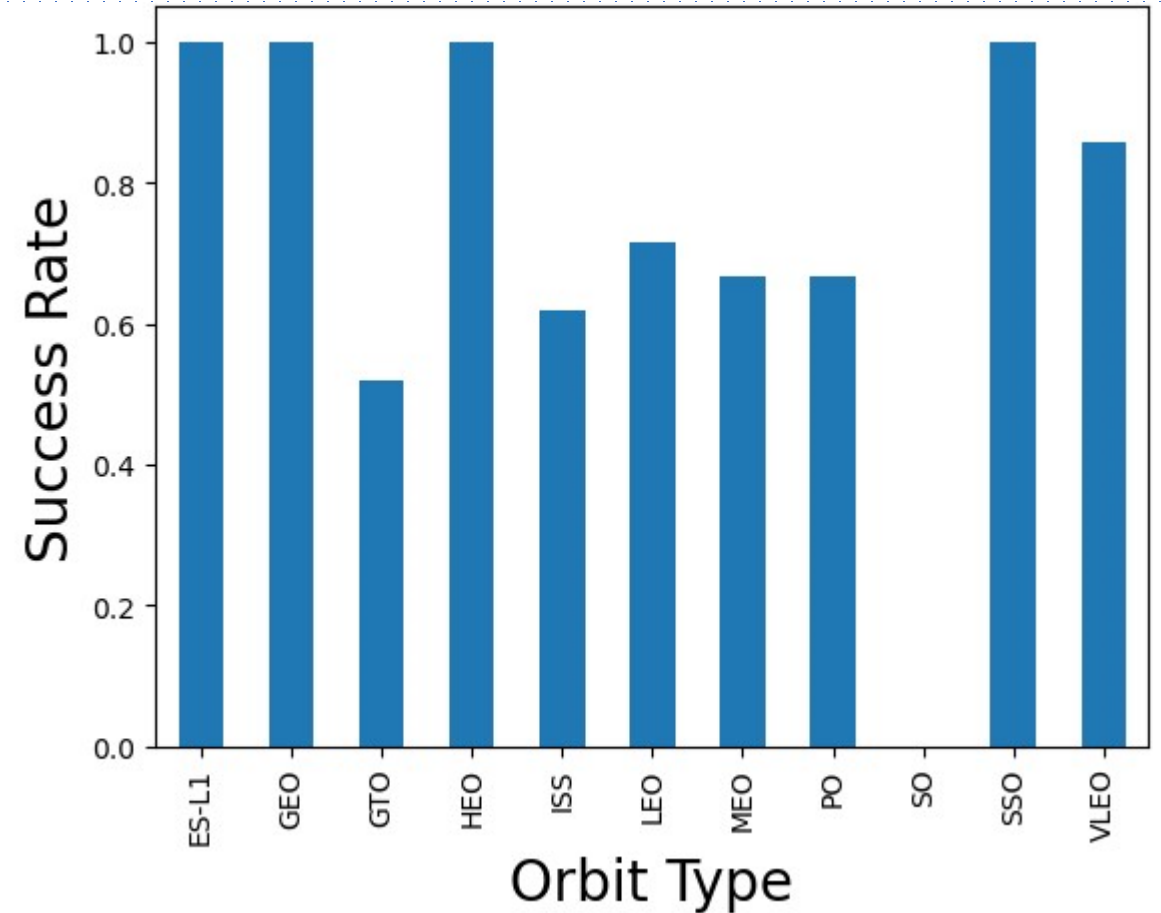
- We performed exploratory data analysis and determined the training labels.
- We calculated the number of launches at each site, and the number and occurrence of each orbits
- We created landing outcome label from outcome column and exported the results to csv.
- The link to the notebook is [https://github.com/Hadiya29/Data\\_Science\\_Capstone\\_Project/blob/main/labs-jupyter-spacex-Data%20wrangling.ipynb](https://github.com/Hadiya29/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/Hadiya29/Data Science Capstone Project/blob/main/jupyter-labs-eda-dataviz.ipynb.jupyterlite.ipynb](https://github.com/Hadiya29/Data_Science_Capstone_Project/blob/main/jupyter-labs-eda-dataviz.ipynb.jupyterlite.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/Hadiya29/Data\\_Science\\_Capstone\\_Project/blob/main/jupyter-labs-eda-sql-coursera\\_sqllite.ipynb](https://github.com/Hadiya29/Data_Science_Capstone_Project/blob/main/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.

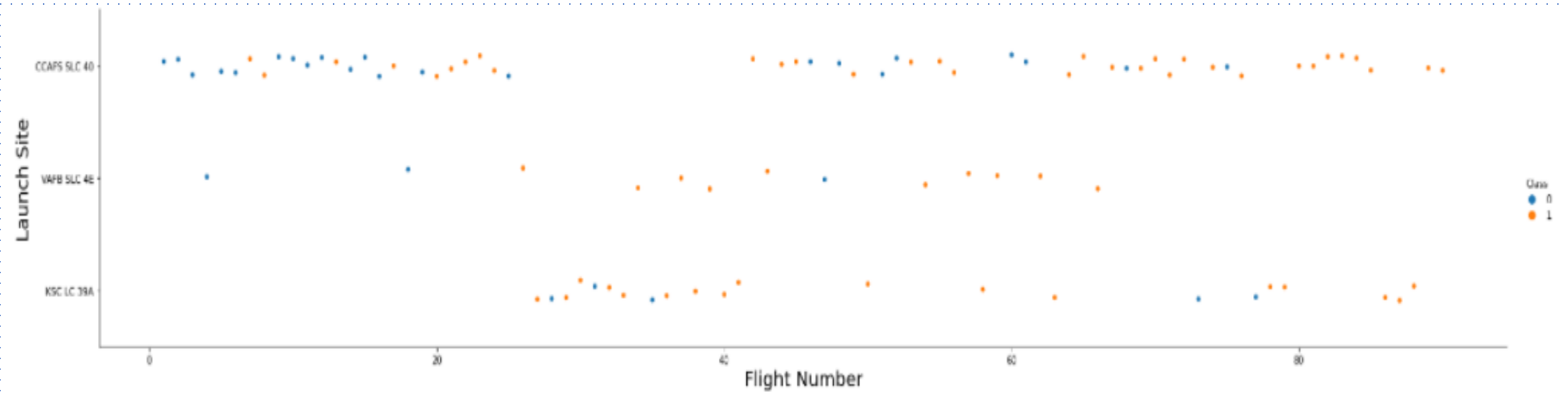
# Predictive Analysis (Classification)

- 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/Hadiya29/Data\\_Science\\_Capstone\\_Project/blob/main/SpaceX\\_Machine\\_Learning\\_Prediction\\_Part\\_5.jupyterlite%20\(1\).ipynb](https://github.com/Hadiya29/Data_Science_Capstone_Project/blob/main/SpaceX_Machine_Learning_Prediction_Part_5.jupyterlite%20(1).ipynb)

# RESULTS

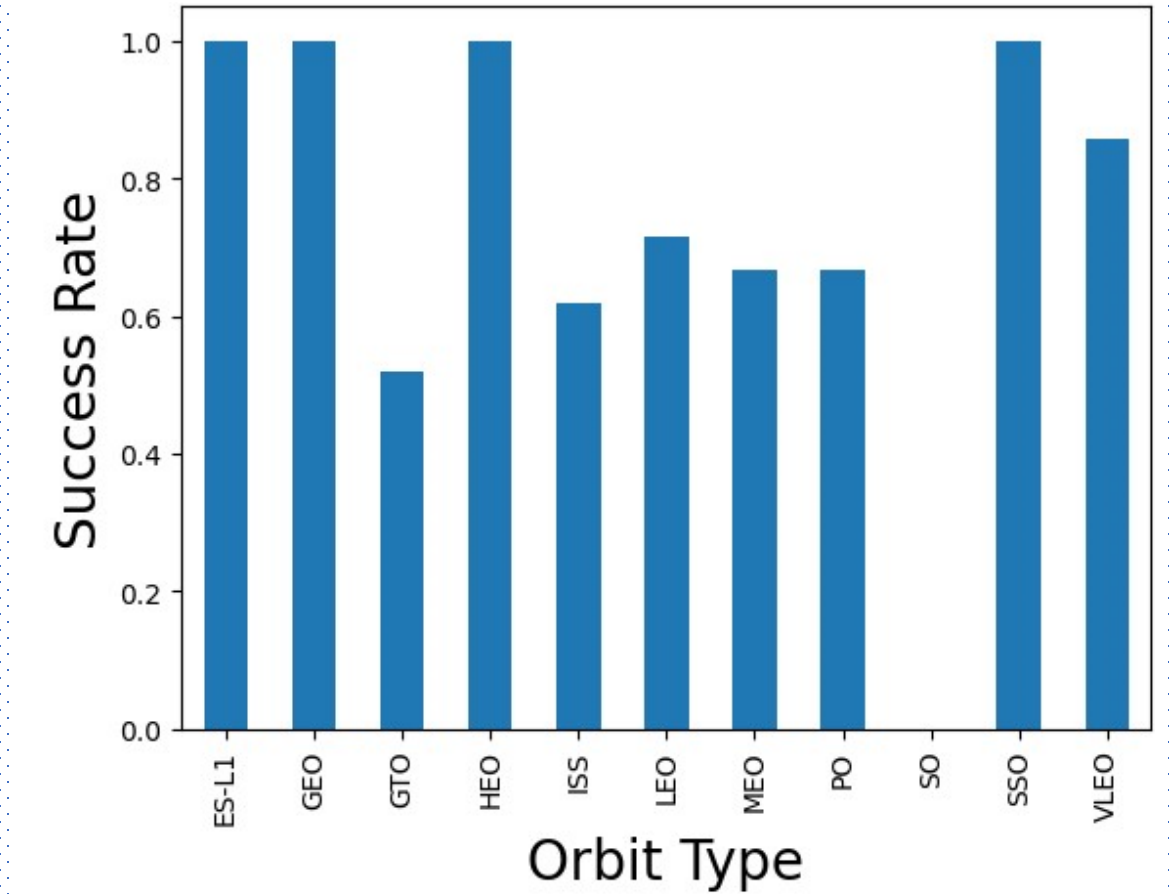
## Flight Number vs. Launch Site

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



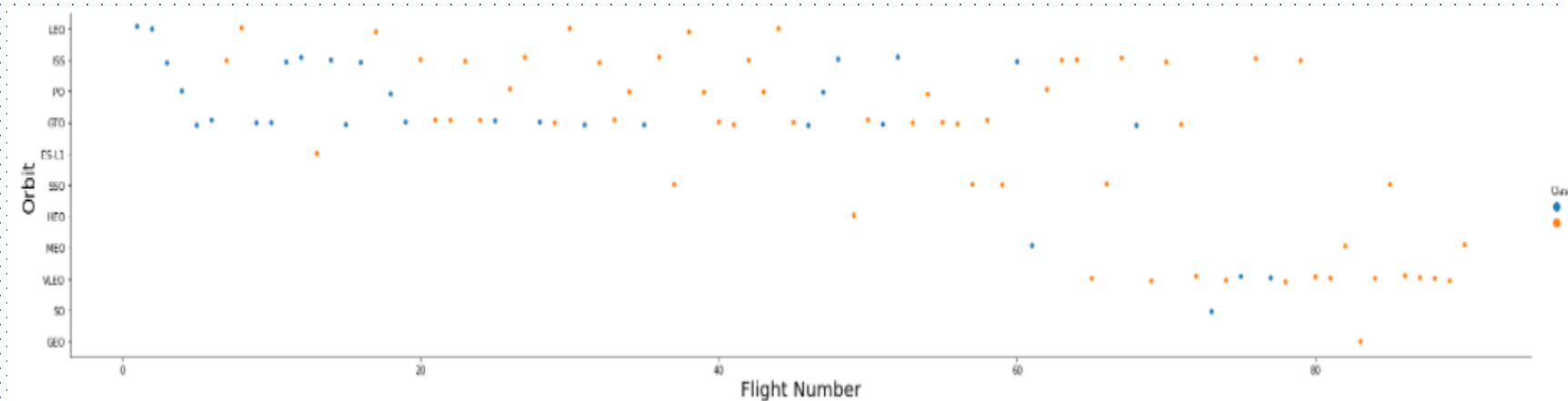
# Success Rate vs. Orbit Type

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



# Flight Number vs. Orbit Type

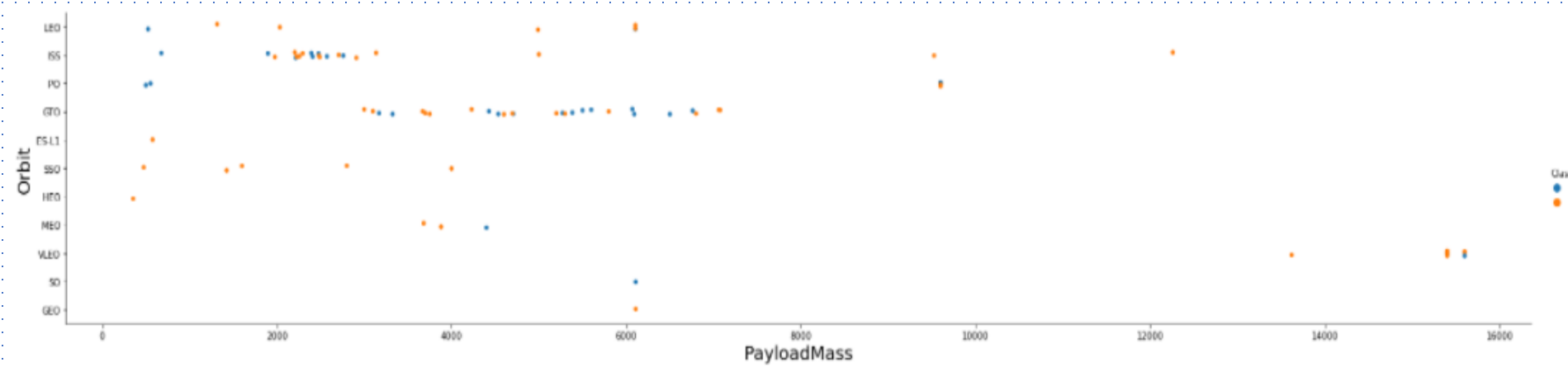
- The plot below shows the Flight Number vs. Orbit type. We observe that in the LEO orbit, success is related to the number of flights whereas in the GTO orbit, there is no relationship between flight number and the orbit.





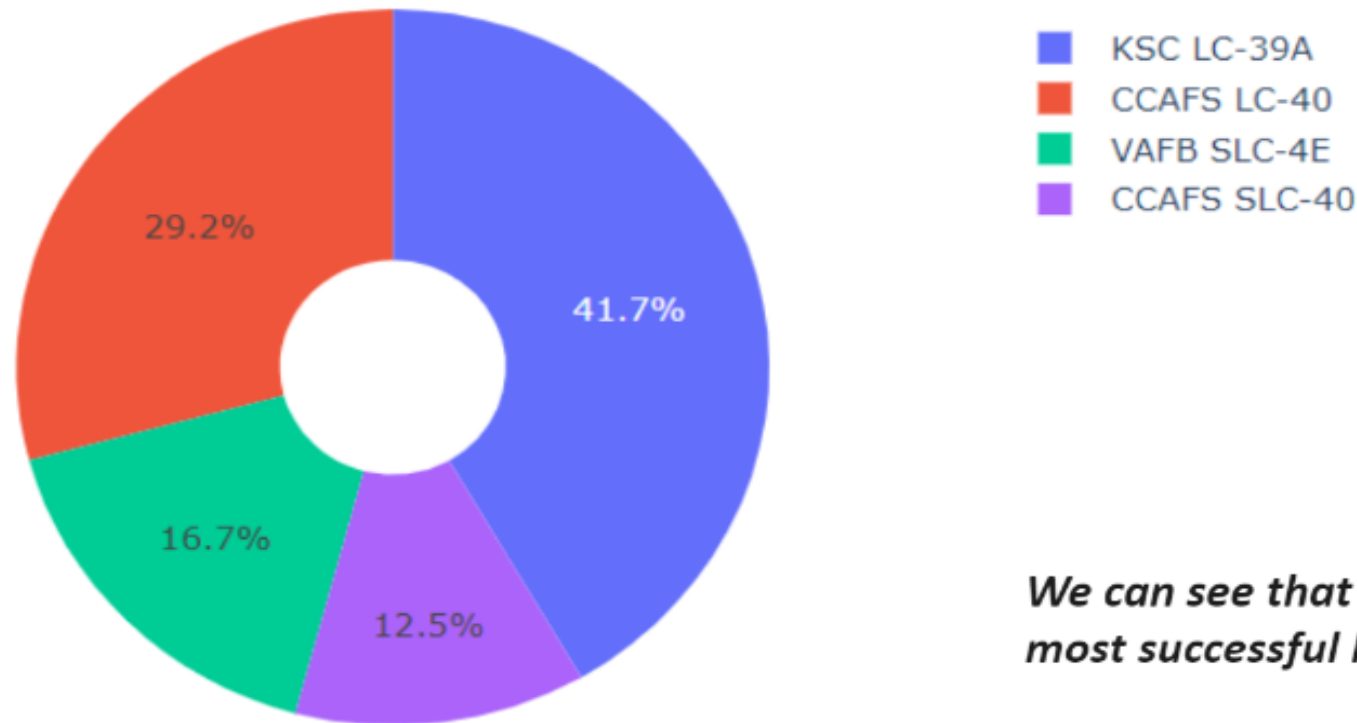
# Payload vs. Orbit Type

- We can observe that with heavy payloads, the successful landing are more for PO, LEO and ISS orbits.



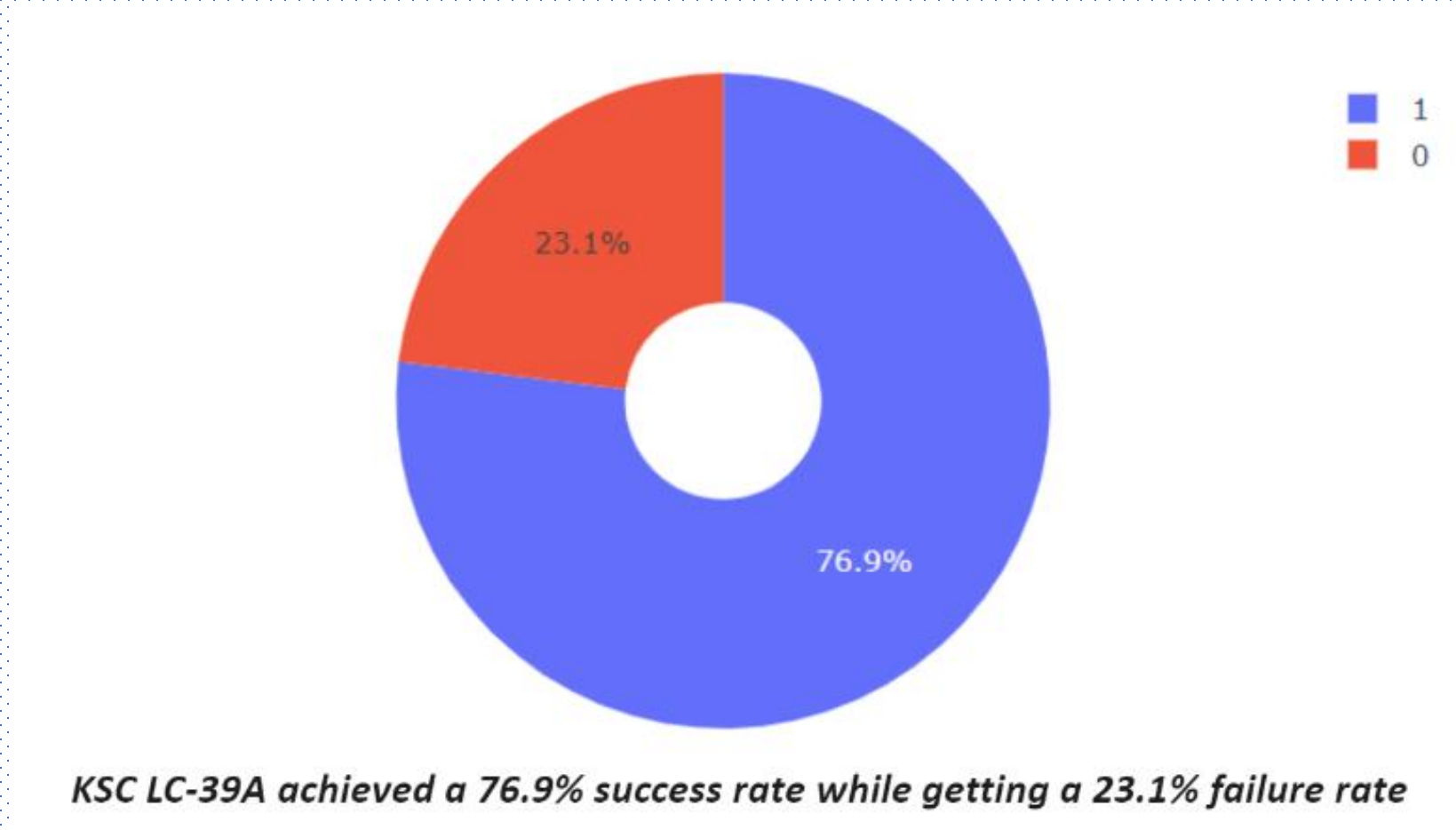
# Pie chart showing the success percentage achieved by each launch site

Total Success Launches By all sites



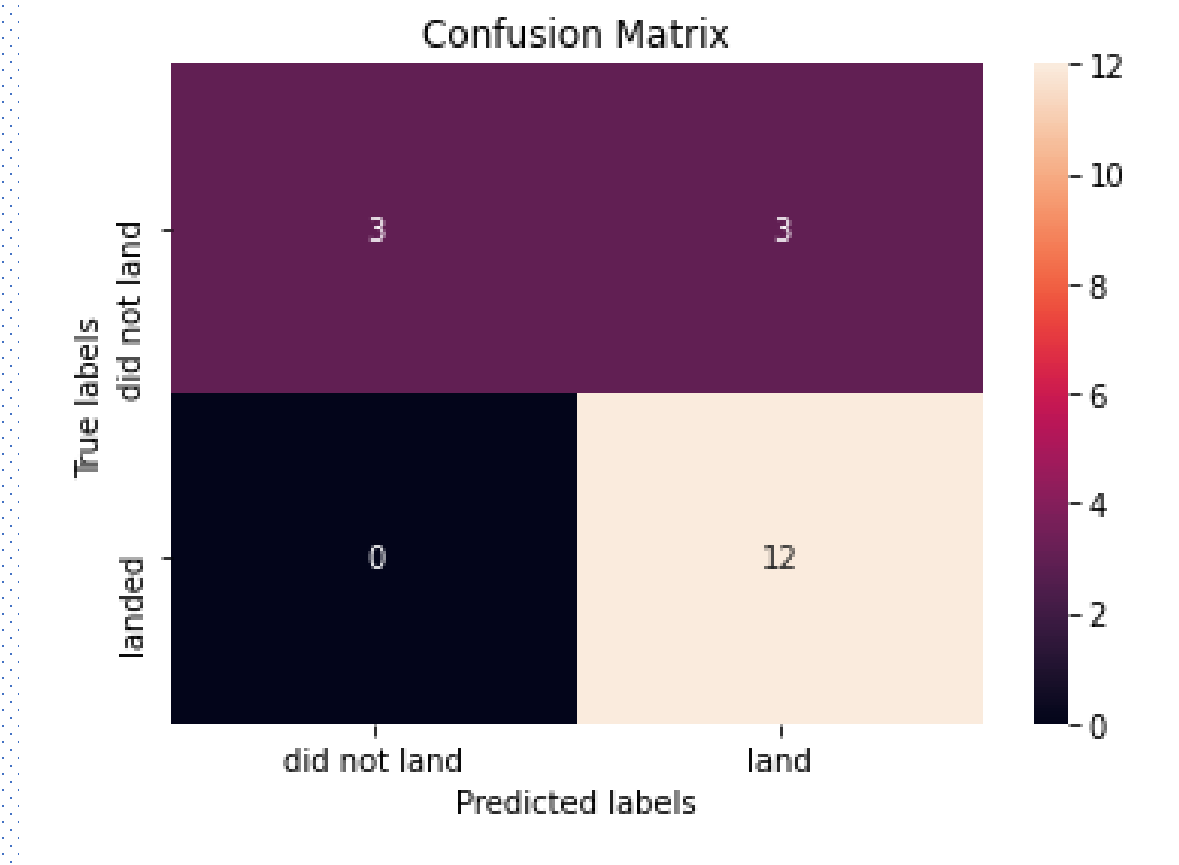
*We can see that KSC LC-39A had the most successful launches from all the sites*

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



# 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.

THANK YOU!!