



IBM Developer
SKILLS NETWORK

Winning Space Race with Data Science

Ganraj Puyad
21/05/25



Outline

- Executive Summary
- Introduction
- Methodology
- Results
- Discussion
- Conclusion

Executive Summary

- In this capstone project, we will predict if the SpaceX Falcon 9 first stage will land successfully using several machine learning classification algorithms.
- The main steps in this project include:
 - Data collection, wrangling, and formatting
 - Exploratory data analysis
 - Interactive data visualization
 - Machine learning prediction
- Our graphs show that some features of the rocket launches have a correlation with the outcome of the launches, i.e., success or failure.
- It is also concluded that decision tree may be the best machine learning algorithm to predict if the Falcon 9 first stage will land successfully.

Introduction

- In this capstone, we will predict if the Falcon 9 first stage will land successfully. SpaceX 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. This information can be used if an alternate company wants to bid against SpaceX for a rocket launch.
- Most unsuccessful landings are planned. Sometimes, SpaceX will perform a controlled landing in the ocean.
- The main question that we are trying to answer is, for a given set of features about a Falcon 9 rocket launch which include its payload mass, orbit type, launch site, and so on, will the first stage of the rocket land successfully?

Section 1

Methodology

Methodology

- The overall methodology includes:
 1. Data collection, wrangling, and formatting, using:
 - SpaceX API
 - Web scraping
 2. Exploratory data analysis (EDA), using:
 - Pandas and NumPy
 - SQL
 3. Data visualization, using:
 - Matplotlib and Seaborn
 - Folium
 - Dash
 4. Machine learning prediction, using
 - Logistic regression
 - Support vector machine (SVM)
 - Decision tree
 - K-nearest neighbors (KNN)

METHODOLOGY

Data collection, wrangling, and formatting

- SpaceX API
 - The API used is <https://api.spacexdata.com/v4/rockets/>.
 - The API provides data about many types of rocket launches done by SpaceX, the data is therefore filtered to include only Falcon 9 launches.
 - Every missing value in the data is replaced the mean the column that the missing value belongs to.
 - We end up with 90 rows or instances and 17 columns or features. The picture below shows the first few rows of the data:

	FlightNumber	Date	BoosterVersion	PayloadMass	Orbit	LaunchSite	Outcome	Flights	GridFins	Reused	Legs		LandingPad	Block	ReusedCount	Serial	Longitude	Latitude
4	1	2010-06-04	Falcon 9	NaN	LEO	CCSFS SLC 40	None None	1	False	False	False		None	1.0	0	B0003	-80.577366	28.561857
5	2	2012-05-22	Falcon 9	525.0	LEO	CCSFS SLC 40	None None	1	False	False	False		None	1.0	0	B0005	-80.577366	28.561857
6	3	2013-03-01	Falcon 9	677.0	ISS	CCSFS SLC 40	None None	1	False	False	False		None	1.0	0	B0007	-80.577366	28.561857
7	4	2013-09-29	Falcon 9	500.0	PO	VAFB SLC 4E	False Ocean	1	False	False	False		None	1.0	0	B1003	-120.610829	34.632093
8	5	2013-12-03	Falcon 9	3170.0	GTO	CCSFS SLC 40	None None	1	False	False	False		None	1.0	0	B1004	-80.577366	28.561857

METHODOLOGY

Data collection, wrangling, and formatting

- Web scraping
 - The data is scraped from [https://en.wikipedia.org/w/index.php?title=List of Falcon 9 and Falcon Heavy launches&oldid=1027686922](https://en.wikipedia.org/w/index.php?title=List_of_Falcon_9_and_Falcon_Heavy_launches&oldid=1027686922)
 - The website contains only the data about Falcon 9 launches.
 - We end up with 121 rows or instances and 11 columns or features. The picture below shows the first few rows of the data:

	Flight No.	Launch site	Payload	Payload mass	Orbit	Customer	Launch outcome	Version Booster	Booster landing	Date	Time
0	1	CCAFS	Dragon Spacecraft Qualification Unit	0	LEO	SpaceX	Success\n	F9 v1.0B0003.1	Failure	4 June 2010	18:45
1	2	CCAFS	Dragon	0	LEO	NASA	Success	F9 v1.0B0004.1	Failure	8 December 2010	15:43
2	3	CCAFS	Dragon	525 kg	LEO	NASA	Success	F9 v1.0B0005.1	No attempt\n	22 May 2012	07:44
3	4	CCAFS	SpaceX CRS-1	4,700 kg	LEO	NASA	Success\n	F9 v1.0B0006.1	No attempt	8 October 2012	00:35
4	5	CCAFS	SpaceX CRS-2	4,877 kg	LEO	NASA	Success\n	F9 v1.0B0007.1	No attempt\n	1 March 2013	15:10

METHODOLOGY

Data collection, wrangling, and formatting

- The data is later processed so that there are no missing entries and categorical features are encoded using one-hot encoding.
- An extra column called 'Class' is also added to the data frame. The column 'Class' contains 0 if a given launch is failed and 1 if it is successful.
- In the end, we end up with 90 rows or instances and 83 columns or features.

METHODOLOGY

Exploratory Data Analysis (EDA)



- Pandas and NumPy

- Functions from the Pandas and NumPy libraries are used to derive basic information about the data collected, which includes:
 - The number of launches on each launch site
 - The number of occurrence of each orbit
 - The number and occurrence of each mission outcome



- SQL

- The data is queried using SQL to answer several questions about the data such as:
 - The names of the 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

METHODOLOGY

Data Visualization

- Matplotlib and Seaborn



- Functions from the Matplotlib and Seaborn libraries are used to visualize the data through scatterplots, bar charts, and line charts.
- The plots and charts are used to understand more about the relationships between several features, such as:
 - The relationship between flight number and launch site
 - The relationship between payload mass and launch site
 - The relationship between success rate and orbit type

- Folium



- Functions from the Folium libraries are used to visualize the data through interactive maps.
- The Folium library is used to:
 - Mark all launch sites on a map
 - Mark the succeeded launches and failed launches for each site on the map
 - Mark the distances between a launch site to its proximities such as the nearest city, railway, or highway

METHODOLOGY

Data Visualization



- Dash
 - Functions from Dash are used to generate an interactive site where we can toggle the input using a dropdown menu and a range slider.
 - Using a pie chart and a scatterplot, the interactive site shows:
 - The total success launches from each launch site
 - The correlation between payload mass and mission outcome (success or failure) for each launch site

METHODOLOGY

Machine Learning Prediction

- Functions from the Scikit-learn library are used to create our machine learning models.
- The machine learning prediction phase include the following steps:
 - Standardizing the data
 - Splitting the data into training and test data
 - Creating machine learning models, which include:
 - Logistic regression
 - Support vector machine (SVM)
 - Decision tree
 - K nearest neighbors (KNN)
 - Fit the models on the training set
 - Find the best combination of hyperparameters for each model
 - Evaluate the models based on their accuracy scores and confusion matrix



Results

- The results are split into 5 sections:
 - SQL (EDA with SQL)
 - Matplotlib and Seaborn (EDA with Visualization)
 - Folium
 - Dash
 - Predictive Analysis
- In all of the graphs that follow, class 0 represents a failed launch outcome while class 1 represents a successful launch outcome.

The background of the slide is an abstract composition. It features a dark blue base color. Overlaid on this are numerous diagonal streaks in shades of red and cyan. A faint, light blue grid pattern is also visible, particularly in the lower-left quadrant. The overall effect is dynamic and technological.

Section 2

Insights drawn from EDA

RESULTS

SQL (EDA with SQL)

- The names of the unique launch sites in the space mission
 - CCAFS LC-40
 - VAFB SLC-4E
 - KSC LC-39A
 - CCAFS SLC-40
- 5 records where launch sites begin with 'CCA'

Records:

```
('06/04/2010', '18:45:00', 'F9 v1.0 B0003', 'CAFS LC-40', 'Dragon Spacecraft Qualification Unit', 0.0, 'LEO', 'SpaceX', 'Success', 'Failure (parachute)')
('12/08/2010', '15:43:00', 'F9 v1.0 B0004', 'CAFS LC-40', 'Dragon demo flight C1, two CubeSats, barrel of Brouere cheese', 0.0, 'LEO (ISS)', 'NASA (COTS) NRO', 'Success', 'Failure (parachute)')
('22/05/2012', '7:44:00', 'F9 v1.0 B0005', 'CAFS LC-40', 'Dragon demo flight C2', 525.0, 'LEO (ISS)', 'NASA (COTS)', 'Success', 'No attempt')
('10/08/2012', '0:35:00', 'F9 v1.0 B0006', 'CAFS LC-40', 'SpaceX CRS-1', 500.0, 'LEO (ISS)', 'NASA (CRS)', 'Success', 'No attempt')
('03/01/2013', '15:10:00', 'F9 v1.0 B0007', 'CAFS LC-40', 'SpaceX CRS-2', 677.0, 'LEO (ISS)', 'NASA (CRS)', 'Success', 'No attempt')
```

RESULTS

SQL (EDA with SQL)

- The total payload mass carried by boosters launched by NASA (CRS)

Total Payload Mass carried by boosters launched by NASA (CRS): 45596.0 kg

- The average payload mass carried by booster version F9 v1.1

Average Payload Mass Carried by F9 v1.1 Boosters: 2928.4 kg

- The date when the first successful landing outcome in ground pad was achieved

Date of the First Successful Landing in a Ground Pad: 01/08/2018

RESULTS

SQL (EDA with SQL)

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

Boosters with a Successful Landing on a Drone Ship and Payload Mass between 4000 kg and 6000 kg:

F9 FT B1022

F9 FT B1026

F9 FT B1021.2

F9 FT B1031.2

- The total number of successful and failure mission outcomes

Total Number of Successful Mission Outcomes: 98

Total Number of Failure Mission Outcomes: 0

RESULTS

SQL (EDA with SQL)

- The names of the booster versions which have carried the maximum payload mass

Booster Versions with Maximum Payload Mass:

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

RESULTS

SQL (EDA with SQL)

- The failed landing outcomes in drone ship, their booster versions, and launch site names for in year 2015

Records:

Month: October

Landing Outcome in Drone Ship: Failure (drone ship)

Booster Version: F9 v1.1 B1012

Launch Site: CCAFS LC-40

Month: November

Landing Outcome in Drone Ship: Controlled (ocean)

Booster Version: F9 v1.1 B1013

Launch Site: CCAFS LC-40

Month: February

Landing Outcome in Drone Ship: No attempt

Booster Version: F9 v1.1 B1014

Launch Site: CCAFS LC-40

RESULTS

SQL (EDA with SQL)

Month: April

Landing Outcome in Drone Ship: Failure (drone ship)

Booster Version: F9 v1.1 B1015

Launch Site: CCAFS LC-40

Month: April

Landing Outcome in Drone Ship: No attempt

Booster Version: F9 v1.1 B1016

Launch Site: CCAFS LC-40

Month: June

Landing Outcome in Drone Ship: Precluded (drone ship)

Booster Version: F9 v1.1 B1018

Launch Site: CCAFS LC-40

RESULTS

SQL (EDA with SQL)

Month: December

Landing Outcome in Drone Ship: Success (ground pad)

Booster Version: F9 FT B1019

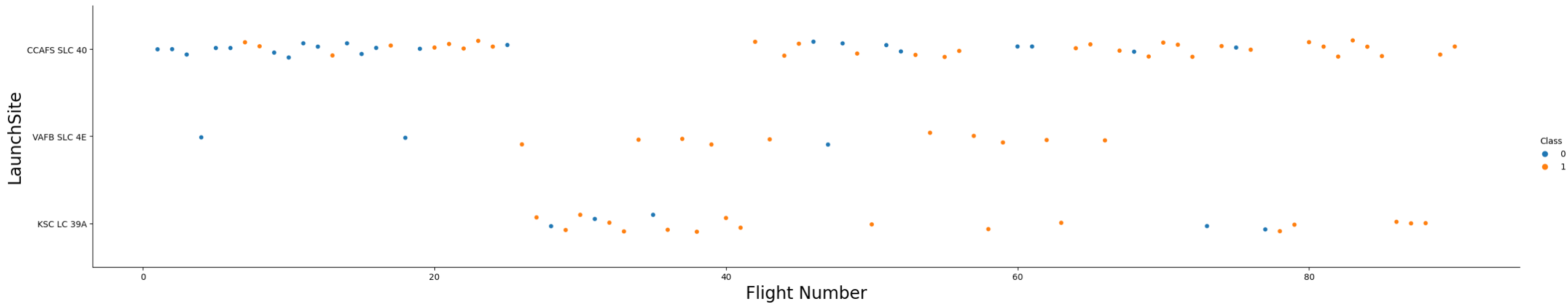
Launch Site: CCAFS LC-40

-
- The count of landing outcomes between the date 2010-06-04 and 2017-03-20, in descending order

RESULTS

Matplotlib and Seaborn (EDA with Visualization)

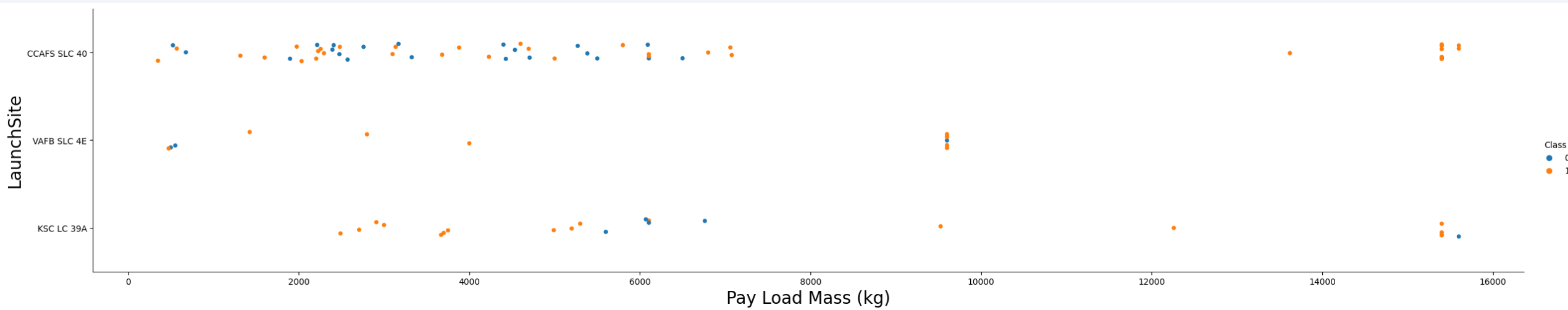
- The relationship between flight number and launch site



RESULTS

Matplotlib and Seaborn (EDA with Visualization)

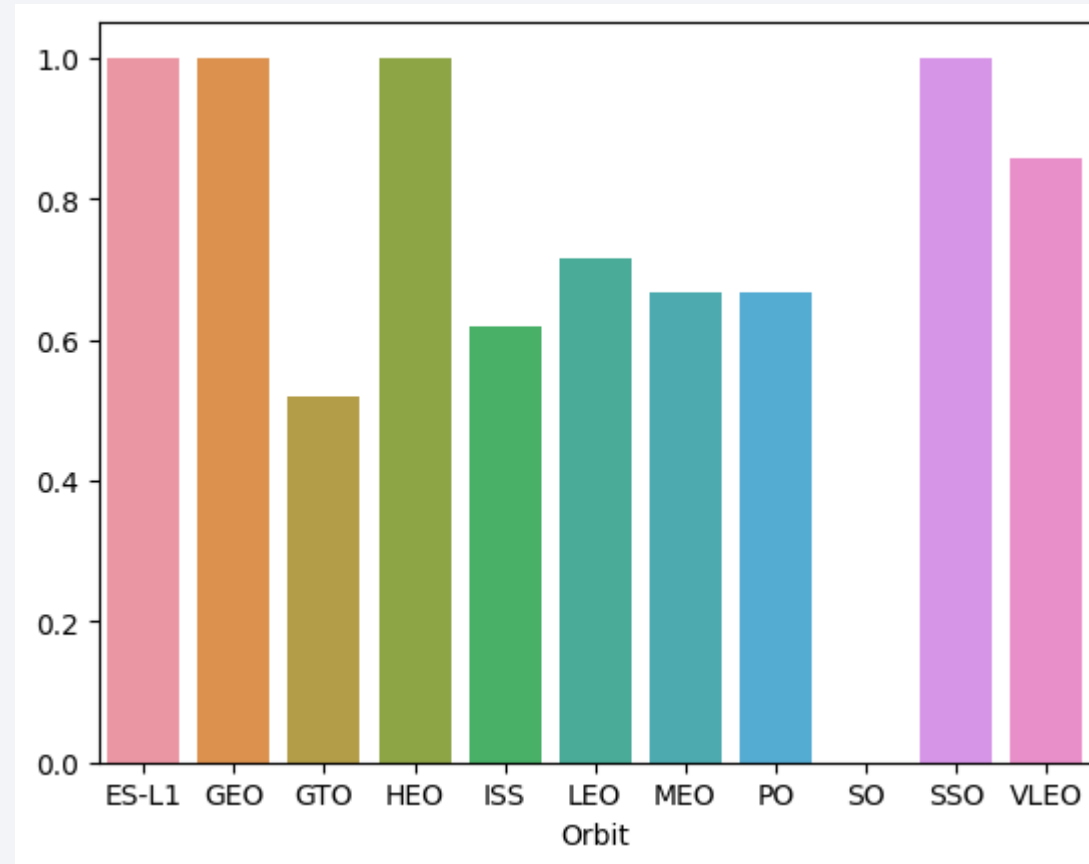
- The relationship between payload mass and launch site



RESULTS

Matplotlib and Seaborn (EDA with Visualization)

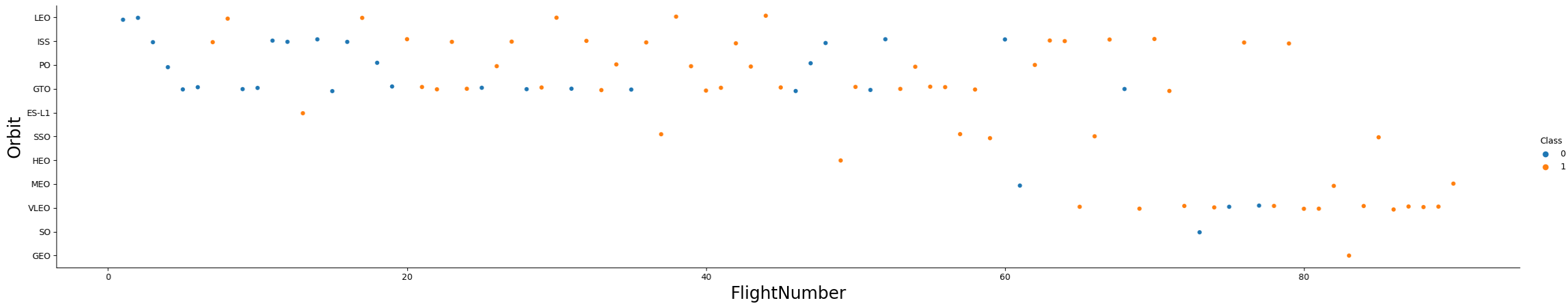
- The relationship between success rate and orbit type



RESULTS

Matplotlib and Seaborn (EDA with Visualization)

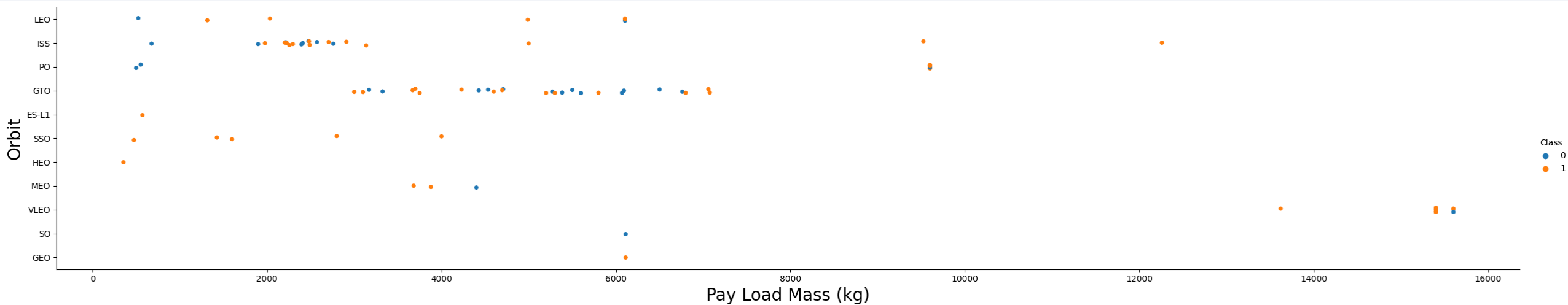
- The relationship between flight number and orbit type



RESULTS

Matplotlib and Seaborn (EDA with Visualization)

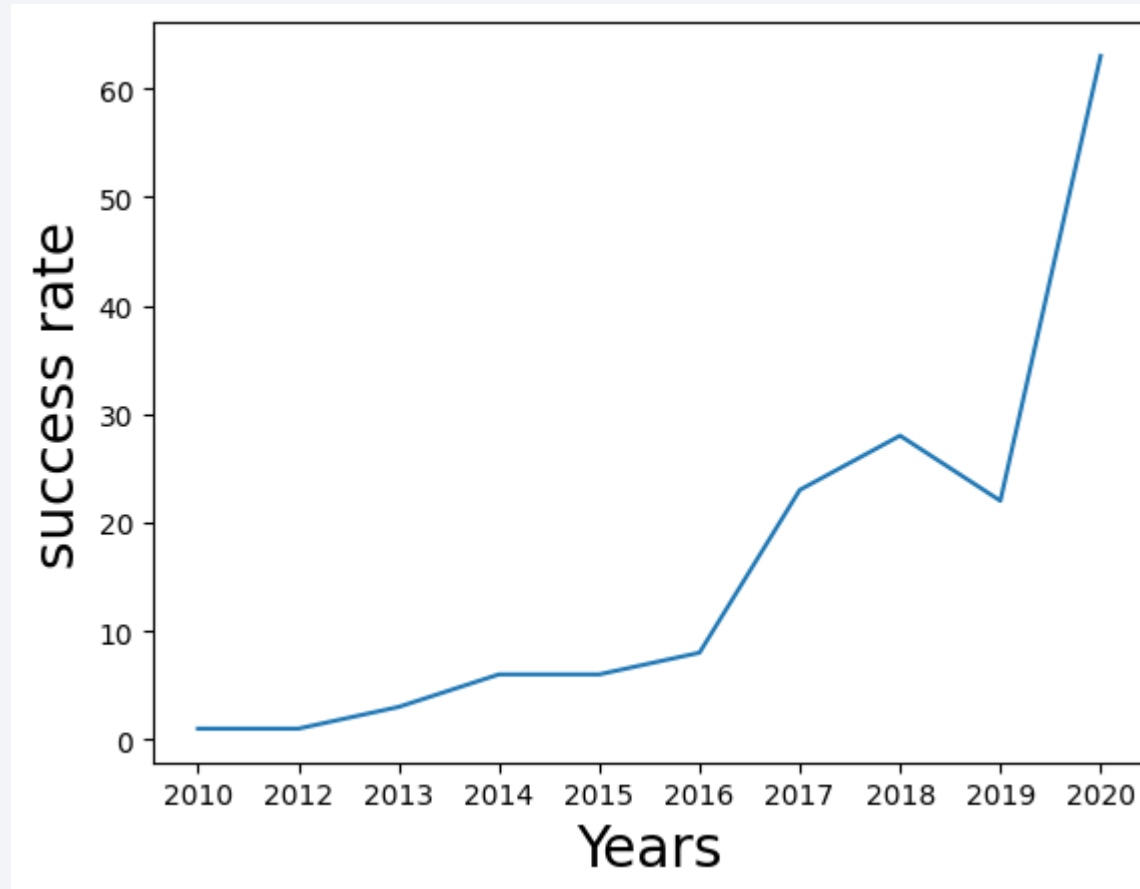
- The relationship between payload mass and orbit type



RESULTS

Matplotlib and Seaborn (EDA with Visualization)

- The launch success yearly trend



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 rectangle on the left and a satellite photograph of Earth on the right. The Earth's surface is dark, with numerous bright yellow and orange lights representing cities and urban areas. The horizon of the Earth is visible, separating the dark surface from the deep blue of the atmosphere and the blackness of space.

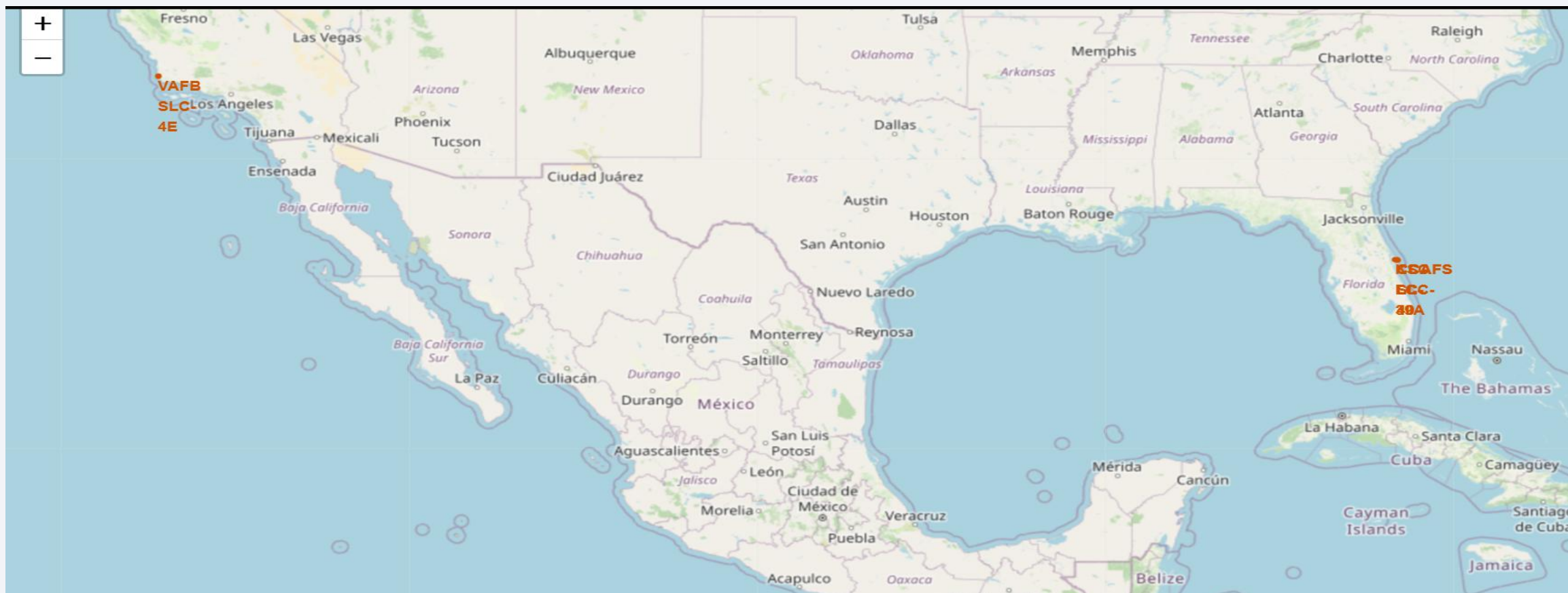
Section 3

Launch Sites Proximities Analysis

RESULTS

Folium

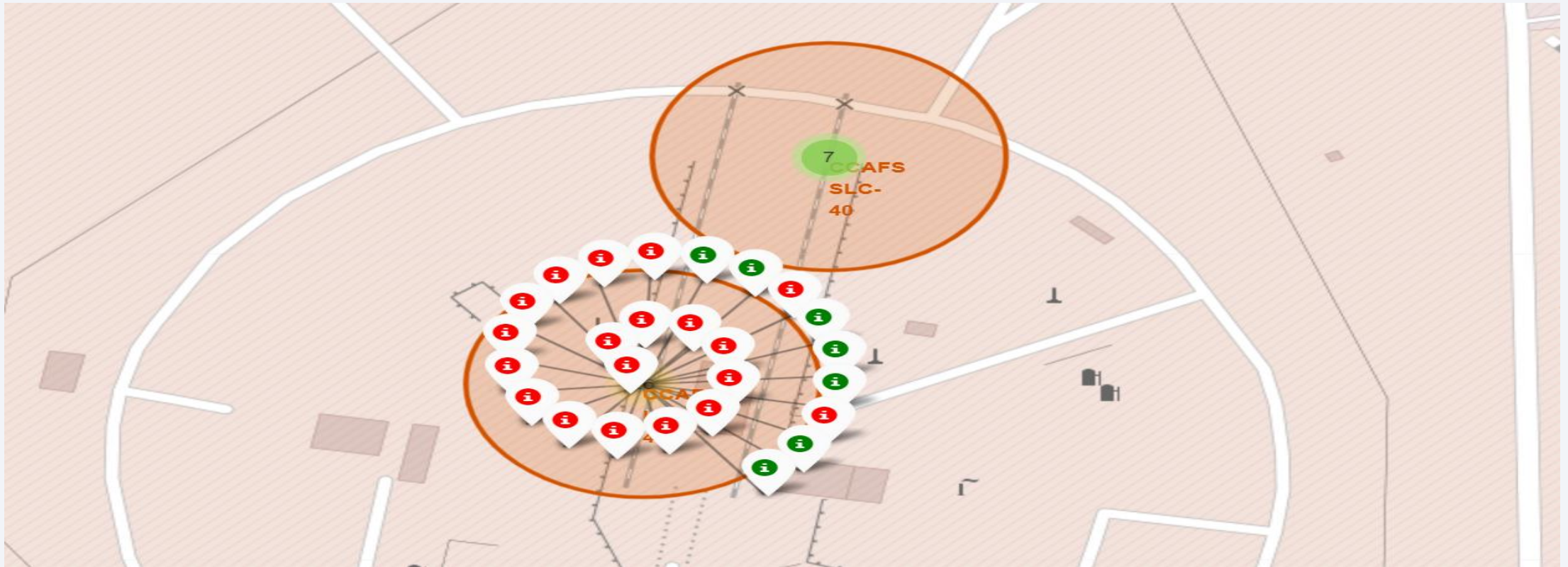
- All launch sites on map



RESULTS

Folium

- The succeeded launches and failed launches for each site on map
 - If we zoom in on one of the launch site, we can see green and red tags. Each green tag represents a successful launch while each red tag represents a failed launch



RESULTS

Folium

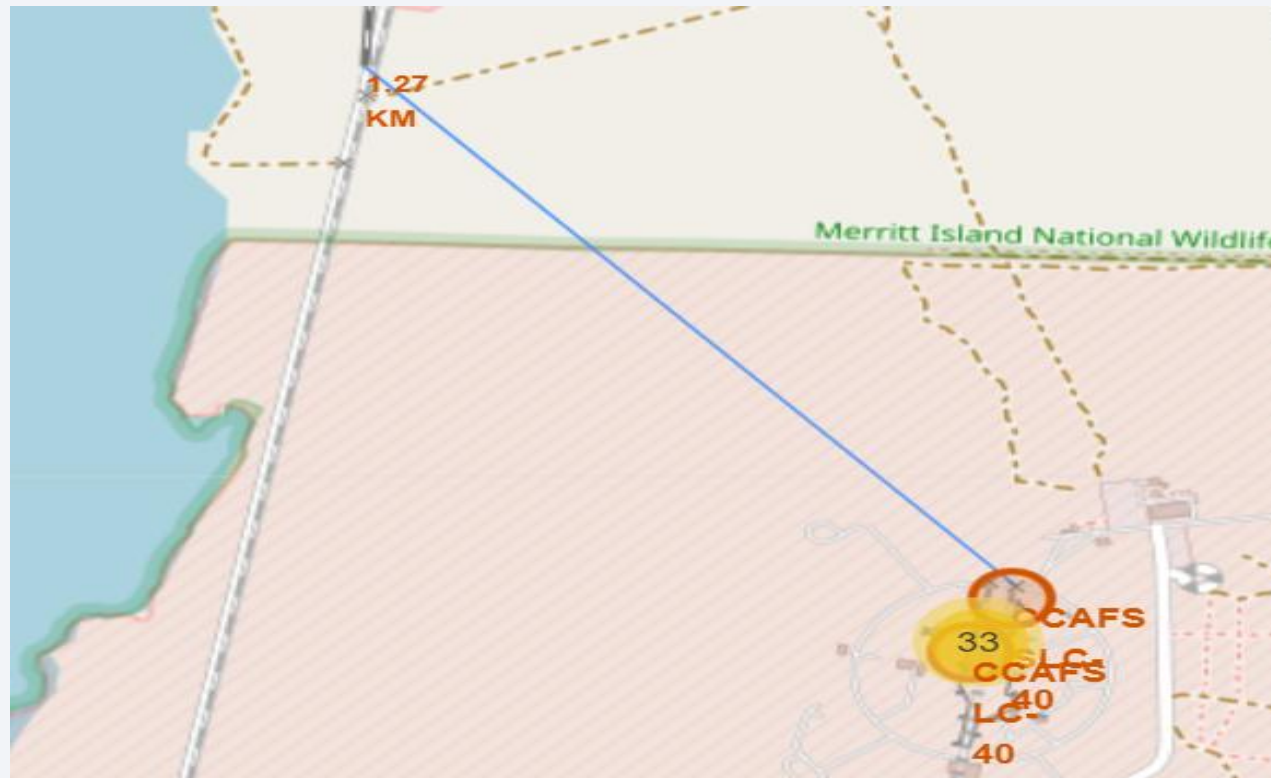
- The distances between a launch site to its proximities such as the nearest city, railway, **coastline** or highway



RESULTS

Folium

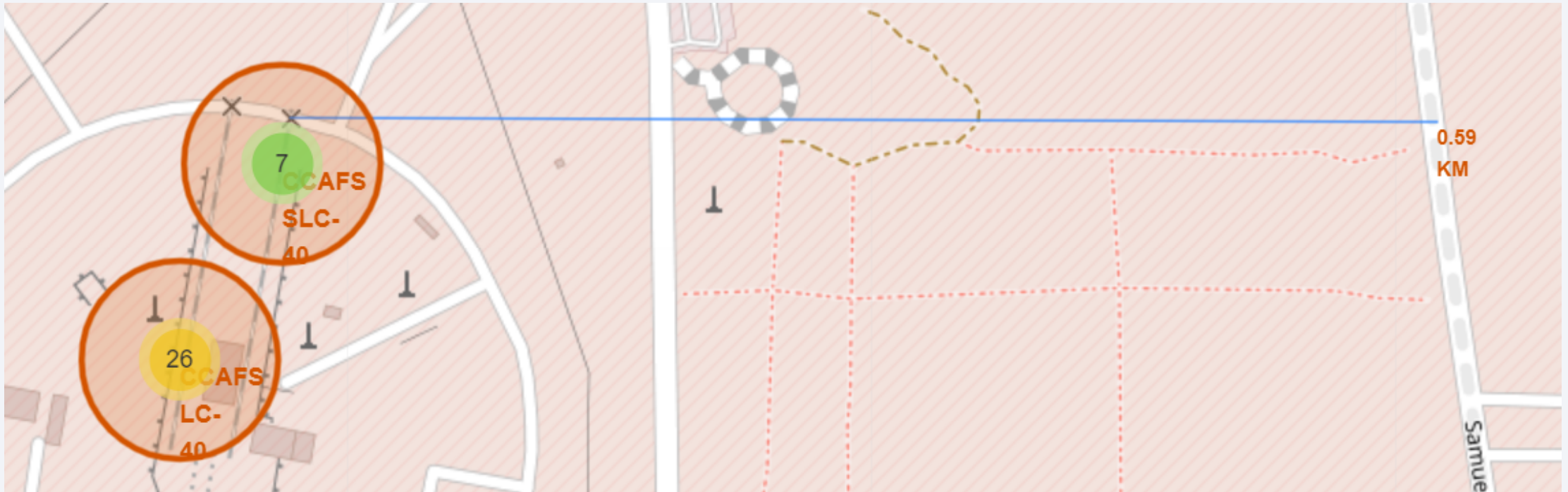
- The distances between a launch site to its proximities such as the nearest city, **railway**, or highway



RESULTS

Folium

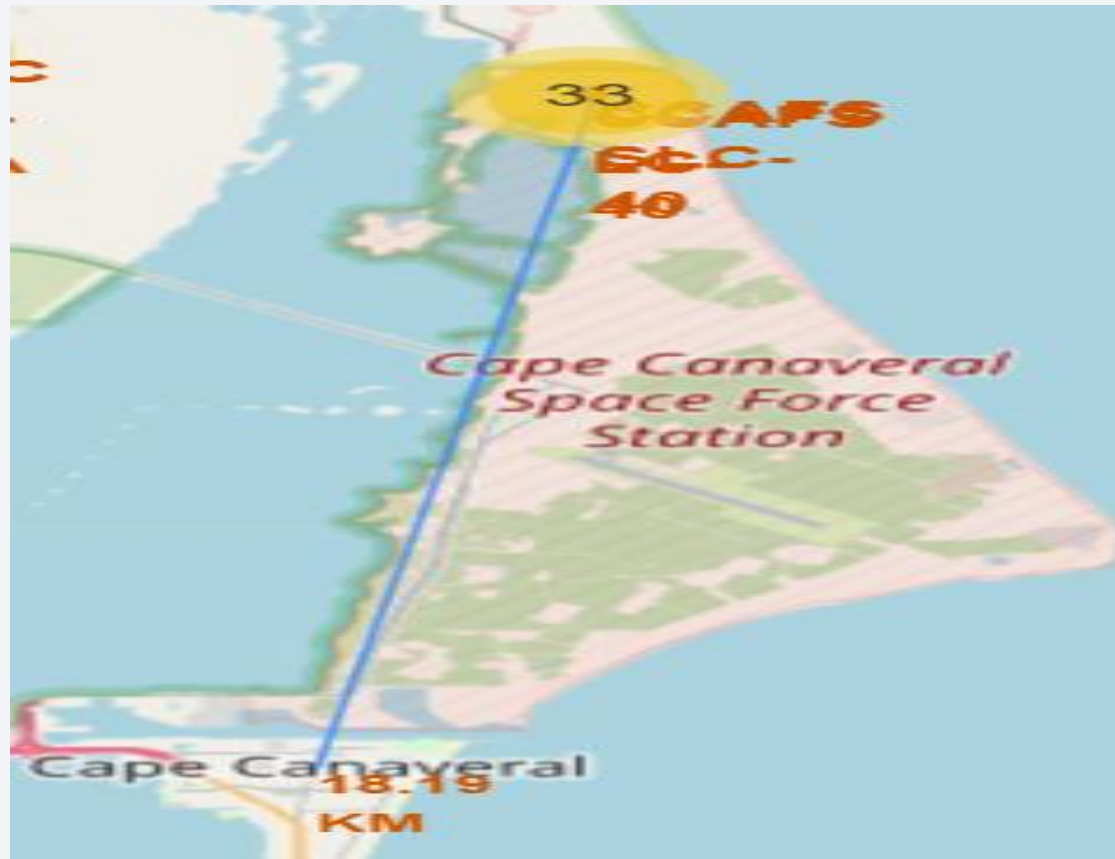
- The distances between a launch site to its proximities such as the nearest city, railway, or **highway**



RESULTS

Folium

- The distances between a launch site to its proximities such as the nearest **city**, railway, or highway





Section 4

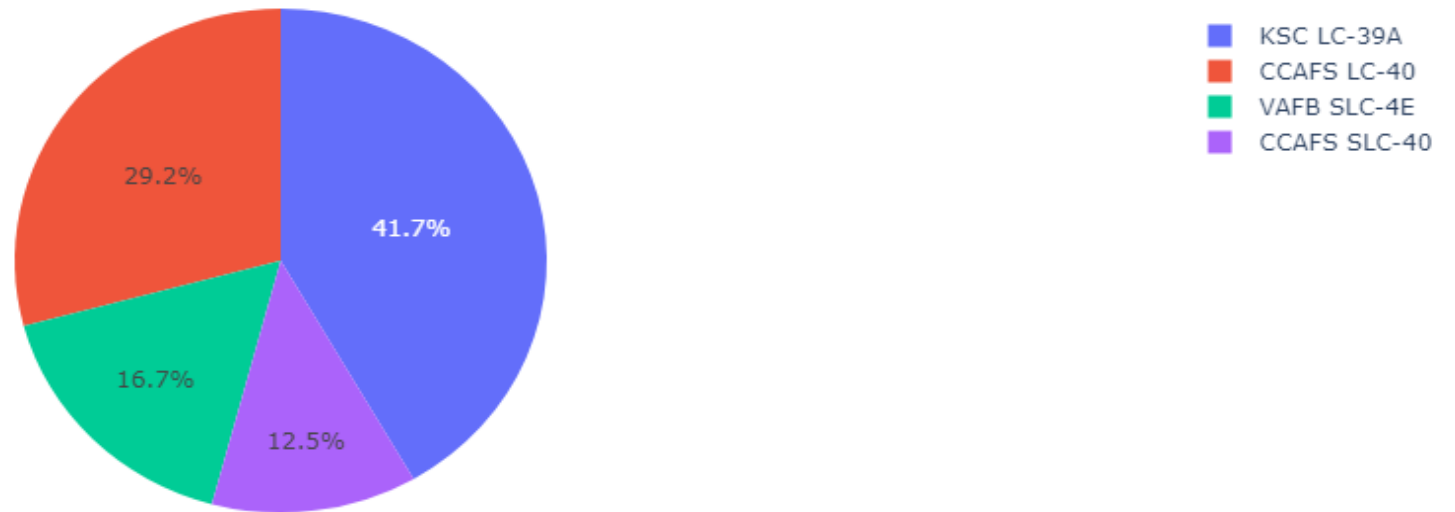
Build a Dashboard with Plotly Dash

RESULTS

Dash

- The picture below shows a pie chart when all launch site is chosen.

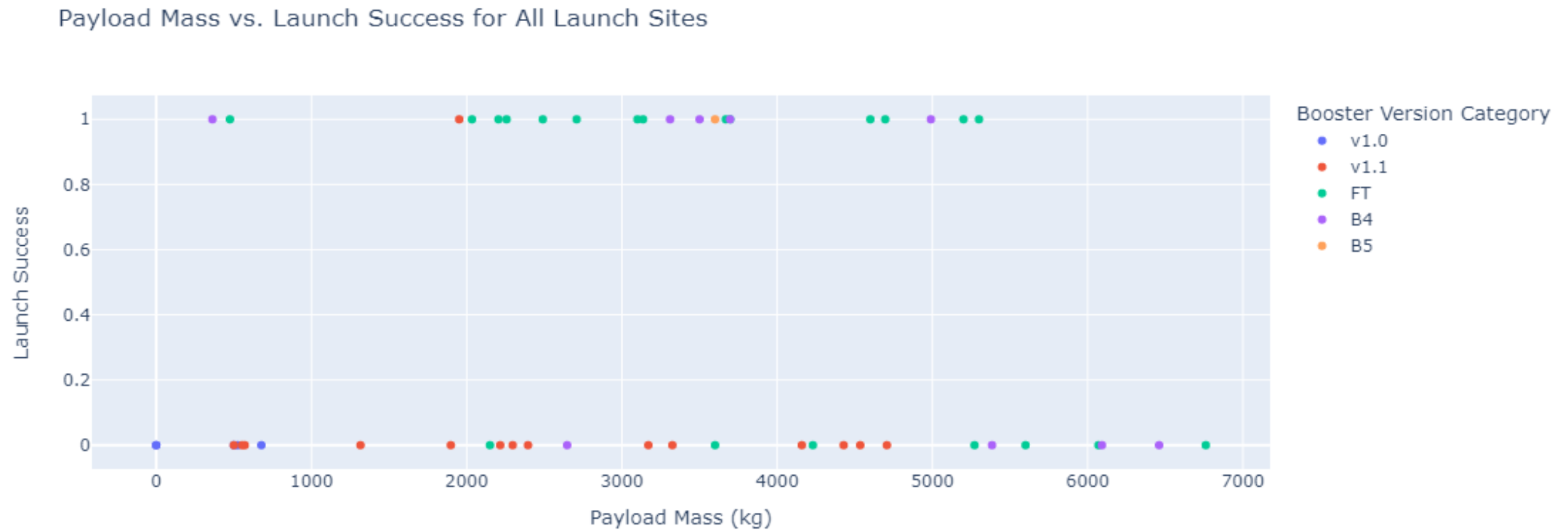
Success Rate for All Launch Sites



RESULTS

Dash

- The picture below shows a scatterplot when the payload mass range is set to be from 2000kg to 8000kg.



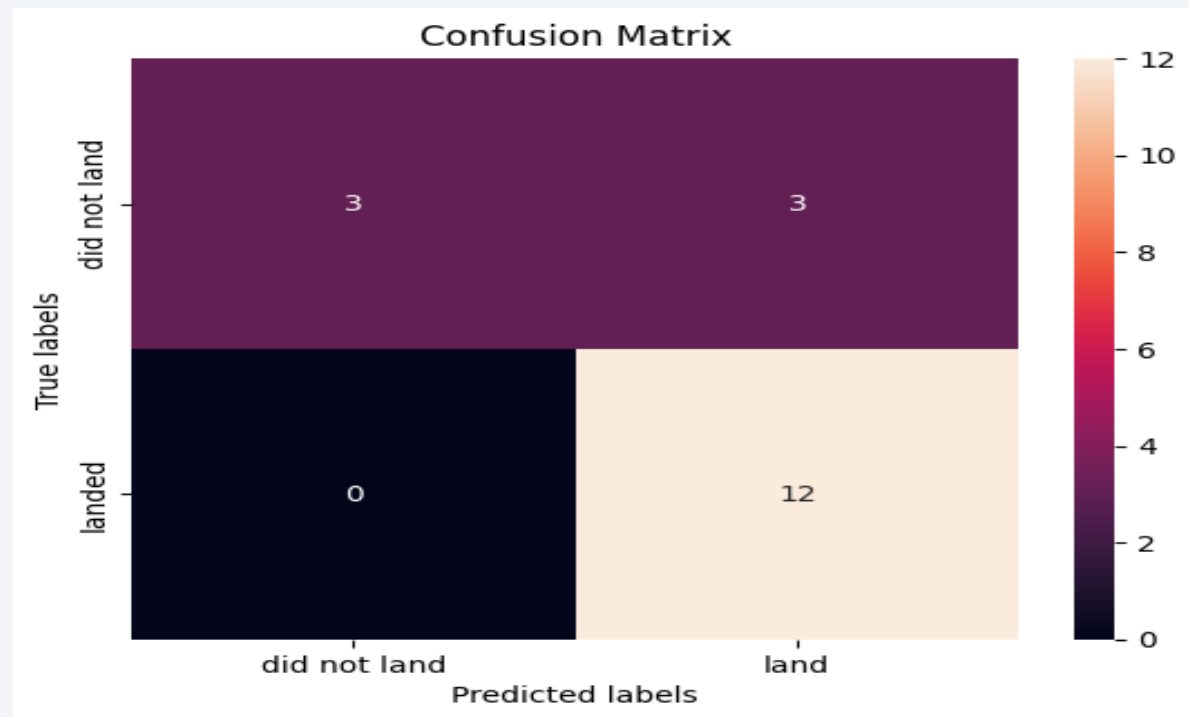
Section 5

Predictive Analysis (Classification)

RESULTS

Predictive Analysis

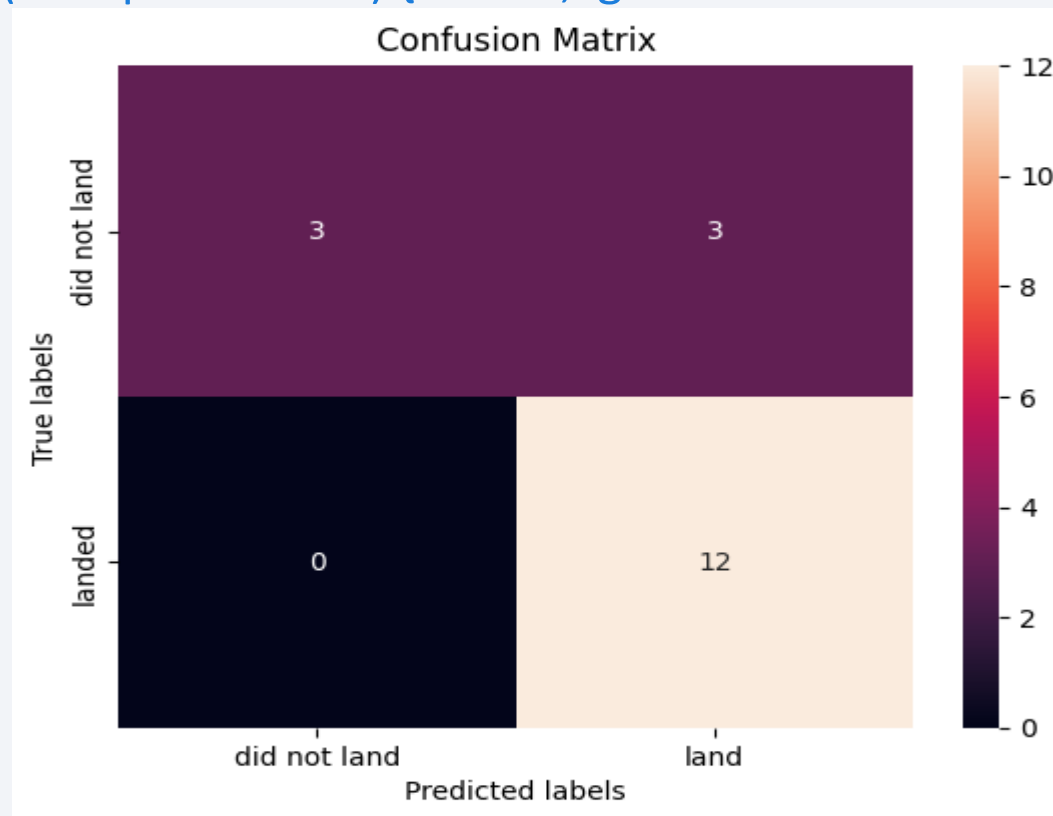
- Logistic regression
 - GridSearchCV best score: 0.8464285714285713
 - Accuracy score on test set: 0.8333333333333334
 - tuned hpyerparameters :(best parameters) {'C': 0.01, 'penalty': 'l2', 'solver': 'lbfgs'}
 - Confusion matrix:



RESULTS

Predictive Analysis

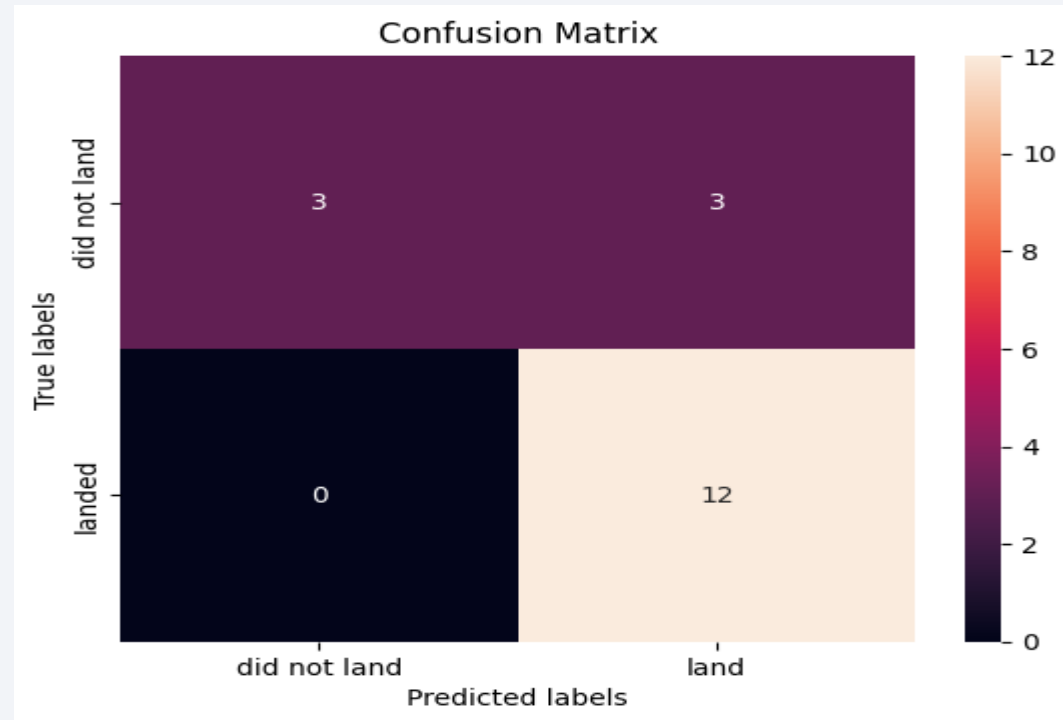
- Support vector machine (SVM)
 - GridSearchCV best score: 0.8482142857142856
 - Accuracy score on test set: 0.8333333333333334
 - tuned hpyerparameters :(best parameters) {'C': 1.0, 'gamma': 0.03162277660168379, 'kernel': 'sigmoid'}
 - Confusion matrix:



RESULTS

Predictive Analysis

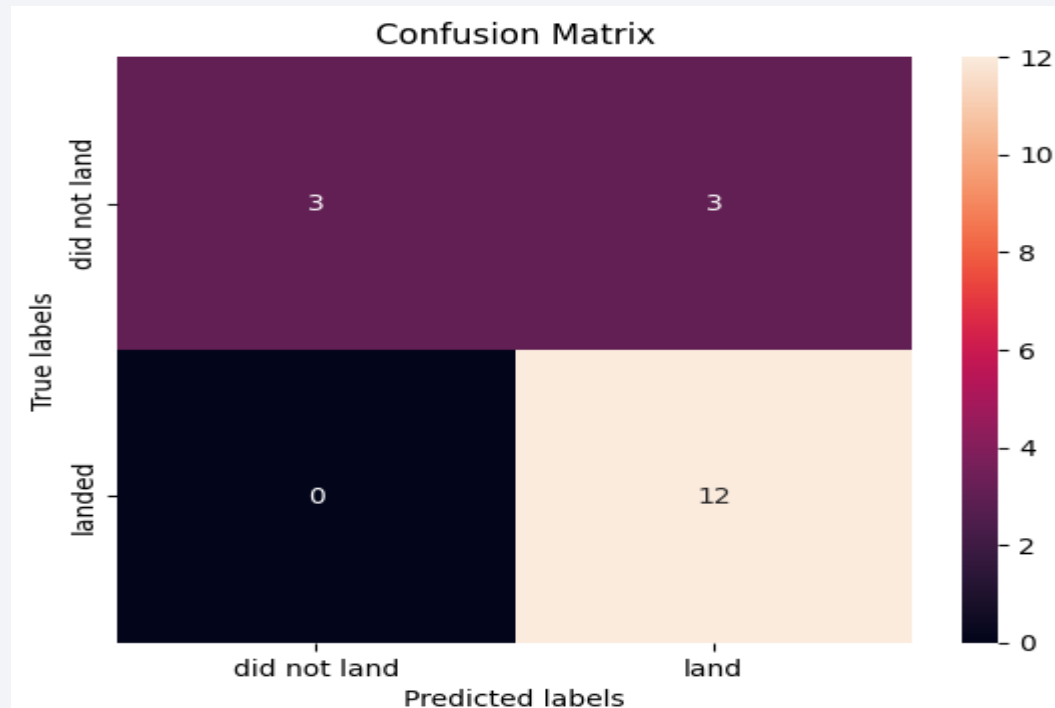
- K nearest neighbors (KNN)
 - GridSearchCV best score: 0.8482142857142858
 - Accuracy score on test set: 0.8333333333333334
 - tuned hpyerparameters :(best parameters) {'algorithm': 'auto', 'n_neighbors': 10, 'p': 1}
 - Confusion matrix:



RESULTS

Predictive Analysis

- DecisionTree
 - GridSearchCV best score: 0.8642857142857142
 - Accuracy score on test set: 0.8333333333333334
 - tuned hpyerparameters :(best parameters) {'criterion': 'gini', 'max_depth': 18, 'max_features': 'sqrt', 'min_samples_leaf': 4, 'min_samples_split': 10, 'splitter': 'random'}
 - Confusion matrix:



RESULTS

Predictive Analysis

- Putting the results of all 4 models side by side, we can see that they all share the same accuracy score and confusion matrix when tested on the test set.
- Therefore, their GridSearchCV best scores are used to rank them instead. Based on the GridSearchCV best scores, the models are ranked in the following order with the first being the best and the last one being the worst:
 1. Decision tree (GridSearchCV best score: 0.8642857142857142)
 2. K nearest neighbors, KNN (GridSearchCV best score: 0.8482142857142858)
 3. Support vector machine, SVM (GridSearchCV best score: 0.8482142857142856)
 4. Logistic regression (GridSearchCV best score: 0.8464285714285713)

DISCUSSION

- From the data visualization section, we can see that some features may have correlation with the mission outcome in several ways. For example, with heavy payloads the successful landing or positive landing rate are more for orbit types Polar, LEO and ISS. However, for GTO, we cannot distinguish this well as both positive landing rate and negative landing(unsuccesful mission) are both there here.
- Therefore, each feature may have a certain impact on the final mission outcome. The exact ways of how each of these features impact the mission outcome are difficult to decipher. However, we can use some machine learning algorithms to learn the pattern of the past data and predict whether a mission will be successful or not based on the given features.

Conclusions

- In this project, we try to predict if the first stage of a given Falcon 9 launch will land in order to determine the cost of a launch.
- Each feature of a Falcon 9 launch, such as its payload mass or orbit type, may affect the mission outcome in a certain way.
- Several machine learning algorithms are employed to learn the patterns of past Falcon 9 launch data to produce predictive models that can be used to predict the outcome of a Falcon 9 launch.
- The predictive model produced by decision tree algorithm performed the best among the 4 machine learning algorithms employed.

Thank you!

