# Capstone Introduction and Understanding the Datasets

# Outline

- Interactive visual analytics and Dashboard

- Wrangling the data

- Predictive Analysis (Classification)

- Data-Driven Insights

▶ 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

# Outline

- Executive Summary

- Introduction

- Methodology

- Results

- Conclusion

- Appendix

# 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. This goal of the project is to create a machine learning pipeline to predict if the first stage will land successfully.

- Problems you want to find answers

  - What factors determine if the rocket will land successfully?

  - The interaction amongst various features that determine the success rate of a successful landing.

  - What operating conditions needs to be in place to ensure a successful landing program.

Section 1

# Methodology

# Methodology

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

- 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

► The data was collected using various methods

- Data collection was done using get request to the SpaceX API.

- Next, we decoded the response content as a Json using .json() function call and turn 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.

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

► The link to the notebook is

► https://github.com/devipriyak/CapStone-ML-Project/blob/main/DataWrangling.ipynb

# Importing Libraries

```
[1] !pip3 install beautifulsoup4
    !pip3 install requests
```

```
Requirement already satisfied: beautifulsoup4 in /usr/local/lib/python3.10/dist-packages (4.12.3)
Requirement already satisfied: soupsieve>1.2 in /usr/local/lib/python3.10/dist-packages (from beautifulsoup4) (2.5)
Requirement already satisfied: requests in /usr/local/lib/python3.10/dist-packages (2.31.0)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from requests) (3.3.2)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests) (3.7)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.10/dist-packages (from requests) (2.0.7)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests) (2024.2.2)
```

```
[2] import sys
    import requests
    from bs4 import BeautifulSoup
    import re
    import unicodedata
    import pandas as pd
```

# Accessing data

```
[4]  static_url = "https://en.wikipedia.org/w/index.php?title=List_of_Falcon_9_and_Falcon_Heavy_launches&oldid=1027686922"
```

```
[5]  # use requests.get() method with the provided static_url

     data  = requests.get(static_url).text
```

```
[6]  # Use BeautifulSoup() to create a BeautifulSoup object from a response text content
     soup = BeautifulSoup(data, 'html5lib')
```

```
[7]  ## Use soup.title attribute
     print(soup.title)
```

```
    <title>List of Falcon 9 and Falcon Heavy launches - Wikipedia</title>
```

```
[8]  html_tables1 = soup.find_all('table')
```

```
[9]  first_launch_table = html_tables1[2]
     print(first_launch_table)
```

✓ 0s    completed at 9:27 PM

# Printing Details

```
[7]  ## Use soup.title attribute
     print(soup.title)

     <title>List of Falcon 9 and Falcon Heavy launches - Wikipedia</title>

[8]  html_tables1 = soup.find_all('table')

[9]  first_launch_table = html_tables1[2]
     print(first_launch_table)
```

# Print Feature Details

```
print(column_names)
```

```
['Flight No.', 'Date and time ( )', 'Launch site', 'Payload', 'Payload mass', 'Orbit', 'Customer', 'Launch outcome']
```

# Code in Git Hub

- https://github.com/devipriyak/CapStone-ML-Project/blob/main/webscrap.ipynb

# Data Wrangling

- ▶ Identify duplicate values in the dataset.
- ▶ Remove duplicate values from the dataset.
- ▶ Identify missing values in the dataset.
- ▶ Impute the missing values in the dataset.
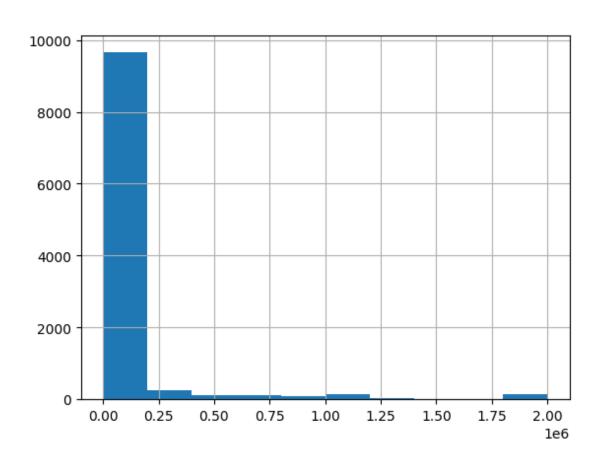- ▶ Normalize data in the dataset.

# Github Id

- https://github.com/devipriyak/CapStone-ML-Project/blob/main/DataWrangling.ipynb
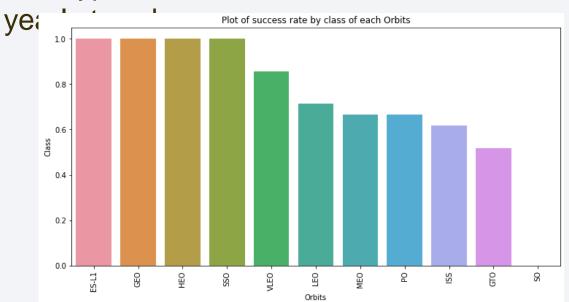
# Module 2 - Exploratory Data Analysis (EDA)

▶ In this module, you will collect data on the Falcon 9 first-stage landings. You will use a RESTful API and web scraping. You will also convert the data into a dataframe and then perform some data wrangling.

▶ **Learning Objectives**

▶ Load a dataset into a database.

▶ Write and execute SQL queries to select and sort data.

▶ Write Python code to conduct exploratory data analysis by manipulating data in a Pandas data frame.

▶ Visualize the data and extract meaningful patterns to guide the modeling process.

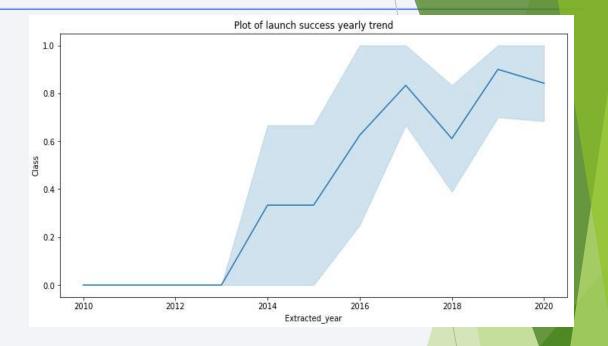▶ Create scatter plots and bar charts to analyze data in a Pandas data frame.

# Visualizing data

# EDA with Data Visualization https://github.com/devipriyak/CapStone-ML-Project/blob/main/Data_Explo.ipynb

▶ 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





https://github.com/devipriyak/CapStone-ML-Project/blob/main/Data_Explo.ipynb

# EDA with SQL

▶ We loaded the SpaceX dataset into a PostgreSQL database without leaving the jupyter notebook.

▶ We applied EDA with SQL to get insight from the data. We wrote queries to find out for instance:

- 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.

The link to the notebook

https://github.com/devipriyak/CapStone-ML-Project/blob/main/Data_Explo.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.

- https://github.com/devipriyak/CapStone-ML-Project/blob/main/Interactive_Visual_Analytics_with_Folium.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.

▶ The link to the notebook is https://github.com/devipriyak/CapStone-ML-Project/

# 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.

▶ The link to the notebook https://github.com/devipriyak/CapStone-ML-Project/blob/main/Machine_Learning_Prediction.ipynb
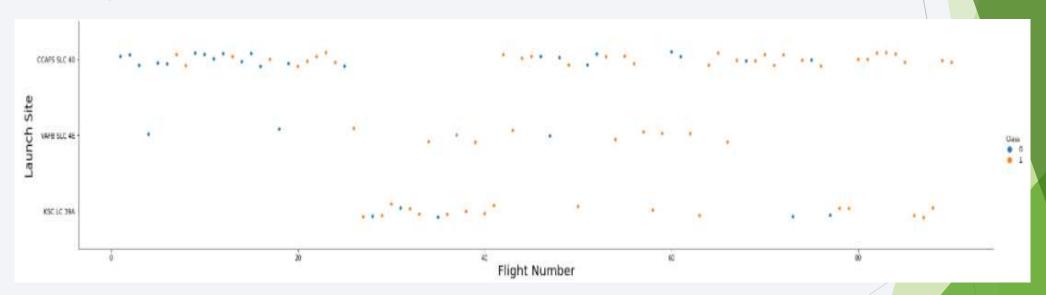
# Results

- Exploratory data analysis results

- Interactive analytics demo in screenshots

- Predictive analysis results

Section 2

# Insights drawn from EDA

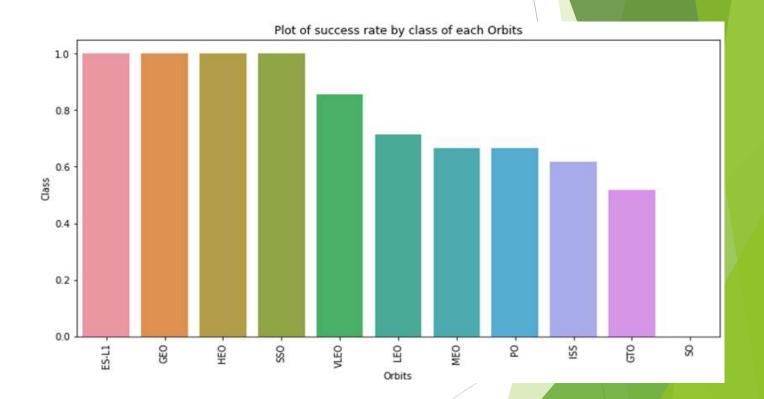# Flight Number vs. Launch Site

► From the plot, we found that the larger the flight amount at a launch site, the greater the success rate at a launch site.
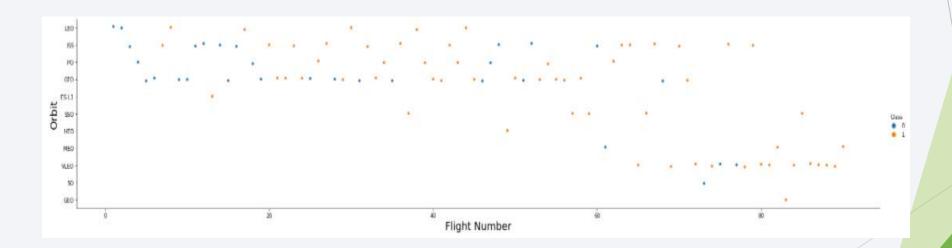
# Success Rate vs. Orbit Type

▶ 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
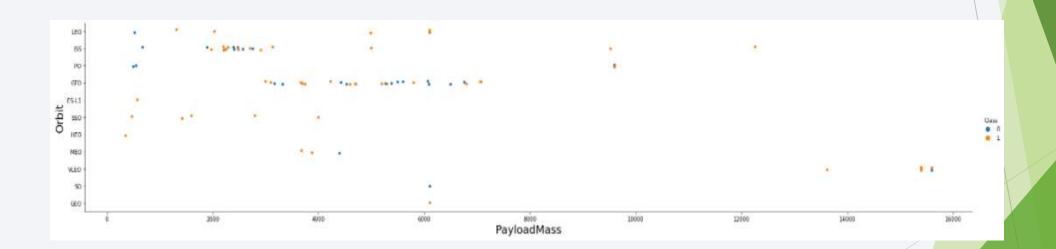
# Flight Number vs. Orbit Type

► 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.
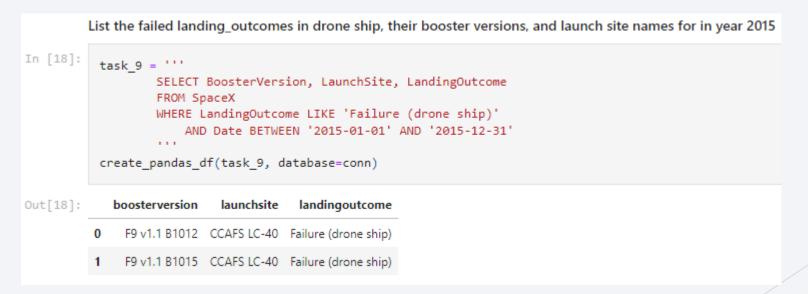
# Boosters Carried Maximum Payload

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

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

```
In [17]:   task_8 = '''
               SELECT BoosterVersion, PayloadMassKG
               FROM SpaceX
               WHERE PayloadMassKG = (
                                       SELECT MAX(PayloadMassKG)
                                       FROM SpaceX
                                       )
               ORDER BY BoosterVersion
               '''

           create_pandas_df(task_8, database=conn)
```

Out[17]:

|    | boosterversion | payloadmasskg |
|----|----------------|---------------|
| 0  | F9 B5 B1048.4  | 15600         |
| 1  | F9 B5 B1048.5  | 15600         |
| 2  | F9 B5 B1049.4  | 15600         |
| 3  | F9 B5 B1049.5  | 15600         |
| 4  | F9 B5 B1049.7  | 15600         |
| 5  | F9 B5 B1051.3  | 15600         |
| 6  | F9 B5 B1051.4  | 15600         |
| 7  | F9 B5 B1051.6  | 15600         |
| 8  | F9 B5 B1056.4  | 15600         |
| 9  | F9 B5 B1058.3  | 15600         |
| 10 | F9 B5 B1060.2  | 15600         |
| 11 | F9 B5 B1060.3  | 15600         |

# 2015 Launch Records

▶ 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

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)
```

Out[18]:

|   | 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 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 |

▶ 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 the grouped landing outcome in descending order.

Section 4

# Launch Sites Proximities Analysis

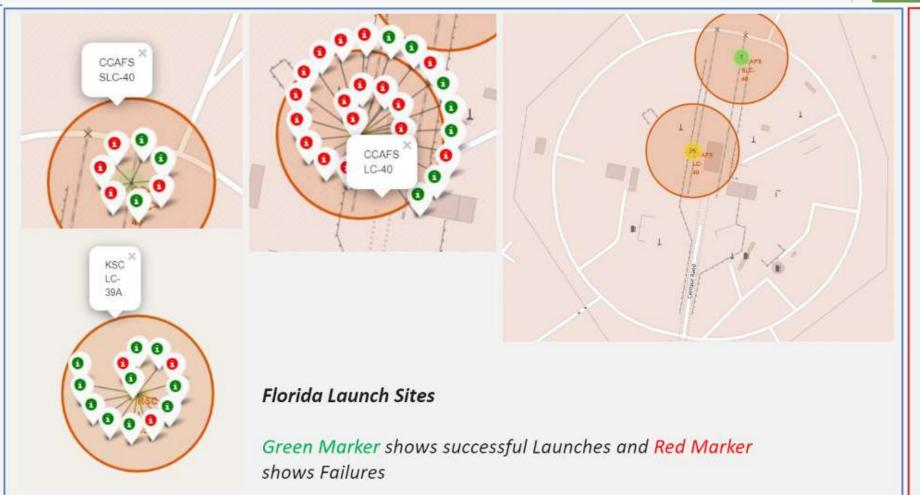# All launch sites global map markers



We can see that the SpaceX launch sites are in the United States of America coasts. Florida and California
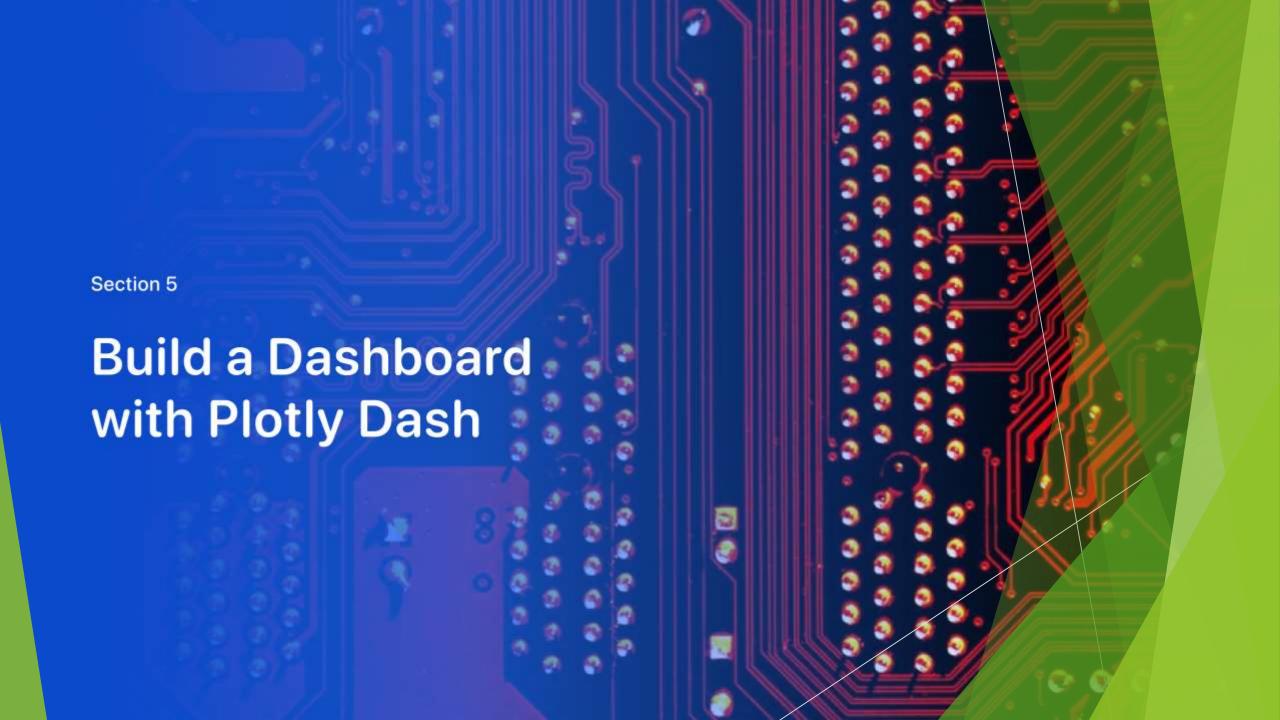
# Markers showing launch sites with color labels



Florida Launch Sites

Green Marker shows successful Launches and Red Marker shows Failures

California Launch Site

# Launch Site distance to landmarks



Distance to Railway Station

Distance to closest Highway

Distance to Coastline

Distance to City

Distance to coast

- 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 5

# Build a Dashboard
# with Plotly Dash

# Pie chart showing the success percentage achieved by each launch site



Total Success Launches By all sites
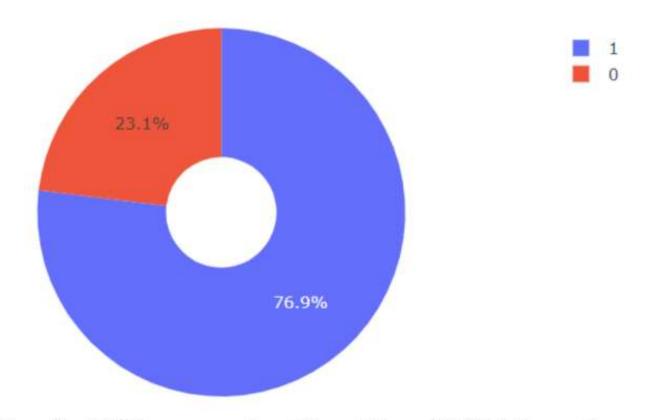
29.2%

41.7%

16.7%
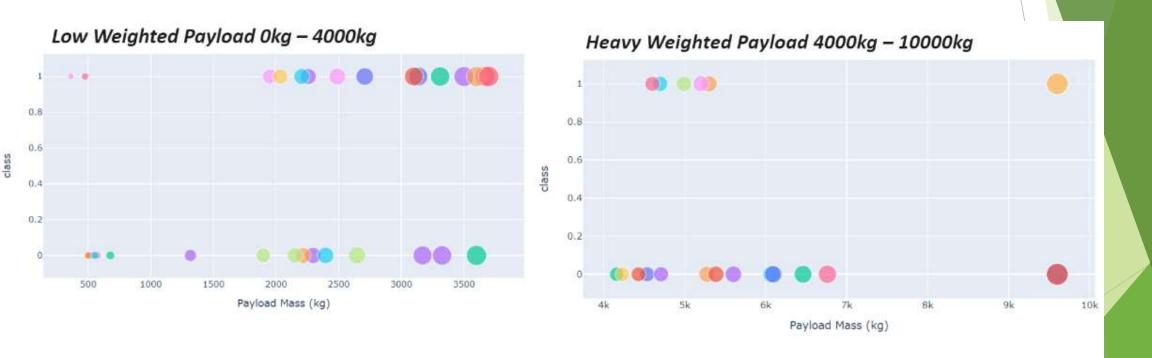
12.5%

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

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

# Pie chart showing the Launch site with the highest launch success ratio



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, with different payload selected in the range slider



We can see the success rates for low weighted payloads is higher than the heavy weighted payloads
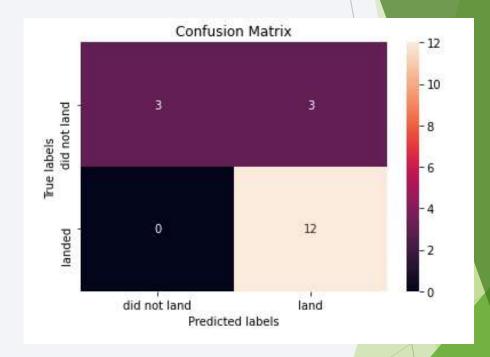
Section 6

# Predictive Analysis (Classification)

# Classification Accuracy

▶ The decision tree classifier is the model with the highest 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 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

We can conclude that:

▶ The larger the flight amount at a launch site, the greater the success rate at a launch site.

▶ Launch success rate started to increase in 2013 till 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.

Thank you!