Space X Falcon 9 Landing Analysis

## IBM Data Science Capstone Project

## OUTLINE

**Executive Summary** 

Introduction

Methodology

Result

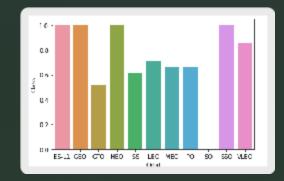
Conclusion

Appendix

#### Summary of Methodologies:

This project follows these steps:

- Data Collection
- Data Wrangling
- Exploratory Data Analysis
- Interactive Visual Analytics
- Predictive Analysis (Classification)

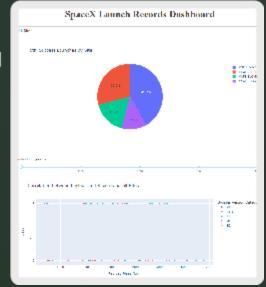


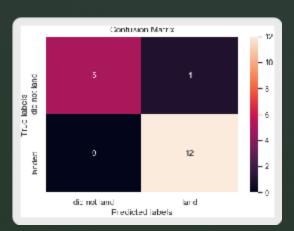


#### Summary of Results:

This project produced the following outputs and visualizations:

- 1. Exploratory Data Analysis (EDA) results
- Geospatial analytics
- Interactive dashboard
- 4. Predictive analysis of classification models





## Introduction

- SpaceX launches Falcon 9 rockets at a cost of around \$62m. This is considerably cheaper than other providers (which usually cost upwards of \$165m), and much of the savings are because SpaceX can land, and then re-use the first stage of the rocket.
- If we can make predictions on whether the first stage will land, we can determine the cost of a launch, and use this information to assess whether or not an alternate company should bid and SpaceX for a rocket launch.
- This project will ultimately predict if the Space X Falcon
   9 first stage will land successfully.

## Methodology summary

#### 1. Data Collection

- Making GET requests to the SpaceX REST API
- Web Scraping

#### 2. Data Wrangling

- Using the .fillna() method to remove NaN values
- Using the .value\_counts() method to determine the following:
  - Number of launches on each site
  - Number and occurrence of each orbit
- Number and occurrence of mission outcome per orbit type
- Creating a landing outcome label that shows the following:
- 0 when the booster did not land successfully
- 1 when the booster did land successfully

#### 3. Exploratory Data Analysis

- Using SQL queries to manipulate and evaluate the SpaceX dataset
- Using Pandas and Matplotlib to visualize relationships between variables, and determine patterns

#### 4. Interactive Visual Analytics

- Geospatial analytics using Folium
- Creating an interactive dashboard using Plotly Dash

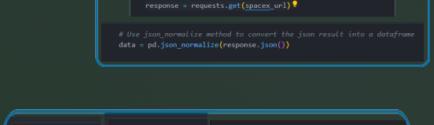
#### 5. Data Modelling and Evaluation

- Using Scikit-Learn to:
- Pre-process (standardize) the data
- Split the data into training and testing data using train\_test\_split
- Train different classification models
- Find hyperparameters using GridSearchCV
- Plotting confusion matrices for each classification model
- Assessing the accuracy of each classification model

## DATA COLLECTION – space x REST api

Using the SpaceX API to retrieve data about launches, including information about the rocket used, payload delivered, launch specifications, landing specifications, and landing outcome.

- Make a GET response to the SpaceX REST API
  - Convert the response to a .json file then to a Pandas DataFrame
- Use custom logic to clean the data (see Appendix)
  - Define lists for data to be stored in
  - Call custom functions (see Appendix) to retrieve data and fill the lists
  - Use these lists as values in a dictionary and construct the dataset
- Create a Pandas DataFrame from the constructed dictionary dataset
- Filter the DataFrame to only include Falcon 9 launches
  - Reset the FlightNumber column
  - Replace missing values of PayloadMass with the mean PayloadMass value



spacex\_url="https://api.spacexdata.com/v4/launches/past" ?

```
# Call getBoosterVersion | launch_dict = ('FlightNumber': list(data['flight number'])
                                                         'Date': list(data['date']),
BoosterVersion = []
                                                         'BoosterVersion':BoosterVersion,
PayloadMass = |
                                                         'PayloadMass':PayloadMass,
                                                         'Orbit':Orbit,
LaunchSite = []
                                                         'LaunchSite':LaunchSite,
Outcome | []
                             getLaunchSite(data)
                                                         'Outcome':Outcome,
                                                         'Flights':Flights,
                                                         'GridFins': GridFins,
                                                         "Reused": Reused,
                                                         'legs':legs,
LandingPad = []
                             getPayloadData(data)
                                                         'LandingPad':LandingPad,
Block = []
                                                         'Block':Block,
ReusedCount = []
                                                         'ReusedCount':ReusedCount,
                                                         'Serial':Serial,
Longitude = []
                                                         'Longitude': Longitude,
Latitude = []
                                                         'Latitude': Latitude
```

## Data Manipulation

#### Context:

- The SpaceX dataset contains several Space X launch facilities, and each location is in the LaunchSite column.
- Each launch aims to a dedicated orbit, and some of the common orbit types are shown in the figure below. The orbit type is in the Orbit column.

# 35768 km 10000 km MEO 1000 km

#### Initial Data Exploration:

- Using the .value counts() method to determine the following:
  - Number of launches on each site
  - Number and occurrence of each orbit
  - Number and occurrence of landing outcome per orbit type

## DATA MANIPULATION/WRANGLING – PANDAS

#### Context:

- The landing outcome is shown in the Outcome column:
  - True Ocean the mission outcome was successfully landed to a specific region of the ocean
- False Ocean the mission outcome was unsuccessfully landed to a specific region of the ocean.
- True RTLS the mission outcome was successfully landed to a ground pad
- False RTLS the mission outcome was unsuccessfully landed to a ground pad.
- True ASDS the mission outcome was successfully landed to a drone ship
- False ASDS the mission outcome was unsuccessfully landed to a drone ship.
- None ASDS and None None these represent a failure to land.

#### Data Wrangling:

- To determine whether a booster will successfully land, it is best to have a binary column, i.e., where the value is 1 or 0, representing the success of the landing.
- This is done by:
  - Defining a set of unsuccessful (bad) outcomes, bad\_outcome
  - 2. Creating a list, landing\_class, where the element is 0 if the corresponding row in Outcome is in the set bad\_outcome, otherwise, it's 1.
  - 3. Create a Class column that contains the values from the list landing class
- Export the DataFrame as a .csv file.

```
bed outcomes=set(landing outcomes.keys()||1,3,5,6,7||)
bed_outcomes

['Felse ASDS', 'Felse Ocean', 'Felse RTLS', 'None ASDS', 'None None']
```

```
# Landing_class = 0 if had_outcome
# landing_class = []

for outcome in df['Outcome']:
    if outcome in had_outcomes:
        landing_class.append(0)
    else:
        landing_class.append(1)
```

```
df.to_csv("dataset_part\_2.csv", index=False)
```

## Exploratory data analysis (eda) – SQL

To gather some information about the dataset, some SQL queries were performed.

#### The SQL gueries performed on the data set were used to:

- 1. Display the names of the unique launch sites in the space mission
- 2. Display 5 records where launch sites begin with the string 'CCA'
- 3. Display the total payload mass carried by boosters launched by NASA (CRS)
- 4. Display the average payload mass carried by booster version F9 v1.1
- 5. List the date when the first successful landing outcome on a ground pad was achieved
- 6. List the names of the boosters which had success on a drone ship and a payload mass between 4000 and 6000 kg
- 7. List the total number of successful and failed mission outcomes
- B. List the names of the booster versions which have carried the maximum payload mass
- 9. List the failed landing outcomes on drone ships, their booster versions, and launch site names for 2015
- 10. Rank 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

## Geospatial analysis – folium

The following steps were taken to visualize the launch data on an interactive map:

#### 1. Mark all launch sites on a map

- Initialise the map using a Folium Map object
- Add a folium.Circle and folium.Marker for each launch site on the launch map

#### 2. Mark the success/failed launches for each site on a map

- As many launches have the same coordinates, it makes sense to cluster them together.
- Before clustering them, assign a marker colour of successful (class = 1) as green, and failed (class = 0) as red.
- To put the launches into clusters, for each launch, add a folium.Marker to the MarkerCluster() object.
- Create an icon as a text label, assigning the icon\_color as the marker\_colour determined previously.

#### 3. Calculate the distances between a launch site to its proximities

- To explore the proximities of launch sites, calculations of distances between points can be made using the Lat and Long values.
- After marking a point using the Lat and Long values, create a folium. Marker object to show the distance.
- To display the distance line between two points, draw a folium. PolyLine and add this to the map.

## Interactive dashboard – plotly dash

The following plots were added to a Plotly Dash dashboard to have an interactive visualisation of the data:

- 1. Pie chart (px.pie()) showing the total successful launches per site
- This makes it clear to see which sites are most successful
- The chart could also be filtered (using a dcc.Dropdown () object) to see the success/failure ratio for an individual site
- 2. Scatter graph (px.scatter()) to show the correlation between outcome (success or not) and payload mass (kg)
- This could be filtered (using a RangeSlider() object) by ranges of payload masses
- It could also be filtered by booster version

## Predictive Analysis - Classification

### Model Development 🗾





### Model Evaluation [







- To prepare the dataset for model development:
  - Load dataset
  - Perform necessary data transformations (standardise and pre-process)
  - Split data into training and test data sets, using train\_test\_split()
  - Decide which type of machine learning algorithms are most appropriate
- For each chosen algorithm:
  - Create a GridSearchCV object and a dictionary of parameters
  - Fit the object to the parameters
  - Use the training data set to train the model

- For each chosen algorithm:
  - Using the output GridSearchCV object:
    - Check the tuned hyperparameters (best\_params\_)
    - Check the accuracy (score and best\_score\_)
  - Plot and examine the Confusion Matrix

- Review the accuracy scores for all chosen algorithms
- The model with the highest accuracy score is determined as the best performing model

## results

Exploratory Data Analysis

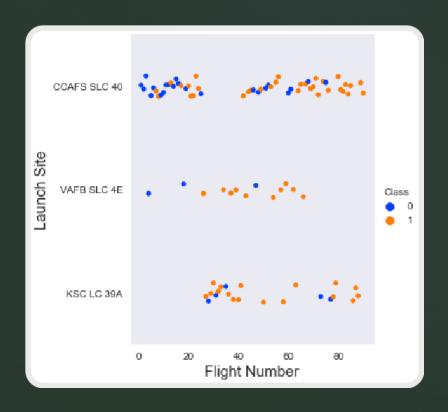
Interactive Analytics

Predictive Analysis

## Launch Site VS. FLIGHT NUMBER

The scatter plot of Launch Site vs. Flight Number shows that:

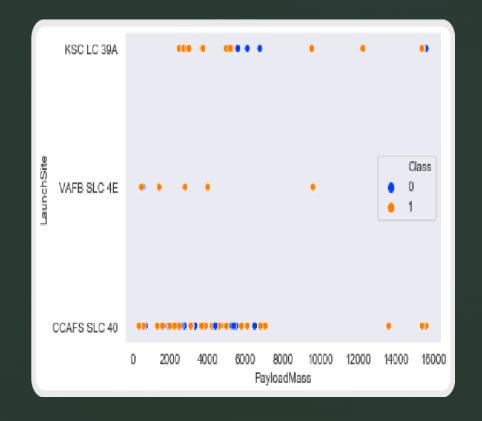
- As the number of flights increases, the rate of success at a launch site increases.
- Most of the early flights (flight numbers < 30) were launched from CCAFS SLC 40, and were generally unsuccessful.
- The flights from VAFB SLC 4E also show this trend, that earlier flights were less successful.
- No early flights were launched from KSC LC 39A, so the launches from this site are more successful.
- Above a flight number of around 30, there are significantly more successful landings (Class = 1).



## LAUNCH SITE vs. PAYLOAD MASS

The scatter plot of Launch Site vs. Payload Mass shows that:

- Above a payload mass of around 7000 kg, there are very few unsuccessful landings, but there is also far less data for these heavier launches.
- There is no clear correlation between payload mass and success rate for a given launch site.
- All sites launched a variety of payload masses, with most of the launches from CCAFS SLC 40 being comparatively lighter payloads (with some outliers).



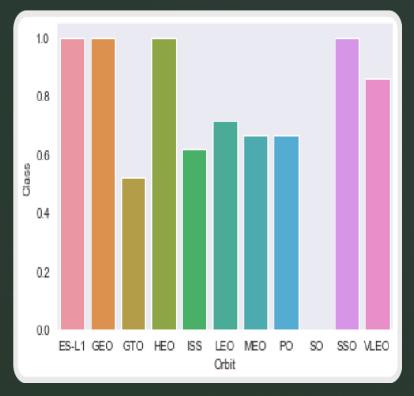
## Success Rate vs. Orbit Type

The bar chart of Success Rate vs. Orbit Type shows that the following orbits have the highest (100%) success rate:

- ES-L1 (Earth-Sun First Lagrangian Point)
- GEO (Geostationary Orbit)
- HEO (High Earth Orbit)
- SSO (Sun-synchronous Orbit)

The orbit with the lowest (0%) success rate is:

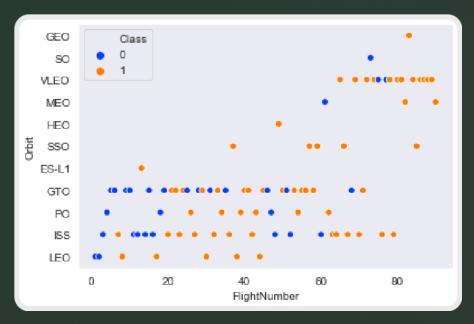
• SO (Heliocentric Orbit)



## Orbit Type vs. flight number

This scatter plot of Orbit Type vs. Flight number shows a few useful things that the previous plots did not, such as:

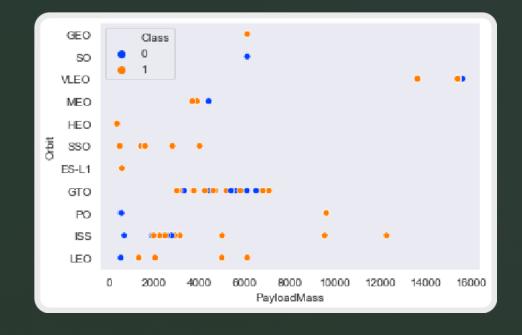
- The 100% success rate of GEO, HEO, and ES-L1 orbits can be explained by only having 1 flight into the respective orbits.
- The 100% success rate in SSO is more impressive, with 5 successful flights.
- There is little relationship between Flight Number and Success Rate for GTO.
- Generally, as Flight Number increases, the success rate increases. This is most extreme for LEO, where unsuccessful landings only occurred for the low flight numbers (early launches).



## ORBIT TYPE VS. PAYLOAD MASS

This scatter plot of Orbit Type vs. Payload Mass shows that:

- The following orbit types have more success with heavy payloads:
  - PO (although the number of data points is small)
  - ISS
  - LEO
- For GTO, the relationship between payload mass and success rate is unclear.
- VLEO (Very Low Earth Orbit) launches are associated with heavier payloads, which makes intuitive sense.

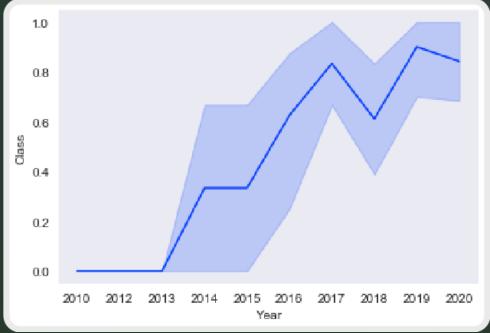


## Launch Success Yearly Trend

The line chart of yearly average success rate shows that:

- Between 2010 and 2013, all landings were unsuccessful (as the success rate is 0).
- After 2013, the success rate generally increased, despite small dips in 2018 and 2020.

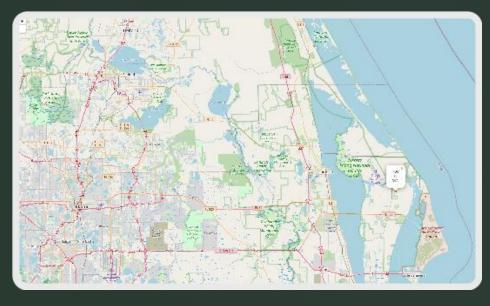
After 2016, there was always a greater than 50% chance of success.



LAUNCH SITES
 PROXIMITY ANALYSIS –
 FOLIUM INTERACTIVE
 MAP

## ALL LAUNCH SITES ON A MAP



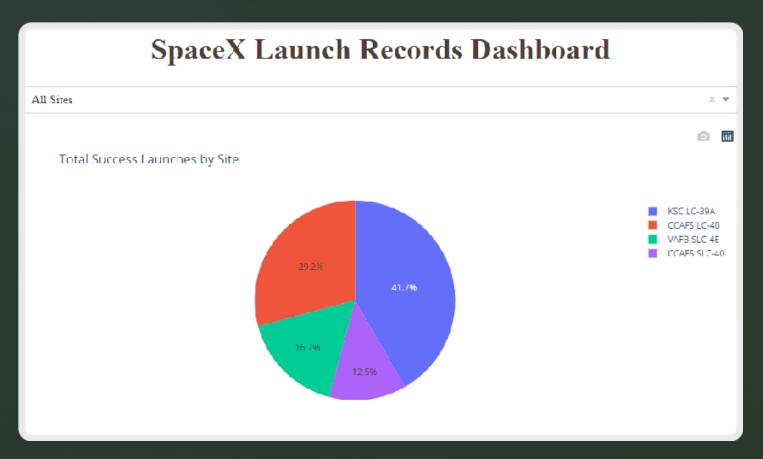






interactivedashboard - PlotlyDash

## launch success count for all sites



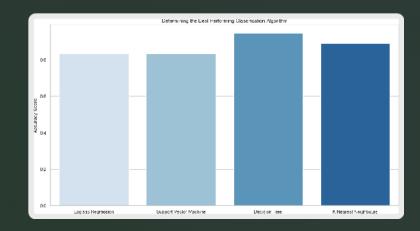
The launch site KSC LC-39 A had the most successful launches, with 41.7% of the total successful launches.

# PREDICTIVE ANALYSIS CLASSIFICATION

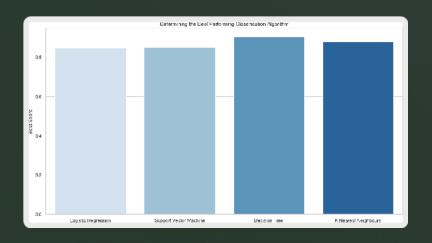
## CLASSIFICATION ACCURACY

Plotting the Accuracy Score and Best Score for each classification algorithm produces the following result:

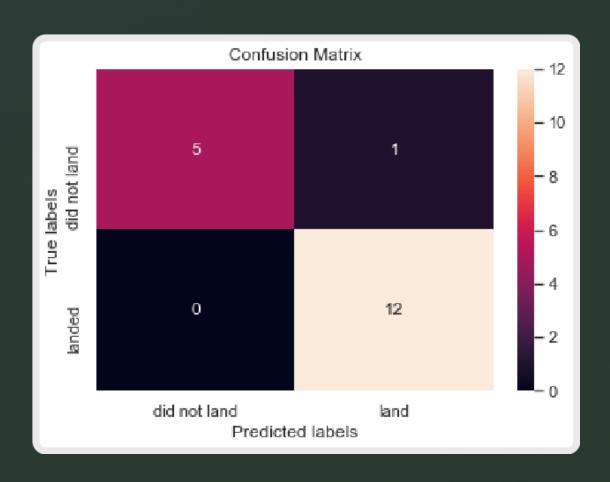
- The <u>Decision Tree</u> model has the highest classification accuracy
  - The Accuracy Score is 94.44%
  - The Best Score is 90.36%



	Algorithm	Accuracy Score	Best Score
ı	Logistic Regression	0.833333	0.846429
ı	Support Vector Machine	0.833333	0.848214
	Decision Tree	0.944444	0.903571
ı	K Nearest Neighbours	0.888889	0.876786



## **Confusion Matrix**



- As shown previously, best performing classification model is the Decision Tree model, with an accuracy of 94.44%.
- This is explained by the confusion matrix, which shows only 1 out of 18 total results classified incorrectly (a false positive, shown in the top-right corner).
- The other 17 results are correctly classified (5 did not land, 12 did land).

## CONCLUSIONS

## CONCLUSIONS

- As the number of flights increases, the rate of success at a launch site increases, with most early flights being unsuccessful. I.e. with more experience, the success rate increases.
  - Between 2010 and 2013, all landings were unsuccessful (as the success rate is 0).
  - After 2013, the success rate generally increased, despite small dips in 2018 and 2020.
  - After 2016, there was always a greater than 50% chance of success.
- Orbit types ES-L1, GEO, HEO, and SSO, have the highest (100%) success rate.
  - The 100% success rate of GEO, HEO, and ES-L1 orbits can be explained by only having 1 flight into the respective orbits.
  - The 100% success rate in SSO is more impressive, with 5 successful flights.
  - The orbit types PO, ISS, and LEO, have more success with heavy payloads:
  - VLEO (Very Low Earth Orbit) launches are associated with heavier payloads, which makes intuitive sense.
- The launch site KSC LC-39 A had the most successful launches, with 41.7% of the total successful launches, and also the highest rate of successful launches, with a 76.9% success rate.
- The success for massive payloads (over 4000kg) is lower than that for low payloads.
- The best performing classification model is the Decision Tree model, with an accuracy of 94.44%.

