# 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 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 background and context

  SpaceX has been a revolutionary organization since its inception. It has changed the space exploration industry by reducing cost of its rocket launches. According to SpaceX, it costs them $62 million to launch its Falcon 9 rocket. In comparison with other providers who have costs of $162 million each. The information that we have gathered is substantial and we can use it to ascertain whether we can compete with our own launches. We are creating a machine learning pipeline to see if the first launch will land successfully.

- Problems you want to find answers
  - What factors determine if the rocket will land successfully?
  - The interaction between various features that determine the success rate of a successful landing.
  - What operating conditions need to be in place to ensure a successful landing program.

Section 1

# **Methodology**

# Methodology

## Executive Summary

- Data collection methodology:

  - We collected the data using SpaceX API and web scraping from Wikipedia.

- Perform data wrangling

  - We use One-hot encoding 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

- Describe how data sets were collected.
  - We used get request to the SpaceX API to collect data.

  - Then we decoded the response content as a Json using .json() function  and we turned 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.

  - Our objective was to extract the launch records as a HTML table, parse the table and convert it to a pandas dataframe for analysis.

# Data Collection – SpaceX API

- With the use of the get request to the SpaceX API we collected data, cleaned the requested data and conducted basic data wrangling and formatting.

- The Link to my Notebook

https://github.com/SmatCreative/SpaceX-Final-Project/blob/main/jupyter-labs-spacex-data-collection-api.ipynb

**Task 1: Request and parse the SpaceX launch data using the GET request**

To make the requested JSON results more consistent, we will use the following static response object for this project:

```
In [9]:   static_json_url='https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DS0321EN-SkillsNetwork/datasets/API_
```

We should see that the request was successfull with the 200 status response code

```
In [10]:  response.status_code
```

Out[10]: 200

Now we decode the response content as a Json using `.json()` and turn it into a Pandas dataframe using `.json_normalize()`

```
In [11]:  # Use json_normalize meethod to convert the json result into a dataframe
          data = pd.json_normalize(response.json())
```

Using the dataframe `data` print the first 5 rows

```
In [12]:  # Get the head of the dataframe
          data.head(5)
```

Out[12]:

| | static_fire_date_utc | static_fire_date_unix | net | window | rocket | success | failures | details | crew | ships | cap |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 2006-03-17T00:00:00.000Z | 1.142554e+09 | False | 0.0 | 5e9d0d95eda69955f709d1eb | False | [{'time': 33, 'altitude': None, 'reason': | Engine failure at 33 seconds | [] | [] | |

# Data Collection - Scraping

- For web scrapping we used BeautifulSoup on the Falcon 9 launch records. Then we parsed the table converting it into a Pandas dataframe.

- GitHub URL of the completed web scraping notebook.

https://github.com/SmatCreative/SpaceX-Final-Project/blob/main/jupyter-labs-webscraping.ipynb

## TASK 1: Request the Falcon9 Launch Wiki page from its URL

First, let's perform an HTTP GET method to request the Falcon9 Launch HTML page, as an HTTP response.

```
In [5]:    # use requests.get() method with the provided static_url
           # assign the response to a object
           html_data = requests.get(static_url)
           html_data.status_code
```

```
Out[5]:    200
```

Create a `BeautifulSoup` object from the HTML `response`

```
In [6]:    # Use BeautifulSoup() to create a BeautifulSoup object from a response text content
           soup = BeautifulSoup(html_data.text)
```

Print the page title to verify if the `BeautifulSoup` object was created properly

```
In [7]:    # Use soup.title attribute
           soup.title
```

```
Out[7]:    <title>List of Falcon 9 and Falcon Heavy launches - Wikipedia</title>
```
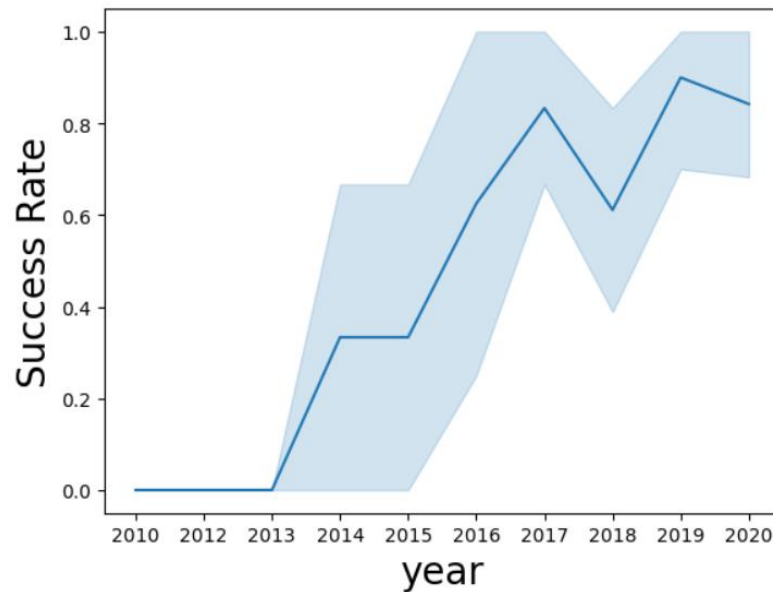
# Data Wrangling

- We conducted an exploratory data analysis and determined the training labels.

- Then we calculated the number of launches at each site, and the number and occurrence of each orbits. Thereafter we created landing outcome label from outcome column and exported the results to csv.

- GitHub URL of completed data wrangling notebooks.

https://github.com/SmatCreative/SpaceX-Final-Project/blob/main/IBM-DS0321EN-SkillsNetwork_labs_module_1_L3_labs-jupyter-spacex-data_wrangling_jupyterlite.jupyterlite.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.



you can observe that the sucess rate since 2013 kept increasing till 2020

- GitHub URL of completed EDA with data visualization notebook.

https://github.com/SmatCreative/SpaceX-Final-Project/blob/main/IBM-DS0321EN-SkillsNetwork_labs_module_2_jupyter-labs-eda-dataviz.ipynb.jupyterlite.ipynb

# EDA with SQL

We did EDA with SQL to collect insights from the data. We wrote queries to find

o Names of unique launch sites.

o Total payload mass carried by boosters launched by NASA (CRS)

o Average payload mass carried by booster version F9 v1.1

o Total number of successful and failure mission outcomes

o Failed landing outcomes in drone ship, their booster version and launch site names.

• GitHub URL of completed EDA with SQL notebook.

https://github.com/SmatCreative/SpaceX-Final-Project/blob/main/jupyter-labs-eda-sql-coursera_sqllite.ipynb

# Build an Interactive Map with Folium

- 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.

- Assigned the feature launch outcomes (failure or success) to class 0 (failure) and 1(success)

- With the use of color-labeled marker clusters, we identified launch sites that have relatively high success rate.

- Calculated the distance between launch sites to their proximities.

- GitHub URL of completed interactive map with Folium map.

https://github.com/SmatCreative/SpaceX-Final-Project/blob/main/IBM-DS0321EN-SkillsNetwork_labs_module_3_lab_jupyter_launch_site_location.jupyterlite.ipynb

# Build a Dashboard with Plotly Dash

- Plotted pie charts showing the total launches by certain sites

- Plotted scatter graph showing the relationship with Outcome and Payload Mass (Kg) for the different booster version.

- GitHub URL of completed Plotly Dash lab.

- https://github.com/SmatCreative/SpaceX-Final-Project/blob/main/IBM-DS0321EN-SkillsNetwork_labs_module_3_lab_jupyter_launch_site_location.jupyterlite.ipynb

# Predictive Analysis (Classification)

- We used numpy and pandas to load the data then transformed the data, split the data into training and testing.

- We built different machine learning models and tuned different hyperparameters using GridSearchCV.

- We used accuracy as the metric for our model.

- Improved the model using feature engineering and algorithm tuning.

- GitHub URL of completed predictive analysis lab.

- https://github.com/SmatCreative/SpaceX-Final-Project/blob/main/IBM-DS0321EN-SkillsNetwork_labs_module_4_SpaceX_Machine_Learning_Prediction_Part_5.jupyterlite%20(1).ipynb

# Results

```
In [62]:  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 0.9035714285714287
Best params is : {'criterion': 'entropy', 'max_depth': 2, 'max_features': 'sqrt', 'min_samples_leaf': 1, 'min_samples_spli
t': 5, 'splitter': 'random'}
```
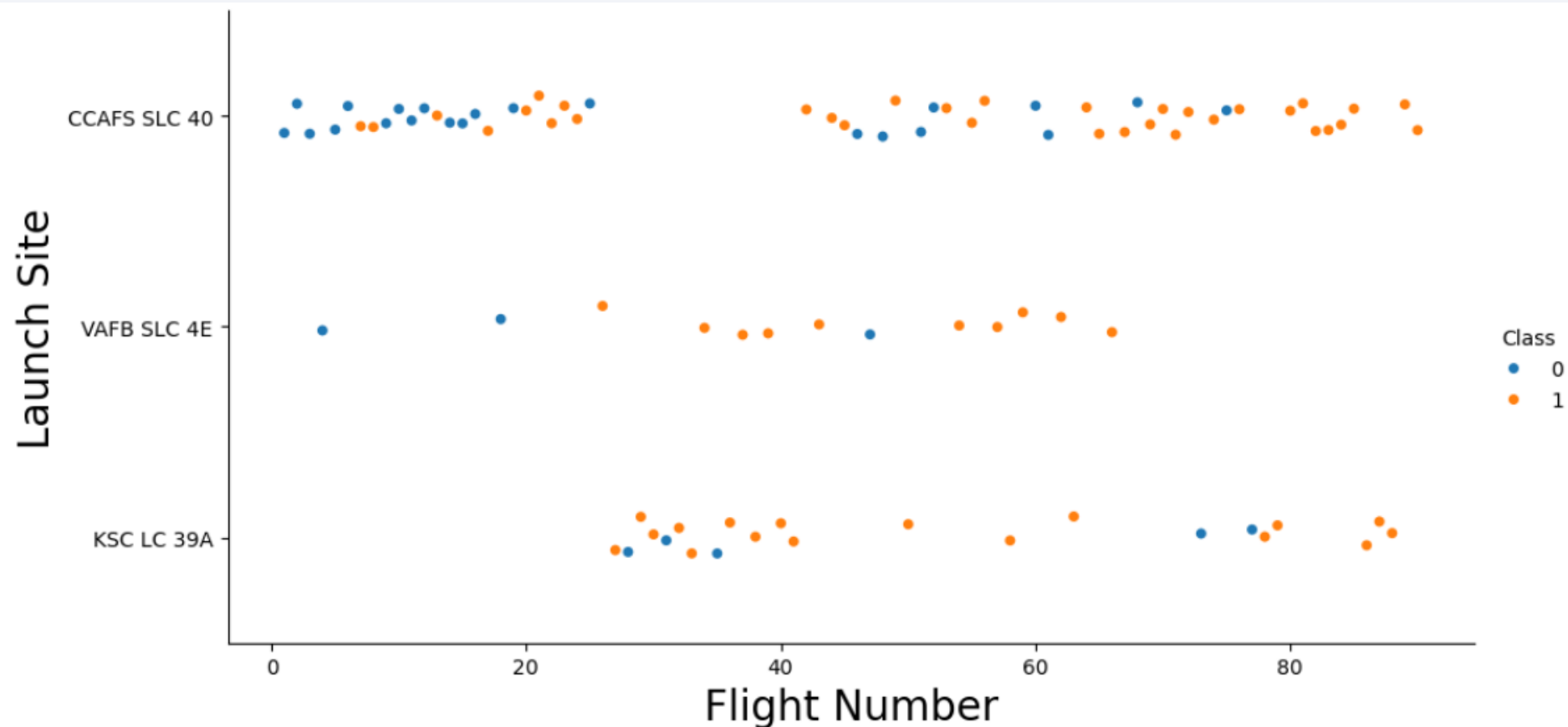
16

Section 2

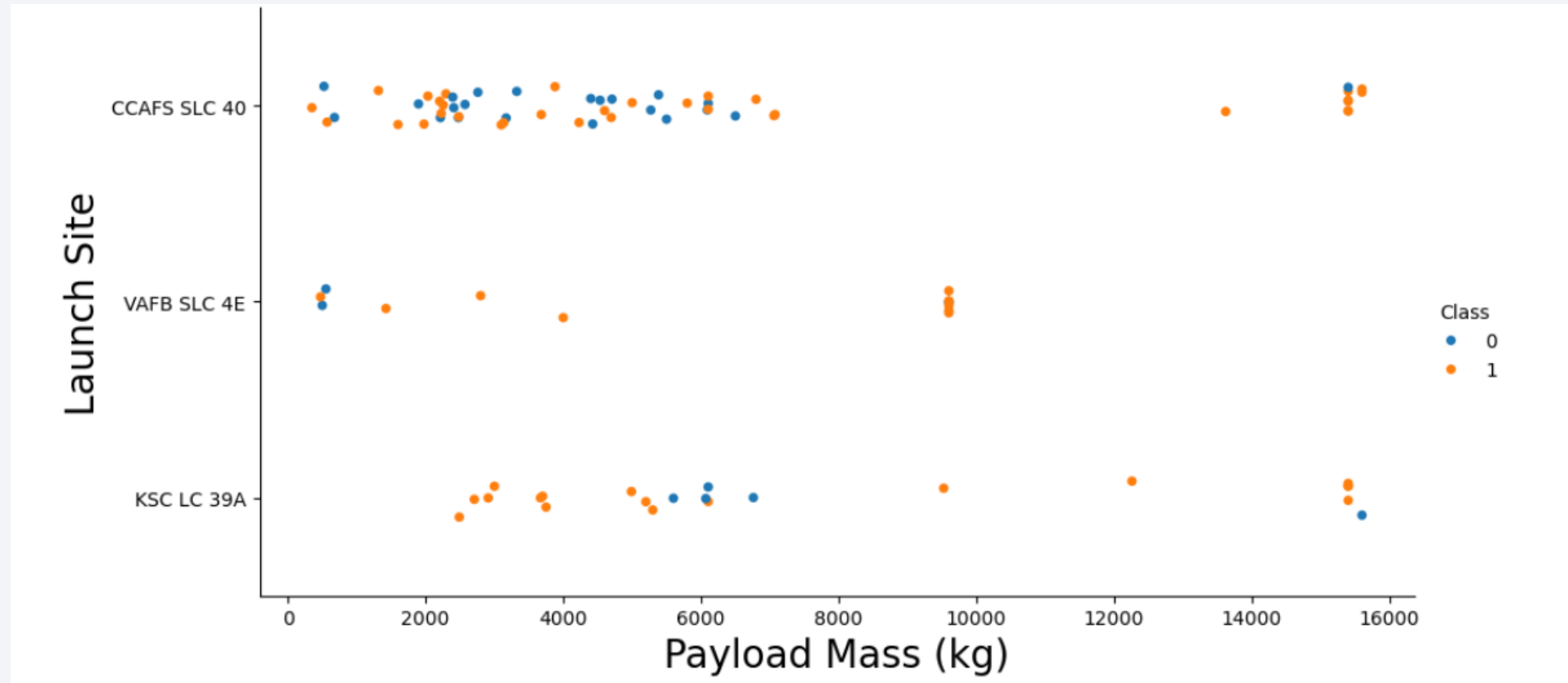# Insights drawn from EDA

# Flight Number vs. Launch Site

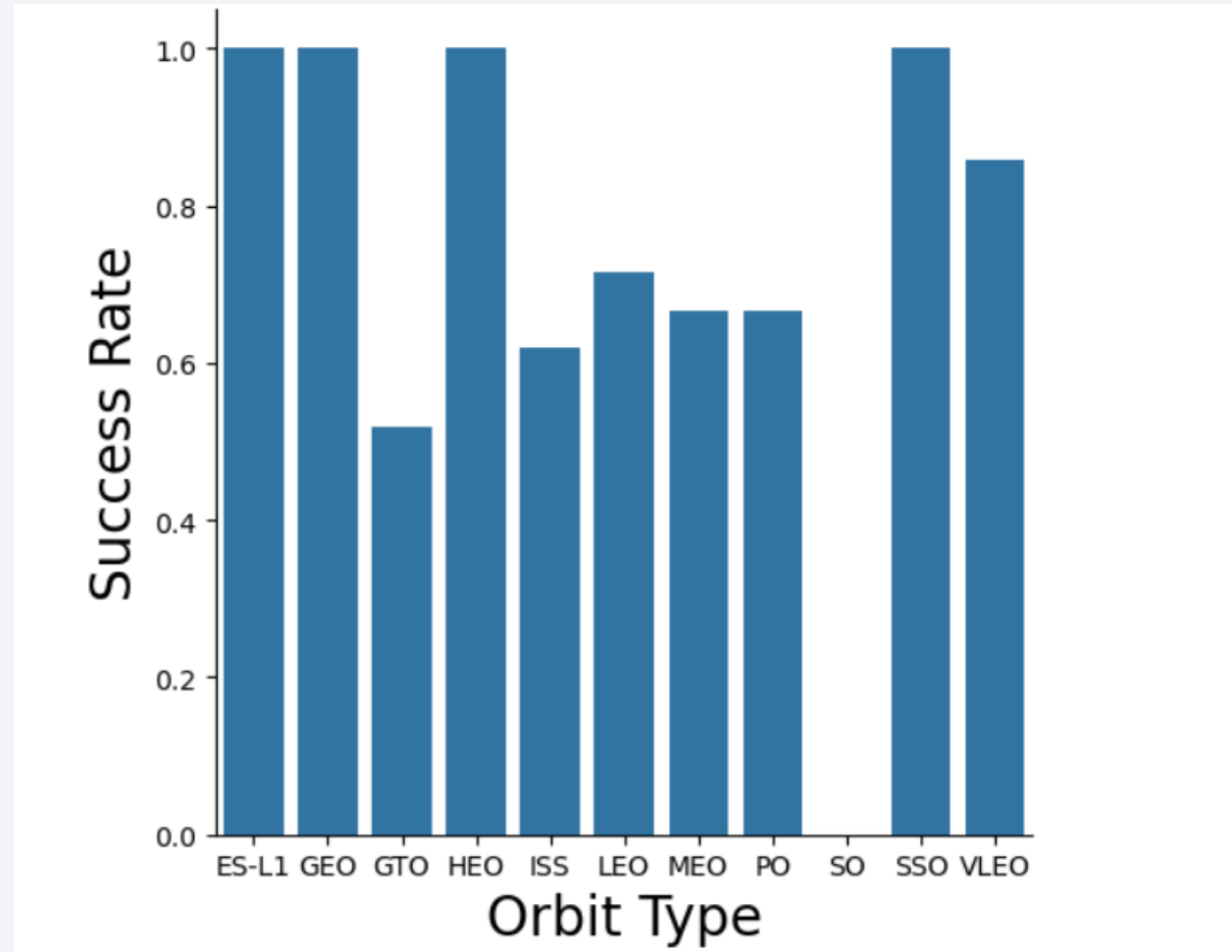- The greater the flight number at a launch site, the greater the success rate.

# Payload vs. Launch Site

- The VAFB-SLC launch site there are no rockets launched for heavy payload mass(greater than 10000).
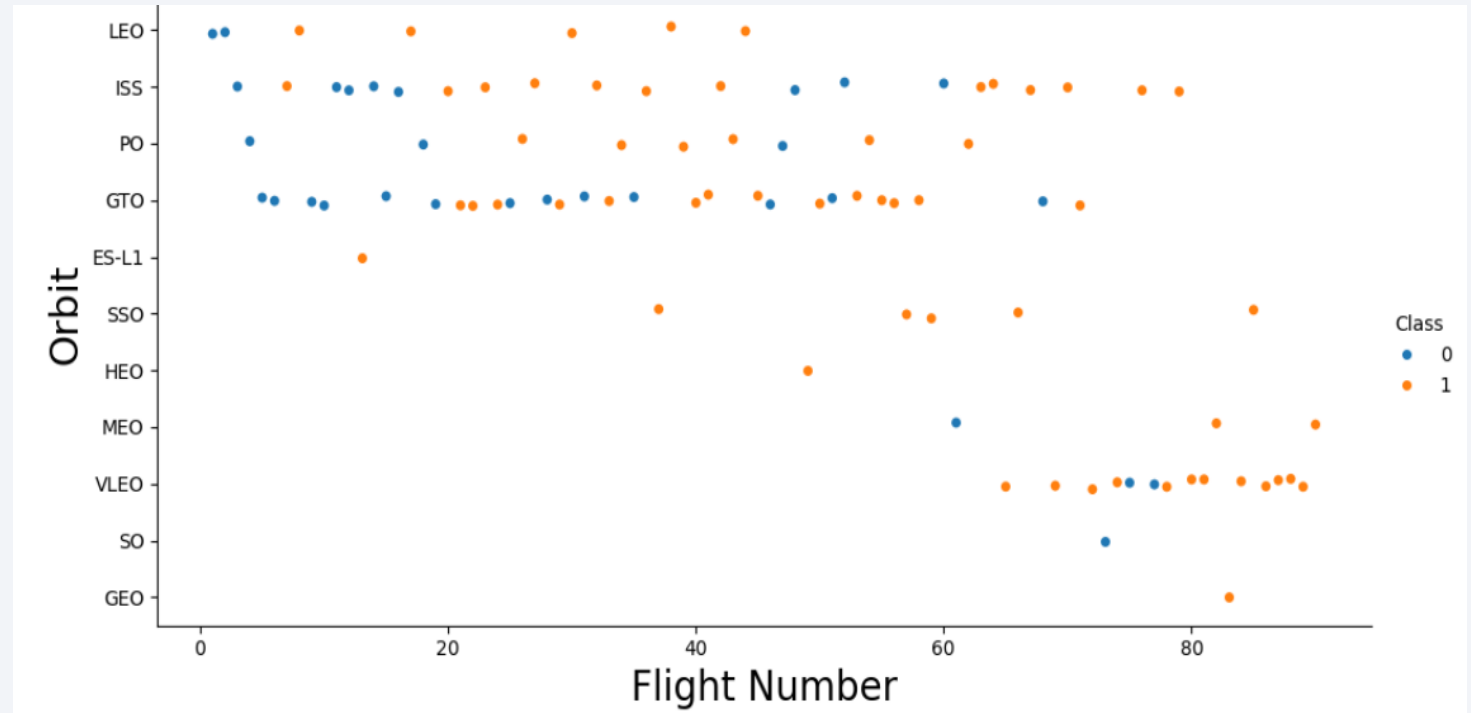
# Success Rate vs. Orbit Type

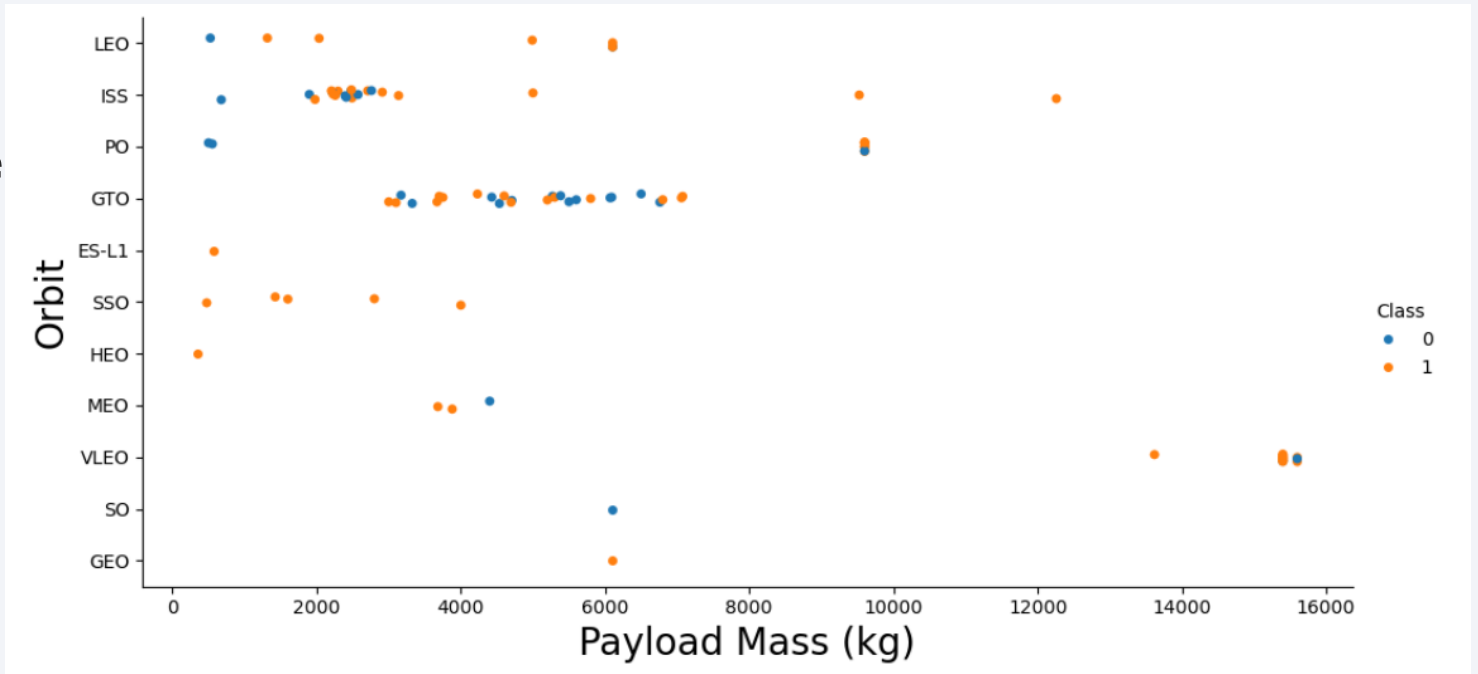- SSO, HEO, GEO and ES-L1 have highly successful rate.

# Flight Number vs. Orbit Type

- LEO orbit the Success appears related to the number of flights; on the other hand, there seems to be no relationship between flight number when in GTO orbit.
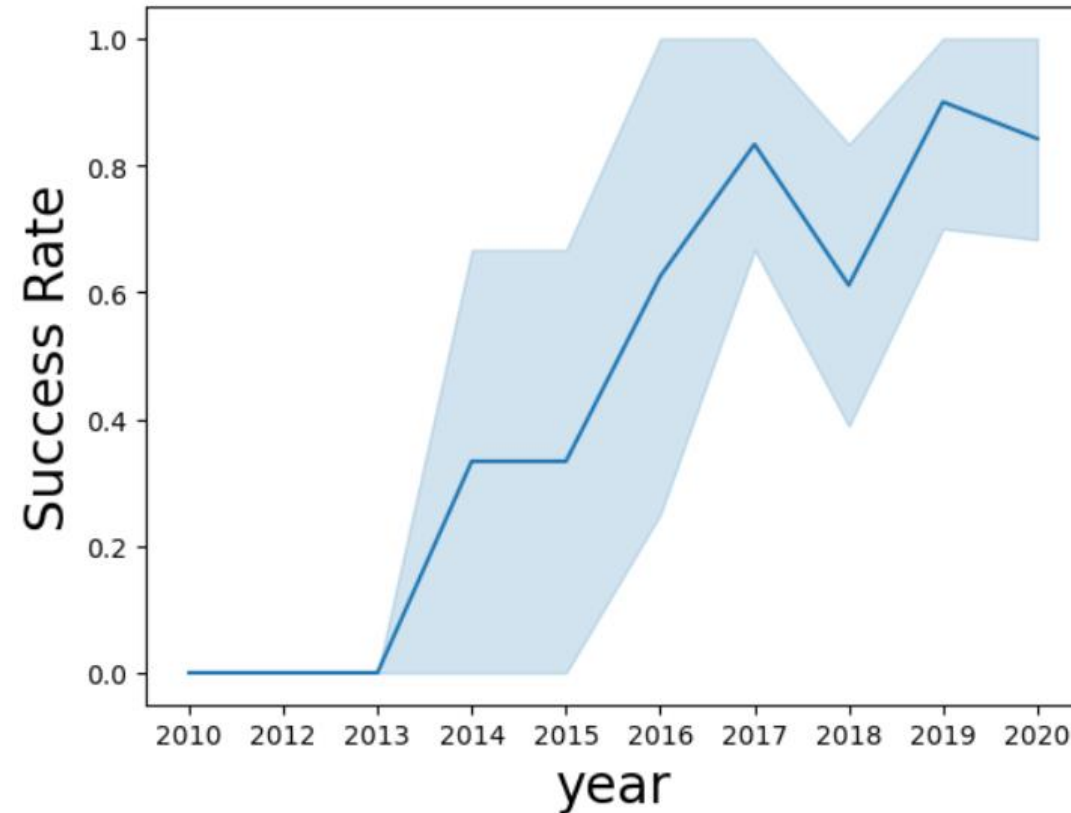
# Payload vs. Orbit Type

- With heavy payloads the successful landing or positive landing rate are more for Polar,LEO and ISS.

- However for GTO we cannot distinguish this well as both positive landing rate and negative landing(unsuccessful mission) are both there here.

# Launch Success Yearly Trend

- Line chart shows yearly average success rate.

- The success rate since 2013 kept increasing till 2020



you can observe that the success rate since 2013 kept increasing till 2020

# All Launch Site Names

- We used the key word **DISTINCT** to show only unique launch sites from the SpaceX data.

Display the names of the unique launch sites in the space mission

```
In [10]:    task_1 = '''
                SELECT DISTINCT LaunchSite
                FROM SpaceX
            '''
            create_pandas_df(task_1, database=conn)
```

Out[10]:
|   | launchsite |
|---|---|
| 0 | KSC LC-39A |
| 1 | CCAFS LC-40 |
| 2 | CCAFS SLC-40 |
| 3 | VAFB SLC-4E |

# Launch Site Names Begin with 'CCA'

- We use the LIMIT to limit the result to 5. We use WHERE to specify the condition of launchsite and LIKE for the 'CCA%'

Display 5 records where launch sites begin with the string 'CCA'

```
In [11]:   task_2 = '''
              SELECT *
              FROM SpaceX
              WHERE LaunchSite LIKE 'CCA%'
              LIMIT 5
           '''
           create_pandas_df(task_2, database=conn)
```

Out[11]:

| | date | time | boosterversion | launchsite | payload | payloadmasskg | orbit | customer | missionoutcome | landingoutcome |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 2010-04-06 | 18:45:00 | F9 v1.0 B0003 | CCAFS LC-40 | Dragon Spacecraft Qualification Unit | 0 | LEO | SpaceX | Success | Failure (parachute) |
| 1 | 2010-08-12 | 15:43:00 | F9 v1.0 B0004 | CCAFS LC-40 | Dragon demo flight C1, two CubeSats, barrel of... | 0 | LEO (ISS) | NASA (COTS) NRO | Success | Failure (parachute) |
| 2 | 2012-05-22 | 07:44:00 | F9 v1.0 B0005 | CCAFS LC-40 | Dragon demo flight C2 | 525 | LEO (ISS) | NASA (COTS) | Success | No attempt |
| 3 | 2012-08-10 | 00:35:00 | F9 v1.0 B0006 | CCAFS LC-40 | SpaceX CRS-1 | 500 | LEO (ISS) | NASA (CRS) | Success | No attempt |
| 4 | 2013-01-03 | 15:10:00 | F9 v1.0 B0007 | CCAFS LC-40 | SpaceX CRS-2 | 677 | LEO (ISS) | NASA (CRS) | Success | No attempt |

# Total Payload Mass

- We used SUM to specify the addition of all the PAYLOAD_MASS_KG and we used FROM to specify where the data is being gotten and WHERE CUSTOMER is Nasa.

## Task 3

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

```
In [10]:  %sql SELECT SUM(PAYLOAD_MASS__KG_) \
              FROM SPACEXTBL \
              WHERE CUSTOMER = 'NASA (CRS)';
```

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

Out[10]:  **SUM(PAYLOAD_MASS__KG_)**

45596

# Average Payload Mass by F9 v1.1

- The average payload mass carried by the booster version F9 v1.1 is 2928.4

## Task 4

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

```
In [11]:   %sql SELECT AVG(PAYLOAD_MASS__KG_) \
               FROM SPACEXTBL \
               WHERE BOOSTER_VERSION = 'F9 v1.1';

           * sqlite:///my_data1.db
           Done.

Out[11]:   AVG(PAYLOAD_MASS__KG_)

                      2928.4
```

# First Successful Ground Landing Date

- The date of the first successful landing outcome on ground pad is 2015-12-22

```
In [13]:   %sql SELECT MIN(DATE) \
           FROM SPACEXTBL \
           WHERE LANDING_OUTCOME = 'Success (ground pad)'

            * sqlite:///my_data1.db
           Done.

Out[13]:   MIN(DATE)

           2015-12-22
```

# Successful Drone Ship Landing with Payload between 4000 and 6000

- The names of boosters which have successfully landed on drone ship with payload mass greater than 4000 but less than 6000.

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

```
In [14]:   %sql SELECT PAYLOAD \
           FROM SPACEXTBL \
           WHERE LANDING_OUTCOME = 'Success (drone ship)' \
           AND PAYLOAD_MASS__KG_ BETWEEN 4000 AND 6000;
```

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

Out[14]:

| Payload |
| --- |
| JCSAT-14 |
| JCSAT-16 |
| SES-10 |
| SES-11 / EchoStar 105 |

# Total Number of Successful and Failure Mission Outcomes

- The total number of successful and failure mission outcomes.

List the total number of successful and failure mission outcomes

```
In [15]:  %sql SELECT MISSION_OUTCOME, COUNT(*) as total_number \
          FROM SPACEXTBL \
          GROUP BY MISSION_OUTCOME;
```

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

Out[15]:

| Mission_Outcome | total_number |
|---|---|
| Failure (in flight) | 1 |
| Success | 98 |
| Success | 1 |
| Success (payload status unclear) | 1 |

# Boosters Carried Maximum Payload

- A list of the names of the booster which have carried the maximum payload mass

List the names of the booster_versions which have carried the maximum payload mass. Use a subquery

```
In [17]:   %sql SELECT BOOSTER_VERSION \
           FROM SPACEXTBL \
           WHERE PAYLOAD_MASS__KG_ = (SELECT MAX(PAYLOAD_MASS__KG_) FROM SPACEXTBL);
```

* sqlite:///my_data1.db
Done.

Out[17]:

| Booster_Version |
| --- |
| 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 |

# 2015 Launch Records

- A list of the failed landing outcomes in drone ship, their booster versions, and launch site names for in year 2015

List the failed landing_outcomes in drone ship, their booster versions, and launch site names for in year 2015

```
In [18]:  task_9 = '''
                SELECT BoosterVersion, LaunchSite, LandingOutcome
                FROM SpaceX
                WHERE LandingOutcome LIKE 'Failure (drone ship)'
                    AND Date BETWEEN '2015-01-01' AND '2015-12-31'
                '''
          create_pandas_df(task_9, database=conn)
```

| | boosterversion | launchsite | landingoutcome |
|---|---|---|---|
| 0 | F9 v1.1 B1012 | CCAFS LC-40 | Failure (drone ship) |
| 1 | F9 v1.1 B1015 | CCAFS LC-40 | Failure (drone ship) |

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

- Rank of 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.

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

```
In [19]:    task_10 = '''
            SELECT LandingOutcome, COUNT(LandingOutcome)
            FROM SpaceX
            WHERE DATE BETWEEN '2010-06-04' AND '2017-03-20'
            GROUP BY LandingOutcome
            ORDER BY COUNT(LandingOutcome) DESC
            '''
            create_pandas_df(task_10, database=conn)
```

Out[19]:

| | landingoutcome | count |
|---|---|---|
| 0 | No attempt | 10 |
| 1 | Success (drone ship) | 6 |
| 2 | Failure (drone ship) | 5 |
| 3 | Success (ground pad) | 5 |
| 4 | Controlled (ocean) | 3 |
| 5 | Uncontrolled (ocean) | 2 |
| 6 | Precluded (drone ship) | 1 |
| 7 | Failure (parachute) | 1 |

# Launch Sites Proximities Analysis

# All Launch Sites

# Launch Sites wit Color Labels



**Florida Launch Sites**

*Green Marker* shows successful Launches and *Red Marker* shows Failures

**California Launch Site**

# Launch Site Proximity to Landmarks



Distance to Railway Station

Distance to closest Highway

Distance to coast

Distance to Coastline

Distance to City

- 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

Section 4

# Build a Dashboard
# with Plotly Dash

# Pie Chart of Success Launches By All Sites



Total Success Launches By all sites

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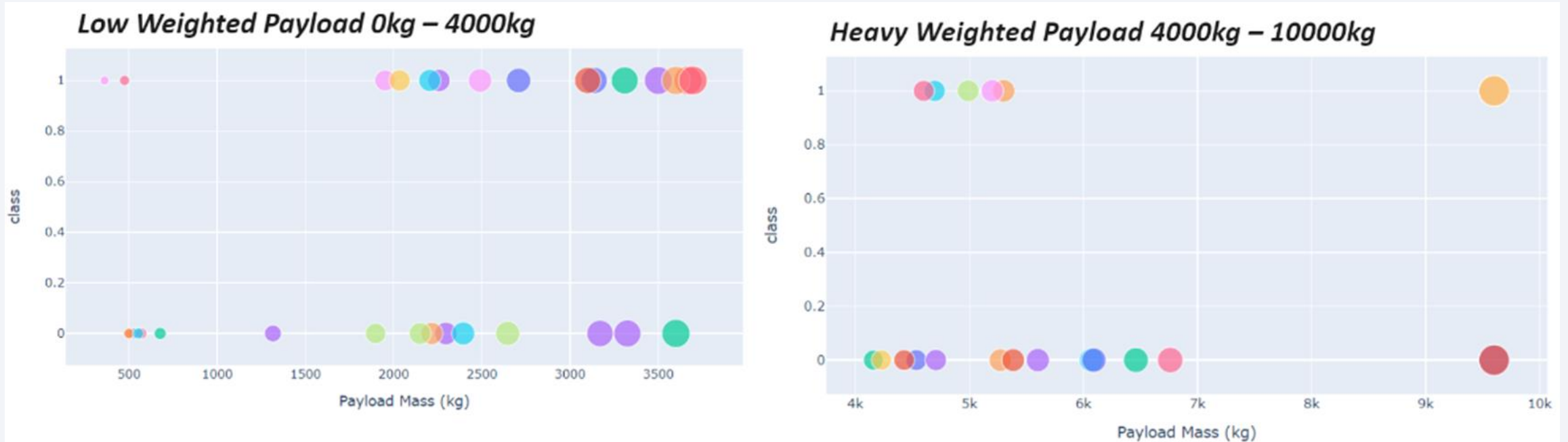CCAFS SLC-40

41.7%
29.2%
16.7%
12.5%

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

# Pie Chart Showing Highest Lauch Success Rate on KSC LC-39A



KSC LC-39A achieved a 76.9% success rate while getting a 23.1% failure rate

# Scatter plot of Payload vs Launch Outcome for all sites



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_)
```
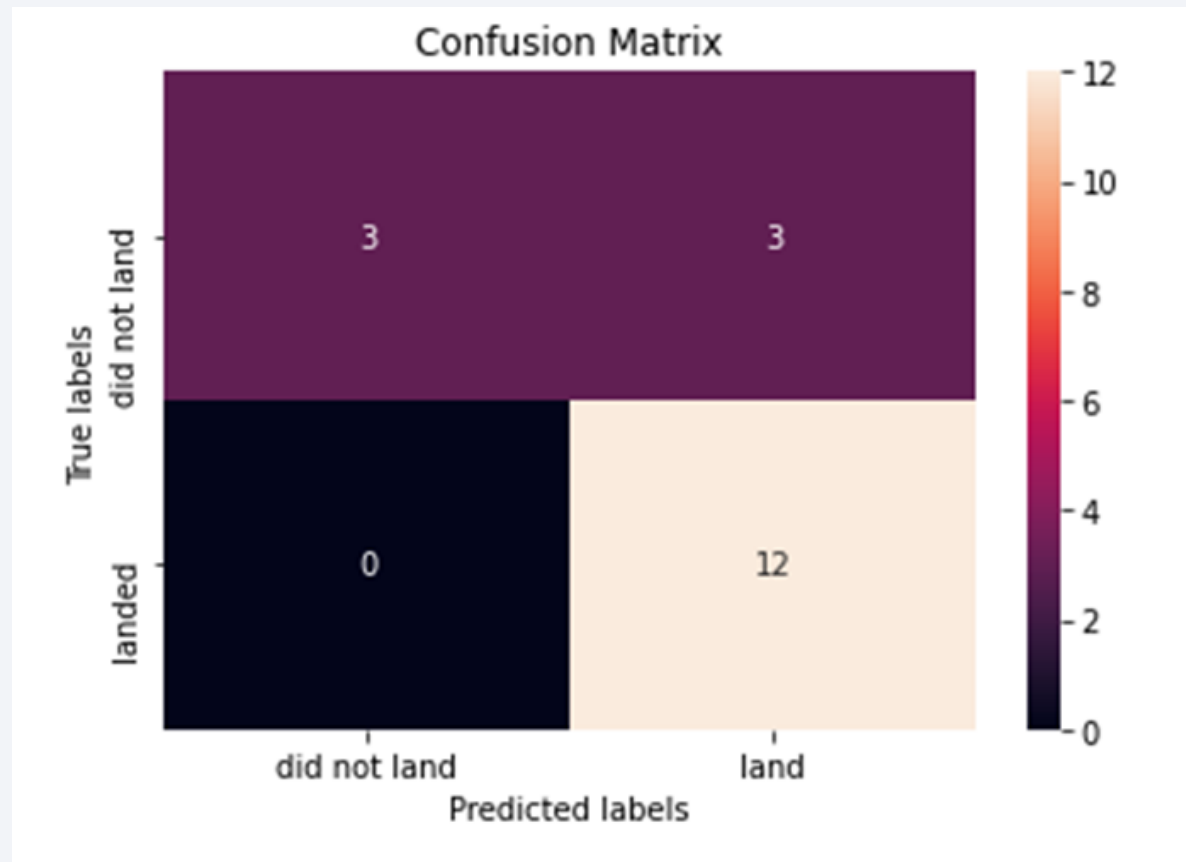
```
Best model is DecisionTree with a score of 0.8732142857142856
Best params is : {'criterion': 'gini', 'max_depth': 6, 'max_features': 'auto', 'min_samples_leaf': 2, 'min_samples_split': 5, 'splitter': 'random'}
```

# Confusion Matrix

- The confusion matrix of the best performing model

# Conclusions

- We found that the greater the flight number at a launch site, the greater the success rate.

- The success rate of the launches started to increase in 2013 to 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 with a score of 0.87321428571.

Thank you!