# Outline

- Executive Summary

- Introduction

- Methodology

- Results

- Conclusion

- Appendix

# Executive Summary

- Summary of methodologies
  - ✓ Data Collection Through API
  - ✓ Data Collection with Web-Scraping
  - ✓ Data Wrangling
  - ✓ Exploratory Data Analysis (EDA) with SQL
  - ✓ Perform interactive visual analytics using Folium and Plotly Dash
  - ✓ Perform predictive analysis using classification models

- Summary of all results
  - ✓ Exploratory Data Analysis result
  - ✓ Interactive analytics in screenshots
  - ✓ Predictive Analytics result

# 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. The primary challenge addressed in this machine learning project is to predict whether the first stage of a Falcon 9 rocket will successfully land or not. By accurately predicting the landing outcome, potential clients or competing companies can make informed decisions when considering SpaceX as a launch provider or bidding against SpaceX for a rocket launch contract.

- Problems you want to find answers

  ❑ Can we accurately predict whether the first stage of a Falcon 9 rocket will successfully land or not?

  ❑ How can the predictions from the machine learning model be utilized in the decision-making process for potential clients or competing companies considering SpaceX for a rocket launch?

  ❑ What insights can be gained from the exploratory data analysis regarding the distribution of key features and potential patterns related to landing outcomes?

  ❑ How accurate, precise, and reliable is the model in predicting first stage landing outcomes?

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 and Data Standardization

- 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

- Data Collection was initiated the data collection process by utilizing a GET request to the SpaceX API.

- Following this, we decoded the response content using the .json() function and transformed it into a pandas dataframe using the .json_normalize() method.

- Subsequently, we conducted data cleaning, examined for any missing values, and addressed them as needed.

- Additionally, we engaged in web scraping from Wikipedia to retrieve Falcon 9 launch records using BeautifulSoup, in order to extract the launch records presented in an HTML table, parse the table, and convert the information into a pandas dataframe for subsequent 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.

- Link Notebook: https://github.com/kevinrioharris/IBM-Courses-Data-Science/blob/main/Applied%20Data%20Science%20Capstone/Data%20Collection%20API_FinalPresentation.ipynb

# 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.

- Link Notebook: https://github.com/kevinrioharris/IBM-Courses-Data-Science/blob/main/Applied%20Data%20Science%20Capstone/Data%20Collection%20with%20Web%20Scraping_FinalPresentation.ipynb

# Data Wrangling

- We conducted exploratory data analysis to establish training labels, computed launch counts per site, and analyzed orbit frequencies.

- Additionally, we derived landing outcome labels from the outcome column and exported the findings to a CSV file.

- Link notebook: https://github.com/kevinrioharris/IBM-Courses-Data-Science/blob/main/Applied%20Data%20Science%20Capstone/Data%20Wrangling_FinalPresentation.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.

- And As we can see that the sucess rate since 2013 kept increasing till 2020

- Link Notebook:
  https://github.com/kevinrioharris/IBM-Courses-Data-Science/blob/main/Applied%20Data%20Science%20Capstone/EDA%20with%20Data%20Visualization_FinalPresentation.ipynb

# EDA with SQL

- We seamlessly imported the SpaceX dataset into a PostgreSQL database within the Jupyter notebook.

- Employing EDA with SQL for exploratory data analysis, we extracted key insights, including:

  - 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.

- Link Notebook: https://github.com/kevinrioharris/IBM-Courses-Data-Science/blob/main/Applied%20Data%20Science%20Capstone/EDA%20with%20SQL_FinalPresentation.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.

- Link notebook: https://github.com/kevinrioharris/IBM-Courses-Data-Science/blob/main/Applied%20Data%20Science%20Capstone/Interactive%20Map%20with%20Folium_FinalP.ipynb

# 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.

- Link Code: https://github.com/kevinrioharris/IBM-Courses-Data-Science/blob/main/Applied%20Data%20Science%20Capstone/PlotlyApps.py

# 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.

- Link Notebook: https://github.com/kevinrioharris/IBM-Courses-Data-Science/blob/main/Applied%20Data%20Science%20Capstone/Machine%20Learning%20Prediction_FinalP.ipynb

# Results

- Best classification model is **DecisionTree** with a accuracy score of 0.8732142857142856
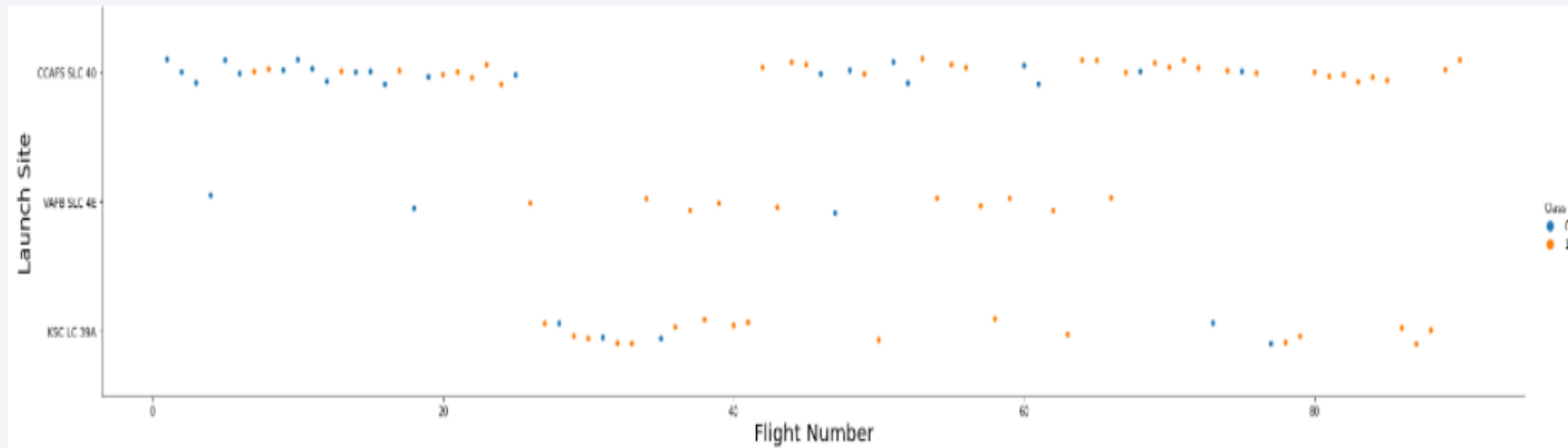
Section 2

# 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.
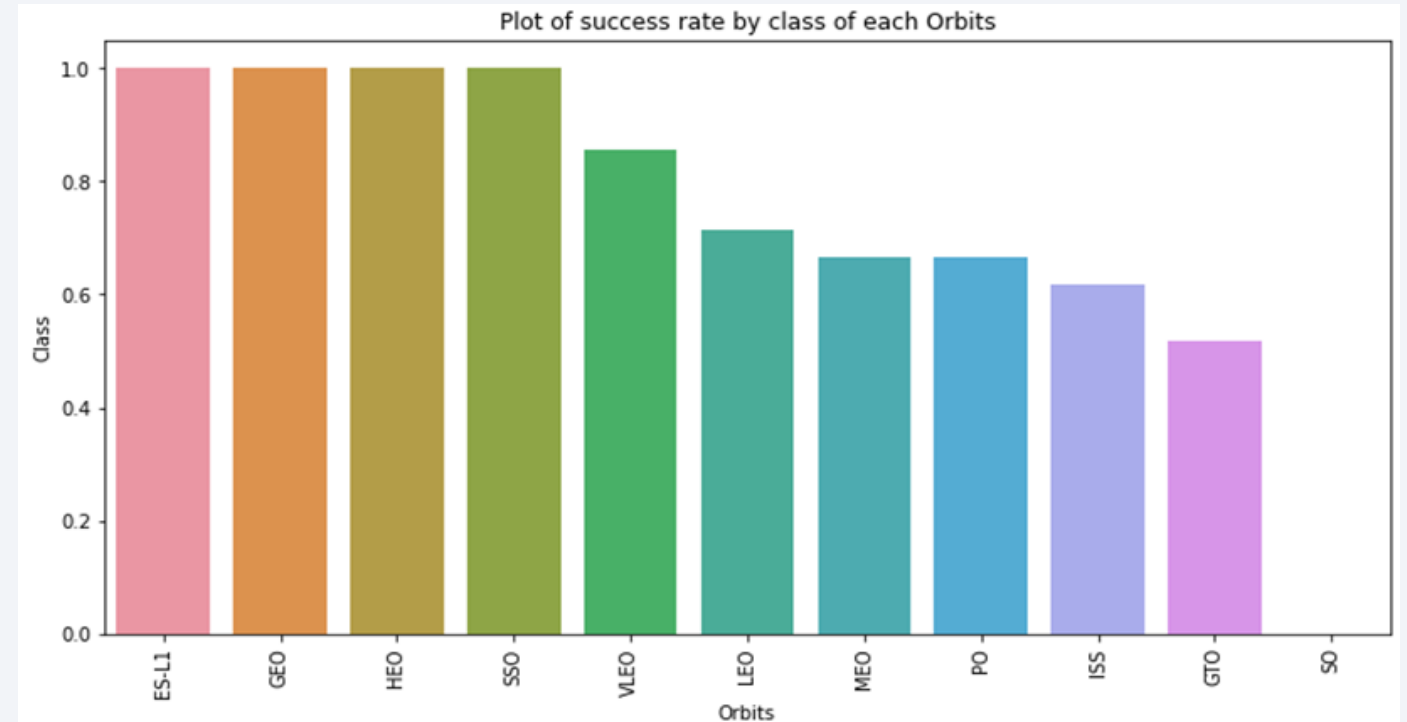
# Payload vs. Launch Site

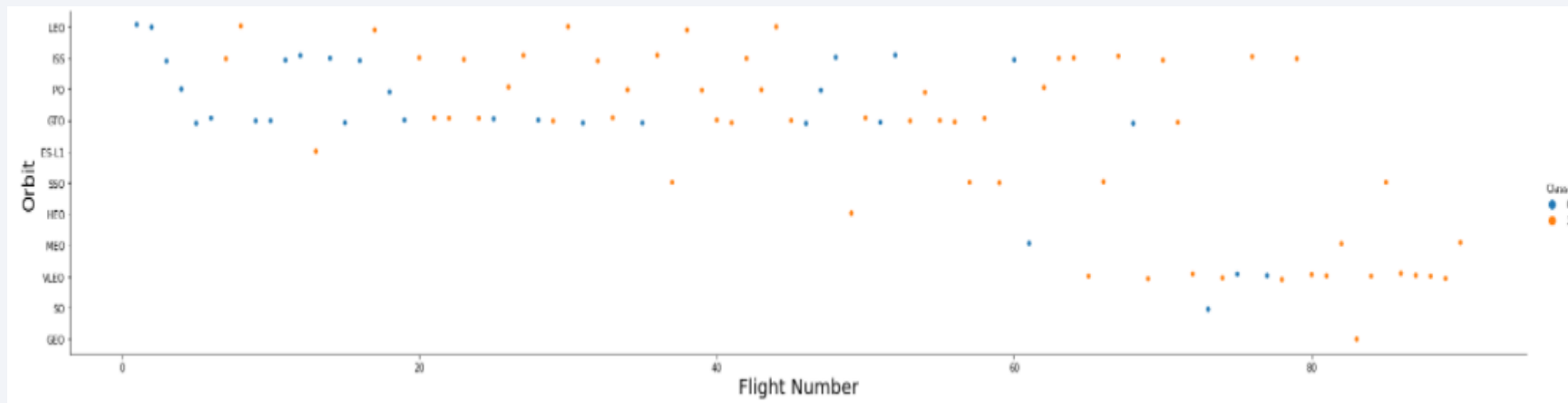The Greate the Payload mass for launch site CCAFS SLC 40, The higher the success rate for the rocket.

# Success Rate vs. Orbit Type

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



Plot of success rate by class of each Orbits

# 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.
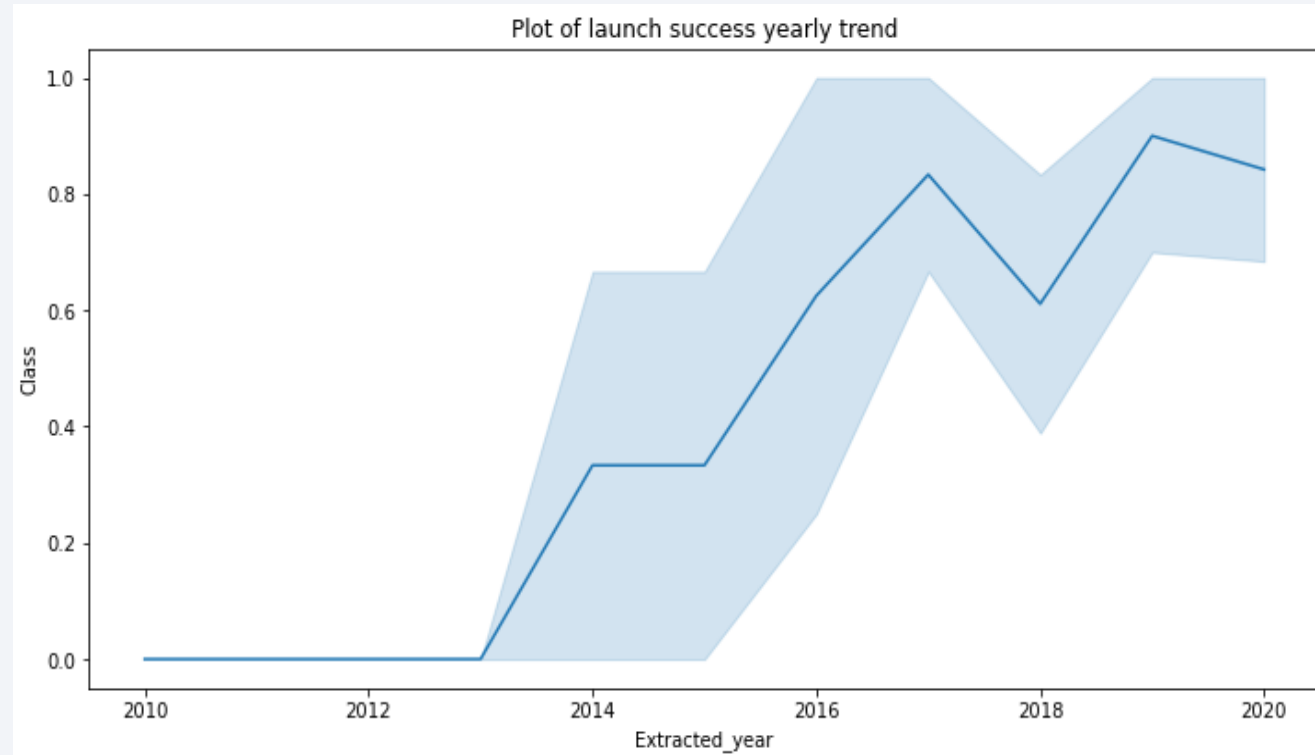
# Payload vs. Orbit Type

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.

# All Launch Site Names

We used the keyword **DISTINCT** on the LaunchSite to show only unique launch sites from the SpaceX data.

It shows

- CCAFS LC-40
- VAFB SLC-4E
- KSC LC-39A
- CCAFS SLC-40

# Launch Site Names Begin with 'CCA'



We used the query above to display 5 records where launch sites begin with `CCA` by using LIKE function and LIMIT for displaying only 5

# Total Payload Mass

- We calculated the total payload carried by boosters launched by NASA is **45596** using the query below

# Average Payload Mass by F9 v1.1

We calculated the average payload mass carried by booster version F9 v1.1 is **2928.4**

Display average payload mass carried by booster version F9 v1.1

```
[18]: %%sql
SELECT AVG(PAYLOAD_MASS__KG_) AS Avg_PayloadMass
FROM SPACEXTABLE
WHERE Booster_Version = 'F9 v1.1'
```

 * sqlite:///my_data1.db
Done.

[18]: **Avg_PayloadMass**

2928.4

# First Successful Ground Landing Date

We observed that the dates of the first successful landing outcome on ground pad was 22$^{nd}$ December 2015

Task 5

List the date when the first succesful landing outcome in ground pad was acheived.

*Hint:Use min function*

```sql
%%sql
SELECT MIN(DATE) AS First_Successfull_Landing_Date
FROM SPACEXT
WHERE Landing_Outcome LIKE 'Success (ground pad)'
```

 * sqlite:///my_data1.db
Done.

**First_Successfull_Landing_Date**

2015-12-22

# Successful Drone Ship Landing with Payload between 4000 and 6000



We used the **WHERE** clause to filter for boosters which have successfully landed on drone ship and applied the **AND** condition to determine successful landing with payload mass greater than 4000 but less than 6000

# Total Number of Successful and Failure Mission Outcomes

Task 7

List the total number of successful and failure mission outcomes

```sql
[63]: %%sql
SELECT
  COUNT(CASE WHEN Mission_Outcome LIKE 'Success%' THEN 1 END) AS Success_Outcome,
  COUNT(CASE WHEN Mission_Outcome LIKE 'Failure%' THEN 1 END) AS Failure_Outcome
FROM SPACEXTABLE;
```

 * sqlite:///my_data1.db
Done.

[63]:

| Success_Outcome | Failure_Outcome |
| --- | --- |
| 100 | 1 |

We used wildcard like '%' to filter for **WHERE** MissionOutcome was a success or a failure.

We used **CASE** function to make it shorter code for doing both Mission_Outcome : **Success** or **Failure**

# Boosters Carried Maximum Payload

We determined the booster that have carried the maximum payload using a subquery in the **WHERE** clause and the **MAX()** function.

Where The maximum payload is 15600 KG

# 2015 Launch Records



Task 9

List the records which will display the month names, failure landing_outcomes in drone ship ,booster versions, launch_site for the months in year 2015.

**Note: SQLLite does not support monthnames. So you need to use substr(Date, 6,2) as month to get the months and substr(Date,0,5)='2015' for year.**

```sql
[68]: %%sql
SELECT Booster_Version, Launch_Site, Landing_Outcome
FROM SPACEXTABLE
WHERE Landing_Outcome = 'Failure (drone ship)' AND DATE BETWEEN '2015-01-01' AND '2015-12-31'
```

 * sqlite:///my_data1.db
Done.

[68]:

| Booster_Version | Launch_Site | Landing_Outcome |
|---|---|---|
| F9 v1.1 B1012 | CCAFS LC-40 | Failure (drone ship) |
| F9 v1.1 B1015 | CCAFS LC-40 | Failure (drone ship) |

We used a combinations of the **WHERE** clause, **LIKE**, **AND**, and **BETWEEN** conditions to filter for failed landing outcomes in drone ship, their booster versions, and launch site names for year 2015

# Rank Landing Outcomes Between 2010-06-04 and 2017-03-20

- We selected Landing outcomes and the **COUNT** of landing outcomes from the data and used the **WHERE** clause to filter for landing outcomes **BETWEEN** 2010-06-04 to 2010-03-20.

- We applied the **GROUP BY** clause to group the landing outcomes and the **ORDER BY** clause to order based on date in descending order.

## Task 10

Rank the count of landing outcomes (such as Failure (drone ship) or Success (ground pad)) between

```
[70]:  %%sql
       SELECT Date, Landing_Outcome, COUNT(Landing_Outcome) AS Total_Failure
       FROM SPACEXTABLE
       WHERE Date BETWEEN '2010-06-04' AND '2017-03-20'
       GROUP BY Landing_Outcome
       ORDER BY Date DESC
```

* sqlite:///my_data1.db
Done.

[70]:

| Date | Landing_Outcome | Total_Failure |
|---|---|---|
| 2016-04-08 | Success (drone ship) | 5 |
| 2015-12-22 | Success (ground pad) | 3 |
| 2015-06-28 | Precluded (drone ship) | 1 |
| 2015-01-10 | Failure (drone ship) | 5 |
| 2014-04-18 | Controlled (ocean) | 3 |
| 2013-09-29 | Uncontrolled (ocean) | 2 |
| 2012-05-22 | No attempt | 10 |
| 2010-06-04 | Failure (parachute) | 2 |

# Launch Sites Proximities Analysis

# All Launch Sites Global Map Markers

We can see that the SpaceX Launch Site are in the United States of America coasts: Florida and California
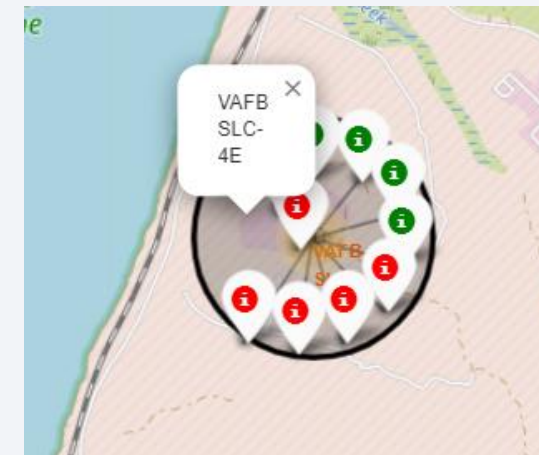
# Markers showing successful launch using color labels

Florida Launch Sites

California Launch Sites



Green Marker Shows successful Launch

Red Marker Shows Failures Launch

# Launch Site (CCAFS SLC-40) Distance to Landmarks



distance_highway = 0.5834695366934144  km
distance_railroad = 1.2845344718142522  km
distance_city = 51.434169995172326  km

Are launch sites in close proximity to railways? **No**
Are launch sites in close proximity to highways? **No**
Are launch sites in close proximity to coastline? **Yes**
Do launch sites keep certain distance away from cities? **Yes**

37

Section 4

# Build a Dashboard
# with Plotly Dash

# Pie Chart Showing The Success Percentage Achieved from Each Launch Site

We can see that KSC LC-39A had the most successful launches from all Launch Sites.



Total Success Launches By all sites

Legend:
- KSC LC-39A
- CCAFS LC-40
- VAFB SLC-4E
- CCAFS SLC-40

41.7%, 29.2%, 16.7%, 12.5%

# Pie Chart Showing The Launch Site with The Highest Launch Success Ratio



KSC LC-39A Achieved 76,9% success rate while getting 23,1% failure rate

1 means Success

0 means Failure

# Scatter Plot of Payload vs Launch Outcome for All Launch Sites (with Different Payloads)



Based on the graph above, We can see the success rates for low weighted payloads
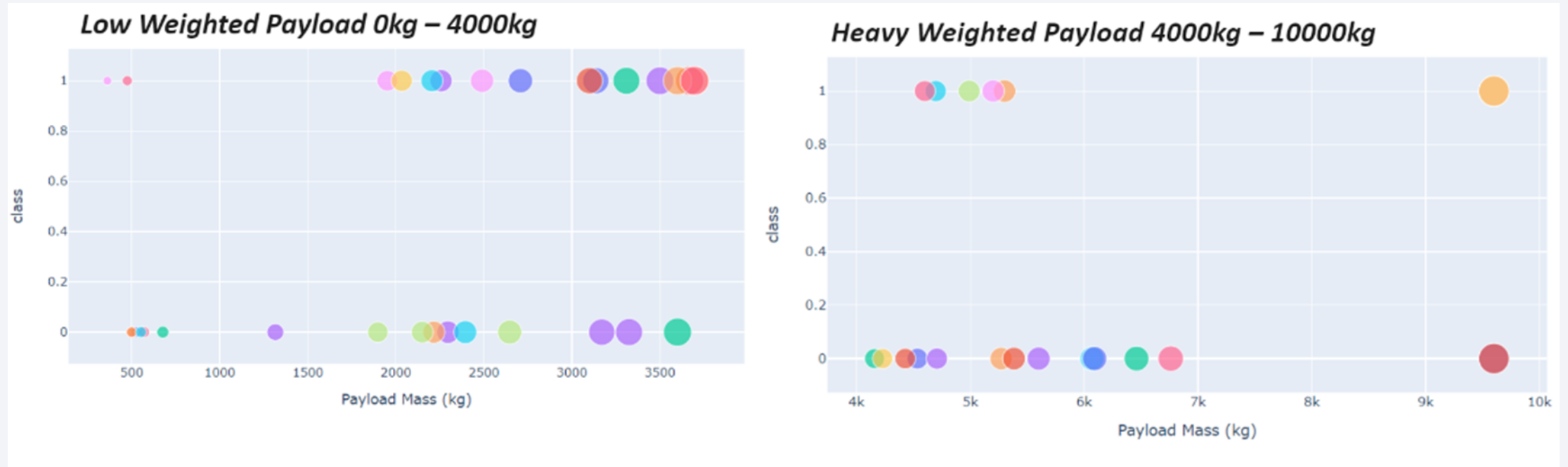is higher than the heavy weighted payloads

Section 5

# Predictive Analysis (Classification)

# Classification Accuracy

```python
models = {'KNeighbors':knn_cv.best_score_,
          'DecisionTree':tree_cv.best_score_,
          'LogisticRegression':logreg_cv.best_score_,
          'SupportVector': svm_cv.best_score_}

bestalgorithm = max(models, key=models.get)
print('Best model is', bestalgorithm,'with a score of', models[bestalgorithm])
if bestalgorithm == 'DecisionTree':
    print('Best params is :', tree_cv.best_params_)
if bestalgorithm == 'KNeighbors':
    print('Best params is :', knn_cv.best_params_)
if bestalgorithm == 'LogisticRegression':
    print('Best params is :', logreg_cv.best_params_)
if bestalgorithm == 'SupportVector':
    print('Best params is :', svm_cv.best_params_)
```
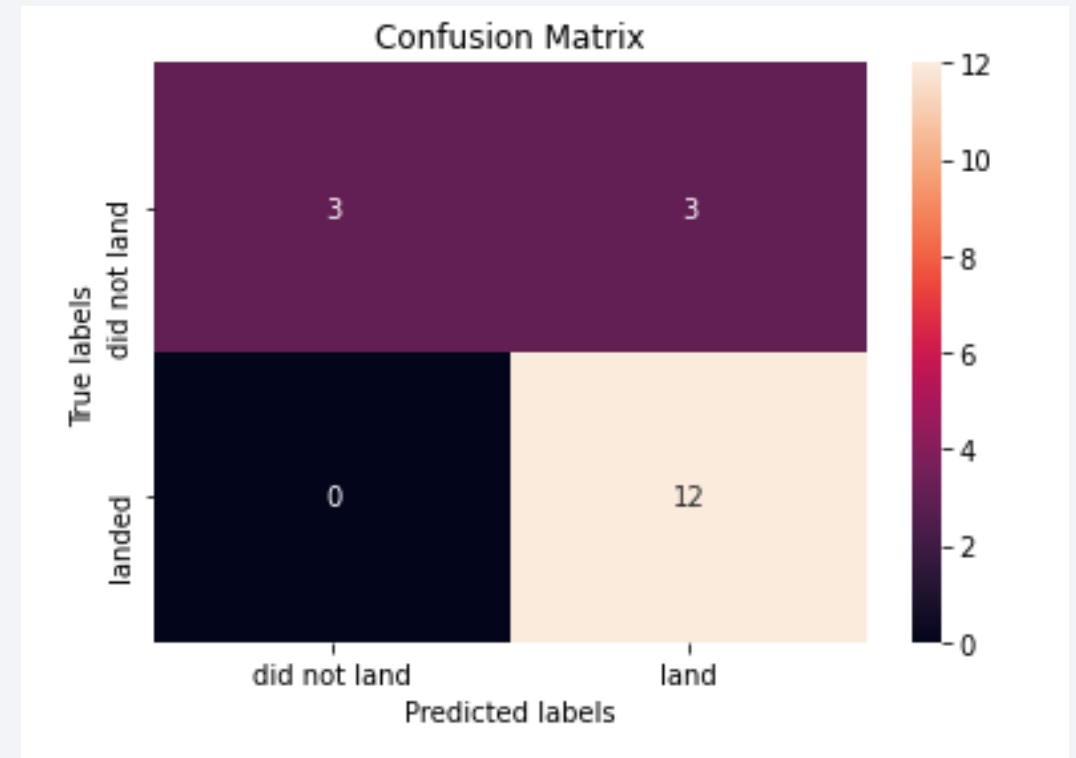
```
Best model is DecisionTree with a score of 0.8732142857142856
Best params is : {'criterion': 'gini', 'max_depth': 6, 'max_features': 'auto', 'min_samples_leaf': 2, 'min_samples_split': 5, 'splitter': 'random'}
```

Based on the code above, The **decision tree classifier** is the model with the highest classification accuracy

# 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

In summary:

- A launch site's success rate is positively correlated with the magnitude of flights conducted there.

- The launch success rate witnessed a continuous rise from 2013 to 2020.

- Orbits ES-L1, GEO, HEO, SSO, and VLEO exhibited the highest success rates.

- Among all launch sites, KSC LC-39A boasted the highest number of successful launches.

- The Decision tree classifier stands out as the most effective machine learning algorithm for this analysis.

Thank you!