# IBM Applied Data Science Capstone Project

Presented By:
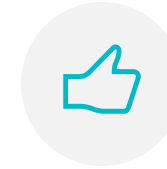
Soumesh khuntia

10th Nov, 2023

# Outline

Executive Summary

Introduction

Methodology

Results

Conclusion

Appendix

# Executive Summary

## Summary of methodologies

- Data Collection
- Data Wrangling
- Exploratory Data Analysis with Data Visualization
- Exploratory Data Analysis with SQL
- Building an interactive map with Folium
- Building a Dashboard with Plotly & Dash
- Predictive Analysis (Classification)

## Summary of all results

- Exploratory Data Analysis
- Interactive Analytics Demo in screenshots
- Predictive Analysis results

# Introduction

SpaceX is the most successful company of the commercial space age, making space travel affordable. The company 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 SpaceX can reuse the first stage. Therefore, if we can determine if the first stage will land, we can determine the cost of a launch. Based on public information and machine learning models, we are going to predict if SpaceX will reuse the first stage.

# Methodology

Data Collection Methodology

- Using SpaceX Rest API

- Using Web Scraping from Wikipedia

Performed data wrangling

- Filtering the data

- Dealing with missing values

- Using One Hot Encoding to prepare the data to a binary classification

Performed exploratory data analysis (EDA) using visualization and SQL Performed interactive visual analytics using Folium and Plotly Dash Performed predictive analysis using classification models - Building, tuning and evaluation of classification models to ensure the best results

# Data Collection

This process involved a combination of API requests from SpaceX **REST API** and **Web Scrapping** data from a table in SpaceX's Wikipedia entry.

We had to use both data collection methods to get complete information about the launches for a more detailed analysis.

Data Columns obtained from **SpaceX REST API**:

FlightNumber, Date, BoosterVersion, PayloadMass, Orbit, LaunchSite, Outcome, Flights, GridFins, Reused, Legs, LandingPad, Block, ReusedCount, Serial, Longitude, Latitude

Data Columns obtained by using **Wikipedia Web Scrapping**:

Flight No., Launch site, Payload, PayloadMass, Orbit, Customer, Launch outcome, Version Booster, Booster landing, Date, Time

# [Data Collection - SpaceX API](#)

- 🚀 Requesting rocket launch date from SpaceX API

- 🗄 Decoding the response content using .json() and turning it into a dataframe using .json_normalize()

- 👤 Requesting needed information about the launches from SpaceX API by applying custom functions

- 📚 Constructing data we have obtained into a dictionary

- 📖 Creating a dataframe from the dictionary

- 🐦 Filtering the dataframe to only include Falcon-9 launches

- 🔍 Replacing missing values of payload mass columns with calculated .mean()

- ⬇ Exporting the data to CSV

# Data Collection – Web Scrapping

| 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|
| Requesting Falcon-9 launch from Wikipedia | Creating a **BeautifulSoup** object from the HTML response | Extracting all column names from the HTML table header | Collecting the data by parsing HTML tables | Constructing data we've obtained into a dictionary | Creating a dataframe from the dictionary | Exporting the data into CSV |

# Data Wrangling

- Perform **Exploratory Data Analysis** and determine training labels

- Calculate the number of launches on each site

- Calculate the number and occurrence of each orbit

- Calculate the number and occurrence of mission outcome per orbit type

- Creating a landing outcome label from **Outcome column**

- Exporting the data to CSV

# EDA with Data Visualization

Charts were plotted: **Flight Number vs. Payload Mass**, **Flight Number vs. Launch Site**, **Payload Mass vs. Launch Site**, **Orbit Type vs. Success Rate**, **Flight Number vs. Orbit Type**, **Payload Mass vs Orbit Type** and Success Rate Yearly Trend Scatter plots show the relationship between variables.

If a relationship exists, they could be used in machine learning model. Bar charts show comparisons among discrete categories. The goal is to show the relationship between the specific categories being compared and a measured value. Line charts show trends in data over time (time series)

# EDA with SQL

Performed SQL queries:

• Displaying the names of the unique launch sites in the space mission

• Displaying 5 records where launch sites begin with the string 'CCA'

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

• Displaying average payload mass carried by booster version F9 v1.1

• Listing the date when the first successful landing outcome in ground pad was achieved

• Listing the names of the boosters which have success in drone ship and have payload mass greater than 4000 but less than 6000

• Listing the total number of successful and failure mission outcomes

• Listing the names of the booster versions which have carried the maximum payload mass

• Listing the failed landing outcomes in drone ship, their booster versions and launch site names for the months in year 2015

• Ranking 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
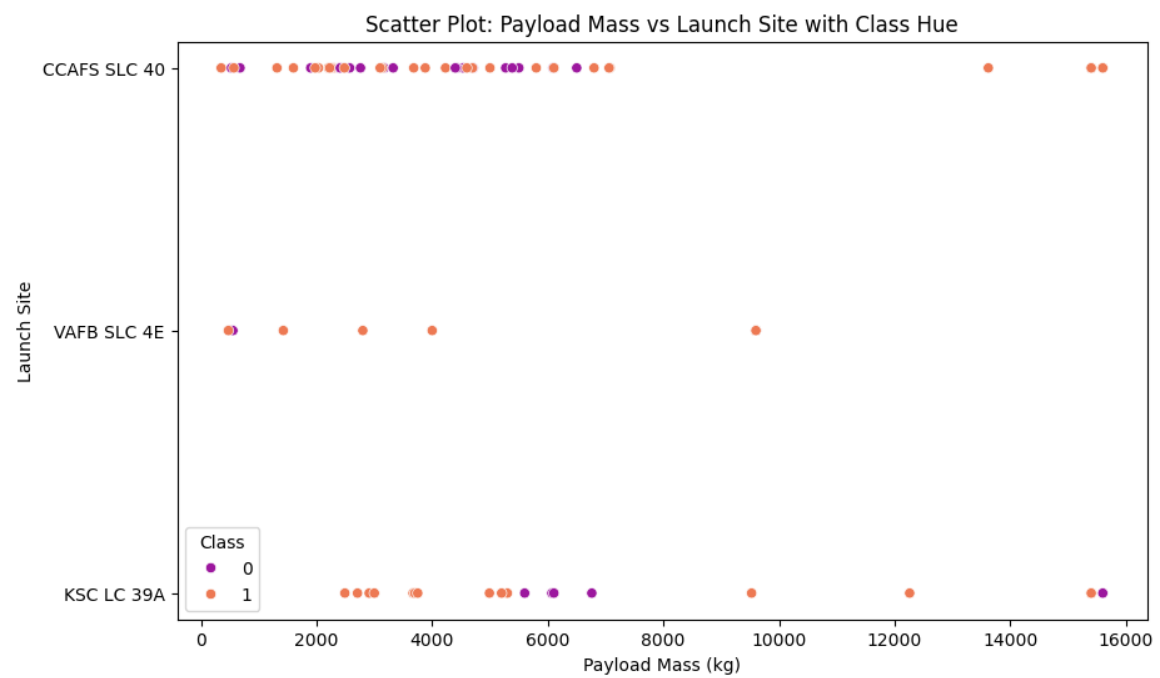
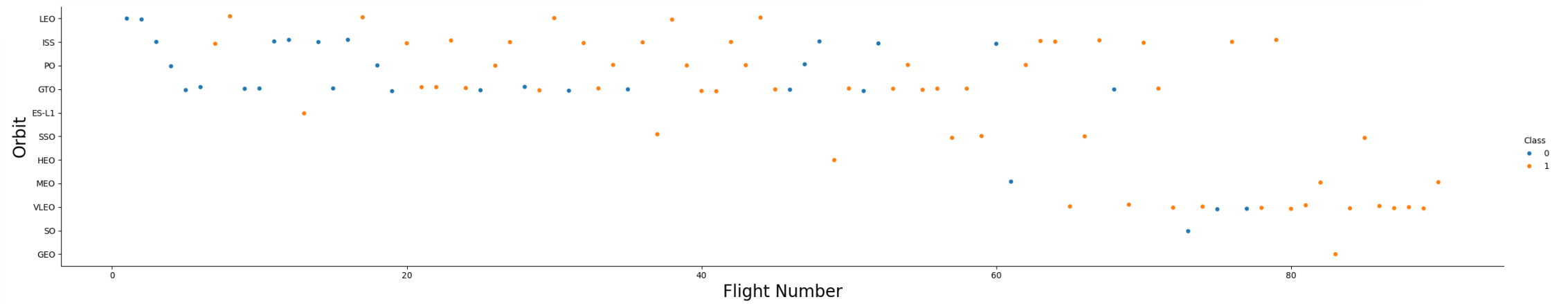# Interactive Visual Analytics using Folium, Plotly and Dash

Markers of all Launch Sites: - Added Marker with Circle, Popup Label and Text Label of NASA Johnson Space Center using its latitude and longitude coordinates as a start location. - Added Markers with Circle, Popup Label and Text Label of all Launch Sites using their latitude and longitude coordinates to show their geographical locations and proximity to Equator and coasts. Coloured Markers of the launch outcomes for each Launch Site: - Added coloured Markers of success (Green) and failed (Red) launches using Marker Cluster to identify which launch sites have relatively high success rates. Distances between a Launch Site to its proximities: - Added coloured Lines to show distances between the Launch Site KSC LC-39A (as an example) and its proximities like Railway, Highway, Coastline and Closest City
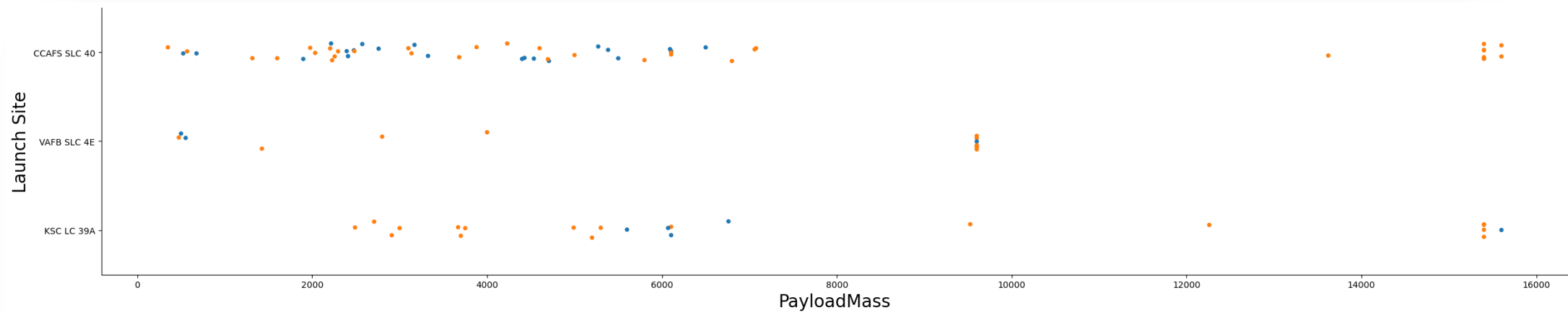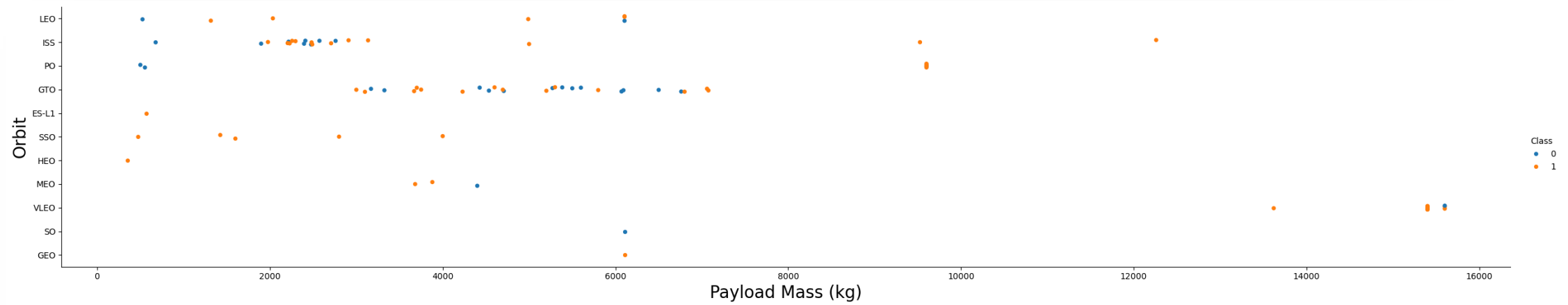
# Results

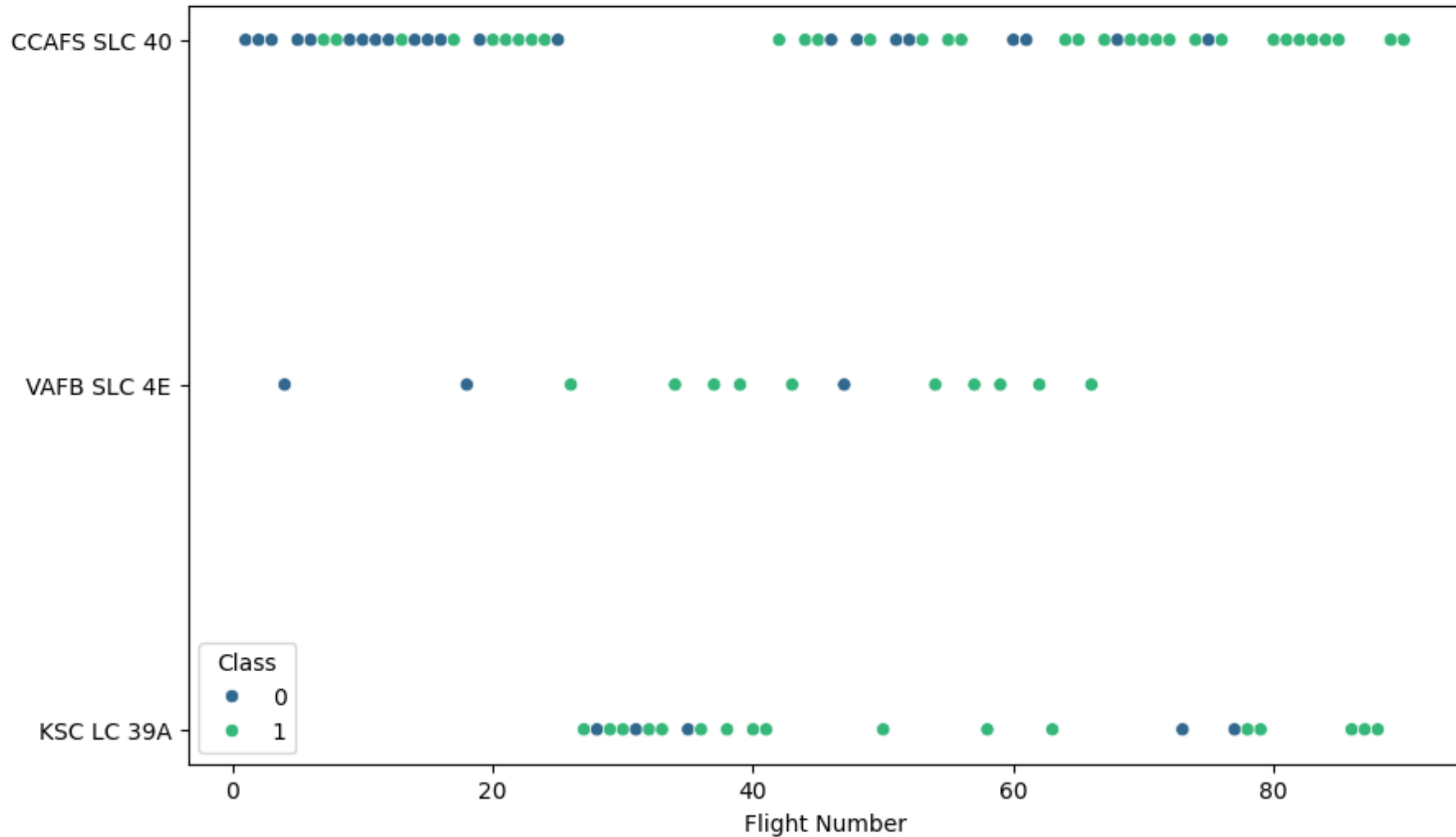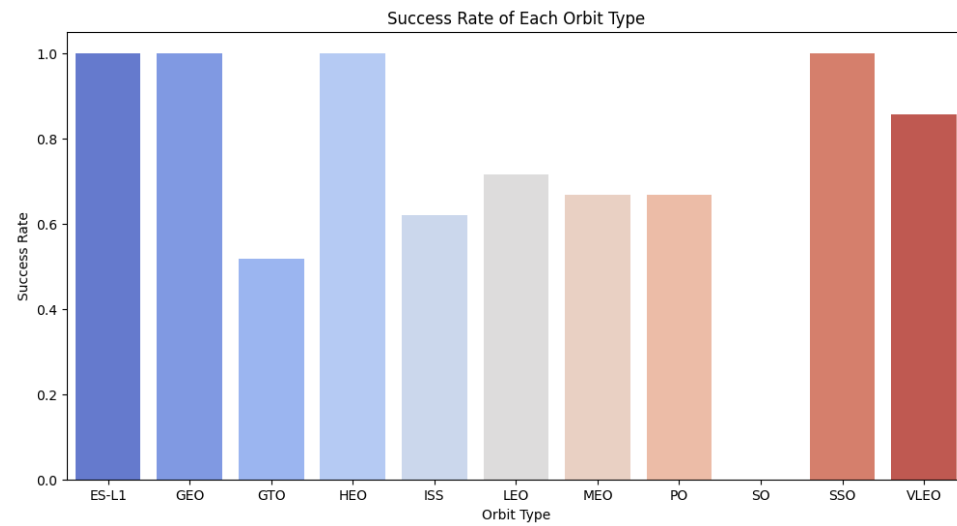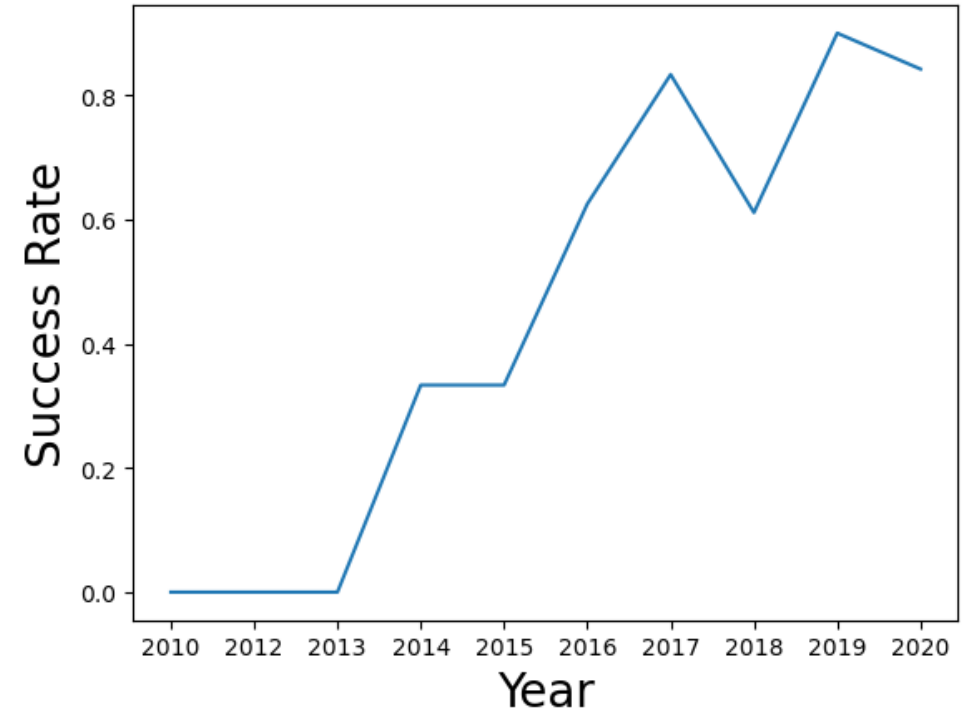Scatter Plot: Payload Mass vs Launch Site with Class Hue

Scatter Plot: Flight Number vs Launch Site with Class Hue
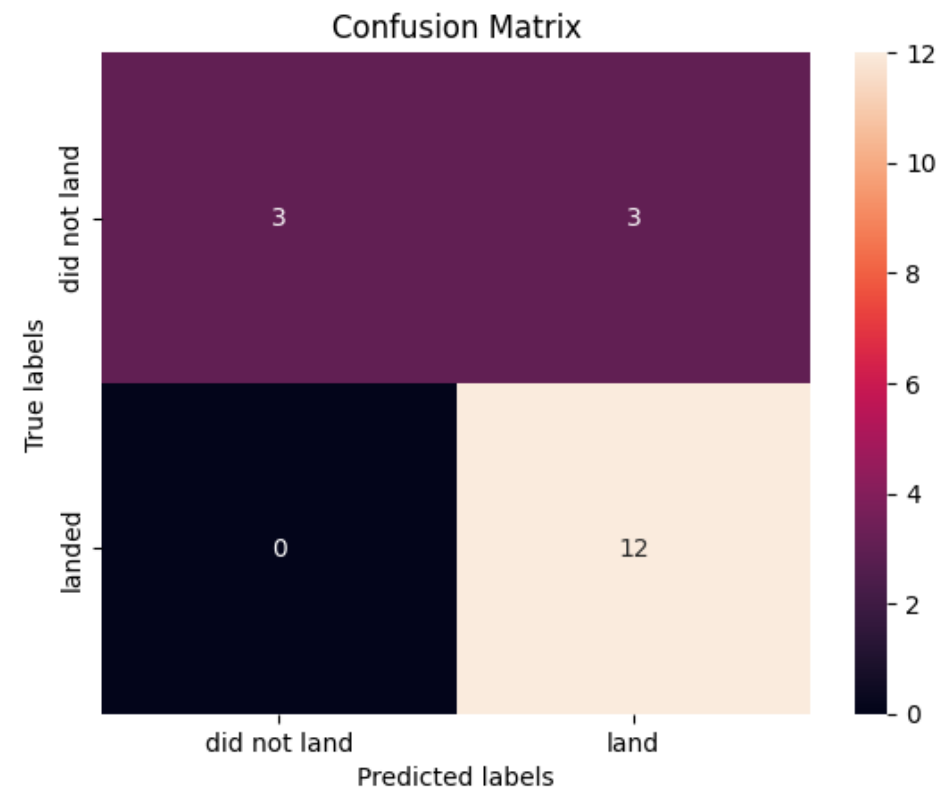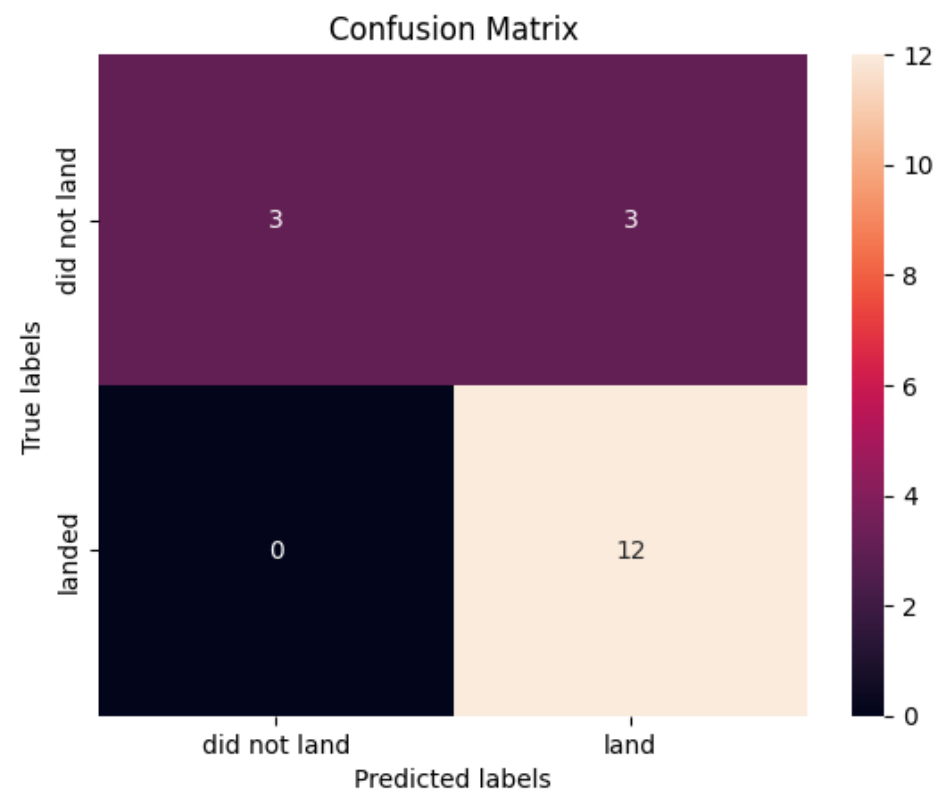
Success Rate of Each Orbit Type

# Success Rate

# Confusion Matrix

SVM

# Confusion Matrix

Decision Tree
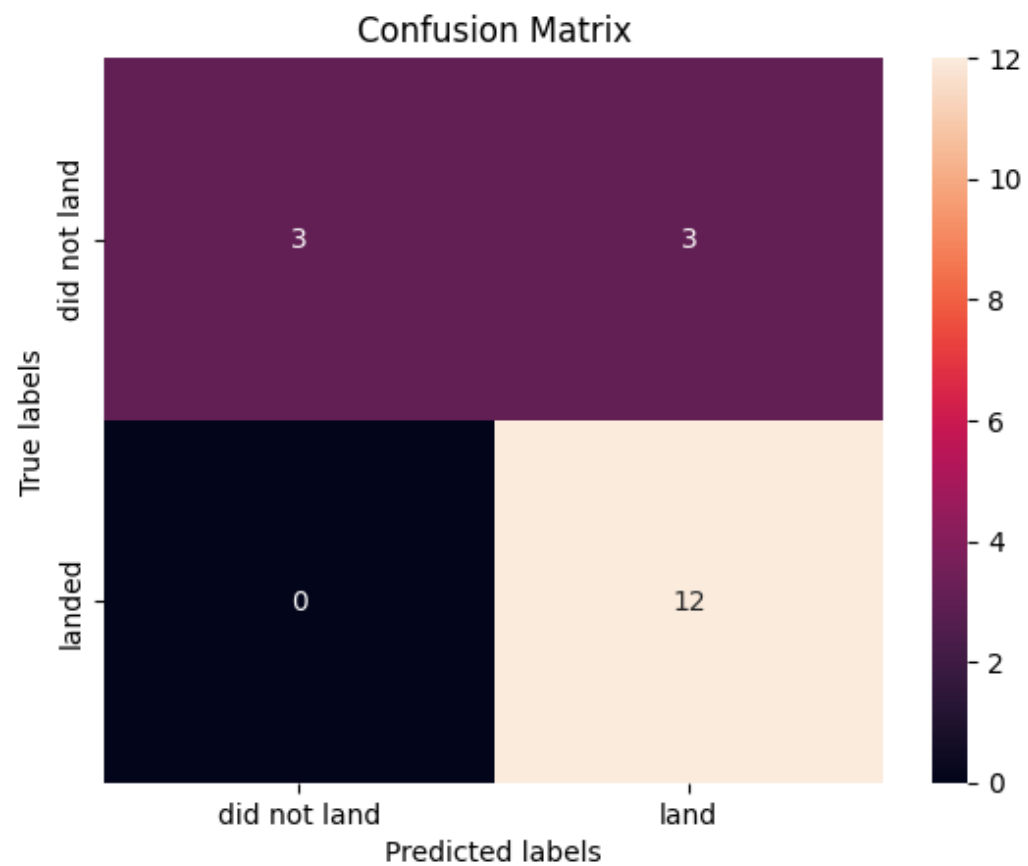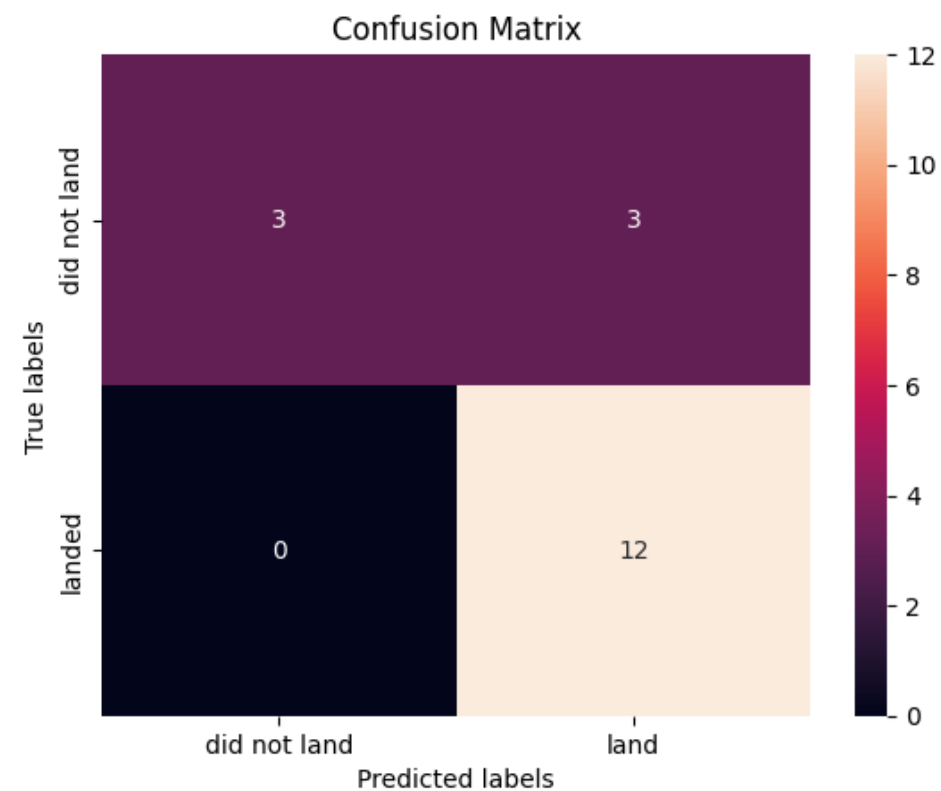


Confusion Matrix

# Confusion Matrix

KNN

# Confusion Matrix

Logistic Regression

# Conclusion

- Decision Tree Model is the best algorithm for this dataset.

- Launches with a low payload mass show better results than launches with a larger payload mass.

- Most of launch sites are in proximity to the Equator line and all the sites are in very close proximity to the coast.

- The success rate of launches increases over the years.

- KSC LC-39A has the highest success rate of the launches from all the sites.

- Orbits ES-L1, GEO, HEO and SSO have 100% success rate