# Data Science and Machine Learning Capstone Project

Stephen Malcolm

05/ 06/ 2023

# Outline

- Executive Summary

- Introduction

- Methodology

- Results
  - EDA with Visualization
  - EDA with SQL
  - Interactive Maps with Folium
  - Interactive Maps with Plotly
  - Machine Learning Prediction

- Conclusion

- Appendix

# Executive Summary

- Summary of methodologies
    - Data collection via API & Web Scraping
    - Exploratory Data Analysis (EDA) with Data Visualization
    - EDA with SQL
    - Building an Interactive Map with Folium
    - Building a Dashboard with Plotly Dash
    - Predictive Analysis (Classification)
- Summary of all results
    - EDA results
    - Interactive maps and dashboard
    - Predictive results

3

# Introduction

### Background

SpaceX is the first private company to develop a liquid-propellent rocket that has reached orbit; to launch, orbit, and recover a spacecraft; to send a spacecraft to the International Space Station; and to send astronauts to orbit and to the International Space Station. It is also the first organization of any type to achieve a vertical propulsive landing of an orbital rocket booster and the first to reuse such a booster. SpaceX's Falcon 9 rockets have landed and re-flown more than 150 times [1].

### Objective

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.

### Questions

- How do variables such as payload mass, launch site, number of flights, and orbits affect the success of the first stage landing?

- Does the rate of successful landings increase over time?

- What is the best model that can be used for predicting a successful landing?

# Methodology

# Methodology - Steps

Data Collection:

- Using SpaceX Rest API
- Using web scraping from Wikipedia

Data Wrangling:

- Filter data, handle missing values and apply one hot encoding (to prepare data for analysis and modelling)

EDA:

- Perform EDA with visualization and SQL
- Perform interactive visual analytics using Folium and Plotly Dash

Build Model:

- To predict landing outcomes using classification models. Tune and evaluate for model optimization

# Data Collection – API

## Steps

- **Request data** from SpaceX API (rocket launch data)

- **Decode response** using .json() and convert to a dataframe using .jason_normalize()

- **Request information** about the launches from SpaceX API using custom functions

- **Create dictionary** from data

- **Create dataframe** from the dictionary

- **Filter dataframe** to contain only Falcon 9 launches

- **Replace missing values** of Payload Mass with calculated.mean()

- **Export data** to csv file

# Data Collection – Web Scraping

Steps

- **Request data** from Wikipedia (Falcon 9 launch data)

- **Create BeautifulSoup (webscraping library) object** from HTML response

- **Extract column names** from HTML table header

- **Collect data** from parsing HTML tables

- **Create dictionary** from the data

- **Create dataframe** from the dictionary

- **Export data** to csv file

8

# Data Wrangling – Part 1

- Calculate number of launches for each site

- Calculate the number occurrence of each object

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

- Create landing outcome column (target variable/ 0 or 1)

- Export data to csv file

# Data Wrangling – Part 2

In the dataset, there are varied cases where the booster has successful and unsuccessful landings.

Examples of Mission Outcomes:

- True Ocean, True RTLS, True ASDS means the mission/ landing was successful (i.e. landed in specific region of Ocean, landed to a ground pad and landed to a drone ship respectively )
- False Ocean, False RTLS, False ASDS means the mission/ landing was unsuccessful (i.e. landed in specific region of Ocean, landed to a ground pad and landed to a drone ship respectively )

Convert these mission outcomes into Training Labels:

- String variables require to be transformed into categorical variables, where 1 means successful landing of booster and 0 means an unsuccessful landing

# EDA with Data Visualization

- Scatter Plots - normally used to observe and visually display the relationship/ correlation between variables
  - Flight Number v Payload Mass
  - Flight Number v Launch Site
  - Payload Mass v Launch Site
  - Orbit Type v Flight Number
  - Payload Mass v Orbit Type
  - Orbit Type v Payload Mass

- Bar Charts - useful to compare different categorical or discrete variables. Easy to understand relations between variables
  - Success Rate v Orbit Type

- Line Chart - A line chart supports monitoring behaviour in a set of data. These charts are useful for tracking change over time
  - Success Rate v Year

# EDA with SQL

Performed SQL queries to retrieve data from dataset, in IBM cloud database

Displaying:

- Names of the unique launch sites in the space mission

- 5 records where launch sites begin with string, 'CCA'

- Total Payload Mass carried by boosters launched by NASA (CRS)

- Average Payload Mass carried by the booster version F9 v1.1

Listing:

- Date when the first successful landing outcome in the drone ship was achieved

- Names of the boosters which have success in the ground pad and have payload mass greater than 4000, but less than 6000kg

- Total number of of successful and failure mission outcomes

- Names of the booster versions which have carried the maximum payload mass (using subquery)

- Successful landing outcomes in ground pad, their booster versions and launch site names for the month in Year 2017

- Count of successful  landing outcomes between 04-06-2010 and 20-03-2017 in descending order

12

# Build an Interactive Map with Folium

Markers to Indicate Launch Sites

- Blue circle and pop up/ text label created, to indicate NASA Johnson Space Center's coordinates as a start location

- More markers with same features created for all Launch Sites, using their coordinates to show geographical locations and proximity to equator and coast line

Markers of Launch Outcomes

- Markers of successful (green) and failed (red) launches implemented, using Marker Cluster to identify which launch sites have high success rates

Distances Between a Launch Site to Proximities

- Colored lines were added to to show distances between launch site, KSC LC-39A (our example) and it's proximity to nearest coastline, railway, highway and city

# Build a Dashboard with Plotly Dash

Dropdown list with Launch Sites

- Allows user to select launch site or all launch sites

Pie Chart Showing Successful Launches (All Sites/ Specific Sites)

- Shows total successful launches count for all sites and the Success v Failed counts for the site, using dropdown component

Slider of Payload Mass Range

- Allows user to select payload mass range

Scatter Chart Showing Payload Mass v Success Rate by Booster Version

- Allows user to view the relationship/ correlation between Payload Mass and Launch Success

14

# Predictive Analytics (Classification)

Data Pre-Processing

- Import and load dataset

- Create NumPy array from Class column

- Standardize data

- Split data into training and test sets

Model Preparation and Evaluation

- Create a GridSearchCV object, for parameter optimization

- Apply GridSearch CV on various ML algorithms, to find best parameters for each model
    - Logistic Regression
    - Support Vector Machine
    - Decision Tree
    - K-Nearest Neighbor

- Calculate the accuracy on the test data using the .score method

- Assess the confusion matrix for all models

- Find the model/ method that performs best. Take other evaluation methods into consideration

# Results

# Results - Summary

- EDA results

- Interactive analytics demo in screenshots
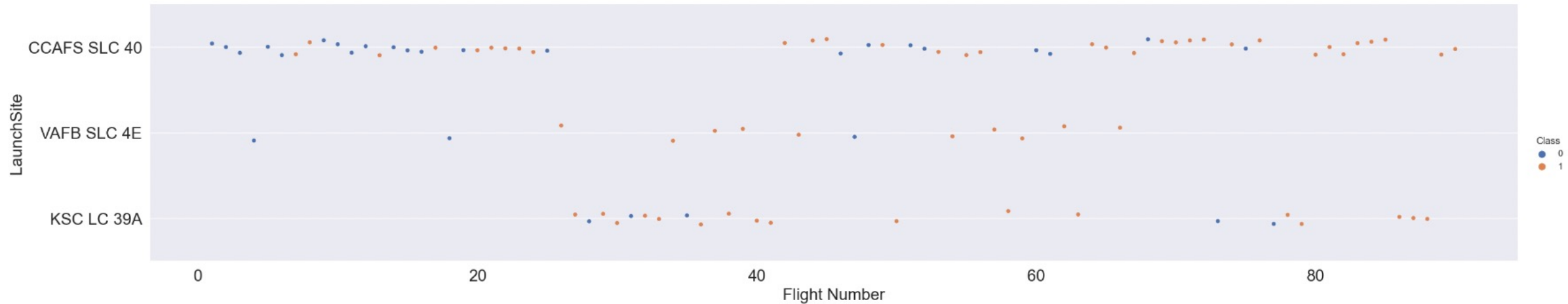
- Predictive analysis results
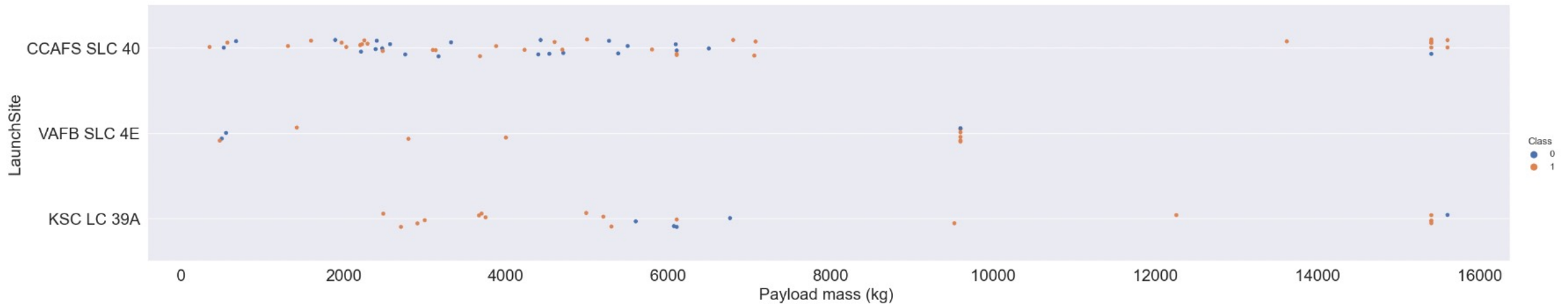
# EDA with Visualization

# Flight Number v Launch Site



Inferences (Note; blue=failed launch and orange=successful launch):
- The earliest flights all failed, whilst the latest flights all succeeded
- Around half of all launches were from the CCAFS SLC 40 launch site
- VAFB SLC 4E and KSC LC 39A have higher success rates
- We finally infer that for each site, the success rate increases over time (later flight numbers)
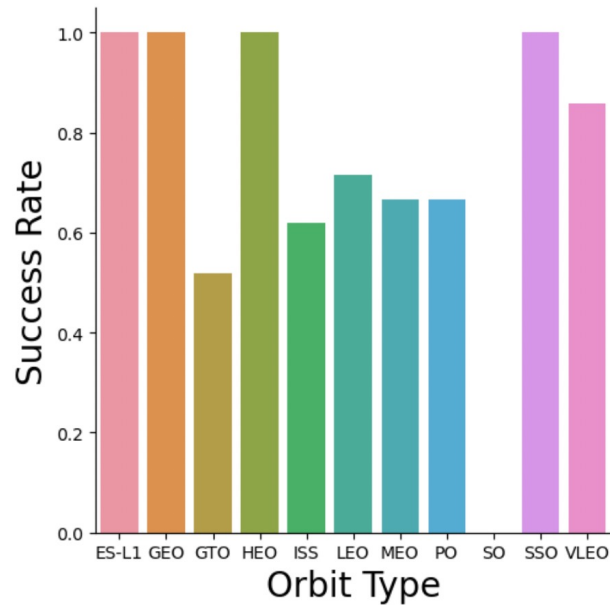
# Payload v Launch Site



Inferences (Note; blue=failed launch and orange=successful launch):
- The majority of launches with payload mass over 7000kg were successful
- KSC LC 39A has a 100% success rate for launches less than 5,500kg
- VAFB SKC 4E has not launched anything greater than 10,000kg
- Our analysis has led us to a significant inference regarding the impact of payload mass on successful landings, depending on the launch site. It appears that a heavier payload mass may be a crucial factor to consider for a successful landing. However, it's important to note that there is a potential risk of failure if the payload mass exceeds a certain threshold. When the payload mass becomes excessively heavy, it can potentially jeopardize the landing outcome.
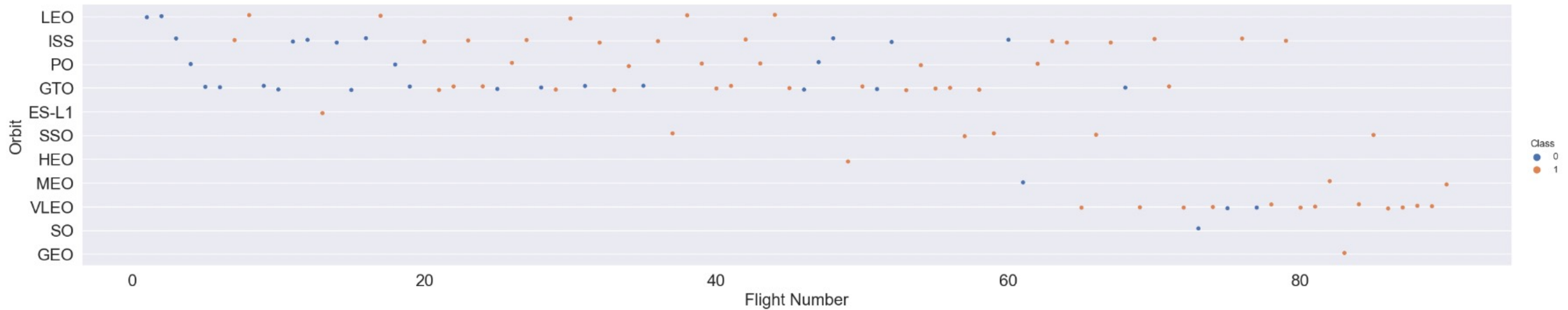
# Success Rate v Orbit Type



Inferences:
- Orbit types with 100% success rate; ES-L1, GEO, HEO and SSO
- Orbit types with success rate between 50% and 85%; GTO, ISS, LEO, MEO and PO
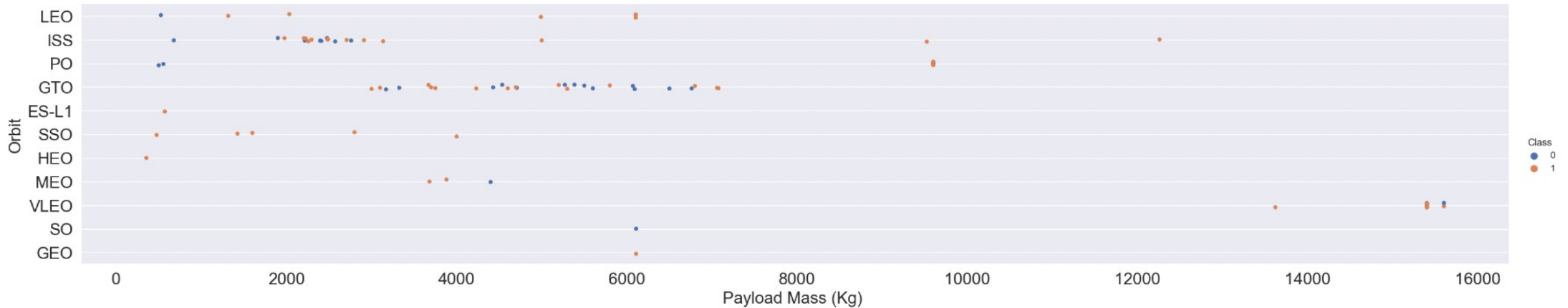- Orbit types with 0% success rate; SO

# Flights Number v Orbit Type



Inferences (Note; blue=failed launch and orange=successful launch):
- We have observed a notable correlation between the success rate and the number of flights for satellites in the Low Earth Orbit (LEO). This suggests that an increase in the number of flights in LEO orbit positively influences the success rate.
- However, when examining satellites in the Geostationary Transfer Orbit (GTO), we did not find any significant relationship between the flight number and the success rate. It seems that the number of Flights does not play a significant role in determining the success rate for satellites in GTO orbit.
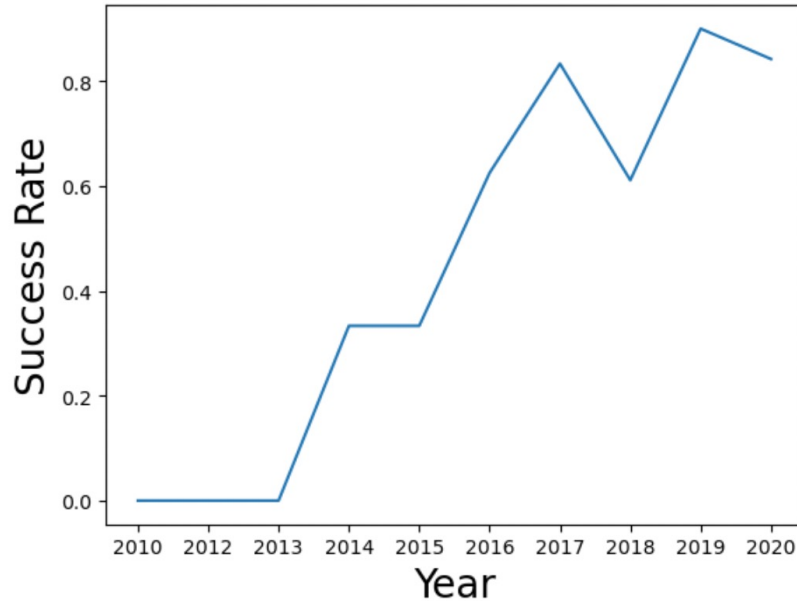
# Payload Mass v Orbit Type



Inferences (Note; blue=failed launch and orange=successful launch):

- Our analysis indicates that the weight of payloads plays a significant role in determining the outcome of launches. Specifically, in LEO, ISS and PO we have observed a positive correlation between heavier payloads and improved success rates. This suggests that increasing the weight of payloads in LEO, ISS and PO orbit tends to enhance the likelihood of a successful launch.
- Conversely, when considering the Geostationary Transfer Orbit (GTO), we have uncovered a mixed trend. GTO orbit seems to depict no correlation between payload and success rate.
- Finally, SO, GEO and HEO orbit requires more data points to make inferences/ observe trends.

# Launch Success Yearly Trend



Inferences:
- The success rate improved from 2013 to 2017 and from 2018 to 2019
- However, the success rate deteriorated from 2017 to 2018 and from 2019 to 2020
- Overall, the SpaceX rocket launch success has improved, since 2013.

# EDA with SQL

# All Launch Site Names

SQL Query and Results:

```
In [30]:    1  %sql SELECT DISTINCT LAUNCH_SITE as "Launch_Sites" FROM SPACEXTBL;

             * sqlite:///my_data1.db
            Done.
```

Out[30]:

| Launch_Sites |
| --- |
| CCAFS LC-40 |
| VAFB SLC-4E |
| KSC LC-39A |
| CCAFS SLC-40 |
| None |

Overview:
Displays the names of unique launch sites, in the space mission

# Launch Site Begins With 'KSC'

SQL Query and Results:

```
In [31]:   1  %sql SELECT * \
           2      FROM SPACEXTBL \
           3      WHERE LAUNCH_SITE LIKE'KSC%' LIMIT 5;

 * sqlite:///my_data1.db
Done.
```

Out[31]:

| Date | Time (UTC) | Booster_Version | Launch_Site | Payload | PAYLOAD_MASS__KG_ | Orbit | Customer | Mission_Outcome | Landing_Outcome |
|------|-----------|-----------------|-------------|---------|-------------------|-------|----------|-----------------|-----------------|
| 19/02/2017 | 14:39:00 | F9 FT B1031.1 | KSC LC-39A | SpaceX CRS-10 | 2490.0 | LEO (ISS) | NASA (CRS) | Success | Success (ground pad) |
| 16/03/2017 | 6:00:00 | F9 FT B1030 | KSC LC-39A | EchoStar 23 | 5600.0 | GTO | EchoStar | Success | No attempt |
| 30/03/2017 | 22:27:00 | F9 FT B1021.2 | KSC LC-39A | SES-10 | 5300.0 | GTO | SES | Success | Success (drone ship) |
| 05/01/2017 | 11:15:00 | F9 FT B1032.1 | KSC LC-39A | NROL-76 | 5300.0 | LEO | NRO | Success | Success (ground pad) |
| 15/05/2017 | 23:21:00 | F9 FT B1034 | KSC LC-39A | Inmarsat-5 F4 | 6070.0 | GTO | Inmarsat | Success | No attempt |

Overview:
Displays 5 records where launch sites begin with the string, 'KSC' (Kennedy Space Center)

# Total Payload Mass

SQL Query and Results:

```
In [32]:   1  %sql SELECT SUM(PAYLOAD_MASS__KG_) \
           2      FROM SPACEXTBL \
           3      WHERE CUSTOMER = 'NASA (CRS)';

           * sqlite:///my_data1.db
          Done.

Out[32]:
          SUM(PAYLOAD_MASS__KG_)

                        45596.0
```

Overview:
Displays the total payload mass carried by boosters, launched by NASA (CRS/ Commercial Resupply Services)

# Average Payload Mass by F9 v1.1

SQL Query and Results:

```
In [33]:    1  %sql SELECT AVG(PAYLOAD_MASS__KG_) \
            2     FROM SPACEXTBL \
            3     WHERE BOOSTER_VERSION = 'F9 v1.1';

          * sqlite:///my_data1.db
         Done.

Out[33]:
```

| AVG(PAYLOAD_MASS__KG_) |
| --- |
| 2928.4 |

Overview:
Displays average payload mass carried by booster version F9 v1.1

# First Successful Drone Ship Date

SQL Query and Results:

```
In [12]:   1  %sql SELECT MIN(DATE) AS "First Succesful Landing Outcome in Drone Ship" FROM SPACEXTBL \
           2  WHERE LANDING_OUTCOME = 'Success (drone ship)';

            * sqlite:///my_data1.db
           Done.
```

Out[12]:

| First Succesful Landing Outcome in Drone Ship |
| :--- |
| 04/08/2016 |

Overview:
Lists the date when the first successful landing outcome for a drone ship was achieved

# Successful Ground Pad Landing With Payload Between 4000 and 6000kg

SQL Query and Results:

```
In [37]:    1  %sql SELECT BOOSTER_VERSION FROM SPACEXTBL WHERE LANDING_OUTCOME = 'Success (ground pad)' \
            2  AND PAYLOAD_MASS__KG_ > 4000 AND PAYLOAD_MASS__KG_ < 6000;

  * sqlite:///my_data1.db
Done.

Out[37]:
```

| Booster_Version |
|---|
| F9 FT B1032.1 |
| F9 B4 B1040.1 |
| F9 B4 B1043.1 |

Overview:
Lists the names of all boosters that had success in ground pad landings, with payload between 4000 and 6000kg.

# Total Number of Successful and Failed Mission Outcomes

SQL Query and Results:

```
In [25]:   1  %sql SELECT sum(case when MISSION_OUTCOME LIKE '%Success%' then 1 else 0 end) AS "Successful Mission", \
           2      sum(case when MISSION_OUTCOME LIKE '%Failure%' then 1 else 0 end) AS "Failed Mission" \
           3  FROM SPACEXTBL;
```

```
 * sqlite:///my_data1.db
Done.
```

Out[25]:

| Successful Mission | Failed Mission |
| --- | --- |
| 100 | 1 |

Overview:
Lists the total number of successful and failed mission outcomes

# Boosters That Carried Maximum Payload Mass

SQL Query and Results:

```
In [17]:    1  %sql SELECT DISTINCT BOOSTER_VERSION AS "Booster Versions which carried the Maximum Payload Mass" FROM SPACEXTBL \
            2  WHERE PAYLOAD_MASS__KG_ =(SELECT MAX(PAYLOAD_MASS__KG_) FROM SPACEXTBL);
            3
```

```
 * sqlite:///my_data1.db
Done.
```

Out[17]:

| Booster Versions which carried the 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 |

Overview:
Lists the names of the booster versions that carried the maximum payload mass

# 2017 Launch Records

SQL Query and Results:

```
In [18]:  1  %sql SELECT substr(Date,4,2) as month, DATE,BOOSTER_VERSION, LAUNCH_SITE, [Landing_Outcome] \
          2  FROM SPACEXTBL \
          3  where [Landing_Outcome] = 'Success (ground pad)' and substr(Date,7,4)='2017';
```

 * sqlite:///my_data1.db
Done.

Out[18]:

| month | Date | Booster_Version | Launch_Site | Landing_Outcome |
|---|---|---|---|---|
| 02 | 19/02/2017 | F9 FT B1031.1 | KSC LC-39A | Success (ground pad) |
| 01 | 05/01/2017 | F9 FT B1032.1 | KSC LC-39A | Success (ground pad) |
| 03 | 06/03/2017 | F9 FT B1035.1 | KSC LC-39A | Success (ground pad) |
| 08 | 14/08/2017 | F9 B4 B1039.1 | KSC LC-39A | Success (ground pad) |
| 07 | 09/07/2017 | F9 B4 B1040.1 | KSC LC-39A | Success (ground pad) |
| 12 | 15/12/2017 | F9 FT B1035.2 | CCAFS SLC-40 | Success (ground pad) |

Overview:
Lists the successful landing outcomes for ground pad, their booster versions and launch site names for all months in year, 2017.

# Rank the Count of Successful Landing Outcomes Between 04/06/10 to 20/03/17

SQL Query and Results:

```
In [19]:  1 %sql SELECT [Landing_Outcome], count(*) as count_outcomes \
          2 FROM SPACEXTBL \
          3 WHERE DATE between '04-06-2010' and '20-03-2017' group by [Landing_Outcome] order by count_outcomes DESC;
          4
```

```
 * sqlite:///my_data1.db
Done.
```

Out[19]:

| Landing_Outcome | count_outcomes |
| --- | --- |
| Success | 20 |
| No attempt | 10 |
| Success (drone ship) | 8 |
| Success (ground pad) | 7 |
| Failure (drone ship) | 3 |
| Failure | 3 |
| Failure (parachute) | 2 |
| Controlled (ocean) | 2 |
| No attempt | 1 |

Overview:
Ranks the count of landing outcomes between dates, 04/06/10 to 20/03/17, in descending order

# Launch Site Location Analysis

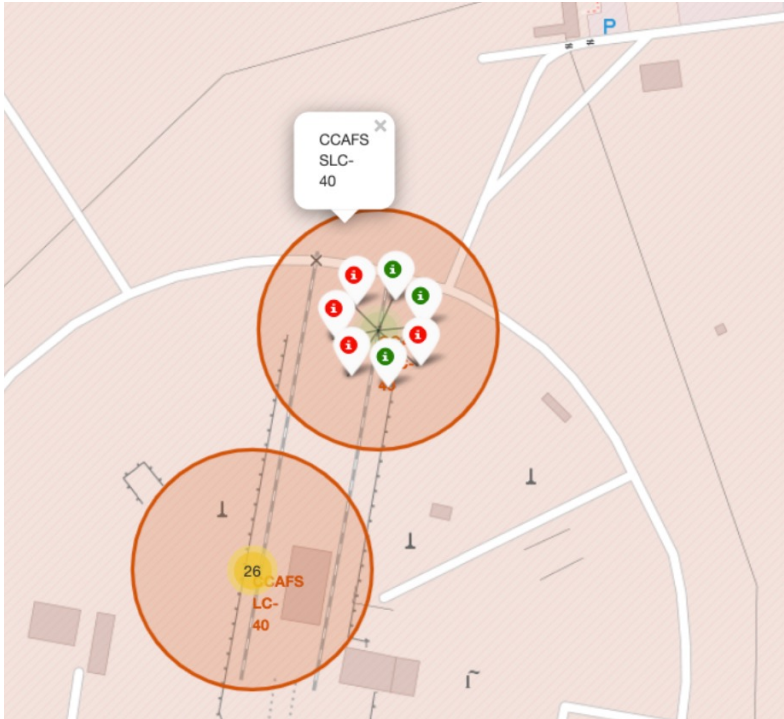Interactive Map with Folium

# Launch Sites (with Markers)



**Equator:**

Most of the launch sites are in proximity to the equator line.  The closer the launch site to the equator, the easier it is to launch to equatorial orbit. Also, if a spacecraft is launched from a site near Earth's equator, it can take optimum advantage of the Earth's substantial rotational speed. Sitting on the launch pad near the equator, it is already moving at a speed of over 1650 km per hour relative to Earth's centre [2]. This naturally helps save costs towards extra fuel and boosters.

**Coast:**

Reviewing the interactive map, all launch sites are in very close proximity to the coast. If anything goes wrong during their launch, the debris would essentially fall into the ocean, which acts as a safety measure i.e. far away from densely populated areas.
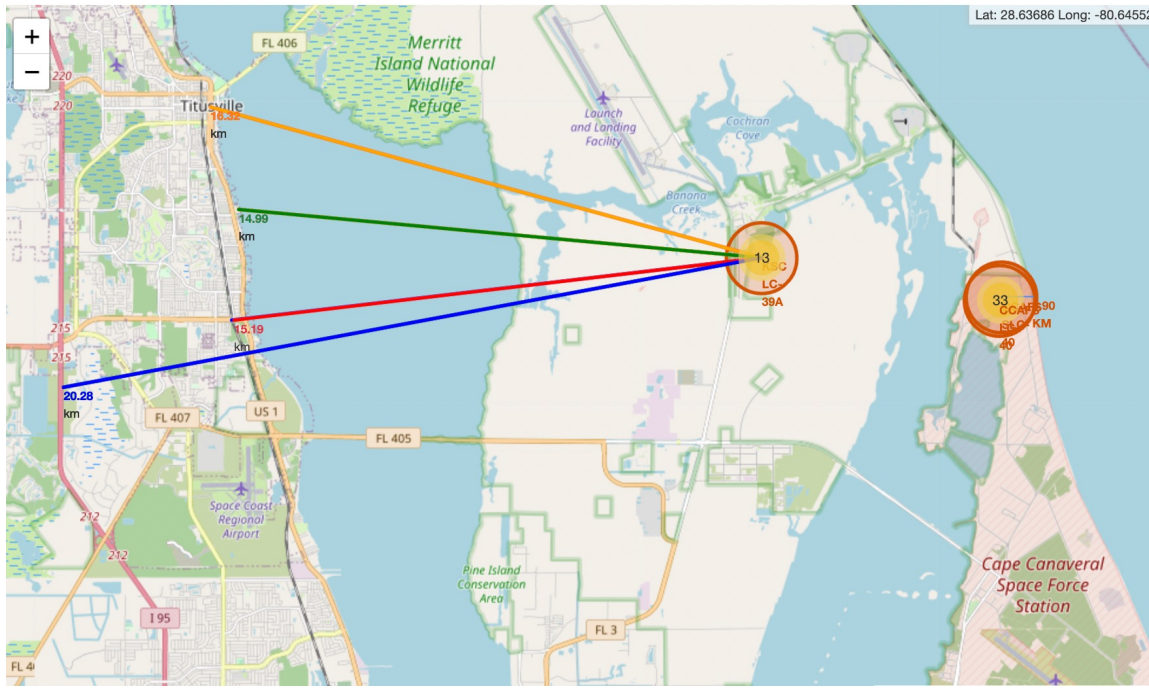
# Launch Outcomes for Each Site



From the color labeled markers above, we can identify which launch sites have a high success rate.

Green Marker = successful launches

Red Marker = unsuccessful launches

Launch site CCAFS SLC-40 has a 3/7 (43%) success rate (refer to pop up above)

# Distance from Launch Site to Proximity



If we take launch site 13 from the map above (KSC-LC-39A), we can make some inferences:

•Launch site is in close proximity of railway i.e. 15.19km
•Launch site is in close proximity to the highways i.e. 20.28km
•Launch site is in close proximity to the coastline i.e. 14.99km
•Launch site is in close proximity to Titusville (closest city on map) i.e.16.32km

A failed rocket launch could potentially be dangerous to these points of interest above. Risk assessment should then be carried out with each launch site, for precaution.

# Launch Success Count for All Sites

Total Success Launches by Site



Inferences:
From the pie chart, it can be clearly seen that KSC LC-39A has the most successful launches (41.2%), amongst other launch sites

# Launch Site with Highest Launch Success Ratio

Total Success Launches for Site KSC LC-39A



Inferences:
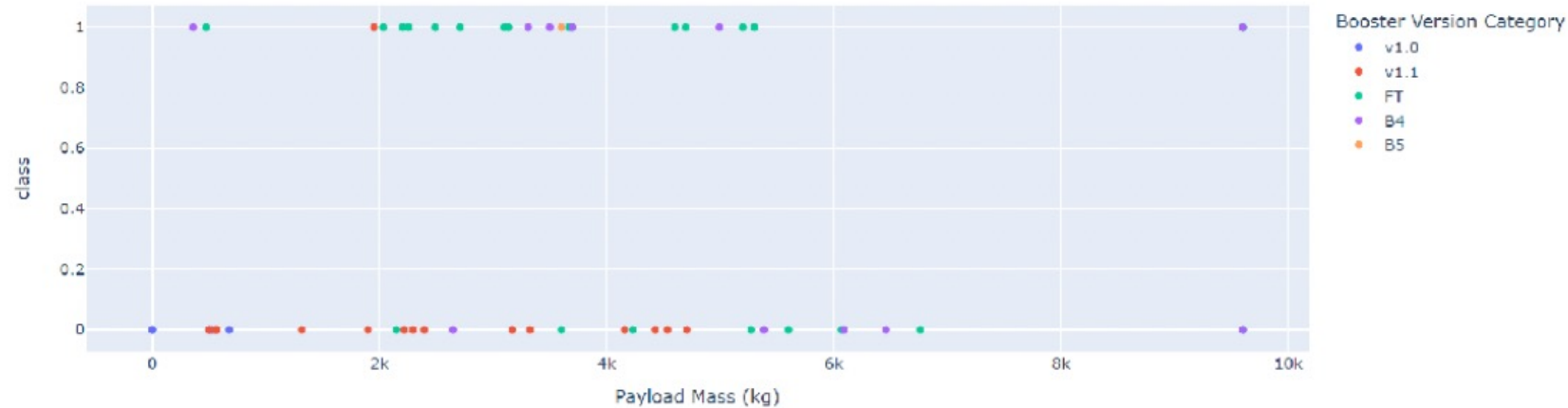KSC LC-39A also has the highest launch success rate (76.9%), amongst other launch sites.
10 successful launches and 3 failed launches

# Payload Mass v Launch Outcome for All Sites



Inferences:
The chart indicates that payload mass between 2000kg and 5500kg have the highest success rate

# Predictive Analytics (Binary Classification)

# Results

All models produced the same test accuracy, apart from the decision tree (DT) model. Indicating that they perform reasonably well in terms of predictive accuracy on the test/ unseen data. We can also consider other evaluation metrics or criteria to determine the best model i.e. F1 Score, as opposed to test accuracy:
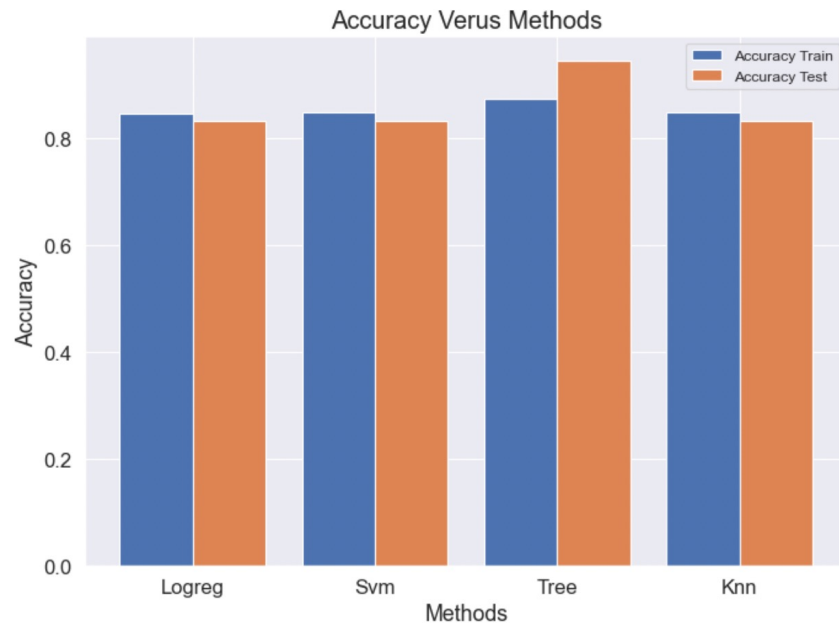
| | Accuracy Train | Accuracy Test |
|---|---|---|
| **Logreg** | 0.846429 | 0.833333 |
| **Svm** | 0.848214 | 0.833333 |
| **Tree** | 0.875000 | 0.944444 |
| **Knn** | 0.848214 | 0.833333 |

| | LogReg | SVM | Tree | KNN |
|---|---|---|---|---|
| **F1_Score** | 0.888889 | 0.888889 | 0.960000 | 0.888889 |
| **Accuracy** | 0.833333 | 0.833333 | 0.944444 | 0.833333 |

We can clearly see the best performance is a F1 Score of 0.96, from the Decision Tree. F1 score is usually more useful than accuracy, especially if we have an uneven class distribution. Accuracy works best if false positives and false negatives have similar cost. However, if the cost of false positives and false negatives are different (such as our case, refer to confusion Matrices on slide 47) it's better to look at both Precision and Recall. And since the model demonstrates strong performance on both the training data and the test/evaluation data, it indicates that our model is not overfitting. This is a positive outcome, indicating that the model generalizes well and can effectively handle unseen data.
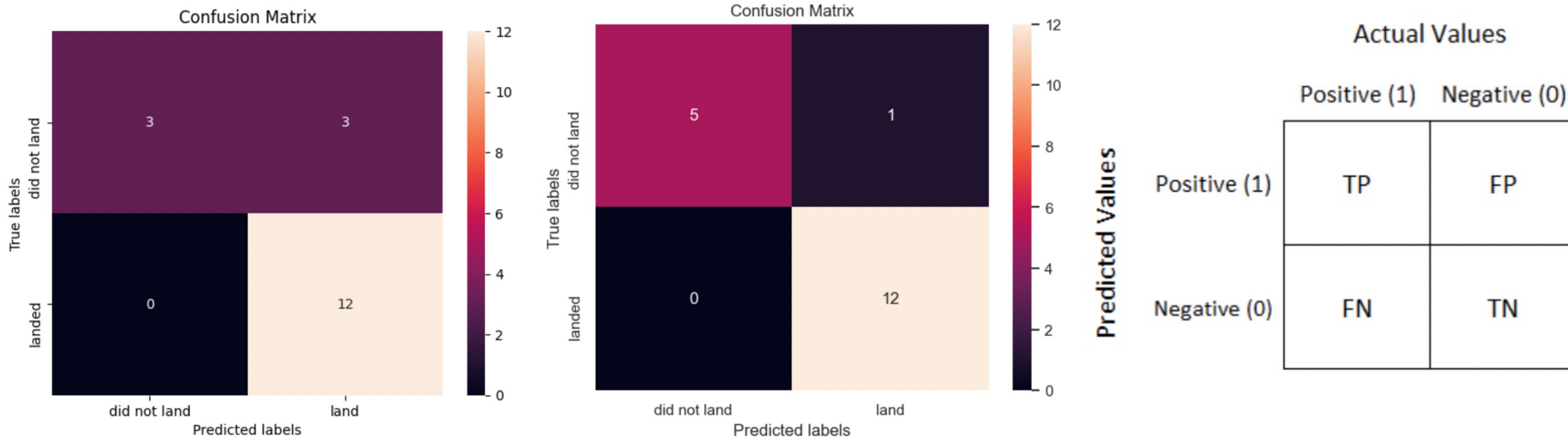
**Interpretability:** We can also consider the interpretability of the models. Some models, such as decision trees or linear models, provide more interpretable results compared to complex models like neural networks or ensemble methods. In our case, the Logistic Regression and Decision tree are considered interpretable, unlike SVM and KNN. And the higher the interpretability of a machine learning model, the easier it is for someone to comprehend why certain decisions or predictions have been made.

# Evaluation

**Data:** In general we lack test data, as we only have 18 test samples to work with. Hence, the test dataset is an unrepresentative sample of data from the domain. Possibly a reason why all the test accuracies are similar (apart from the DT model). Overall, the evaluation metrics and criteria may vary depending on our specific problem, data, and objectives of our project. Therefore, it is recommended to consider a combination of metrics and criteria to make an informed decision about the best model, taking into account our specific requirements and constraints. In our case, the DT method/ model performs the best, when we take into consideration all evaluation metrics or criteria.

# Evaluation (Cont.) – Confusion Matrix



In the context of predicting a successful landing, false positives (FP) refer to cases where the model predicts a successful landing when it actually resulted in a crash or failure. Having 3 false positives (all models, except the DT with only 1 FP) is a major problem in this scenario because it can lead to false confidence in the model's predictions, potentially putting human lives at risk.

Consider a situation where the model predicts a successful landing, but in reality, the landing is unsuccessful. If these false positives are not caught and appropriate precautions are not taken, the consequences could be catastrophic. For example, if the rocket is deemed safe to land based on the model's prediction but is actually not in a suitable condition, it could result in a crash.

# Evaluation (Cont.)

To ensure the safety and reliability of the landing process, it is crucial to minimize false positives. The DT model performed best here, where it only produced 1 FP. The model should be highly accurate in identifying actual successful landings to avoid any erroneous decisions or actions based on incorrect predictions. Therefore, false positives are a major problem because they can compromise safety and undermine the trustworthiness of the predictive model.

Finally, here are some strategies that could reduce the FP occurrence and mitigate their impact:

•   Gathering more training data: Increasing the quantity and diversity of training data can help the model learn better decision boundaries, potentially reducing false positives. Ensure that the training data includes a representative distribution of both positive and negative instances.
•   Ensemble methods: Employing ensemble methods, such as bagging or boosting, can help reduce false positives by combining multiple models' predictions. Ensemble methods often result in more robust and accurate predictions than a single model.
•   Feature engineering: Carefully selecting and engineering features can improve the model's ability to discriminate between positive and negative instances. Consider incorporating additional relevant features that capture meaningful information for the problem at hand.
•   Using different evaluation metrics: Accuracy alone might not be an appropriate metric when dealing with imbalanced datasets. Consider using evaluation metrics that are more sensitive to false positives, such as precision, recall, F1 score, or area under the precision-recall curve.

# Conclusion

- The decision tree is the best model/ algorithm for this particular dataset
- Most launch sites are located in close proximity to the equator, which offers an inherent advantage of providing an additional natural boost[2]. This geographical advantage results in cost savings by reducing the need for extra fuel and boosters.
- The success rate of launches increases over the years
- Launches with a low payload mass show better results than launches with a larger payload mass
- KSC-LC39A has the highest success rate of all launches, amongst all launch sites
- Orbits ES-L1, GEO, HEO, and SSO have 100% success rate

# Appendix

References

[1] SpaceX Wiki

https://en.wikipedia.org/wiki/SpaceX

[2] Basics of Space Flight/ Launch

https://solarsystem.nasa.gov/basics/chapter14-1/#:~:text=If%20a%20spacecraft%20is%20launched,hour%20relative%20to%20Earth's%20center.